

©Copyright 2020

Shiqing Yu

# Non-Gaussian Graphical Models: Estimation with Score Matching and Causal Discovery under Zero-Inflation

Shiqing Yu

A dissertation  
submitted in partial fulfillment of the  
requirements for the degree of

Doctor of Philosophy

University of Washington

2020

Reading Committee:

Mathias Drton, Chair

Ali Shojaie, Chair

Yen-Chi Chen

Program Authorized to Offer Degree:  
Statistics

University of Washington

**Abstract**

Non-Gaussian Graphical Models: Estimation with Score Matching and Causal Discovery  
under Zero-Inflation

Shiqing Yu

Co-Chairs of the Supervisory Committee:

Professor Mathias Drton

Department of Mathematics, Technische Universität München

Associate Professor Ali Shojaie

Department of Biostatistics

Graphical models specify conditional independence relations between variables. These include undirected graphical models and directed graphical models, the latter of which also capture causal relationships. This dissertation considers challenges but also opportunities brought about by non-Gaussianity in undirected and directed graphical models.

A common challenge in estimating parameters of probability density functions is the intractability of the normalizing constant, which hinders the use of maximum likelihood estimation, especially for non-Gaussian graphical models or when the density is supported on a subspace of  $\mathbb{R}^m$ . This dissertation presents a generalized form of score matching for distributions supported on general domains. The proposed approach generalizes the original forms proposed in Hyvärinen (2005) for  $\mathbb{R}^m$  and Hyvärinen (2007) for  $\mathbb{R}_+^m$ , while avoiding direct calculation of the normalizing constant, and yielding closed-form estimates for exponential families of continuous distributions. We generalize the regularized score matching method of Lin et al. (2016), and apply it to a general class of pairwise interaction graphical models, establishing strong theoretical guarantees.

Motivated by modern RNA sequencing technologies that provide gene expression measurements from single cells, this dissertation also addresses the challenge of causal discovery

from zero-inflated expression patterns in single cell data. In particular, we propose directed graphical models based on Hurdle conditional distributions to explore cause-effect relationships among the genes. We show that, under a natural and weak assumption, the exact directed acyclic graph for our model can be identified. We propose methods for graph recovery and show simulated experiments that validate the identifiability and graph estimation methods in practice.

# TABLE OF CONTENTS

	Page
List of Figures . . . . .	v
List of Tables . . . . .	viii
Glossary . . . . .	ix
Chapter 1: Introduction and Background . . . . .	1
1.1 Score Matching . . . . .	2
1.2 Undirected Graphical Models . . . . .	4
1.3 Directed Graphical Models and Causal Discovery . . . . .	5
1.4 Score Matching Applied to Pairwise Interaction Graphical Models . . . . .	6
1.5 Notation . . . . .	7
Chapter 2: Generalized Score Matching for Non-Negative Data . . . . .	9
2.1 Structure of Chapter . . . . .	10
2.2 Score Matching . . . . .	10
2.2.1 Original Score Matching . . . . .	10
2.2.2 Generalized Score Matching for Non-Negative Data . . . . .	11
2.3 Exponential Families . . . . .	14
2.4 Regularized Generalized Score Matching . . . . .	18
2.5 Score Matching for Graphical Models for Non-negative Data . . . . .	19
2.5.1 A General Framework of Pairwise Interaction Models . . . . .	20
2.5.2 Implementation for Different Models . . . . .	21
2.5.3 Computational Details . . . . .	24
2.5.4 Choice of the Function $\mathbf{h}$ . . . . .	25
2.5.5 Tuning Parameter Selection . . . . .	28
2.6 Theory for Graphical Models . . . . .	28

2.6.1	Theory for Pairwise Interaction Models . . . . .	29
2.6.2	Revisiting Gaussian Score Matching . . . . .	31
2.6.3	Generalized Score Matching for Truncated GGMs . . . . .	31
2.7	Numerical Experiments . . . . .	34
2.7.1	Structure of $\mathbf{K}$ . . . . .	35
2.7.2	Truncated GGMs . . . . .	35
2.7.3	Other $a/b$ Models . . . . .	42
2.7.4	RNAseq Data . . . . .	50
2.8	Discussion . . . . .	54
Chapter 3:	Generalized Score Matching for General Domain . . . . .	56
3.1	Introduction . . . . .	57
3.2	Preliminaries . . . . .	58
3.2.1	Original Score Matching on $\mathbb{R}^m$ (Hyvärinen, 2005) . . . . .	58
3.2.2	Score Matching on $\mathbb{R}_+^m$ (Hyvärinen, 2007; Yu et al., 2018, 2019b) . . . . .	59
3.2.3	Truncated Score Matching on Bounded Open Subsets of $\mathbb{R}^m$ (Liu and Kanamori, 2019) . . . . .	60
3.3	Generalized Score Matching for General Domains . . . . .	60
3.3.1	Assumption on the Domain . . . . .	60
3.3.2	Generalized Score Matching Loss for General Domains . . . . .	61
3.3.3	Comparison to Previous Work . . . . .	62
3.3.4	Form of the Truncated Componentwise Distance $\varphi$ . . . . .	64
3.3.5	The Empirical Generalized Score Matching Loss . . . . .	65
3.3.6	Extension of $g_0$ in Liu and Kanamori (2019) . . . . .	67
3.4	Exponential Families and $a$ - $b$ Models . . . . .	68
3.4.1	Exponential Families . . . . .	68
3.4.2	Pairwise Interaction Power $a$ - $b$ Models . . . . .	69
3.4.3	Regularized Score Matching . . . . .	70
3.5	$a$ - $b$ Models on Domains with Positive Measure . . . . .	70
3.5.1	Finite Normalizing Constant and Validity of Score Matching . . . . .	71
3.5.2	Estimation . . . . .	73
3.5.3	Univariate Examples . . . . .	74
3.6	$a$ - $b$ Models on Standard Simplices . . . . .	77

3.6.1	Estimation for General $a$ and $b$ . . . . .	78
3.6.2	log–log Models on Standard Simplices . . . . .	81
3.7	Theoretical Properties . . . . .	85
3.7.1	Truncated Gaussian Graphical Models on A Finite Disjoint Union of Convex Sets with Positive Measure . . . . .	86
3.7.2	Bounded Domains in $\mathbb{R}_+^m$ with Positive Measure . . . . .	87
3.7.3	Unbounded Domains in $\mathbb{R}_+^m$ with Positive Measure . . . . .	89
3.7.4	Models on the Standard Simplices . . . . .	90
3.8	Numerical Experiments . . . . .	91
3.8.1	Estimation — Choice of $\mathbf{h}$ and $\mathbf{C}$ . . . . .	91
3.8.2	Numerical Experiments for Domains with Positive Measure . . . . .	92
3.8.3	Numerical Experiments for $A^{m-1}$ Models on the Simplex . . . . .	103
3.9	DNA Methylation Data . . . . .	104
3.10	Discussion . . . . .	111
Chapter 4:	Directed Graphical Models and Causal Discovery for Zero-Inflated Data	113
4.1	Introduction . . . . .	114
4.2	Directed Graphical Models for Zero-Inflated Data . . . . .	117
4.2.1	Hurdle Joint Distributions for Zero-Inflated Continuous Observations	118
4.2.2	Hurdle Conditionals . . . . .	118
4.2.3	Directed Graphical Models for Zero-Inflation Data . . . . .	120
4.3	Identifiability . . . . .	121
4.3.1	Strong Identifiability . . . . .	121
4.3.2	Weak Identifiability . . . . .	122
4.4	Estimation of DAGs from Zero-Inflated Data . . . . .	125
4.4.1	Fitting Hurdle Conditionals . . . . .	126
4.4.2	Graph Search . . . . .	127
4.4.3	Stability Selection . . . . .	128
4.5	Numerical Experiments . . . . .	128
4.5.1	True Underlying DAGs and Distributions . . . . .	129
4.5.2	Estimation . . . . .	132
4.5.3	Results on Graph Recovery . . . . .	132
4.6	T Helper Cell Data . . . . .	137

4.7 Discussion . . . . .	138
Chapter 5: Discussion and Future Work . . . . .	142
Bibliography . . . . .	145
Appendix A: Appendices to Chapter 2 . . . . .	154
A.1 Proofs . . . . .	155
A.1.1 Proof of Theorem 2 . . . . .	155
A.1.2 Proof of Theorems and Examples in Section 2.3 . . . . .	157
A.1.3 Proof of Theorems in Section 2.5 . . . . .	161
A.1.4 Proof of Theorems in Section 2.6 . . . . .	167
A.2 Auxiliary Lemmas and Definitions . . . . .	172
A.3 Simulation Results for Erdős-Rényi Graphs . . . . .	179
A.3.1 Truncated GGMs . . . . .	180
A.3.2 Other $a/b$ Models . . . . .	185
Appendix B: Appendix to Chapter 3 . . . . .	189
Appendix C: Appendix to Chapter 4 . . . . .	211

## LIST OF FIGURES

Figure Number	Page
2.1 Log of asymptotic variance and efficiency with respect to the Cramér-Rao bound for $\hat{\mu}_h$ ( $\sigma^2 = 1$ known). . . . .	17
2.2 Log of asymptotic variance and efficiency with respect to the Cramér-Rao bound for $\hat{\sigma}_h^2$ ( $\mu = 0.5$ known). . . . .	17
2.3 Average ROC curves for the centered estimators for GGMs on $\mathbb{R}_+^m$ . . . . .	39
2.4 Average ROC curves for the non-centered profiled estimators for GGMs on $\mathbb{R}_+^m$ . . . . .	39
2.5 Performance of $\min(x, 3)$ for truncated centered GGMs on $\mathbb{R}_+^m$ using different multipliers. . . . .	41
2.6 Performance of the non-centered estimator with $h(x) = \min(x, 3)$ for GGMs on $\mathbb{R}_+^m$ with different choices of $\lambda_{\mathbf{K}}/\lambda_{\boldsymbol{\eta}}$ . . . . .	43
2.7 Average AUCs for edge recovery for the exponential models on $\mathbb{R}_+^m$ . . . . .	46
2.8 Average AUCs for edge recovery for the gamma models on $\mathbb{R}_+^m$ . . . . .	47
2.9 Average AUCs for edge recovery for $a = 3/2, b = 1/2$ on $\mathbb{R}_+^m$ . . . . .	48
2.10 Average AUCs for edge recovery for $a = 3/2, b = 0$ on $\mathbb{R}_+^m$ . . . . .	49
2.11 Graphs for RNAseq data estimated by $\min(x, 3)$ and by Lin et al. (2016) with isolated nodes removed. . . . .	51
2.12 Graphs for RNAseq data estimated by $\min(x, 3)$ and by Lin et al. (2016) with isolated nodes. . . . .	52
3.1 Comparison of $g_0, \varphi_{1_2, \mathcal{D}, 1}$ and $\varphi_{1_2, \mathcal{D}, 2}$ on $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}^2 : \ \mathbf{x}\ _2 < 1\}$ . . . . .	63
3.2 Comparison of $g_0, \varphi_{1_2, \mathcal{D}, 1}$ and $\varphi_{1_2, \mathcal{D}, 2}$ on $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}_+^2 : \ \mathbf{x}\ _2 < 1\}$ . . . . .	63
3.3 Univariate Gaussian example. . . . .	76
3.4 Average AUCs for the log models ( $a = 0$ ) on different types of domains. . . . .	96
3.5 Average AUCs for the exponential square-root models ( $a = 1/2$ ) on different types of domains. . . . .	97
3.6 Average AUCs for the Gaussian models ( $a = 1$ ) on different types of subsets of $\mathbb{R}_+^m$ . . . . .	98

3.7	Average AUCs for the Gaussian models ( $a = 1$ ) on different types of subsets of $\mathbb{R}^m$ . . . . .	99
3.8	Average AUCs for the $a = 3/2$ models on different types of domains. . . . .	100
3.9	Average AUCs for the $a = 2$ models on different types of domains. . . . .	101
3.10	Average AUCs for the $a = 3$ models on different types of domains. . . . .	102
3.11	Average AUCs for the $A^{m-1}$ models on the simplex. . . . .	104
3.12	Average errors for the $A^{m-1}$ models on the simplex. . . . .	105
3.13	Graphs for CpG sites estimated by regularized generalized score matching estimator. . . . .	108
3.14	Graphs for CpG sites aggregated by the genes. . . . .	109
3.15	Interlaced histogram (left) and Q-Q plot (right) showing the node degree distributions for both site graphs. . . . .	110
4.1	Pairwise scatter plots and kernel densities on four genes from the T helper cell data. . . . .	116
4.2	DAG structures used in our experiments. . . . .	129
4.3	Pairwise scatter plots of zero-inflated data generated using chain graphs and complete graphs. . . . .	131
4.4	Results of DAG recovery for chain graphs with $m = 5$ . . . . .	134
4.5	Results of DAG recovery for complete graphs with $m = 5$ . . . . .	135
4.6	Results of DAG recovery for lattice graphs with $m = 9$ . . . . .	136
4.7	Graphs estimated for for T helper cell data. . . . .	139
A.1	Average ROC curves for the centered estimators for GGMs on $\mathbb{R}_+^m$ under Erdős-Rényi graphs. . . . .	180
A.2	Average ROC curves for the non-centered profiled estimators for GGMs on $\mathbb{R}_+^m$ under Erdős-Rényi graphs. . . . .	183
A.3	Performance of $\min(x, 3)$ for truncated centered GGMs on $\mathbb{R}_+^m$ using different multipliers under Erdős-Rényi graphs. . . . .	183
A.4	Performance of the non-centered estimator for GGMs on $\mathbb{R}_+^m$ with $h(x) = \min(x, 3)$ with different choices of $\lambda_{\mathbf{K}}/\lambda_{\boldsymbol{\eta}}$ under Erdős-Rényi graphs. . . . .	184
A.5	Average AUCs for edge recovery for exponential models on $\mathbb{R}_+^m$ under Erdős-Rényi Graphs. . . . .	185
A.6	Average AUCs for edge recovery for gamma models on $\mathbb{R}_+^m$ under Erdős-Rényi Graphs. . . . .	186

A.7	Average AUCs for edge recovery for $a = 3/2$ and $b = 1/2$ on $\mathbb{R}_+^m$ under Erdős-Rényi Graphs. . . . .	187
A.8	Average AUCs for edge recovery for $a = 3/2$ and $b = 0$ on $\mathbb{R}_+^m$ under Erdős-Rényi Graphs. . . . .	188

## LIST OF TABLES

Table Number	Page
2.1 Mean and standard deviation of AUCs for the centered estimators. . . . .	37
2.2 Mean and standard deviation of AUCs for the profiled non-centered estimators.	40
2.3 List of genes with the highest node degrees in each estimated graph. . . . .	51
3.1 List of sites with the highest node degrees in each estimated graph. . . . .	107
A.1 Mean and standard deviation of AUCs for the centered estimators under Erdős-Rényi graphs. . . . .	181
A.2 Mean and standard deviation of AUCs for the profiled non-centered estimators under Erdős-Rényi graphs. . . . .	182

## GLOSSARY

AUC: The area under the ROC curve.

BIC: The Bayesian information criterion.

DAGS: Directed acyclic graphs; directed graphs with no directed cycles.

EBIC: The extended Bayesian information criterion.

FDR: The false discovery rate.

FPR: The false positive rate.

GGMS: Gaussian graphical models.

HURDLE CONDITIONALS: See Definition 10.

HURDLE POLYNOMIALS: See Definition 12.

ROC CURVE: Receiver operating characteristic curve.

TPR: The true positive rate.

## ACKNOWLEDGMENTS

First of all, I would like to thank my PhD advisors, Mathias Drton and Ali Shojaie, whose support and guidance have made this dissertation possible. I would also like to thank Yen-Chi Chen, Steve Mooney and Dmitriy Drusvyatskiy for their inputs and for serving on my supervisory committee.

Next, I would like to thank the faculty, the department staff, and my fellow PhDs in the Department of Statistics at the University of Washington for their great help and assistance during my time here. I would also like to thank the Department of Mathematical Sciences at the University of Copenhagen for their hospitality during my visit.

Furthermore, I would like to thank Divakar Viswanath at the University of Michigan for his help and guidance during my graduate school application process. I would also like to thank Xiang Yu for guiding me through my undergraduate research project at the University of Michigan.

In addition, I would like to thank Prof. Zhang, who had been treating me like his grandson since I was born, and encouraged me to pursue a PhD when I was in high school. May his soul rest in peace.

Last but not the least, I would like to thank my parents, my grandfather and my deceased grandmother for bringing me up and giving me support over the last 25 years. Many thanks to Carmen for her emotional support and patience during my PhD. I would also like to thank my friends and other family relatives for their help and encouragement. Finally, I would like to express my gratitude to all my deceased family and friends who unfortunately are no longer able to see me reaching this stage of my life, for leaving me with good memories in my life, and may they rest in peace.

# DEDICATION

to my family

Chapter 1

**INTRODUCTION AND BACKGROUND**

Graphical models specify conditional independence relations among variables. Non-Gaussian models provide the flexibility to express more complicated interactions among the variables, especially for distributions that are naturally supported on some proper subset of the  $m$ -dimensional real space  $\mathbb{R}^m$ , or those that are not absolutely continuous with respect to the Lebesgue measure. This comes with both opportunities and significant challenges. A challenge for undirected graphical models is the fact that non-Gaussian models as well as Gaussian models that are restricted to a proper subset of  $\mathbb{R}^m$  generally have intractable normalizing constants, making maximum likelihood estimation computationally difficult. For directed graphical models, which capture causal relationships among the variables, non-Gaussian modeling is again more involved, but here it also offers an opportunity for causal discovery. Indeed, for many Gaussian settings, they may only be identifiable up to an equivalence class and thus preventing causal discovery. In this dissertation, we address these opportunities and challenges in detail.

### 1.1 Score Matching

Probability density functions, especially those for multivariate graphical models, are often defined only up to a normalizing constant. In high-dimensional settings, calculating the normalizing constant is often computationally intensive and intractable except in some special cases such as Gaussian densities (Gaussian graphical models) on the  $m$ -dimensional real space  $\mathbb{R}^m$ . The situation becomes worse when the distributions are defined only on a proper subset of  $\mathbb{R}^m$ . For example, even the Gaussian densities on  $\mathbb{R}_+^m$  have intractable normalizing constants except for special cases e.g. with diagonal covariance matrices. This inability to calculate normalizing constants makes density estimation challenging and maximum likelihood estimation becomes less ideal.

Fortunately, Hyvärinen (2005) introduced the method of *score matching* that comes to our rescue. Given a continuous distribution  $P_0$  supported on  $\mathbb{R}^m$  with Lebesgue density  $p_0$ , let  $\mathcal{P}$  be a family of distributions with twice continuously differentiable densities. The score matching estimator of  $p_0$  using  $\mathcal{P}$  as a model is the minimizer of the expected  $\ell_2$

distance between the gradients of the true log density  $\log p_0$  on  $\mathbb{R}^m$  and a proposed log density  $\log p$  from  $\mathcal{P}$ , namely  $\int_{\mathbb{R}^m} p_0(\mathbf{x}) \|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_0(\mathbf{x})\|_2^2 d\mathbf{x}$ . The loss depends on  $p_0$ , but integration by parts can be used to rewrite it in a form that can be approximated by averaging over the sample without knowing  $p_0$ . A key feature of score matching is that normalizing constants cancel in gradients of log-densities, allowing for simple treatment of models with intractable normalizing constants. For exponential families, the loss is quadratic in the canonical parameter, making optimization straightforward.

If the considered distributions are supported on a proper subset of  $\mathbb{R}^m$ , then the integration by parts arguments underlying the score matching estimator may fail due to discontinuities at the boundary of the support. For data supported on the non-negative orthant  $\mathcal{D} \equiv \mathbb{R}_+^m$ , Hyvärinen (2007) addresses this problem by modifying the loss to  $\int_{\mathbb{R}^m} p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) \odot \mathbf{x} - \nabla \log p_0(\mathbf{x}) \odot \mathbf{x}\|_2^2 d\mathbf{x}$ , where  $\odot$  denotes entrywise multiplication. In this loss, boundary effects are dampened by multiplying gradients element-wise with the identity functions  $x_j$ .

In Chapter 2, we propose *generalized score matching* methods for densities supported on  $\mathcal{D} \equiv \mathbb{R}_+^m$  based on element-wise multiplication with slowly growing/bounded functions  $\mathbf{h}^{1/2}(\mathbf{x})$  in place of  $\mathbf{x}$ . We show that this modification drastically improves estimation accuracy, both theoretically and empirically. This work is published in Yu et al. (2018, 2019b).

In Chapter 3, we further extend the generalized score matching methods to densities supported on more general domains  $\mathcal{D} \subseteq \mathbb{R}^m$  that may be either bounded or unbounded. This includes densities supported on simplex domains, for which density estimation has been difficult due to their sum-to-one constraint, with important applications such as compositional microbiome data. Our generalization completes the framework for estimation of continuous distributions supported on various types of domains with intractable normalizing constants.

Liu and Kanamori (2019), on the other hand, extended our results in Chapter 2 to bounded open subspaces  $\mathcal{D}$  of  $\mathbb{R}^m$  that have piecewise smooth boundaries by minimizing  $\sup_{g \in \mathcal{G}} \int_{\mathcal{D}} p_0(\mathbf{x}) g(\mathbf{x}) \|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \nabla_{\mathbf{x}} \log p_0(\mathbf{x})\|_2^2 d\mathbf{x}$ , where  $\mathcal{G} \equiv \{g | g(\mathbf{x}) = 0 \forall \mathbf{x} \in \partial\mathcal{D} \text{ and } g \text{ is 1-Lipschitz continuous}\}$ ,  $\partial\mathcal{D}$  being the boundary of  $\mathcal{D}$ . The  $\sup_{g \in \mathcal{G}}$  in the loss is obtained at  $g_0(\mathbf{x}) \equiv \min_{\mathbf{x}' \in \partial\mathcal{D}} \|\mathbf{x} - \mathbf{x}'\|$ , the distance of  $\mathbf{x}$  to  $\partial\mathcal{D}$ . In Chapter 3, we extend

their slightly different approach to unbounded subsets of  $\mathbb{R}^m$ , but we empirically show that their approach and its extension seem inferior to our generalized score matching estimators in many settings.

## 1.2 Undirected Graphical Models

Graphical models specify conditional independence relations among variables in a random vector  $\mathbf{Y}$  indexed by the nodes  $\mathcal{V}$  of a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with edge set  $\mathcal{E}$  (Maathuis et al., 2019). Models based on undirected graphs may be used to explore conditional independence between any two variables  $Y_V$  and  $Y_U$  given all others  $(Y_W)_{W \neq U, V}$ , as represented by the absence of an edge between  $V$  and  $U$  in  $\mathcal{E}$ .

Among undirected graphical models, largely due to their tractability, Gaussian graphical models (GGMs) have gained great popularity. The conditional independence graph of a multivariate normal vector  $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is determined by the *inverse covariance matrix*  $\mathbf{K} \equiv \boldsymbol{\Sigma}^{-1}$ , also termed *concentration* or *precision matrix*. Specifically,  $X_i$  and  $X_j$  are conditionally independent given all other variables if and only if the  $(i, j)$ -th and the  $(j, i)$ -th entries of  $\mathbf{K}$  are both zero. This simple relation underlies a rich literature including Drton and Perlman (2004), Meinshausen and Bühlmann (2006), Yuan and Lin (2007) and Friedman et al. (2008), among others.

More recent work has provided tractable procedures also for non-Gaussian undirected graphical models. This includes Gaussian copula models (Liu et al., 2009; Dobra and Lenkoski, 2011; Liu et al., 2012), Ising models (Ravikumar et al., 2010), other exponential family models (Chen et al., 2015; Yang et al., 2015), as well as semi- or non-parametric estimation techniques (Fellinghauer et al., 2013; Voorman et al., 2014). In this dissertation, we apply our generalized score matching method to a class of pairwise interaction power models, introduced in Section 1.4, that generalizes Gaussian and non-negative Gaussian random variables, recently considered by Lin et al. (2016) and Yu et al. (2016), as well as square root graphical models proposed by Inouye et al. (2016) when the sufficient statistic function is a pure power.

### 1.3 Directed Graphical Models and Causal Discovery

Models based on directed acyclic graphs (DAGs), for which  $\mathcal{E}$  is comprised of directed edges, capture conditional independence structure that naturally arises from cause-effect relationships between the variables. In biology and genetics, graphical models have been applied to infer the structure of gene regulatory networks based on measurements of gene expression (Maathuis et al., 2019, Sections 20-21). Traditional technologies produce expression levels aggregated over hundreds or thousands of individual cells, and these bulk measurements are frequently modeled using the assumption of Gaussianity. In directed GGMs the exact structure of the underlying DAG cannot be identified from purely observational data, and the target of inference becomes an equivalence class of DAGs. For instance, one cannot differentiate between  $V \rightarrow U$  and  $U \rightarrow V$  when the variables are assumed bivariate normal. In the Gaussian case, directed graphical models posit linear functional relationships between the variables coupled with additive Gaussian noise. A more recent line of work emphasizes that directed graphical models that alter this assumption to nonlinear functional relationships and additive noise (Peters et al., 2014), or linear relations and non-Gaussian noise (Shimizu et al., 2006; Wang and Drton, 2020), or linear relations with homoscedastic Gaussian noise (Peters and Bühlmann, 2013; Chen et al., 2019) are amenable to causal discovery in the sense that different DAGs are no longer equivalent.

More recent technology obtains sequencing measurements of mRNA present in single cells. This new technology, as well as the larger sample sizes it provides, promise to give more information than bulk measurements, but at the same time bring in a unique new challenge. At the single cell level, genes appear as “on” with positive single cell gene expression levels, or as “off” with the recorded measurements zero or negligible (McDavid et al., 2019). A novel undirected graphical model that deals with this zero-inflation was introduced by McDavid et al. (2019). Since more information can be inferred from single-cell sequencing data, one would hope that the data can also be analyzed using more informative directed graphical models, and that we can infer which variables (genes) are the causes of change

in other variables (expression levels of other genes). In Chapter 4, we formulate directed graphical models for zero-inflated data, and prove that under a weak assumption one can recover the exact DAG from the joint distribution. We use simulation studies to support our identifiability theory and justify the use of two methods for estimating the DAG in practice.

#### 1.4 Score Matching Applied to Pairwise Interaction Graphical Models

In this section we introduce the *pairwise interaction power models*, also called *a-b models*, that we examine in Chapters 2 and 3. In particular, we assume the models have (Lebesgue) probability density functions proportional to

$$\exp \left\{ -\frac{1}{2a} \mathbf{x}^a \mathbf{K} \mathbf{x}^a + \frac{1}{b} \boldsymbol{\eta}^\top \mathbf{x}^b \right\} \quad (1.1)$$

supported on some general domain  $\mathcal{D} \subseteq \mathbb{R}^m$ . Here  $a \geq 0$  and  $b \geq 0$  are known constants, and  $\mathbf{K} \in \mathbb{R}^{m \times m}$  and  $\boldsymbol{\eta} \in \mathbb{R}^m$  are unknown parameters of interest. For  $a = 0$  we define  $\mathbf{x}^{a\top} \mathbf{K} \mathbf{x}^a / a \equiv (\log \mathbf{x})^\top \mathbf{K} (\log \mathbf{x})$  and for  $b = 0$  we define  $\boldsymbol{\eta}^\top \mathbf{x}^b / b \equiv \boldsymbol{\eta}^\top (\log \mathbf{x})$ . This class of models is motivated by the form of important univariate distributions, including gamma and truncated normal distributions.

This model class provides a framework for pairwise interaction that is concrete yet rich enough to capture key differences in how densities may behave; in particular, if  $\mathcal{D}$  is a product set of the form  $\mathcal{D}_1 \times \cdots \times \mathcal{D}_m$  with one-dimensional sets  $\mathcal{D}_1, \dots, \mathcal{D}_m$ ,  $X_i$  and  $X_j$  are conditionally independent given all others if and only if  $\kappa_{ij} = \kappa_{ji} = 0$  in the interaction matrix  $\mathbf{K}$ , just as for Gaussian graphical models. We will develop estimators of  $(\boldsymbol{\eta}, \mathbf{K})$  and the associated conditional independence graph using the proposed *generalized score matching*.

A special case of (1.1) are (truncated) Gaussian graphical models, with  $a = b = 1$ . Let  $\boldsymbol{\mu} \in \mathbb{R}^m$ , and let  $\mathbf{K}$  be a positive definite matrix. Then a random vector  $\mathbf{X} \in \mathcal{D}$  follows a (truncated) normal distribution for mean parameter  $\boldsymbol{\mu}$  and inverse covariance parameter  $\mathbf{K}$ , in symbols  $\mathbf{X} \sim \text{TN}_{\mathcal{D}}(\boldsymbol{\mu}, \mathbf{K})$ , if it has density proportional to

$$\exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{K} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (1.2)$$

on domain  $\mathcal{D} \subseteq \mathbb{R}^m$ . We refer to  $\boldsymbol{\Sigma} = \mathbf{K}^{-1}$  as the covariance parameter of the distribution, and note that the  $\boldsymbol{\eta}$  parameter in (1.1) is  $\mathbf{K}\boldsymbol{\mu}$ . Another special case of (1.1) is the exponential square root graphical models in Inouye et al. (2016), where  $a = b = 1/2$  and  $\mathcal{D} \equiv \mathbb{R}_+^m$ .

Lin et al. (2016) proposed estimating the truncated GGMs supported on  $\mathcal{D} \equiv \mathbb{R}_+^m$  using the score matching loss from Hyvärinen (2007) plus an  $\ell_1$  loss on the off-diagonal elements of  $\mathbf{K}$ . We propose a similar regularized loss using our generalized score matching loss for estimating the  $a$ - $b$  models above. The loss in Lin et al. (2016) may be unbounded from below, and thus infinitely many minimizers exist; in Chapter 2 we fix this problem by introducing *diagonal multipliers*, effectively imposing an elastic net-type penalty (Zou and Hastie, 2005) on  $\mathbf{K}$ . In addition, Lin et al. (2016) requires a sample size  $n = \Omega((\log m)^8)$  for GGMs on  $\mathbb{R}_+^m$ , while for our generalized score matching estimator the sample complexity can be reduced to  $n = \Omega(\log m)$  for any domain  $\mathcal{D}$  that is a finite disjoint union of convex sets. In Chapter 3, we present similar theoretical guarantees for general  $a$ - $b$  models on a general domain  $\mathcal{D}$ .

## 1.5 Notation

Constant scalars, vectors, and functions are written in lower-case (e.g.,  $a$ ,  $\mathbf{a}$ ), random scalars and vectors in upper-case (e.g.,  $X$ ,  $\mathbf{X}$ ). Regular font is used for scalars (e.g.  $a$ ,  $X$ ), and boldface for vectors (e.g.  $\mathbf{a}$ ,  $\mathbf{X}$ ). Graph nodes are written in unbolded upper-case (e.g.  $U$ ,  $V$ ), and domains, sets of nodes, vectors of functions in calligraphic font (e.g.  $\mathcal{D}$ ,  $\mathcal{V}$ ,  $\mathcal{H}$ ). Matrices are in upright bold, with constant matrices in upper-case ( $\mathbf{K}$ ,  $\mathbf{M}$ ) and random matrices holding observations in lower-case ( $\mathbf{x}$ ,  $\mathbf{y}$ ).

Subscripts refer to entries in vectors and columns in matrices. When used as a subscript of a vector, a set of nodes/indices selects the corresponding entries from the vector, e.g.  $\mathbf{y}_{\mathcal{V}} = (y_V)_{V \in \mathcal{V}}$ . Superscripts refer to rows in matrices. So  $X_j$  is the  $j$ -th component of a random vector  $\mathbf{X}$ . For a data matrix  $\mathbf{x} \in \mathbb{R}^{n \times m}$ , each row comprising one observation of  $m$  variables/features,  $X_j^{(i)}$  is the  $j$ -th feature for the  $i$ -th observation. For a matrix  $\mathbf{K} = [\kappa_{ij}]_{i,j} \in \mathbb{R}^{q \times q}$ ,  $\text{diag}(\mathbf{K}) = (\kappa_{11}, \dots, \kappa_{qq})^\top$  denotes its diagonal, and for a vector  $\mathbf{v} \in \mathbb{R}^q$ ,  $\text{diag}(\mathbf{v})$  is the  $q \times q$  diagonal matrix with diagonal entries  $v_1, \dots, v_q$ . Stacking the columns

of  $\mathbf{K}$  gives its vectorization  $\text{vec}(\mathbf{K}) = (\kappa_{11}, \dots, \kappa_{q1}, \kappa_{12}, \dots, \kappa_{q2}, \dots, \kappa_{1r}, \dots, \kappa_{qr})^\top$ .

For  $a \geq 1$ , the  $\ell_a$ -norm of a vector  $\mathbf{v} \in \mathbb{R}^q$  is denoted

$$\|\mathbf{v}\|_a = \left( \sum_{j=1}^q |v_j|^a \right)^{1/a},$$

with  $\|\mathbf{v}\|_\infty = \max_{j=1, \dots, q} |v_j|$ . A matrix  $\mathbf{K} = [\kappa_{ij}]_{i,j} \in \mathbb{R}^{q \times r}$  has Frobenius norm

$$\|\mathbf{K}\|_F \equiv \|\text{vec}(\mathbf{K})\|_2 \equiv \sqrt{\sum_{i=1}^q \sum_{j=1}^r \kappa_{ij}^2},$$

and max norm  $\|\mathbf{K}\|_\infty \equiv \|\text{vec}(\mathbf{K})\|_\infty \equiv \max_{i,j} |\kappa_{ij}|$ . Its  $\ell_a$ - $\ell_b$  operator norm is

$$\|\mathbf{K}\|_{a,b} \equiv \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{K}\mathbf{x}\|_b}{\|\mathbf{x}\|_a}$$

with shorthand notation  $\|\mathbf{K}\|_a \equiv \|\mathbf{K}\|_{a,a}$ ; for instance,  $\|\mathbf{K}\|_\infty \equiv \max_{i=1, \dots, q} \sum_{j=1}^r |\kappa_{ij}|$ .

For a function  $f : \mathbb{R}^m \rightarrow \mathbb{R}$ , we define  $\partial_j f(\mathbf{x})$  as the partial derivative with respect to  $x_j$ , and  $\partial_{jj} f(\mathbf{x}) = \partial_j \partial_j f(\mathbf{x})$ . For vector-valued  $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ ,  $\mathbf{x} \mapsto (f_1(\mathbf{x}), \dots, f_m(\mathbf{x}))$ , we let  $\mathbf{f}'(\mathbf{x}) = (\partial_1 f_1(\mathbf{x}), \dots, \partial_m f_m(\mathbf{x}))^\top$  be the vector of partial derivatives. Likewise  $\mathbf{f}''(\mathbf{x})$  is used for second partial derivatives. For  $a \in \mathbb{R}$ , let  $\mathbf{f}^a(\mathbf{x}) \equiv (f_1^a(\mathbf{x}), \dots, f_m^a(\mathbf{x}))^\top$ . For two compatible functions  $f$  and  $g$ ,  $f \circ g$  denotes their function composition. For  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^m$ ,  $\mathbf{a} \odot \mathbf{b} \equiv (a_1 b_1, \dots, a_m b_m)^\top$ .

Finally, the symbol  $\mathbf{1}_A(\cdot)$  denotes the indicator function of a set  $A$ , while  $\mathbf{1}_n \in \mathbb{R}^n$  is the vector of all 1's. When it is clear from the context,  $\mathbb{E}_0$  denotes the expectation under a true distribution  $P_0$ .

Chapter 2

**GENERALIZED SCORE MATCHING FOR NON-NEGATIVE  
DATA**

## 2.1 Structure of Chapter

In this chapter we discuss our *generalized score matching* estimator for continuous distributions supported on  $\mathbb{R}_+^m$ , with applications of *regularized general score matching* to the pairwise interaction power models (see (1.1)). The chapter is organized as follows. Section 2.2 introduces score matching and our proposed *generalized score matching* for distributions on  $\mathbb{R}_+^m$ . In Section 2.3, we apply generalized score matching to exponential families, with univariate truncated normal distributions as an example. *Regularized generalized score matching* for graphical models and the notion of *diagonal multipliers* are formulated in Section 2.4. The estimators for pairwise interaction power models are shown in Section 2.5, while theoretical consistency results are presented in Section 2.6, where we treat the probabilistically most tractable case of truncated GGMs. Simulation results and applications to RNAseq data are given in Section 2.7. Proofs for theorems in Sections 2.2–2.6 are presented in Appendices A.1 and A.2. Additional experimental results are presented in Appendix A.3. The work in this chapter is published in Yu et al. (2018) and Yu et al. (2019b).

## 2.2 Score Matching

In this section, we review the original score matching and develop our generalized score matching estimators.

### 2.2.1 Original Score Matching

Let  $\mathbf{X}$  be a random vector taking values in  $\mathbb{R}^m$  with distribution  $P_0$  and density  $p_0$ . Let  $\mathcal{P}$  be a family of distributions of interest with twice continuously differentiable densities supported on  $\mathbb{R}^m$ . Suppose  $P_0 \in \mathcal{P}$ . The *score matching loss* for  $P \in \mathcal{P}$ , with density  $p$ , is given by

$$J(P) = \int_{\mathbb{R}^m} p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) - \nabla \log p_0(\mathbf{x})\|_2^2 d\mathbf{x}. \quad (2.1)$$

The gradients in (2.1) can be thought of as gradients with respect to a hypothetical location parameter, evaluated at the origin (Hyvärinen, 2005). The loss  $J(P)$  is minimized if and only

if  $P = P_0$ , which forms the basis for estimation of  $P_0$ . Importantly, since the loss depends on  $p$  only through its log-gradient, it suffices to know  $p$  up to a normalizing constant. Under mild conditions, (2.1) can be rewritten as

$$J(P) = \int_{\mathbb{R}^m} p_0(\mathbf{x}) \sum_{j=1}^m \left[ \partial_{jj} \log p(\mathbf{x}) + \frac{(\partial_j \log p(\mathbf{x}))^2}{2} \right] d\mathbf{x}, \quad (2.2)$$

plus a constant independent of  $p$ . The integral in (2.2) can be approximated by a sample average; this alleviates the need for knowing the true density  $p_0$ , and provides a way to estimate  $p_0$ .

### 2.2.2 Generalized Score Matching for Non-Negative Data

When the true density  $p_0$  is supported on a proper subset of  $\mathbb{R}^m$ , the integration by parts underlying the equivalence of (2.1) and (2.2) may fail due to discontinuity at the boundary. For distributions supported on the non-negative orthant,  $\mathbb{R}_+^m$ , Hyvärinen (2007) addressed this issue by instead minimizing the *non-negative score matching loss*

$$J_+(P) = \int_{\mathbb{R}_+^m} p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) \odot \mathbf{x} - \nabla \log p_0(\mathbf{x}) \odot \mathbf{x}\|_2^2 d\mathbf{x}. \quad (2.3)$$

This loss can be motivated via gradients with respect to a hypothetical scale parameter (Hyvärinen, 2007). Under mild conditions,  $J_+(P)$  can again be rewritten in terms of an expectation of a function independent of  $p_0$ , thus allowing one to form a sample loss.

In this work, we consider generalizing the non-negative score matching loss as follows.

**Definition 1.** Let  $\mathcal{P}_+$  be the family of distributions of interest, and assume every  $P \in \mathcal{P}_+$  has a twice continuously differentiable density supported on  $\mathbb{R}_+^m$ . Suppose the  $m$ -variate random vector  $\mathbf{X}$  has true distribution  $P_0 \in \mathcal{P}_+$ , and let  $p_0$  be its twice continuously differentiable density. Let  $h_1, \dots, h_m : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  be a.s. positive functions that are absolutely continuous in every bounded sub-interval of  $\mathbb{R}_+$ , and set  $\mathbf{h}(\mathbf{x}) = (h_1(x_1), \dots, h_m(x_m))^\top$ . For  $P \in \mathcal{P}_+$  with density  $p$ , the generalized  $\mathbf{h}$ -score matching loss is

$$J_{\mathbf{h}}(P) = \int_{\mathbb{R}_+^m} \frac{1}{2} p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) \odot \mathbf{h}(\mathbf{x})^{1/2} - \nabla \log p_0(\mathbf{x}) \odot \mathbf{h}(\mathbf{x})^{1/2}\|_2^2 d\mathbf{x}, \quad (2.4)$$

where  $\mathbf{h}^{1/2}(\mathbf{x}) \equiv (h_1^{1/2}(x_1), \dots, h_m^{1/2}(x_m))^\top$ .

**Proposition 1.** *The distribution  $P_0$  is the unique minimizer of  $J_{\mathbf{h}}(P)$  for  $P \in \mathcal{P}_+$ .*

*Proof.* First, observe that  $J_{\mathbf{h}}(P) \geq 0$  and  $J_{\mathbf{h}}(P_0) = 0$ . For uniqueness, suppose  $J_{\mathbf{h}}(P_1) = 0$  for some  $P_1 \in \mathcal{P}_+$ . Let  $p_0$  and  $p_1$  be the respective densities. By assumption  $p_0(\mathbf{x}) > 0$  a.s. and  $h_j^{1/2}(\mathbf{x}) > 0$  a.s. for all  $j = 1, \dots, m$ . Therefore, we must have  $\nabla \log p_1(\mathbf{x}) = \nabla \log p_0(\mathbf{x})$  a.s., or equivalently,  $p_1(\mathbf{x}) = \text{const} \times p_0(\mathbf{x})$  almost surely in  $\mathbb{R}_+^m$ . Since  $p_1$  and  $p_0$  are continuous densities supported on  $\mathbb{R}_+^m$ , it follows that  $p_1(\mathbf{x}) = p_0(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{R}_+^m$ .  $\square$

Choosing all  $h_j(x) = x^2$  recovers the loss from (2.3). In our generalization, we will focus on using functions  $h_j$  that are increasing but are bounded or grow rather slowly. This will alleviate the need to estimate higher moments, leading to better practical performance and improved theoretical guarantees.

We will consider the following assumptions:

$$(A1) \quad p_0(\mathbf{x})h_j(x_j)\partial_j \log p(\mathbf{x}) \Big|_{x_j \searrow 0^+}^{x_j \nearrow +\infty} = 0, \quad \forall \mathbf{x}_{-j} \in \mathbb{R}_+^{m-1}, \quad \forall p \in \mathcal{P}_+;$$

$$(A2) \quad \mathbb{E}_{p_0} \|\nabla \log p(\mathbf{X}) \odot \mathbf{h}^{1/2}(\mathbf{X})\|_2^2 < +\infty, \quad \mathbb{E}_{p_0} \|(\nabla \log p(\mathbf{X}) \odot \mathbf{h}(\mathbf{X}))'\|_1 < +\infty, \quad \forall p \in \mathcal{P}_+,$$

where  $\partial_j \log p(\mathbf{x}) \equiv \frac{\partial \log p(\mathbf{y})}{\partial y_j} \Big|_{\mathbf{y}=\mathbf{x}}$ ,  $f(\mathbf{x}) \Big|_{x_j \searrow 0^+}^{x_j \nearrow +\infty} \equiv \lim_{x_j \nearrow +\infty} f(\mathbf{x}) - \lim_{x_j \searrow 0} f(\mathbf{x})$ , “ $\forall p \in \mathcal{P}_+$ ” is a shorthand for “for all  $p$  being the density of some  $P \in \mathcal{P}_+$ ”, and the prime symbol denotes component-wise differentiation. While the second half of (A2) was not made explicit in Hyvärinen (2005, 2007), (A1)-(A2) were both required for integration by parts and Fubini-Tonelli to apply.

Once the forms of  $p_0$  and  $p$  are given, sufficient conditions for  $\mathbf{h}$  for Assumptions (A1)-(A2) to hold are easy to find. In particular, (A1) and (A2) are easily satisfied and verified for exponential families.

Integration by parts yields the following theorem which shows that  $J_{\mathbf{h}}$  from (2.4) is an expectation (under  $P_0$ ) of a function that does not depend on  $p_0$ , similar to (2.2). The proof is given in Appendix A.1.1.

**Theorem 2.** *Under (A1) and (A2), the loss from (2.4) equals*

$$J_{\mathbf{h}}(P) = \int_{\mathbb{R}_+^m} p_0(\mathbf{x}) \sum_{j=1}^m \left[ h'_j(x_j) \partial_j(\log p(\mathbf{x})) + h_j(x_j) \partial_{jj}(\log p(\mathbf{x})) + \frac{1}{2} h_j(x_j) (\partial_j(\log p(\mathbf{x})))^2 \right] d\mathbf{x} \quad (2.5)$$

plus a constant independent of  $p$ .

Given a data matrix  $\mathbf{x} \in \mathbb{R}^{n \times m}$  with rows  $\mathbf{X}^{(i)}$ , we define the sample version of (2.5) as

$$\hat{J}_{\mathbf{h}}(P) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \left\{ h'_j(X_j^{(i)}) \partial_j(\log p(\mathbf{X}^{(i)})) + h_j(X_j^{(i)}) \left[ \partial_{jj}(\log p(\mathbf{X}^{(i)})) + \frac{1}{2} (\partial_j(\log p(\mathbf{X}^{(i)})))^2 \right] \right\}. \quad (2.6)$$

Subsequently, for a distribution  $P$  with density  $p$ , we let  $J_{\mathbf{h}}(p) \equiv J_{\mathbf{h}}(P)$ . Similarly, when a distribution  $P_{\boldsymbol{\theta}}$  with density  $p_{\boldsymbol{\theta}}$  is associated to a parameter vector  $\boldsymbol{\theta}$ , we write  $J_{\mathbf{h}}(\boldsymbol{\theta}) \equiv J_{\mathbf{h}}(p_{\boldsymbol{\theta}}) \equiv J_{\mathbf{h}}(P_{\boldsymbol{\theta}})$ . We apply similar conventions to the sample version  $\hat{J}_{\mathbf{h}}(P)$ . We note that this type of loss is also treated in slightly different settings in Parry (2016) and Almeida and Gidas (1993).

**Remark 1.** In the one-dimensional case, using the notation in Parry et al. (2012),  $J_{\mathbf{h}}(P)$  and  $\hat{J}_{\mathbf{h}}(P)$  correspond to  $d(P_0, P)$  and  $S(x, P)$ , respectively, and can be generated by  $\phi(x, p, p_1) \equiv -h(x)p_1^2/(2p)$  (c.f. Equations (39), (51), (53) and Section 10.1 therein). Thus Theorem 2 follows from this correspondence. While (A1) is equivalent to the condition implied by the boundary divergence  $d_b = 0$  in that paper, (A2), which we assume for invoking Fubini-Tonelli due to multi-dimensionality, is not present. On the other hand, while Parry (2016) treats the multivariate case, it does not cover the connection between our  $J_{\mathbf{h}}$  and  $\hat{J}_{\mathbf{h}}$ . Since  $\phi$  is concave but not strictly concave in  $(p, p_1)$ , the results in Parry (2016) only imply that  $P_0$  is a minimizer, a weaker conclusion than Proposition 1.

### 2.3 Exponential Families

In this section, we study the case where  $\mathcal{P}_+ \equiv \{p_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \Theta\}$  is an exponential family comprising continuous distributions with support  $\mathbb{R}_+^m$ . More specifically, we consider densities that are indexed by the canonical parameter  $\boldsymbol{\theta} \in \mathbb{R}^r$  and have the form

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{t}(\mathbf{x}) - \psi(\boldsymbol{\theta}) + b(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}_+^m, \quad (2.7)$$

where  $\mathbf{t}(\mathbf{x}) \in \mathbb{R}_+^r$  comprises the sufficient statistics,  $\psi(\boldsymbol{\theta})$  is a normalizing constant depending on  $\boldsymbol{\theta}$  only, and  $b(\mathbf{x})$  is the base measure, with  $\mathbf{t}$  and  $b$  a.s. differentiable with respect to each component. Define  $\mathbf{t}'_j(\mathbf{x}) \equiv (\partial_j t_1(\mathbf{x}), \dots, \partial_j t_r(\mathbf{x}))^\top$  and  $b'_j(\mathbf{x}) \equiv \partial_j b(\mathbf{x})$ .

**Theorem 3.** *Under Assumptions (A1)-(A2) from Section 2.2.2, the empirical generalized  $\mathbf{h}$ -score matching loss (2.6) can be rewritten as a quadratic function in  $\boldsymbol{\theta} \in \mathbb{R}^r$ :*

$$\hat{J}_{\mathbf{h}}(p_{\boldsymbol{\theta}}) = \frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{\Gamma}(\mathbf{x}) \boldsymbol{\theta} - \mathbf{g}(\mathbf{x})^\top \boldsymbol{\theta} + \text{const}, \quad \text{where} \quad (2.8)$$

$$\boldsymbol{\Gamma}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m h_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)})^\top \quad \text{and} \quad (2.9)$$

$$\mathbf{g}(\mathbf{x}) = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \left[ h_j(X_j^{(i)}) b'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) + h_j(X_j^{(i)}) \mathbf{t}''_j(\mathbf{X}^{(i)}) + h'_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \right] \quad (2.10)$$

are sample averages of functions of the data matrix  $\mathbf{x}$  only.

Define  $\boldsymbol{\Gamma}_0 \equiv \mathbb{E}_{p_0} \boldsymbol{\Gamma}(\mathbf{x})$ ,  $\mathbf{g}_0 \equiv \mathbb{E}_{p_0} \mathbf{g}(\mathbf{x})$ , and  $\boldsymbol{\Sigma}_0 \equiv \mathbb{E}_{p_0} [(\boldsymbol{\Gamma}(\mathbf{x}) \boldsymbol{\theta}_0 - \mathbf{g}(\mathbf{x}))(\boldsymbol{\Gamma}(\mathbf{x}) \boldsymbol{\theta}_0 - \mathbf{g}(\mathbf{x}))^\top]$ .

**Theorem 4.** *Suppose that*

(C1)  $\boldsymbol{\Gamma}$  is a.s. invertible, and

(C2)  $\boldsymbol{\Gamma}_0$ ,  $\boldsymbol{\Gamma}_0^{-1}$ ,  $\mathbf{g}_0$  and  $\boldsymbol{\Sigma}_0$  exist and are entry-wise finite.

Then the minimizer of (2.8) is a.s. unique with closed-form solution  $\hat{\boldsymbol{\theta}} \equiv \boldsymbol{\Gamma}(\mathbf{x})^{-1} \mathbf{g}(\mathbf{x})$ . Moreover,

$$\hat{\boldsymbol{\theta}} \xrightarrow{\text{a.s.}} \boldsymbol{\theta}_0 \quad \text{and} \quad \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \rightarrow_d \mathcal{N}_r(\mathbf{0}, \boldsymbol{\Gamma}_0^{-1} \boldsymbol{\Sigma}_0 \boldsymbol{\Gamma}_0^{-1}) \quad \text{as } n \rightarrow \infty.$$

Theorems 3 and 4 are proved in Appendix A.1.2. Theorem 3 clarifies the quadratic nature of the loss, and Theorem 4 provides a basis for asymptotically valid tests and confidence intervals for the parameter  $\boldsymbol{\theta}$ . Note that Condition (C1) holds if and only if  $h_j(X_j) > 0$  a.s. and  $[\mathbf{t}'_j(\mathbf{X}^{(1)}), \dots, \mathbf{t}'_j(\mathbf{X}^{(n)})] \in \mathbb{R}^{r \times n}$  has rank  $r$  a.s. for some  $j = 1, \dots, m$ .

The conclusion in Theorem 4 indicates that, similar to the estimator in Hyvärinen (2007) with  $h_j(x) = x^2$ , the closed-form solution for our generalized  $\hat{\boldsymbol{\theta}}$  allows one to consistently estimate the canonical parameter in an exponential family distribution without needing to calculate the often complicated normalizing constant  $\psi(\boldsymbol{\theta})$  or resort to numerical methods. Computational details are explicated in Section 2.5.3.

Below we illustrate the estimator  $\hat{\boldsymbol{\theta}}$  in the case of univariate truncated normal distributions. We assume (A1)-(A2) and (C1)-(C2) throughout.

**Example 3.1.** *Univariate ( $m = r = 1$ ) truncated normal distributions for mean parameter  $\mu$  and variance parameter  $\sigma^2$  have density*

$$p_{\mu, \sigma^2}(x) \propto \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}, \quad x \in \mathbb{R}_+. \quad (2.11)$$

If  $\sigma^2$  is known but  $\mu$  unknown, then writing the density in canonical form as in (2.7) yields

$$p_{\theta}(x) \propto \exp \{ \theta t(x) + b(x) \}, \quad \theta \equiv \frac{\mu}{\sigma^2}, \quad t(x) \equiv x, \quad b(x) = -\frac{x^2}{2\sigma^2}.$$

Given an i.i.d. sample  $X_1, \dots, X_n \sim p_{\mu_0, \sigma^2}$ , the generalized  $h$ -score matching estimator of  $\mu$  is

$$\hat{\mu}_h \equiv \frac{\sum_{i=1}^n h(X_i) X_i - \sigma^2 h'(X_i)}{\sum_{i=1}^n h(X_i)}.$$

If  $\lim_{x \searrow 0^+} h(x) = 0$ ,  $\lim_{x \nearrow +\infty} h^2(x)(x - \mu_0)p_{\mu_0, \sigma^2}(x) = 0$  and the expectations are finite (for example, when  $h(x) = o(\exp(Mx^2))$  for  $M < \frac{1}{4\sigma^2}$ ), then

$$\sqrt{n}(\hat{\mu}_h - \mu_0) \rightarrow_d \mathcal{N} \left( 0, \frac{\mathbb{E}_0[\sigma^2 h^2(X) + \sigma^4 h'^2(X)]}{\mathbb{E}_0^2[h(X)]} \right).$$

We recall that the Cramér-Rao lower bound (i.e. the lower bound on the variance of any unbiased estimator) for estimating  $\mu$  is

$$\frac{\sigma^4}{\text{var}(X - \mu_0)}.$$

**Example 3.2.** Consider the univariate truncated normal distributions from (2.11) in the setting where the mean parameter  $\mu$  is known but the variance parameter  $\sigma^2 > 0$  is unknown. In canonical form as in (2.7), we write

$$p_\theta(x) \propto \exp \{ \theta t(x) + b(x) \}, \quad \theta \equiv \frac{1}{\sigma^2}, \quad t(x) \equiv -(x - \mu)^2/2, \quad b(x) = 0.$$

Given an i.i.d. sample  $X_1, \dots, X_n \sim p_{\mu, \sigma_0^2}$ , the generalized  $h$ -score matching estimator of  $\sigma^2$  is

$$\hat{\sigma}_h^2 \equiv \frac{\sum_{i=1}^n h(X_i)(X_i - \mu)^2}{\sum_{i=1}^n h(X_i) + h'(X_i)(X_i - \mu)}.$$

If, in addition to the assumptions in Example 3.1,  $\lim_{x \nearrow +\infty} h^2(x)(x - \mu)^3 p_{\mu, \sigma_0^2}(x) = 0$ , then

$$\sqrt{n}(\hat{\sigma}_h^2 - \sigma_0^2) \rightarrow_d \mathcal{N} \left( 0, \frac{2\sigma_0^6 \mathbb{E}_0[h^2(X)(X - \mu)^2] + \sigma_0^8 \mathbb{E}_0[h'^2(X)(X - \mu)^2]}{\mathbb{E}_0^2[h(X)(X - \mu)^2]} \right).$$

Moreover, the Cramér-Rao lower bound for estimating  $\sigma^2$  is

$$\frac{4\sigma_0^8}{\text{var}(X - \mu)^2}.$$

**Remark 2.** In Example 3.2, if  $\mu_0 = 0$ , then  $h(x) \equiv 1$  also satisfies (A1)-(A2) and (C1)-(C2) and one recovers the sample variance  $\frac{1}{n} \sum_i X_i^2$ , which obtains the Cramér-Rao lower bound.

In these examples, there is a benefit in using a bounded function  $h$ , which can be explained as follows. When  $\mu \gg \sigma$ , there is effectively no truncation to the Gaussian distribution, and our method adapts to using low moments in (2.4), since a bounded and increasing  $h(x)$  becomes almost constant as it reaches its asymptote for  $x$  large. Hence, we effectively revert to the original score matching (recall Section 2.2.1). In the other cases, the truncation effect is significant and our estimator uses higher moments accordingly.

Figure 2.1 plots the asymptotic variance of  $\hat{\mu}_h$  from Example 3.1, with  $\sigma = 1$  known. Efficiency as measured by the Cramér-Rao lower bound divided by the asymptotic variance is also shown. We see that two truncated versions of  $\log(1 + x)$  have asymptotic variance close to the Cramér-Rao bound. This asymptotic variance is also reflective of the variance for smaller finite samples.

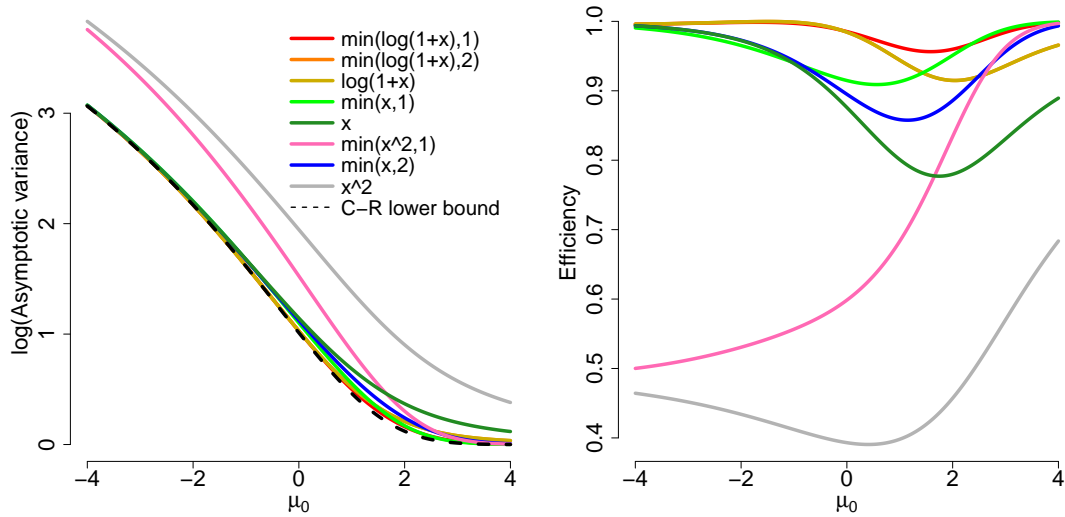


Figure 2.1: Log of asymptotic variance and efficiency with respect to the Cramér-Rao bound for  $\hat{\mu}_h$  ( $\sigma^2 = 1$  known).

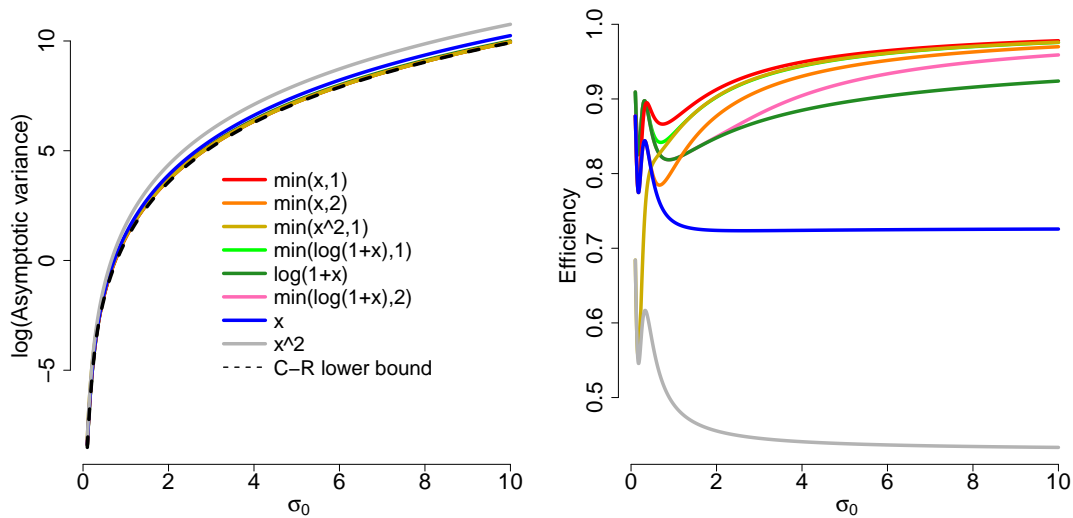


Figure 2.2: Log of asymptotic variance and efficiency with respect to the Cramér-Rao bound for  $\hat{\sigma}_h^2$  ( $\mu = 0.5$  known).

Figure 2.2 is the analog of Figure 2.1 for  $\hat{\sigma}_h^2$  from Example 3.2 with  $\mu = 0.5$  known. While the specifics are a bit different the benefits of using bounded or slowly growing  $h$  are again clear. We note that when  $\sigma$  is small, the effect of truncation to the positive part of the real line is small.

In both plots we order/color the curves based on their overall efficiency, so they have different colors in one from the other, although the same functions are presented. For all functions presented here (A1)–(A2) and (C1)–(C2) are satisfied.

## 2.4 Regularized Generalized Score Matching

In high-dimensional settings, when the number  $r$  of parameters to estimate may be larger than the sample size  $n$ , it is hard, if not impossible, to estimate the parameters consistently without turning to some form of regularization. More specifically, for exponential families, condition (C1) in Section 2.3 fails when  $r > n$ . A popular approach is then the use of  $\ell_1$  regularization to exploit possible sparsity.

Let the data matrix  $\mathbf{x} \in \mathbb{R}^{n \times m}$  comprise  $n$  i.i.d. samples from distribution  $P_0$ . Assume  $P_0$  has density  $p_0$  belonging to an exponential family  $\mathcal{P}_+ \equiv \{p_\theta : \theta \in \Theta\}$ , where  $\Theta \subseteq \mathbb{R}^r$ . Adding an  $\ell_1$  penalty to (2.8), we obtain the regularized generalized score matching loss

$$\frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{\Gamma}(\mathbf{x}) \boldsymbol{\theta} - \mathbf{g}(\mathbf{x})^\top \boldsymbol{\theta} + \lambda \|\boldsymbol{\theta}\|_1 \quad (2.12)$$

as in Lin et al. (2016). The loss in (2.12) involves a quadratic smooth part as in the familiar lasso loss for linear regression. However, although the matrix  $\boldsymbol{\Gamma}$  is positive semidefinite, the regularized loss in (2.12) is not guaranteed to be bounded unless the tuning parameter  $\lambda$  is sufficiently large—a problem that does not occur in lasso. We note that here, and throughout, we suppress the dependence on the data  $\mathbf{x}$  for  $\boldsymbol{\Gamma}(\mathbf{x})$ ,  $\mathbf{g}(\mathbf{x})$  and derived quantities.

For a more detailed explanation, note that by (2.9),  $\boldsymbol{\Gamma} = \mathbf{H}^\top \mathbf{H}$  for some  $\mathbf{H} \in \mathbb{R}^{nm \times r}$ . In the high-dimensional case, the rank of  $\boldsymbol{\Gamma}$ , or equivalently that of  $\mathbf{H}$ , is at most  $nm < r$ . Hence,  $\boldsymbol{\Gamma}$  is not invertible and  $\mathbf{g}$  does not necessarily lie in the column span of  $\boldsymbol{\Gamma}$ . Let  $\text{Ker}(\boldsymbol{\Gamma})$

be the kernel of  $\mathbf{\Gamma}$ . Then there may exist  $\boldsymbol{\nu} \in \text{Ker}(\mathbf{\Gamma})$  with  $\mathbf{g}^\top \boldsymbol{\nu} \neq 0$ . In this case, if

$$0 \leq \lambda < \sup_{\boldsymbol{\nu} \in \text{Ker}(\mathbf{\Gamma})} |\mathbf{g}^\top \boldsymbol{\nu}| / \|\boldsymbol{\nu}\|_1,$$

there exists  $\boldsymbol{\nu} \in \text{Ker}(\mathbf{\Gamma})$  with  $\frac{1}{2} \boldsymbol{\nu}^\top \mathbf{\Gamma} \boldsymbol{\nu} = 0$  and  $-\mathbf{g}^\top \boldsymbol{\nu} + \lambda \|\boldsymbol{\nu}\|_1 < 0$ . Evaluating at  $\boldsymbol{\theta}(a) = a \cdot \boldsymbol{\nu}$  for scalar  $a > 0$ , the loss becomes  $a(-\mathbf{g}^\top \boldsymbol{\nu} + \lambda \|\boldsymbol{\nu}\|_1)$ , which is negative and linear in  $a$ , and thus unbounded below. In this case no minimizer of (2.12) exists for small values of  $\lambda$ . This issue also exists for the estimators from Zhang and Zou (2014) and Liu and Luo (2015), which correspond to score matching for GGMs. We note that in the context of estimating the interaction matrix in pairwise models,  $r = m^2$ ; thus, the condition  $nm < r$  reduces to  $n < m$ , or  $n < m + 1$  when both  $\mathbf{K}$  and  $\boldsymbol{\eta}$  in (2.14) below are estimated.

To circumvent the unboundedness problem, we add small values  $\gamma_\ell > 0$  to the diagonal entries of  $\mathbf{\Gamma}$ , which become  $\mathbf{\Gamma}_{\ell,\ell} + \gamma_\ell$ ,  $\ell = 1, \dots, r$ . This is in the spirit of work such as Ledoit and Wolf (2004) and corresponds to an elastic net-type penalty (Zou and Hastie, 2005) with weighted  $\ell_2$  penalty  $\sum_{\ell=1}^r \gamma_\ell \theta_\ell^2$ . After this modification,  $\mathbf{\Gamma}$  is positive definite, our regularized loss is strongly convex in  $\boldsymbol{\theta}$ , and a unique minimizer exists for all  $\lambda \geq 0$ . For the special case of truncated GGMs, we will show that a result on consistent estimation holds if we choose  $\gamma_\ell = \delta_0 \mathbf{\Gamma}_{\ell,\ell}$  for a suitably small constant  $\delta_0 > 0$ , for which we propose a particular choice to avoid tuning. This choice of  $\gamma_\ell$  depends on the data through  $\mathbf{\Gamma}_{\ell,\ell}$ .

**Definition 2.** For  $\boldsymbol{\gamma} \in \mathbb{R}_+^r \setminus \{\mathbf{0}\}$ , let  $\mathbf{\Gamma}_\boldsymbol{\gamma} \equiv \mathbf{\Gamma} + \text{diag}(\boldsymbol{\gamma})$ . The regularized generalized  $\mathbf{h}$ -score matching estimator with tuning parameter  $\lambda \geq 0$  and amplifier  $\boldsymbol{\gamma}$  is the estimator

$$\hat{\boldsymbol{\theta}} \in \underset{\boldsymbol{\theta} \in \Theta}{\text{argmin}} \hat{J}_{\mathbf{h},\lambda,\boldsymbol{\gamma}}(\boldsymbol{\theta}) \equiv \underset{\boldsymbol{\theta} \in \Theta}{\text{argmin}} \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{\Gamma}_\boldsymbol{\gamma}(\mathbf{x}) \boldsymbol{\theta} - \mathbf{g}(\mathbf{x})^\top \boldsymbol{\theta} + \lambda \|\boldsymbol{\theta}\|_1. \quad (2.13)$$

In the case where  $\boldsymbol{\gamma} = (\delta - 1)\text{diag}(\mathbf{\Gamma})$  for some  $\delta > 1$ , we also call  $\delta$  the *multiplier*. We note that  $\hat{\boldsymbol{\theta}}$  from (2.13) is a *piecewise linear* function of  $\lambda$  (Lin et al., 2016).

## 2.5 Score Matching for Graphical Models for Non-negative Data

In this section we apply our generalized score matching estimator to a general class of graphical models for non-negative data.

### 2.5.1 A General Framework of Pairwise Interaction Models

We consider the class of pairwise interaction power models with density introduced in (1.1). We recall the form of the density:

$$p_{\boldsymbol{\eta}, \mathbf{K}}(\mathbf{x}) \propto \exp\left(-\frac{1}{2a} \mathbf{x}^{a\top} \mathbf{K} \mathbf{x}^a + \boldsymbol{\eta}^\top \frac{\mathbf{x}^b - \mathbf{1}_m}{b}\right) \mathbb{1}_{\mathbb{R}_+^m}(\mathbf{x}), \quad (2.14)$$

where  $a$  and  $b$  are known constants, and the interaction matrix  $\mathbf{K}$  and the vector  $\boldsymbol{\eta}$  are parameters. When  $b = 0$ , we use the convention that  $\frac{x^0 - 1}{0} \equiv \log x$  and apply the logarithm element-wise. Our focus will be on the interaction matrix  $\mathbf{K}$  that determines the conditional independence graph through its support  $S(\mathbf{K}) \equiv \{(i, j) : \kappa_{ij} \neq 0\}$ . However, unless  $\boldsymbol{\eta}$  is known or assumed to be zero, we also need to estimate  $\boldsymbol{\eta}$  as a nuisance parameter. In the case where we assume  $\boldsymbol{\eta} \equiv \mathbf{0}$  is known (i.e. the linear part  $(\mathbf{x}^b - \mathbf{1}_m)/b$  is not present), we call the distribution (and the corresponding estimator) a *centered* distribution (estimator), in contrast to the general case termed *non-centered* when we assume  $\boldsymbol{\eta} \neq \mathbf{0}$  or unknown.

We first give a set of sufficient conditions for the density to be valid, i.e., the right-hand side of (2.14) to be integrable. The proof is given in Appendix A.1.3.

**Theorem 5.** *Define conditions*

(CC1)  $\mathbf{K}$  is strictly co-positive, i.e.,  $\mathbf{v}^\top \mathbf{K} \mathbf{v} > 0$  for all  $\mathbf{v} \in \mathbb{R}_+^m \setminus \{\mathbf{0}\}$ ;

(CC2)  $2a > b > 0$ ;

(CC3)  $a > 0$ ,  $b = 0$ , and  $\eta_j > -1$  for  $j = 1, \dots, m$  ( $\boldsymbol{\eta} \succ -\mathbf{1}_m$ ).

*In the non-centered case, if (CC1) and one of (CC2) and (CC3) holds, then the function on the right-hand side of (2.14) is integrable over  $\mathbb{R}_+^m$ . In the centered case, (CC1) and  $a > 0$  are sufficient.*

We emphasize that (CC1) is a weaker condition than positive definiteness. Criteria for strict co-positivity are discussed in Väliäho (1986).

### 2.5.2 Implementation for Different Models

In this section we give some implementation details for the regularized generalized  $\mathbf{h}$ -score matching estimator defined in (2.13) applied to the pairwise interaction models from (2.14). We let  $\Psi \equiv (\mathbf{K}^\top, \boldsymbol{\eta})^\top \in \mathbb{R}^{(m+1) \times m}$ . The unregularized loss is then

$$\hat{J}_{\mathbf{h}}(P) = \frac{1}{2} \text{vec}(\Psi)^\top \Gamma(\mathbf{x}) \text{vec}(\Psi) - \mathbf{g}(\mathbf{x})^\top \text{vec}(\Psi).$$

The general form of the matrix  $\Gamma$  and the vector  $\mathbf{g}$  in the loss were given in equations (2.8)–(2.10). Here  $\Gamma \in \mathbb{R}^{(m+1)m \times (m+1)m}$  is block-diagonal, with the  $j$ -th  $\mathbb{R}^{(m+1) \times (m+1)}$  block

$$\begin{aligned} \Gamma_j(\mathbf{x}) &\equiv \begin{bmatrix} \Gamma_{11,j} & \Gamma_{12,j} \\ \Gamma_{12,j}^\top & \Gamma_{22,j} \end{bmatrix} \\ &\equiv \frac{1}{n} \sum_{i=1}^n \begin{bmatrix} h_j(X_j^{(i)}) X_j^{(i)2a-2} \mathbf{X}^{(i)a} \mathbf{X}^{(i)a\top} & -h_j(X_j^{(i)}) X_j^{(i)a+b-2} \mathbf{X}^{(i)a} \\ -h_j(X_j^{(i)}) X_j^{(i)a+b-2} \mathbf{X}^{(i)a\top} & h_j(X_j^{(i)}) X_j^{(i)2b-2} \end{bmatrix} \\ &= \frac{1}{n} \mathbf{y}^\top \mathbf{y}, \quad \mathbf{y} \equiv \left[ -\left(\sqrt{\mathbf{h}_j(\mathbf{X}_j)} \odot \mathbf{X}_j^{a-1}\right) \odot \mathbf{x}^a \quad \sqrt{\mathbf{h}_j(\mathbf{X}_j)} \odot \mathbf{X}_j^{b-1} \right] \in \mathbb{R}^{n,m+1}, \end{aligned} \quad (2.15)$$

where the  $\odot$  product between a vector and a matrix means an elementwise multiplication of the vector with each *column* of the matrix, and  $\mathbf{h}_j(\mathbf{X}_j) \equiv [h_j(X_j^{(1)}), \dots, h_j(X_j^{(n)})]^\top \in \mathbb{R}^m$ .

Furthermore,  $\mathbf{g} \equiv \begin{bmatrix} \text{vec}(\mathbf{g}_1) \\ \mathbf{g}_2 \end{bmatrix} \in \mathbb{R}^{(m+1)m}$ , where  $\mathbf{g}_1$  and  $\mathbf{g}_2$  correspond to each entry of  $\mathbf{K}$  and  $\boldsymbol{\eta}$ , respectively. The  $j$ -th column of  $\mathbf{g}_1 \in \mathbb{R}^{m \times m}$ , written as  $\mathbf{g}_{1,j}(\mathbf{x})$ , is

$$\frac{1}{n} \sum_{i=1}^n \left( h'_j(X_j^{(i)}) X_j^{(i)a-1} + (a-1)h_j(X_j^{(i)}) X_j^{(i)a-2} \right) \mathbf{X}^{(i)a} + ah_j(X_j^{(i)}) X_j^{(i)2a-2} \mathbf{e}_{j,m},$$

where  $\mathbf{e}_{j,m}$  is the  $m$ -vector with 1 at the  $j$ -th position and 0 elsewhere, and the  $j$ -th entry of  $\mathbf{g}_2 \in \mathbb{R}^m$  is

$$g_{2,j} = \frac{1}{n} \sum_{i=1}^n -h'_j(X_j^{(i)}) X_j^{(i)b-1} - (b-1)h_j(X_j^{(i)}) X_j^{(i)b-2}.$$

These formulae also hold for  $b = 0$  since  $\Gamma$  and  $\mathbf{g}$  only depend on the gradient of the log density, and  $\frac{d(x^b-1)/b}{dx} = x^{b-1}$  also holds for  $b = 0$ . In the centered case where we know

$\boldsymbol{\eta}_0 \equiv \mathbf{0}$ , we only estimate  $\mathbf{K} \in \mathbb{R}^{m \times m}$ , and  $\boldsymbol{\Gamma} \in \mathbb{R}^{m^2 \times m^2}$  is still block-diagonal, with the  $j$ -th block being the  $\boldsymbol{\Gamma}_{11,j}$  submatrix in (2.15), while  $\mathbf{g}$  is just  $\text{vec}(\mathbf{g}_1)$ . Since  $b$  only appears in the  $\boldsymbol{\eta}$  part of the density, the formulae only depend on  $a$  in the centered case.

We emphasize that it is indeed necessary to introduce amplifiers  $\boldsymbol{\gamma} \succ \mathbf{0}$  or a multiplier  $\delta > 1$  in addition to the  $\ell_1$  penalty. It is clear from (2.16) that  $\text{rank}(\boldsymbol{\Gamma}_j) \leq \min\{n, m + 1\}$  (or  $\min\{n, m\}$  if centered). Thus,  $\boldsymbol{\Gamma}$  is non-invertible when  $n \leq m$  (or  $n < m$  if centered) and  $\mathbf{g}$  need not lie in its column span.

We claim that including amplifiers/multipliers for the submatrices  $\boldsymbol{\Gamma}_{11,j}$  only is sufficient for unique existence of a solution for all penalty parameters  $\lambda \geq 0$ . To see this, consider any nonzero vector  $\boldsymbol{\nu} \in \mathbb{R}^{m+1}$ . Partition it as  $\boldsymbol{\nu} \equiv (\boldsymbol{\nu}_1, \boldsymbol{\nu}_2)$  with  $\boldsymbol{\nu}_1 \in \mathbb{R}^m$ . Let  $\boldsymbol{\Gamma}_{j,\boldsymbol{\gamma}}$  be our amplified version of the matrix  $\boldsymbol{\Gamma}_j$  from (2.19), so

$$\boldsymbol{\Gamma}_{j,\boldsymbol{\gamma}} = \begin{pmatrix} \boldsymbol{\Gamma}_{11,j} + \text{diag}(\gamma_1, \dots, \gamma_m) & \boldsymbol{\Gamma}_{12,j} \\ \boldsymbol{\Gamma}_{12,j}^\top & \boldsymbol{\Gamma}_{22,j} \end{pmatrix}.$$

As  $\boldsymbol{\Gamma}_j$  itself is positive semidefinite, we find that if at least one of the first  $m$  entries of  $\boldsymbol{\nu}$  is nonzero then

$$\boldsymbol{\nu}^\top \boldsymbol{\Gamma}_{j,\boldsymbol{\gamma}} \boldsymbol{\nu} \geq \boldsymbol{\nu}^\top \boldsymbol{\Gamma}_j \boldsymbol{\nu} + \sum_{k=1}^m \nu_k^2 \gamma_k \geq \sum_{k=1}^m \nu_k^2 \gamma_k > 0.$$

If only the last entry of  $\boldsymbol{\nu}$  is nonzero then

$$\boldsymbol{\nu}^\top \boldsymbol{\Gamma}_{j,\boldsymbol{\gamma}} \boldsymbol{\nu} = \nu_{m+1}^2 \boldsymbol{\Gamma}_{22,j} > 0$$

almost surely; recall that  $\boldsymbol{\Gamma}_{22,j} = \frac{1}{n} \sum_{i=1}^n h_j \left( X_j^{(i)} \right) X_j^{2b-2}$ . We conclude that  $\boldsymbol{\Gamma}_{j,\boldsymbol{\gamma}}$  (and thus the entire amplified  $\boldsymbol{\Gamma}$ ) is a.s. positive definite, which ensures unique existence of the loss minimizer.

Given the formulae for  $\boldsymbol{\Gamma}$  and  $\mathbf{g}$ , one adds the  $\ell_1$  penalty on  $\boldsymbol{\Psi}$  to get the regularized loss (2.22). Our methodology readily accommodates two different choices of the penalty parameter  $\lambda$  for  $\mathbf{K}$  and  $\boldsymbol{\eta}$ . This is also theoretically supported for truncated GGMs, since if the ratio of the respective values  $\lambda_{\mathbf{K}}$  and  $\lambda_{\boldsymbol{\eta}}$  is fixed, the proof of the theorems in Section 2.6 can be easily modified by replacing  $\boldsymbol{\eta}$  by  $(\lambda_{\boldsymbol{\eta}}/\lambda_{\mathbf{K}})\boldsymbol{\eta}$ . To avoid picking two tuning parameters,

one may also choose to remove the penalty on  $\boldsymbol{\eta}$  altogether by profiling out  $\boldsymbol{\eta}$  and solve for  $\hat{\boldsymbol{\eta}} \equiv \boldsymbol{\Gamma}_{22}^{-1} \left( \mathbf{g}_2 - \boldsymbol{\Gamma}_{12}^\top \text{vec}(\hat{\mathbf{K}}) \right)$ , with  $\hat{\mathbf{K}}$  the minimizer of the profiled loss

$$\hat{J}_{\mathbf{h}, \lambda, \boldsymbol{\gamma}, \text{profile}}(\mathbf{K}) \equiv \frac{1}{2} \text{vec}(\mathbf{K})^\top \boldsymbol{\Gamma}_{\boldsymbol{\gamma}, 11.2} \text{vec}(\mathbf{K}) - (\mathbf{g}_1 - \boldsymbol{\Gamma}_{12} \boldsymbol{\Gamma}_{22}^{-1} \mathbf{g}_2)^\top \text{vec}(\mathbf{K}) + \lambda \|\mathbf{K}\|_1, \quad (2.17)$$

where the Schur complement  $\boldsymbol{\Gamma}_{\boldsymbol{\gamma}, 11.2} \equiv \boldsymbol{\Gamma}_{\boldsymbol{\gamma}, 11} - \boldsymbol{\Gamma}_{12} \boldsymbol{\Gamma}_{22}^{-1} \boldsymbol{\Gamma}_{12}^\top$  is a.s. positive definite such that the profiled estimator exists a.s. for all  $\lambda \geq 0$ . This profiled approach corresponds to choosing  $\lambda_\eta / \lambda_{\mathbf{K}} = 0$ . A detailed theoretical analysis of the profiled estimator is beyond the scope of this chapter, however. We note that in the other extreme, with  $\lambda_\eta / \lambda_{\mathbf{K}} = +\infty$ , the non-centered estimator reduces to the estimator from the centered case.

**Example 5.3.** *The truncated normal model comprises the density*

$$p_{\boldsymbol{\mu}, \mathbf{K}}(\mathbf{x}) \propto \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{K} (\mathbf{x} - \boldsymbol{\mu}) \right\} \mathbb{1}_{[0, \infty)^m}(\mathbf{x}). \quad (2.18)$$

*This corresponds to (2.14) with  $a = b = 1$ , and  $\boldsymbol{\eta} = \mathbf{K}\boldsymbol{\mu}$ . The  $j$ -th  $(m+1) \times (m+1)$  block of  $\boldsymbol{\Gamma}(\mathbf{x})$  is*

$$\frac{1}{n} \begin{bmatrix} \mathbf{x}^\top \text{diag}(\mathbf{h}_j(\mathbf{X}_j)) \mathbf{x} & -\mathbf{x}^\top \mathbf{h}_j(\mathbf{X}_j) \\ -\mathbf{h}_j(\mathbf{X}_j)^\top \mathbf{x} & \mathbf{h}_j(\mathbf{X}_j)^\top \mathbf{1}_n \end{bmatrix}. \quad (2.19)$$

*Partitioning the vector  $\mathbf{g}(\mathbf{x})$  into  $m$  subvectors  $\mathbf{g}_j(\mathbf{x}) \in \mathbb{R}^{m+1}$ , where the entries of  $\mathbf{g}_j(\mathbf{x})$  correspond to column  $\boldsymbol{\Psi}_j$ , the  $k$ -th entry of  $\mathbf{g}_j(\mathbf{x})$  is*

$$g_{jk}(\mathbf{x}) \equiv \begin{cases} \frac{1}{n} \sum_{i=1}^n h'_j \left( X_j^{(i)} \right) X_k^{(i)} & \text{if } k \leq m, k \neq j, \\ \frac{1}{n} \sum_{i=1}^n h'_j \left( X_j^{(i)} \right) X_k^{(i)} + h_j \left( X_j^{(i)} \right) & \text{if } k = j, \\ -\frac{1}{n} \sum_{i=1}^n h'_j \left( X_j^{(i)} \right) & \text{if } k = m+1. \end{cases} \quad (2.20)$$

**Example 5.4.** *The exponential square-root graphical model in Inouye et al. (2016) has*

$$p_{\boldsymbol{\eta}, \mathbf{K}}(\mathbf{x}) \propto \exp \left( -\sqrt{\mathbf{x}}^\top \mathbf{K} \sqrt{\mathbf{x}} + 2\boldsymbol{\eta}^\top \sqrt{\mathbf{x}} \right) \mathbb{1}_{[0, \infty)^m}(\mathbf{x}),$$

*which corresponds to (2.14) with  $a = b = 1/2$ . We refer to this as the exponential model. In this case, the  $j$ -th  $\mathbb{R}^{(m+1) \times (m+1)}$  block of  $\boldsymbol{\Gamma}$  is*

$$\boldsymbol{\Gamma}_j(\mathbf{x}) \equiv \frac{1}{n} \sum_{i=1}^n \frac{h_j \left( X_j^{(i)} \right)}{X_j^{(i)}} \begin{pmatrix} -\sqrt{\mathbf{X}^{(i)}} \\ 1 \end{pmatrix} \left( -\sqrt{\mathbf{X}^{(i)}}^\top, 1 \right)$$

and  $\mathbf{g} = \text{vec}(\mathbf{g}_0)$ , where the  $j$ -th column of  $\mathbf{g}_0 \in \mathbb{R}^{(m+1) \times m}$  is

$$\mathbf{g}_j(\mathbf{x}) \equiv \frac{1}{n} \sum_{i=1}^n \frac{2h'_j(X_j^{(i)}) X_j^{(i)} - h_j(X_j^{(i)})}{2X_j^{(i)3/2}} \begin{pmatrix} \sqrt{\mathbf{X}^{(i)}} \\ -1 \end{pmatrix} + \frac{h_j(X_j^{(i)})}{2X_j^{(i)}} \mathbf{e}_{j,m+1}.$$

**Example 5.5.** If  $a = 1/2$  and  $b = 0$ , then (2.14) becomes

$$p_{\boldsymbol{\eta}, \mathbf{K}}(\mathbf{x}) \propto \exp\left(-\sqrt{\mathbf{x}}^\top \mathbf{K} \sqrt{\mathbf{x}} + \boldsymbol{\eta}^\top \log(\mathbf{x})\right) \mathbf{1}_{(0, \infty)^m}(\mathbf{x}). \quad (2.21)$$

If  $\mathbf{K}$  is diagonal in this case, then  $\mathbf{X} \sim p_{\boldsymbol{\eta}, \mathbf{K}}$  has independent entries with  $X_j$  following the gamma distribution with rate  $\kappa_{jj}$  and shape  $\eta_j + 1$ , which gives an intuition for condition (CC3)  $\eta_j > -1$  in Theorem 5. We can thus view (2.21) as a multivariate gamma distribution with pairwise interactions among the covariates, and call this the gamma model. For this model, the  $j$ -th block of  $\boldsymbol{\Gamma}$  is

$$\boldsymbol{\Gamma}_j(\mathbf{x}) \equiv \frac{1}{n} \sum_{i=1}^n \frac{h_j(X_j^{(i)})}{X_j^{(i)2}} \begin{pmatrix} -\sqrt{X_j^{(i)} \mathbf{X}^{(i)}} \\ 1 \end{pmatrix} \begin{pmatrix} -\sqrt{X_j^{(i)} \mathbf{X}^{(i)}}^\top \\ 1 \end{pmatrix}$$

and the part of  $\mathbf{g}$  corresponding to  $\mathbf{K}_j$  is

$$\mathbf{g}_{1,j}(\mathbf{x}) \equiv \frac{1}{n} \sum_{i=1}^n \frac{2h'_j(X_j^{(i)}) X_j^{(i)} - h_j(X_j^{(i)})}{2X_j^{(i)3/2}} \sqrt{\mathbf{X}^{(i)}} + \frac{h_j(X_j^{(i)})}{2X_j^{(i)}} \mathbf{e}_{j,m},$$

while the part for  $\eta_j$  is

$$g_{2,j}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{h_j(X_j^{(i)})}{X_j^{(i)2}} - \frac{h'_j(X_j^{(i)})}{X_j^{(i)}}.$$

We note that the  $\boldsymbol{\Gamma}_{11,j}$  sub-matrix of  $\boldsymbol{\Gamma}_j$  and the  $\mathbf{g}_{1,j}$  sub-vector of  $\mathbf{g}_j$  for the gamma model are the same as those for the exponential model, since  $a = 1/2$  in both cases and the parts involving  $\mathbf{K}$  in the densities are the same.

### 2.5.3 Computational Details

In the most general exponential family setting, as in Eq. (2.8)–(2.10) in Theorem 3, the time complexity for forming  $\boldsymbol{\Gamma} \in \mathbb{R}^{r \times r}$  and  $\mathbf{g} \in \mathbb{R}^r$  is  $\mathcal{O}(nm(f_{b'}(m) + r^2 + r(f_{\nu'}(m) + f_{\nu''}(m))))$ .

Here  $f_b(m)$  is the average time complexity for calculating  $\partial_j b(\mathbf{x})$  over  $j = 1, \dots, m$ , and similarly  $f_{t'}(m)$  for  $\partial_j t_\ell(\mathbf{x})$  and  $f_{t''}(m)$  for  $\partial_{jj} t_\ell(\mathbf{x})$  over  $j = 1, \dots, m$  and  $\ell = 1 \dots, r$ . In many applications, however, these three functions would be constant in  $m$ , thus giving an  $\mathcal{O}(nmr^2)$  computational complexity, with the dominating term coming from the operations for  $\mathbf{t}'_j \mathbf{t}'_j{}^\top$  in  $\mathbf{\Gamma}$  since  $\mathbf{\Gamma}$  is of dimension  $r \times r$ .

For pairwise interaction power models,  $r = m^2$  and the formula above becomes  $\mathcal{O}(nm^5)$ . However, since  $\mathbf{\Gamma}$  is block-diagonal with only  $m^3$  nonzero entries and by the special form of  $\mathbf{t}(\mathbf{x}) = \mathbf{x}^a \mathbf{x}^{a\top}$ , the true complexity is in fact  $\mathcal{O}(nm^3)$ .

While the introduction of the  $\ell_1$  penalty inevitably precludes the estimator from having a closed-form solution and introduces non-differentiability, state-of-art numerical optimization algorithms, such as coordinate-descent (Friedman et al., 2007), can be applied for fast estimation. To speed up estimation, one can usually use warm starts using the solution from the previous  $\lambda$ 's, as well as lasso-type strong screening rules (Tibshirani et al., 2012) to eliminate components of  $\hat{\boldsymbol{\theta}}$  that are known a priori to have zero estimates.

In our implementation for pairwise interaction models of Section 2.5.1 (available in the `genscore` R package), we optimize our loss functions with respect to a symmetric matrix  $\hat{\mathbf{K}}$ ; in the non-centered case the vector  $\hat{\boldsymbol{\eta}}$  is also included. We use a coordinate-descent method analogous to Algorithm 2 in Lin et al. (2016), where in each step we update each element of  $\hat{\mathbf{K}}$  and  $\hat{\boldsymbol{\eta}}$  based on the other entries from the previous steps, while maintaining symmetry. In our simulations in Section 2.7 we always scale the data matrix by column  $\ell_2$  norms before proceeding to estimation. Note that estimation of  $\hat{\mathbf{K}}$  without symmetry can be parallelized as the loss can be decomposed into a sum over the columns.

#### 2.5.4 Choice of the Function $\mathbf{h}$

In this subsection we discuss the requirements on the function  $\mathbf{h}$  as well as some reasonable choices of  $\mathbf{h}$ .

*Requirements on  $\mathbf{h}$*

In Section 2.2.2, we presented two assumptions (A1) and (A2) under which the generalized score-matching loss is valid, i.e., the integration by parts is justified and Theorem 2 holds. In this section, we present some sufficient (and nearly necessary) requirements on  $\mathbf{h}$  such that (A1) and (A2) are satisfied.

**Definition 3.** *Suppose  $\mathbf{h} : \mathbb{R}_+^m \rightarrow \mathbb{R}_+^m$  with  $\mathbf{h}(\mathbf{x}) = (h_1(x_1), \dots, h_m(x_m))^\top$ . We write that  $\mathbf{h} \in \mathcal{H}_{a,b}$  (for simplicity we omit the dependency on  $m$ ) if for all  $j = 1, \dots, m$ :*

i)  $h_j$  is absolutely continuous in every bounded sub-interval of  $\mathbb{R}_+$ , and thus has derivative  $h'_j$  a.s.;

ii)  $h_j(x) > 0$  a.s. on  $\mathbb{R}_+$ ;

iii)  $h_j$  and  $h'_j$  are both bounded by some piecewise powers of  $x$  a.s. on  $\mathbb{R}_+$ ;

iv)  $\lim_{x \searrow 0^+} h_j(x)/x_j^q = 0$ , where  $q \equiv \begin{cases} \max\{1 - a, 1 - b\} & \text{if } b > 0, \\ 1 - \eta_{0,j} & \text{if } b = 0. \end{cases}$

**Theorem 6.** *Assume every  $P$  in the family of distribution  $\mathcal{P}_+$  satisfies (CC1)–(CC3) and thus has finite normalizing constants. If  $\mathbf{h} \in \mathcal{H}_{a,b}$ , then (A1) and (A2) are satisfied.*

In centered models, where  $\boldsymbol{\eta} \equiv \mathbf{0}$ , we can assume  $b = 2a$  and iv) in the definition of  $\mathcal{H}_{a,2a}$  has  $q = 1 - a$ . For truncated GGMs,  $a = b = 1$ , so iv) in Definition 3 is simply  $\lim_{x_j \searrow 0^+} h_j(x_j) = 0$ .

In the case of  $b = 0$ ,  $\boldsymbol{\eta}$  is an unknown parameter, and (CC3) requires each of its component to be greater than  $-1$ . If one has prior information on  $\boldsymbol{\eta}$  or restricts the parameter space for  $\boldsymbol{\eta}$ , the requirement reduces to  $h_j(x_j) = o(x_j^{1-\eta_{0,j}})$  as  $x_j \searrow 0^+$ . Otherwise, it suffices to require  $h_j(x_j) = o(x_j^2)$ . Note that this is only a condition for  $x_j \searrow 0^+$ , and the globally quadratic behavior of  $h_j(x_j) = x_j^2$  from the original score matching is not needed on the entire  $\mathbb{R}_+$ , leaving opportunities for improvements.

*Reasonable Choices of  $\mathbf{h}$*

Assume a common univariate  $h$  for all components in  $\mathbf{h}$ . Inspired by Theorem 6, we consider  $h$  that behaves like a power of  $x$  both as  $x \nearrow +\infty$  and as  $x \searrow 0^+$ . Since the requirements on the two tails are separate, we can choose  $h$  to be a piecewise defined function that joins two powers with possibly different degrees. In other words,  $h(x) = \min(x^{p_1}, cx^{p_2})$  for some powers  $p_1 \geq p_2 \geq 0$  and constant  $c > 0$ . Only one constant  $c$  is required since generalized score matching is invariant to scaling of  $h$ . In determining the exact power of  $p_1$  we have the following considerations:

a) In the centered case:

- (i) (A1) and (A2): Theorem 6 requires that  $p_1 \geq 1 - a$ .
- (ii) “Controlled  $\mathbf{\Gamma}$  and  $\mathbf{g}$  for  $\mathbf{x}^a$ ”: We propose avoiding poles at the origin for the entries of  $\mathbf{\Gamma}$  and  $\mathbf{g}$ . The formula for  $\mathbf{\Gamma}_{11}$  in (2.16) shows that to this end  $\sqrt{h(x)}x^{a-1}$  needs to have a non-negative degree. This requires  $p_1 \geq 2 - 2a$ . The formula for  $\mathbf{g}_1$  similarly shows that  $h'(x)x^{a-1}$ ,  $h(x)x^{a-2}$  and  $h(x)x^{2a-2}$  all need to have a non-negative degree for small  $x$ . This requires  $p_1 \geq 2 - a$ .

b) In the non-centered case, in addition to (i) and (ii),

- (iii) (A1) and (A2): Theorem 6 requires  $p_1 \geq \max\{1 - a, 1 - b\}$  for  $b > 0$ , or  $1 - \min_j \eta_{0,j}$  for  $b = 0$ .
- (iv) “Controlled  $\mathbf{\Gamma}$  and  $\mathbf{g}$  for  $\mathbf{x}^b$ ”: From the definition of  $\mathbf{\Gamma}_{22}$  and  $\mathbf{g}_2$  and by the same reasoning as above,  $\sqrt{h(x)}x^{b-1}$ ,  $h'(x)x^{b-1}$  and  $h(x)x^{b-2}$  need to be non-negative powers of  $x$ , thus requiring  $p_1 \geq \max\{2 - b, 2 - 2b\} = 2 - b$ .

The choice of  $p_2$ , is only relevant for large data points. Our main consideration is then merely how well  $\mathbf{\Gamma}$  and  $\mathbf{g}$  concentrate on their true population values (Theorem 7). From this perspective, our intuition is that  $p_2$  should be chosen small so that the tails of the

distributions of the entries of  $\mathbf{\Gamma}$  and  $\mathbf{g}$  are well-behaved. Thus, we can choose  $p_2 = 0$ , in which case  $h(x) = \min(x^{p_1}, c)$  is a truncated power.

### 2.5.5 Tuning Parameter Selection

By treating the unpenalized loss (i.e.,  $\lambda = 0$ ,  $\gamma = 0$ ) as a negative log-likelihood, we may use the extended Bayesian Information Criterion (eBIC) to choose the tuning parameter (Chen and Chen, 2008; Foygel and Drton, 2010). Consider the centered case as an example. Let  $\hat{S}^\lambda \equiv \{(i, j) : \hat{\kappa}_{ij}^\lambda \neq 0, i < j\}$ , where  $\hat{\mathbf{K}}^\lambda$  be the estimate associated with tuning parameter  $\lambda$ . The eBIC is then

$$\text{eBIC}(\lambda) = -n \text{vec}(\hat{\mathbf{K}})^\top \mathbf{\Gamma}(\mathbf{x}) \text{vec}(\hat{\mathbf{K}}) + 2n \mathbf{g}(\mathbf{x})^\top \text{vec}(\hat{\mathbf{K}}) + |\hat{S}^\lambda| \log n + 2 \log \binom{p(p-1)/2}{|\hat{S}^\lambda|},$$

where  $\hat{\mathbf{K}}$  can be either the original estimate associated with  $\lambda$ , or a refitted solution obtained by restricting the support to  $\hat{S}^\lambda$ .

We use the eBIC instead of the ordinary BIC (Bayesian Information Criterion) since the BIC tends to choose an overly complex model when the model space is large, as encountered in the high-dimensional setting. The extension in eBIC comes from the last term in the above display which can be motivated by a prior distribution under which the number of edges in the conditional independence graph is uniformly distributed; see also Żak-Szatkowska and Bogdan (2011) and Barber and Drton (2015).

## 2.6 Theory for Graphical Models

In our regularized generalized score matching framework, we introduced the amplifiers and the multipliers to address the inexistence problem. We also proposed using a general function  $\mathbf{h}$  in place of  $\mathbf{x}^2$  as a means to improve estimation accuracy. This section provides a theoretical analysis of these two aspects.

In Section 2.6.1, we present the theory for our regularized generalized score matching estimators for general pairwise interaction models before going into the details for the special

cases of (truncated) GGMs. Next, we show that a specific choice of amplifiers/multipliers yields consistent estimation without the need for tuning. This point is important even in the case of Gaussian models on all of  $\mathbb{R}^m$ . Therefore, in Section 2.6.2 we digress from non-negative data and consider the original score matching of Hyvärinen (2005) for centered Gaussian distributions. Finally, in Section 2.6.3, we derive probabilistic results for  $\hat{\Psi}$  based on Theorem 7, justifying the benefits of using a general bounded  $\mathbf{h}$  over  $\mathbf{x}^2$  in the non-negative setting. As the most important models from the class of pairwise interaction power models over  $\mathbb{R}_+^m$ , we only treat truncated GGMs since they have the most tractable concentration bounds; this case also provides a comparison to Corollary 2 in Lin et al. (2016), which uses  $\mathbf{x}^2$ .

### 2.6.1 Theory for Pairwise Interaction Models

The graphical models we treat are parametrized by the interaction matrix  $\mathbf{K}$  and the coefficients  $\boldsymbol{\eta}$  on  $(\mathbf{x}^b - \mathbf{1}_m)/b$ . It is convenient to accommodate this setting with a matrix-valued parameter  $\Psi \in \mathbb{R}^{r_1 \times r_2}$  (in place of  $\boldsymbol{\theta}$ ) and specify our regularized  $\mathbf{h}$ -score matching loss as

$$\hat{J}_{\mathbf{h}, \lambda, \gamma}(\Psi) \equiv \operatorname{argmin}_{\Psi \in \mathbb{R}^{r_1 \times r_2}} \frac{1}{2} \operatorname{vec}(\Psi)^\top \Gamma_\gamma(\mathbf{x}) \operatorname{vec}(\Psi) - \mathbf{g}(\mathbf{x})^\top \operatorname{vec}(\Psi) + \lambda \|\Psi\|_1. \quad (2.22)$$

In the non-centered case we thus take  $\Psi = [\mathbf{K}, \boldsymbol{\eta}]^\top \in \mathbb{R}^{m(m+1) \times m}$ . In the centered case,  $\Psi$  is simply the  $m \times m$  interaction matrix  $\mathbf{K}$ . Following related prior work such as Lin et al. (2016), for ease of proof we allow the matrix  $\mathbf{K}$  to be nonsymmetric, which allows us to decouple optimization over the different columns of  $\mathbf{K}$  or  $\Psi$ , while in our implementations we ensure that  $\mathbf{K}$  is symmetric.

**Definition 4.** Let  $\Gamma_0 \equiv \mathbb{E}_0 \Gamma(\mathbf{x})$  and  $\mathbf{g}_0 \equiv \mathbb{E}_0 \mathbf{g}(\mathbf{x})$  be the population versions of  $\Gamma(\mathbf{x})$  and  $\mathbf{g}(\mathbf{x})$  under the distribution given by a true parameter matrix  $\Psi_0$ . The support of a matrix  $\Psi$  is  $S(\Psi) \equiv \{(i, j) : \psi_{ij} \neq 0\}$ , and we let  $S_0 = S(\Psi_0)$ . For a matrix  $\Psi_0$ , we define  $d_{\Psi_0}$  to be the maximum number of non-zero entries in any column, and  $c_{\Psi_0} \equiv \|\Psi_0\|_{\infty, \infty}$ . Writing  $\Gamma_{0, AB}$  for the  $A \times B$  submatrix of  $\Gamma_0$ , we define

$$c_{\Gamma_0} \equiv \|\|(\Gamma_{0, S_0 S_0})^{-1}\|_{\infty, \infty}. \quad (2.23)$$

Finally,  $\mathbf{\Gamma}_0$  satisfies the irrepresentability condition with incoherence parameter  $\alpha \in (0, 1]$  and edge set  $S_0$  if

$$\|\mathbf{\Gamma}_{0, S_0^c S_0} (\mathbf{\Gamma}_{0, S_0 S_0})^{-1}\|_{\infty, \infty} \leq (1 - \alpha). \quad (2.24)$$

Our analysis of the regularized generalized  $\mathbf{h}$ -score matching estimator builds on the following theorem taken from Lin et al. (2016, Theorem 1).

**Theorem 7.** *Suppose  $\mathbf{\Gamma}_0$  has  $\mathbf{\Gamma}_{0, S_0 S_0}$  invertible and satisfies the irrepresentability condition (2.24) with incoherence parameter  $\alpha \in (0, 1]$ . Assume that*

$$\|\mathbf{\Gamma}_\gamma(\mathbf{x}) - \mathbf{\Gamma}_0\|_\infty < \epsilon_1, \quad \|\mathbf{g}(\mathbf{x}) - \mathbf{g}_0\|_\infty < \epsilon_2, \quad (2.25)$$

with  $d_{\Psi_0} \epsilon_1 \leq \alpha / (6c_{\mathbf{\Gamma}_0})$ . If

$$\lambda > \frac{3(2 - \alpha)}{\alpha} \max\{c_{\Psi_0} \epsilon_1, \epsilon_2\},$$

then the following holds:

- (a) *The regularized generalized  $\mathbf{h}$ -score matching estimator  $\hat{\Psi}$  minimizing (2.22) is unique, with support  $\hat{S} \equiv S(\hat{\Psi}) \subseteq S_0$ , and satisfies*

$$\|\hat{\Psi} - \Psi_0\|_\infty \leq \frac{c_{\mathbf{\Gamma}_0}}{2 - \alpha} \lambda.$$

- (b) *If*

$$\min_{1 \leq j < k \leq m} |\Psi_{0, jk}| > \frac{c_{\mathbf{\Gamma}_0}}{2 - \alpha} \lambda,$$

*then  $\hat{S} = S_0$  and  $\text{sign}(\hat{\Psi}_{jk}) = \text{sign}(\Psi_{0, jk})$  for all  $(j, k) \in S_0$ .*

This result is deterministic, and the improvement of our generalized estimator over the one in Lin et al. (2016) is in its probabilistic guarantees, as shown for truncated GGMs in Theorems 10 and 11 in Section 2.6.3. Before going into these examples, we state a general corollary.

**Corollary 8.** *Under the assumptions of Theorem 7, the matrix  $\hat{\Psi}$  minimizing (2.22) satisfies*

$$\begin{aligned} \|\hat{\Psi} - \Psi_0\|_F &\leq \frac{c_{\mathbf{\Gamma}_0}}{2 - \alpha} \lambda \sqrt{|S_0|} \leq \frac{c_{\mathbf{\Gamma}_0}}{2 - \alpha} \lambda \sqrt{d_{\Psi_0} m}, \\ \|\hat{\Psi} - \Psi_0\|_2 &\leq \frac{c_{\mathbf{\Gamma}_0}}{2 - \alpha} \lambda \min(\sqrt{|S_0|}, d_{\Psi_0}). \end{aligned}$$

### 2.6.2 Revisiting Gaussian Score Matching

In this section we consider estimating the inverse covariance matrix  $\mathbf{K}$  of a centered Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{K})$ , which has density proportional to (1.2) on all of  $\mathbb{R}^m$ . As shown, e.g., in Example 1 of Lin et al. (2016), the  $\ell_1$ -regularized score matching loss then takes the form

$$\frac{1}{2}\text{tr}(\mathbf{K}\mathbf{K}\mathbf{x}\mathbf{x}^\top) - \text{tr}(\mathbf{K}) + \lambda\|\mathbf{K}\|_1, \quad (2.26)$$

which can be written as (2.12) with  $\boldsymbol{\theta} = \text{vec}(\mathbf{K})$ ,  $\boldsymbol{\Gamma} = \text{diag}(\mathbf{x}\mathbf{x}^\top, \dots, \mathbf{x}\mathbf{x}^\top)$  and  $\mathbf{g} = \text{vec}(\mathbf{I}_m)$ . Thus, in general, the kernel of  $\boldsymbol{\Gamma}$  need not be orthogonal to  $\mathbf{g}$ , and for  $\lambda$  small the loss can be unbounded from below as discussed above. Hence, an amplifier/multiplier on the diagonals of  $\boldsymbol{\Gamma}$  is needed. We have the following theorem on the estimator using the amplification.

**Theorem 9.** *Suppose the data matrix  $\mathbf{x}$  holds  $n$  i.i.d. copies of  $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_0)$ . Adopt the amplifying in Section 2.4 and redefine the loss in (2.26) as*

$$\frac{1}{2}\text{tr}(\mathbf{K}\mathbf{K}\mathbf{G}) - \text{tr}(\mathbf{K}) + \lambda\|\mathbf{K}\|_1, \quad \mathbf{G}_{jk} = (\mathbf{x}\mathbf{x}^\top)_{jk} (\mathbf{1}_{\{j \neq k\}} + (\delta - 1)\mathbf{1}_{\{j=k\}}), \quad (2.27)$$

where  $1 < \delta < 2 - \left(1 + 80\sqrt{\log m/n}\right)^{-1}$ . Let  $\hat{\mathbf{K}}$  be the resulting estimator. Let  $c^* \equiv 12800 (\max_j \boldsymbol{\Sigma}_{0,jj})^2$  and  $c_1 = 4c_{\mathbf{r}_0}/\alpha$ . If for some  $\tau > 2$ , the regularization parameter and the sample size satisfy

$$\begin{aligned} \lambda &> 2c_{\mathbf{K}_0}(2 - \alpha)\sqrt{c^*(\tau \log m + \log 4)/n/\alpha}, \\ n &> \max(c^*c_1^2d_{\mathbf{K}_0}^2, 2)(\tau \log m + \log 4), \end{aligned}$$

then  $\|\hat{\mathbf{K}} - \mathbf{K}_0\|_\infty \leq \frac{c_{\mathbf{r}_0}}{2-\alpha}\lambda$  with probability  $1 - m^{2-\tau}$ .

In Corollary 1 of Lin et al. (2016) the same results were shown with  $c^* \equiv 3200 (\max_j \boldsymbol{\Sigma}_{0,jj})^2$  when a unique minimizer exists, but the existence was not guaranteed.

### 2.6.3 Generalized Score Matching for Truncated GGMS

Next, we provide theory for the regularized generalized  $\mathbf{h}$ -score matching estimator  $\hat{\boldsymbol{\Psi}}$  in the special case of truncated GGMS. Again, assume a common  $h$  for all components in  $\mathbf{h}$ .

**Theorem 10.** *Suppose the data matrix  $\mathbf{x}$  holds  $n$  i.i.d. copies of  $\mathbf{X} \sim \text{TN}(\mathbf{0}, \mathbf{K}_0)$ , where the mean parameter is known to be zero. Assume that  $\mathbf{h} \in \mathcal{H}_{1,1}$  and that  $0 \leq h \leq M$ ,  $0 \leq h' \leq M'$  a.s. for constants  $M, M'$ , and choose  $\boldsymbol{\gamma} = (\delta - 1)\text{diag}(\boldsymbol{\Gamma})$  with*

$$1 < \delta < C(n, m) \equiv 2 - \left(1 + 4e \max\{6 \log m/n, \sqrt{6 \log m/n}\}\right)^{-1}.$$

*Suppose that the  $\boldsymbol{\Gamma}_{0, S_0 S_0}$  block of  $\boldsymbol{\Gamma}_0$  is invertible and  $\boldsymbol{\Gamma}_0$  satisfies the irrepresentability condition (2.24) with  $\alpha \in (0, 1]$  and true edge set  $S_0$ . Define  $c_{\mathbf{X}} \equiv 2 \max_j \left(2\sqrt{(\mathbf{K}_0^{-1})_{jj}} + \sqrt{e} \mathbb{E}_0 X_j\right)$ . If for  $\tau > 3$  the sample size and the regularization parameter satisfy*

$$n > \mathcal{O} \left( \tau \log m \max \left\{ \frac{M^2 c_{\boldsymbol{\Gamma}_0}^2 c_{\mathbf{X}}^4 d_{\mathbf{K}_0}^2}{\alpha^2}, \frac{M c_{\boldsymbol{\Gamma}_0} c_{\mathbf{X}}^2 d_{\mathbf{K}_0}}{\alpha} \right\} \right), \quad (2.28)$$

$$\lambda > \mathcal{O} \left[ (M c_{\mathbf{K}_0} c_{\mathbf{X}}^2 + M' c_{\mathbf{X}} + M) \left( \sqrt{\frac{\tau \log m}{n}} + \frac{\tau \log m}{n} \right) \right], \quad (2.29)$$

*then the following statements hold with probability  $1 - m^{3-\tau}$ :*

- (a) *The regularized generalized  $\mathbf{h}$ -score matching estimator  $\hat{\mathbf{K}}$  that minimizes (2.22) is unique, has its support included in the true support,  $\hat{S} \equiv S(\hat{\mathbf{K}}) \subseteq S_0$ , and satisfies*

$$\begin{aligned} \|\hat{\mathbf{K}} - \mathbf{K}_0\|_{\infty} &\leq \frac{c_{\boldsymbol{\Gamma}_0}}{2 - \alpha} \lambda, \\ \|\hat{\mathbf{K}} - \mathbf{K}_0\|_F &\leq \frac{c_{\boldsymbol{\Gamma}_0}}{2 - \alpha} \lambda \sqrt{|S_0|}, \\ \|\hat{\mathbf{K}} - \mathbf{K}_0\|_2 &\leq \frac{c_{\boldsymbol{\Gamma}_0}}{2 - \alpha} \lambda \min(\sqrt{|S_0|}, d_{\mathbf{K}_0}), \end{aligned}$$

*where  $c_{\boldsymbol{\Gamma}_0}$  is defined in (2.23).*

- (b) *Moreover, if*

$$\min_{j,k:(j,k) \in S_0} |\kappa_{0,jk}| > \frac{c_{\boldsymbol{\Gamma}_0}}{2 - \alpha} \lambda,$$

*then  $\hat{S} = S_0$  and  $\text{sign}(\hat{\kappa}_{jk}) = \text{sign}(\kappa_{0,jk})$  for all  $(j, k) \in S_0$ .*

The theorem is proved in Appendix A.1.4, where details on the dependencies on constants are provided. A key ingredient of the proof is a tail bound on  $\|\boldsymbol{\Gamma}_{\boldsymbol{\gamma}} - \boldsymbol{\Gamma}_0\|_{\infty}$ , which

features products of the  $X_j^{(i)}$ 's. In Lin et al. (2016), the products are up to fourth order. Using bounded  $\mathbf{h}$ , our products automatically calibrates to a quadratic polynomial when the observed values are large, and resort to higher moments only when they are small. This leads to improved bounds and convergence rates, underscored in the new requirement on the sample size  $n$ , which should be compared to  $n \geq \mathcal{O}(d_{\mathbf{K}_0}^2 (\log m^\tau)^8)$  in Lin et al. (2016).

For the non-centered case, by definition,  $c_{\Psi_0} \equiv \|\Psi_0^\top\|_{\infty, \infty} \leq c_{\mathbf{K}_0} + \|\boldsymbol{\eta}_0\|_\infty$ ,  $d_{\Psi_0} \leq d_{\mathbf{K}_0} + 1$ . The proof given for Theorem 10 goes through again here, and we have the following consistency results.

**Theorem 11.** *Suppose the data matrix holds  $n$  i.i.d. copies of  $\mathbf{X} \sim \text{TN}(\boldsymbol{\mu}_0, \mathbf{K}_0)$ . Assume that  $\mathbf{h} \in \mathcal{H}_{1,1}$  and that  $0 \leq h \leq M$ ,  $0 \leq h' \leq M'$  a.s. for constants  $M, M'$ . Let  $\boldsymbol{\gamma}$  be a vector of amplifiers that are non-zero only for the diagonal entries of the matrices  $\boldsymbol{\Gamma}_{11,j}$ , amplifying those by  $(\delta - 1)\text{diag}(\boldsymbol{\Gamma}_{11,j})$  with*

$$1 < \delta < C(n, m) \equiv 2 - \left(1 + 4e \max\{6 \log m/n, \sqrt{6 \log m/n}\}\right)^{-1}.$$

*Suppose further that  $\boldsymbol{\Gamma}_{0, S_0 S_0}$  is invertible and satisfies the irrepresentability condition (2.24) with  $\alpha \in (0, 1]$ . Define  $c_{\mathbf{X}} \equiv 2 \max_j \left(2\sqrt{(\mathbf{K}_0^{-1})_{jj}} + \sqrt{e} \mathbb{E}_0 X_j\right)$ . Suppose for  $\tau > 3$  the sample size and the regularization parameter satisfy*

$$n > \mathcal{O} \left( \tau \log m \max \left\{ \frac{M^2 c_{\boldsymbol{\Gamma}_0, \Psi_0}^2 c_{\mathbf{X}}^4 d_{\Psi_0}^2}{\alpha^2}, \frac{M c_{\boldsymbol{\Gamma}_0, \Psi_0} c_{\mathbf{X}}^2 d_{\Psi_0}}{\alpha} \right\} \right), \quad (2.30)$$

$$\lambda > \mathcal{O} \left[ (M c_{\Psi_0} c_{\mathbf{X}}^2 + M' c_{\mathbf{X}} + M) \left( \sqrt{\frac{\tau \log m}{n}} + \frac{\tau \log m}{n} \right) \right], \quad (2.31)$$

*where  $c_{\boldsymbol{\Gamma}_0, \Psi_0}$  is  $c_{\boldsymbol{\Gamma}_0}$  as in (2.23) but with notation  $\Psi_0$  to differentiate it from the centered case. Then the following statements hold with probability  $1 - m^{3-\tau}$ :*

- (a) *The regularized generalized  $\mathbf{h}$ -score matching estimator  $\hat{\Psi}$  that minimizes (2.22) is unique, has its support included in the true support,  $\hat{S} \equiv S(\hat{\Psi}) \subseteq S_0$ , and satisfies*

$$\|\hat{\mathbf{K}} - \mathbf{K}_0\|_\infty \leq \frac{c_{\boldsymbol{\Gamma}_0, \Psi_0}}{2 - \alpha} \lambda, \quad \|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\|_\infty \leq \frac{c_{\boldsymbol{\Gamma}_0, \Psi_0}}{2 - \alpha} \lambda,$$

$$\begin{aligned} \|\hat{\mathbf{K}} - \mathbf{K}_0\|_F &\leq \frac{c_{\Gamma_0, \Psi_0}}{2 - \alpha} \lambda \sqrt{|S_0|}, & \|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\|_F &\leq \frac{c_{\Gamma_0, \Psi_0}}{2 - \alpha} \lambda \sqrt{|S_0|}, \\ \|\hat{\mathbf{K}} - \mathbf{K}_0\|_2 &\leq \frac{c_{\Gamma_0, \Psi_0}}{2 - \alpha} \lambda \min\left(\sqrt{|S_0|}, d_{\Psi_0}\right), & \|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\|_2 &\leq \frac{c_{\Gamma_0, \Psi_0}}{2 - \alpha} \lambda \min\left(\sqrt{|S_0|}, d_{\Psi_0}\right). \end{aligned}$$

(b) Moreover, if

$$\min_{j,k:(j,k) \in S_0} |\kappa_{0,jk}| > \frac{c_{\Gamma_0}}{2 - \alpha} \lambda \quad \text{and} \quad \min_{j:(m+1,j) \in S_0} |\eta_{0,j}| > \frac{c_{\Gamma_0}}{2 - \alpha} \lambda,$$

then  $\hat{S} = S_0$  and  $\text{sign}(\hat{\kappa}_{jk}) = \text{sign}(\kappa_{0,jk})$  for all  $(j, k) \in S_0$  and  $\text{sign}(\hat{\eta}_j) = \text{sign}(\eta_{0j})$  for  $(m + 1, j) \in S_0$ .

**Remark 3.** The quantity  $c_{\mathbf{X}}$  in Theorem 11 depends on  $\mathbb{E}_0 X_j$ , which in turn depends on the structure of both  $\boldsymbol{\mu}_0$  and  $\mathbf{K}_0$ . If  $\mu_{0,j}$  is large compared to  $(\mathbf{K}_0)_{jj}^{-1}$ , then  $c_{\mathbf{X}}$  seems to scale with  $\boldsymbol{\mu}_0$ , which negatively impacts the guarantees stated in Theorem 11. However, as in the one-dimensional case for estimation of  $\mu_0$  (Example 3.1), our estimator should automatically adapt to the large mean parameter. This suggests that it might be possible to improve our analysis involving  $c_{\mathbf{X}}$ .

## 2.7 Numerical Experiments

In this section, we compare the performance of our estimator with different choices of  $\mathbf{h}$  to the existing approaches for pairwise interaction power models. In our simulation experiments, we consider  $m = 100$  variables and  $n = 80$  and  $n = 1000$  samples, corresponding to high- and low-dimensional settings. We also tried intermediate sample sizes between these two extremes, but found no interesting result worth reporting. For  $n = 80$ , amplification is necessary. Except in Section 2.7.2, the amplifier is set based on Theorem 10 to  $\delta = C(n, m) = 1.8647$  for truncated GGMs. The same amplifier is also used for settings with other  $a$  and  $b$ . For  $n = 1000$ , we consider  $\delta = 1$ , i.e., no amplification, and  $\delta = C(n, m) = 1.6438$  (again, based on Theorem 10). Throughout, we assume a common univariate  $h$  for all components in  $\mathbf{h}$ .

### 2.7.1 Structure of $\mathbf{K}$

The underlying interaction matrices are selected as follows: Proceeding as in Section 4.2 of Lin et al. (2016), the graph is chosen to have 10 disconnected subgraphs, each containing  $m/10$  nodes. Thus,  $\mathbf{K}_0$  is block-diagonal. In each block, each lower-triangular element is set to 0 with probability  $1 - \pi$  for some  $\pi \in (0, 1)$ , and is otherwise drawn from  $\text{Uniform}[0.5, 1]$ . The upper triangular elements are determined by symmetry. The diagonal elements of  $\mathbf{K}_0$  are chosen as a common positive value such that the minimum eigenvalue of  $\mathbf{K}_0$  is 0.1.

We generate 5 different true precision matrices  $\mathbf{K}_0$ , and run 10 trials with each of these precision matrices. For  $n = 1000$ , we choose  $\pi = 0.8$ , which is in accordance with Lin et al. (2016). For  $n = 80$ , we set  $\pi = 0.2$ . This way  $n/(d_{\mathbf{K}_0}^2 \log m)$  is roughly constant; recall Theorems 10 and 11 for truncated GGMs.

In Appendix A.3, we report results on *Erdős-Rényi graphs*, which lead to similar conclusions.

### 2.7.2 Truncated GGMs

Given our focus on truncated GGMs and their relevance in graphical modeling applications, we start with experiments for these models.

#### *Choice of $\mathbf{h}$*

Our estimator requires choosing a function  $\mathbf{h} : \mathbb{R}_+^m \rightarrow \mathbb{R}_+^m$ . For simplicity, we will always specify  $\mathbf{h}(\mathbf{x}) = (h(x_1), \dots, h(x_m))$  for a single non-decreasing univariate function  $h : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ , i.e. all coordinates share the same  $h$  function.

As previously explained,  $\mathbf{h} \in \mathcal{H}_{a,b}$  is a sufficient condition for assumptions (A1)-(A2), as well as (C1)-(C2) in the case of unregularized estimators. Only in the proofs of our theoretical guarantees in Section 2.6 for truncated GGMs, did we require  $h$  to be bounded and to have bounded derivatives. As motivated by the discussion in Section 2.5.4, we consider truncated and untruncated powers,  $\min(x, c)$  and  $x$  (since  $2 - a = 2 - b = 1$ ); we evaluate this choice

by contrasting them with powers  $x^{1.5}$  and  $x^2$ . We also explore functions like  $\log(1+x)$  that seem natural and are linear near 0. In particular, we make a further comparison to functions linear near 0 with a finite asymptote as  $x \nearrow +\infty$  but differentiable everywhere: MCP- (Fan and Li, 2001) and SCAD-like (Zhang, 2010) functions defined below. The results we report are based on selections of best performing choices of  $h$ .

$$\text{SCAD}(x; \lambda, \gamma) \equiv \begin{cases} \lambda x & \text{if } 0 \leq x \leq \lambda, \\ \frac{2\gamma\lambda x - x^2 - \lambda^2}{2(\gamma-1)} & \text{if } \lambda < x < \gamma\lambda, \\ \frac{\lambda^2(\gamma+1)}{2} & \text{if } x \geq \gamma\lambda; \end{cases} \quad \text{MCP}(x; \lambda, \gamma) \equiv \begin{cases} \lambda x - \frac{x^2}{2\gamma} & \text{if } 0 \leq x \leq \gamma\lambda, \\ \frac{1}{2}\gamma\lambda^2 & \text{if } x > \gamma\lambda. \end{cases}$$

We do not observe any clear relationship between features such as convexity, differentiability or the slope of  $h$  at 0, and performance of the estimator. Nonetheless, for many choices of rather simple functions  $h$ , our estimator provides a significant improvement over existing methods. In particular, most  $h$  functions that behave linearly for small  $x$ , namely  $\log(1+x)$  and  $x$  and their truncations, and additionally MCP and SCAD, always perform better than  $x^{1.5}$  and  $x^2$ . This agrees with our discussion in Section 2.5.4, where  $2-a=1$  is a reasonable choice of the power for small  $x$ ; also see Section 2.7.3. However, we conclude that there is no real gain from making the function smoother by using MCP or SCAD.

*Truncated Centered GGMS:* For data from a truncated centered Gaussian distribution, we compare our generalized score matching estimator with various choices of  $h$ , to *SpaCE JAM* (SJ, Voorman et al., 2014), which estimates graphs using additive models for conditional means, a pseudo-likelihood method *SPACE* (Peng et al., 2009) in the reformulation of Khare et al. (2015), *graphical lasso* (GLASSO, Yuan and Lin, 2007; Friedman et al., 2008), the *neighborhood selection* estimator (NS) of Meinshausen and Bühlmann (2006), and *non-paranormal SKEPTIC* (Liu et al., 2012) with Kendall's  $\tau$ . Recall that the choice of  $h(x) = x^2$  corresponds to the estimator from Lin et al. (2016).

Centered, $n = 80$ , multiplier 1.8647											
min(log(1 + x), c)						min(x, c)					
c	Mean	sd	c	Mean	sd	c	Mean	sd	c	Mean	sd
$\infty$	0.694	0.033	$\infty$	0.702	0.031	$\infty$	0.702	0.031	$\infty$	0.702	0.031
2	0.694	0.033	3	0.702	0.031	3	0.702	0.031	3	0.702	0.031
1	0.692	0.033	2	0.698	0.033	2	0.698	0.033	2	0.698	0.033
0.5	0.664	0.038	1	0.686	0.030	1	0.686	0.030	1	0.686	0.030
MCP(1, c)						SCAD(1, c)					
c	Mean	sd	c	Mean	sd	c	Mean	sd	c	Mean	sd
10	0.701	0.032	10	0.702	0.031	10	0.702	0.031	10	0.702	0.031
5	0.700	0.032	5	0.701	0.032	5	0.701	0.032	5	0.701	0.032
1	0.672	0.036	2	0.696	0.033	2	0.696	0.033	2	0.696	0.033
$x^{1.5}$ : (0.683, 0.030)						$x^2$ : (0.630, 0.029)					
GLASSO (0.600,0.032)						SPACE: (0.587, 0.031)					
NS: (0.587,0.031)						SJ: (0.540,0.036)					

Centered, $n = 1000$ , multiplier 1						Centered, $n = 1000$ , multiplier 1.6438					
min(log(1 + x), c)			min(x, c)			min(log(1 + x), c)			min(x, c)		
c	Mean	sd	c	Mean	sd	c	Mean	sd	c	Mean	sd
2	0.826	0.015	2	0.820	0.014	$\infty$	0.857	0.011	3	0.855	0.011
$\infty$	0.826	0.015	3	0.820	0.015	2	0.857	0.011	$\infty$	0.855	0.011
1	0.824	0.014	$\infty$	0.819	0.015	1	0.855	0.011	2	0.854	0.011
0.5	0.804	0.015	1	0.817	0.014	0.5	0.833	0.012	1	0.847	0.011
MCP(1, c)			SCAD(1, c)			MCP(1, c)			SCAD(1, c)		
c	Mean	sd	c	Mean	sd	c	Mean	sd	c	Mean	sd
5	0.824	0.015	2	0.823	0.014	5	0.857	0.011	5	0.856	0.011
10	0.822	0.015	5	0.822	0.015	10	0.856	0.011	10	0.855	0.011
1	0.810	0.015	10	0.821	0.015	1	0.840	0.012	2	0.855	0.011
$x^{1.5}$ : (0.782,0.014)			$x^2$ : (0.732,0.015)			$x^{1.5}$ : (0.812,0.011)			$x^2$ : (0.736,0.011)		
SPACE: (0.780,0.015)			NS: (0.779,0.015)			SPACE: (0.780,0.015)			NS: (0.779,0.015)		
GLASSO (0.764,0.014)			SJ: (0.703,0.015)			GLASSO (0.764,0.014)			SJ: (0.703,0.015)		

Table 2.1: Mean and standard deviation of areas under the ROC curves (AUC) using different estimators in the centered setting, with  $n = 80$  and multiplier 1.8647, or  $n = 1000$  and multiplier 1 and 1.6438. Methods include our estimator with different choices of  $h$ , GLASSO, SPACE, neighborhood selection (NS), and Space JAM (SJ).

The ROC (*receiver operating characteristic*) curves for different estimators are shown in

Figure 2.3 on Page 39. Each plotted curve corresponds to the average of 50 ROC curves, where the averaging is based on the vertical averaging from Algorithm 3 in Fawcett (2006), and is mean AUC-preserving. The  $x$  and  $y$  axes of each ROC curve represent the false positive and true positive rates at varying levels of penalty parameter  $\lambda$ , defined as

$$\text{FPR} \equiv \frac{|\hat{S}_{\text{off}} \setminus S_{0,\text{off}}|}{m(m-1) - |S_{0,\text{off}}|} \quad \text{and} \quad \text{TPR} \equiv \frac{|\hat{S}_{\text{off}} \cap S_{0,\text{off}}|}{|S_{0,\text{off}}|},$$

where  $S_{0,\text{off}} \equiv \{(i, j) : i \neq j \wedge \kappa_{0,ij} \neq 0\}$ , and  $\hat{S}_{\text{off}} \equiv \{(i, j) : i \neq j \wedge \hat{\kappa}_{ij} \neq 0\}$ .

To reduce clutter, we only report the results for the top performing competing methods. In particular, results for nonparanormal SKEPTIC are omitted, as the method always performs the worst in our experiments. The corresponding means and standard deviations of AUCs (*areas under the curves*) over 50 curves are given in Table 2.1.

Looking at the mean AUCs, with the standard deviations in mind, all choices of  $h$  considered here perform better than  $h(x) = x^2$  from Hyvärinen (2007) and Lin et al. (2016) and the competing methods. The results for  $n = 1000$  in Table 2.1 also show that the multiplier does help improve the AUCs, a matter to be discussed in Section 2.7.2.

*Truncated Non-Centered GGMS:* We generate data from a truncated non-centered Gaussian distribution with both parameters  $\boldsymbol{\mu}$  and  $\mathbf{K}$  unknown. In each trial, we form the true  $\mathbf{K}_0$  as in Section 2.7.1, and generate each component of  $\boldsymbol{\mu}_0$  independently from the normal distribution with mean 0 and standard deviation 0.5.

We compare the performance of our *profiled* estimator based on (2.17), with different  $h$  functions, but with no penalty on  $\boldsymbol{\eta} \equiv \mathbf{K}\boldsymbol{\mu}$ , to SPACE, SpaCE JAM (SJ), GLASSO, and neighborhood selection (NS). As before, we consider 50 trials. Representative ROC curves are plotted in Figure 2.4, and the corresponding AUCs are summarized in Table 2.2.

Even without tuning the extra penalty parameter on  $\boldsymbol{\eta} \equiv \mathbf{K}\boldsymbol{\mu}$ , our profiled estimator beats the competing methods by a large margin when  $n = 80$ . With multipliers 1 and  $n = 1000$ , our estimators still do better than Space JAM and GLASSO, and have performance comparable to other competing methods. It might appear that the performance of our

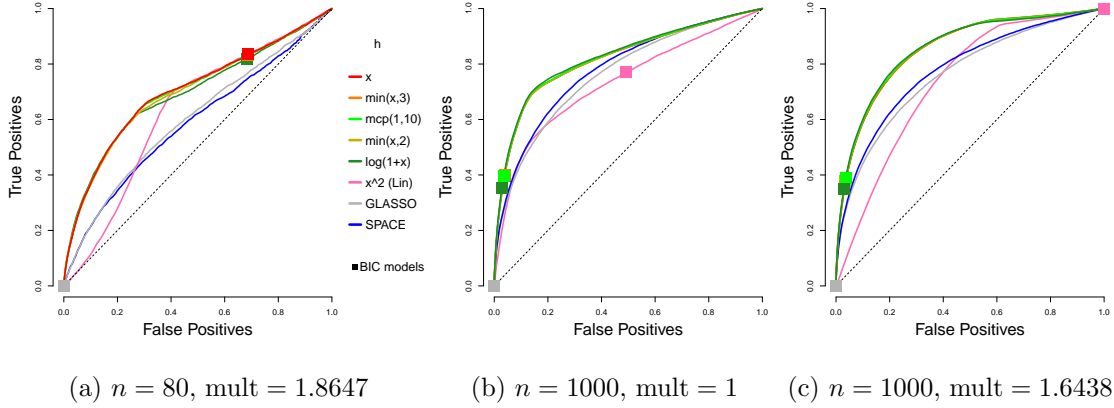


Figure 2.3: Average ROC curves of our centered estimator with various choices of  $h$ , compared to SPACE and GLASSO, for the truncated centered GGM case;  $m = 100$  variables and  $n = 80$  or 1000 samples are considered. Squares indicate average true positive rate (TPR) and false positive rate (FPR) of models picked by eBIC with refitting for the estimator in the same color.

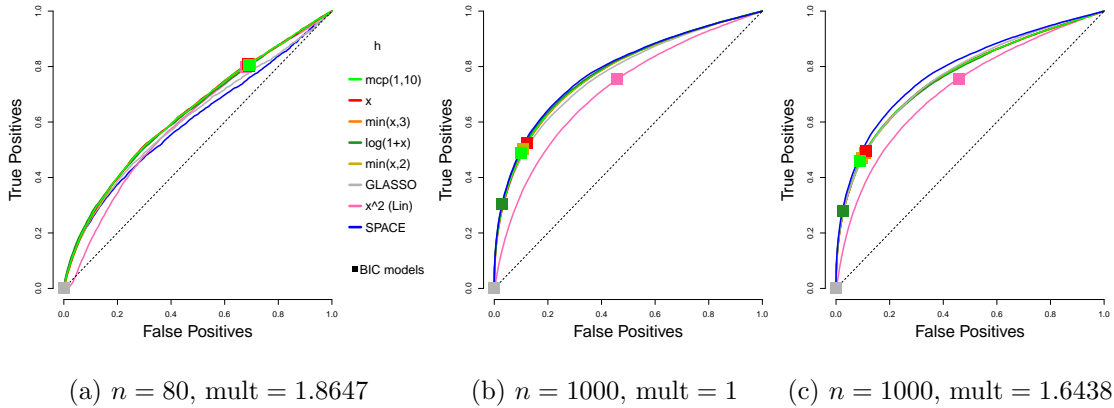


Figure 2.4: Average ROC curves of our non-centered profiled estimator with various choices of  $h$ , compared to SPACE and GLASSO, for the truncated non-centered GGM case;  $m = 100$  variables and  $n = 80$  or 1000 samples are considered.

Non-centered profiled, $n = 80$ , multiplier 1.8647					
$\min(\log(1+x), c)$			$\min(x, c)$		
$c$	Mean	sd	$c$	Mean	sd
$\infty$	0.632	0.032	$\infty$	0.634	0.032
2	0.632	0.032	3	0.634	0.032
1	0.631	0.032	2	0.632	0.032
0.5	0.619	0.033	1	0.628	0.032
MCP(1, $c$ )			SCAD(1, $c$ )		
$c$	Mean	sd	$c$	Mean	sd
10	0.634	0.032	5	0.634	0.032
5	0.634	0.032	10	0.634	0.032
1	0.622	0.032	2	0.634	0.032
$x^{1.5}$ : (0.623,0.031)			$x^2$ : (0.607,0.030)		
GLASSO: (0.614,0.029)			NS: (0.604,0.028)		
SPACE: (0.602,0.029)			SJ: (0.561,0.036)		

Non-centered profiled, $n = 1000$ , multiplier 1						Non-centered profiled, $n = 1000$ , multiplier 1.6438					
$\min(\log(1+x), c)$			$\min(x, c)$			$\min(\log(1+x), c)$			$\min(x, c)$		
$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd
$\infty$	0.783	0.020	2	0.779	0.020	$\infty$	0.764	0.018	$\infty$	0.766	0.019
2	0.783	0.020	$\infty$	0.779	0.020	2	0.764	0.018	3	0.765	0.019
1	0.782	0.020	3	0.779	0.020	1	0.762	0.018	2	0.764	0.018
0.5	0.767	0.021	0.5	0.758	0.020	0.5	0.738	0.018	1	0.753	0.018
MCP(1, $c$ )			SCAD(1, $c$ )			MCP(1, $c$ )			SCAD(1, $c$ )		
$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd
5	0.782	0.020	2	0.780	0.020	10	0.766	0.019	10	0.766	0.019
10	0.780	0.020	5	0.780	0.020	5	0.766	0.019	5	0.766	0.019
1	0.771	0.021	10	0.779	0.020	1	0.745	0.018	2	0.763	0.018
$x^{1.5}$ : (0.751,0.019)			$x^2$ : (0.713,0.018)			$x^{1.5}$ : (0.748,0.018)			$x^2$ : (0.718,0.017)		
SPACE: (0.786,0.020)			NS: (0.785,0.02)			SPACE: (0.786,0.020)			NS: (0.785,0.020)		
GLASSO (0.770,0.019)			SJ: (0.720,0.019)			GLASSO (0.770,0.019)			SJ: (0.720,0.019)		

Table 2.2: Mean and standard deviation of AUC using different profiled estimators in the non-centered setting, with  $n = 80$  and multiplier 1.8647, or  $n = 1000$  and multipliers 1 and 1.6438. Methods include our estimator with different choices of  $h$ , GLASSO, SPACE, neighborhood selection (NS), and Space JAM (SJ).

estimators deteriorate with a multiplier larger than 1; however, as we will see, there can be significant improvement in AUCs if we tune an additional parameter for the multiplier. As in the centered case, the leading  $h$  functions in each category perform similarly, and the exact choice is not crucial. Subsequently, we will simply use  $h(x) = \min(x, 3)$ .

### Choice of multiplier

*Truncated Centered GGMs:* In Figure 2.5, the ROC curves for GLASSO, SPACE, and our estimator with  $h(x) = \min(x, 3)$ , but with different levels of amplification, via different choices of multipliers  $\delta$ , are compared for the centered case of Section 2.7.2.

While Theorem 10 guarantees consistency only for  $\delta < C(n, m)$ , we observe that there can be a gain from going beyond the *upper-bound multiplier*  $C(n, m)$ , which is 1.8647 for  $n = 80$  and 1.6438 for  $n = 1000$  (when  $n = 1000$ ,  $C(n, m)$  turns out to be the best-performing multiplier). However, the effect deteriorates fast as the multiplier grows larger. The figure suggests that while some additional gains are possible by tuning over the choice of multiplier, the *upper-bound multiplier* is a good default.

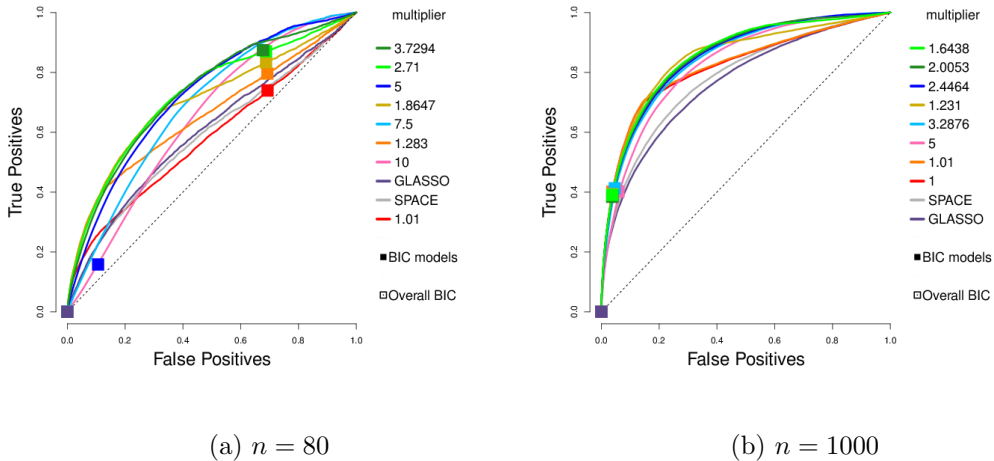


Figure 2.5: Performance of  $\min(x, 3)$  for truncated centered GGMs using different multipliers, compared to GLASSO and SPACE, in the centered setting,  $n = 80$  or 1000.

*Truncated Non-Centered GGMs* : In Figure 2.6, we consider the non-centered case of Section 2.7.2, and use the non-profiled estimator; that is, the non-centered estimator with  $\ell_1$  penalty on both  $\mathbf{K}$  and  $\boldsymbol{\eta} \equiv \mathbf{K}\boldsymbol{\mu}$ . The ROC curves are compared to competing methods GLASSO and SPACE. For the choice of amplification in our estimator, we consider the upper-bound multiplier  $C(n, m)$  from Theorem 11 as the default. We refer to this as *high* amplification. We also consider lower amplification, with  $\delta = 2 - (1 + 24e \log m/n)^{-1}$ , referred to as *medium*. For  $n = 1000$ , we also consider a *low* multiplier 1, which corresponds to no amplification. We compare these possible defaults to a finer grid of multipliers of which we show some representatives in the plots.

We see that among our defaults, the upper-bound choice  $C(n, m)$  performs best. Some additional gains are possible by tuning the multiplier over a grid of values containing this choice. Moreover, we see that it can be beneficial to tune over both  $\lambda_{\mathbf{K}}$  and  $\lambda_{\mathbf{K}}/\lambda_{\boldsymbol{\eta}}$ .

We remark that while for each run, the best model picked by BIC falls on the ROC curve, a few squares are off the curve in Figure 2.6 (c). This is because these squares correspond to the average of the true and false positive rates of the chosen BIC models over 50 runs, potentially due to multimodality of the distribution of the models. Nonetheless, in all cases, the average of the models picked by BIC tuned over both  $\lambda_{\mathbf{K}}$  and  $\lambda_{\mathbf{K}}/\lambda_{\boldsymbol{\eta}}$  looks reasonable.

### 2.7.3 Other a/b Models

We now turn to the non-Gaussian ( $a \neq 1$  or  $b \neq 1$ ) setting. Based on the observations in Section 2.5.4, we focus on functions of type  $\min(x^p, c)$  for some power  $p > 0$  and truncation point  $c > 0$ . For simplicity, for the non-centered models we use the profiled estimator (2.17) (i.e.,  $\lambda_{\boldsymbol{\eta}} = 0$ ) and use the multiplier  $C(n, m)$  in Theorem 10 for truncated GGMs as a guidance. We note that tuning over the  $\lambda_{\boldsymbol{\eta}}$  parameter and the multiplier can potentially give a significant improvement as seen in Section 2.7.2.

These simulations suggest that among the class of functions of the form  $\min(x^p, c)$ ,  $x^{2-a}$  or  $\min(x^{2-a}, c)$  with a moderately large  $c$  can be used as the default choice of  $h(x)$ . This agrees with our findings in Section 2.7.2. We note that bounded  $h$  functions were only

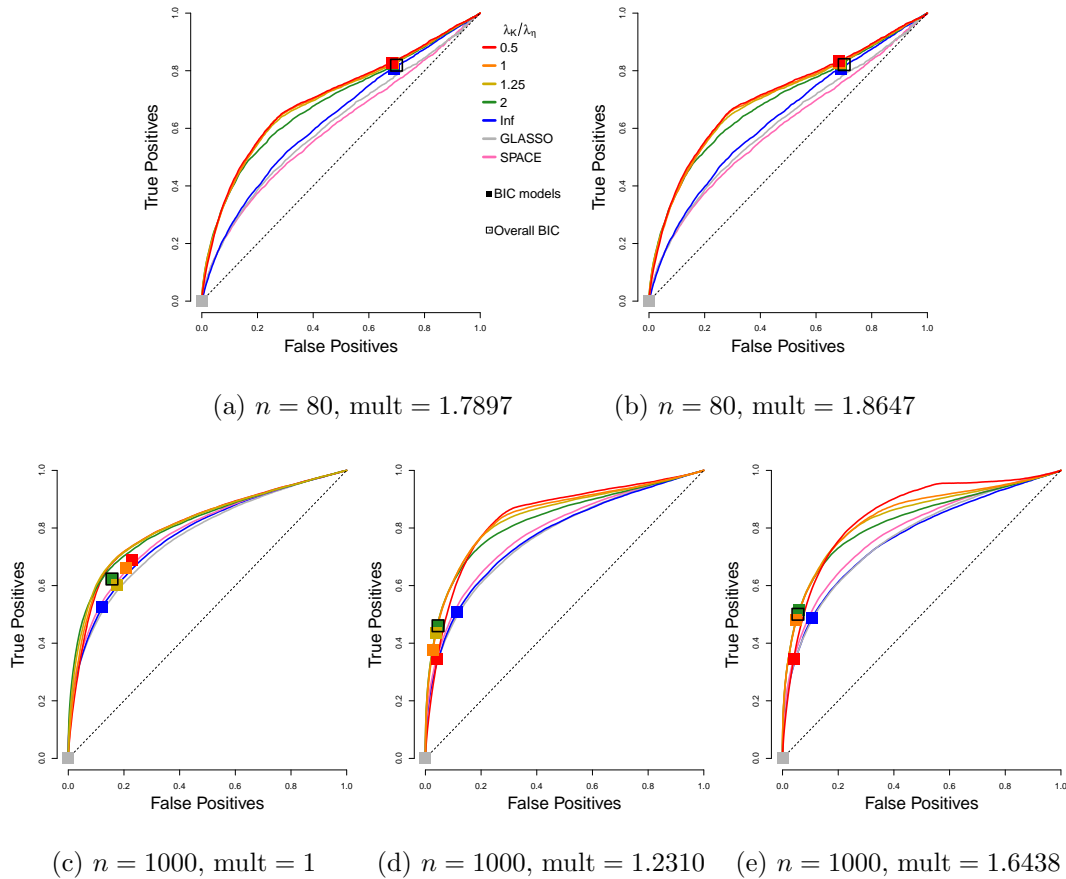


Figure 2.6: Performance of the non-centered estimator with  $h(x) = \min(x, 3)$ . Each curve corresponds to a different choice of  $\lambda_{\mathbf{K}}/\lambda_{\eta}$ . Squares indicate models picked by eBIC with refit. The square with black outline has the highest eBIC among all models (combinations of  $\lambda_{\mathbf{K}}$ ,  $\lambda_{\eta}$ ). Multipliers correspond to medium or high for  $n = 80$ , and low, medium or high for  $n = 1000$ , respectively.

used in the proof for truncated GGMs, and picking a moderately large truncation point can correspond to having an untruncated power.

### *Exponential Setting*

For the exponential models,  $a = b = 1/2$ . Since  $a = b$ , for both centered and non-centered settings, based on the principle in Section 2.5.4, choosing  $h(x) = \min(x^{3/2}, c)$  satisfies (A1) and (A2) and also ensures that entries in  $\mathbf{\Gamma}$  and  $\mathbf{g}$  are bounded (for small  $x$ ), while choosing  $h(x) = \min(\sqrt{x}, c)$  only guarantees (A1) and (A2).

In Figure 2.7, we present the AUCs for the ROC curves of edge recovery with different choices of  $h(x) = \min(x^{\text{pow}}, c)$ . As before, we set  $n = 80$  or  $1000$  and  $m = 100$ , but we use an  $\boldsymbol{\eta}_0$  with each component uniformly equal to  $-0.5$ ,  $0$  or  $0.5$ ; for  $\boldsymbol{\eta}_0 \equiv \mathbf{0}$ , we assume this information is known and use the centered estimator. The results suggest that  $\text{pow} = 3/2 = 2 - a$  is the best choice of power. For this optimal choice, the performance improves with larger  $c$ , so  $x^{2-a}$  gives the best results. For sub-optimal powers, including truncation gives better results.

### *Gamma Setting*

The centered gamma models reduce to the centered exponential models. Thus, in this section, we only consider the non-centered settings, with  $a = 1/2$ ,  $b = 0$ . From Section 2.5.4, we have the following choices:

- $\min(x^2, c)$  both satisfies (A1)–(A2) and ensures  $\mathbf{\Gamma}$  and  $\mathbf{g}$  are bounded;
- $\min(x^{\max\{3/2, 1 - \min_j \eta_{0,j}\}}, c)$  ensures (A1)–(A2) and bounds  $\mathbf{\Gamma}_{11}$  and  $\mathbf{g}_1$ ; by default without prior information on  $\boldsymbol{\eta}_0$  this is  $\min(x^2, c)$ ;
- $\min(x^{3/2}, c)$  satisfies both conditions on the interaction part only ( $\mathbf{x}^a$ ), but does not guarantee (A1)–(A2);

- $\min(x^{1/2}, c)$  satisfies the sufficient conditions for (A1)–(A2) on the interaction only.

The results are shown in Figure 2.8, where we consider  $n = 80, 1000$ , and  $\boldsymbol{\eta} = \pm 0.5\mathbf{1}_{100}$ . They suggest that  $\text{pow} = 2 - a = 1.5$  works consistently well, although slightly outperformed by 1 and 1.25 in one case. As in the exponential case, with the optimal power it is beneficial to choose a large truncation point, or work with an untruncated power  $x^{1.5}$ . We conclude that the performance is likely only dependent on the  $(2 - a)$  power requirement for the  $\mathbf{x}^{a\top} \mathbf{K} \mathbf{x}^a$  part or  $2 - \min_j \eta_{0,j}$ ; simulations in the next section rule out the possibility of the latter.

#### *Other Choices of $a$ and $b$*

In this section, we consider other choices of  $a$  and  $b$ . Specifically,  $a = 3/2$  and  $b = 1/2$  or  $0$ . These combinations are chosen to confirm, in a more extreme setting, that the performance is mainly determined by requirements on the power based on  $a$ , which correspond to choosing a power of  $1 - a$  or  $2 - a$ , but not those on  $b$  (or on  $\boldsymbol{\eta}$  when  $b = 0$ ) that correspond to  $1 - b$  and  $2 - b$ . The relationship between these two settings is analogous to that between the exponential and gamma models (same  $a, b$  nonzero/zero).

The results are shown in Figures 2.9 and 2.10, and indeed confirm that  $x^{2-a} = x^{0.5}$  consistently gives the optimal results, even though  $\boldsymbol{\eta}^\top \mathbf{x}^b$  is in favor of  $x^{2-b} = x^{1.5}$  for  $b = 0.5$ , and  $\boldsymbol{\eta}^\top \log(\mathbf{x})$  is in favor of  $x^2$  or at least  $x^{1-\min_j \eta_{0,j}}$  when  $b = 0$ . There are two possible explanations for the optimality of  $2 - a$  over  $\max\{2 - a, 2 - b\}$  or  $\max\{2 - a, 1 - \min_j \eta_{0,j}\}$ : (1) The AUC metric is measured only on our interest, edge recovery for the interaction matrix, which only depends on  $\mathbf{x}^a$ ; (2) using the profiled estimator weakens the effect of  $b$ .

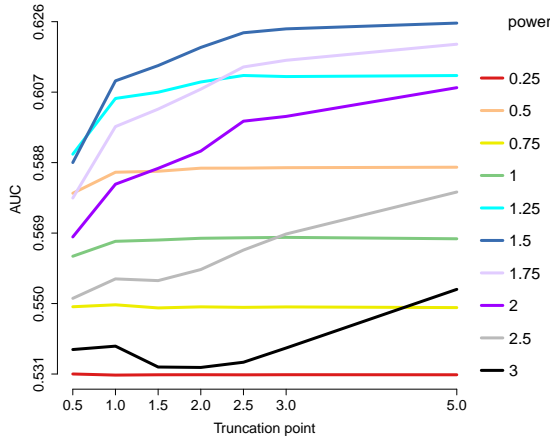
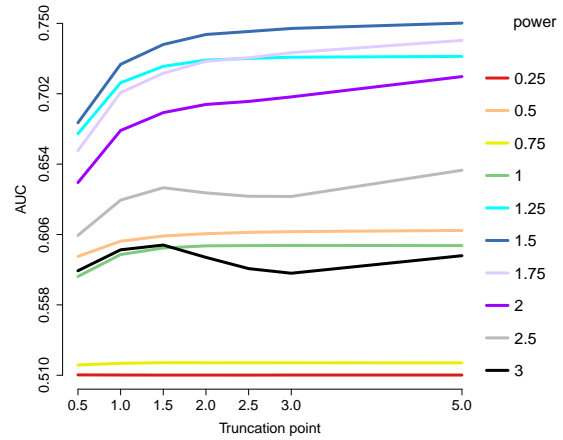
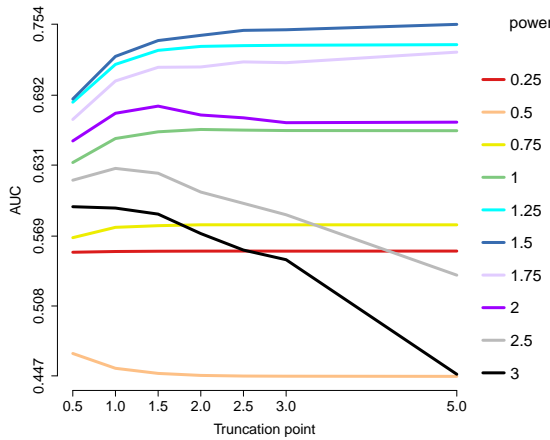
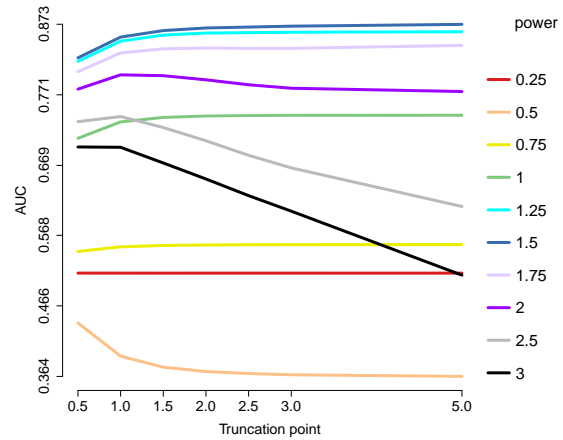
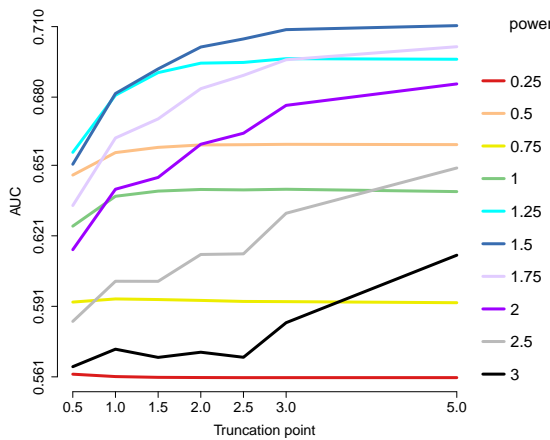
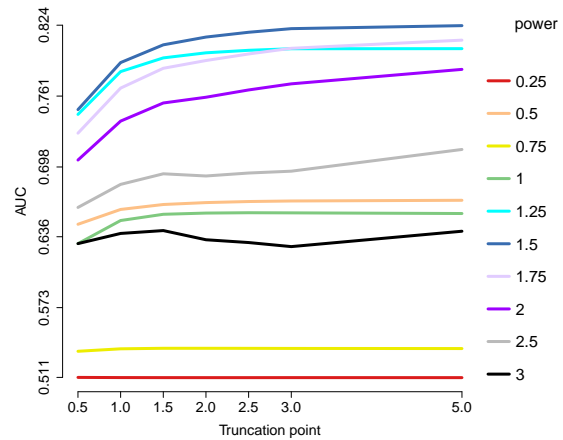
(a)  $n = 80$ ,  $\eta = -0.5\mathbf{1}_{100}$ , profiled estimator(b)  $n = 1000$ ,  $\eta = -0.5\mathbf{1}_{100}$ , profiled estimator(c)  $n = 80$ ,  $\eta = \mathbf{0}$ , centered estimator(d)  $n = 1000$ ,  $\eta = \mathbf{0}$ , centered estimator(e)  $n = 80$ ,  $\eta = 0.5\mathbf{1}_{100}$ , profiled estimator(f)  $n = 1000$ ,  $\eta = 0.5\mathbf{1}_{100}$ , profiled estimator

Figure 2.7: AUCs for edge recovery using generalized score matching for the exponential models. Each curve represents a different choice of power  $p$  in  $h(x) = \min(x^p, c)$ , and the  $x$  axis marks the truncation point  $c$ . Colors are sorted by  $p$ .

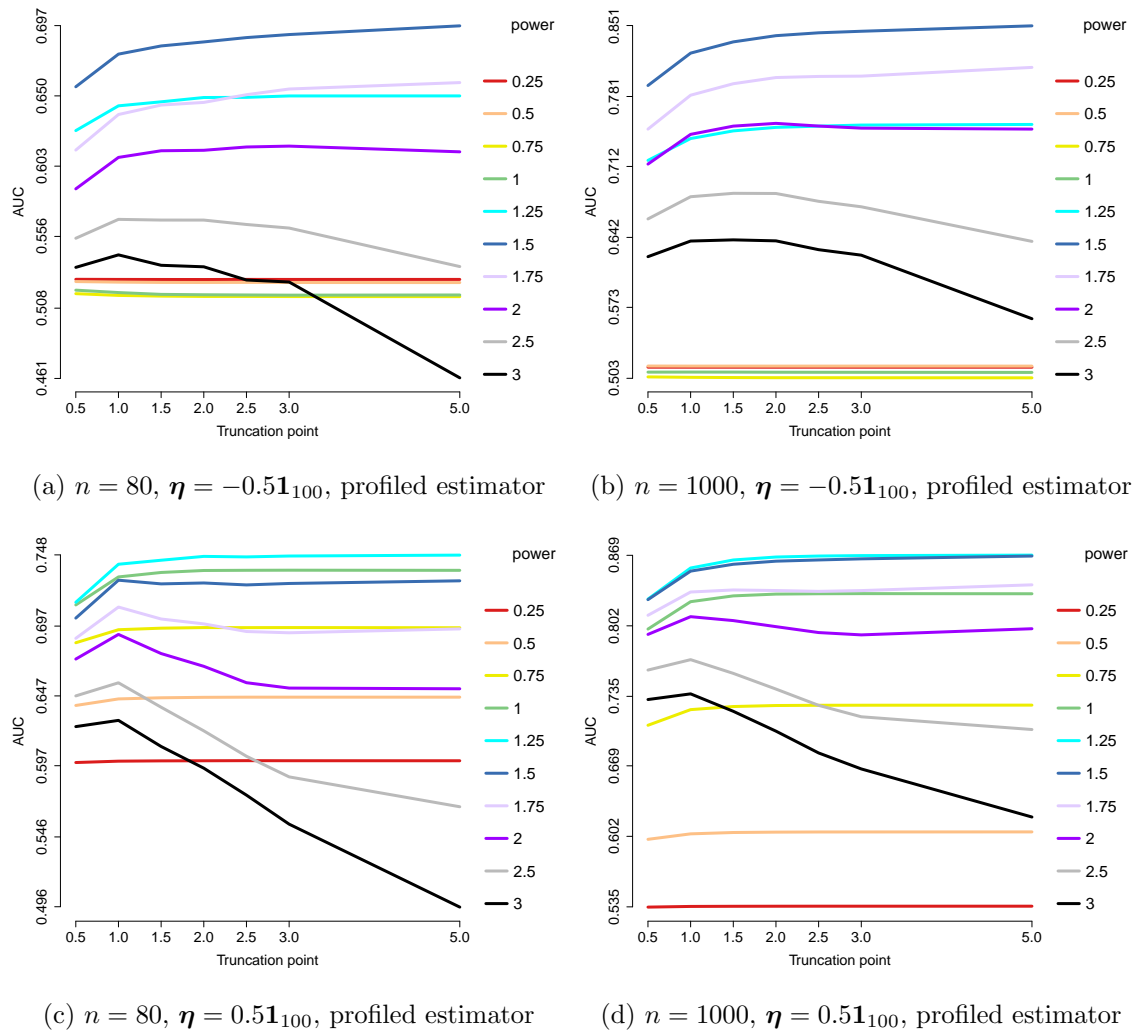
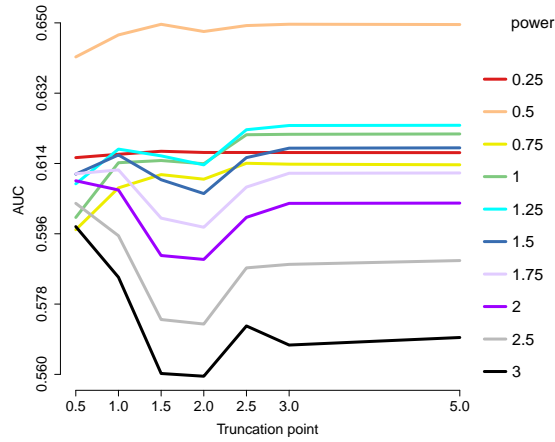
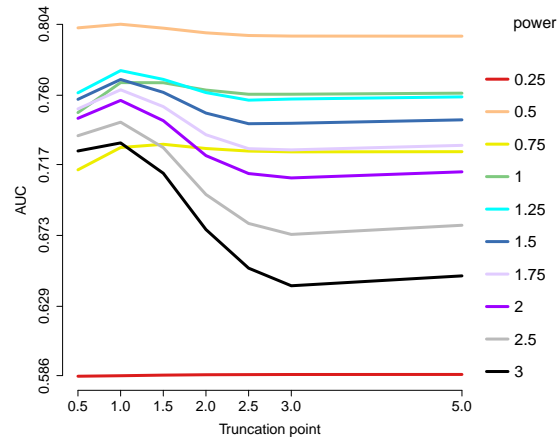
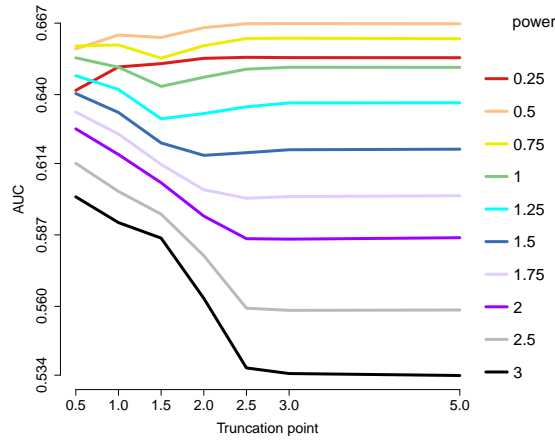
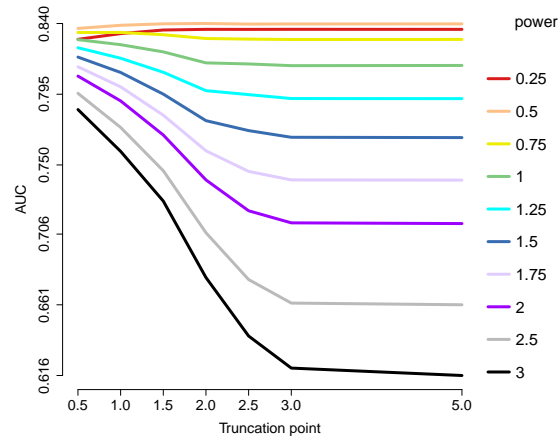
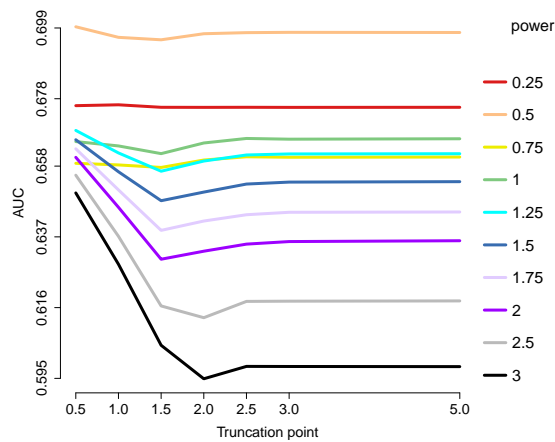
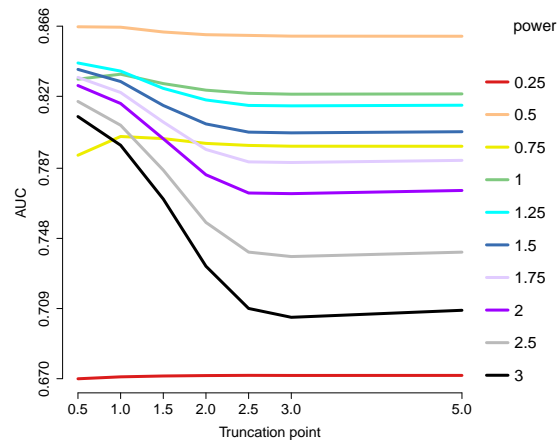
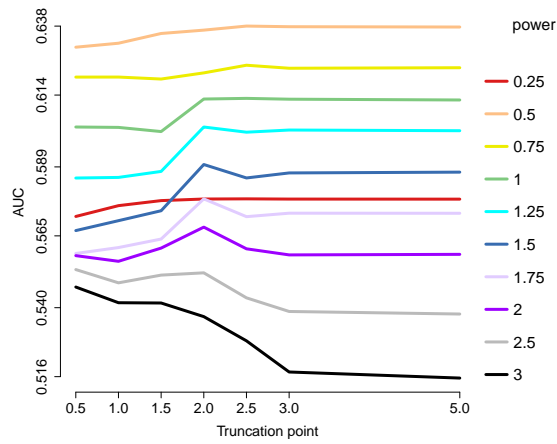
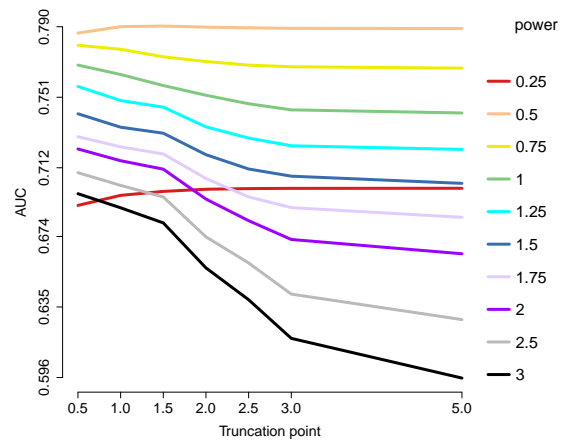
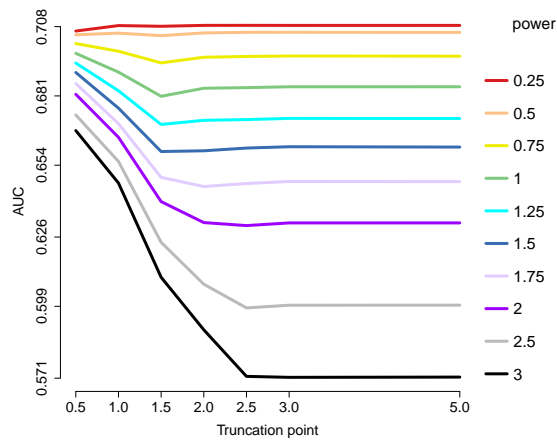
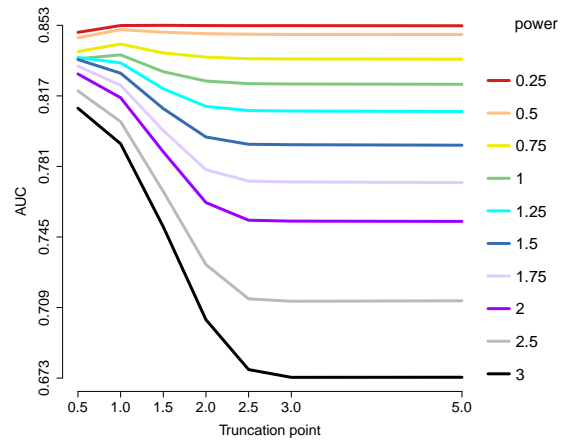


Figure 2.8: AUCs for edge recovery using generalized score matching for the gamma models.

(a)  $n = 80$ ,  $\eta = -0.51_{100}$ , profiled estimator(b)  $n = 1000$ ,  $\eta = -0.51_{100}$ , profiled estimator(c)  $n = 80$ ,  $\eta = \mathbf{0}$ , centered estimator(d)  $n = 1000$ ,  $\eta = \mathbf{0}$ , centered estimator(e)  $n = 80$ ,  $\eta = 0.51_{100}$ , profiled estimator(f)  $n = 1000$ ,  $\eta = 0.51_{100}$ , profiled estimatorFigure 2.9: AUCs for edge recovery using generalized score matching for  $a = 3/2$ ,  $b = 1/2$ .

(a)  $n = 80$ ,  $\boldsymbol{\eta} = -0.5\mathbf{1}_{100}$ , profiled estimator(b)  $n = 1000$ ,  $\boldsymbol{\eta} = -0.5\mathbf{1}_{100}$ , profiled estimator(c)  $n = 80$ ,  $\boldsymbol{\eta} = 0.5\mathbf{1}_{100}$ , profiled estimator(d)  $n = 1000$ ,  $\boldsymbol{\eta} = 0.5\mathbf{1}_{100}$ , profiled estimatorFigure 2.10: AUCs for edge recovery using generalized score matching for  $a = 3/2$ ,  $b = 0$ .

### 2.7.4 RNAseq Data

In this section we apply our regularized generalized  $h$ -score matching estimator for truncated non-centered GGMs to RNAseq data also studied in Lin et al. (2016), since the same model is considered therein. The data consists of  $n = 487$  prostate adenocarcinoma samples from The Cancer Genome Atlas (TCGA) data set. Following Lin et al. (2016), we focus on  $m = 333$  genes that belong to the known cancer pathways in the Kyoto Encyclopedia of Genes and Genomes (KEGG) and that have no more than 10% missing values. Missing values are set to 0. We choose  $h(x) = \min(x, 3)$  and use the upper-bound multiplier (*high*), as discussed in Section 2.7.2. For simplicity, we use the profiled estimator, and choose the regularization parameter  $\lambda$  so that the estimated graph has exactly  $m = 333$  edges, all these choices being as in Lin et al. (2016).

We compare our graph to the one in Lin et al. (2016), which corresponds to  $h(x) = x^2$  with no multiplier. Shown in Figure 2.11 are the estimated graphs, with their intersection in the middle. To improve visualization, isolated nodes are removed and the layouts are optimized for each plot. Red-colored points are the “hub nodes”, namely nodes with degree at least 10. In Figure 2.12, we plot the same graphs in a fixed layout optimized for the graph corresponding to  $h = \min(x, 3)$ , and include the isolated nodes.

Out of 333 edges, the two estimated graphs share 117 edges in common. Assuming that edges are placed at random between nodes and the two graphs are independent, the distribution of the number  $R$  of common edges follow a hypergeometric distribution, so  $P(R = r) = \frac{\binom{m}{r} \binom{m(m-1)/2 - m}{m-r}}{\binom{m(m-1)/2}{m}}$ . For  $m = 333$  the probability of at least 117 common edges is essentially zero. The large number of shared edges between the two methods can be explained by the fact that they both minimize the same underlying score-matching loss.

The graph using  $h(x) = \min(x, 3)$  has much more isolated nodes (204) than the other (108), and has a slightly smaller max degree (16 versus 19). Table 2.3 provides another way of comparing between the two graphs by listing the genes with the highest node degrees.

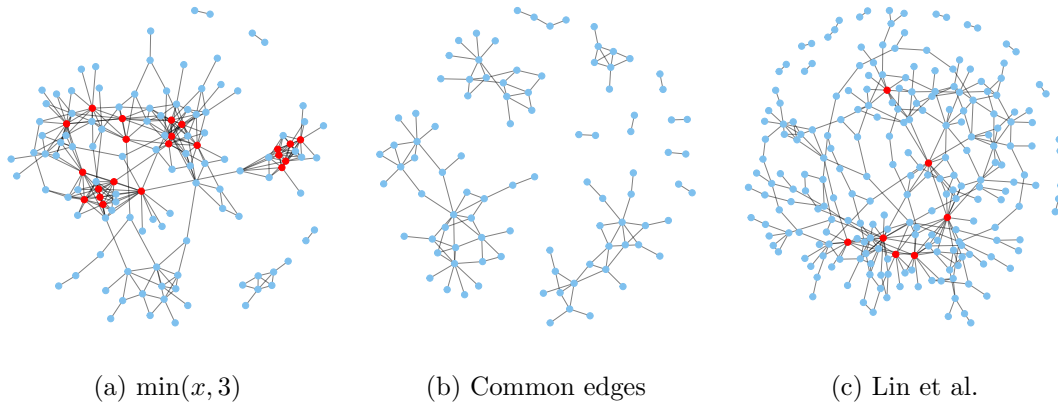


Figure 2.11: Graphs estimated by regularized generalized score matching estimator with  $h(x) = \min(x, 3)$  with upper-bound multiplier (left) and  $h(x) = x^2$  with no multiplier (Lin et al., 2016, right), and their intersection graph (middle). Isolated nodes with no edges are removed, and the layout is optimized for each plot. In (a) and (c), red points indicate nodes with degree at least 10 (“hub nodes”).

$\min(x, 3)$ with multiplier 1.63	Lin et al.
<b>LAMB3 (16)</b>	CCNE2 (19)
<b>PIK3CG (16)</b>	<b>PIK3CG (16)</b>
MMP2 (15)	BRCA2 (13)
GLI2 (13)	<b>BIRC5 (12)</b>
LAMA4 (13)	<b>LAMB3 (10)</b>
<b>PDGFRB (13)</b>	<b>PIK3CD (10)</b>
<b>PIK3CD (13)</b>	SKP2 (10)
RASSF5 (13)	HRAS (9)
<b>BIRC5 (12)</b>	STAT5B (9)
FLT3 (12)	<b>GSTP1(8)</b>
<b>GSTP1 (12)</b>	<b>PDGFRB (8)</b>
LAMA2 (12)	
RAC2 (12)	

Table 2.3: List of genes with the highest node degrees in each estimated graph.

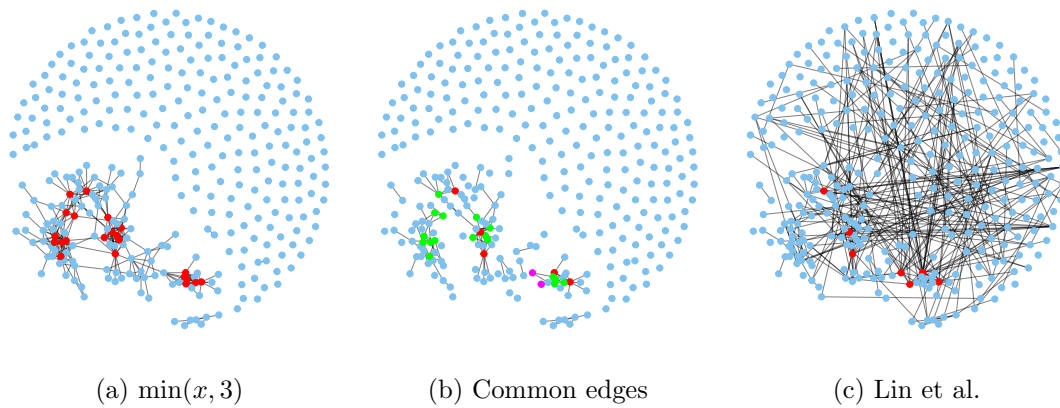


Figure 2.12: Graphs estimated by regularized generalized score matching estimator with  $h(x) = \min(x, 3)$  with upper-bound multiplier (left) and  $h(x) = x^2$  with no multiplier (Lin et al., 2016, right), and their intersection graph (middle). Isolated nodes are included and the layout is fixed across plots and optimized for graph (a). In (b) the red nodes are hub nodes shared by both graphs, the green ones are hub nodes in graph (a) only, and the magenta ones are hub nodes in graph (c) only.

In Table 2.3 we list the top ten genes in terms of node degree for both estimated graphs. Due to ties, 13 genes are listed for  $h(x) = \min(x, 3)$  and 11 for Lin et al. (2016). As noted in Lin et al. (2016), genes with high node-degrees are known to be important in biological networks (Carter et al., 2004; Jeong et al., 2001; Han et al., 2004). Among these top genes, six are common in both graphs, and are discussed in Lin et al. (2016). We next elaborate on the evidence supporting the first four of the newly discovered genes.

- MMP2 (Matrix metalloproteinase 2): According to Trudel et al. (2003), increased MMP-2 expression is an independent predictor of decreased prostate cancer disease-free survival. Morgia et al. (2005) state that activity of MMP-2 can be useful in diagnosis, therapy, and assessment of malignant progression in prostate cancer.
- GLI2 (GLI family zinc finger 2): GLI2 is a primary mediator of the hedgehog signaling pathway, which has been reported in prostate cancer, and plays a critical role in the malignant phenotype of prostate cancer cells (Thiyagarajan et al., 2007). Its increased level of expression is also related to AI prostate cancer, and may be a therapeutic target in castrate-resistant prostate cancer (Narita et al., 2008).
- LAMA4 (Laminin subunit alpha 4): LAMA4 is consistently upregulated in benign prostatic hyperplasia when compared to normal prostate tissues (Luo et al., 2002).
- RASSF5 (RAS association domain family member 5): The combination of RASSF5 along with four other DNA methylation markers can effectively differentiate between benign prostate biopsy cores from non-cancer patients and cancer cores, and can be used to identify patients at risk without repeat biopsies (Brikun et al., 2014).

We note that the two methods indeed use different estimators (different  $h$  functions and multipliers), and it is thus not surprising to see that some of the top genes by one method are not among those for the other. In particular, CCNE2, BRCA2, SKP2 and STAT5B, while previously reported as newly discovered in Lin et al. (2016), are dropped by our new

analysis. Testing and inference (potentially using bootstrapping) is an important problem but is beyond the scope of this chapter.

## 2.8 Discussion

In this chapter, we proposed a generalization of the score matching estimator of Hyvärinen (2007), based on scaling the log-gradients to be matched with a suitably chosen function  $\mathbf{h}$ . The generalization retains the advantages of Hyvärinen’s method: Estimates can be computed without knowledge of normalizing constants, and for canonical parameters of exponential families, the estimation loss is a quadratic function.

For high-dimensional exponential family graphical models, following Lin et al. (2016), we add an  $\ell_1$  penalty to regularize the generalized score matching loss. One practical issue that is overlooked in Lin et al. (2016) is the fact that the score matching loss can be unbounded from below for a small tuning parameter, when the dimension  $m$  exceeds the sample size  $n$ . We fix this issue by amplifying the diagonal entries in the quadratic matrix in the definition of the generalized score matching loss by a factor/multiplier, and we give an upper bound on that multiplier that guarantees consistency.

As examples we consider *pairwise interaction power models* on the non-negative orthant  $\mathbb{R}_+^m$ . Specifically, the considered models are exponential families in which the log density is the sum of pairwise interactions between entries of powers  $\mathbf{x}^a$  plus linearly weighted effects  $\mathbf{x}^b$ , or  $\log(\mathbf{x})$  when  $b = 0$ . Our main interest is in the matrix of interaction parameters whose support determines the distributions’ conditional independence graph. The considered framework covers truncated normal distributions ( $a = b = 1$ ), exponential square root graphical models ( $a = b = 1/2$ ) from Inouye et al. (2016), as well as a class of multivariate gamma distributions ( $a = 1/2, b = 0$ ).

In the case of multivariate truncated normal distributions, where the conditional independence graph is given by the underlying Gaussian inverse covariance matrix, the sample size required for the consistency of our method using bounded  $\mathbf{h}$  is  $\Omega(d^2 \log m)$ , where  $d$  is the degree of the graph. This matches the rates for Gaussian graphical models in Ravikumar

et al. (2011) and Lin et al. (2016). In contrast, the sample complexity for truncated Gaussian models given in Lin et al. (2016) is  $\Omega(d^2 \log^8 m)$ .

For the considered class of pairwise interaction models, we recommend using the function  $\mathbf{h}$  with coordinates  $h_j(x) = \min(x^{2-a}, c)$  for some moderately large  $c$ , or simply  $h_j(x) = x^{2-a}$ . While this choice is effective, it would be an interesting problem for future work to develop a method that adaptively chooses an optimized function  $\mathbf{h}$  from data.

## Chapter 3

**GENERALIZED SCORE MATCHING FOR GENERAL  
DOMAIN**

### 3.1 Introduction

In this chapter, we aim to further extend the generalized score matching loss to complicated and possibly unbounded domains  $\mathcal{D}$  not limited to  $\mathbb{R}^m$  and  $\mathbb{R}_+^m$ , and apply the method to estimation of graphical models supported on  $\mathcal{D}$ . In particular, we allow  $\mathcal{D}$  to be any *componentwise countable* domain defined in Definition 5, which requires that the induced domain of one component  $x_j$  fixing all other  $\mathbf{x}_{-j}$  is a union of at most countably many intervals in  $\mathbb{R}$ ; this should cover most but pathological cases in practice. We propose composing the  $\mathbf{h}$  function in the loss defined in Chapter 2 (Yu et al., 2019b) with a distance function  $\boldsymbol{\psi} = (\psi_1, \dots, \psi_m) : \mathcal{D} \rightarrow \mathbb{R}_+^m$ , with  $\psi_j(\mathbf{x})$  the distance of  $x_j$  to its induced boundary while fixing the other components  $\mathbf{x}_{-j}$ , truncated from above by some user-selected constant  $C_j$ . We then use  $(\mathbf{h} \circ \boldsymbol{\psi})$  in place of  $\mathbf{h}$  and show that the generalized score matching loss can again be approximated by an empirical loss, which is quadratic in the canonical parameters for exponential family distributions.

Since the density corresponding to the minimizer of the expected (true) loss is equal to the true density a.e.  $\mathbf{x} \in \mathcal{D}$ , the method above only makes sense for  $\mathcal{D}$  with positive Lebesgue measure in  $\mathbb{R}^m$ . As an important example of domains with zero Lebesgue measure in  $\mathbb{R}^m$ , we also present results for the  $m$ -simplex domain  $\{\mathbf{x} \in \mathbb{R}_+^m \mid \mathbf{x} \succ \mathbf{0}, \mathbf{1}_m^\top \mathbf{x} = 1\}$  by profiling out  $x_m \equiv 1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m}$ . Notably, we revisit the  $A^{m-1}$  models for simplex domains introduced in Aitchison (1985), which corresponds to  $a = b = 0$  with  $\mathbf{K}\mathbf{1}_m = \mathbf{K}^\top \mathbf{1}_m = \mathbf{0}_m$  in the pairwise interaction power models introduced in Section 1.4. We show that estimation for models on the simplex can be done with little additional effort compared to general  $a$ - $b$  models on  $\mathcal{D}$  with positive Lebesgue measure.

In Chapter 2 we proved theoretical consistency results for high-dimensional settings requiring sample size  $n = \Omega(\log m)$  for Gaussian graphical models (GGMs) restricted to  $\mathbb{R}_+^m$ . In this chapter we show that this sample complexity also holds for GGMs restricted to any domain  $\mathcal{D}$  that is a finite disjoint union of convex sets, as well as general  $a$ - $b$  models on bounded subsets of  $\mathbb{R}_+^m$  with positive measure for  $a \geq 0$  and on simplex domains for  $a > 0$ .

On unbounded domains for  $a > 0$  and on simplex domains for  $a = 0$  (notably, the  $A^{m-1}$  models) we require the sample size  $n$  to be  $\Omega(\log m)$  times an unknown constant that may weakly depend on  $m$ .

The structure of the chapter is as follows. In Section 3.2 we provide an introduction to score matching and in Section 3.3 we detail our new methodology and theory necessary for the method to be usable in practice. In Section 3.4 we revisit the pairwise interaction power  $a$ - $b$  models and the regularized generalized score matching estimator for exponential families, while in Sections 3.5 and 3.6 we focus on application of our method to graphical models on domains with positive Lebesgue measure and simplex domains, respectively. Theoretical results and numerical experiments are demonstrated in Sections 3.7 and 3.8, respectively. We apply our method to a DNA methylation dataset in Section 3.9. Longer proofs are included in Appendix B. An implementation that incorporates various types of domain  $\mathcal{D}$  is available in the `genscore` R package.

## 3.2 Preliminaries

Suppose  $\mathbf{X} \in \mathbb{R}^m$  is a random vector with distribution function  $P_0$  supported on domain  $\mathcal{D} \subseteq \mathbb{R}^m$  and a twice continuously differentiable probability density function  $p_0$  with respect to the Lebesgue measure restricted to  $\mathcal{D}$ . Let  $\mathcal{P}(\mathcal{D})$  be a family of distributions of interest with twice continuously differentiable densities on  $\mathcal{D}$ . The goal is to estimate  $p_0$  by picking the distribution  $P$  from  $\mathcal{P}(\mathcal{D})$  with density  $p$  that minimizes some empirical loss that measures the distance between  $p$  and  $p_0$ .

### 3.2.1 Original Score Matching on $\mathbb{R}^m$ (Hyvärinen, 2005)

The original *score matching* loss proposed by Hyvärinen (2005) for  $\mathcal{D} \equiv \mathbb{R}^m$  is given by

$$J_{\mathbb{R}^m}(P) \equiv \frac{1}{2} \int_{\mathbb{R}^m} p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) - \nabla \log p_0(\mathbf{x})\|_2^2 d\mathbf{x}, \quad (3.1)$$

in which the gradients can be thought of as gradients with respect to a hypothetical location parameter and evaluated at the origin (Hyvärinen, 2005). The log densities enable estimation

without calculating the normalizing constants of  $p$  and  $p_0$ . Under mild conditions, using integration by parts the loss can be rewritten as

$$J_{\mathbb{R}^m}(P) \equiv \int_{\mathbb{R}^m} p_0(\mathbf{x}) \sum_{j=1}^m \left[ \partial_{jj} \log p(\mathbf{x}) + \frac{1}{2} (\partial_j \log p(\mathbf{x}))^2 \right] d\mathbf{x}$$

plus a constant independent of  $p$ . One can thus use a sample average to approximate the loss without knowing the true density  $p_0$ .

### 3.2.2 Score Matching on $\mathbb{R}_+^m$ (Hyvärinen, 2007; Yu et al., 2018, 2019b)

Consider  $\mathcal{D} \equiv \mathbb{R}_+^m$ . Let  $\mathbf{h} : \mathbb{R}_+^m \rightarrow \mathbb{R}_+^m$ ,  $\mathbf{x} \mapsto (h_1(x_1), \dots, h_m(x_m))^\top$ , where  $h_1, \dots, h_m : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  are almost surely positive functions that are absolutely continuous in every bounded sub-interval of  $\mathbb{R}_+$ . The *generalized  $\mathbf{h}$ -score matching loss* proposed in Chapter 2 (Yu et al., 2018, 2019b) is

$$J_{\mathbf{h}, \mathbb{R}_+^m}(P) \equiv \frac{1}{2} \int_{\mathbb{R}_+^m} p_0(\mathbf{x}) \left\| \nabla \log p(\mathbf{x}) \odot \mathbf{h}^{1/2}(\mathbf{x}) - \nabla \log p_0(\mathbf{x}) \odot \mathbf{h}^{1/2}(\mathbf{x}) \right\|_2^2 d\mathbf{x}. \quad (3.2)$$

The score matching loss for  $\mathbb{R}_+^m$  originally proposed by Hyvärinen (2007) is a special case of (3.2) with  $\mathbf{h}(\mathbf{x}) = \mathbf{x}^2$ . In Chapter 2 (Yu et al., 2018, 2019b) we proved that by choosing  $h_1, \dots, h_m$  that grow slowly and preferably bounded the estimation efficiency can be significantly improved. Under assumptions that for all  $P \in \mathcal{P}(\mathbb{R}_+^m)$  with density  $p$ ,

$$(A0.1) \quad p_0(x_j; \mathbf{x}_{-j}) h_j(x_j) \partial_j \log p(x_j; \mathbf{x}_{-j}) \Big|_{x_j \searrow 0^+}^{x_j \nearrow +\infty} = 0, \quad \forall \mathbf{x}_{-j} \in \mathbb{R}_+^{m-1} \quad \forall j;$$

$$(A0.2) \quad \mathbb{E}_{p_0} \left\| \nabla \log p(\mathbf{X}) \odot \mathbf{h}^{1/2}(\mathbf{X}) \right\|_2^2 < +\infty, \quad \mathbb{E}_{p_0} \left\| (\nabla \log p(\mathbf{X}) \odot \mathbf{h}(\mathbf{X}))' \right\|_1 < +\infty,$$

where  $f(\mathbf{x}) \Big|_{x_j \searrow 0^+}^{x_j \nearrow +\infty} \equiv \lim_{x_j \nearrow +\infty} f(\mathbf{x}) - \lim_{x_j \searrow 0^+} f(\mathbf{x})$ , the loss (3.2) can be rewritten as

$$J_{\mathbf{h}, \mathbb{R}_+^m}(P) \equiv \int_{\mathbb{R}_+^m} p_0(\mathbf{x}) \sum_{j=1}^m \left[ h_j'(x_j) \partial_j (\log p(\mathbf{x})) + h_j(x_j) \partial_{jj} (\log p(\mathbf{x})) + \frac{1}{2} h_j(x_j) [\partial_j (\log p(\mathbf{x}))]^2 \right] d\mathbf{x} \quad (3.3)$$

plus a constant independent of  $p$ . One can thus use the minimizer of the empirical loss (3.3) as an estimator of  $p_0$ .

### 3.2.3 Truncated Score Matching on Bounded Open Subsets of $\mathbb{R}^m$ (Liu and Kanamori, 2019)

Liu and Kanamori (2019) approach the problem of density estimation on a bounded open subset  $\mathcal{D} \subsetneq \mathbb{R}^m$  with a piecewise smooth boundary  $\partial\mathcal{D}$  by minimizing the following “maximally weighted score matching” loss

$$J_{g_0, \mathcal{D}}(P) \equiv \sup_{g \in \mathcal{G}} \frac{1}{2} \int_{\mathbb{R}_+^m} g(\mathbf{x}) p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) - \nabla \log p_0(\mathbf{x})\|_2^2 d\mathbf{x}. \quad (3.4)$$

with  $\mathcal{G} \equiv \{g | g(\mathbf{x}) = 0, \forall \mathbf{x} \in \partial\mathcal{D} \text{ and } g \text{ is } L\text{-Lipschitz continuous}\}$  for some constant  $L > 0$ . They prove that the maximum is obtained with  $g_0(\mathbf{x}) \equiv L \inf_{\mathbf{x}' \in \partial\mathcal{D}} \|\mathbf{x} - \mathbf{x}'\|_2$ , i.e. the  $\ell_2$  distance of  $\mathbf{x}$  to the boundary of  $\mathcal{D}$ , and with integration by parts similar to the previous methods, (3.4) can be estimated using the empirical loss which can be calculated with a closed form.

## 3.3 Generalized Score Matching for General Domains

### 3.3.1 Assumption on the Domain

For any index  $j = 1, \dots, m$ , write  $\mathcal{C}_{j, \mathcal{D}}(\mathbf{x}_{-j}) \equiv \{y \in \mathbb{R} : (y; \mathbf{x}_{-j}) \in \mathcal{D}\}$  as a set-valued function of values  $\mathbf{x}_{-j}$  of the other components. This is the  $j$ -th coordinate of the intersection between  $\mathcal{D}$  and the line  $\{(y; \mathbf{x}_{-j}) : y \in \mathbb{R}\}$ . Also write  $\mathcal{S}_{-j, \mathcal{D}} \equiv \{\mathbf{x}_{-j} : \mathcal{C}_{j, \mathcal{D}}(\mathbf{x}_{-j}) \neq \emptyset\} \subset \mathbb{R}^{m-1}$ , a projection of  $\mathcal{D}$  onto  $\mathbb{R}^{m-1}$ . For notational simplicity we may drop the dependency of  $\mathcal{C}_j$  and  $\mathcal{S}_{-j}$  on  $\mathcal{D}$ .

**Definition 5.** We call a domain  $\mathcal{D} \subseteq \mathbb{R}^m$  componentwise countable if for any index  $j = 1, \dots, m$  and any  $\mathbf{x}_{-j} \in \mathcal{S}_{-j, \mathcal{D}}$ , we can write

$$\mathcal{C}_{j, \mathcal{D}}(\mathbf{x}_{-j}) \equiv \bigcup_{k=1}^{K_j(\mathbf{x}_{-j})} [a_{k,j}(\mathbf{x}_{-j}), b_{k,j}(\mathbf{x}_{-j})], \quad (3.5)$$

where  $1 \leq K_j \leq +\infty$ ,  $-\infty \leq a_{k,j} \leq b_{k,j} \leq +\infty$  for all  $1 \leq k \leq K_j$ , and further if the boundary set

$$\partial\mathcal{D} \equiv \left\{ \mathbf{x} \in \mathcal{D} : \exists j = 1, \dots, m, x_j \in \bigcup_{k=1}^{K_j(\mathbf{x}_{-j})} \{a_{k,j}(\mathbf{x}_{-j}), b_{k,j}(\mathbf{x}_{-j})\} \setminus \{\pm\infty\} \right\} \quad (3.6)$$

is a Lebesgue-null set in  $\mathbb{R}^m$ ; here each bounded interval in (3.5) can be open/closed/half-open half-closed and must be non-overlapping.

That is, for any  $j$  and any  $\mathbf{x}_{-j}$  such that there exists some  $x_j$  with  $(x_j; \mathbf{x}_{-j}) \in \mathcal{D}$ , the intersection between  $\mathcal{D}$  and the line  $\{(y_j; \mathbf{x}_{-j}) : \mathbf{y} \in \mathbb{R}\}$ , when projected to the  $j$ -th dimension, is a union of at most countably many disjoint (bounded or unbounded) intervals in  $\mathbb{R}$ , and its “boundary points” must form a Lebesgue-null set. Note that  $K_j(\mathbf{x}_{-j})$  cannot be 0 since  $\mathcal{C}_{j,\mathcal{D}}(\mathbf{x}_{-j}) \neq \emptyset$  by definition of  $\mathcal{S}_{-j,\mathcal{D}}$ .

### 3.3.2 Generalized Score Matching Loss for General Domains

We first define the notion of the *truncated componentwise distance*.

**Definition 6.** Fix positive constants  $\mathbf{C} \succ \mathbf{0}_m$ . Given a non-empty componentwise countable domain  $\mathcal{D} \subseteq \mathbb{R}^m$  as in Definition 5 and a vector  $\mathbf{x} \in \mathcal{D}$ , define the truncated componentwise distance of  $\mathbf{x}$  to the boundary of  $\mathcal{D}$  as

$$\boldsymbol{\varphi}_{\mathbf{C},\mathcal{D}}(\mathbf{x}) \equiv (\varphi_{C_1,\mathcal{D},1}(\mathbf{x}), \dots, \varphi_{C_m,\mathcal{D},m}(\mathbf{x})) \in \mathbb{R}_+^m, \quad (3.7)$$

$$\varphi_{C_j,\mathcal{D},j}(\mathbf{x}) \equiv \begin{cases} C_j, & \text{if } a_{k,j} = -\infty \text{ and } b_{k,j} = +\infty, \\ \min(C_j, b_{k,j} - x_j), & \text{if } a_{k,j} = -\infty \text{ and } x_j \leq b_{k,j} < +\infty, \\ \min(C_j, x_j - a_{k,j}, b_{k,j} - x_j) & \text{if } -\infty < a_{k,j} \leq x_j \leq b_{k,j} < +\infty, \\ \min(C_j, x_j - a_{k,j}) & \text{if } -\infty < a_{k,j} \leq x_j \text{ and } b_{k,j} = +\infty, \end{cases} \quad (3.8)$$

where  $k$  is such that  $a_{k,j} \leq x_j \leq b_{k,j}$  as in (3.5).

We thus propose the *generalized  $(\mathbf{h}, \mathbf{C}, \mathcal{D})$ -score matching loss* for a general domain  $\mathcal{D}$  similar to (3.2) with  $\mathbf{h}$  applied to  $\boldsymbol{\varphi}_{\mathbf{C},\mathcal{D}}(\mathbf{x})$  instead of to  $\mathbf{x}$ :

**Definition 7.** Suppose the true distribution  $P_0$  has a twice continuously differentiable density  $p_0$  supported on  $\mathcal{D} \subseteq \mathbb{R}^m$ , a non-empty componentwise countable domain (Definition 5). Given positive constants  $\mathbf{C} \succ \mathbf{0}$ , and  $\mathbf{h} : \mathbb{R}_+^m \rightarrow \mathbb{R}_+^m$ ,  $\mathbf{y} \mapsto (h_1(y_1), \dots, h_m(y_m))$  with

$h_1, \dots, h_m : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ , the generalized  $(\mathbf{h}, \mathbf{C}, \mathcal{D})$ -score matching loss in  $P \in \mathcal{P}(\mathcal{D})$  with density  $p$  is defined as

$$J_{\mathbf{h}, \mathbf{C}, \mathcal{D}}(P) \equiv \frac{1}{2} \int_{\mathcal{D}} p_0(\mathbf{x}) \left\| \nabla \log p(\mathbf{x}) \odot (\mathbf{h} \circ \varphi_{\mathbf{C}, \mathcal{D}})^{1/2}(\mathbf{x}) - \nabla \log p_0(\mathbf{x}) \odot (\mathbf{h} \circ \varphi_{\mathbf{C}, \mathcal{D}})^{1/2}(\mathbf{x}) \right\|_2^2 d\mathbf{x}. \quad (3.9)$$

### 3.3.3 Comparison to Previous Work

The intuition of (3.9) is applying the loss from (3.2) to  $(\mathbf{h} \circ \varphi_{\mathbf{C}, \mathcal{D}})$  in place of  $\mathbf{h}$ . Thus, the generalized score matching loss (3.2) in Chapter 2 becomes a special case with  $\mathcal{D} = \mathbb{R}_+^m$  and  $\mathbf{C} = +\infty^m$ , since  $\varphi_{+\infty^m, \mathbb{R}_+^m}(\mathbf{x}) = \mathbf{x}$ . In Chapter 2 we suggested using an  $\mathbf{h}$  with each component a bounded function, which is now incorporated into the definition of  $\varphi$  with finite truncation points  $\mathbf{C}$ . When  $\mathbf{h}(\mathbf{x}) = \mathbf{1}_m$ ,  $(\mathbf{h} \circ \varphi_{\mathbf{C}, \mathcal{D}})^{1/2}(\mathbf{x}) \equiv \mathbf{1}_m$  and it corresponds to the original score-matching for  $\mathbb{R}^m$  in Hyvärinen (2005). If  $\mathbf{h}(\mathbf{x}) = \mathbf{x}^2$ ,  $\mathbf{C} = +\infty^m$  and  $\mathcal{D} \equiv \mathbb{R}_+^m$ ,  $(\mathbf{h} \circ \varphi_{\mathbf{C}, \mathcal{D}})^{1/2}(\mathbf{x}) \equiv x^2$  corresponds to the estimator in Hyvärinen (2007) and Lin et al. (2016).

The  $\varphi_{\mathbf{C}, \mathcal{D}}$  function transforms an  $\mathbf{x} \in \mathcal{D}$  into the component-wise distance vector taking values in  $\mathbb{R}_+^m$ , and the rest remains the same as in Chapter 2. It is thus a natural extension to the work and requires little extra work to be done as  $\varphi_{\mathbf{C}, \mathcal{D}}$  usually has closed-form solution and can be computationally trivial.

We note that for a bounded  $\mathcal{D}$ , our approach is different than using a uniform  $g_0(\mathbf{x})$  the  $\ell_2$  distance of  $\mathbf{x}$  to the boundary of  $\mathcal{D}$  as in Liu and Kanamori (2019), in that we decompose the distance for each component and apply an extra  $\mathbf{h}$  function. In Figure 3.1 we illustrate the two components of our  $\varphi$  function on the 2-d unit disk  $x_1^2 + x_2^2 \leq 1$ , considered below in Example 3.8, compared to the  $g_0$  function in Liu and Kanamori (2019). While  $g_0(\mathbf{x}) = 1 - \sqrt{x_1^2 + x_2^2}$ , choosing  $\mathbf{C} \succeq \mathbf{1}_2$ , we have  $\varphi_1(\mathbf{x}) = \sqrt{1 - x_2^2} - |x_1|$  and  $\varphi_2(\mathbf{x}) = \sqrt{1 - x_1^2} - |x_2|$ . Similarly, in Figure 3.2 we plot these functions for the 2-d unit disk restricted to  $\mathbb{R}_+^2$  considered in Example 3.9, where  $\varphi_1(\mathbf{x}) = \min\{x_1, \sqrt{1 - x_2^2} - x_1\}$ ,  $\varphi_2(\mathbf{x}) = \min\{x_2, \sqrt{1 - x_1^2} - x_2\}$ ,  $g_0(\mathbf{x}) = \min\{x_1, x_2, 1 - \sqrt{x_1^2 + x_2^2}\}$ .

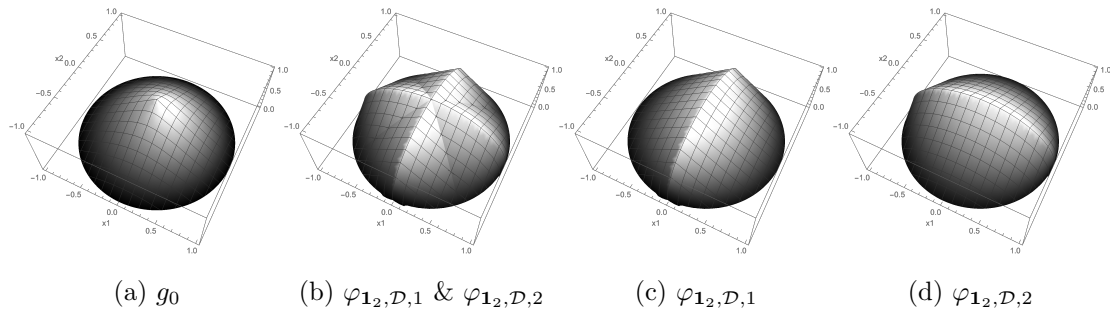


Figure 3.1: Comparison of  $g_0$ ,  $\varphi_{1,2,\mathcal{D},1}$  and  $\varphi_{1,2,\mathcal{D},2}$  on  $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\|_2 < 1\}$ .

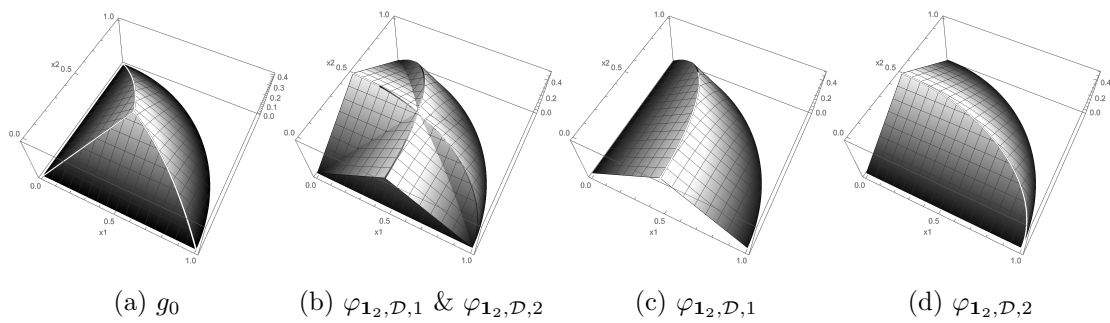


Figure 3.2: Comparison of  $g_0$ ,  $\varphi_{1,2,\mathcal{D},1}$  and  $\varphi_{1,2,\mathcal{D},2}$  on  $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}_+^2 : \|\mathbf{x}\|_2 < 1\}$ .

### 3.3.4 Form of the Truncated Componentwise Distance $\varphi$

**Example 3.6** ( $\mathbb{R}^m$  and  $\mathbb{R}_+^m$ ).  $\varphi_{\mathbf{C}, \mathbb{R}^m}(\mathbf{x}) = \mathbf{C}$  and  $\varphi_{\mathbf{C}, \mathbb{R}_+^m}(\mathbf{x}) = (\min(C_1, x_1), \dots, \min(C_m, x_m))$  by definition.

**Example 3.7** (Unit hypercube). Consider the unit hypercube  $\mathcal{D}$  as the convex hull of  $2^m$  points  $\{\mathbf{x} \in \mathbb{R}^m : |x_1| = \dots = |x_m| = 1/2\}$ . The domain of  $x_j$  is  $[-1/2, 1/2]$  independent of the other components and so  $\varphi_{C_j, \mathcal{D}, j}(\mathbf{x}) = \min\{C_j, 1/2 - |x_j|\}$ . Since the hypercube is bounded by nature, it is natural to drop the truncation by  $C_j$  and simply use  $\varphi_{\mathcal{D}}(\mathbf{x}) = \mathbf{1}_m/2 - |\mathbf{x}|$ .

**Example 3.8** ( $\mathcal{L}^q$  balls). Consider the  $\mathcal{L}^q$  ball with radius  $r > 0$  and  $q \geq 1$  defined by  $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}^m : \|\mathbf{x}\|_q \leq r\}$ . Given a point  $\mathbf{x} \in \mathcal{D}$  and for  $j = 1, \dots, m$ , the domain for  $x_j$  while fixing  $\mathbf{x}_{-j}$  is simply  $[-(r^q - \mathbf{1}^\top |\mathbf{x}_{-j}|^q)^{1/q}, (r^q - \mathbf{1}^\top |\mathbf{x}_{-j}|^q)^{1/q}]$ , and so  $\varphi_{C_j, \mathcal{D}, j}(\mathbf{x}) = \min\{C_j, (r^q - \mathbf{1}^\top |\mathbf{x}_{-j}|^q)^{1/q} - |x_j|\}$ .

**Example 3.9** ( $\mathcal{L}^q$  balls restricted to  $\mathbb{R}_+^m$ ). Consider the intersection between  $\mathbb{R}_+^m$  and the  $\mathcal{L}^q$  ball with radius  $r > 0$  and  $q \geq 1$  defined by  $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}_+^m : \|\mathbf{x}\|_q \leq r\}$ . Given a point  $\mathbf{x} \in \mathcal{D}$  and for  $j = 1, \dots, m$ , the domain for  $x_j$  while fixing  $\mathbf{x}_{-j}$  is simply  $[0, (r^q - \mathbf{1}^\top |\mathbf{x}_{-j}|^q)^{1/q}]$ , and so  $\varphi_{C_j, \mathcal{D}, j}(\mathbf{x}) = \min\{C_j, x_j, (r^q - \mathbf{1}^\top |\mathbf{x}_{-j}|^q)^{1/q} - x_j\}$ .

**Example 3.10** (Complement of  $\mathcal{L}^q$  balls). Consider the complement of the  $\mathcal{L}^q$  ball with radius  $r > 0$  and  $q \geq 1$  defined by  $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}^m : \|\mathbf{x}\|_q > r\}$ . Given an  $\mathbf{x} \in \mathcal{D}$  and  $j$ , the domain for  $x_j$  while fixing  $\mathbf{x}_{-j}$  is  $\mathbb{R}$  if  $\mathbf{1}_{m-1}^\top |\mathbf{x}_{-j}|^q > r^q$ , or  $(-\infty, -(r^q - \mathbf{1}_{m-1}^\top |\mathbf{x}_{-j}|^q)^{1/q}) \cup ((r^q - \mathbf{1}_{m-1}^\top |\mathbf{x}_{-j}|^q)^{1/q}, +\infty)$  otherwise. So

$$\varphi_{C_j, \mathcal{D}, j}(\mathbf{x}) = \begin{cases} C_j & \text{if } \mathbf{1}_{m-1}^\top |\mathbf{x}_{-j}|^q > r^q, \\ \min\{C_j, |x_j| - (r^q - \mathbf{1}_{m-1}^\top |\mathbf{x}_{-j}|^q)^{1/q}\} & \text{otherwise.} \end{cases}$$

**Example 3.11** (Complement of  $\mathcal{L}^q$  balls restricted to  $\mathbb{R}_+^m$ ). Consider the intersection between  $\mathbb{R}_+^m$  and the complement of the  $\mathcal{L}^q$  ball with radius  $r > 0$  and  $q \geq 1$ :  $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}_+^m : \|\mathbf{x}\|_q > r\}$ . Given an  $\mathbf{x} \in \mathcal{D}$  and  $j$ , the domain for  $x_j$  while fixing  $\mathbf{x}_{-j}$  is  $\mathbb{R}_+$  if  $\mathbf{1}_{m-1}^\top |\mathbf{x}_{-j}|^q > r^q$ , or

$((r^q - \mathbf{1}_{m-1}^\top |\mathbf{x}_{-j}|^q)^{1/q}, +\infty)$  otherwise. So

$$\varphi_{C_j, \mathcal{D}, j}(\mathbf{x}) = \begin{cases} \min\{C_j, x_j\} & \text{if } \mathbf{1}_{m-1}^\top |\mathbf{x}_{-j}|^q > r^q, \\ \min\{C_j, x_j - (r^q - \mathbf{1}_{m-1}^\top |\mathbf{x}_{-j}|^q)^{1/q}\} & \text{otherwise.} \end{cases}$$

**Example 3.12** (Complicated domains defined by inequality constraints). A domain  $\mathcal{D}$  may be determined by a series of intersections/unions of regions determined by inequality constraints, e.g.  $\mathcal{D} = \{\mathbf{x} \in \mathbb{R}^m : (f_1(\mathbf{x}) \leq c_1 \wedge f_2(\mathbf{x}) \leq c_2) \vee f_3(\mathbf{x}) \geq c_3\}$ . In this case, for example, to calculate  $\varphi_{\mathcal{C}, \mathcal{D}}$  one plug in  $\mathbf{x}_{-j}$  as given and solve numerically  $f_i(x_j; \mathbf{x}_{-j}) = c_i$  for  $i = 1, 2, 3$ , and obtain the boundary points for  $x_j$  using simple algorithms for interval unions/intersections. This is implemented in the package `genscore` for some types of polynomial  $f_i$  and arbitrary intersections/unions.

**Example 3.13** (Standard Simplices). Consider  $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}_+^m \mid \mathbf{x} \succ \mathbf{0}, \mathbf{1}_m^\top \mathbf{x} = 1\}$ , the standard simplex. Since this has zero Lebesgue measure in  $\mathbb{R}^m$ , we drop the last coordinate and consider instead  $\mathcal{D}_{-m} \equiv \{\mathbf{x}_{-m} \in \mathbb{R}_+^{m-1} \mid \mathbf{x}_{-m} \succ \mathbf{0}, \mathbf{1}_{m-1}^\top \mathbf{x}_{-m} < 1\}$  as a set in  $\mathbb{R}^{m-1}$ . For  $j = 1, \dots, m-1$ , the domain of  $x_j$  given  $\mathbf{x}_{-j, -m}$  for  $\mathbf{x}_{-m} \in \mathcal{D}_{-m}$  is thus  $(0, 1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j, -m})$ , and so  $\varphi_{C_j, \mathcal{D}_{-m}, j}(\mathbf{x}) = \min\{C_j, x_j, 1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m}\} = \min\{C_j, x_j, x_m\}$ . In practice since the standard simplex is naturally bounded by the unit cube, it is natural to drop the truncation by  $C_j$  and simply use  $\varphi_{\mathcal{D}_{-m}, j}(\mathbf{x}) = \min\{x_j, 1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m}\} = \min\{x_j, x_m\}$ .

### 3.3.5 The Empirical Generalized Score Matching Loss

From this section, we drop the dependence of  $\varphi$  on  $\mathbf{C}$  and  $\mathcal{D}$  in notation for simplicity.

**Lemma 12.** Suppose  $\mathbf{C} \succ \mathbf{0}$ ,  $p_0(\mathbf{x}) > 0$  for almost every  $\mathbf{x} \in \mathcal{D}$  and  $h_1(y), \dots, h_m(y) > 0$  for all  $y > 0$ . Then  $J_{\mathbf{h}, \mathbf{C}, \mathcal{D}}(P) = 0$  if and only if  $p_0 = p$  for a.e.  $\mathbf{x} \in \mathcal{D}$ .

*Proof of Lemma 12.* By the assumption on  $\mathcal{D}$  and definition of  $\varphi$ ,  $\varphi(\mathbf{x}) \succ \mathbf{0}$  for almost every  $\mathbf{x} \in \mathcal{D}$ , and thus  $(\mathbf{h} \circ \varphi)(\mathbf{x}) \succ \mathbf{0}$  for a.e.  $\mathbf{x} \in \mathcal{D}$ . So we have  $J_{\mathbf{h}, \mathbf{C}, \mathcal{D}}(P) = 0$  if and only if  $\nabla \log p_0(\mathbf{x}) = \nabla \log p(\mathbf{x})$  for a.e.  $\mathbf{x} \in \mathcal{D}$ , i.e.  $\log p_0(\mathbf{x}) = \log p(\mathbf{x}) + c_0$  for a.e.  $\mathbf{x} \in \mathcal{D}$  for some constant  $c_0$ , or  $p_0(\mathbf{x}) = c_1 \cdot p(\mathbf{x})$  for a.e.  $\mathbf{x} \in \mathcal{D}$  for some non-zero constant  $c_1 \equiv \exp(c_0)$ . Since  $p_0$  and  $p$  both integrate to 1 over  $\mathcal{D}$ , we have  $c_1 = 1$  and  $p_0 = p$  for a.e.  $\mathbf{x} \in \mathcal{D}$ .  $\square$

**Lemma 13.** *Similar to (A0.1) and (A0.2) in Section 3.2.2, assume the following assumptions hold for all  $p \in \mathcal{P}(\mathcal{D})$ ,*

$$(A.1) \quad p_0(x_j; \mathbf{x}_{-j}) h_j(\varphi_j(\mathbf{x})) \partial_j \log p(x_j; \mathbf{x}_{-j}) \Big|_{x_j \searrow a_k(\mathbf{x}_{-j})^+}^{x_j \nearrow b_k(\mathbf{x}_{-j})^-} = 0 \text{ for all } k = 1, \dots, K_j(\mathbf{x}_{-j}) \text{ and } \mathbf{x}_{-j} \in \mathcal{S}_{-j} \text{ for all } j;$$

$$(A.2) \quad \int_{\mathcal{D}} p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) \odot (\mathbf{h} \circ \varphi)^{1/2}(\mathbf{x})\|_2^2 d\mathbf{x} < +\infty, \\ \int_{\mathcal{D}} p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) \odot (\mathbf{h} \circ \varphi)(\mathbf{x})\|_1 d\mathbf{x} < +\infty.$$

Also assume that

$$(A.3) \quad \forall j = 1, \dots, m \text{ and a.e. } \mathbf{x}_{-j} \in \mathcal{S}_{-j}, h_j \text{ is absolutely continuous in any bounded subinterval of } \mathcal{C}_j(\mathbf{x}_{-j}).$$

(This implies the same for  $(h_j \circ \varphi_j)$  and also that  $\partial_j(h_j \circ \varphi_j)$  exists a.e.) Then

$$J_{\mathbf{h}, \mathcal{C}, \mathcal{D}}(P) \equiv \frac{1}{2} \sum_{j=1}^m \int_{\mathcal{D}} p_0(\mathbf{x}) \cdot (h_j \circ \varphi_j)(\mathbf{x}) \cdot [\partial_j \log p(\mathbf{x})]^2 d\mathbf{x} \\ + \sum_{j=1}^m \int_{\mathcal{D}} p_0(\mathbf{x}) \cdot \partial_j [(h_j \circ \varphi_j)(\mathbf{x}) \cdot \partial_j \log p(\mathbf{x})] d\mathbf{x} \quad (3.10)$$

plus a constant depending on  $p_0$  only and independent of  $p$ .

The proof of the lemma is in Appendix. The lemma enables us to estimate the population loss using the empirical loss

$$\hat{J}_{\mathbf{h}, \mathcal{C}, \mathcal{D}}(P) = \frac{1}{2} \sum_{j=1}^m \sum_{i=1}^n \frac{1}{2} (h_j \circ \varphi_j)(\mathbf{x}^{(i)}) \cdot [\partial_j \log p(\mathbf{x}^{(i)})]^2 + \partial_j [(h_j \circ \varphi_j)(\mathbf{x}^{(i)}) \cdot \partial_j \log p(\mathbf{x}^{(i)})]. \quad (3.11)$$

As the canonical choices of  $\mathbf{h}$  are power functions in  $\mathbf{x}$ , we give the following sufficient conditions for the assumptions in the lemma.

**Proposition 14.** *Suppose for all  $j = 1, \dots, m$ ,  $h_j(x_j) = x_j^{\alpha_j}$  for some  $\alpha_j > 0$ . Suppose in addition that for all  $j$  and  $\mathbf{x}_{-j} \in \mathcal{S}_{-j}$  and all  $p \in \mathcal{P}$  we have*

- (1)  $p_0(x_j; \mathbf{x}_{-j}) \partial_j \log p(x_j; \mathbf{x}_{-j}) = o(1/(x_j - c_{k,j})^\alpha)$  as  $x_j \nearrow c_{k,j} \equiv b_{k,j} \neq +\infty$  or as  $x_j \searrow c_{k,j} \equiv a_{k,j} \neq -\infty$  for all  $k$ , and
- (2)  $p_0(x_j; \mathbf{x}_{-j}) \partial_j \log p(x_j; \mathbf{x}_{-j}) \rightarrow 0$  as  $x_j \nearrow +\infty$  if  $\mathcal{C}_j(\mathbf{x}_{-j})$  is unbounded from above, and as  $x_j \searrow -\infty$  if  $\mathcal{C}_j(\mathbf{x}_{-j})$  is unbounded from below.

Then (A.1) and (A.3) are satisfied.

*Proof of Proposition 14.* (A.3) is satisfied by the property of  $h_j$ . (A.1) is trivial since by construction  $(h_j \circ \varphi_j)(\mathbf{x})$  becomes  $|x_j - c_{k,j}|^\alpha$  as  $x_j \rightarrow c_{k,j} \in \cup_{k=1}^{K_j} \{a_{k,j}, b_{k,j}\}$  and also  $(h_j \circ \varphi_j)$  is bounded by  $C^\alpha$  as  $x_j \nearrow +\infty$  or  $x_j \searrow -\infty$ , if applicable.  $\square$

According to Lemma 12 the method only makes sense if the domain  $\mathcal{D}$  has a positive Lebesgue measure in  $\mathbb{R}^m$ . However, for some Lebesgue null sets one may apply appropriate transforms under which  $\mathcal{D}$  becomes a set with positive Lebesgue measure in lower dimensions with dimensionality reduction. One important example is the  $(m-1)$ -dimensional standard simplex  $\{\mathbf{x} \in \mathbb{R}_+^m \mid \mathbf{x} \succ \mathbf{0}, \mathbf{1}_m^\top \mathbf{x} = 1\}$ , a null set since it is a subset of the  $(m-1)$ -dimensional hyperplane  $\{\mathbf{1}_m^\top \mathbf{x} = 1\}$ . Simply dropping one of the coordinate, e.g.  $x_m$ , and substituting  $x_m = 1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m}$  in  $p_0$  and  $p$  we can thus apply the method to  $\{\mathbf{x}_{-m} \in \mathbb{R}_+^{m-1} \mid \mathbf{x}_{-m} \succ \mathbf{0}, \mathbf{1}_{m-1}^\top \mathbf{x}_{-m} < 1\}$ , a set with positive measure in  $\mathbb{R}^{m-1}$ . We further discuss this example in Section 3.6.

### 3.3.6 Extension of $g_0$ in Liu and Kanamori (2019)

We extend the method in Liu and Kanamori (2019) introduced in Section 3.2.3 from bounded domains to any *componentwise countable* domain  $\mathcal{D}$  defined in Definition 5 by defining the loss of the same form as (3.4):

$$J_{g_0, C, \mathcal{D}}(P) \equiv \sup_{g \in \mathcal{G}} \frac{1}{2} \int_{\mathbb{R}_+^m} g(\mathbf{x}) p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) - \nabla \log p_0(\mathbf{x})\|_2^2 d\mathbf{x},$$

but with  $\mathcal{G} \equiv \{g \mid g(\mathbf{x}) = 0, \forall \mathbf{x} \in \partial \mathcal{D}, g \text{ is } L\text{-Lipschitz continuous and } g \leq C\}$  instead, for the same Lipschitz constant  $L > 0$  but with an extra truncation constant  $C > 0$ . Here the

boundary set  $\partial\mathcal{D}$  is as defined in (3.6). Following the same proof as for their Proposition 1 and dropping the Lipschitz constant  $L$  by replacing  $C$  with  $C/L$  (or equivalently choosing  $L = 1$  without loss of generality), it is easy to see that the maximum is obtained at

$$g_0(\mathbf{x}) \equiv \min \left\{ C, \inf_{\mathbf{x}' \in \partial\mathcal{D}} \|\mathbf{x} - \mathbf{x}'\|_2 \right\}, \quad (3.12)$$

the  $\ell_2$  distance of  $\mathbf{x}$  to  $\partial\mathcal{D}$  truncated above by  $C$ , which naturally extends the method in Liu and Kanamori (2019) to unbounded domains. In the special case where  $\partial\mathcal{D} = \emptyset$ , we must have  $\mathcal{D} \equiv \mathbb{R}^m$ , and (3.12) becomes constant  $C$  by the convention  $\inf \emptyset = +\infty$ , which coincides with the original score matching in Hyvärinen (2005).

Assuming assumptions (A.1) and (A.2) in Lemma 13 hold when replacing  $(\mathbf{h} \circ \boldsymbol{\varphi})(\mathbf{x})$  by  $g_0(\mathbf{x})\mathbf{1}_m$  and  $h_j(\varphi_j(\mathbf{x}))$  by  $g_0(\mathbf{x})$ , the same conclusion there applies, i.e.

$$J_{g_0, C, \mathcal{D}}(P) \equiv \frac{1}{2} \sum_{j=1}^m \int_{\mathcal{D}} p_0(\mathbf{x}) g_0(\mathbf{x}) [\partial_j \log p(\mathbf{x})]^2 d\mathbf{x} + \sum_{j=1}^m \int_{\mathcal{D}} p_0(\mathbf{x}) \partial_j [g_0(\mathbf{x}) \partial_j \log p(\mathbf{x})] d\mathbf{x}$$

plus a constant depending on  $p_0$  only and independent of  $p$ , the same empirical loss as (13) in Liu and Kanamori (2019) with an additional truncation from above applied to  $g_0$ . The proof is in the same spirit as that for Lemma 13 and is thus omitted.

### 3.4 Exponential Families and $a$ - $b$ Models

#### 3.4.1 Exponential Families

Consider the case where  $\mathcal{P}(\mathcal{D}) \equiv \{p_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \Theta \subset \mathbb{R}^r\}$  for some  $r = 1, 2, \dots$  is an exponential family with continuous distributions supported on  $\mathcal{D}$  with densities of the form

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{t}(\mathbf{x}) - \psi(\boldsymbol{\theta}) + b(\mathbf{x}), \quad \mathbf{x} \in \mathcal{D}, \quad (3.13)$$

where  $\boldsymbol{\theta} \in \mathbb{R}^r$  is the unknown canonical parameter of interest,  $\mathbf{t}(\mathbf{x}) \in \mathbb{R}^r$  are the sufficient statistics,  $\psi(\boldsymbol{\theta})$  is the normalizing constant, and  $b(\mathbf{x})$  is the base measure. The empirical loss  $\hat{J}_{\mathbf{h}, C, \mathcal{D}}$  (3.11) can then be written as a quadratic function in the canonical parameter:

$$\hat{J}_{\mathbf{h}, C, \mathcal{D}}(p_{\boldsymbol{\theta}}) = \frac{1}{2} \boldsymbol{\theta}^\top \mathbf{\Gamma}(\mathbf{x}) \boldsymbol{\theta} - \mathbf{g}(\mathbf{x})^\top \boldsymbol{\theta} + \text{const}, \quad \text{with} \quad (3.14)$$

$$\mathbf{\Gamma}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m (h_j \circ \varphi_j) (\mathbf{X}^{(i)}) \partial_j \mathbf{t}(\mathbf{X}^{(i)}) (\partial_j \mathbf{t}(\mathbf{X}^{(i)}))^\top \quad \text{and} \quad (3.15)$$

$$\begin{aligned} \mathbf{g}(\mathbf{x}) = & -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m [(h_j \circ \varphi_j) (\mathbf{X}^{(i)}) \partial_j b(\mathbf{X}^{(i)}) \partial_j \mathbf{t}(\mathbf{X}^{(i)}) \\ & + (h_j \circ \varphi_j) (\mathbf{X}^{(i)}) \partial_{jj} \mathbf{t}(\mathbf{X}^{(i)}) + \partial_j (h_j \circ \varphi_j) (\mathbf{X}^{(i)}) \partial_j \mathbf{t}(\mathbf{X}^{(i)})], \end{aligned} \quad (3.16)$$

where  $\partial_j \mathbf{t}(\mathbf{x}) = (\partial_j t_1(\mathbf{x}), \dots, \partial_j t_r(\mathbf{x}))^\top \in \mathbb{R}^r$ . Note that (3.15) and (3.16) are sample averages of functions in the data matrix  $\mathbf{x}$  only. These forms are an exact analog of those in Theorem 3. One thus have the following consistency result similar to 4:

**Theorem 15.** *Suppose the true density is  $p_0 \equiv p_{\boldsymbol{\theta}_0}$  and that*

(C1)  $\mathbf{\Gamma}$  is almost surely invertible, and

(C2)  $\mathbf{\Gamma}_0 \equiv \mathbb{E}_{p_0} \mathbf{\Gamma}(\mathbf{x})$ ,  $\mathbf{\Gamma}_0^{-1}$ ,  $\mathbf{g}_0 \equiv \mathbb{E}_{p_0} \mathbf{g}(\mathbf{x})$ , and  $\mathbf{\Sigma}_0 \equiv \mathbb{E}_{p_0} [(\mathbf{\Gamma}(\mathbf{x})\boldsymbol{\theta}_0 - \mathbf{g}(\mathbf{x})) (\mathbf{\Gamma}(\mathbf{x})\boldsymbol{\theta}_0 - \mathbf{g}(\mathbf{x}))^\top]$  exist and are component-wise finite.

Then the minimizer of (3.14) is a.s. unique with closed form solution  $\hat{\boldsymbol{\theta}} \equiv \mathbf{\Gamma}(\mathbf{x})^{-1} \mathbf{g}(\mathbf{x})$  with

$$\hat{\boldsymbol{\theta}} \rightarrow_{a.s.} \boldsymbol{\theta}_0 \quad \text{and} \quad \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \rightarrow_d \mathcal{N}_r(\mathbf{0}, \mathbf{\Gamma}_0^{-1} \mathbf{\Sigma}_0 \mathbf{\Gamma}_0^{-1}) \quad \text{as } n \rightarrow \infty.$$

### 3.4.2 Pairwise Interaction Power a-b Models

As an important example of graphical models from exponential family distributions, in this chapter we focus on the class of the pairwise interaction power models introduced in (1.1):

$$p_{\boldsymbol{\eta}, \mathbf{K}}(\mathbf{x}) \propto \exp \left( -\frac{1}{2a} \mathbf{x}^{a\top} \mathbf{K} \mathbf{x}^a + \frac{1}{b} \boldsymbol{\eta}^\top \mathbf{x}^b \right) \mathbf{1}_{\mathcal{D}}(\mathbf{x}) \quad (3.17)$$

for which we treat  $\mathbf{x}^0 \equiv \log \mathbf{x}$  and  $1/0 \equiv 1$ , on a domain  $\mathcal{D} \subseteq \mathbb{R}^m$  with a positive measure. Here  $a \geq 0$  and  $b \geq 0$  are known constants and the interaction matrix  $\mathbf{K} \in \mathbb{R}^{m \times m}$  and the linear vector  $\boldsymbol{\eta} \in \mathbb{R}^m$  are unknown parameters of interest. As in Chapter 2, our focus will be on the support of  $\mathbf{K}$ ,  $S(\mathbf{K}) = \{(i, j) : \kappa_{ij} \neq 0\}$ , that defines the conditional independence

graph of  $\mathbf{X} \sim p_{\boldsymbol{\eta}, \mathbf{K}}$ , while simultaneously estimating the nuisance parameter  $\boldsymbol{\eta}$  unless it is known or assumed to be  $\mathbf{0}$ .

When  $a = b = 1$  (3.17) is simply a truncated Gaussian model. When  $a = b = 1/2$  we obtain the exponential square-root graphical model in Inouye et al. (2016). The gamma model considered in Chapter 2 corresponds to  $a = 1/2$  and  $b = 0$ . As we discuss in Section 3.6.2,  $A^{m-1}$  models in Aitchison (1985) are covered by  $a = b = 0$  with a simplex domain  $\mathcal{D}$ .

### 3.4.3 Regularized Score Matching

In high-dimensional settings where the number of parameters  $r$  to be estimated is larger than the sample size  $n$ , one usually add some form of regularization for consistent estimation. For exponential family distributions, as in Chapter 2, we add an  $\ell_1$  penalty on  $\boldsymbol{\theta}$  to the loss in (3.14), while multiplying the diagonals of the  $\boldsymbol{\Gamma}$  by a *diagonal multiplier*  $\delta > 1$ :

**Definition 8.** Let  $\delta > 1$ , and  $\boldsymbol{\Gamma}_\delta(\mathbf{x})$  be  $\boldsymbol{\Gamma}(\mathbf{x})$  with diagonal entries multiplied by  $\delta$ . For exponential family distributions in (3.13), the regularized generalized  $(\mathbf{h}, \mathbf{C}, \mathcal{D})$ -score matching loss is defined as

$$\hat{J}_{\mathbf{h}, \mathbf{C}, \mathcal{D}, \lambda, \delta}(p_{\boldsymbol{\theta}}) \equiv \frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{\Gamma}_\delta(\mathbf{x}) \boldsymbol{\theta} - \mathbf{g}(\mathbf{x})^\top \boldsymbol{\theta} + \lambda \|\boldsymbol{\theta}\|_1. \quad (3.18)$$

The multiplier  $\delta > 1$ , together with the  $\ell_1$  penalty, resembles an elastic net penalty and prevents the loss in (3.18) from being unbounded from below for smaller  $\lambda$ , in which case there can be infinitely many minimizers. For a more detailed discussion on this, refer to Section 2.4. The minimization of (3.18) can be efficiently done using coordinate-descent with warm starts, along with other computational details discussed in Section 2.5.3.

## 3.5 $a$ - $b$ Models on Domains with Positive Measure

Throughout this section, we assume that  $\mathcal{D}$  has positive Lebesgue measure in  $\mathbb{R}^m$ .

### 3.5.1 Finite Normalizing Constant and Validity of Score Matching

The following theorem gives detailed sufficient conditions for the  $a$ - $b$  density  $p_{\boldsymbol{\eta}, \mathbf{K}}$  (3.17) to be a proper density on  $\mathcal{D} \subseteq \mathbb{R}^m$  with positive Lebesgue measure.

**Theorem 16** (Sufficient conditions for finite normalizing constant). *Denote the closure of the range of  $x_j$  in the domain  $\mathcal{D}$  as  $\rho_j(\mathcal{D}) \equiv \overline{\{x_j : \mathbf{x} \in \mathcal{D}\}}$ . If any of the following conditions holds, the density in (3.17) is a proper density, i.e. the right-hand of (3.17) is integrable over  $\mathcal{D}$ :*

(CC1)  $a > 0, b > 0, \mathcal{D}$  bounded;

(CC2)  $a > 0, b > 0, \mathbf{v}^a \top \mathbf{K} \mathbf{v}^a > 0 \forall \mathbf{v} \in \mathcal{D} \setminus \{\mathbf{0}\}$ , and either  $2a > b$  or  $\boldsymbol{\eta} \top \mathbf{v}^b \leq 0 \forall \mathbf{v} \in \mathcal{D}$ ;

(CC3)  $a > 0, b = 0, \eta_j > -1$  for all  $j$  s.t.  $0 \in \rho_j(\mathcal{D})$ , and one of the following:

(i)  $\mathcal{D}$  bounded;

(ii)  $\mathcal{D}$  is unbounded and  $\mathbf{v}^a \top \mathbf{K} \mathbf{v}^a > 0 \forall \mathbf{v} \in \mathcal{D} \setminus \{\mathbf{0}\}$ ;

(iii)  $\mathcal{D}$  is unbounded,  $\mathbf{v}^a \top \mathbf{K} \mathbf{v}^a \geq 0 \forall \mathbf{v} \in \mathcal{D}$  and  $\eta_j < -1$  for all  $j$  s.t.  $\rho_j(\mathcal{D})$  is unbounded  
( $\Rightarrow \rho_j(\mathcal{D}) = (0, +\infty)$  is not allowed for any  $j$ );

(CC4)  $a = 0, \mathcal{D}$  is bounded and  $0 \notin \rho_j(\mathcal{D})$  for all  $j$ ;

(CC5)  $a = 0, b = 0, \log(\mathbf{x}) \top \mathbf{K} \log(\mathbf{x}) > 0 \forall \mathbf{x} \in \mathcal{D}$ ;

(CC6)  $a = 0, b > 0, \log(\mathbf{x}) \top \mathbf{K} \log(\mathbf{x}) > 0 \forall \mathbf{x} \in \mathcal{D}$  and  $\eta_j \leq 0$  for all  $j$  s.t.  $\rho_j(\mathcal{D})$  is unbounded;

(CC7)  $a = 0, b > 0, \log(\mathbf{x}) \top \mathbf{K} \log(\mathbf{x}) \geq 0 \forall \mathbf{x} \in \mathcal{D}$  and  $\eta_j < 0$  for all  $j$  s.t.  $\rho_j(\mathcal{D})$  is unbounded.

In the centered case where  $\boldsymbol{\eta} = \mathbf{0}$  is known, any condition in terms of  $b$  and  $\boldsymbol{\eta}$  can be ignored.

To simplify our discussion, the following corollary gives a simpler set of sufficient conditions for integrability of the density.

**Corollary 17** (Sufficient conditions for finite normalizing constant; simplified). *Suppose*

*(CC0\*)  $\mathbf{K}$  is positive definite*

*and one of the following conditions holds, then the right-hand side of (3.17) is integrable over  $\mathcal{D}$ .*

*(CC1\*)  $2a > b > 0$  or  $a = b = 0$ ;*

*(CC2\*)  $a > 0, b = 0, \boldsymbol{\eta} \succ -\mathbf{1}_m$ ;*

*(CC3\*)  $a = 0, b > 0, \eta_j \leq 0$  for any  $j$  such that  $x_j$  is unbounded in  $\mathcal{D}$ .*

*In the non-centered case where  $\boldsymbol{\eta} \equiv \mathbf{0}$ , (CC0\*) is sufficient.*

For simplicity, we use the set of conditions (CC0\*) through (CC3\*) throughout the rest of the chapter. The following theorem gives sufficient conditions on  $\mathbf{h}$  such that conditions (A.1)–(A.3) in Lemma 13 for score matching are satisfied.

**Theorem 18** (Sufficient conditions that satisfy assumptions for score matching). *Suppose (CC0\*) and one of (CC1\*) through (CC3\*) holds, and  $\mathbf{h}(\mathbf{x}) = (x_1^{\alpha_1}, \dots, x_m^{\alpha_m})$ , where*

*(1) if  $a > 0$  and  $b > 0$ ,  $\alpha_j > \max\{0, 1 - a, 1 - b\}$ ;*

*(2) if  $a > 0$  and  $b = 0$ ,  $\alpha_j > 1 - \eta_{0,j}$ ;*

*(3) if  $a = 0$ ,  $\alpha_j \geq 0$ .*

*Then conditions (A.1), (A.2) and (A.3) in Lemma 13 are satisfied and the equivalent form of the generalized score matching loss (3.10) holds, and the empirical loss (3.11) is valid. In the centered case with  $\boldsymbol{\eta} \equiv \mathbf{0}$ , it suffices to have  $a > 0$  and  $\alpha_j > \max\{0, 1 - a\}$  or  $a = 0$  and  $\alpha_j \geq 0$ .*

### 3.5.2 Estimation

Let  $\Psi \equiv [\mathbf{K}^\top \boldsymbol{\eta}]^\top \in \mathbb{R}^{(m+1) \times m}$ . In this section, we give the form of  $\mathbf{\Gamma} \in \mathbb{R}^{(m+1)m \times (m+1)m}$  and  $\mathbf{g} \in \mathbb{R}^{(m+1)m}$  in the unpenalized loss  $\frac{1}{2} \text{vec}(\Psi)^\top \mathbf{\Gamma} \text{vec}(\Psi) - \mathbf{g}^\top \text{vec}(\Psi)$  following (3.15)–(3.16).  $\mathbf{\Gamma}$  is block-diagonal, with the  $j$ -th  $\mathbb{R}^{(m+1) \times (m+1)}$  block

$$\begin{aligned} \mathbf{\Gamma}_j(\mathbf{x}) &\equiv \begin{bmatrix} \mathbf{\Gamma}_{\mathbf{K},j} & \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},j} \\ \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},j}^\top & \boldsymbol{\gamma}_{\boldsymbol{\eta},j} \end{bmatrix} \\ &\equiv \frac{1}{n} \sum_{i=1}^n \begin{bmatrix} (h_j \circ \varphi_j)(\mathbf{X}^{(i)}) X_j^{(i)2a-2} \mathbf{X}^{(i)a} \mathbf{X}^{(i)a\top} & -(h_j \circ \varphi_j)(\mathbf{X}^{(i)}) X_j^{(i)a+b-2} \mathbf{X}^{(i)a} \\ -(h_j \circ \varphi_j)(\mathbf{X}^{(i)}) X_j^{(i)a+b-2} \mathbf{X}^{(i)a\top} & (h_j \circ \varphi_j)(\mathbf{X}^{(i)}) X_j^{(i)2b-2} \end{bmatrix}, \end{aligned}$$

and  $\mathbf{g} \equiv \text{vec}([\mathbf{g}_{\mathbf{K}}^\top \mathbf{g}_{\boldsymbol{\eta}}]^\top) \in \mathbb{R}^{(m+1)m}$ , where  $\mathbf{g}_{\mathbf{K}} \in \mathbb{R}^{m \times m}$  and  $\mathbf{g}_{\boldsymbol{\eta}} \in \mathbb{R}^m$  correspond to  $\mathbf{K}$  and  $\boldsymbol{\eta}$ , respectively. The  $j$ -th column of  $\mathbf{g}_{\mathbf{K}}$  is

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left( \partial_j (h_j \circ \varphi_j)(\mathbf{X}^{(i)}) X_j^{(i)a-1} + (a-1)(h_j \circ \varphi_j)(\mathbf{X}^{(i)}) X_j^{(i)a-2} \right) \mathbf{X}^{(i)a} \\ + a(h_j \circ \varphi_j)(\mathbf{X}^{(i)}) X_j^{(i)2a-2} \mathbf{e}_{j,m}, \end{aligned} \quad (3.19)$$

where  $\mathbf{e}_{j,m} \in \mathbb{R}^m$  has 1 at the  $j$ -th composition and 0 elsewhere, and the  $j$ -th entry of  $\mathbf{g}_2$  is

$$\frac{1}{n} \sum_{i=1}^n -\partial_j (h_j \circ \varphi_j)(\mathbf{X}^{(i)}) X_j^{(i)b-1} - (b-1)(h_j \circ \varphi_j)(\mathbf{X}^{(i)}) X_j^{(i)b-2}. \quad (3.20)$$

(If  $a = 0$ , set the coefficients  $(a-1) \leftarrow -1$  and  $a \leftarrow 1$  in (3.19); for  $b = 0$  let  $(b-1) \leftarrow -1$  in the second term of (3.20).)

As in Chapter 2 we only apply the diagonal multiplier  $\delta$  to the diagonals of  $\mathbf{\Gamma}_{\mathbf{K},j} \in \mathbb{R}^{m \times m}$ , not  $\boldsymbol{\gamma}_{\boldsymbol{\eta},j} \in \mathbb{R}$ . Note that each block of  $\mathbf{\Gamma}$  and  $\mathbf{g}$  correspond to each column of  $\Psi$ , i.e.  $(\boldsymbol{\kappa}_j, \eta_j) \in \mathbb{R}^{m+1}$ . In the penalized generalized score-matching loss (3.18), the penalty  $\lambda$  for  $\mathbf{K}$  and  $\boldsymbol{\eta}$  can be different,  $\lambda_{\mathbf{K}}$  and  $\lambda_{\boldsymbol{\eta}}$ , respectively, as long as the ratio  $\lambda_{\boldsymbol{\eta}}/\lambda_{\mathbf{K}}$  is fixed. For  $\mathbf{K}$  we follow the convention that we penalize its off-diagonal entries only. That is,

$$\frac{1}{2} \text{vec}(\Psi)^\top \mathbf{\Gamma}_\delta(\mathbf{x}) \text{vec}(\Psi) - \mathbf{g}(\mathbf{x})^\top \text{vec}(\Psi) + \lambda_{\mathbf{K}} \|\mathbf{K}_{\text{off}}\|_1 + \lambda_{\boldsymbol{\eta}} \|\boldsymbol{\eta}\|_1. \quad (3.21)$$

In the case where we do not penalize  $\boldsymbol{\eta}$ , i.e.  $\lambda_\eta = 0$ , as in Chapter 2 we can simply profile out  $\boldsymbol{\eta}$ , solve for  $\hat{\boldsymbol{\eta}} = \boldsymbol{\Gamma}_\eta^{-1} \left( \mathbf{g}_\eta - \boldsymbol{\Gamma}_{\mathbf{K},\eta}^\top \text{vec}(\hat{\mathbf{K}}) \right)$ , plug this back in and rewrite the loss in  $\mathbf{K}$  only. Let  $\boldsymbol{\Gamma}_{\delta,\mathbf{K}} \in \mathbb{R}^{m^2 \times m^2}$  be the block-diagonal matrix with blocks  $\boldsymbol{\Gamma}_{\mathbf{K},j}$  and diagonal multiplier  $\delta$ , and let  $\boldsymbol{\Gamma}_{\mathbf{K},\eta} \in \mathbb{R}^{m^2 \times m}$  and  $\boldsymbol{\Gamma}_\eta \in \mathbb{R}^{m \times m}$  be the (block-)diagonal matrices with blocks  $\boldsymbol{\gamma}_{\mathbf{K},\eta,j}$  and  $\gamma_{\eta,j}$ , respectively. Further denote  $\boldsymbol{\Gamma}_{\delta,*}$  as the Schur complement of  $\boldsymbol{\Gamma}_{\delta,\mathbf{K}}$  of  $\begin{bmatrix} \boldsymbol{\Gamma}_{\delta,\mathbf{K}} & \boldsymbol{\Gamma}_{\mathbf{K},\eta} \\ \boldsymbol{\Gamma}_{\mathbf{K},\eta}^\top & \boldsymbol{\Gamma}_\eta \end{bmatrix}$ , i.e.  $\boldsymbol{\Gamma}_{\delta,\mathbf{K}} - \boldsymbol{\Gamma}_{\mathbf{K},\eta} \boldsymbol{\Gamma}_\eta^{-1} \boldsymbol{\Gamma}_{\mathbf{K},\eta}^\top$ , which is guaranteed to be positive definite for  $\delta > 1$ . Then the profiled loss is

$$\hat{J}_{h,\mathcal{C},\mathcal{D},\lambda,\delta,*}(p_{\mathbf{K}}) \equiv \frac{1}{2} \text{vec}(\mathbf{K})^\top \boldsymbol{\Gamma}_{\delta,*} \text{vec}(\mathbf{K}) - (\mathbf{g}_{\mathbf{K}} - \boldsymbol{\Gamma}_{\mathbf{K},\eta} \boldsymbol{\Gamma}_\eta^{-1} \mathbf{g}_\eta)^\top \text{vec}(\mathbf{K}) + \lambda_{\mathbf{K}} \|\mathbf{K}\|_1. \quad (3.22)$$

### 3.5.3 Univariate Examples

In this section we present univariate Gaussian models on general domains as an example of our generalized score matching estimator. In particular, we estimate one of  $\mu_0$  and  $\sigma_0^2$  assuming the other is known, given that the true density is

$$p_{\mu_0,\sigma_0^2}(x) \propto \exp \left\{ -\frac{(x - \mu_0)^2}{2\sigma_0^2} \right\}, \quad x \in \mathcal{D}$$

with  $\mu_0 \in \mathbb{R}$ ,  $\sigma_0^2 > 0$  and  $\mathcal{D} \subset \mathbb{R}$  *piecewise decomposable*, i.e. a union of at most countably many intervals and with Lebesgue positive measure.

If we do not impose  $\ell_1$  or  $\ell_2$  penalty and assume the true  $\sigma_0^2$  is known, given i.i.d. samples  $X^{(1)}, \dots, X^{(n)} \sim p_{\mu_0,\sigma_0^2}$ , similar to Example 3.1 in Section 2.3, we have the estimator for  $\mu_0$

$$\hat{\mu} \equiv \frac{\sum_{i=1}^n (h \circ \varphi_{\mathcal{C}})(X^{(i)}) \cdot X^{(i)} - \sigma_0^2 (h \circ \varphi_{\mathcal{C}})'(X^{(i)})}{\sum_{i=1}^n (h \circ \varphi_{\mathcal{C}})(X^{(i)})}.$$

By Theorem 18, it suffices to choose  $h(x) = x^\alpha$  with  $\alpha > 0$ . Similar to Section 2.3, we have

$$\sqrt{n}(\hat{\mu} - \mu_0) \rightarrow_d \mathcal{N} \left( 0, \frac{\mathbb{E}_0 [\sigma_0^2 (h \circ \varphi_{\mathcal{C}})^2(X) + \sigma_0^4 (h \circ \varphi_{\mathcal{C}})'^2(X)]}{\mathbb{E}_0^2 [(h \circ \varphi_{\mathcal{C}})(X)]} \right), \quad (3.23)$$

if the expectations exist. On the other hand, similar to Example 3.2 in Section 2.3, assuming the true  $\mu_0$  is known, the estimator for  $\sigma^2$  is,

$$\hat{\sigma}^2 \equiv \frac{\sum_{i=1}^n (h \circ \varphi_{\mathcal{C}})(X^{(i)}) \cdot (X^{(i)} - \mu_0)^2}{\sum_{i=1}^n (h \circ \varphi_{\mathcal{C}})(X^{(i)}) + (h \circ \varphi_{\mathcal{C}})'(X^{(i)}) \cdot (X^{(i)} - \mu_0)}$$

with limiting distribution

$$\sqrt{n} (\hat{\sigma}^2 - \sigma_0^2) \rightarrow_d \mathcal{N} \left( 0, \frac{\mathbb{E}_0 [2\sigma_0^6 (h \circ \varphi_C)^2(X) \cdot (X - \mu_0)^2 + \sigma_0^8 (h \circ \varphi_C)'{}^2(X) \cdot (X - \mu_0)^2]}{\mathbb{E}_0^2 [(h \circ \varphi_C)(X) \cdot (X - \mu_0)^2]} \right). \quad (3.24)$$

In Figure 3.3 we show an example of standard normal  $\mathcal{N}(0, 1)$  restricted to three univariate domains, namely  $\mathcal{D}_2 \equiv (-\infty, -3/2] \cup [3/2, +\infty)$ ,  $\mathcal{D}_3 \equiv [-1, -3/4] \cup [3/4, 1]$ , and their union  $\mathcal{D}_1 \equiv (-\infty, -3/2] \cup [-1, -3/4] \cup [3/4, 1] \cup [3/2, +\infty)$ . Note that  $\mathcal{D}_1$  is the union of  $\mathcal{D}_2$  and  $\mathcal{D}_3$ , where  $\mathcal{D}_2$  comprises two unbounded intervals, and  $\mathcal{D}_3$  consists of two bounded intervals. The endpoints are chosen so that the probability of the variable lying in each interval is roughly the same ( $\mathbb{P}(\mathcal{N}(0, 1) \in (-\infty, -3/2]) \approx 0.0668$ ,  $\mathbb{P}(\mathcal{N}(0, 1) \in [-1, -3/4]) \approx 0.0680$ ). To pick the truncation point  $C$  in the definition of the truncated distance  $\varphi_{C, \mathcal{D}}$ , we choose a  $\pi \in (0, 1]$ , and let  $C$  be the  $\pi$  quantile of  $\varphi_{+\infty, \mathcal{D}}(X)$  with  $X \sim \mathcal{N}(0, 1)$  on  $\mathcal{D}$ , that is, the  $C$  such  $\mathbb{P}(\varphi_{+\infty, \mathcal{D}}(X) \leq C) = \pi$ . Here,  $\varphi_{+\infty, \mathcal{D}_1}(X) = |X| - 3/2$  if  $|X| > 3/2$ , or  $\min(|X| - 3/4, 1 - |X|)$  otherwise,  $\varphi_{+\infty, \mathcal{D}_2}(X) = |X| - 3/2$ , and  $\varphi_{+\infty, \mathcal{D}_3}(X) = \min(|X| - 3/4, 1 - |X|)$ .

The first subfigure in each row of Figure 3.3 shows the density on each domain, along with the corresponding  $\varphi_{+\infty}(X)$  in red, whose  $y$  axis is on the right. The second in each row shows the log of the asymptotic variance for the corresponding  $\hat{\mu}$ , as on the right-hand side of (3.23), and the third shows that for  $\hat{\sigma}^2$  as in (3.24). Each curve represents a different power  $\alpha$  in  $h(x) = x^\alpha$ , and the  $x$  axis represents the quantiles  $\pi$  associated with the truncation point  $C$ , as mentioned above. Finally, the red dotted curve shows the truncation point  $C$  versus  $\pi$  for each domain. The ‘‘bumps’’ in the variance for  $x^{0.5}$  are due to numerical instability in integration.

As we show in Section 3.7, for the purpose of edge recovery for graphical models, we recommend using  $\mathbf{h}(\mathbf{x}) = (x_1^{\alpha_1}, \dots, x_m^{\alpha_m})$  with  $\alpha \geq 1$  for  $\mathcal{D}$  that is a finite disjoint union of convex subsets of  $\mathbb{R}^m$ . Although minimizing the asymptotic variance in the univariate case is a different task,  $\alpha = 1$  also seems to be consistently the best performing choice.

For  $\mathcal{D}_2$  and  $\mathcal{D}_3$ , all variance curves are U-shaped, while for  $\mathcal{D}_1 = \mathcal{D}_2 \cup \mathcal{D}_3$  each curve is visually piecewise connected by two parabolae, and the turning point corresponds to

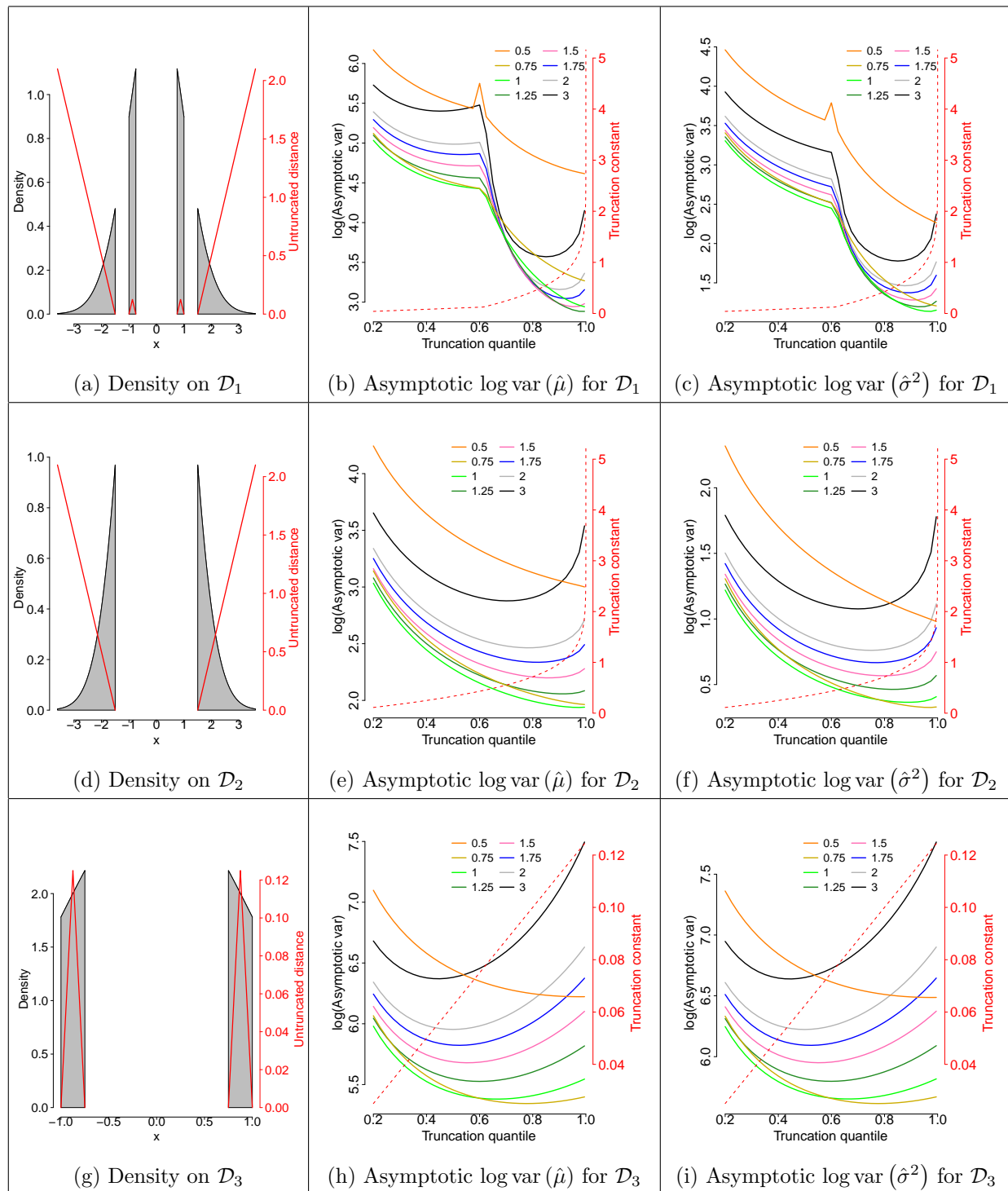


Figure 3.3: Univariate Gaussian example.

$C_0 = \max \varphi_{+\infty, \mathcal{D}_3}(x) = (1 - 3/4)/2 = 0.125$ . To the right of this  $C_0$ , the truncation is applied to the two unbounded intervals (i.e.  $\mathcal{D}_2$ ) only. The first segment of most  $\text{var}(\hat{\mu})$  curves for  $\mathcal{D}_1$ , as well as most curves for  $\mathcal{D}_2$  suggest there may still be benefit from truncating the distances  $\varphi$  within the bounded intervals, although the  $\text{var}(\hat{\sigma}^2)$  curves for  $\mathcal{D}_1$  as well as both curves for  $x^{0.75}$  on  $\mathcal{D}_2$  suggest otherwise. On the other hand, the curves for  $\mathcal{D}_1$  and  $\mathcal{D}_2$  imply that a truncation constant larger than  $C_0$  is favorable; by referring to the ticks on the right-hand side, note that the curves for  $\mathcal{D}_2$  reach their minimum at  $C \geq 0.5$ . This suggests the choice of a separate truncation point  $C$  for each connected component of  $\mathcal{D}$ , especially for unbounded versus bounded ones. Although a further generalization of our proposed method, this becomes infeasible for  $m \gg 1$  and involves tuning multiple parameters, which we do not further examine in this paper.

### 3.6 $a$ - $b$ Models on Standard Simplices

In this section we consider  $a$ - $b$  models on the most important example of  $\mathcal{D}$  with zero Lebesgue measure in  $\mathbb{R}^m$ , the  $(m-1)$ -dimensional standard simplex  $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}_+^m \mid \mathbf{x} \succ \mathbf{0}, \mathbf{1}_m^\top \mathbf{x} = 1\}$ . As in Example 3.13, in the density we profile out the last component of  $\mathbf{x}$  using  $x_m = 1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m}$ :

$$p_{\boldsymbol{\eta}, \mathbf{K}}(\mathbf{x}_{-m}) \propto \exp \left[ -\frac{1}{2a} \mathbf{x}_{-m}^a \top \mathbf{K}_{-m, -m} \mathbf{x}_{-m}^a - \frac{1}{a} \mathbf{x}_{-m}^a \top \boldsymbol{\kappa}_{-m, m} (1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m})^a - \frac{1}{2a} \kappa_{m, m} (1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m})^{2a} + \frac{1}{b} \boldsymbol{\eta}_{-m}^\top \mathbf{x}_{-m}^b + \frac{\eta_m}{b} (1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m})^b \right] \quad (3.25)$$

on  $\mathcal{D}_{-m} \equiv \{\mathbf{x}_{-m} \in \mathbb{R}_+^{m-1} \mid \mathbf{x}_{-m} \succ \mathbf{0}, \mathbf{1}_{m-1}^\top \mathbf{x}_{-m} < 1\} \subseteq \mathbb{R}_+^{m-1}$ , where for simplifying the notation of this density we again define that  $1/0 \equiv 1$  and  $\mathbf{x}_{-m}^0 \equiv \log(\mathbf{x}_{-m})$ .

Sufficient conditions on  $a$  and  $b$  for the density to be proper as discussed in Theorem 16 and Corollary 17 also hold in the simplex cases. Additional sufficient conditions for the case where  $a = b = 0$  are extensively studied in Section 3.6.2. The following theorem states that the sufficient conditions on  $h_j$  for a general  $\mathcal{D} \subseteq \mathbb{R}^m$  are also sufficient for  $\mathcal{D}_{-m} \subseteq \mathbb{R}^{m-1}$ , and the equivalent definition (3.10) of the generalized score matching loss and its empirical loss (3.11) are again valid.

**Theorem 19.** *Suppose the conditions in Theorem 18 hold, then (A.1), (A.2) and (A.3) in Lemma 13 are also satisfied for  $a$ - $b$  models on  $\mathcal{D}_{-m}$  in (3.25).*

By the nature of the simplex domain, sparsity pattern of  $\mathbf{K}$  does not directly relate to the conditional independence between two components, as any two components are perfectly correlated when all others are fixed, due to the sum-to-one condition. However, one may still be interested in the pattern as a measure of how two components interact, as well as estimation of  $\mathbf{K}$  and  $\boldsymbol{\eta}$  with minimized error for the sole purpose of recovering the density itself. As a prerequisite, the following theorem gives the conditions for the identifiability of  $\mathbf{K}$  and  $\boldsymbol{\eta}$  from a given  $a$ - $b$  density on the simplex. In short, the parameters are identifiable if  $2a = b > 0$  does not hold and  $a \neq 1$ , and in these cases the goal of recovering the underlying true parameters  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  is well-defined.

**Theorem 20.** *Suppose there exist  $\mathbf{K}_1, \mathbf{K}_2, \boldsymbol{\eta}_1, \boldsymbol{\eta}_2$  such that*

$$\exp\left(-\frac{1}{2a}\mathbf{x}^{a\top}\mathbf{K}_1\mathbf{x}^a + \frac{1}{b}\boldsymbol{\eta}_1^\top\mathbf{x}^b\right) = \exp\left(-\frac{1}{2a}\mathbf{x}^{a\top}\mathbf{K}_2\mathbf{x}^a + \frac{1}{b}\boldsymbol{\eta}_2^\top\mathbf{x}^b\right)$$

for all  $\mathbf{x} \in \mathcal{D}$ , where by  $\mathbf{x}^0$  we mean  $\log(\mathbf{x})$  and  $1/0 \equiv 1$ . Then  $\mathbf{K}_1 = \mathbf{K}_2$  and  $\boldsymbol{\eta}_1 = \boldsymbol{\eta}_2$ , or else one of the following must hold: (1)  $a = b = 1$ , (2)  $a = 1, b = 2$ , (3)  $a = 1$  and  $\boldsymbol{\eta}_1 = \boldsymbol{\eta}_2$ , (4)  $2a = b > 0$  and  $\mathbf{K}_1 - \mathbf{K}_2 = 2\boldsymbol{\eta}_1 - 2\boldsymbol{\eta}_2$ .

### 3.6.1 Estimation for General $a$ and $b$

In the equations in this section, for  $a = 0$  (and similarly for  $b = 0$ ) we substitute the coefficients “ $a$ ” with “1” and “ $(a - 1)$ ” with “ $-1$ ”, and let  $x^a \equiv \log x$ .

By substituting  $x_m \equiv 1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m}$  and working on  $\mathcal{D}_{-m}$  instead,  $\partial_j \log p(\mathbf{x}_{-m})$  now depends on both  $(\boldsymbol{\kappa}_{\cdot j}, \eta_j)$  and  $(\boldsymbol{\kappa}_{\cdot m}, \eta_m)$ . Thus, unlike in Section 3.5.2,  $(\boldsymbol{\kappa}_{\cdot j}, \eta_j)$  and  $(\boldsymbol{\kappa}_{\cdot m}, \eta_m)$  are no longer isolated in the score-matching loss. In particular,

$$\partial_j \log p(\mathbf{x}_{-m}) = -(\boldsymbol{\kappa}_{\cdot j}^\top \mathbf{x}^a) x_j^{a-1} + (\boldsymbol{\kappa}_{\cdot m}^\top \mathbf{x}^a) x_m^{a-1} + \eta_j x_j^{b-1} - \eta_m x_m^{b-1} \quad (3.26)$$

$$\partial_{jj} \log p(\mathbf{x}_{-m}) = -(a-1)(\boldsymbol{\kappa}_{\cdot j}^\top \mathbf{x}^a) x_j^{a-2} - (a-1)(\boldsymbol{\kappa}_{\cdot m}^\top \mathbf{x}^a) x_m^{a-2} - a\kappa_{jj} x_j^{2a-2} - a\kappa_{mm} x_m^{2a-2} \quad (3.27)$$

$$+ (b-1)\eta_j x_j^{b-2} + (b-1)\eta_m x_m^{b-2} + 2a\kappa_{jm} x_j^{a-1} x_m^{a-1}$$

and the unpenalized loss becomes

$$\frac{1}{2} (\text{vec}(\mathbf{K}), \boldsymbol{\eta})^\top \boldsymbol{\Gamma} (\text{vec}(\mathbf{K}), \boldsymbol{\eta}) - \mathbf{g}^\top (\text{vec}(\mathbf{K}), \boldsymbol{\eta}), \quad \text{where}$$

$$\boldsymbol{\Gamma} \equiv \begin{bmatrix} \boldsymbol{\Gamma}_{\mathbf{K}} & \boldsymbol{\Gamma}_{\mathbf{K},\boldsymbol{\eta}} \\ \boldsymbol{\Gamma}_{\mathbf{K},\boldsymbol{\eta}}^\top & \boldsymbol{\Gamma}_{\boldsymbol{\eta}} \end{bmatrix} \in \mathbb{R}^{m(m+1) \times m(m+1)}, \quad \boldsymbol{\Gamma}_{\mathbf{K}} \in \mathbb{R}^{m^2 \times m^2}, \quad \boldsymbol{\Gamma}_{\mathbf{K},\boldsymbol{\eta}} \in \mathbb{R}^{m^2 \times m}, \quad \boldsymbol{\Gamma}_{\boldsymbol{\eta}} \in \mathbb{R}^{m \times m}, \quad (3.28)$$

$$\mathbf{g} \equiv (\text{vec}(\mathbf{g}_{\mathbf{K}}), \mathbf{g}_{\boldsymbol{\eta}}) \in \mathbb{R}^{m(m+1)}, \quad \mathbf{g}_{\mathbf{K}} \in \mathbb{R}^{m \times m}, \quad \mathbf{g}_{\boldsymbol{\eta}} \in \mathbb{R}^m. \quad (3.29)$$

The components  $\boldsymbol{\Gamma}_{\mathbf{K}} \in \mathbb{R}^{m^2 \times m^2}$ ,  $\boldsymbol{\Gamma}_{\mathbf{K},\boldsymbol{\eta}} \in \mathbb{R}^{m(m+1)}$  and  $\boldsymbol{\Gamma}_{\boldsymbol{\eta}} \in \mathbb{R}^{m \times m}$  are no longer block-diagonal with  $m$  blocks as for general  $\mathcal{D} \subseteq \mathbb{R}^m$  in Section 3.4.2. Instead, each of them can be partitioned into  $m \times m$  blocks with  $m$  “block-columns” and  $m$  “block-rows”, and each matrix not only has  $m$  blocks on the diagonal, but also have the entire right-most “block-column” and bottom-most “block-row” filled, representing the interaction between the  $(\boldsymbol{\kappa}_j, \eta_j)$  and  $(\boldsymbol{\kappa}_m, \eta_m)$ . In particular,

$$\boldsymbol{\Gamma}_{\mathbf{K}} \equiv \begin{bmatrix} \boldsymbol{\Gamma}_{\mathbf{K},1} & \mathbf{0} & \cdots & \mathbf{0} & \boldsymbol{\Gamma}_{\mathbf{K},(1,m)} \\ \mathbf{0} & \boldsymbol{\Gamma}_{\mathbf{K},2} & \cdots & \mathbf{0} & \boldsymbol{\Gamma}_{\mathbf{K},(2,m)} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \boldsymbol{\Gamma}_{\mathbf{K},m-1} & \boldsymbol{\Gamma}_{\mathbf{K},(m-1,m)} \\ \boldsymbol{\Gamma}_{\mathbf{K},(1,m)}^\top & \boldsymbol{\Gamma}_{\mathbf{K},(2,m)}^\top & \cdots & \boldsymbol{\Gamma}_{\mathbf{K},(m-1,m)}^\top & \boldsymbol{\Gamma}_{\mathbf{K},m} \end{bmatrix} \in \mathbb{R}^{m^2 \times m^2}, \quad \text{each block} \in \mathbb{R}^{m \times m},$$

$$\boldsymbol{\Gamma}_{\mathbf{K},\boldsymbol{\eta}} \equiv \begin{bmatrix} \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},1} & \mathbf{0} & \cdots & \mathbf{0} & \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},(1,m)} \\ \mathbf{0} & \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},2} & \cdots & \mathbf{0} & \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},(2,m)} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},m-1} & \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},(m-1,m)} \\ \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},(m,1)} & \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},(m,2)} & \cdots & \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},(m,m-1)} & \boldsymbol{\gamma}_{\mathbf{K},\boldsymbol{\eta},m} \end{bmatrix} \in \mathbb{R}^{m^2 \times m}, \quad \text{each} \in \mathbb{R}^m,$$

$$\mathbf{\Gamma}_\eta \equiv \begin{bmatrix} \gamma_{\eta,1} & 0 & \cdots & 0 & \gamma_{\eta,(1,m)} \\ 0 & \gamma_{\eta,2} & \cdots & 0 & \gamma_{\eta,(2,m)} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & \gamma_{\eta,m-1} & \gamma_{\eta,(m-1,m)} \\ \gamma_{\eta,(1,m)} & \gamma_{\eta,(2,m)} & \cdots & \gamma_{\eta,(m-1,m)} & \gamma_{\eta,m} \end{bmatrix} \in \mathbb{R}^{m \times m}, \text{ each block } \in \mathbb{R},$$

where for  $j = 1, \dots, m-1$ ,

$$\begin{aligned} \mathbf{\Gamma}_j &\equiv \begin{bmatrix} \mathbf{\Gamma}_{\mathbf{K},j} & \gamma_{\mathbf{K},\eta,j} \\ \gamma_{\mathbf{K},\eta,j}^\top & \gamma_{\eta,j} \end{bmatrix} \\ &\equiv \frac{1}{n} \sum_{i=1}^n (h_j \circ \varphi_j) (\mathbf{X}^{(i)}) \begin{bmatrix} X_j^{(i)a-1} \mathbf{X}^{(i)a} \\ -X_j^{(i)b-1} \end{bmatrix} \begin{bmatrix} X_j^{(i)a-1} \mathbf{X}^{(i)a\top} & -X_j^{(i)b-1} \end{bmatrix}, \\ \mathbf{\Gamma}_m &\equiv \begin{bmatrix} \mathbf{\Gamma}_{\mathbf{K},m} & \gamma_{\mathbf{K},\eta,m} \\ \gamma_{\mathbf{K},\eta,m}^\top & \gamma_{\eta,m} \end{bmatrix} \\ &\equiv \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^{m-1} (h_k \circ \varphi_k) (\mathbf{X}^{(i)}) \begin{bmatrix} X_m^{(i)a-1} \mathbf{X}^{(i)a} \\ -X_m^{(i)b-1} \end{bmatrix} \begin{bmatrix} X_m^{(i)a-1} \mathbf{X}^{(i)a\top} & -X_m^{(i)b-1} \end{bmatrix}, \\ \mathbf{\Gamma}_{(j,m)} &\equiv \begin{bmatrix} \mathbf{\Gamma}_{\mathbf{K},(j,m)} & \gamma_{\mathbf{K},\eta,(j,m)} \\ \gamma_{\mathbf{K},\eta,(m,j)} & \gamma_{\eta,(j,m)} \end{bmatrix} \\ &\equiv -\frac{1}{n} \sum_{i=1}^n (h_j \circ \varphi_j) (\mathbf{X}^{(i)}) \begin{bmatrix} X_j^{(i)a-1} \mathbf{X}^{(i)a} \\ -X_j^{(i)b-1} \end{bmatrix} \begin{bmatrix} X_m^{(i)a-1} \mathbf{X}^{(i)a\top} & -X_m^{(i)b-1} \end{bmatrix}. \end{aligned}$$

In addition,

$$\begin{aligned} \mathbf{g}_{\mathbf{K},j} &\equiv \frac{1}{n} \sum_{i=1}^n \left[ \partial_j (h_j \circ \varphi_j) (\mathbf{X}^{(i)}) X_j^{(i)a-1} + (a-1)(h_j \circ \varphi_j) (\mathbf{X}^{(i)}) X_j^{(i)a-2} \right] \mathbf{X}^{(i)a} \\ &\quad + a(h_j \circ \varphi_j) (\mathbf{X}^{(i)}) X_j^{(i)2a-2} \mathbf{e}_{j,m} - a(h_j \circ \varphi_j) (\mathbf{X}^{(i)}) X_j^{(i)a-1} X_m^{(i)a-1} \mathbf{e}_{m,m} \\ \mathbf{g}_{\mathbf{K},m} &\equiv \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^{m-1} \left[ -\partial_k (h_k \circ \varphi_k) (\mathbf{X}^{(i)}) X_m^{(i)a-1} + (a-1)(h_k \circ \varphi_k) (\mathbf{X}^{(i)}) X_m^{(i)a-2} \right] \mathbf{X}^{(i)a} \\ &\quad + a(h_k \circ \varphi_k) (\mathbf{X}^{(i)}) X_m^{(i)2a-2} \mathbf{e}_{m,m} - a(h_k \circ \varphi_k) (\mathbf{X}^{(i)}) X_k^{(i)a-1} X_m^{(i)a-1} \mathbf{e}_{k,m}, \\ \mathbf{g}_{\eta,j} &\equiv \frac{1}{n} \sum_{i=1}^n -\partial_j (h_j \circ \varphi_j) (\mathbf{X}^{(i)}) X_j^{(i)b-1} - (b-1)(h_j \circ \varphi_j) (\mathbf{X}^{(i)}) X_j^{(i)b-2}, \end{aligned}$$

$$\mathbf{g}_{\boldsymbol{\eta}, m} \equiv \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^{m-1} \partial_k (h_k \circ \varphi_k) (\mathbf{X}^{(i)}) X_m^{(i)b-1} - (b-1)(h_k \circ \varphi_k) (\mathbf{X}^{(i)}) X_m^{(i)b-2}.$$

Instead of  $m$  diagonal blocks for a general  $\mathcal{D} \subseteq \mathbb{R}^m$  of positive measure, in the simplex case each matrix has  $(3m-2) = \mathcal{O}(m)$  blocks, so the computational complexity is the same as the general case.

We note that the profiled estimator in (3.22) for  $\lambda_{\boldsymbol{\eta}} = 0$  is not practical in this case: The profiled quadratic matrix  $\boldsymbol{\Gamma}_{\delta, \star} = \boldsymbol{\Gamma}_{\delta, \mathbf{K}} - \boldsymbol{\Gamma}_{\mathbf{K}, \boldsymbol{\eta}} \boldsymbol{\Gamma}_{\boldsymbol{\eta}}^{-1} \boldsymbol{\Gamma}_{\mathbf{K}, \boldsymbol{\eta}}^{\top}$  is no longer block-diagonal; in fact, all of its  $m \times m$  blocks can be nonzero due to the new structure of  $\boldsymbol{\Gamma}_{\delta, \mathbf{K}}$ ,  $\boldsymbol{\Gamma}_{\mathbf{K}, \boldsymbol{\eta}}$  and  $\boldsymbol{\Gamma}_{\boldsymbol{\eta}}$ . The inversion of  $\boldsymbol{\Gamma}_{\boldsymbol{\eta}}$  also increases the computational burden and introduces more numerical error. Instead of profiling out  $\boldsymbol{\eta}$  and solve for  $\mathbf{K}$  only, it is thus recommended to directly perform coordinate-descent on both  $\mathbf{K}$  and  $\boldsymbol{\eta}$  without any penalty on  $\boldsymbol{\eta}$ .

### 3.6.2 log-log Models on Standard Simplices

In this section we consider the special case of  $a = 0$  and  $b = 0$ , namely models with density proportional to  $\exp(-\log \mathbf{x}^{\top} \mathbf{K} \log \mathbf{x} / 2 + \boldsymbol{\eta}^{\top} \log \mathbf{x})$  supported on the  $m$ -dimensional standard simplex. Central to these models, consider the  $A^{m-1}$  class of distributions formulated in (2.7) of Aitchison (1985) with parameters  $\boldsymbol{\beta} \equiv (\beta_j)_{j=1, \dots, m}$  and  $(\gamma_{jk})_{1 \leq j \neq k \leq m}$ , where  $\gamma_{jk} = \gamma_{kj}$ , and density proportional to

$$\exp \left( -\frac{1}{2} \sum_{\substack{j \neq k \\ j, k=1, \dots, m}} \gamma_{jk} (\log x_j - \log x_k)^2 + (\boldsymbol{\beta} - \mathbf{1}_m)^{\top} \log \mathbf{x} \right) \mathbf{1}_{\mathcal{D}}(\mathbf{x}).$$

Expanding the expression this can be rewritten as

$$\exp \left( -\sum_{j=1}^m (\log x_j)^2 \left( \sum_{k \in \{1, \dots, m\} \setminus \{j\}} \gamma_{jk} \right) + \sum_{\substack{j \neq k \\ j, k=1, \dots, m}} \gamma_{jk} \log x_j \log x_k + (\boldsymbol{\beta} - \mathbf{1}_m)^{\top} \log \mathbf{x} \right)$$

and thus belong to our  $a$ - $b$  model

$$\exp \left( -\frac{1}{2} \log \mathbf{x}^{\top} \mathbf{K} \log \mathbf{x} + \boldsymbol{\eta}^{\top} \log \mathbf{x} \right)$$

with  $\boldsymbol{\eta} \equiv \boldsymbol{\beta} - \mathbf{1}_m$  as well as  $\kappa_{jj} = 2 \sum_{i \in \{1, \dots, m\} \setminus \{j\}} \gamma_{ji}$  and  $\kappa_{kj} = \kappa_{jk} = -2\gamma_{kj}$  for  $j \neq k = 1, \dots, m$ . The  $A^{m-1}$  model thus translates to the  $a$ - $b$  model with  $a = b = 0$  and  $\mathbf{K}\mathbf{1}_m = \mathbf{0}_m$  and  $\mathbf{K} = \mathbf{K}^\top$ .

In Theorem 21 we show that under simple conditions the density is again proper and the equivalent definition (3.10) of the generalized score matching loss is still valid. For estimation, we again consider  $\mathcal{D}_{-m} \equiv \{\mathbf{x}_{-m} \in \mathbb{R}_+^{m-1} \mid \mathbf{x}_{-m} \succ \mathbf{0}, \mathbf{1}_{m-1}^\top \mathbf{x}_{-m} < 1\} \subseteq \mathbb{R}_+^{m-1}$  and profile out the last coordinate  $x_m$  in the density.

**Theorem 21.** *Suppose  $\mathbf{K}$  is symmetric, and one of the following holds:*

- (1)  $\mathbf{K}$  is positive definite, or
- (2)  $\mathbf{K}\mathbf{1}_m = \mathbf{0}_m$ ,  $\mathbf{K}_{-k,-k}$  is positive definite for some  $k = 1, \dots, m$ , and  $\mathbf{1}_m^\top \boldsymbol{\eta} + m \geq 0$ , or
- (3)  $\mathbf{K}\mathbf{1}_m = \mathbf{0}_m$ ,  $\mathbf{K}$  is positive semi-definite, and  $\boldsymbol{\eta} \succ -\mathbf{1}_m$ .

*Then the density has a finite normalizing constant. Note that (2) implies that  $\mathbf{K}$  is positive semi-definite and (3) implies that  $\mathbf{K}_{-k,-k}$  is positive semi-definite for all  $k$  (but not necessarily positive definite).*

*For all  $j = 1, \dots, m-1$ , assume  $h_j(x) = x^{\alpha_j}$  with  $\alpha_j > 0$  and for (3) also assume  $\alpha_j > \max\{1 - \eta_{0,j}, 1 - \eta_{0,m}\}$ , then (A.1)–(A.3) in Lemma 13 are satisfied.*

### *Inference*

As discussed in the beginning of Section 3.6, conditional dependence between two components cannot be inferred from  $\mathbf{K}$  due to the nature of the simplex domain. However, one can still be interested in recovering the sparsity of  $\mathbf{K}$  to see how different components interact, and recovering  $\mathbf{K}$  and  $\boldsymbol{\eta}$  to estimate the density itself. For this purpose, we have the following corollary of Theorem 20, i.e.  $\mathbf{K}$  and  $\boldsymbol{\eta}$  are exactly identifiable from the density (assuming  $\mathbf{K}\mathbf{1}_m = \mathbf{0}_m$  or not):

**Corollary 22.** *Suppose there exist  $\mathbf{K}_1, \mathbf{K}_2, \boldsymbol{\eta}_1, \boldsymbol{\eta}_2$  such that*

$$\exp\left(-\frac{1}{2}\log(\mathbf{x})^\top \mathbf{K}_1 \log(\mathbf{x}) + \boldsymbol{\eta}_1^\top \log(\mathbf{x})\right) = \exp\left(-\frac{1}{2}\log(\mathbf{x})^\top \mathbf{K}_2 \log(\mathbf{x}) + \boldsymbol{\eta}_2^\top \log(\mathbf{x})\right)$$

for all  $\mathbf{x} \in \mathcal{D}$ , the standard simplex. Then  $\mathbf{K}_1 = \mathbf{K}_2$  and  $\boldsymbol{\eta}_1 = \boldsymbol{\eta}_2$ .

Now define the *additive log-ratio transformation*  $\mathbf{y}_{-m} \equiv \text{alt}(\mathbf{x}) \equiv \log \mathbf{x}_{-m} - (\log x_m) \mathbf{1}_{m-1} = (\log(x_1/x_m), \dots, \log(x_{m-1}/x_m))$ . The  $A^{m-1}$  model, corresponding to (2) and (3) in Theorem 21, is proposed in Aitchison (1985) as a generalization of both the *additive logistic normal model* with a normal density in  $\mathbf{y}_{-m}$  with inverse covariance  $\mathbf{K}_{-m,-m}$  and mean  $\mathbf{K}_{-m,-m}^{-1} \boldsymbol{\eta}_{-m}$ , corresponding to (2) with  $\mathbf{1}_m^\top \boldsymbol{\eta} = -m$ , and the Dirichlet distribution with parameters  $\boldsymbol{\eta} + \mathbf{1}_m$ , corresponding to (3) with  $\mathbf{K} = \mathbf{0}$ . This generalization requires only one additional parameter than the additive logistic normal model, namely  $\mathbf{1}_m^\top \boldsymbol{\eta}$  is no longer assumed to be  $-m$ . This brings in more flexibility to the two basic models over the simplex.

By nature of the simplex domains, any  $X_j$  and  $X_k$  are perfectly conditionally correlated given all other  $\mathbf{X}_{-j,-k}$ . On the other hand, under the additive logistic normal model,  $Y_j = \log(X_j/X_m)$  and  $Y_k = \log(X_k/X_m)$  are conditionally independent given all other  $\log(X_\ell/X_m)$ ,  $\ell \neq j, k, m$  if and only if  $\kappa_{jk} = \kappa_{kj} = 0$ ; conversely, as we make clear in the proof of Theorem 21, this is true only for the additive logistic normal model ( $\mathbf{1}_m^\top \boldsymbol{\eta} = -m$ ). We note that the statement in this paragraph still holds when we replace the ratios with respect to  $X_m$  by those w.r.t. any  $X_\ell$ .

Assuming the more general  $A^{m-1}$  model, we can thus perform the following two-step test of conditional independence between  $X_j/X_\ell$  and  $X_k/X_\ell$  for triplet  $(j, k, \ell)$  with  $j, k, \ell$  all different. First test if  $\mathbf{1}_m^\top \boldsymbol{\eta} + m$  is significantly different from 0 by comparing the BIC/eBIC for fitted  $A^{m-1}$  and additive logistic normal models with regularization discussed in Section 3.4.3; if we reject the null, then we cannot establish conditional independence for any such triplets. Otherwise, we claim that for all  $\ell \neq j, k$ ,  $X_j/X_\ell$  and  $X_k/X_\ell$  are conditionally independent given all other  $X_i/X_\ell$  ( $i \neq j, k, \ell$ ) if and only if  $\kappa_{jk} = \kappa_{kj} = 0$  in either fitted model. While the additive logistic normal model can be fitted as a Gaussian to the additive

log-ratio transformed data, the  $A^{m-1}$  model on the simplex can also be easily fitted with the help of generalized score matching, as discussed below.

### *Estimation for $A^{m-1}$ Models*

In Section 3.6.1 we formulate the elements for estimating  $\mathbf{K}$  and  $\boldsymbol{\eta}$  under the log-log model ( $a = b = 0$ ) assuming (1)  $\mathbf{K}$  is positive definite only. But for the  $A^{m-1}$  models (Aitchison, 1985) which assume the additional constraint  $\mathbf{K}\mathbf{1}_m = \mathbf{0}_m$  and  $\mathbf{K} = \mathbf{K}^\top$  as in (2) and (3) of Theorem 21, we need the following modification. We marginalize out the diagonals of  $\mathbf{K}$  with  $\kappa_{jj} = -\boldsymbol{\kappa}_{-j,j}^\top \mathbf{1}_{m-1}$  and estimate all off-diagonal elements. Under the assumption, for any matrices  $\mathbf{A}$  with  $m$  rows and  $\mathbf{B}$  with  $m$  columns, we can write

$$\begin{aligned}\boldsymbol{\kappa}_{\cdot j}^\top \mathbf{A} &= \boldsymbol{\kappa}_{-j,j}^\top \mathbf{A}_{-j} + \kappa_{jj} \mathbf{a}_j^\top = \boldsymbol{\kappa}_{-j,j}^\top (\mathbf{A}_{-j} - \mathbf{1}_{m-1} \mathbf{a}_j^\top) = \boldsymbol{\kappa}_{-j,j}^\top (\mathbf{C}^j \mathbf{A}), \\ \mathbf{B} \boldsymbol{\kappa}_{\cdot j} &= \mathbf{B}_{\cdot, -j} \boldsymbol{\kappa}_{-j,j} + \mathbf{b}_{\cdot j} \kappa_{jj} = (\mathbf{B}_{\cdot, -j} - \mathbf{b}_{\cdot j} \mathbf{1}_{m-1}^\top) \boldsymbol{\kappa}_{-j,j} = (\mathbf{B} \mathbf{C}^j)^\top \boldsymbol{\kappa}_{-j,j},\end{aligned}$$

where  $\mathbf{C}^j \in \mathbb{R}^{(m-1) \times m}$  has its  $j$ -th column all equal to  $-1$ , and  $(1, 1), \dots, (j-1, j-1), (j, j+1), \dots, (m-1, m)$ -th positions equal to  $1$ , and  $0$  everywhere else. Let  $\mathbf{C} \in \mathbb{R}^{(m-1)m \times m^2}$  be the block-diagonal matrix with blocks  $\mathbf{C}^1, \dots, \mathbf{C}^m$ . The unpenalized generalized score-matching loss  $\frac{1}{2} \text{vec}([\mathbf{K}^\top \boldsymbol{\eta}])^\top \boldsymbol{\Gamma} \text{vec}([\mathbf{K}^\top \boldsymbol{\eta}]) - \mathbf{g}^\top \text{vec}([\mathbf{K}^\top \boldsymbol{\eta}])$  thus becomes

$$\frac{1}{2} \text{vec} \left( \begin{bmatrix} \mathbf{K}_{\text{off}} \\ \boldsymbol{\eta}^\top \end{bmatrix} \right)^\top \begin{bmatrix} \mathbf{C} \boldsymbol{\Gamma}_{\mathbf{K}} \mathbf{C}^\top & \mathbf{C} \boldsymbol{\Gamma}_{\mathbf{K}, \boldsymbol{\eta}} \\ \boldsymbol{\Gamma}_{\mathbf{K}, \boldsymbol{\eta}}^\top \mathbf{C}^\top & \boldsymbol{\Gamma}_{\boldsymbol{\eta}} \end{bmatrix} \text{vec} \left( \begin{bmatrix} \mathbf{K}_{\text{off}} \\ \boldsymbol{\eta}^\top \end{bmatrix} \right) - \begin{bmatrix} \mathbf{C} \text{vec}(\mathbf{g}_{\mathbf{K}}) \\ \mathbf{g}_{\boldsymbol{\eta}} \end{bmatrix}^\top \text{vec} \left( \begin{bmatrix} \mathbf{K}_{\text{off}} \\ \boldsymbol{\eta}^\top \end{bmatrix} \right). \quad (3.30)$$

Denote the new components as  $\tilde{\boldsymbol{\Gamma}} \in \mathbb{R}^{m^2 \times m^2}$  and  $\tilde{\mathbf{g}} \in \mathbb{R}^{m^2}$ . For any  $m^2$ -dimensional nonzero  $\mathbf{v} \equiv (\mathbf{v}_{\mathbf{K},1}, \dots, \mathbf{v}_{\mathbf{K},m}, \mathbf{v}_{\boldsymbol{\eta}}) \in \mathbb{R}^{m^2}$  with  $\mathbf{v}_{\mathbf{K},j} \in \mathbb{R}^{m-1}$  and  $\mathbf{v}_{\boldsymbol{\eta}} \in \mathbb{R}^m$ , form  $\tilde{\mathbf{V}} \in \mathbb{R}^{m \times m}$  whose  $j$ -th column has  $\tilde{\mathbf{v}}_{-j,j} = \mathbf{v}_{\mathbf{K},j}$  and  $\tilde{v}_{jj} = 1 - \mathbf{1}_{m-1}^\top \mathbf{v}_{\mathbf{K},j}$  (analogous to the inverse operation of  $\mathbf{K} \mapsto \text{vec}(\mathbf{K}_{\text{off}})$ ), and let  $\tilde{\mathbf{v}} \equiv (\text{vec}(\tilde{\mathbf{V}}), \mathbf{v}_{\boldsymbol{\eta}})$ . Then by definition  $\mathbf{v}^\top \tilde{\boldsymbol{\Gamma}} \mathbf{v} = \tilde{\mathbf{v}}^\top \boldsymbol{\Gamma} \tilde{\mathbf{v}}$ . Since  $\tilde{\mathbf{v}} \neq \mathbf{0}_{m(m+1)}$  if and only if  $\mathbf{v} \neq \mathbf{0}_{m^2}$ ,  $\tilde{\boldsymbol{\Gamma}}$  is positive (semi-)definite if and only if  $\boldsymbol{\Gamma}$  is positive (semi-)definite.

Since  $\tilde{\boldsymbol{\Gamma}}$  is merely a linear transformation of  $\boldsymbol{\Gamma}$ , we may apply diagonal multipliers directly to  $\tilde{\boldsymbol{\Gamma}}$  for simplicity, and the concentration properties of the original  $\boldsymbol{\Gamma}$  (high probability bounds

for deviation of  $\mathbf{\Gamma}$  from its expectation) will be inherited, subject to different constants. The penalized generalized score-matching loss for the  $A^{m-1}$  models is thus defined as

$$\begin{aligned} & \hat{J}_{\mathbf{h}, \mathcal{C}, \mathcal{D}, \lambda, \delta}(p_{\theta}) \\ & \equiv \frac{1}{2} \text{vec} \left( \begin{bmatrix} \mathbf{K}_{\text{off}} \\ \boldsymbol{\eta}^{\top} \end{bmatrix} \right)^{\top} \tilde{\mathbf{\Gamma}}_{\delta}(\mathbf{x}) \text{vec} \left( \begin{bmatrix} \mathbf{K}_{\text{off}} \\ \boldsymbol{\eta}^{\top} \end{bmatrix} \right) - \tilde{\mathbf{g}}(\mathbf{x})^{\top} \text{vec} \left( \begin{bmatrix} \mathbf{K}_{\text{off}} \\ \boldsymbol{\eta}^{\top} \end{bmatrix} \right) + \lambda_{\mathbf{K}} \|\mathbf{K}_{\text{off}}\|_1 + \lambda_{\boldsymbol{\eta}} \|\boldsymbol{\eta}\|_1, \end{aligned} \quad (3.31)$$

$$\tilde{\mathbf{\Gamma}}_{\delta} \equiv \begin{bmatrix} (\mathbf{C}\mathbf{\Gamma}_{\mathbf{K}}\mathbf{C}^{\top})_{\delta} & \mathbf{C}\mathbf{\Gamma}_{\mathbf{K}, \boldsymbol{\eta}} \\ \mathbf{\Gamma}_{\mathbf{K}, \boldsymbol{\eta}}^{\top} \mathbf{C}^{\top} & \mathbf{\Gamma}_{\boldsymbol{\eta}} \end{bmatrix}, \quad \tilde{\mathbf{g}} \equiv \begin{bmatrix} \mathbf{C} \text{vec}(\mathbf{g}_{\mathbf{K}}) \\ \mathbf{g}_{\boldsymbol{\eta}} \end{bmatrix},$$

where  $(\cdot)_{\delta}$  denotes the operation of multiplying the diagonals of a matrix by  $\delta > 1$ .

### 3.7 Theoretical Properties

In this section we present theoretical guarantees for our generalized method applied to the pairwise interaction power  $a$ - $b$  models. In particular, similar to Section 2.6, with high probability we bound the deviation of our estimates  $\hat{\mathbf{K}}$  and  $\hat{\boldsymbol{\eta}}$  from their true values  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$ . We first restate Definition 4 as below.

**Definition 9.** Let  $\mathbf{\Gamma}_0 \equiv \mathbb{E}_0 \mathbf{\Gamma}(\mathbf{x})$  and  $\mathbf{g}_0 \equiv \mathbb{E}_0 \mathbf{g}(\mathbf{x})$  be the population versions of  $\mathbf{\Gamma}(\mathbf{x})$  and  $\mathbf{g}(\mathbf{x})$  under the distribution given by a true parameter matrix  $\boldsymbol{\Psi}_0 \equiv [\mathbf{K}_0, \boldsymbol{\eta}_0]^{\top} \in \mathbb{R}^{m(m+1)}$ , or  $\boldsymbol{\Psi}_0 \equiv \mathbf{K}_0 \in \mathbb{R}^{m^2}$  in the centered case if we assume  $\boldsymbol{\eta}_0 \equiv \mathbf{0}$ . The support of a matrix  $\boldsymbol{\Psi}$  is  $S(\boldsymbol{\Psi}) \equiv \{(i, j) : \psi_{ij} \neq 0\}$ , and we let  $S_0 = S(\boldsymbol{\Psi}_0)$ . For a matrix  $\boldsymbol{\Psi}_0$ , we define  $d_{\boldsymbol{\Psi}_0}$  to be the maximum number of non-zero entries in any column, and  $c_{\boldsymbol{\Psi}_0} \equiv \|\boldsymbol{\Psi}_0\|_{\infty, \infty}$ . Writing  $\mathbf{\Gamma}_{0, AB}$  for the  $A \times B$  submatrix of  $\mathbf{\Gamma}_0$ , we define

$$c_{\mathbf{\Gamma}_0} \equiv \|\mathbf{\Gamma}_{0, S_0 S_0}^{-1}\|_{\infty, \infty}. \quad (3.32)$$

Finally,  $\mathbf{\Gamma}_0$  satisfies the irrepresentability condition with incoherence parameter  $\omega \in (0, 1]$  and edge set  $S_0$  if

$$\|\mathbf{\Gamma}_{0, S_0^c S_0} (\mathbf{\Gamma}_{0, S_0 S_0})^{-1}\|_{\infty, \infty} \leq (1 - \omega). \quad (3.33)$$

### 3.7.1 Truncated Gaussian Graphical Models on A Finite Disjoint Union of Convex Sets with Positive Measure

Truncated Gaussian graphical models are covered by our  $a$ - $b$  models described in Section 3.4.2 with  $a = b = 1$ . When the domain  $\mathcal{D}$  is a finite disjoint union of convex sets with a positive Lebesgue measure, we have the following theorem similar to Theorem 10 in Chapter 2, which bounds the errors as long as one uses finite truncation points  $\mathbf{C}$  for  $\varphi_{\mathbf{C},\mathcal{D}}$  and each component in  $\mathbf{h}(\mathbf{x})$  is a power function with a positive exponent.

**Theorem 23.** *Let  $\mathcal{D} \subset \mathbb{R}^m$  be a componentwise decomposable set (Def 5) with positive Lebesgue measure, and assume it is a finite disjoint union of convex sets  $\Delta \equiv \{\mathcal{D}_1, \dots, \mathcal{D}_{|\Delta|}\}$ , i.e.  $\mathcal{D} \equiv \mathcal{D}_1 \sqcup \dots \sqcup \mathcal{D}_{|\Delta|}$ . Suppose the data matrix holds  $n$  i.i.d. copies of  $\mathbf{X}$  following a truncated Gaussian distribution on  $\mathcal{D}$  with inverse covariance parameter  $\mathbf{K}_0 \in \mathbb{R}^{m \times m}$  and mean parameter  $\boldsymbol{\mu}_0$ , namely with density*

$$p_{\boldsymbol{\eta}_0, \mathbf{K}_0}(\mathbf{x}) \propto \exp\left(-\frac{1}{2}\mathbf{x}^\top \mathbf{K}_0 \mathbf{x} + \boldsymbol{\eta}_0^\top \mathbf{x}\right) \mathbb{1}_{\mathcal{D}}(\mathbf{x})$$

with  $\mathbf{K}_0$  positive definite and  $\boldsymbol{\eta}_0 \equiv \mathbf{K}_0 \boldsymbol{\mu}_0$ . Let  $\boldsymbol{\Psi}_0 \equiv [\mathbf{K}_0, \boldsymbol{\eta}_0]^\top$ .

Assume that (A.1)–(A.3) in Lemma 13 hold, and in addition that  $\mathbf{h}$  and the truncation points  $\mathbf{C}$  in the truncated componentwise distance  $\varphi_{\mathbf{C},\mathcal{D}}$  satisfy  $0 \leq (h_j \circ \varphi_{C_j, \mathcal{D}, j})(\mathbf{x}) \leq M$  and  $0 \leq \partial_j (h_j \circ \varphi_{C_j, \mathcal{D}, j})(\mathbf{x}) \leq M'$  almost surely for all  $j = 1, \dots, m$  for some constants  $0 < M, M' < +\infty$ . Note that  $\mathbf{h}(\mathbf{x}) = (x_1^{\alpha_1}, \dots, x_m^{\alpha_m})$  with  $\alpha_1, \dots, \alpha_m \geq 1$  satisfies all these assumptions, according to Theorem 18.

Let the diagonal multiplier  $\delta$  introduced in Section 3.4.3 satisfy

$$1 < \delta < C(n, m) \equiv 2 - \left(1 + 4e \max\left\{(6 \log m + 2 \log |\Delta|) / n, \sqrt{(6 \log m + 2 \log |\Delta|) / n}\right\}\right)^{-1}$$

and suppose further that  $\boldsymbol{\Gamma}_{0, S_0 S_0}$  is invertible and satisfies the irrepresentability condition (2.24) with  $\omega \in (0, 1]$ . Define  $c_{\mathbf{X}} \equiv 2 \max_{\mathcal{D}' \in \Delta} \max_j \left| 2\sqrt{(\mathbf{K}_0^{-1})_{jj}} + \sqrt{e} \mathbb{E}_0 X_j \mathbb{1}_{\mathcal{D}'}(\mathbf{X}) \right|$ . Suppose for  $\tau > 3$  the sample size and the regularization parameter satisfy

$$n > \mathcal{O}\left((\tau \log m + \log |\Delta|) \max\left\{\frac{M^2 c_{\boldsymbol{\Gamma}_0}^2 c_{\mathbf{X}}^4 d_{\boldsymbol{\Psi}_0}^2}{\omega^2}, \frac{M c_{\boldsymbol{\Gamma}_0} c_{\mathbf{X}}^2 d_{\boldsymbol{\Psi}_0}}{\omega}\right\}\right), \quad (3.34)$$

$$\lambda > \mathcal{O} \left[ (Mc_{\Psi_0}c_{\mathbf{X}}^2 + M'c_{\mathbf{X}} + M) \left( \sqrt{\frac{\tau \log m + \log |\Delta|}{n}} + \frac{\tau \log m + \log |\Delta|}{n} \right) \right]. \quad (3.35)$$

Then the following statements hold with probability  $1 - m^{3-\tau}$ :

(a) The regularized generalized  $\mathbf{h}$ -score matching estimator  $\hat{\Psi}$  that minimizes (3.18) is unique, has its support included in the true support,  $\hat{S} \equiv S(\hat{\Psi}) \subseteq S_0$ , and satisfies

$$\begin{aligned} \|\hat{\mathbf{K}} - \mathbf{K}_0\|_{\infty} &\leq \frac{c_{\Gamma_0}}{2-\omega} \lambda, & \|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\|_{\infty} &\leq \frac{c_{\Gamma_0}}{2-\omega} \lambda, \\ \|\hat{\mathbf{K}} - \mathbf{K}_0\|_F &\leq \frac{c_{\Gamma_0}}{2-\omega} \lambda \sqrt{|S_0|}, & \|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\|_F &\leq \frac{c_{\Gamma_0}}{2-\omega} \lambda \sqrt{|S_0|}, \\ \|\hat{\mathbf{K}} - \mathbf{K}_0\|_2 &\leq \frac{c_{\Gamma_0}}{2-\omega} \lambda \min \left( \sqrt{|S_0|}, d_{\Psi_0} \right), \\ \|\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}_0\|_2 &\leq \frac{c_{\Gamma_0}}{2-\omega} \lambda \min \left( \sqrt{|S_0|}, d_{\Psi_0} \right). \end{aligned}$$

(b) Moreover, if

$$\min_{j,k:(j,k) \in S_0} |\kappa_{0,jk}| > \frac{c_{\Gamma_0}}{2-\omega} \lambda \quad \text{and} \quad \min_{j:(m+1,j) \in S_0} |\eta_{0,j}| > \frac{c_{\Gamma_0}}{2-\omega} \lambda,$$

then  $\hat{S} = S_0$  and  $\text{sign}(\hat{\kappa}_{jk}) = \text{sign}(\kappa_{0,jk})$  for all  $(j, k) \in S_0$  and  $\text{sign}(\hat{\eta}_j) = \text{sign}(\eta_{0j})$  for  $(m+1, j) \in S_0$ .

In the centered setting, the same bounds hold by removing the dependencies on  $\hat{\boldsymbol{\eta}}$  and  $\boldsymbol{\eta}_0$ .

The detailed proof of the theorem is omitted, since it is analogous to that for Theorem 10 of Chapter 2 with two modifications: (i) using the triangle inequality, split the concentration bounds (A.6) (A.10) (A.11) in the Appendix for Chapter 2 into one for each  $\mathcal{D}_1, \dots, \mathcal{D}_{|\Delta|}$  and use a union bound; (ii) in the proof of Lemma 36.1, replace  $\mathcal{D} \equiv \mathbb{R}_+^m$  by any convex  $\mathcal{D}' = \mathcal{D}_1, \dots, \mathcal{D}_{|\Delta|}$  and replace  $X_1$  by  $X_1 \mathbb{1}_{\mathcal{D}'}(\mathbf{X})$ , as the proof there only uses the convexity of the domain.

### 3.7.2 Bounded Domains in $\mathbb{R}_+^m$ with Positive Measure

In this section we present results for general  $a$ - $b$  models on bounded domains with positive measure.

**Theorem 24.** (1) Suppose  $a > 0$  and  $b \geq 0$ . Let  $\mathcal{D}$  be a bounded subset of  $\mathbb{R}_+^m$  with positive Lebesgue measure with  $\mathcal{D} \subseteq [u_1, v_1] \times \cdots \times [u_m, v_m]$  for finite nonnegative constants  $u_1, v_1, \dots, u_m, v_m$ , and suppose that the true parameters  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  satisfy the conditions in Corollary 17 (so that the density is proper). Assume  $\mathbf{h}(\mathbf{x}) \equiv (x_1^{\alpha_1}, \dots, x_m^{\alpha_m})$  with  $\alpha_1, \dots, \alpha_m \geq \max\{1, 2-a, 2-b\}$ , and suppose  $\varphi_{\mathbf{C}, \mathcal{D}}$  has truncation points  $\mathbf{C} = (C_1, \dots, C_m)$  with  $0 < C_j < +\infty$  for  $j = 1, \dots, m$ . Define

$$\zeta(C_j, u_j, v_j, \alpha_j, p_j) \equiv \begin{cases} \min\{C_j, (v_j - u_j)/2\}^{\alpha_j} (u_j + v_j)^{p_j}/2^{p_j}, & p_j < 0, v_j - u_j \leq 2C_j, \\ \min\{C_j, (v_j - u_j)/2\}^{\alpha_j} (u_j + C_j)^{p_j}, & p_j < 0, v_j - u_j > 2C_j, \\ \min\{C_j, (v_j - u_j)/2\}^{\alpha_j} v_j^{p_j}, & p_j \geq 0, \end{cases}$$

$$\varsigma_{\mathbf{r}} \equiv \max_{j,k=1,\dots,m} \max\{\zeta(C_j, u_j, v_j, \alpha_j, 2a-2)v_k^{2a}, \zeta(C_j, u_j, v_j, \alpha_j, 2b-2)\},$$

$$\varsigma_{\mathbf{g}} \equiv \max_{j,k=1,\dots,m} \max\{\alpha_j \zeta(C_j, u_j, v_j, \alpha_j - 1, a-1)v_k^a$$

$$+ |a-1|\zeta(C_j, u_j, v_j, \alpha_j, a-2)v_k^a + a\zeta(C_j, u_j, v_j, \alpha_j, 2a-2),$$

$$\alpha_j \zeta(C_j, u_j, v_j, \alpha_j - 1, b-1) + |b-1|\zeta(C_j, u_j, v_j, \alpha_j, b-2)\}.$$

Suppose that  $\mathbf{\Gamma}_{0, S_0 S_0}$  is invertible and satisfies the irrepresentability condition (3.33) with  $\omega \in (0, 1]$ . Suppose for  $\tau > 0$  the sample size, the regularization parameter and the diagonal multiplier  $\delta$  from Section 3.4.3 satisfy

$$n > 72c_{\mathbf{r}}^2 d_{\Psi_0}^2 \varsigma_{\mathbf{r}}^2 (\tau \log m + \log 2)/\omega^2, \quad (3.36)$$

$$\lambda > \frac{3(2-\omega)}{\omega} \max\left\{c_{\Psi_0} \varsigma_{\mathbf{r}} \sqrt{2(\tau \log m + \log 4)/n}, \varsigma_{\mathbf{g}} \sqrt{(\tau \log m + \log 4)/(2n)}\right\}, \quad (3.37)$$

$$1 < \delta < C_{\text{bounded}}(n, m, \tau) \equiv 1 + \sqrt{(\tau \log m + \log 4)/(2n)}. \quad (3.38)$$

Then the statements (a) and (b) in Theorem 23 hold with probability at least  $1 - m^{-\tau}$ .

(2) For  $a = 0$  and  $b \geq 0$ , if  $u_j > 0$  for all  $k = 1, \dots, m$ , the above holds with

$$\varsigma_{\mathbf{r}} \equiv \max_{j,k=1,\dots,m} \max\{\zeta(C_j, u_j, v_j, \alpha_j, -2) \max\{|\log u_k|, |\log v_k|\}^2, \zeta(C_j, u_j, v_j, \alpha_j, 2b-2)\},$$

$$\varsigma_{\mathbf{g}} \equiv \max_{j,k=1,\dots,m} \max\{\alpha_j \zeta(C_j, u_j, v_j, \alpha_j - 1, -1) \max\{|\log u_k|, |\log v_k|\}$$

$$\begin{aligned}
& + |a - 1| \zeta(C_j, u_j, v_j, \alpha_j, -2) \max\{|\log u_k|, |\log v_k|\}, \\
& \alpha_j \zeta(C_j, u_j, v_j, \alpha_j - 1, b - 1) + |b - 1| \zeta(C_j, u_j, v_j, \alpha_j, b - 2)\}.
\end{aligned}$$

We note that the requirement on  $\alpha_j \geq 1$  is only for bounding the two  $\partial_j(h_j \circ \varphi_j)$  terms in  $\mathbf{g}(\mathbf{x})$  and might not be necessary in practice as we see in the simulation studies.

### 3.7.3 Unbounded Domains in $\mathbb{R}_+^m$ with Positive Measure

For unbounded domains  $\mathcal{D}$  in the non-negative orthant, we also have consistency results, albeit introducing an additional unknown constant factor that may depend on  $m$  to the sample complexity. For simplicity we only show the results for  $a > 0$  and those for  $a = 0$  should be similar. The following lemma enables us to bound each row of the data matrix  $\mathbf{x}$  by a finite cube with high probability and then apply the same proof as for Theorem 24.

**Lemma 25.** *Suppose  $\mathcal{D}$  has positive measure and the true parameters  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  satisfy the conditions in Corollary 17. Then for all  $j = 1, \dots, m$ ,  $X_j^{2a}$  is sub-exponential if  $a > 0$  and  $\log X_j$  is sub-exponential if  $a = 0$ .*

We have the following corollary of Theorem 24. The results involve an unknown constant, namely the sub-exponential norm  $\|\cdot\|_{\psi_1}$  of  $X_j^{2a}$ , and we have roughly and asymptotically  $n = \Omega(\log m) \max_j \mathcal{O}\left(\|X_j^{2a}\|_{\psi_1}\right)^{(\alpha_j + \max\{4a, 2b\} - 2)/a}$ . We suspect the sub-exponential norm to scale like  $\Omega((\log m)^c)$  for some  $c$  small, and the exact dependency on  $m$  is left for further research.

**Corollary 26.** *Suppose  $a > 0$  and  $\mathcal{D}$  is a subset of  $\mathbb{R}_+^m$  with positive measure and suppose that the true parameters  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  satisfy the conditions in Corollary 17. Let  $\rho_j(\mathcal{D}) \equiv \overline{\{x_j : \mathbf{x} \in \mathcal{D}\}}$  and  $\rho_{\mathcal{D}}^* \equiv \{j = 1, \dots, m : \sup \rho_j(\mathcal{D}) < +\infty\}$ , and suppose  $\rho_j(\mathcal{D}) \subseteq [u_j, v_j]$  for  $j \in \rho_{\mathcal{D}}^*$ . Then Theorem 24 holds with  $\log 4$  replaced by  $\log 6$  in (3.36)–(3.38), and  $u_j = \max\{\mathbb{E}_0 X_j^{2a} - \epsilon_{3,j}, 0\}^{1/(2a)}$  and  $v_j = (\mathbb{E}_0 X_j^{2a} + \epsilon_{3,j})^{1/(2a)}$  for  $j \notin \rho_{\mathcal{D}}^*$ , where*

$$\epsilon_{3,j} \equiv \max \left\{ 2\sqrt{2}e \|X_j^{2a}\|_{\psi_1} \sqrt{\log 3 + \log n + \tau \log m + \log(m - |\rho_{\mathcal{D}}^*|)}, \right.$$

$$4e \left\| X_j^{2a} \right\|_{\psi_1} (\log 3 + \log n + \tau \log m + \log (m - |\rho_{\mathcal{D}}^*|)) \Big\},$$

$$\left\| X_j^{2a} \right\|_{\psi_1} \equiv \sup_{q \geq 1} (\mathbb{E}_0 |X_j|^{2aq})^{1/q} / q \geq \mathbb{E}_0 X_j^{2a}.$$

### 3.7.4 Models on the Standard Simplices

For models with  $a > 0$  on the simplex, since each coordinate is in  $[0, 1]$ , we have the following corollary of Theorem 24.

**Corollary 27.** *Suppose  $a > 0$  and  $b \geq 0$ , and let  $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}_+^m \mid \mathbf{x} \succ \mathbf{0}, \mathbf{1}_m^\top \mathbf{x} = 1\}$  be the simplex domain in  $\mathbb{R}_+^m$ . Suppose that the true parameters  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  satisfy the conditions in Corollary 17. Then for the minimizer of (3.18) using the  $\boldsymbol{\Gamma}$  and  $\mathbf{g}$  formulated in Section 3.6.1, Theorem 24 holds with  $\varsigma_{\boldsymbol{\Gamma}}$  and  $\varsigma_{\mathbf{g}}$  replaced by  $\varsigma_{\boldsymbol{\Gamma}} \equiv 1$  and  $\varsigma_{\mathbf{g}} \equiv \max_{j=1, \dots, m} \alpha_j + \max\{|a - 1| + 2a, |b - 1|\}$ .*

In the proof of Corollary 27 we show that we can give tighter constant bounds  $\varsigma_{\boldsymbol{\Gamma}}$  and  $\varsigma_{\mathbf{g}}$  for entries in  $\boldsymbol{\Gamma}$  and  $\mathbf{g}$ , respectively, which may be a lot smaller but have rather complicated forms.

For models with  $a = 0$  on simplex domains, including the  $A^{m-1}$  models discussed in Section 3.6.2 we first use the following lemma to bound  $\log X_j$  with high probability, just as in Lemma 25.

**Lemma 28.** *Suppose  $\mathcal{D}$  is the standard simplex, and the true parameters  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  satisfy the conditions in Corollary 17 for  $a > 0$  or  $b > 0$ , or in Theorem 21 for  $a = b = 0$ . Then for all  $j = 1, \dots, m$ ,  $X_j^{2a}$  is sub-exponential for  $a > 0$ , and  $\log X_j$  is sub-exponential for  $a = 0$ .*

We then have the following corollary of Theorem 24.

**Corollary 29.** *Suppose  $\mathcal{D}$  is the simplex domain in  $\mathbb{R}_+^m$  and  $a = 0$ . Suppose  $b = 0$  and the conditions for  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  in Theorem 21 hold, or  $b > 0$  and the condition in Corollary 17 hold. Let  $\mathbf{h}(\mathbf{x}) \equiv (x_1^{\alpha_1}, \dots, x_m^{\alpha_m})$  with  $\alpha_1, \dots, \alpha_m \geq 2$ . Then Theorem 24 holds with  $\log 4$  replaced by  $\log 6$  in (3.36)–(3.38), and*

$$\varsigma_{\boldsymbol{\Gamma}} \equiv \max \left\{ 1, c_{\log, \mathbf{K}_0, \boldsymbol{\eta}_0}^2 \right\},$$

$$\begin{aligned} \varsigma_{\mathbf{g}} &\equiv \max \left\{ \left( \max_{j=1,\dots,m} \alpha_j + 1 \right) c_{\log, \mathbf{K}_0, \boldsymbol{\eta}_0} + 2, \max_j \alpha_j + |b - 1| \right\}, \quad \text{where} \\ c_{\log, \mathbf{K}_0, \boldsymbol{\eta}_0} &\equiv \max_j \mathbb{E}_0 \log X_j \\ &\quad + \max \left\{ 2\sqrt{2}e \max_{j=1,\dots,m} \|\log X_j\|_{\psi_1} \sqrt{\log 3 + \log n + (\tau + 1) \log m}, \right. \\ &\quad \left. 4e \max_{j=1,\dots,m} \|\log X_j\|_{\psi_1} (\log 3 + \log n + (\tau + 1) \log m) \right\}, \\ \|\log X_j\|_{\psi_1} &\equiv \sup_{q \geq 1} (\mathbb{E}_0 |\log X_j|^q)^{1/q} / q \geq -\mathbb{E}_0 \log X_j. \end{aligned}$$

The results are written in terms of the max of the sub-exponential norms of  $\log X_1, \dots, \log X_m$ , an unknown constant, and we have roughly  $n = \Omega(\log m) \mathcal{O}\left(\max_j \|\log X_j\|_{\psi_1}\right)^2$ . We expect the sub-exponential norm to scale like  $\Omega((\log m)^c)$  for some  $c$  small, and the exact dependency on  $m$  is left for further research.

### 3.8 Numerical Experiments

In this section we present results from numerical experiments using our method formulated in Section 3.3.2 and 3.3.5, as well as our extension to Liu and Kanamori (2019) presented in Section 3.3.6. The plots we report are in a similar fashion to those in Section 2.7.

#### 3.8.1 Estimation — Choice of $\mathbf{h}$ and $\mathbf{C}$

Recall that multiplication of  $\nabla \log p(\mathbf{x})$  with an  $(\mathbf{h} \circ \boldsymbol{\varphi}_{\mathbf{C}, \mathcal{D}})^{1/2}(\mathbf{x})$  is key to our method, in which the  $j$ -th component of  $\boldsymbol{\varphi}_{\mathbf{C}, \mathcal{D}}(\mathbf{x}) = (\varphi_{C_1, \mathcal{D}, 1}(\mathbf{x}), \dots, \varphi_{C_m, \mathcal{D}, m}(\mathbf{x}))$  is the distance of  $x_j$  to the boundary of its domain holding  $\mathbf{x}_{-j}$  fixed, with this distance truncated from above by some constant  $C_j > 0$ .

We use a uniform  $h$  for all components ( $\mathbf{h}(\mathbf{x}) = (h(x_1), \dots, h(x_m))$ ) and compare the performance of our method using various  $h(x)$  of the form  $x^c$  for some power  $c \geq 0$  along with various truncation points  $\mathbf{C}$ . In particular, we choose  $c = i/4$  for  $i = 0, 1, \dots, 8$ . Instead of prespecifying constants  $\mathbf{C}$ , given  $0 < \pi \leq 1$  we choose each  $C_j$  to be the  $\pi$  sample quantiles of  $\varphi_{+\infty, \mathcal{D}, j}$  applied to each row of the data matrix  $\mathbf{x}$ , namely the  $\pi$  quantile of

$\{\varphi_{+\infty, \mathcal{D}, j}(\mathbf{x}^{(1)}), \dots, \varphi_{+\infty, \mathcal{D}, j}(\mathbf{x}^{(n)})\}$ , assuming there are  $n$  samples in the data; infinite  $\varphi_{+\infty, \mathcal{D}, j}$  values are ignored. If all untruncated distances for the  $j$ -th component in the sample are infinite, i.e.  $\varphi_{+\infty, \mathcal{D}, j}(\mathbf{x}^{(1)}) = \dots = \varphi_{+\infty, \mathcal{D}, j}(\mathbf{x}^{(n)}) = +\infty$ , we set  $\varphi_j \equiv 1$ . This dynamic way of choosing the truncation points allows us to automatically adapt to the scale of data, and is more informative than fixing the constant to a grid from 0.5 to 3 as done in Section 2.7. In our experiments, we choose  $\pi = 0.2, 0.4, 0.6, 0.8, 1$ , where  $\pi = 1$  means no truncation for all finite  $\varphi_j$  values.

We remind that with power  $c = 0$ ,  $(\mathbf{h} \circ \varphi_{\mathbf{C}, \mathcal{D}})(\mathbf{x}) \equiv 1$  and it corresponds to the original score-matching for  $\mathbb{R}^m$  of Hyvärinen (2005), while with  $c = 2$ ,  $\mathbf{C} = +\infty^m$  and  $\mathcal{D} \equiv \mathbb{R}_+^m$ ,  $(\mathbf{h} \circ \varphi_{\mathbf{C}, \mathcal{D}})(\mathbf{x}) \equiv \mathbf{x}^2$  corresponds to the estimator of Hyvärinen (2007) and Lin et al. (2016). The case where  $\mathcal{D} \equiv \mathbb{R}_+^m$  corresponds to Section 2.7.

Finally, we also include results for our extension to the method proposed in Liu and Kanamori (2019) using  $g_0(\mathbf{x})$  as opposed to  $(\mathbf{h} \circ \varphi_{\mathbf{C}, \mathcal{D}})^{1/2}(\mathbf{x})$ , taking the  $\ell_2$  distance to the boundary  $\partial\mathcal{D}$  upper truncated by a constant  $C$ ; see Section 3.3.6. The constant  $C$  in this case is also determined using quantiles of the untruncated  $\ell_2$  distances of the given data sample to  $\partial\mathcal{D}$ . For  $C = +\infty$  ( $\pi = 1$ ) there is no truncation and the estimator corresponds to Liu and Kanamori (2019).

### 3.8.2 Numerical Experiments for Domains with Positive Measure

In this section we present results for general  $a$ - $b$  models restricted to domains with positive Lebesgue measure.

#### *Experiment Setup*

Throughout our experiments we choose dimension  $m = 100$  and sample sizes  $n = 80$  and 1000. For simplicity and given the length of this dissertation, we only present results for the centered case (assuming  $\boldsymbol{\eta} \equiv \mathbf{0}$ ) where the  $b$  power does not come into play in the distribution, i.e. the density is proportional to  $\exp\{-\mathbf{x}^a \mathbf{K} \mathbf{x}^a / (2a)\}$  for  $a > 0$  or  $\exp(-\log \mathbf{x}^\top \mathbf{K} \log \mathbf{x} / 2)$  for  $a = 0$ . Supported by the exhaustive experiments of the same kind carried out in Section

2.7, we expect that the results for non-centered settings are similar, meaning that the best choice of  $h$  mainly depends on  $a$  but not  $b$ .

We consider six settings for  $a$ : (1)  $a = 0$  (log), (2)  $a = 1/2$  (exponential square root; Inouye et al. (2016)), (3)  $a = 1$  (Gaussian), (4)  $a = 3/2$  as well as some extreme cases (5)  $a = 2$  and (6)  $a = 3$ . For all settings, we consider the following subsets of  $\mathbb{R}_+^m$  as our domain  $\mathcal{D}$ :

- i) non-negative  $\ell_2$  ball  $\{\mathbf{x} \in \mathbb{R}_+^m : \|\mathbf{x}\|_2 \leq c_1\}$ , which we call  $\ell_2$ -nn (“non-negative”),
- ii) complement of  $\ell_2$  ball in  $\mathbb{R}_+^m$ :  $\{\mathbf{x} \in \mathbb{R}_+^m : \|\mathbf{x}\|_2 \geq c_1\}$ , which we call  $\ell_2^c$ -nn, and
- iii)  $[c_1, +\infty)^m$ , which we call *unif-nn*,

for some  $c_1 > 0$ . For the Gaussian ( $a = 1$ ) case consider in addition the following subsets of  $\mathbb{R}^m$ :

- iv) the entire  $\ell_2$  ball  $\{\mathbf{x} \in \mathbb{R}^m : \|\mathbf{x}\|_2 \leq c_1\}$ , which we call  $\ell_2$ ,
- v) the complement of  $\ell_2$  ball in  $\mathbb{R}^m$ :  $\{\mathbf{x} \in \mathbb{R}^m : \|\mathbf{x}\|_2 \geq c_1\}$ , which we call  $\ell_2^c$ , and
- vi)  $((-\infty, c_1] \cup [c_1, +\infty))^m$ , which we call *unif*.

The constant  $c_1$  in each setting above is determined in the following way. We first generate  $n$  samples from the corresponding untruncated distribution on  $\mathbb{R}_+^m$  for i)–iii) or  $\mathbb{R}^m$  for iv)–vi), then determine the  $c_1$  so that exactly half of the samples would fall inside the truncated boundary. We believe that the  $c_1$  chosen adaptively in this way would be more meaningful and the truncation would not be too aggressive, similar to data one might see in practice.

The underlying true interaction matrices  $\mathbf{K}_0$  are selected as in Lin et al. (2016) and Section 2.7. Each graph has 10 disconnected subgraphs, each containing  $m/10 = 10$  nodes, making  $\mathbf{K}_0$  block-diagonal. In each block, each lower-triangular element is set to 0 with probability  $1 - \rho$  for some  $\rho \in (0, 1)$ , and is otherwise drawn from  $\text{Uniform}[0.5, 1]$ . The upper

triangular elements are determined by symmetry. The diagonal elements of  $\mathbf{K}_0$  are chosen as a common positive value such that the minimum eigenvalue of  $\mathbf{K}_0$  is 0.1. We generate 5 different true precision matrices  $\mathbf{K}_0$ , and run 10 trials using each of them. We choose  $\rho = 0.8$  for  $n = 1000$  and  $\rho = 0.2$  for  $n = 80$ . This way  $n/(d_{\mathbf{K}_0}^2 \log m)$  is roughly constant, coinciding with our theory in Section 3.7.

### Results

We use the area under the ROC curve (AUC) as the measure of performance of edge recovery of  $\mathbf{K}_0$ . Specifically, writing the estimate of  $\mathbf{K}_0$  as  $\hat{\mathbf{K}}$ , the ROC curve plots the true positive rate (TPR) against the false positive rate (FPR), with

$$\text{FPR} \equiv \frac{|\hat{S}_{\text{off}} \setminus S_{0,\text{off}}|}{m(m-1) - |S_{0,\text{off}}|} \quad \text{and} \quad \text{TPR} \equiv \frac{|\hat{S}_{\text{off}} \cap S_{0,\text{off}}|}{|S_{0,\text{off}}|},$$

where  $S_{0,\text{off}} \equiv \{(i, j) : i \neq j \wedge \kappa_{0,i,j} \neq 0\}$ , and  $\hat{S}_{\text{off}} \equiv \{(i, j) : i \neq j \wedge \hat{\kappa}_{i,j} \neq 0\}$ . We plot the AUC averaged over 50 trials for each setting against  $\pi = 0.2, 0.4, 0.6, 0.8, 1$ , whose sample quantiles are used as the truncation points  $\mathbf{C}$  of  $\varphi_{\mathbf{C},\mathcal{D}}$  (c.f. Section 3.8.1), while each curve represents  $h(x) = 1$  (Hyvärinen, 2005),  $g_0(\mathbf{x})$  from Section 3.3.6, or  $h(x) = x^c$  with  $c = 1/4, 1/2, \dots, 2$ . The  $y$ -ticks on the right-hand side are the original AUC values, whereas those on the left are the AUCs divided by the AUC for  $h(x) = 1$ , measuring the relative performance of each method compared to the original score matching in Hyvärinen (2005). (Note that  $h(x) = 1$  is invariant to the truncation points and is thus constant in each plot.)

From the plots we conclude that in most settings our method using  $h(x) = x^c$  with a  $c$  close to  $\max\{2 - a, 0\}$  works the best, as we observed in Section 2.7, and in most settings the truncated  $g_0$  function does not work well (Liu and Kanamori (2019) corresponds to  $\pi = 1$ ). The only notable exception is the domains iv)–vi) above, namely Gaussian models on subsets of  $\mathbb{R}^m$  not restricted to  $\mathbb{R}_+^m$ , shown in Figure 3.7. In particular, the original score matching in Hyvärinen (2005) designed for densities in the entire  $\mathbb{R}^m$  seems to work the best in these settings, suggesting that estimation of GGMs for domains with such truncations might not be challenging enough for one to switch to the more complex generalized methods. The

reason behind this is left for future research. On the other hand, by reading the  $y$  ticks on the left side of the plots in Figure 3.7, one should however note that for the iv)  $\ell_2$  and v)  $\ell_2^{\mathcal{C}}$  domains, the difference in the performance of all estimators is insignificant.

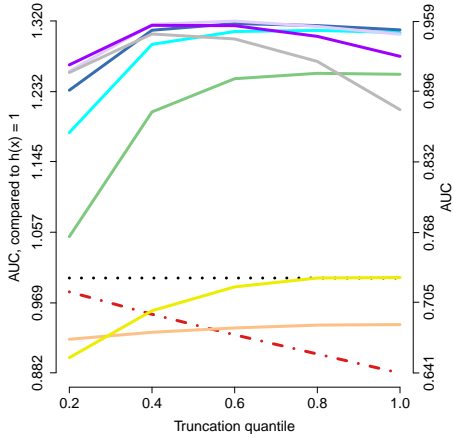
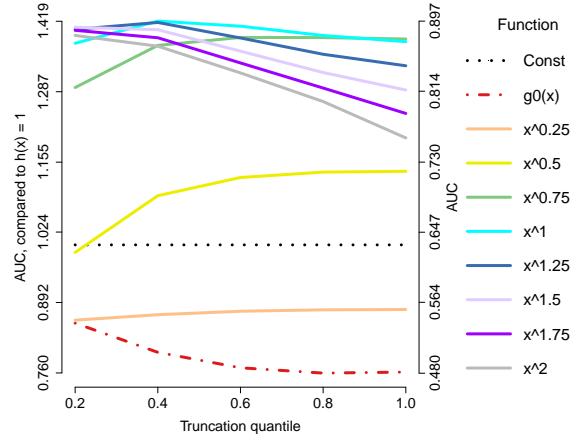
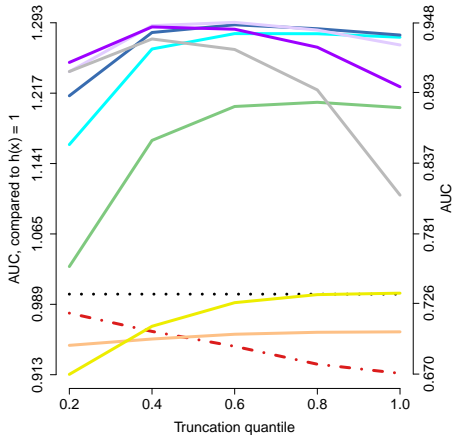
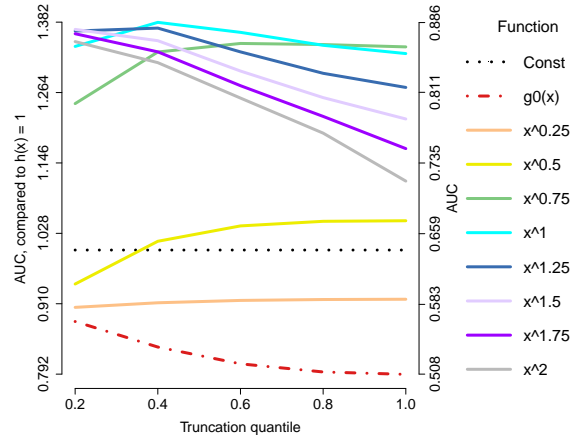
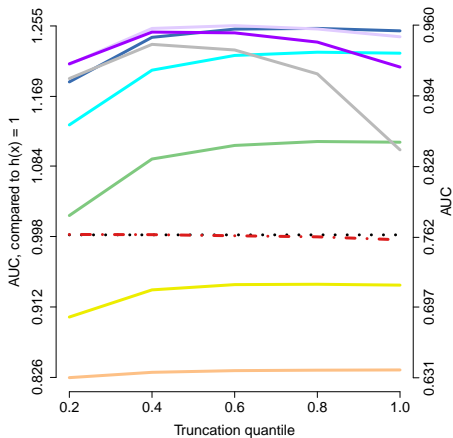
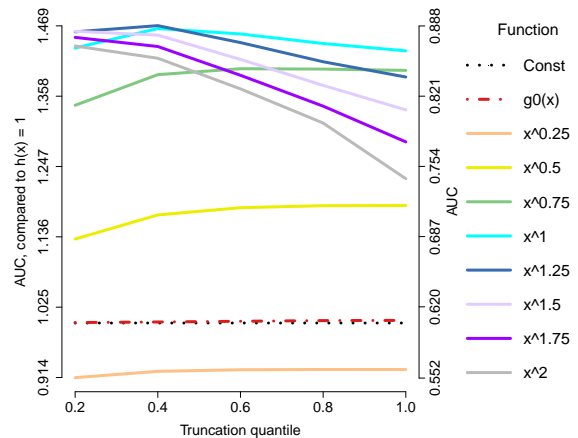
(a)  $n = 80$ ,  $\ell_2$ -nn domain(b)  $n = 1000$ ,  $\ell_2$ -nn domain(c)  $n = 80$ ,  $\ell_2^c$ -nn domain(d)  $n = 1000$ ,  $\ell_2^c$ -nn domain(e)  $n = 80$ , unif-nn domain(f)  $n = 1000$ , unif-nn domain

Figure 3.4: AUCs averaged over 50 trials for edge recovery using generalized score matching for the log models ( $a = 0$ ). Each curve represents using either our extension to  $g_0(\mathbf{x})$  proposed in Liu and Kanamori (2019) or a different choice of power function  $h(x) = x^c$ , and the  $x$  axis marks the  $\pi$  values whose columnwise sample quantiles are used as the truncation points  $\mathbf{C}$  for the truncated componentwise distances. The colors are sorted by the power  $c$ .

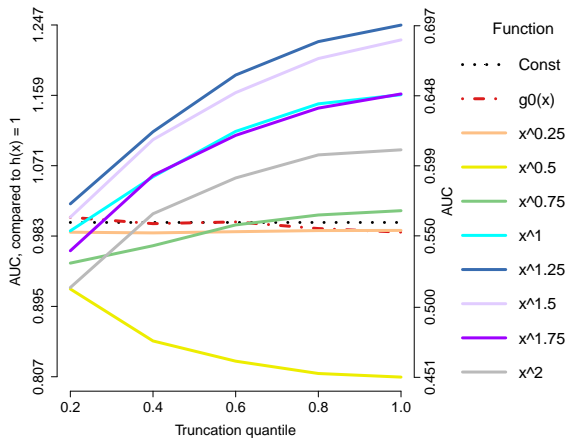
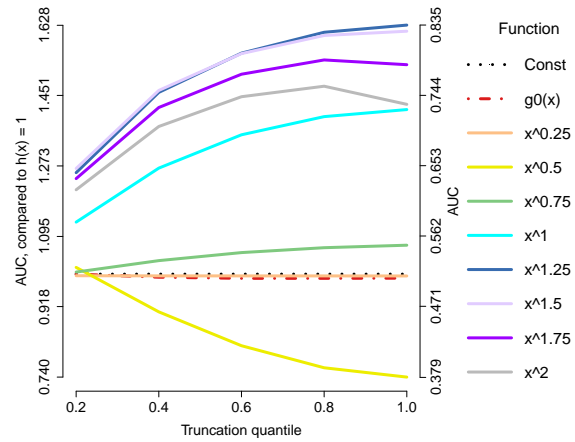
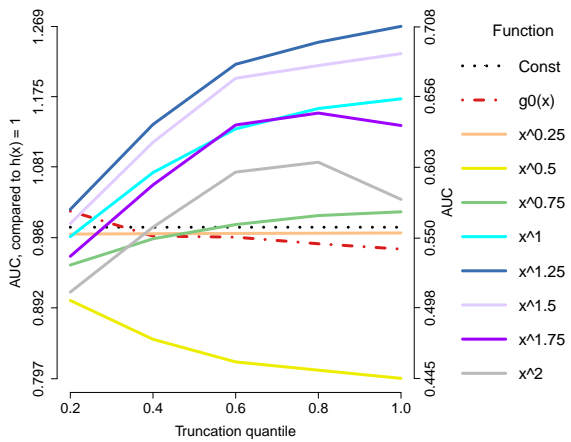
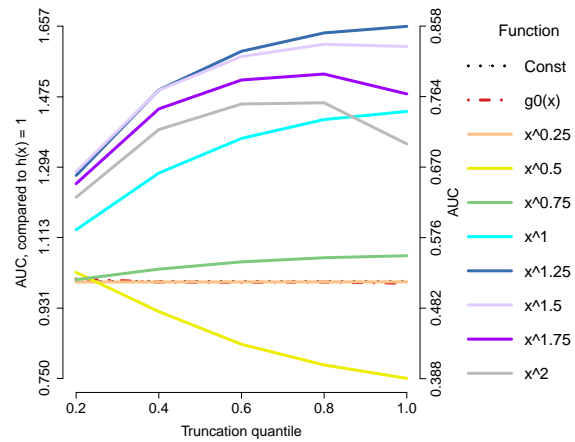
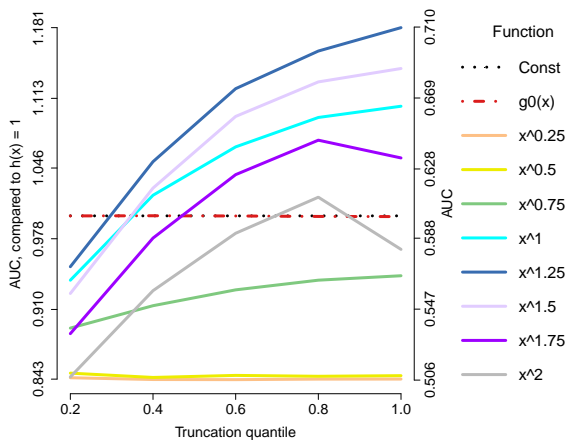
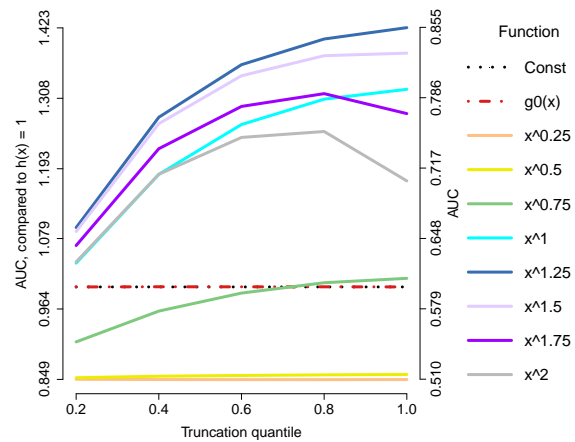
(a)  $n = 80$ ,  $\ell_2$ -nn domain(b)  $n = 1000$ ,  $\ell_2$ -nn domain(c)  $n = 80$ ,  $\ell_2^c$ -nn domain(d)  $n = 1000$ ,  $\ell_2^c$ -nn domain(e)  $n = 80$ , unif-nn domain(f)  $n = 1000$ , unif-nn domain

Figure 3.5: AUCs averaged over 50 trials for edge recovery using generalized score matching for the exponential square-root models ( $a = 1/2$ ).

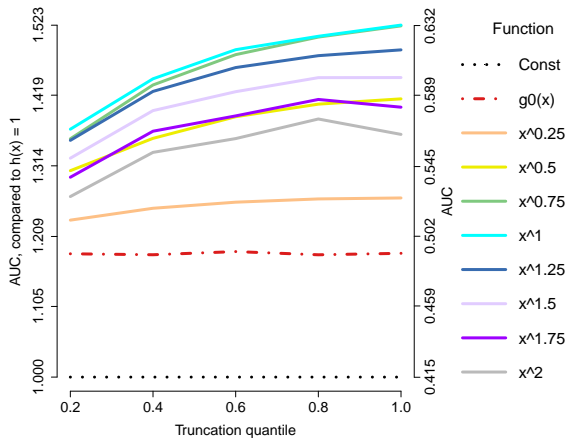
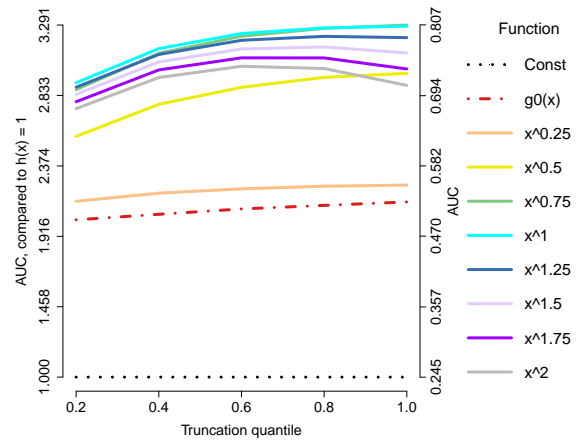
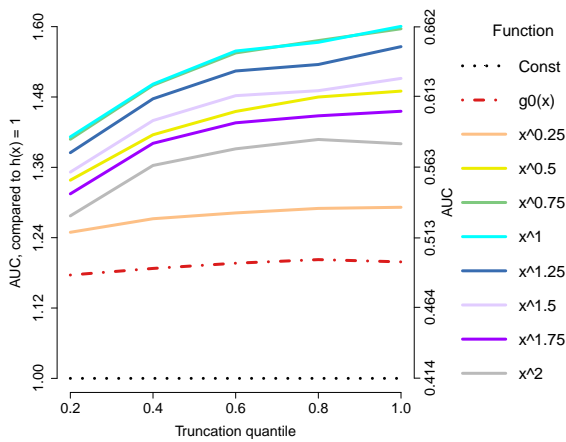
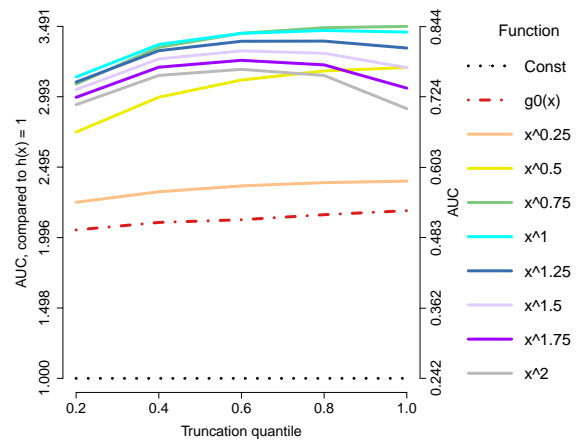
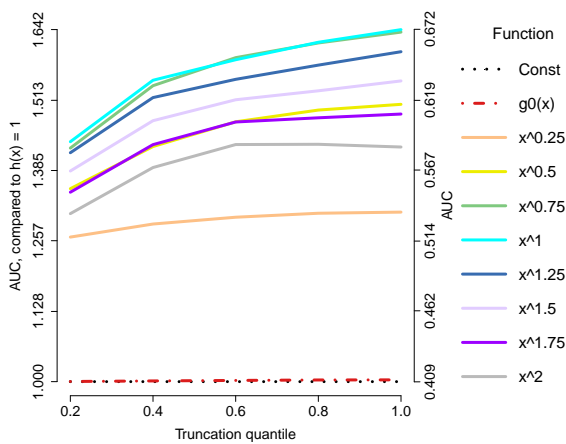
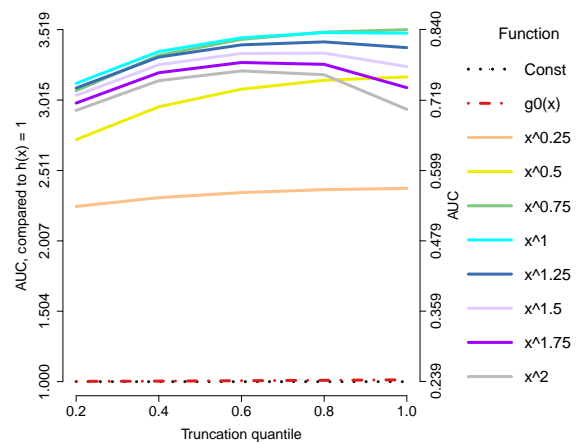
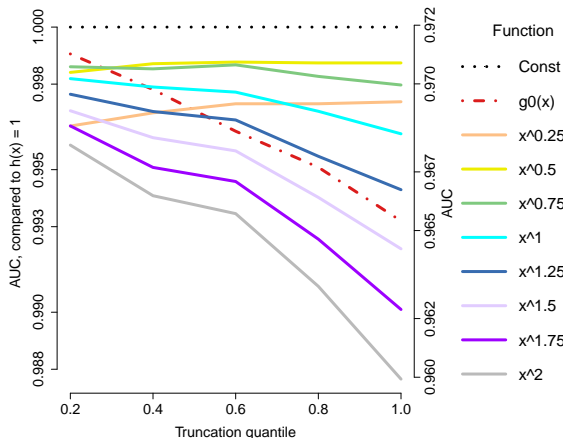
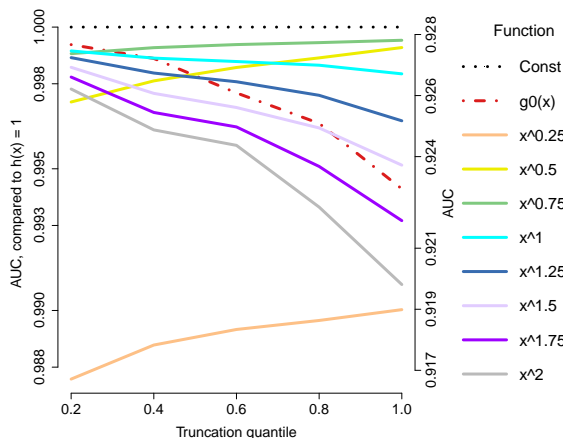
(a)  $n = 80$ ,  $\ell_2$ -nn domain(b)  $n = 1000$ ,  $\ell_2$ -nn domain(c)  $n = 80$ ,  $\ell_2^c$ -nn domain(d)  $n = 1000$ ,  $\ell_2^c$ -nn domain(e)  $n = 80$ , unif-nn domain(f)  $n = 1000$ , unif-nn domain

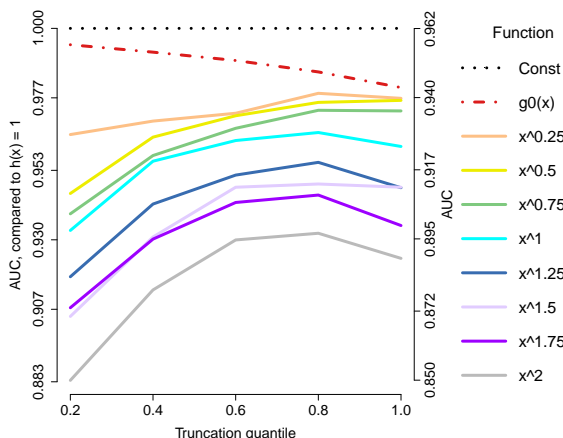
Figure 3.6: AUCs averaged over 50 trials for edge recovery using generalized score matching for the Gaussian models ( $a = 1$ ) on domains being subsets of  $\mathbb{R}_+^m$ .



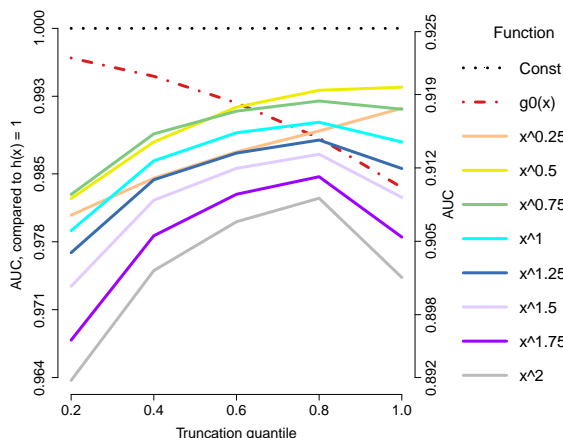
(a)  $n = 80, \ell_2$  domain



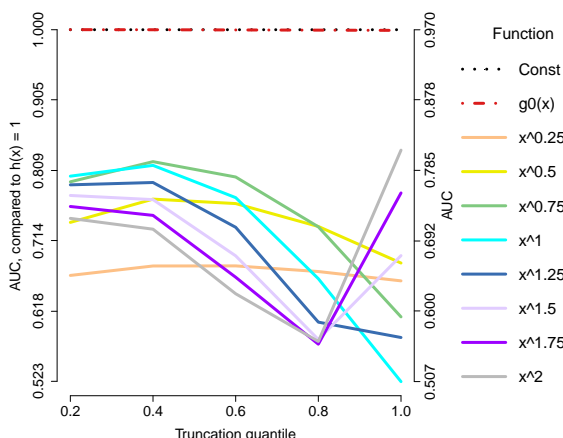
(b)  $n = 1000, \ell_2$  domain



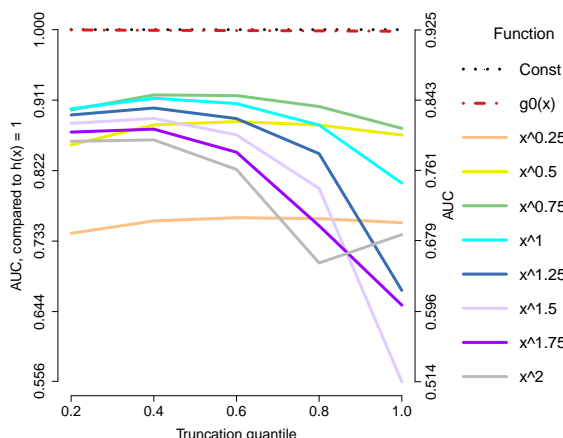
(c)  $n = 80, \ell_2^G$  domain



(d)  $n = 1000, \ell_2^G$  domain



(e)  $n = 80, \text{unif domain}$



(f)  $n = 1000, \text{unif domain}$

Figure 3.7: AUCs averaged over 50 trials for edge recovery using generalized score matching for the Gaussian models ( $a = 1$ ) on domains being subsets of  $\mathbb{R}^m$  (not restricted to  $\mathbb{R}_+^m$ ).

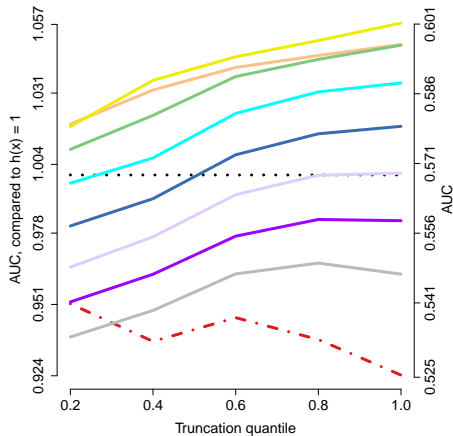
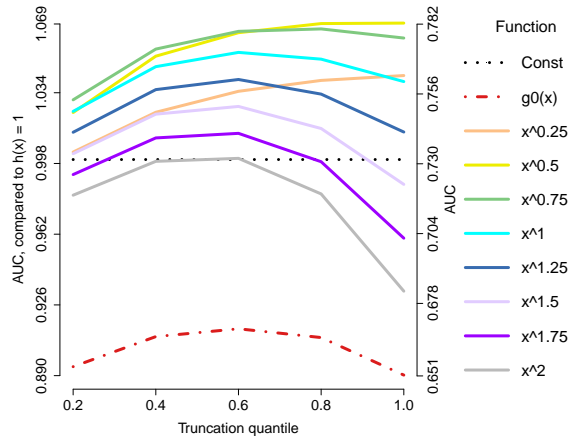
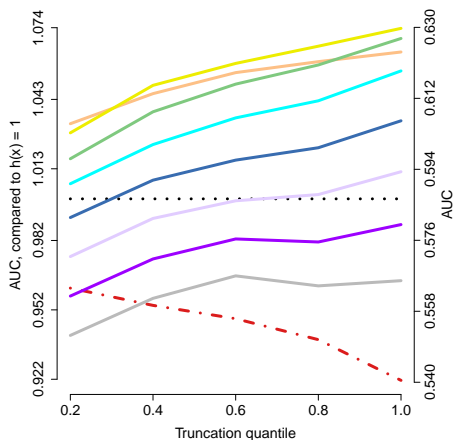
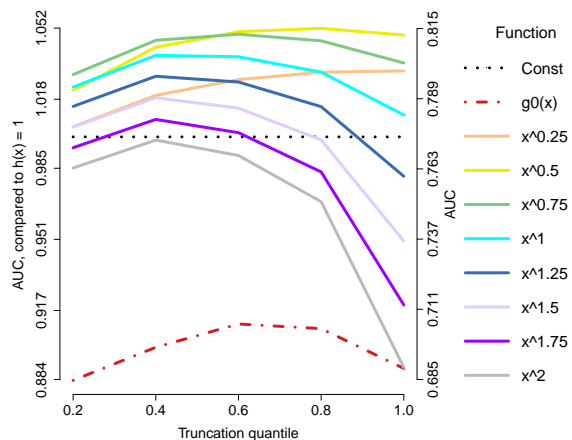
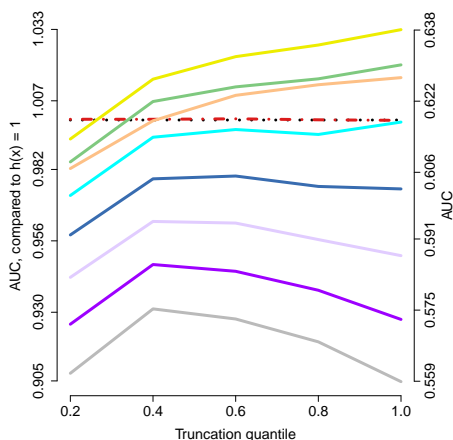
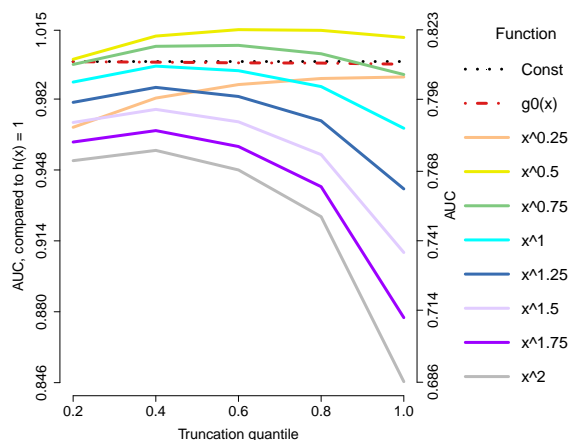
(a)  $n = 80$ ,  $\ell_2$ -nn domain(b)  $n = 1000$ ,  $\ell_2$ -nn domain(c)  $n = 80$ ,  $\ell_2^c$ -nn domain(d)  $n = 1000$ ,  $\ell_2^c$ -nn domain(e)  $n = 80$ , unif-nn domain(f)  $n = 1000$ , unif-nn domain

Figure 3.8: AUCs averaged over 50 trials for edge recovery using generalized score matching for the  $a = 3/2$  models.

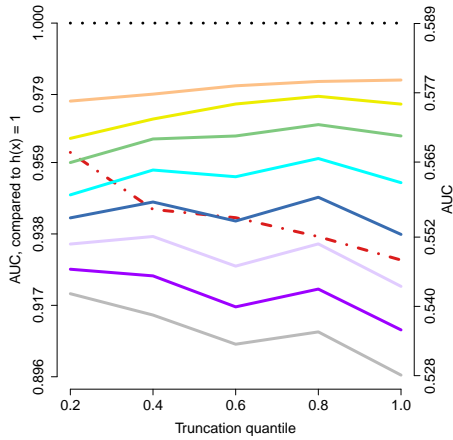
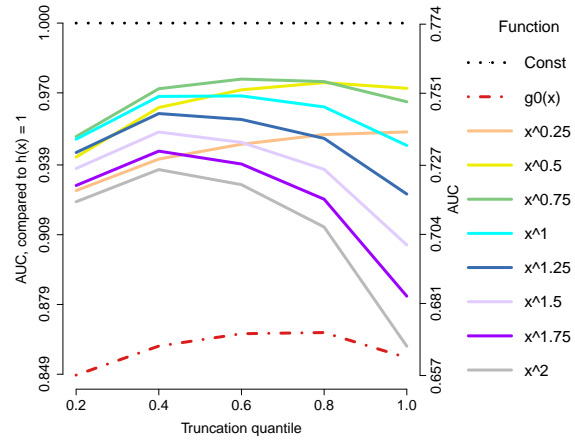
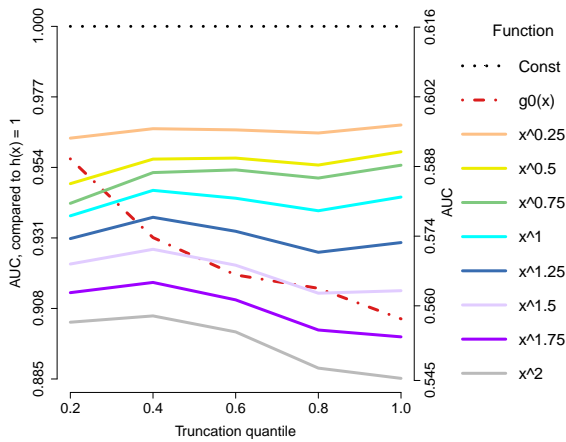
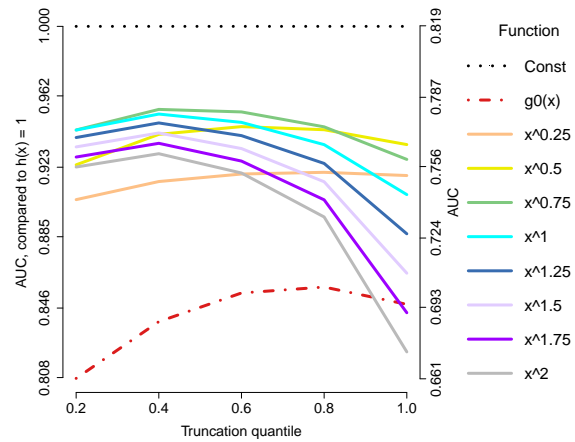
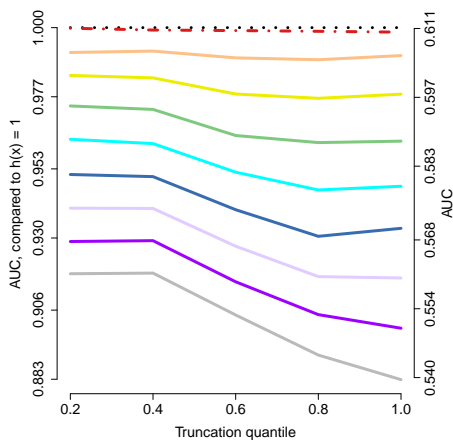
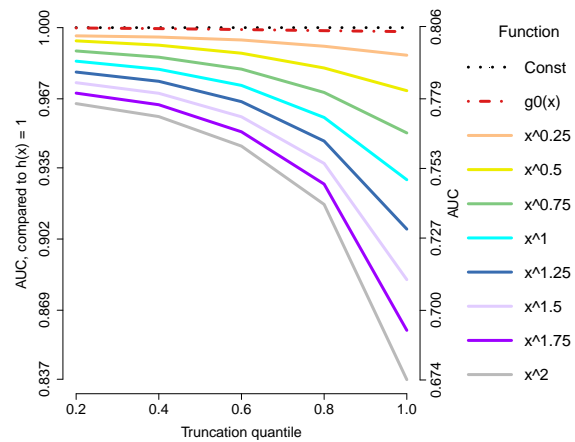
(a)  $n = 80$ ,  $\ell_2$ -nn domain(b)  $n = 1000$ ,  $\ell_2$ -nn domain(c)  $n = 80$ ,  $\ell_2^c$ -nn domain(d)  $n = 1000$ ,  $\ell_2^c$ -nn domain(e)  $n = 80$ , unif-nn domain(f)  $n = 1000$ , unif-nn domain

Figure 3.9: AUCs averaged over 50 trials for edge recovery using generalized score matching for the  $a = 2$  models.

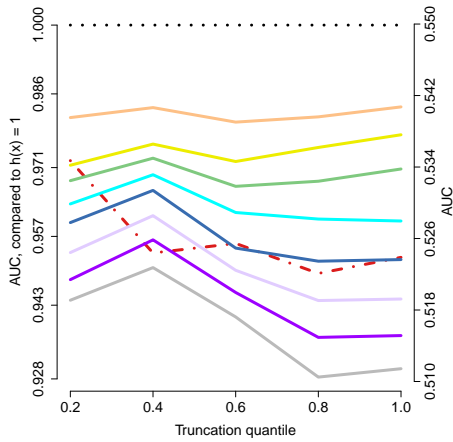
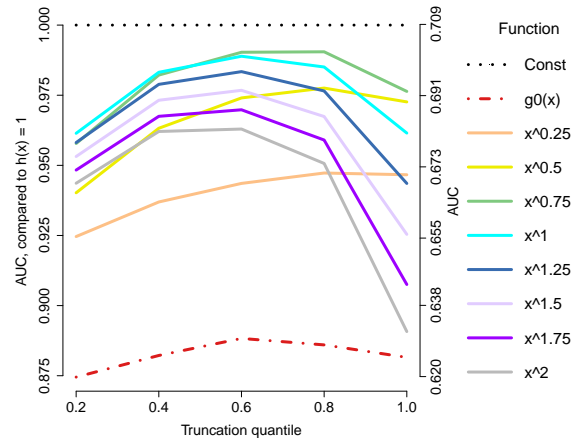
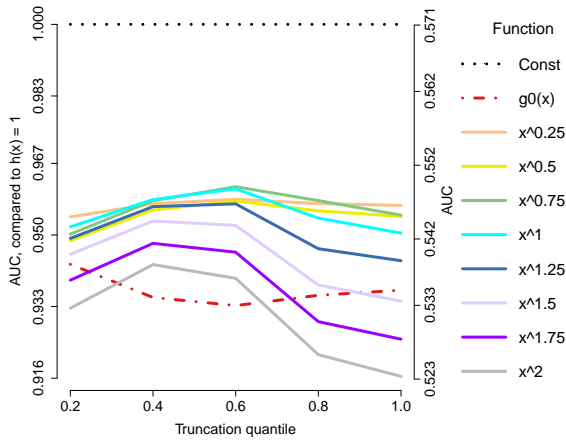
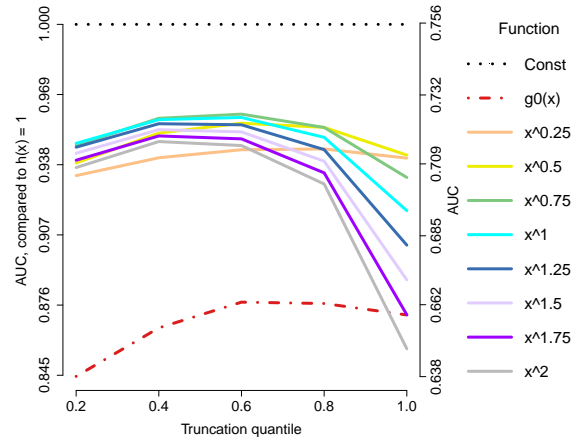
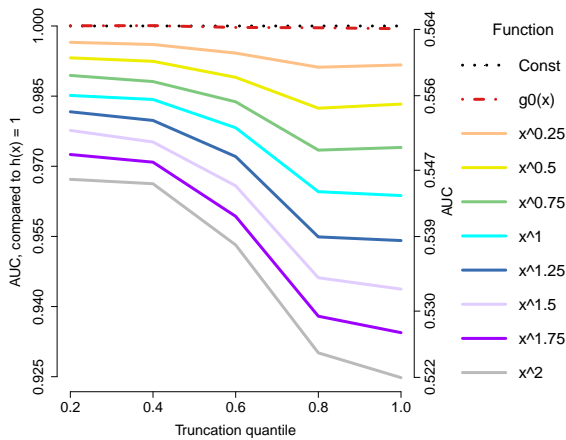
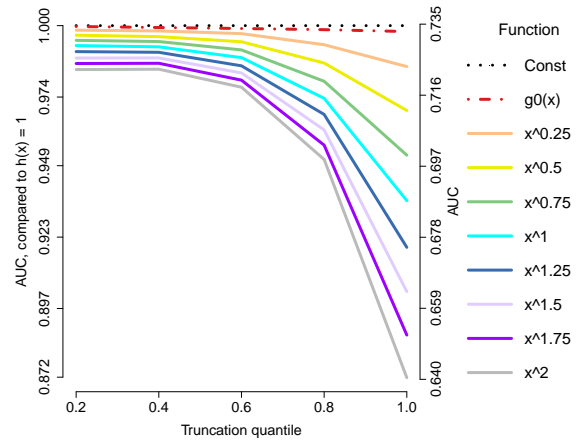
(a)  $n = 80$ ,  $\ell_2$ -nn domain(b)  $n = 1000$ ,  $\ell_2$ -nn domain(c)  $n = 80$ ,  $\ell_2^c$ -nn domain(d)  $n = 1000$ ,  $\ell_2^c$ -nn domain(e)  $n = 80$ , unif-nn domain(f)  $n = 1000$ , unif-nn domain

Figure 3.10: AUCs averaged over 50 trials for edge recovery using generalized score matching for the  $a = 3$  models.

### 3.8.3 Numerical Experiments for $A^{m-1}$ Models on the Simplex

#### Experiment Setup

As the most important example of models on the simplex  $\mathcal{D} \equiv \{\mathbf{x} \in \mathbb{R}_+^m \mid \mathbf{x} \succ \mathbf{0}, \mathbf{1}_m^\top \mathbf{x} = 1\}$ , in this section, we consider  $A^{m-1}$  models discussed in Section 3.6.2, i.e. when  $a = b = 0$  and  $\mathbf{K}\mathbf{1}_m = \mathbf{0}_m$  with  $\mathbf{K} = \mathbf{K}^\top$ . As in Section 3.8.2 we consider dimension  $m = 100$ , sample sizes  $n = 80$  and  $n = 1000$ , and assume  $\boldsymbol{\eta}_0 \equiv \mathbf{0}$  is known for simplicity; the density is proportional to  $\exp(-\log \mathbf{x}^\top \mathbf{K} \log \mathbf{x}/2)$ .

For the true interaction matrix  $\mathbf{K}_0$  we use band matrices with bandwidths  $s = 7$  for  $n = 1000$  and  $s = 2$  for  $n = 80$ , where the bandwidth is defined as  $\max\{|i - j| : \kappa_{0,i,j} > 0\}$ . We set  $\kappa_{0,i,j}$  to  $1 - |i - j|/(s + 1)$  for  $1 \leq |i - j| \leq s$ , and set the diagonals so that  $\mathbf{K}_0\mathbf{1}_m = \mathbf{0}_m$ . Each  $n$  is thus associated with only one graph, for which we run 50 trials. Note that the bandwidth is chosen so that  $n/(d_{\mathbf{K}_0}^2 \log m)$  is roughly constant, where  $d$  is the maximum node degree, as this quantity is suggested to be linked to the probability of successful edge recovery, according to our consistency theory in Section 3.7.

#### Results

As discussed in Section 3.6.2, although conditional dependence between two components does not correspond to the zero/nonzero pattern of the entries in  $\mathbf{K}$ , one may still be interested in successful recovery of such pattern as well as minimization of the estimation error of  $\mathbf{K}$  and  $\boldsymbol{\eta}$ . This is backed up by Corollary 22, according to which  $\mathbf{K}$  and  $\boldsymbol{\eta}$  are exactly identifiable from the distribution.

As in Section 3.8.2, in Figure 3.11 we plot the AUC averaged over 50 trials against the  $\pi = 0.2, 0.4, 0.6, 0.8, 1$  whose columnwise sample quantiles are used as the truncation points  $\mathbf{C}$  of  $\boldsymbol{\varphi}_{\mathbf{C},\mathcal{D}}$  (c.f. Section 3.8.1), and each curve represents  $h(x) = 1$  (Hyvärinen, 2005),  $g_0(\mathbf{x})$  from Section 3.3.6, or  $h(x) = x^c$  with  $c = 1/4, 1/2, \dots, 2$ . In addition, in Figure 3.12 we plot the estimation error in spectral and Frobenius norms, namely  $\|\hat{\mathbf{K}} - \mathbf{K}_0\|_2$  and  $\|\hat{\mathbf{K}} - \mathbf{K}_0\|_F$ , against  $\pi$ . The  $y$ -ticks on the right-hand side are the errors, and those on the left are the

errors divided by the error for  $h(x) = 1$ , measuring the relative performance of each method compared to Hyvärinen (2005). Opposite to the earlier plots, curves closer to the bottom of each plot indicate better performance.

Based on the plots, for both edge recovery and error measured by the Frobenius norm,  $h(x) = x^2$  is among the best performers, supporting our previous conclusion of the choice of  $h(x) = x^{\max\{2-a, 0\}}$  for general  $a$ - $b$  models. When the error is measured in the spectral norm,  $h(x) = x^2$  has the largest error for  $n = 80$  but shows better improvements over other estimators as  $n$  is increased.

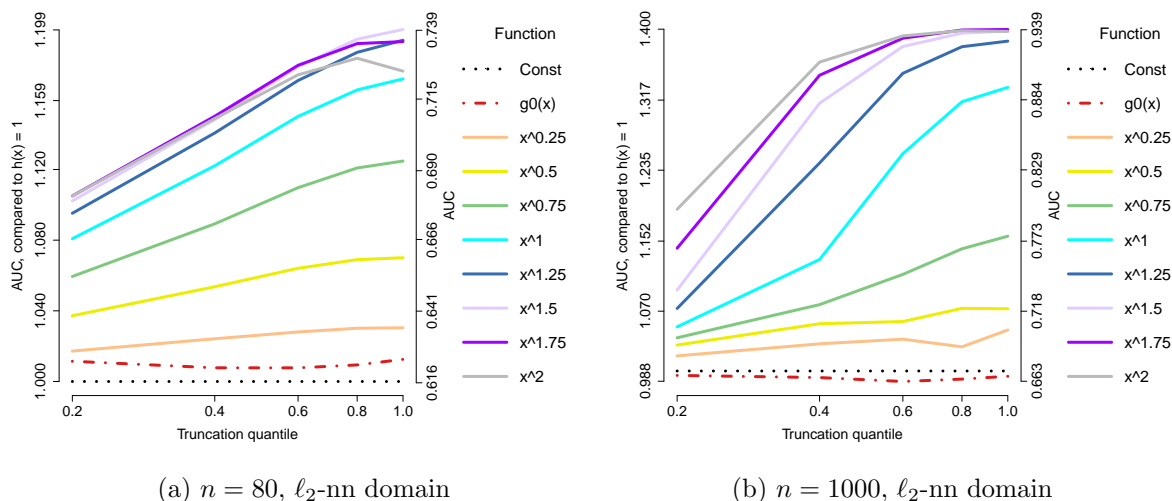
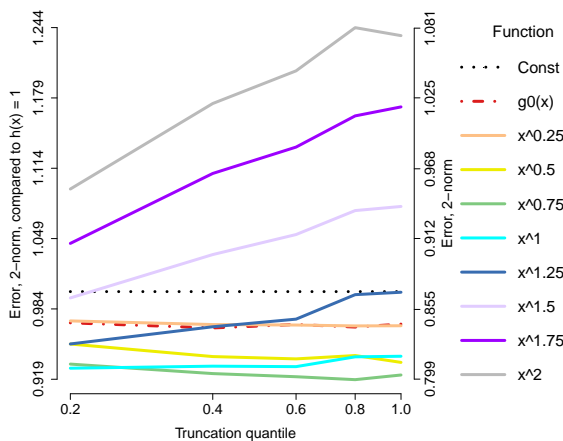


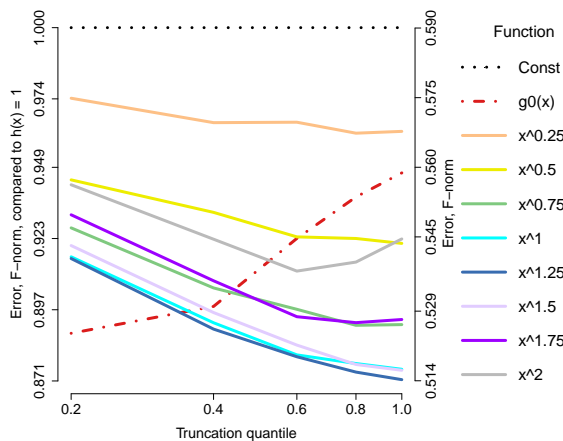
Figure 3.11: AUCs averaged over 50 trials for edge recovery for the  $A^{m-1}$  models on the simplex.

### 3.9 DNA Methylation Data

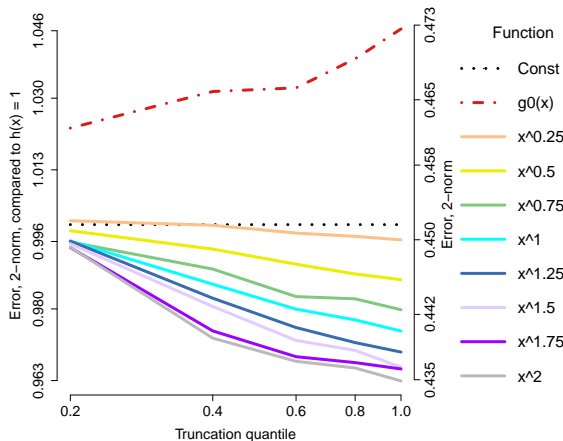
In this section, we apply our generalized score matching to a DNA methylation dataset. The dataset contains methylation levels of CpG islands for 500 patients associated with head and neck cancer from the Cancer Genome Atlas (TCGA), where CpG islands are regions of DNA where a guanine nucleotide very frequently follows a cytosine nucleotide along its 5'  $\rightarrow$  3' direction. Methylation levels are associated with epigenetic regulation of genes (Du



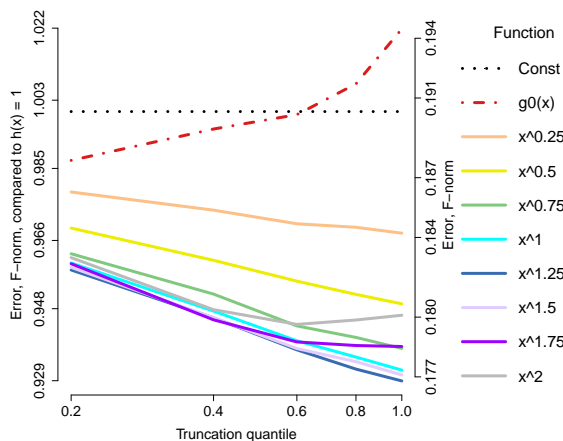
(a)  $n = 80$ , error in matrix 2-norm



(b)  $n = 1000$ , error in matrix 2-norm



(c)  $n = 80$ , error in matrix  $F$ -norm



(d)  $n = 1000$ , error in matrix  $F$ -norm

Figure 3.12: Averaged error in spectral and  $F$  norms over 50 trials normalized by the corresponding norms of the true  $\mathbf{K}_0$ ; sparsity chosen by cross validation for the  $A^{m-1}$  models on the simplex.

et al., 2010). According to Du et al. (2010), to measure the methylation level, one can use the Beta-values, a value in  $[0, 1]$  as the ratio of the methylated probe intensity and the sum of methylated and unmethylated probe intensities, or the M-values, which is defined as  $\log_2(\text{Beta}/1 - \text{Beta})$  taking values in  $\mathbb{R}$ , the logit of the Beta-values with base 2. We are interested in applying our method and the  $a$ - $b$  model framework to a subset of the data using Beta-values and M-values to estimate the network/undirected graph of the CpG sites or the corresponding genes.

Since the dataset contains 500 samples and 93715 sites, we choose a subset of sites corresponding to genes known to belong to the pathway for Thyroid cancer according to the Kyoto Encyclopedia of Genes and Genomes (KEGG). With the  $a$ - $b$  model assumption, we remove the sites that are clearly bimodal; specifically, for each site, we apply `Mclust` to its non-outlier Beta-value observations and test if it is a Gaussian mixture having more than one component. This results in 500 samples and 478 sites, corresponding to 36 genes.

To estimate the graph using the M-values, we assume a Gaussian model on  $\mathbb{R}^m$ , i.e.  $a$ - $b$  model with  $a = b = 1$ , and use the profiled estimator in (3.22), and choose the upper-bound diagonal multiplier  $2 - \left(1 + 80\sqrt{\log m/n}\right)^{-1}$  as suggested in Chapter 2.6.2. Note that the  $(\mathbf{h} \circ \boldsymbol{\varphi})$  is unnecessary in this setting and we simply use the original score matching with  $(\mathbf{h} \circ \boldsymbol{\varphi})(\mathbf{x}) = \mathbf{1}_m$ . For Beta-values, we assume a log-log model ( $a = b = 0$ ) on  $[0, 1]^m$ , and use the profiled estimator with the upper-bound diagonal multiplier  $1 + \sqrt{(\tau \log m + \log 4)/(2n)}$  as in (3.38) with the choice of  $\tau = 3$ . We use  $\mathbf{h}(\mathbf{x}) = \mathbf{x}^2$  as suggested by our theory, and choose the truncation points in  $\boldsymbol{\varphi}$  to be the 40th sample percentile, as suggested by the simulation results in Figure 3.4. As an illustration, the  $\lambda$  parameter that defines the  $\ell_1$  penalty on  $\mathbf{K}$  is chosen so that the number of edges is equal to 478, the number of sites, following Section 2.7.4 and Lin et al. (2016).

The estimated graphs are presented in Figure 3.13, where graphs (a) and (d) are the estimated graphs for Beta values, graphs (c) and (f) are those for M values, and (b) and (e) are the common edges, i.e. intersection graphs. In particular, graphs (a) (b) (c) exclude isolated nodes that have no edges and the layout is optimized for each graph, while for (d)

(e) (f) the layout is optimized for the Beta graph. Figure 3.14 shows graphs in Figure 3.13 aggregated by the genes associated with the sites. In (a), (c), (d), (f), red points indicate nodes with degree at least 10 (i.e. genes connected to at least 10 other genes in the case of Figure 3.14). Sites with highest node degrees are listed in Table 3.1, where those shared by the two graphs are highlighted in bold.

Beta values	M values
CDH1—4 (28)	RXRB—24 (25)
<b>TCF7L1—18 (22)</b>	<b>MAPK3—8 (22)</b>
<b>RXRA—19 (21)</b>	PAX8—6 (21)
RXRA—22 (21)	CCND1—19 (20)
RET—22 (21)	RXRA—10 (20)
RXRB—82 (21)	<b>RXRA—19 (20)</b>
NTRK1—40 (21)	RXRB—18 (20)
LEF1—2 (20)	PAX8—9 (20)
TCF7L1—13 (20)	TCF7—3 (18)
CDKN1A—10 (20)	TCF7L1—9 (18)
CDKN1A—6 (19)	<b>TCF7L1—18 (18)</b>
<b>MAPK3—8 (17)</b>	TCF7L2—63 (18)
PAX8—28 (17)	TPM3—12 (18)
<b>PAX8—29 (17)</b>	<b>PAX8—29 (17)</b>

Table 3.1: List of sites with the highest node degrees in each estimated graph.

To quantify the similarity between the two graphs, we calculate their Hamming distance, namely the sum of absolute differences in their adjacency matrices divided by 2. The Hamming distance is 568, which is a lot smaller than 936, the minimal Hamming distance between the graph for Beta values and 10000 randomly generated graphs with the same number of edges, and 940, that value using the graph for M values. Furthermore, the discrete one-

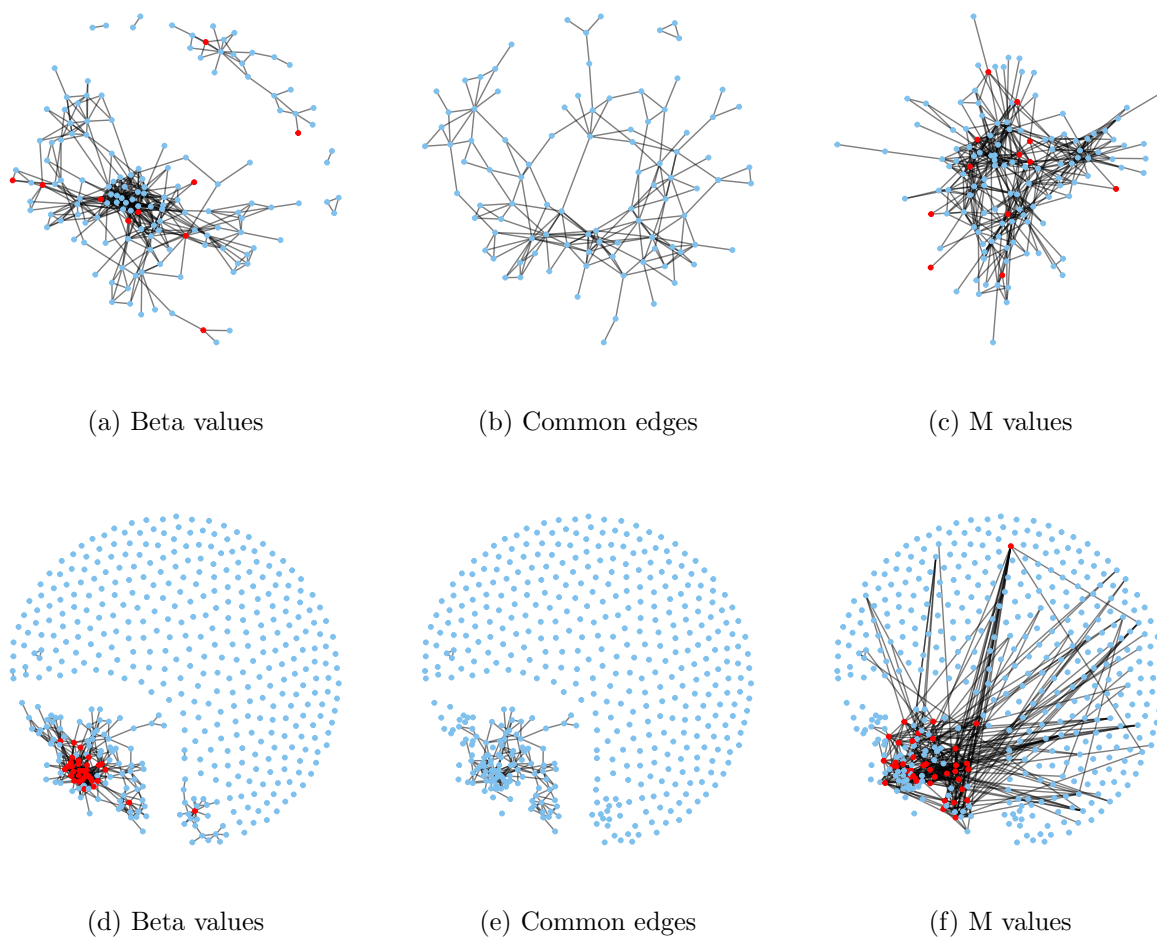


Figure 3.13: Graphs for CpG sites estimated by regularized generalized score matching estimator using Beta values (a, d) and M values (c, f), and their intersection graph (b, e). In (a), (b) and (c), isolated nodes with no edges are removed, and the layout is optimized for each plot. In (d), (e) and (f), isolated nodes are included and the layout is optimized for the graph for Beta values. In (a), (c), (d), (f), red points indicate nodes with degree at least 10 (“hub nodes”).

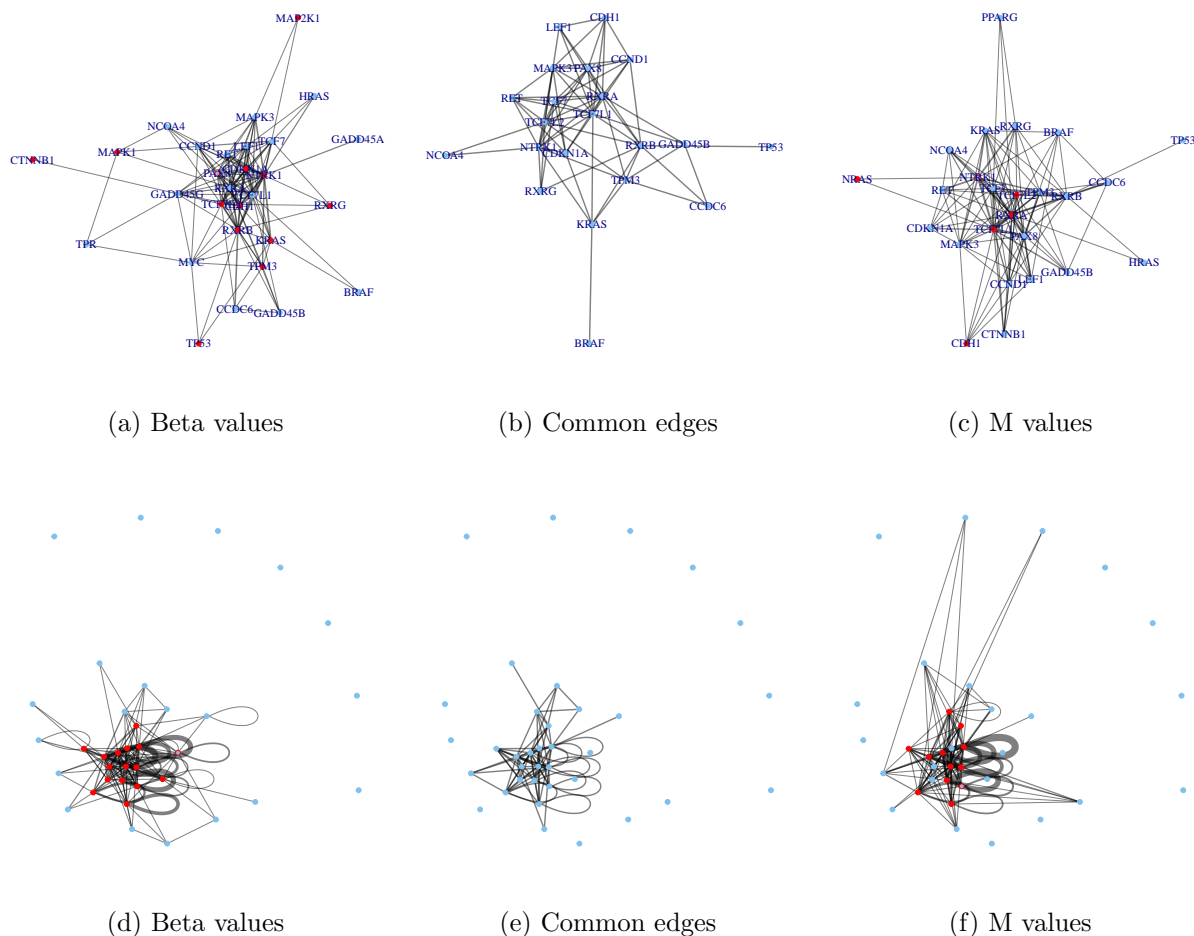


Figure 3.14: Graphs in Figure 3.13 aggregated by the genes associated with the CpG sites, with Beta values (a, d) and M values (c, f), and their intersection graph (b, e). In (a), (b) and (c), isolated nodes and self loops (edges between sites for the same gene) are removed, and the layout is optimized for each plot. In (d), (e) and (f), isolated nodes and self loops are included and the layout is optimized for the graph for Beta values. In (a), (c), (d), (f), red points indicate genes connected to at least 10 other genes. Thickness of edges scales with square root of the total number of edges for sites between each pair of genes (or edges between sites within the same group in the case of self loops).

sample two-sided Kolmogorov-Smirnov test on the sample distribution of node degrees of one graph, while treating that of the other graph as the null distribution, gives  $p$ -values 0.1208 and 0.1397, respectively. In Figure 3.15 we compare the distribution of node degrees for the site graphs estimated using Beta and M values, with interlaced histogram on the left and Q-Q plot on the right. All these results suggest that the two estimated graphs are sufficiently similar to each other.

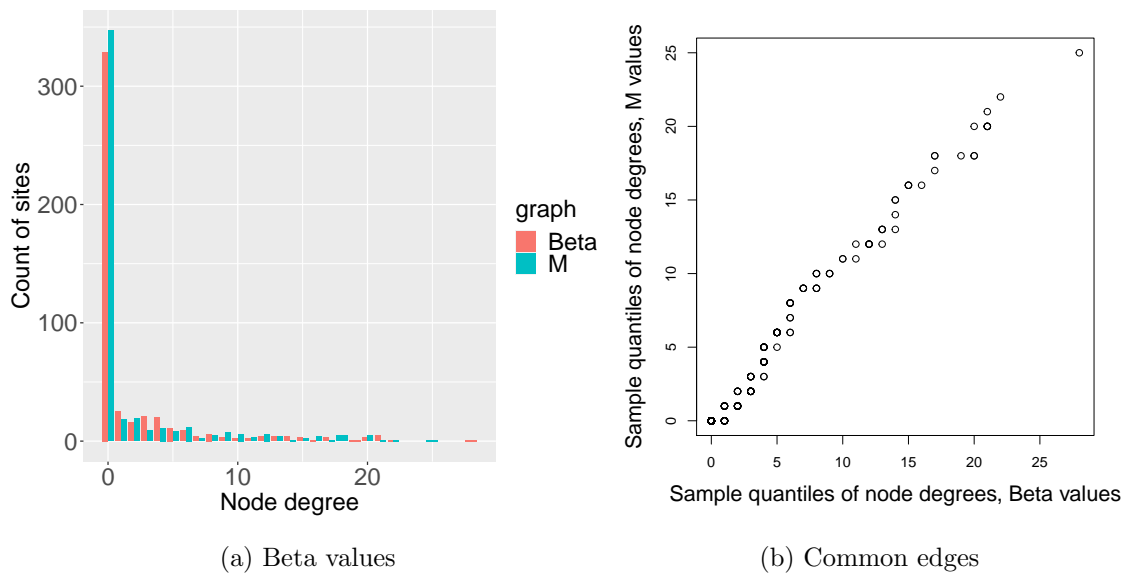


Figure 3.15: Interlaced histogram (left) and Q-Q plot (right) showing the node degree distributions for both site graphs.

### 3.10 Discussion

The generalized score matching loss proposed in Chapter 2 (Yu et al., 2019b) estimates densities supported on  $\mathbb{R}_+^m$  using a generalized score matching loss, whose definition involves componentwise multiplication of  $\nabla \log p(\mathbf{x})$  with an  $\mathbf{h}(\mathbf{x})$  function, generalizing the choice of  $\mathbf{h}(\mathbf{x}) = \mathbf{x}^2$  in Hyvärinen (2007). As in Hyvärinen (2007), the estimator has a closed-form solution in the low-dimensional settings without any regularization and avoids calculating the often computationally intensive normalizing constant.

In this chapter, we extend our work from distributions on  $\mathbb{R}_+^m$  to those on *componentwise countable* domains as defined in Definition 5; informally these are the domains for which, fixing all other components  $\mathbf{x}_{-j}$ , the induced domain of one component  $x_j$  is a union of at most countably many intervals in  $\mathbb{R}$ . We accomplish this by composing the  $\mathbf{h}$  with a distance function  $\boldsymbol{\varphi}_{\mathcal{C}} = (\varphi_{C_{1,1}}, \dots, \varphi_{C_{m,m}}) : \mathcal{D} \rightarrow \mathbb{R}_+^m$ , with  $\varphi_{C_{j,j}}(\mathbf{x})$  being the distance of  $x_j$  to its induced boundary holding  $\mathbf{x}_{-j}$  fixed, truncated from above by some prespecified  $C_j$ . We show that the loss can again be approximated by an empirical loss when we substitute  $\mathbf{h}$  in Chapter 2 with  $(\mathbf{h} \circ \boldsymbol{\varphi}_{\mathcal{C}})$ , which is quadratic in the canonical parameter for exponential families, thus having closed-form solutions in low-dimensional settings. For domains with zero Lebesgue measure, we take the  $m$ -simplex domain as a canonical example and defined the loss by profiling out the last component.

As in Chapter 2, we focus on  $a$ - $b$  pairwise interaction models with density proportional to  $\exp\{-\mathbf{x}^a \top \mathbf{K} \mathbf{x}^a / (2a) + \boldsymbol{\eta} \top \mathbf{x}^b / b\}$ , where for  $a = 0$  we let  $\mathbf{x}^a \top \mathbf{K} \mathbf{x}^a / (2a) \equiv \log \mathbf{x} \top \mathbf{K} \log \mathbf{x} / 2$  and for  $b = 0$ ,  $\boldsymbol{\eta} \top \mathbf{x}^b / b \equiv \boldsymbol{\eta} \top \log \mathbf{x}$ . We formulate our estimators for  $a$ - $b$  models on domains with positive Lebesgue measures as well as on simplex domains, for the latter of which we extensively discuss the  $A^{m-1}$  models (Aitchison, 1985) as the most important example.

We extend the consistency theory for edge recovery in Chapter 2 to Gaussian graphical models on domains that are finite disjoint unions of convex sets, as well as general  $a$ - $b$  models on simplex domains with  $a > 0$  and on bounded domains with positive Lebesgue measure, requiring the sample size to be  $n = \Omega(\log m)$ . For unbounded domains with  $a > 0$  and

the simplex domains with  $a = 0$ , we require an additional multiplicative factor that may weakly depend on  $m$ . Finally, through simulation studies we confirm that the choice of  $\mathbf{h}(\mathbf{x}) = (x_1^c, \dots, x_m^c)$  with  $c = \max\{2 - a, 0\}$  performs the best in most settings in terms of edge recovery, similar to the conclusion in Chapter 2.

In the simulations we adaptively select the truncation points  $\mathbf{C}$  of  $\varphi$  using the sample quantiles of the untruncated distances. Developing a method to choose the best truncation points remain further research. For the theory part, it is also interesting to investigate into the real sample complexity in terms of dimension  $m$  for unbounded domains with  $a > 0$  or simplex domains with  $a = 0$ .

## Chapter 4

**DIRECTED GRAPHICAL MODELS AND CAUSAL  
DISCOVERY FOR ZERO-INFLATED DATA**

## 4.1 Introduction

Graphical models specify conditional independence relations among variables in a random vector  $\mathbf{Y}$  indexed by the nodes  $\mathcal{V}$  of a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with edge set  $\mathcal{E}$  (Maathuis et al., 2019). Models based on undirected graphs may be used to explore conditional independence between any two variables  $Y_V$  and  $Y_U$  given all others  $(Y_W)_{W \neq U, V}$ , as represented by the absence of an edge between  $V$  and  $U$  in  $\mathcal{E}$ . Models based on directed acyclic graphs (DAGs), for which  $\mathcal{E}$  is comprised of directed edges, capture conditional independence structure that naturally arises from cause-effect relationships between the variables.

In biology and genetics, graphical models have been applied to infer the structure of gene regulatory networks based on measurements of gene expression (Maathuis et al., 2019, Sections 20-21). Traditional technologies produce expression levels aggregated over hundreds or thousands of individual cells, and these bulk measurements are frequently modeled using the assumption of Gaussianity. In directed GGMs the exact structure of the underlying DAG cannot be identified from purely observational data, and the target of inference becomes an equivalence class of DAGs. For instance, one cannot differentiate between  $V \rightarrow U$  and  $U \rightarrow V$  when the variables are assumed bivariate normal. In the Gaussian case, directed graphical models posit linear functional relationships between the variables coupled with additive Gaussian noise. A more recent line of work emphasizes that directed graphical models that alter this assumption to nonlinear functional relationships and additive noise (Peters et al., 2014), or linear relations and non-Gaussian noise (Shimizu et al., 2006; Wang and Drton, 2020), or linear relations with homoscedastic Gaussian noise (Peters and Bühlmann, 2013; Chen et al., 2019) are amenable to causal discovery in the sense that different DAGs are no longer equivalent.

More recent technology obtains sequencing measurements of mRNA present in single cells. This new technology, as well as the larger sample sizes it provides, promise to give more information than bulk measurements, but at the same time bring in a unique new challenge. At the single cell level, genes appear as “on” with positive single cell gene expression levels,

or as “off” with the recorded measurements zero or negligible (McDavid et al., 2019).

Figure 4.1 shows pairwise scatter plots of four genes from a T helper single-cell dataset with 1951 measurements from eight healthy donors, which we analyze in Section 4.6. It is a superset of the single-cell T-follicular helper data considered in McDavid et al. (2019), which is similarly plotted in their Figure 1. The lower panels show the pairwise scatter plots along with a fitted linear regression curve, and the diagonal panels show the univariate smoothed kernel density estimates for each gene. As we can see, each gene has a large number of zero values and a linear regression model is not sufficient for modeling the pairwise relationships.

A novel undirected graphical model that deals with this zero-inflation was introduced by McDavid et al. (2019). Their approach considers Hurdle density models, where for a random vector of dimension  $m$ , the joint probability density function has the form

$$f(\mathbf{y}; \mathbf{A}, \mathbf{B}, \mathbf{K}) \propto \exp\left(\mathbf{1}_y^\top \mathbf{A} \mathbf{1}_y + \mathbf{1}_y^\top \mathbf{B} \mathbf{y} - \frac{1}{2} \mathbf{y}^\top \mathbf{K} \mathbf{y}\right), \quad \mathbf{y} \in \mathbb{R}^m, \quad (4.1)$$

with  $\mathbf{1}_y$  being the elementwise indicators of nonzero entries in  $\mathbf{y} \in \mathbb{R}^m$ . The dominating measure for this density is the  $m$ -fold product of the sum of the Lebesgue measure and a point mass at zero. However, since more information can be inferred from single-cell sequencing data, one would hope that the data can also be analyzed using more informative directed graphical models, and that we can infer which variables (genes) are the causes of change in other variables (expression levels of other genes). In this chapter, we formulate such directed graphical models for zero-inflated data, and prove that under a weak assumption one can recover the exact DAG from the joint distribution. In contrast to the setting of McDavid et al. (2019), the distributions in our models are not merely zero-inflated Gaussian as we allow variables that are “on” to be non-linear polynomial functions of other variables and stochastic noise.

In DAG models, the joint distribution can be factorized into the product of conditional distributions of each variable given parent variables. For simplicity we call these conditional distributions the *node conditionals*. In our DAG model for zero-inflated data, we form the node conditionals by taking the conditional distribution of one variable given the others

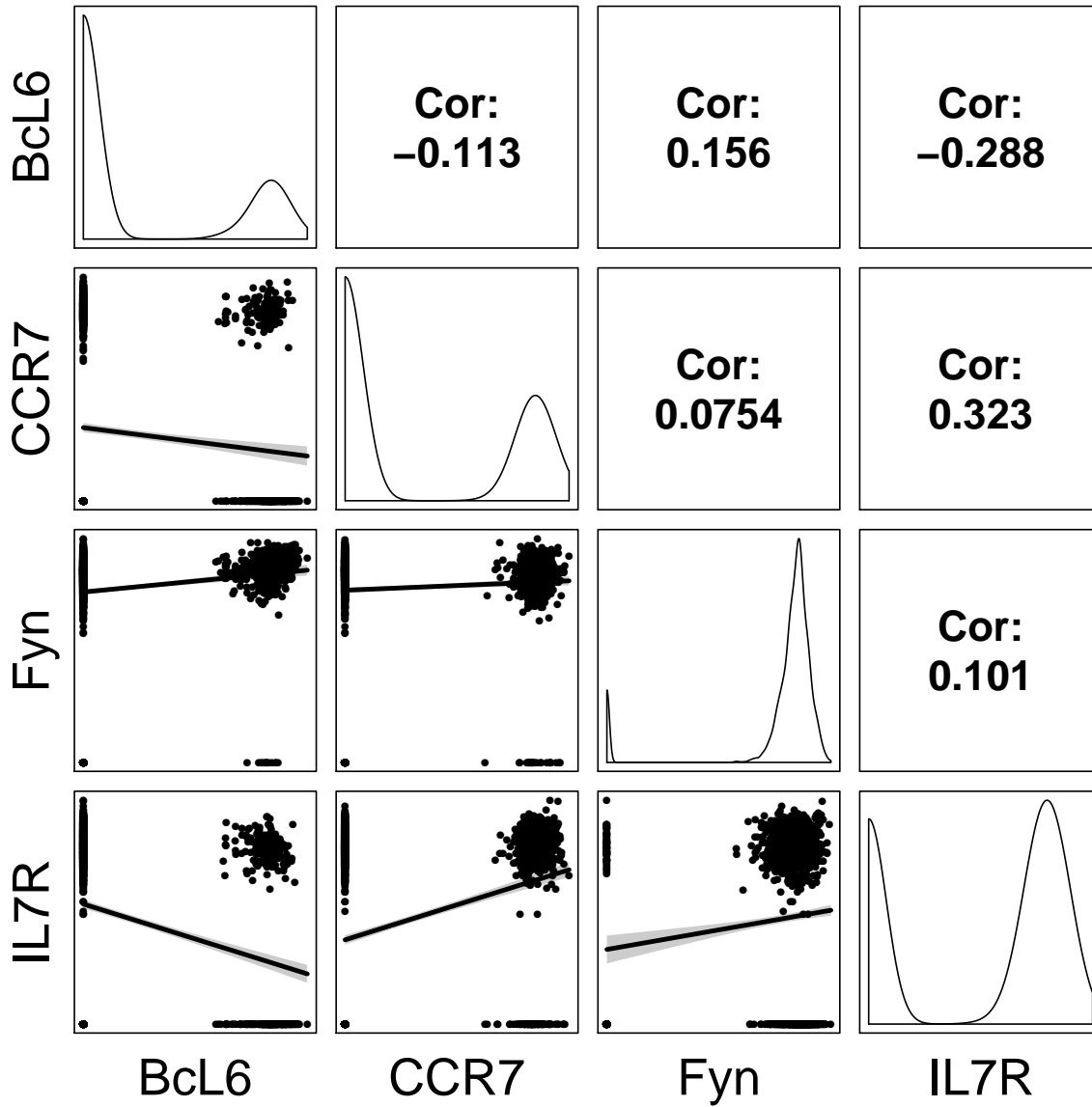


Figure 4.1: Pairwise scatter plots and kernel densities on four genes from the T helper cell data analyzed in Section 4.6.

in the joint model from (4.1). We refer to the resulting graphical model as the model in  $(\alpha, \beta, k)$ -parametrization, or in *canonical* parametrization. Here, the  $\alpha$  and  $\beta$  parameters in the conditional distribution are derived from the matrix parameters  $\mathbf{A}$  and  $\mathbf{B}$  in (4.1). The  $\alpha$  and  $\beta$  are polynomials in the parent variables and their 0/1 indicators being zero/nonzero. An alternative second type of model is obtained by directly specifying each node conditional as a mixture of a point mass at zero and a Gaussian distribution, with the log odds of being nonzero ( $\log(p/(1-p))$ ) and the mean in the Gaussian part being polynomials in the parent variables and their indicators. The Gaussian variance is taken constant in the parents. We call this second formulation the model in  $(p, \mu, \sigma^2)$ - or *moment* parametrization, since the parameters directly correspond to the (conditional) moments. The detailed specification of both model types is developed in Section 4.2.

In Section 4.3, we show that under our models, the distributions that can be represented by two different DAGs must be distributions of *two-Gaussian type* (Definition 15). We then prove that such distributions do not exist for dimension  $m = 2$  and  $m = 3$ ; we also conjecture they do not exist for  $m > 3$ . Moreover, we are able to prove that under a natural and practical assumption, we have full identifiability in the sense of being able to identify the exact DAG underlying the model. This assumption specifies that for each node,  $\alpha + \beta^2/(2k)$  or equivalently  $\log(p/(1-p))$  has a separate univariate term for each parent (e.g.  $y_1 + y_2 + y_1y_2 + y_1^2$  instead of  $y_1 + y_1y_2 + y_1^2$ , which does not have a separate term for  $y_2$ ). Note that this assumption is natural in the context of linear regression.

In Section 4.4, we introduce different methods for estimation of the DAG. Simulation studies supporting the use of these methods are given in Section 4.5, and they are then applied to the T-follicular helper cell dataset (Section 4.6).

The work in this chapter is presented in Yu et al. (2020).

## 4.2 Directed Graphical Models for Zero-Inflated Data

In this section we motivate and formally define our models for zero-inflated data based on directed acyclic graphs (DAGs).

#### 4.2.1 Hurdle Joint Distributions for Zero-Inflated Continuous Observations

McDavid et al. (2019) proposed a *Hurdle joint distribution* with density

$$f(\mathbf{y}; \mathbf{A}, \mathbf{B}, \mathbf{K}) \propto \exp \left( \mathbf{1}_{\mathbf{y}}^{\top} \mathbf{A} \mathbf{1}_{\mathbf{y}} + \mathbf{1}_{\mathbf{y}}^{\top} \mathbf{B} \mathbf{y} - \frac{1}{2} \mathbf{y}^{\top} \mathbf{K} \mathbf{y} \right), \quad \mathbf{y} \in \mathbb{R}^m, \quad (4.2)$$

with respect to  $\lambda^m$ , where  $\lambda$  is the sum of a point mass at 0 and the Lebesgue measure on  $\mathbb{R}$ , and  $\mathbf{A} = (\alpha_{ij})_{i,j}$ ,  $\mathbf{B} = (\beta_{ij})_{i,j}$ ,  $\mathbf{K} = (k_{ij})_{i,j} \in \mathbb{R}^{m \times m}$  are matrices of interaction parameters with  $\mathbf{K}$  positive definite. The indicator vector  $\mathbf{1}_{\mathbf{y}} \equiv (\mathbb{1}_{\{y_1 \neq 0\}}, \dots, \mathbb{1}_{\{y_m \neq 0\}}) \in \{0, 1\}^m$  captures which components of  $\mathbf{y}$  are non-zero.

Consider a random vector  $\mathbf{Y} \in \mathbb{R}^m$  that follows the Hurdle joint distribution. Intuitively, the density in (4.2) is obtained by combining an Ising model for the indicator vector  $\mathbf{1}_{\mathbf{Y}}$  and a conditional normal distribution for  $\mathbf{Y}$  given its nonzero pattern  $\mathbf{1}_{\mathbf{Y}}$ . The Ising model postulates a probability mass function proportional to  $\exp(\mathbf{1}_{\mathbf{y}}^{\top} \mathbf{A} \mathbf{1}_{\mathbf{y}})$ . The conditional normal distribution has density  $p(\mathbf{Y} = \mathbf{y} | \mathbf{1}_{\mathbf{Y}} = \mathbf{1}_{\mathbf{y}}; \mathbf{B}, \mathbf{K}) \propto \exp(\mathbf{1}_{\mathbf{y}}^{\top} \mathbf{B} \mathbf{y} - \frac{1}{2} \mathbf{y}^{\top} \mathbf{K} \mathbf{y})$  with respect to the Lebesgue measure restricted to the subspace of  $\mathbb{R}^m$  compatible with  $\mathbf{1}_{\mathbf{y}}$ .

The exponential specification in (4.2) entails that conditional independence between two variables is equivalent to the corresponding entries in all interaction matrices  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{K}$  being 0. In other words,  $\alpha_{ij} = \alpha_{ji} = \beta_{ij} = \beta_{ji} = k_{ij} = k_{ji} = 0$  if and only if  $Y_i$  and  $Y_j$  are conditionally independent given all other variables. Indeed, it is easy to see that the induced conditional distribution of  $Y_i$  given all other variables  $\mathbf{Y}_{-i}$  in  $\mathbf{Y}$ , has density

$$p(Y_i = y_i | \mathbf{Y}_{-i} = \mathbf{y}_{-i}) = f(y_i; \alpha_{ii} + \boldsymbol{\alpha}_{i,-i}^{\top} \mathbf{1}_{\mathbf{y}_{-i}} + \boldsymbol{\beta}_{i,-i}^{\top} \mathbf{y}_{-i}, \beta_{ii} + \boldsymbol{\beta}_{-i,i}^{\top} \mathbf{1}_{\mathbf{y}_{-i}} - \mathbf{k}_{i,-i}^{\top} \mathbf{y}_{-i}, k_{ii}), \quad (4.3)$$

that is, the distribution is a Hurdle distribution in  $m = 1$  dimension with parameters  $\alpha$ ,  $\beta$ , and  $k$  being linear functions in  $\mathbf{Y}_{-i}$  and  $\mathbf{1}_{\mathbf{Y}_{-i}}$ ; here  $f$  is the univariate version of (4.2).

#### 4.2.2 Hurdle Conditionals

The observation in (4.3) above gives rise to the following definition. Recall that  $\lambda$  is the sum of a point mass at 0 and the Lebesgue measure on  $\mathbb{R}$ .

**Definition 10** ( $(\alpha, \beta, k)$ -Hurdle conditionals). *Given a scalar random variable  $X$  and an  $m$ -dimensional random vector  $\mathbf{Z}$ , we say that the conditional distribution of  $X$  given  $\mathbf{Z}$  is of  $(\alpha, \beta, k)$ -Hurdle type if it admits conditional densities with respect to  $\lambda$  of the form*

$$p(X = x | \mathbf{Z} = \mathbf{z}) = f_{\alpha, \beta, k}^{(m)}(X | \mathbf{Z}) \equiv \frac{\exp(\alpha(\mathbf{z})\mathbb{1}_x + \beta(\mathbf{z})x - kx^2/2)}{\sqrt{2\pi/k} \exp(\alpha(\mathbf{z}) + \beta^2(\mathbf{z})/(2k)) + 1}. \quad (4.4)$$

Here,  $\alpha$  and  $\beta$  are functions of  $\mathbf{Z}$  (and its indicator vector).

Reparametrizing we give another intuitive formulation of Hurdle conditionals that clearly exhibits their nature of a mixture between a point mass at 0 and a conditional Gaussian distribution.

**Definition 11** ( $(p, \mu, \sigma^2)$ -Hurdle conditionals). *Given a scalar random variable  $X$  and an  $m$ -dimensional random vector  $\mathbf{Z}$ , we say that the conditional distribution of  $X$  given  $\mathbf{Z}$  is of  $(p, \mu, \sigma^2)$ -Hurdle type if it admits conditional densities with respect to  $\lambda$  of the form*

$$p(X = x | \mathbf{Z} = \mathbf{z}) = f_{p, \mu, \sigma^2}^{(m)}(X | \mathbf{Z}) \equiv (1 - p(\mathbf{z}))(1 - \mathbb{1}_x) + p(\mathbf{z})\mathbb{1}_x \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu(\mathbf{z}))^2}{2\sigma^2}\right). \quad (4.5)$$

Here,  $p$  and  $\mu$  are functions of  $\mathbf{Z}$  (and its indicator vector).

It is easy to show that the two parametrizations (4.4) and (4.5) are connected through

$$\log \frac{p}{1-p} = \alpha + \frac{\beta^2}{2k} - \frac{1}{2} \log\left(\frac{k}{2\pi}\right), \quad \mu = \frac{\beta}{k}, \quad \sigma^2 = \frac{1}{k}. \quad (4.6)$$

That is, the conditional log odds of being nonzero is linear in  $\alpha$  and quadratic in  $\beta$ , and the conditional Gaussian mean is proportional to  $\beta$ .

We note that while the  $(\alpha, \beta, k)$ -parametrization takes canonical parameters  $\alpha(\mathbf{Z})$ ,  $\beta(\mathbf{Z})$  and  $k$  using a representation as exponential family, the moment parametrization directly models the conditional mixing probability  $p(\mathbf{Z})$ , and the mean  $\mu(\mathbf{Z})$  and variance  $\sigma^2$  parameters of the conditional Gaussian distribution. We thus refer to (4.4) as the *canonical parametrization*, and (4.5) as the *moment parametrization*.

### 4.2.3 Directed Graphical Models for Zero-Inflation Data

Consider an  $m$ -dimensional random vector  $\mathbf{Y}$  whose components are indexed by the vertices of a DAG  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and whose distribution is dominated by a product measure on  $\mathbb{R}^m$ . A graphical model based on  $\mathcal{G}$  requires that the density of the joint distribution admits a factorization as

$$f(\mathbf{y}) = \prod_{V \in \mathcal{V}} f_V(y_V | \mathbf{y}_{\text{pa}(V)}), \quad (4.7)$$

where each factor  $f_V(y_V | \mathbf{y}_{\text{pa}(V)})$  is a conditional density for  $y_V$  given its parent variables  $\mathbf{y}_{\text{pa}(V)}$ . The set of parents is defined to be  $\text{pa}(V) \equiv \{U : U \rightarrow V \in \mathcal{E}\}$ .

In Section 4.2.1 we observed that, for the Hurdle joint distributions from (4.2), the conditional distribution of one variable  $Y_i$  given the others is an  $(\alpha, \beta, k)$ -Hurdle with  $k$  constant, and  $\alpha$  and  $\beta$  linear functions of those variables (and their indicators) that are conditionally dependent on  $Y_i$ ; see (4.3). Motivated by this fact, we specify directed graphical models for zero-inflated data by assuming the conditional densities in the factorization in (4.7) to be  $(\alpha, \beta, k)$ - or  $(p, \mu, \sigma^2)$ -Hurdle conditionals. We then assume the parameters in these conditionals to be *Hurdle polynomials* in its parents, as defined now.

**Definition 12** (Hurdle polynomials). *Let  $\mathbf{Y} = (Y_V)_{V \in \mathcal{V}}$  be an  $m$ -dimensional random vector indexed by a set  $\mathcal{V}$ , and suppose  $\mathcal{S} \subseteq \mathcal{V}$ . If  $\mathcal{S} \neq \emptyset$ , define the space of Hurdle polynomials in  $\mathbf{y}_{\mathcal{S}}$  as*

$$\mathcal{H}(\mathbf{Y}; \mathcal{S}) \equiv \left\{ c_0 + \sum_{j=1}^T c_j \prod_{U \in \mathcal{U}_j} Y_U^{d_{j,U}} \prod_{V \in \mathcal{V}_j} \mathbf{1}_{Y_V}, \quad c_0 \in \mathbb{R}, T \in \mathbb{N}, c_j \neq 0, \right. \\ \left. \mathcal{U}_j \subseteq \mathcal{S}, \mathcal{V}_j \subseteq \mathcal{S} \setminus \mathcal{U}_j, d_{j,U} \in \mathbb{N} \forall U \in \mathcal{U}_j \forall j = 1, \dots, T \right\}, \quad (4.8)$$

where  $\mathbb{N} = \{1, 2, \dots\}$ . This is the set of polynomials in values and indicators of nodes in  $\mathcal{S}$ . If  $\mathcal{S} = \emptyset$ , define  $\mathcal{H}(\mathbf{Y}; \mathcal{S}) \equiv \mathbb{R}$ . The degree of a hurdle polynomial as specified in (4.8) is  $\max_{j=1, \dots, T} \sum_{U \in \mathcal{U}_j} d_{j,U} + |\mathcal{V}_j|$ . Here  $|\cdot|$  denotes the set cardinality.

We are now ready to formally define our models.

**Definition 13** (DAG models for zero-inflated data). *Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a DAG with  $|\mathcal{V}| = m$  nodes. A zero-inflated conditional Gaussian DAG model associated with  $\mathcal{G}$  is a set of joint distributions on  $\mathbb{R}^m$  that admit a density (with respect to  $\lambda^m$ ) that factors as in (4.7) with each conditional density  $f_V(y_V | \mathbf{y}_{\text{pa}(V)})$  being a Hurdle conditional*

- (1) *in the  $(\alpha, \beta, k)$ -parametrization with parameters  $\alpha_V$ ,  $\beta_V$  and  $k_V$ , where  $k_V$  is constant,  $\alpha_V$  and  $\beta_V$  are Hurdle polynomials in  $\mathbf{y}_{\text{pa}(V)}$ ; or*
- (2) *in the  $(p, \mu, \sigma^2)$ -parametrization with parameters  $p_V$ ,  $\mu_V$  and  $\sigma_V^2$ , where  $\sigma_V^2$  is constant,  $\log(p_V/(1-p_V))$  and  $\mu_V$  are Hurdle polynomials in  $\mathbf{y}_{\text{pa}(V)}$ .*

It is apparent from the relationship (4.6) that if we allow the relevant parameters to be Hurdle polynomials of *any* degree, the two parametrizations are equivalent, meaning that given an underlying DAG, they share the same space of all possible joint distributions. However for computational convenience it is useful to bound the degree. In later applications, we will only consider degrees up to three.

### 4.3 Identifiability

#### 4.3.1 Strong Identifiability

As we show next, the directed graphical models from Definition 13 are amenable to causal discovery in the sense that the DAG underlying the model is uniquely identifiable from a given joint distribution. More precisely, we prove identifiability under an explicit mild assumption on the Hurdle conditionals determining the considered joint distribution.

Let  $\pi(\mathbf{y}_S) \in \mathcal{H}(\mathbf{Y}; \mathcal{S})$  be a Hurdle polynomial for a subset  $\mathcal{S} \subseteq \mathcal{V}$ . For  $U \in \mathcal{S}$ , let  $\pi_U(y_U) \equiv \pi(y_U, \mathbf{0})$  be the restriction of  $\pi(\mathbf{y}_S)$  obtained by setting all entries other than  $y_U$  to zero. Then  $\pi_U(y_U) \in \mathcal{H}(\mathbf{Y}; \{U\})$  is a univariate Hurdle polynomial.

**Definition 14** (Strong Hurdle polynomials). *Let  $\pi(\mathbf{y}_S) \in \mathcal{H}(\mathbf{Y}; \mathcal{S})$ . We call  $\pi(\mathbf{y}_S)$  a strong Hurdle polynomial if all of its restrictions  $\pi_U(y_U)$  take at least three different values. In*

other words, for each  $U \in \mathcal{S}$ , the Hurdle polynomial  $\pi(\mathbf{y}_{\mathcal{S}})$  contains at least one term of the form  $c_j y_U^d$  with  $c_j \neq 0$  and  $d \geq 1$ .

**Theorem 30** (DAG identifiability with strong Hurdle polynomials). *Let  $f(\mathbf{y})$  be a joint density with respect to  $\lambda^m$  that factors according to a DAG  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , as in (4.7). Suppose for each  $V \in \mathcal{V}$ , the conditional  $f_V(y_V | \mathbf{y}_{\text{pa}(V)})$  is of Hurdle type with parameters  $(\alpha_V, \beta_V, k_V)$  or  $(p_V, \mu_V, \sigma_V^2)$ . If for each  $V$ ,  $\alpha_V + \beta_V^2 / (2k_V)$ , or equivalently  $\log(p_V / (1 - p_V))$ , is a strong Hurdle polynomial, then  $f(\mathbf{y})$  does not factor with respect to any other DAG  $\mathcal{G}' \neq \mathcal{G}$ .*

In the proof in Appendix C we show that the given restriction on the parameters of the Hurdle conditionals is actually stronger than what we need for identifiability. However, the assumption of *strong* Hurdle polynomials is very natural in that it specifies a weak form of hierarchy among interactions by requiring that the conditional distributions are parametrized to include at least one univariate power term in every parent variable and not just indicators or interaction terms with other parents.

#### 4.3.2 Weak Identifiability

Without assuming the Hurdle polynomials for the conditional distributions to be *strong*, we can still offer a weaker identifiability result that shows that the distributions in the intersection between the models obtained from two Markov equivalent DAGs with Hurdle polynomial parameters always have to be of what we call *two-Gaussian type*. In our definition of this concept, we write  $\phi(\cdot; \mu, \nu)$  for the univariate normal density function with mean  $\mu$  and inverse variance  $\nu$ .

**Definition 15.** *Let  $\mathbf{Y} = (Y_V)_{V \in \mathcal{V}}$  be a random vector, and let  $W, U \in \mathcal{V}$  be the indices for two of its components. Further, let  $\mathcal{P} \subseteq \mathcal{V} \setminus \{W, U\}$  be a set of additional indices. Then the joint distribution of  $\mathbf{Y}$  is of two-Gaussian type w.r.t.  $(W, U, \mathcal{P})$  if the following holds for both  $V = W$  and  $V = U$ : There exists a constant  $\nu_1^V$ , polynomials  $\mu_1^V(\mathbf{y}_{\mathcal{P}})$ ,  $\mu_2^V(\mathbf{y}_{\mathcal{P}})$ ,  $\nu_2^V(\mathbf{y}_{\mathcal{P}})$ , and functions  $c_1^V(\mathbf{y}_{\mathcal{P}})$  and  $c_2^V(\mathbf{y}_{\mathcal{P}})$  such that for almost every  $\mathbf{y}_{\mathcal{P}} \in \mathbb{R}^{|\mathcal{P}|}$ ,  $c_1^V(\mathbf{y}_{\mathcal{P}}) > 0$ ,  $c_2^V(\mathbf{y}_{\mathcal{P}}) > 0$ , either  $\mu_1^V(\mathbf{y}_{\mathcal{P}}) \neq \mu_2^V(\mathbf{y}_{\mathcal{P}})$  or  $\nu_1^V \neq \nu_2^V(\mathbf{y}_{\mathcal{P}})$ , and the conditional density*

$$\mathbb{P}(Y_V = y | Y_V \neq 0, \mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}}) = c_1^V(\mathbf{y}_{\mathcal{P}})\phi(y; \mu_1^V(\mathbf{y}_{\mathcal{P}}), \nu_1^V) + c_2^V(\mathbf{y}_{\mathcal{P}})\phi(y; \mu_2^V(\mathbf{y}_{\mathcal{P}}), \nu_2^V(\mathbf{y}_{\mathcal{P}}))$$

is a mixture of exactly two distinct Gaussian distributions with means polynomial in  $\mathbf{y}_{\mathcal{P}}$ , one with an absolute constant inverse variance parameter and the other polynomial in  $\mathbf{y}_{\mathcal{P}}$ .

If  $\mathcal{P} = \emptyset$ , then two-Gaussian type w.r.t.  $(W, U, \emptyset)$  requires that both  $\mathbb{P}(Y_W | Y_W \neq 0)$  and  $\mathbb{P}(Y_U | Y_U \neq 0)$  are mixtures of exactly two distinct univariate Gaussian distributions with constant parameters, respectively.

We next recall the following observation that appears as Proposition 29(ii) in Peters et al. (2014); see Sections 1.8 and 15.3.2 of Maathuis et al. (2019) for definitions of the Markov property and faithfulness.

**Proposition 31.** *Suppose the distribution of  $\mathbf{Y}$  is Markov and faithful with respect to two distinct Markov equivalent graphs  $\mathcal{G}$  and  $\mathcal{G}'$ . Then, there must exist nodes  $W$  and  $U$  such that  $W \rightarrow U$  in  $\mathcal{G}$  and  $U \rightarrow W$  in  $\mathcal{G}'$ , while  $\mathcal{P} \equiv \text{pa}_{\mathcal{G}}(U) \setminus \{W\} = \text{pa}_{\mathcal{G}'}(W) \setminus \{U\}$ .*

**Remark 4.** *Proposition 31 is at the heart of many proofs of DAG identifiability, which combine it with suitable probabilistic conditioning to reduce the comparison of two DAG models to bivariate problems involving the two graphs  $W \rightarrow U$  and  $W \leftarrow U$ . However, in our setting, a key new challenge arises because the form of the Hurdle conditionals precludes us from applying conditioning to form sets of bivariate distributions that are of the considered Hurdle type. Indeed, conditioning on descendants of the considered variables (i.e., other variables that in the graph can be reached along directed paths) generally gives conditional distributions that are no longer of the Hurdle type used in the definition of our model class.*

We claim that the intersection of sets of joint distributions represented by two distinct Markov equivalent  $\mathcal{G}$  and  $\mathcal{G}'$  must be a subset of 2-Gaussian type distributions with respect to a triplet  $(W, U, \mathcal{P})$  obtained from Proposition 31.

**Theorem 32** (General Identifiability). *Let  $\mathbf{Y}$ ,  $\mathcal{G}$ ,  $\mathcal{G}'$ ,  $W$ ,  $U$ ,  $\mathcal{P}$  be as in Proposition 31. Let  $\mathbf{Y}$  have a  $\lambda^m$ -density that factors w.r.t. both graphs  $\mathcal{G}$  and  $\mathcal{G}'$ . For each  $\mathcal{H} = \mathcal{G}, \mathcal{G}'$ , let the*

node conditionals in the factorization be Hurdle conditionals with the parameters  $(\alpha_V^{\mathcal{H}})_{V \in \mathcal{V}}$  and  $(\beta_V^{\mathcal{H}})_{V \in \mathcal{V}}$  from (4.4), or equivalently  $(p_V^{\mathcal{H}})_{V \in \mathcal{V}}$  and  $(\mu_V^{\mathcal{H}})_{V \in \mathcal{V}}$  from (4.5), that are Hurdle polynomials of the form (4.8), where for  $(V, T, \mathcal{H}) = (U, W, \mathcal{G})$  and  $(V, T, \mathcal{H}) = (W, U, \mathcal{G}')$  it holds that

(i)  $\beta_V^{\mathcal{H}}(y_T, \mathbf{y}_{\mathcal{P}})$  (or  $\mu_V^{\mathcal{H}}(y_T, \mathbf{y}_{\mathcal{P}})$ ) depends on at least one of  $\mathbb{1}_{y_T}$  and  $y_T$ , or

(ii)  $\alpha_V^{\mathcal{H}}(y_T, \mathbf{y}_{\mathcal{P}})$  (or  $p_V^{\mathcal{H}}(y_T, \mathbf{y}_{\mathcal{P}})$ ) depends on the value of  $y_T$  (and maybe additionally on  $\mathbb{1}_{y_T}$ ).

Then the distribution of  $\mathbf{Y}$  must be of two-Gaussian type w.r.t.  $(W, U, \mathcal{P})$ . In this case we also say the distribution is of two-Gaussian type w.r.t.  $\mathcal{G}$  and  $\mathcal{G}'$ .

Note that the assumption of faithfulness in Proposition 31 implies that we have (i) or (ii) or a condition (iii) that states that  $\alpha_V^{\mathcal{H}}(y_T, \mathbf{y}_{\mathcal{P}})$  (or  $p_V^{\mathcal{H}}(y_T, \mathbf{y}_{\mathcal{P}})$ ) depends on  $\mathbb{1}_{y_T}$  only and  $\beta_V^{\mathcal{H}}(y_T, \mathbf{y}_{\mathcal{P}})$  (or  $\mu_V^{\mathcal{H}}(y_T, \mathbf{y}_{\mathcal{P}})$ ) is constant in  $y_T$ . It is case (iii) that we rule out in our assumption of Theorem 32.

The result is proved in Appendix C. It is easy to show that the result also holds if we make modifications such as restricting the maximum degree of the polynomial or excluding interactions between the discrete and continuous components.

In the two- and three-dimensional cases (i.e.,  $m = 2, 3$ ) we show in Appendix C that there does not exist a joint distribution for  $\mathbf{Y}$  that is of two-Gaussian type with respect to two distinct Markov equivalent graphs. We thus have the following result on full identifiability for graphs with two or three nodes.

**Corollary 33** (Identifiability in two and three dimensions). *If  $|\mathcal{V}| \leq 3$ , i.e., in a binary/triary setting, there does not exist a joint distribution that is of two-Gaussian type w.r.t. two distinct Markov equivalent DAGs  $\mathcal{G}$  and  $\mathcal{G}'$ . Thus, strong identifiability is guaranteed as in Theorem 30, meaning that the sets of Markov and faithful distributions associated to  $\mathcal{G}$  and  $\mathcal{G}'$  must be disjoint.*

Theorem 30 and Corollary 33 state that the DAGs are perfectly identifiable from the distributions if  $m = 2, 3$  or if we assume the Hurdle polynomials to be *strong*; Theorem 32 claims that without assuming *strong* Hurdle polynomials, the distributions for  $m > 3$  from which the graph is not identifiable must be a subset of the *two-Gaussian type* distributions. We conjecture that in general, with  $m > 3$ , the set of two-Gaussian type distributions with respect to any two graphs is an empty set.

In Section 4.5.1 we show scatter plots of simulated data that give some indication of how Markov equivalent graphs may be differentiated under our models.

#### 4.4 Estimation of DAGs from Zero-Inflated Data

Suppose now that we are given an i.i.d. sample  $\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)}$  comprised of  $m$ -variate observations. The log-likelihood function  $\ell$  of any DAG model can be decomposed into the sum of conditional (or nodewise) log-likelihood functions  $\ell^V$  for the  $V$ -th variable conditional on its parent variables. Let  $y_V^{(1)}, \dots, y_V^{(n)}$  be the  $n$  observations of the  $V$ -th variable. For the canonical  $(\alpha, \beta, k)$ -parametrization from (4.4), the nodewise log-likelihood function is

$$\begin{aligned} \ell^V(\alpha_V, \beta_V, k_V | \mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)}) &= \sum_{i=1}^n \left( \alpha_V \left( \mathbf{y}_{\text{pa}(V)}^{(i)} \right) \mathbf{1}_{y_V^{(i)}} + \beta_V \left( \mathbf{y}_{\text{pa}(V)}^{(i)} \right) y_V^{(i)} - k_V y_V^{(i)2} / 2 \right. \\ &\quad \left. - \log \left[ \sqrt{2\pi/k_V} \exp \left\{ \alpha_V \left( \mathbf{y}_{\text{pa}(V)}^{(i)} \right) + \beta_V^2 \left( \mathbf{y}_{\text{pa}(V)}^{(i)} \right) / (2k_V) \right\} + 1 \right] \right); \end{aligned}$$

for the moment  $(p, \mu, \sigma^2)$ -parametrization from (4.5) it is

$$\begin{aligned} \ell^V(p_V, \mu_V, \sigma_V^2 | \mathbf{y}^{(1)}, \dots, \mathbf{y}^{(n)}) &= \sum_{i: y_V^{(i)}=0} \log \left\{ 1 - p_V \left( \mathbf{y}_{\text{pa}(V)}^{(i)} \right) \right\} \\ &\quad + \sum_{i: y_V^{(i)} \neq 0} \left[ \log p_V \left( \mathbf{y}_{\text{pa}(V)}^{(i)} \right) - \frac{1}{2} \log(2\pi\sigma_V^2) - \left\{ y_V^{(i)} - \mu_V \left( \mathbf{y}_{\text{pa}(V)}^{(i)} \right) \right\}^2 / (2\sigma_V^2) \right]. \end{aligned}$$

In the latter case, we see the sum of the log-likelihood functions from the logistic regression model for  $p_V$  and the linear regression for  $\mu_V$  restricted to the observations with  $y_V^{(i)} \neq 0$ . Here we recall that the parameters  $\alpha_V, \beta_V, p_V, \mu_V$  are themselves polynomials in  $\mathbf{y}_{\text{pa}(V)}$  and

their indicators, and we are using them as a shorthand notation on the left-hand sides where we really mean  $\ell^V$  as a function of the parameters (i.e., coefficients) in those polynomials.

#### 4.4.1 Fitting Hurdle Conditionals

Estimation of the graphical models amounts to fitting the conditional distribution of one node given a set of others. For the canonical  $(\alpha, \beta, k)$ -parametrization, the log-likelihood function is convex in  $\alpha_V$ ,  $\beta_V$  and  $k_V$ . Moreover,  $\alpha_V$  and  $\beta_V$  are linear in the polynomial coefficients. Therefore, the log-likelihood is convex in the coefficients to estimate and can be maximized by standard methods; e.g., coordinate descent. Estimation for the moment  $(p, \mu, \sigma^2)$ -parametrization (4.5), on the other hand, can be easily solved by separately fitting a logistic regression to  $p_V$  and a linear regression to  $\mu_V$ . Recall again that the two parametrizations, canonical and moment, are equivalent when assuming a full polynomial model, i.e., when the degree and structure of the polynomials is unrestricted. However, when restricting, for instance, the degree the two parametrizations yield different models.

The  $(\alpha, \beta, k)$ -parametrization with linear Hurdle polynomials (i.e., degree 1) is interesting as it naturally comes from conditional distributions of the joint distribution defined for undirected graphical models in McDavid et al. (2019). However, at least for higher degree, the  $(p, \mu, \sigma^2)$ -parametrization may be more intuitive and useful in practice as it leads to a decomposition into a logistic regression and a linear regression. This decomposition enables one to use optimized standard regression solvers for model fitting. The  $(p, \mu, \sigma^2)$ -parametrization also makes it easy to apply available routines to incorporate regularization on the coefficients/parameters into our loss, which is helpful when the number of samples is small compared to the number of parameters. Such higher dimensionality of the models arises in particular when assuming a higher degree for the Hurdle polynomials. The regularization is automatically applied in the implementation in our R package `ZiDAG` available on [GitHub](#).

For estimation of our models, we assume a highest degree of the Hurdle polynomials. To select the degree from data we adopt the Bayesian information criterion (BIC). This

functionality is incorporated in ZiDAG.

#### 4.4.2 Graph Search

For estimation of the DAG underlying the graphical model, we mainly consider two state-of-the-art methods: (A) exhaustive score-based search and (B) greedy search. Both methods rely on a model score which we take to be the BIC defined as  $\nu \log n - 2\ell$ , where  $\nu$  is the total number of parameters in the model,  $n$  is the sample size, and  $\ell$  is the log-likelihood as introduced in the beginning of Section 4.4.

- (A) **Exhaustive search:** Optimizing the BIC over the set of all DAGs is possible for moderately small  $m$  using the dynamic programming algorithm of Silander and Myllymäki (2006). This approach is justified by the asymptotic consistency of the BIC as well as the identifiability of our model (recall Section 4.3). The experiments of Silander and Myllymäki (2006) suggest that for Gaussian models the search is practical for  $m < 32$ . Estimation of our models is computationally more challenging but exhaustive search is feasible at least for  $m < 16$ .
- (B) **Greedy search:** Instead of optimizing BIC over all DAGs, we may apply a greedy search that iteratively improves BIC by moving to a neighboring DAG that provides the largest improvement. The neighborhood is defined using edge additions, deletions, and reversals; compare Chickering (2003). While Chickering (2003) discusses consistency of graph recovery in terms of equivalence classes, in our case the algorithm determines individual graphs. For faster estimation in sparse settings, we consider restricting the maximum node in-degree (i.e., the maximum number of parents).

**Remark 5.** *We have also experimented with a version of the PC algorithm, which is not easily applicable since it relies on a suitable conditional independence test between two variables given any potential parent set. Indeed, by the nature of our models if the potential parent set is misspecified, the conditional distributions may no longer be Hurdle. Another*

*possible approach starts with greedily estimating the topological ordering of nodes by iteratively picking the node that maximizes the conditional likelihood given nodes already chosen, followed by a variable selection problem using, for example, a Wald test or  $\ell_1$  regularization techniques; this method relies on very subtle features of the distributions. Neither the PC algorithm we designed nor the approach focusing on the topological ordering were competitive in our experiments.*

#### 4.4.3 Stability Selection

In our application to single-cell gene expression data, we seek to also achieve some control of the false discovery rate (FDR). To this end, we apply stability selection in graph estimation. In particular, we take up the approach outlined in Shah and Samworth (2013). We randomly choose  $B = 50$  subsets of the data (each of size  $\lfloor n/2 \rfloor$ ), and obtain  $B$  other sets as their complements of equal size, randomly throwing out one sample if  $n$  is odd. We then estimate the graph using  $\lfloor n/2 \rfloor$  subsamples indexed by each of these  $2B$  sets of equal size, and obtain  $2B$  estimated DAGs. Given the desired FDR, we compute a frequency threshold using the formula from Shah and Samworth (2013, Eqn. (8)) with number of total parameters  $m(m-1)/2$ . We then keep all edges that occur more often than the frequency threshold and produce a graph as our final estimate. In our implementation in ZiDAG, if a graph estimated this way is not acyclic, the user can choose to return it as is, or the function will increase the threshold up to the point where the resulting graph is a DAG, even though the resulting graph might be empty in extreme cases.

## 4.5 Numerical Experiments

In this section we present numerical experiments for exact DAG recovery using simulated zero-inflated conditional Gaussian data. The main goal is to verify identifiability using exhaustive search, and examine how accurately greedy search can recover the true graph.

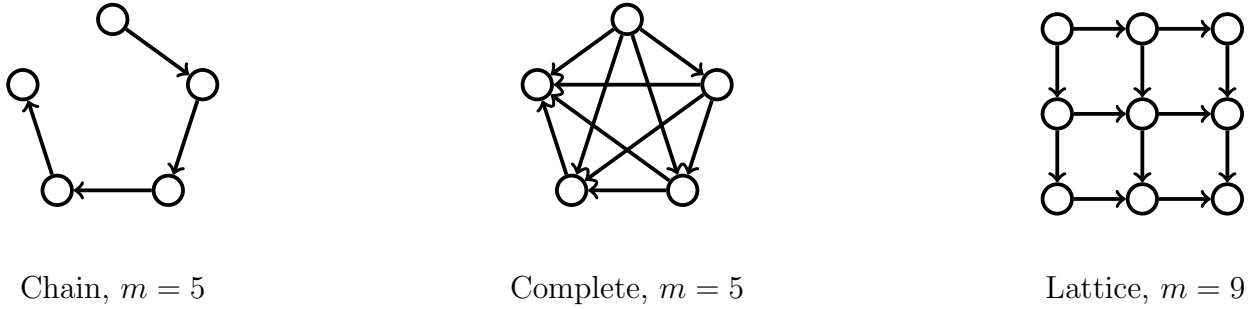


Figure 4.2: Graph structures used in our experiments.

#### 4.5.1 True Underlying DAGs and Distributions

We consider three DAG structures, i) chain graph with  $m = 5$ , ii) complete graph with  $m = 5$ , iii) lattice graph with  $m = 9$ , as illustrated in Figure 4.2. We keep  $m$  small for tractability of repeated application of Silander and Myllymäki (2006) with quadratic Hurdle polynomials. In practice, we suggest applying the directed graph models to the connected components inferred from undirected graphs, estimated using the joint Hurdle distribution (4.2) as in McDavid et al. (2019). This often results in considerably smaller sizes  $m$ . In particular, sizes of  $m$  in this section are similar to the largest component in our data analysis in Section 4.6.

For each structure, we consider true generating conditional distributions using the following parametrizations: a)  $(\alpha, \beta, k)$ -(canonical) parametrization with *linear* Hurdle polynomials, b)  $(p, \mu, \sigma^2)$ -(moment) parametrization with *linear* Hurdle polynomials, and c)  $(p, \mu, \sigma^2)$ -(moment) parametrization with *quadratic* Hurdle polynomials. We note that the distributions represented by c) is a superset of those by a) and b). By (4.6), distributions represented by a) and b) are disjoint because  $\log(p/(1-p))$  is a weighted sum of  $\alpha$  and  $\beta^2$ .

Recall the definition of Hurdle conditionals in (4.4) and (4.5) in Section 4.2.2. In our experiments, whenever  $\text{pa}(V) = \emptyset$ , we generate  $y_V \sim f_0$  such that  $f_0(x) = \frac{1}{2}(1 - \mathbf{1}_x) + \frac{1}{2}\phi(x; 0, 1)$ , where  $\phi$  is the standard normal density. Otherwise, for parametrization a),

we use Hurdle conditionals with parameters  $k_V = 1$ ,  $\alpha_V(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} \mathbb{1}_{y_U} + y_U$  and  $\beta_V(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} (\mathbb{1}_{y_U} - y_U)$ ; similarly for parametrization b) we take  $\sigma_V^2 = 1$ ,  $\log \frac{p_V}{1-p_V}(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} (\mathbb{1}_{y_U} + y_U)$  and  $\mu_V(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} (\mathbb{1}_{y_U} - y_U)$ ; finally, for parametrization c) we take  $\sigma_V^2 = 1$  and

$$\log \frac{p_V}{1-p_V}(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} \left( \mathbb{1}_{y_U} + y_U + \frac{y_U^2}{10} \right) + \frac{1}{10} \sum_{\substack{U, W \in \text{pa}(V) \\ U \neq W}} (\mathbb{1}_{y_U} + y_U)(\mathbb{1}_{y_W} + y_W), \text{ and}$$

$$\mu_V(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} \left( \mathbb{1}_{y_U} - y_U - \frac{y_U^2}{10} \right) + \frac{1}{10} \sum_{\substack{U, W \in \text{pa}(V) \\ U \neq W}} (\mathbb{1}_{y_U} \mathbb{1}_{y_W} - y_U \mathbb{1}_{y_W} - y_W \mathbb{1}_{y_U} - y_U y_W).$$

We then normalize the coefficients in the above expressions  $(\pm 1, \pm 1/10)$  such that  $\alpha_V$ ,  $\beta_V$ ,  $\log p_V/(1-p_V)$  and  $\mu_V$  have means 0 and 1, respectively, across the samples. This normalization ensures that the marginal probability of being nonzero, the marginal mean, and the marginal variance for each node are stabilized, in order to show that the DAGs are truly recovered based on the conditional dependency structure instead of additional signals from these marginal quantities. In fact, in the generated samples the marginal probability is about 0.5 and the marginal mean is about 0 for all nodes, and the marginal variance for the nonzero part only is about the same for all except the source node.

In Figure 4.3, we present pairwise scatter plots of one instance of data generated with the chain graph (upper row) and the complete graph (lower row), respectively, both with  $(p, \mu, k)$ -linear parametrization. Since the true topological ordering is  $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5$ , for clarity we exclude the source and sink nodes (1 and 5) and only include nodes 2, 3 and 4. Plots on the left are plotted in the order 2, 3, 4 and those on the right are reversed. In the histograms on the diagonals we only plot the continuous part.

The scatter plots indicate a slight difference in the respective marginal distributions of nodes 2 and 4 conditioned on node 3 being 0 (and vice versa). This difference intuitively explains how the orientation  $2 \rightarrow 3 \rightarrow 4$  versus  $4 \rightarrow 3 \rightarrow 2$  can be identified. Note that other than this difference, the marginal statistics for the three nodes are indistinguishable and there is little noticeable difference between plots on the left and on the right.

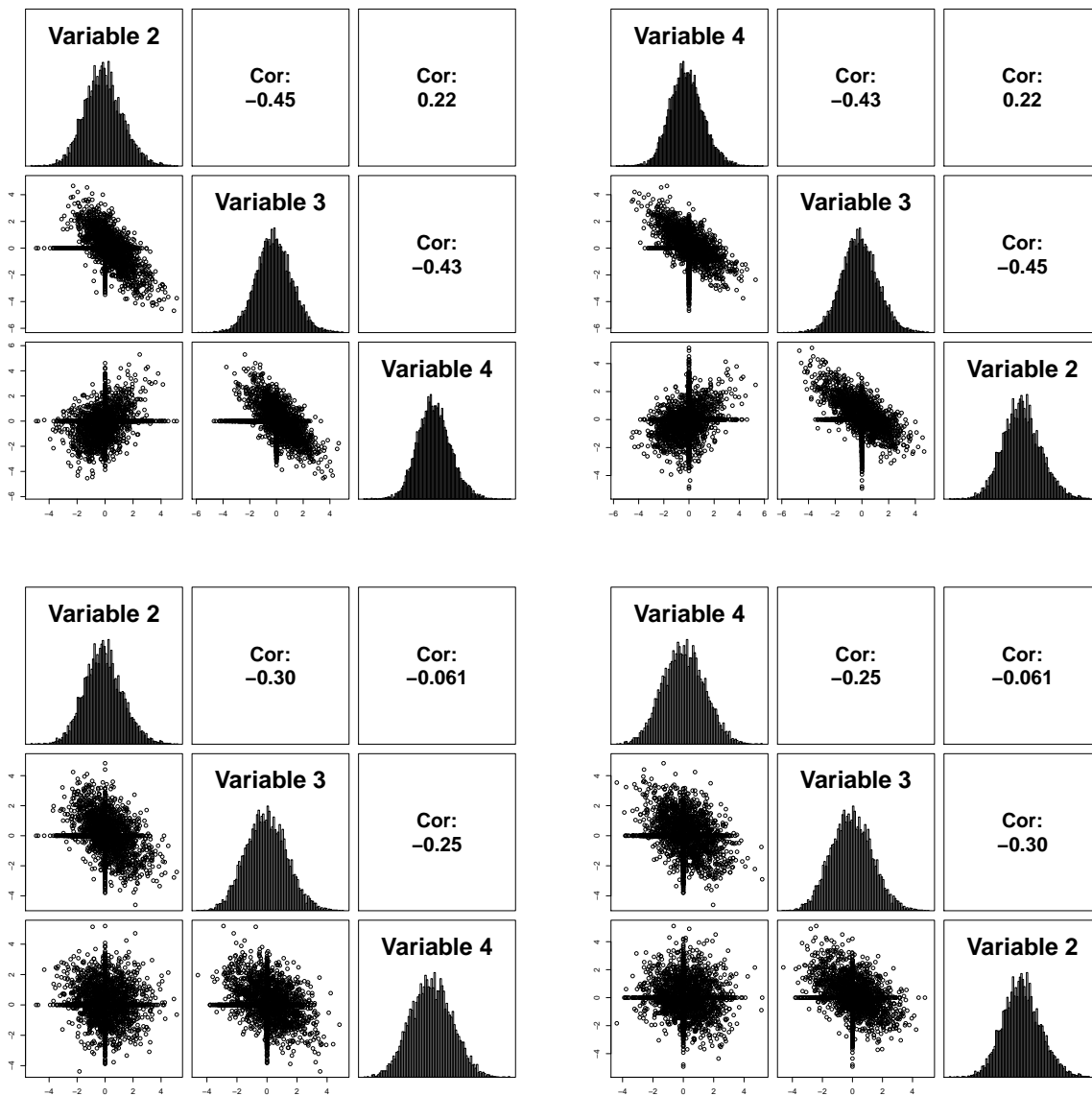


Figure 4.3: Pairwise scatter plots of zero-inflated data generated using chain graphs (upper row) and complete graphs (lower row), both with topological ordering  $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5$ ; only nodes 2, 3 and 4 are plotted. Plots on the left are plotted in the order 2, 3, 4, and 4, 3, 2 on the right. Only the continuous part is plotted in the histograms on the diagonals. There is little noticeable difference between the histograms and scatter plots when we reverse the graph order, yet our model can still determine the correct topological ordering.

### 4.5.2 Estimation

We use the two graph estimation methods described in Section 4.4.2, namely our self-implemented greedy search (GDS) (Chickering, 2003) with BIC score, and exhaustive search with dynamic programming (Silander and Myllymäki, 2006). Details on fitting the Hurdle conditionals themselves were presented in Section 4.4.1.

In our simulation, we aim to assess the performance of the different estimation procedures for correctly specific and misspecified parametrizations. To this end, for each combination of true DAG and true data generating parametrization— $(\alpha, \beta, k)$ -linear and  $(p, \mu, \sigma^2)$ -linear and quadratic—we estimate the DAG using all three parametrizations for generating data. For simplicity and given that the simulation results are presented over  $B = 100$  iterations, stability selection is not used in these experiments.

### 4.5.3 Results on Graph Recovery

Results are shown in Figures 4.4–4.6. Each figure has one true underlying DAG from those shown in Figure 4.2. In all figures, each row indicates one choice of true data generating parametrization— $(\alpha, \beta, k)$ -linear, and  $(p, \mu, \sigma^2)$ -linear and quadratic—and each column shows the results using each estimating parametrization. Thus, plots on the diagonal (with bold titles) correspond to correct parametrizations, where the estimating parametrization agrees with the truth. Off-diagonal plots, in contrast, correspond to cases where the model parametrization is misspecified.

The solid lines with ‘×’ points show the success rates (exact recovery) out of 100 trials versus  $n$  for greedy search (GDS), and the solid lines with ‘o’ points represent exhaustive search (Silander). The gray dotted lines signify success rates measured by recovery of the equivalence class.

Since exhaustive search compares all possible DAGs for  $m$  nodes, for  $n$  large enough it provides an indicator of identifiability. Indeed, the results indicate that in all settings, exhaustive search with correct parametrization almost always identifies the exact DAG for

large  $n$ . In fact, since the  $(p, \mu, \sigma^2)$ -quadratic parametrization covers the other two, in all cases the graphs can be perfectly recovered using the quadratic estimating parametrization. In contrast, when the underlying truth is quadratic, the graph may not be easily identified from estimates that use the other two parametrizations. This is especially the case for the lattice graph. Comparing the linear parametrizations themselves,  $(p, \mu, \sigma^2)$  seems less prone to model misspecification and has the advantage of faster estimation with the help of standard softwares for logistic and linear regressions.

Overall, our simulation studies confirm the identifiability theory (Theorem 30). In particular, our experiments indicate that exhaustive search performs well. They also indicate that GDS works reasonably well for sparse graphs but may require larger samples for recovering the structure of complete, or very dense, graphs. While exhaustive search often succeeds with high probability even with small samples, it may not be scalable for large  $m$ . In such cases, the greedy and faster GDS method, which shows promising results, provides a viable alternative. Utilizing the stability selection method of Section 4.4.3 can further improve the GDS results.

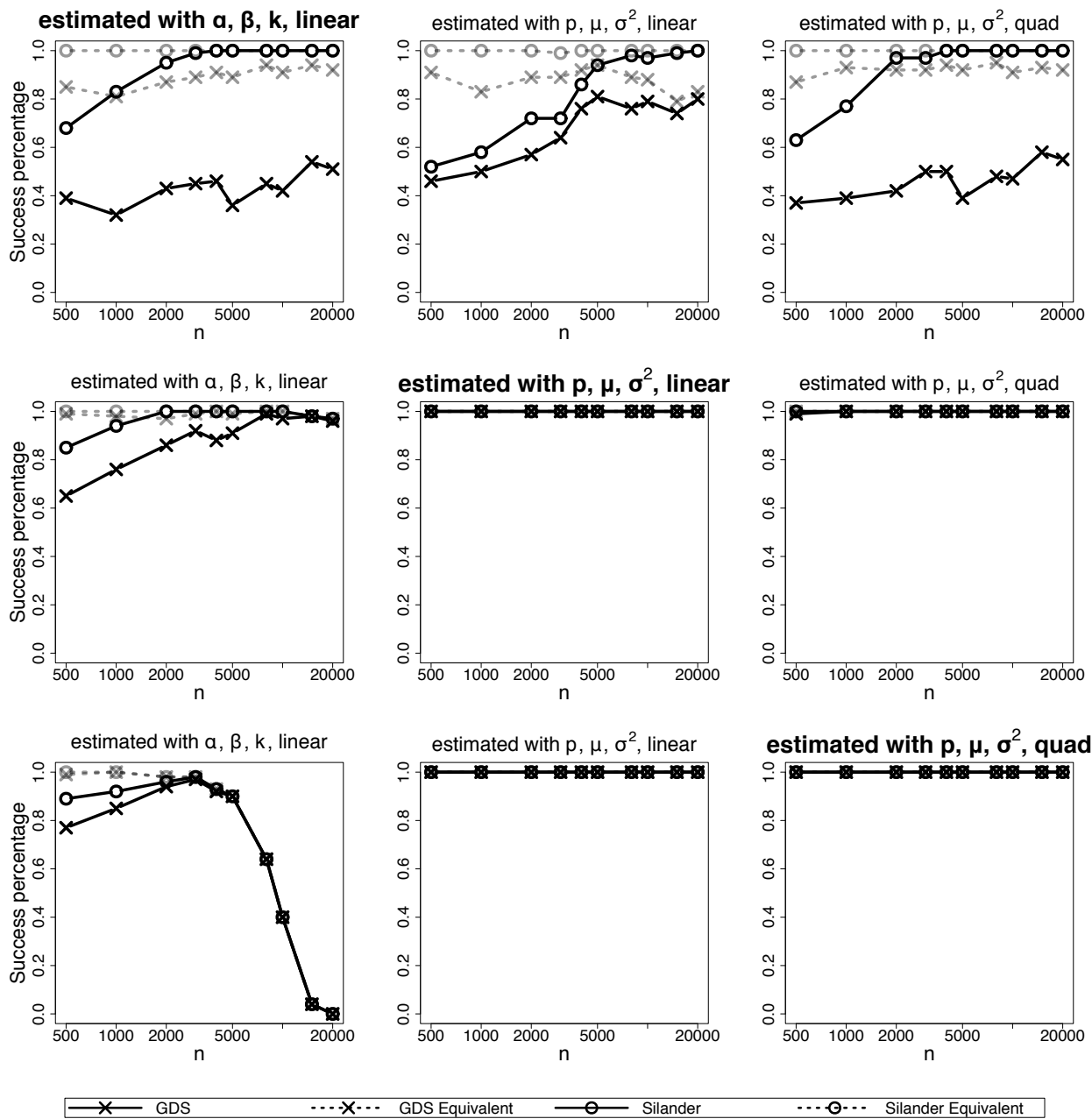


Figure 4.4: Chain graph,  $m = 5$ . Each row corresponds to a different generating parametrization, and each column a different estimating parametrization. Generating and estimating parametrizations agree on the diagonal. Solid 'x': success rates of exact DAG recovery for greedy search; solid 'o': success rates of exact DAG recovery for exhaustive search; gray dotted lines: success rates for recovery of equivalence class.

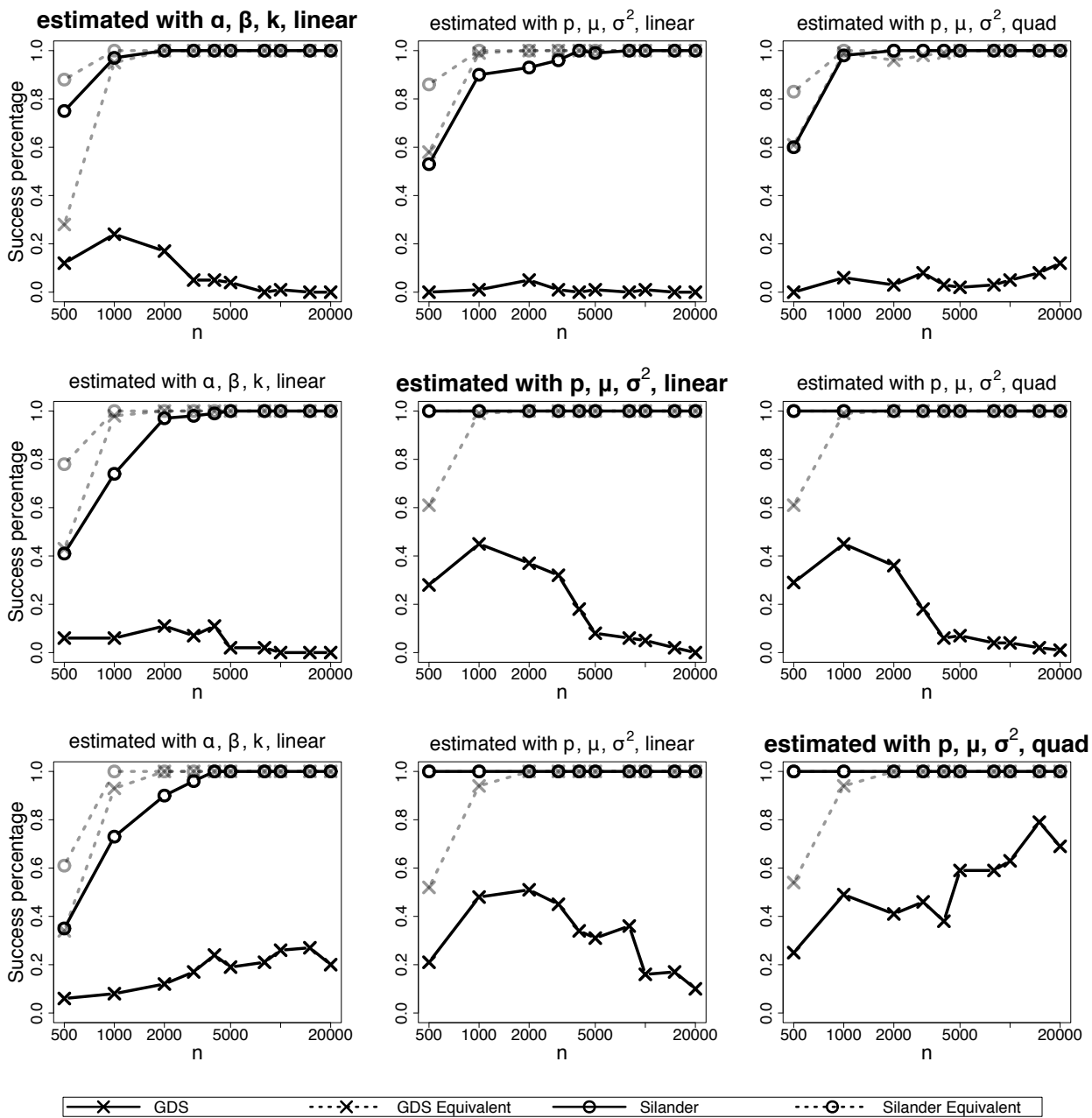


Figure 4.5: Complete graph,  $m = 5$ .

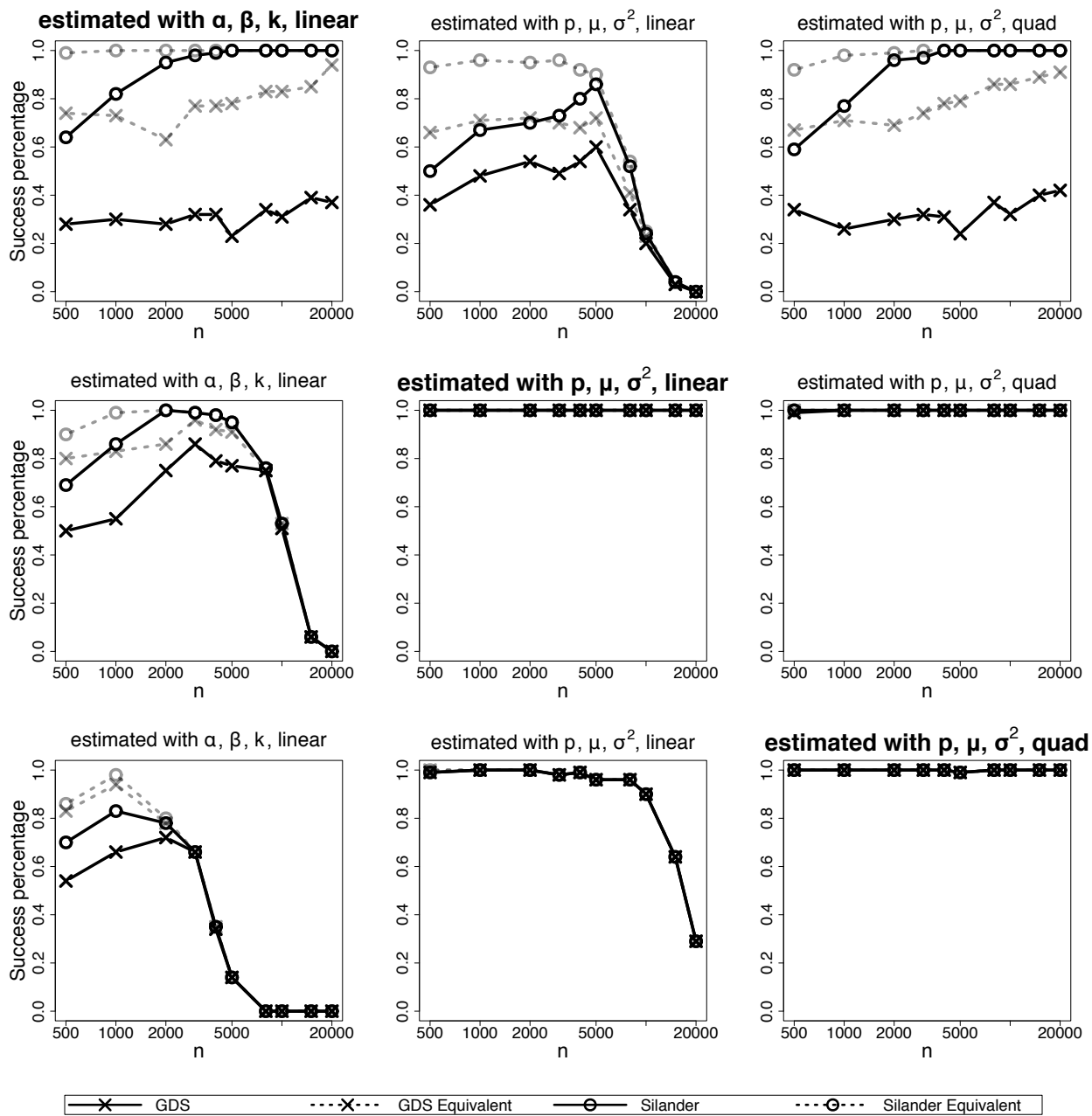


Figure 4.6: Lattice graph,  $m = 9$ .

## 4.6 T Helper Cell Data

In this section we present the results of applying our model to a T helper cell expression dataset. Specifically, the dataset is considered in McDavid et al. (2019) and contains both single cell and 10-cell expression measurements for T helper cells for 80 genes in eight healthy donors. We use all 1951 single cell measurements for these donors (a superset of the 465 measurements in McDavid et al. (2019)) to ensure we have a large enough sample size to produce reliable estimates. In particular, McDavid et al. (2019) consider only the T-follicular (CXCR5<sup>+</sup>PD1<sup>+</sup>) cells that produce high levels of proteins CXCR5 and PD1, while we do not make this restriction. Instead, we add the indicators of CXCR5<sup>+/-</sup> and PD1<sup>+/-</sup> as regressors when fitting the conditional distributions. Following McDavid et al. (2019), we choose the 61 genes that have at least 5% zero and 5% nonzero values.

While the measurements are all nonnegative, the minimum, mean, and standard deviation of the nonzero values in the dataset are 7.89, 18.53, and 1.91, respectively. We thus assume zero-inflated conditional Gaussianity without considering the effect of truncation from below at 0.

To estimate the DAG structure, we first use the procedure of McDavid et al. (2019) to identify the connected components in an estimated undirected graph with the same sparsity as the graph therein. We then estimate the directed edges in each connected component using our method. This procedure is justified by the fact that theoretically the connected components for the underlying true undirected and directed graphs coincide. Thus, we generally recommend this strategy in practice, as the connected components can be much more efficiently obtained from the undirected graph.

We use the  $(p, \mu, \sigma^2)$ -parametrization as it is more flexible than the  $(\alpha, \beta, k)$ , and extra fixed covariates and controlling factors can be easily added, since fitting the conditionals only involves linear and logistic regressions. As discussed in Section 4.5, the  $(p, \mu, \sigma^2)$  is also more robust than  $(\alpha, \beta, k)$ . We use polynomials up to degree three and data-adaptively choose the optimal degree by BIC.

To estimate the DAG, we use the greedy search (GDS) algorithm, which showed promising performance in the simulations of Section 4.5. We also use the stability selection procedure of Section 4.4.3, with the goal of controlling the FDR at 10% for each connected component. For smaller connected components, if controlling the FDR at 10% is not possible, we pick the sparsest graph that maximally maintains the connectivity. Finally, we restrict the node in-degrees to five, in order to both speed up estimation and to constrain the search space. This constraint is motivated by the fact that in gene regulatory networks, each gene is only expected to be regulated by a small number of other genes (Albert, 2005). In contrast, since genetic networks often involve hub genes that regulate many others, we do not restrict the out-degree.

The estimated undirected graph using the procedure of McDavid et al. (2019) is plotted in the upper half of Figure 4.7, with edge width and saturation representing the edge strength. In the lower half of the figure, we plot the estimated DAG using our method; the estimate with stability selection and FDR control is shown on the left and the one without stability selection is on the right. Examining the estimates, we find that CD3E is a hub node with degree five in both estimated DAGs, while it has four neighbors in the estimated undirected graph. On the other hand, in the estimate with stability selection, three genes in the largest connected components, namely CD28, JAK1 and STATS5B, are isolated. This is reasonable as they are each associated with only one weak edge in the undirected graph. Moreover, the undirected and directed graph estimates have different thresholds for determining whether an edge is present. For the other nodes, the estimated DAG structures are very similar to the undirected graph estimate.

#### **4.7 Discussion**

Motivated by the recent advent of single-cell RNA-seq data, in this paper we develop new methods for learning DAGs from zero-inflated data. Our procedures take advantage of two key features of single-cell RNA-seq data, namely, the zero-inflated nature of the data, and the large number of observations from individual samples.

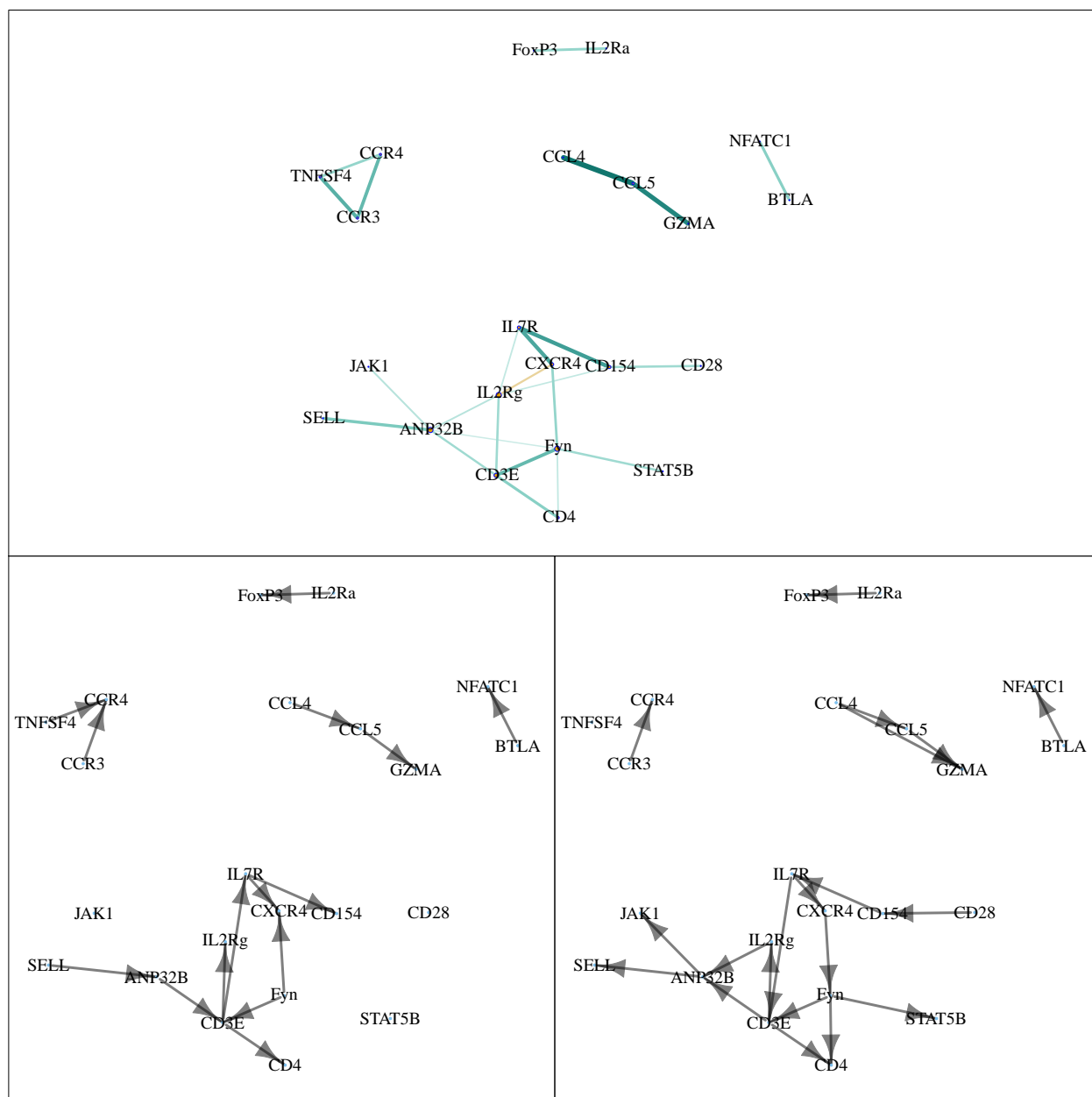


Figure 4.7: Upper: Graph estimated for T helper cell data using undirected zero-inflated graphical models; similar to Figure 6 from McDavid et al. (2019). Lower: Directed graph estimated for T helper cell data with (left) and without (right) stability selection and FDR control.

Our key contribution is establishing identifiability of DAGs from observational zero-inflated data. Specifically, we prove that the exact DAG can be recovered from the joint distribution under reasonable assumptions. We also show that in the most general case, the distributions from which the DAGs are not identifiable only form a small subset, which we prove to be empty in the bivariate and trivariate cases. While our proof uses a very general result on DAGs from Peters et al. (2014) as its first step, our models do not fit into the framework in that paper; we thus take a different approach that considers the zero-inflation and polynomial structures directly.

Our approach is based on factorizing the joint distribution into zero-inflated conditional Gaussian distributions with parameters polynomial in the parents and their indicators of having nonzero values. We present models in terms of two parametrizations, one called  $(\alpha, \beta, k)$  that is linked to the undirected graphs studied in McDavid et al. (2019), and the other called  $(p, \mu, \sigma^2)$  that directly models the conditional moments. Both approaches have computational appeal. In particular, the  $(\alpha, \beta, k)$ -parametrization leads to convex loss functions in the parameters to be estimated, while the  $(p, \mu, \sigma^2)$ -parametrization offers the additional benefit of allowing one to utilize standard software for logistic and linear regression. We combine these models with two state-of-the-art estimation procedures, namely greedy DAG search (GDS) and exhaustive search with dynamic programming. We also validate our identifiability theory using extensive numerical studies. These experiments indicate that the exhaustive search algorithm is effective in correctly identifying DAGs with small number of nodes. For moderate to large DAGs, the GDS algorithm offers a reasonable alternative, with performance comparable to the exhaustive search when the sample size is large enough.

Our work opens the door to multiple future research directions and extensions. The first is to prove our conjecture that the sets of distributions from which the DAG is not identifiable are empty also for graphs with more than 3 nodes. The second direction of future research is proving the consistency and investigating finite sample properties of the proposed estimation procedures. Finally, an interesting extension of our model would be to consider zero-inflated distributions under a truncation to the nonnegative orthant  $\mathbb{R}_+^m$ , which

would be of interest for nonnegative *omics* data. The main challenge in this case would be the normalizing constant as a function of the parents in the conditional distributions, since it would not have a closed-form expression. While this may be resolved by generalizing the *score matching* loss (Hyvärinen, 2005, 2007; Lyu, 2009; Yu et al., 2019b) to data of mixed type, the additional difficulty would lie in proving identifiability and addressing estimation from observational data.

Chapter 5

**DISCUSSION AND FUTURE WORK**

In this dissertation, we discuss both undirected and directed graphical models and their estimation. In Chapter 2 we propose the generalized score matching estimator for estimating continuous distributions supported on  $\mathbb{R}_+^m$  with applications to pairwise interaction power models; the estimator enjoys the advantages of avoiding calculations of normalizing constants as well as being quadratic in the canonical parameters of exponential families, while significantly outperforming the original estimator proposed in Hyvärinen (2007). The generalization is done by replacing the  $\nabla \log p(\mathbf{x}) \otimes \mathbf{x}$  term in Hyvärinen (2007) with  $\nabla \log p(\mathbf{x}) \otimes \mathbf{h}^{1/2}(\mathbf{x})$  for some general  $\mathbf{h} : \mathbb{R}_+^m \rightarrow \mathbb{R}_+^m$ . In Chapter 3 we extend the estimator to distributions supported on more general domains  $\mathcal{D}$ , which should cover most applications in practice, but only involves little extra work from calculating the *truncated componentwise distance*  $\varphi$  of the data to the boundary of  $\mathcal{D}$ , truncated above by prespecified constants  $\mathbf{C}$ . Finally, in Chapter 4, inspired by the recent technology of single-cell RNA-seq data, we propose a class of directed graphical models for zero-inflated data, for which we establish identifiability of directed acyclic graphs (DAGs) under reasonable assumptions, and present graph estimation methods that prove to be successful in practice.

In Chapters 2 and 3, for high-dimensional exponential family graphical models, following Lin et al. (2016) we add an  $\ell_1$  penalty, and in addition we introduce an  $\ell_2$ -type penalty using *diagonal multipliers* to overcome the issue in Lin et al. (2016) that the regularized loss may be unbounded if the  $\ell_1$  penalty is small. We consider pairwise interaction power models, which are exponential families in which the log density is the sum of pairwise interactions between  $\mathbf{x}^a$  (log  $\mathbf{x}$  when  $a = 0$ ) plus linearly weighted effects  $\mathbf{x}^b$  (log  $\mathbf{x}$  when  $b = 0$ ). This covers a wide range of models, such as truncated normal distributions with  $a = b = 1$ , exponential square root models (Inouye et al., 2016) with  $a = b = 1/2$ , a class of multivariate gamma distributions  $a = 1/2$  and  $b = 0$ , as well as  $A^{m-1}$  models (Aitchison, 1985) with  $a = b = 0$  and  $\mathcal{D}$  being the  $m$ -dimensional simplex. We show that for Gaussian graphical models on  $\mathcal{D} \equiv \mathbb{R}_+^m$ , we can reduce the sample size required for consistency from  $\Omega(d^2 \log^8 m)$  in Lin et al. (2016) to  $\Omega(d^2 \log m)$ , where  $d$  is the max degree of the graph. For general domains  $\mathcal{D}$ , we have similar theory requiring  $n = \Omega(d^2 \log m)$  for Gaussian graphical models on domains

that are finite disjoint unions of convex sets, as well as general  $a$ - $b$  models on simplex domains assuming  $a > 0$  or on bounded subsets of  $\mathbb{R}_+^m$  with positive Lebesgue measure and  $a \geq 0$ . For other domains  $\mathcal{D}$  we require an additional multiplicative factor that may weakly depend on  $m$ . Through simulation studies, we recommend the choice of  $\mathbf{h}(\mathbf{x}) = (x_1^c, \dots, x_m^c)$  with  $c = \max\{2 - a, 0\}$ , and adaptively choose the truncation points using sample quantiles.

In Chapter 4, we present our DAG models with two different parametrizations, which take advantage of the zero-inflated nature of the data and the large number of observations from individual samples. We prove that under reasonable assumptions, one can recover the exact DAG from the joint distribution, and that in the most general case, the distributions from which the exact DAGs are not identifiable only form a small subset, which is shown to be empty in bivariate and trivariate cases and conjectured to be empty for  $m > 3$ . The identifiability is supported by simulation studies using exhaustive search with BIC score. We use simulated experiments to show that the greedy DAG search (GDS) also works well when the sample size is large enough, a reasonable assumption for single-cell RNA-seq data.

Possible future extensions to our work in Chapters 2 and 3 could include developing a method that adaptively chooses the  $\mathbf{h}$  function along with the truncation points  $\mathbf{C}$  for the truncated distances  $\varphi$ , as well as probing into the dependence of the additional multiplicative constant factor on  $\log m$  that appears in the consistency results for certain domains  $\mathcal{D}$ . For the DAG models in Chapter 4, one may be interested in proving our conjecture above and the consistency and finite sample properties of the proposed estimation procedures, as well as extending our models to truncated zero-inflated distributions, possibly utilizing our work in Chapters 2 and 3.

## BIBLIOGRAPHY

- John Aitchison. A general class of distributions on the simplex. *Journal of the Royal Statistical Society: Series B (Methodological)*, 47(1):136–146, 1985.
- Reka Albert. Scale-free networks in cell biology. *Journal of Cell Science*, 118(21):4947–4957, 2005.
- Murilo P. Almeida and Basilis Gidas. A variational method for estimating the parameters of MRF from complete or incomplete data. *Ann. Appl. Probab.*, 3(1):103–136, 1993.
- Albert-László Barabási and Réka Albert. Emergence of scaling in random networks. *Science*, 286(5439):509–512, 1999.
- Rina Foygel Barber and Mathias Drton. High-dimensional Ising model selection with Bayesian information criteria. *Electron. J. Stat.*, 9(1):567–607, 2015.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration Inequalities*. Oxford University Press, Oxford, 2013.
- Igor Brikun, Deborah Nusskern, Daniel Gillen, Amy Lynn, Daniel Murtagh, John Feczko, William G Nelson, and Diha Freije. A panel of DNA methylation markers reveals extensive methylation in histologically benign prostate biopsy cores from cancer patients. *Biomarker Research*, 2(1):25, 2014.
- Peter Bühlmann and Sara van de Geer. *Statistics for High-dimensional Data: Methods, Theory and Applications*. Springer Science & Business Media, 2011.
- V. V. Buldygin and Yu. V. Kozachenko. *Metric Characterization of Random Variables and Random Processes*, volume 188 of *Translations of Mathematical Monographs*. American

- Mathematical Society, Providence, RI, 2000. Translated from the 1998 Russian original by V. Zaiats.
- Scott L. Carter, Christian M. Brechbühler, Michael Griffin, and Andrew T. Bond. Gene co-expression network topology provides a framework for molecular characterization of cellular state. *Bioinformatics*, 20(14):2242–2250, 2004.
- Jiahua Chen and Zehua Chen. Extended Bayesian information criteria for model selection with large model spaces. *Biometrika*, 95(3):759–771, 2008.
- Shizhe Chen, Daniela M. Witten, and Ali Shojaie. Selection and estimation for mixed graphical models. *Biometrika*, 102(1):47–64, 2015.
- Wenyu Chen, Mathias Drton, and Y. Samuel Wang. On causal discovery with an equal-variance assumption. *Biometrika*, 106(4):973–980, 2019.
- David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of Machine Learning Research (JMLR)*, 3(3):507–554, 2003.
- Adrian Dobra and Alex Lenkoski. Copula Gaussian graphical models and their application to modeling functional disability data. *Ann. Appl. Stat.*, 5(2A):969–993, 2011.
- Mathias Drton and Marloes H. Maathuis. Structure learning in graphical modeling. *Annual Review of Statistics and Its Application*, 4:365–393, 2017.
- Mathias Drton and Michael D. Perlman. Model selection for Gaussian concentration graphs. *Biometrika*, 91(3):591–602, 2004.
- Pan Du, Xiao Zhang, Chiang-Ching Huang, Nadereh Jafari, Warren A Kibbe, Lifang Hou, and Simon M Lin. Comparison of beta-value and m-value methods for quantifying methylation levels by microarray analysis. *BMC bioinformatics*, 11(1):587, 2010.
- Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.*, 96(456):1348–1360, 2001.

- Tom Fawcett. An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8):861–874, 2006.
- Bernd Fellinghauer, Peter Bühlmann, Martin Ryffel, Michael von Rhein, and Jan D. Reinhardt. Stable graphical model estimation with random forests for discrete, continuous, and mixed variables. *Comput. Statist. Data Anal.*, 64:132–152, 2013.
- Rina Foygel and Mathias Drton. Extended Bayesian information criteria for Gaussian graphical models. In *Advances in Neural Information Processing Systems*, pages 604–612, 2010.
- Jerome Friedman, Trevor Hastie, Holger Höfling, and Robert Tibshirani. Pathwise coordinate optimization. *Ann. Appl. Stat.*, 1(2):302–332, 2007.
- Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.
- Asish Ghoshal and Jean Honorio. Learning identifiable Gaussian Bayesian networks in polynomial time and sample complexity. In *Advances in Neural Information Processing Systems*, pages 6457–6466, 2017.
- Jing-Dong J. Han, Nicolas Bertin, Tong Hao, Debra S. Goldberg, Gabriel F. Berriz, Lan V. Zhang, Denis Dupuy, Albertha JM. Walhout, Michael E. Cusick, Frederick P. Roth, et al. Evidence for dynamically organized modularity in the yeast protein–protein interaction network. *Nature*, 430(6995):88, 2004.
- Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *J. Mach. Learn. Res.*, 6:695–709, 2005.
- Aapo Hyvärinen. Some extensions of score matching. *Comput. Statist. Data Anal.*, 51(5):2499–2512, 2007.
- David Inouye, Pradeep Ravikumar, and Inderjit Dhillon. Square root graphical models: Multivariate generalizations of univariate exponential families that permit positive de-

- dependencies. In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 2445–2453, 2016.
- Hawoong Jeong, Sean P. Mason, A-L. Barabási, and Zoltan N. Oltvai. Lethality and centrality in protein networks. *Nature*, 411(6833):41, 2001.
- Markus Kalisch and Peter Bühlmann. Estimating high-dimensional directed acyclic graphs with the pc-algorithm. *Journal of Machine Learning Research*, 8(Mar):613–636, 2007.
- Kshitij Khare, Sang-Yun Oh, and Bala Rajaratnam. A convex pseudolikelihood framework for high dimensional partial correlation estimation with convergence guarantees. *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, 77(4):803–825, 2015.
- Diederik P. Kingma and Yann L. Cun. Regularized estimation of image statistics by score matching. In *Advances in Neural Information Processing Systems*, pages 1126–1134, 2010.
- Urs Köster and Aapo Hyvärinen. A two-layer ICA-like model estimated by score matching. *Artificial Neural Networks–ICANN 2007*, pages 798–807, 2007.
- Steffen L. Lauritzen. *Graphical Models*, volume 17 of *Oxford Statistical Science Series*. The Clarendon Press, Oxford University Press, New York, 1996.
- Olivier Ledoit and Michael Wolf. A well-conditioned estimator for large-dimensional covariance matrices. *J. Multivariate Anal.*, 88(2):365–411, 2004.
- Lina Lin, Mathias Drton, and Ali Shojaie. Estimation of high-dimensional graphical models using regularized score matching. *Electron. J. Stat.*, 10(1):806–854, 2016.
- Han Liu, John Lafferty, and Larry Wasserman. The nonparanormal: semiparametric estimation of high dimensional undirected graphs. *J. Mach. Learn. Res.*, 10:2295–2328, 2009.
- Han Liu, Fang Han, Ming Yuan, John Lafferty, and Larry Wasserman. High-dimensional semiparametric Gaussian copula graphical models. *Ann. Statist.*, 40(4):2293–2326, 2012.

- Song Liu and Takafumi Kanamori. Estimating density models with complex truncation boundaries. *arXiv preprint arXiv:1910.03834*, 2019.
- Weidong Liu and Xi Luo. Fast and adaptive sparse precision matrix estimation in high dimensions. *J. Multivariate Anal.*, 135:153–162, 2015.
- Jun Luo, Thomas Dunn, Charles Ewing, Jurga Sauvageot, Yidong Chen, Jeffrey Trent, and William Isaacs. Gene expression signature of benign prostatic hyperplasia revealed by cDNA microarray analysis. *The Prostate*, 51(3):189–200, 2002.
- Siwei Lyu. Interpretation and generalization of score matching. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pages 359–366. AUAI Press, 2009.
- Marloes Maathuis, Mathias Drton, Steffen Lauritzen, and Martin Wainwright. *Handbook of graphical models*. Chapman & Hall/CRC Handbooks of Modern Statistical Methods. CRC Press, Boca Raton, FL, 2019. ISBN 978-1-4987-8862-5.
- Andrew McDavid, Raphael Gottardo, Noah Simon, and Mathias Drton. Graphical models for zero-inflated single cell gene expression. *The Annals of Applied Statistics*, 13(2):848–873, 2019.
- Nicolai Meinshausen and Peter Bühlmann. High-dimensional graphs and variable selection with the lasso. *Ann. Statist.*, 34(3):1436–1462, 2006.
- Giuseppe Morgia, Mario Falsaperla, Grazia Malaponte, Massimo Madonia, Manuela Indelicato, Salvatore Travali, and Maria Clorinda Mazzarino. Matrix metalloproteinases as diagnostic (MMP-13) and prognostic (MMP-2, MMP-9) markers of prostate cancer. *Urological Research*, 33(1):44–50, 2005.
- Preetam Nandy, Alain Hauser, Marloes H. Maathuis, et al. High-dimensional consistency in score-based and hybrid structure learning. *The Annals of Statistics*, 46(6A):3151–3183, 2018.

Shintaro Narita, Alan So, Susan Ettinger, Norihiro Hayashi, Mototsugu Muramaki, Ladan Fazli, Youngsoo Kim, and Martin E Gleave. GLI2 knockdown using an antisense oligonucleotide induces apoptosis and chemosensitizes cells to paclitaxel in androgen-independent prostate cancer. *Clinical Cancer Research*, 14(18):5769–5777, 2008.

Masashi Okamoto. Distinctness of the eigenvalues of a quadratic form in a multivariate sample. *The Annals of Statistics*, 1:763–765, 1973.

Matthew Parry. Extensive scoring rules. *Electron. J. Stat.*, 10(1):1098–1108, 2016.

Matthew Parry, A. Philip Dawid, and Steffen Lauritzen. Proper local scoring rules. *Ann. Statist.*, 40(1):561–592, 2012.

Jie Peng, Nengfeng Zhou, and Ji Zhu. Partial correlation estimation by joint sparse regression models. *J. Amer. Statist. Assoc.*, 104(486):735–746, 2009.

Jonas Peters and Peter Bühlmann. Identifiability of Gaussian structural equation models with equal error variances. *Biometrika*, 101(1):219–228, 2013.

Jonas Peters, Joris M. Mooij, Dominik Janzing, and Bernhard Schölkopf. Causal discovery with continuous additive noise models. *The Journal of Machine Learning Research*, 15(1):2009–2053, 2014.

Pradeep Ravikumar, Martin J. Wainwright, and John D. Lafferty. High-dimensional Ising model selection using  $\ell_1$ -regularized logistic regression. *Ann. Statist.*, 38(3):1287–1319, 2010.

Pradeep Ravikumar, Martin J. Wainwright, Garvesh Raskutti, and Bin Yu. High-dimensional covariance estimation by minimizing  $\ell_1$ -penalized log-determinant divergence. *Electron. J. Stat.*, 5:935–980, 2011.

Rajen D. Shah and Richard J. Samworth. Variable selection with error control: another

- look at stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(1):55–80, 2013.
- Shohei Shimizu, Patrik O. Hoyer, Aapo Hyvärinen, and Antti Kerminen. A linear non-Gaussian acyclic model for causal discovery. *Journal of Machine Learning Research*, 7 (Oct):2003–2030, 2006.
- Tomi Silander and Petri Myllymäki. A simple approach for finding the globally optimal bayesian network structure. In *Conference on Uncertainty in Artificial Intelligence*, pages 445–452, 2006.
- Peter Spirtes, Clark Glymour, and Richard Scheines. Causation, prediction, and search. adaptive computation and machine learning, 2000.
- Mona Stearns and Mark E. Stearns. Evidence for increased activated metalloproteinase 2 (MMP-2a) expression associated with human prostate cancer progression. *Oncology Research Featuring Preclinical and Clinical Cancer Therapeutics*, 8(2):69–75, 1996.
- G. M. Tallis. The moment generating function of the truncated multi-normal distribution. *J. Roy. Statist. Soc. Ser. B*, 23:223–229, 1961.
- Saravanan Thiyagarajan, Neehar Bhatia, Shannon Reagan-Shaw, Diana Cozma, Andrei Thomas-Tikhonenko, Nihal Ahmad, and Vladimir S Spiegelman. Role of GLI2 transcription factor in growth and tumorigenicity of prostate cells. *Cancer Research*, 67(22):10642–10646, 2007.
- Robert Tibshirani. Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B*, 58(1):267–288, 1996.
- Robert Tibshirani, Jacob Bien, Jerome Friedman, Trevor Hastie, Noah Simon, Jonathan Taylor, and Ryan J. Tibshirani. Strong rules for discarding predictors in lasso-type problems. *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, 74(2):245–266, 2012.

- Ryan J. Tibshirani. The lasso problem and uniqueness. *Electron. J. Stat.*, 7:1456–1490, 2013.
- Dominique Trudel, Yves Fradet, François Meyer, François Harel, and Bernard Têtu. Significance of MMP-2 expression in prostate cancer. *Cancer Research*, 63(23):8511–8515, 2003.
- Hannu Väliäho. Criteria for copositive matrices. *Linear Algebra Appl.*, 81:19–34, 1986.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. In *Compressed Sensing*, pages 210–268. Cambridge Univ. Press, Cambridge, 2012.
- Arend Voorman, Ali Shojaie, and Daniela Witten. Graph estimation with joint additive models. *Biometrika*, 101(1):85–101, 2014.
- Martin J Wainwright. *High-dimensional statistics: A Non-Asymptotic Viewpoint*, volume 48. Cambridge University Press, 2019.
- Y. Samuel Wang and Mathias Drton. High-dimensional causal discovery under non-Gaussianity. *Biometrika*, 107(1):41–59, 2020.
- Susan R Wilson, Sandra Gallagher, Kate Warpeha, and Susan J Hawthorne. Amplification of MMP-2 and MMP-9 production by prostate cancer cell lines via activation of protease-activated receptors. *The Prostate*, 60(2):168–174, 2004.
- Tiancheng Xie, Binbin Dong, Yangye Yan, Guanghui Hu, and Yunfei Xu. Association between MMP-2 expression and prostate cancer: A meta-analysis. *Biomedical Reports*, 4(2):241–245, 2016.
- Eunho Yang, Pradeep Ravikumar, Genevera I. Allen, and Zhandong Liu. Graphical models via univariate exponential family distributions. *J. Mach. Learn. Res.*, 16:3813–3847, 2015.
- Ming Yu, Mladen Kolar, and Varun Gupta. Statistical inference for pairwise graphical models using score matching. In *Advances in Neural Information Processing Systems*, pages 2829–2837, 2016.

- Ming Yu, Varun Gupta, and Mladen Kolar. Simultaneous inference for pairwise graphical models with generalized score matching. *arXiv preprint arXiv:1905.06261*, 2019a.
- Shiqing Yu, Mathias Drton, and Ali Shojaie. Graphical models for non-negative data using generalized score matching. In *International Conference on Artificial Intelligence and Statistics*, pages 1781–1790, 2018.
- Shiqing Yu, Mathias Drton, and Ali Shojaie. Generalized score matching for non-negative data. *Journal of Machine Learning Research*, 20(76):1–70, 2019b.
- Shiqing Yu, Mathias Drton, and Ali Shojaie. Directed graphical models and causal discovery for zero-inflated data. *arXiv preprint arXiv:2004.04150*, 2020.
- Ming Yuan and Yi Lin. Model selection and estimation in the Gaussian graphical model. *Biometrika*, 94(1):19–35, 2007.
- Małgorzata Żak-Szatkowska and Małgorzata Bogdan. Modified versions of the Bayesian information criterion for sparse generalized linear models. *Comput. Statist. Data Anal.*, 55(11):2908–2924, 2011.
- Cun-Hui Zhang. Nearly unbiased variable selection under minimax concave penalty. *Ann. Statist.*, 38(2):894–942, 2010.
- Teng Zhang and Hui Zou. Sparse precision matrix estimation via lasso penalized D-trace loss. *Biometrika*, 101(1):103–120, 2014.
- Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 67(2):301–320, 2005.

Appendix A

**APPENDICES TO CHAPTER 2**

## A.1 Proofs

### A.1.1 Proof of Theorem 2

The following integration by parts lemma is used in the proof of Theorem 2.

**Lemma 34.** *Let  $f, g : \mathbb{R}_+ \rightarrow \mathbb{R}$  be functions that are absolutely continuous in every bounded sub-interval of  $\mathbb{R}_+$ . Then*

$$\lim_{x \nearrow +\infty} f(x)g(x) - \lim_{x \searrow 0^+} f(x)g(x) = \int_0^\infty f(\mathbf{x}) \frac{dg(x)}{dx} dx + \int_0^\infty g(\mathbf{x}) \frac{df(x)}{dx} dx.$$

*Proof.* This is an analog of Lemma 4 from Hyvärinen (2005) that can be proved by integrating the partial derivatives, and follows from the fundamental theorem of calculus for absolutely continuous functions and the product rule. In particular, we work on integrals in bounded  $[0, c]$ , where the product of two absolutely continuous functions in a bounded interval is again absolutely continuous, and the result is then obtained by letting  $c \nearrow +\infty$ .  $\square$

*Proof of Theorem 2.* Recall the following assumptions from Section 2.2.2:

$$(A1) \quad p_0(\mathbf{x})h_j(x_j)\partial_j \log p(\mathbf{x}) \Big|_{x_j \searrow 0^+}^{x_j \nearrow +\infty} = 0, \quad \forall \mathbf{x}_{-j} \in \mathbb{R}_+^{m-1}, \quad \forall p \in \mathcal{P}_+;$$

$$(A2) \quad \mathbb{E}_{p_0} \|\nabla \log p(\mathbf{X}) \circ \mathbf{h}^{1/2}(\mathbf{X})\|_2^2 < +\infty, \quad \mathbb{E}_{p_0} \|(\nabla \log p(\mathbf{X}) \circ \mathbf{h}(\mathbf{X}))'\|_1 < +\infty, \quad \forall p \in \mathcal{P}_+.$$

Without explicitly writing the domains  $\mathbb{R}_+$  or  $\mathbb{R}_+^m$  in all integrals, by (2.4) we have

$$\begin{aligned} J_{\mathbf{h}}(p) &= \frac{1}{2} \int p_0(\mathbf{x}) \left[ \|\nabla \log p(\mathbf{x}) \circ \mathbf{h}^{1/2}(\mathbf{x})\|_2^2 \right. \\ &\quad \left. - 2(\nabla \log p(\mathbf{x}) \circ \mathbf{h}^{1/2}(\mathbf{x}))^\top (\nabla \log p_0(\mathbf{x}) \circ \mathbf{h}^{1/2}(\mathbf{x})) + \|\nabla \log p_0(\mathbf{x}) \circ \mathbf{h}^{1/2}(\mathbf{x})\|_2^2 \right] d\mathbf{x} \\ &= \frac{1}{2} \underbrace{\int p_0(\mathbf{x}) \sum_{j=1}^m h_j(x_j) \left( \frac{\partial \log p(\mathbf{x})}{\partial x_j} \right)^2 d\mathbf{x}}_{\equiv A} + \frac{1}{2} \underbrace{\int p_0(\mathbf{x}) \sum_{j=1}^m h_j(x_j) \left( \frac{\partial \log p_0(\mathbf{x})}{\partial x_j} \right)^2 d\mathbf{x}}_{\equiv C} \\ &\quad - \underbrace{\int p_0(\mathbf{x}) \sum_{j=1}^m h_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} \frac{\partial \log p_0(\mathbf{x})}{\partial x_j} d\mathbf{x}}_{\equiv B}, \end{aligned}$$

where  $A$  will simply appear in the final display as is,  $C$  is a constant as it only involves the true pdf  $p_0$ , and we wish to simplify  $B$  by integration by parts. We can split the integral into these three parts since  $A$  and  $C$  are assumed finite in the first part of (A2), and the integrand in  $B$  is integrable since  $|2ab| \leq a^2 + b^2$ . Thus, by linearity and Fubini's theorem, we can write

$$\begin{aligned} B &= - \sum_{j=1}^m \int p_0(\mathbf{x}) h_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} \frac{\partial \log p_0(\mathbf{x})}{\partial x_j} d\mathbf{x} \\ &= - \sum_{j=1}^m \int \left[ \int p_0(\mathbf{x}) h_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} \frac{\partial \log p_0(\mathbf{x})}{\partial x_j} dx_j \right] d\mathbf{x}_{-j}. \end{aligned}$$

By the fact that  $\frac{\partial \log p_0(\mathbf{x})}{\partial x_j} = \frac{1}{p_0(\mathbf{x})} \frac{\partial p_0(\mathbf{x})}{\partial x_j}$ , this can be simplified to

$$B = - \sum_{j=1}^m \int \left[ \int \frac{\partial p_0(\mathbf{x})}{\partial x_j} h_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} dx_j \right] d\mathbf{x}_{-j}.$$

But, we assume  $p_0$  and  $p$  are twice continuously differentiable, for every  $j = 1, \dots, m$  and fixed  $\mathbf{x}_{-j} \in \mathbb{R}_+^{m-1}$ . Hence, in every bounded sub-interval of  $\mathbb{R}_+$ ,  $p_0(\mathbf{x}_{-j}; x_j)$  is an absolutely continuous function of  $x_j$ ,  $\partial_j \log p(\mathbf{x}_{-j}; x_j) = \partial_j p(\mathbf{x}_{-j}; x_j) / p(\mathbf{x}_{-j}; x_j)$  is a continuously differentiable (and hence absolutely continuous) function of  $x_j$  by the quotient rule. Thus  $h_j(x_j) \partial_j \log p(\mathbf{x}_{-j}; x_j)$  is also absolutely continuous by the absolute continuity assumption on  $h_j$ . Then, by Lemma 34, where we take  $f \equiv p_0(\mathbf{x}_{-j}; x_j)$  and  $g \equiv h_j(x_j) \partial_j \log p(\mathbf{x}_{-j}; x_j)$  as functions of  $x_j$ , followed by assumption (A1),

$$\begin{aligned} B &= - \sum_{j=1}^m \int \left[ \lim_{a \nearrow +\infty, b \searrow 0^+} [p_0(\mathbf{x}_{-j}; a) h_j(a) \partial_j \log p(\mathbf{x}_{-j}, a) - p_0(\mathbf{x}_{-j}; b) h_j(b) \partial_j \log p(\mathbf{x}_{-j}, b)] \right. \\ &\quad \left. - \int p_0(\mathbf{x}) \frac{\partial (h_j(x_j) \partial_j \log p(\mathbf{x}))}{\partial x_j} dx_j \right] d\mathbf{x}_{-j} \\ &= \sum_{j=1}^m \int \left[ \int p_0(\mathbf{x}) \frac{\partial (h_j(x_j) \partial_j \log p(\mathbf{x}))}{\partial x_j} dx_j \right] d\mathbf{x}_{-j}. \end{aligned}$$

Justified by the second half of (A2), by Fubini-Tonelli and linearity again

$$B = \sum_{j=1}^m \int p_0(\mathbf{x}) \frac{\partial (h_j(x_j) \partial_j \log p(\mathbf{x}))}{\partial x_j} d\mathbf{x},$$

$$= \sum_{j=1}^m \int h'_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} p_0(\mathbf{x}) \, d\mathbf{x} + \sum_{j=1}^m \int h_j(x_j) \frac{\partial^2 \log p(\mathbf{x})}{\partial x_j^2} p_0(\mathbf{x}) \, d\mathbf{x}.$$

Thus,

$$\begin{aligned} & J_{\mathbf{h}}(p) \\ &= B + A + C \\ &= \int_{\mathbb{R}_+^m} p_0(\mathbf{x}) \sum_{j=1}^m \left[ h'_j(x_j) \frac{\partial \log p(\mathbf{x})}{\partial x_j} + h_j(x_j) \frac{\partial^2 \log p(\mathbf{x})}{\partial x_j^2} + \frac{1}{2} h_j(x_j) \left( \frac{\partial \log p(\mathbf{x})}{\partial x_j} \right)^2 \right] d\mathbf{x} + C, \end{aligned}$$

where  $C$  is a constant that does not depend on  $p$ .  $\square$

### A.1.2 Proof of Theorems and Examples in Section 2.3

*Proof of Theorem 3.* For exponential families and under the assumptions, the empirical loss  $\hat{J}_{\mathbf{h}}(p_{\boldsymbol{\theta}})$  in (2.6) becomes (up to an additive constant)

$$\begin{aligned} & \hat{J}_{\mathbf{h}}(p_{\boldsymbol{\theta}}) \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \left[ h'_j(X_j^{(i)}) \frac{\partial \log p_{\boldsymbol{\theta}}(\mathbf{X}^{(i)})}{\partial X_j^{(i)}} + h_j(X_j^{(i)}) \frac{\partial^2 \log p_{\boldsymbol{\theta}}(\mathbf{X}^{(i)})}{\partial (X_j^{(i)})^2} \right. \\ & \qquad \qquad \qquad \left. + \frac{1}{2} h_j(X_j^{(i)}) \left( \frac{\partial \log p_{\boldsymbol{\theta}}(\mathbf{X}^{(i)})}{\partial X_j^{(i)}} \right)^2 \right] \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \left[ h'_j(X_j^{(i)}) (\boldsymbol{\theta}^\top \mathbf{t}'_j(\mathbf{X}^{(i)}) + b'_j(\mathbf{X}^{(i)})) + h_j(X_j^{(i)}) (\boldsymbol{\theta}^\top \mathbf{t}''_j(\mathbf{X}^{(i)}) + b''_j(\mathbf{X}^{(i)})) \right. \\ & \qquad \qquad \qquad \left. + \frac{1}{2} h_j(X_j^{(i)}) (\boldsymbol{\theta}^\top \mathbf{t}'_j(\mathbf{X}^{(i)}) + b'_j(\mathbf{X}^{(i)}))^2 \right] \\ &= \frac{1}{n} \left\{ \frac{1}{2} \boldsymbol{\theta}^\top \left[ \sum_{i=1}^n \sum_{j=1}^m h_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)})^\top \right] \boldsymbol{\theta} + \right. \\ & \qquad \left. \left[ \sum_{i=1}^n \sum_{j=1}^m h_j(X_j^{(i)}) b'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) + h_j(X_j^{(i)}) \mathbf{t}''_j(\mathbf{X}^{(i)}) + h'_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \right]^\top \boldsymbol{\theta} \right\} + \text{const}, \end{aligned}$$

which is quadratic in  $\boldsymbol{\theta}$ . Let

$$\boldsymbol{\Gamma}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m h_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)})^\top, \quad (\text{A.1})$$

$$\mathbf{g}(\mathbf{x}) = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \left[ h_j(X_j^{(i)}) b'_j(\mathbf{X}^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) + h_j(X_j^{(i)}) \mathbf{t}''_j(\mathbf{X}^{(i)}) + h'_j(X_j^{(i)}) \mathbf{t}'_j(\mathbf{X}^{(i)}) \right]. \quad (\text{A.2})$$

Then we can write  $\hat{J}_h(p_\theta) = \frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{\Gamma}(\mathbf{x}) \boldsymbol{\theta} - \mathbf{g}(\mathbf{x})^\top \boldsymbol{\theta} + \text{const}$ .  $\square$

*Proof of Theorem 4.* By Theorem 3,  $\hat{J}_h(p_\theta) = \frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{\Gamma} \boldsymbol{\theta} - \mathbf{g}^\top \boldsymbol{\theta} + \text{const}$ . The minimizer of  $\hat{J}_h(p_\theta)$  is thus available in the unique closed form  $\hat{\boldsymbol{\theta}} \equiv \boldsymbol{\Gamma}(\mathbf{x})^{-1} \mathbf{g}(\mathbf{x})$  as long as  $\boldsymbol{\Gamma}$  is invertible (C1). Since  $\boldsymbol{\Gamma}$  and  $\mathbf{g}$  are sample averages, the weak law of large numbers yields that  $\boldsymbol{\Gamma} \rightarrow_p \mathbb{E}_{p_0} \boldsymbol{\Gamma} \equiv \boldsymbol{\Gamma}_0$  and  $\mathbf{g} \rightarrow_p \mathbb{E}_{p_0} \mathbf{g} \equiv \mathbf{g}_0$ , where existence of  $\boldsymbol{\Gamma}_0$  and  $\mathbf{g}_0$  is assumed in (C2). Since  $J_h(p_\theta) = \mathbb{E}[\hat{J}_h(p_\theta)] = \mathbb{E}[\frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{\Gamma}(\mathbf{x}) \boldsymbol{\theta} - \mathbf{g}(\mathbf{x})^\top \boldsymbol{\theta}] = \frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{\Gamma}_0 \boldsymbol{\theta} - \mathbf{g}_0^\top \boldsymbol{\theta}$  and we know  $\boldsymbol{\theta}_0$  minimizes  $J_h(p_\theta)$  by definition, by the first-order condition we must have  $\boldsymbol{\Gamma}_0 \boldsymbol{\theta}_0 = \mathbf{g}_0$ . Then by the Lindeberg-Lévy central limit theorem,

$$\sqrt{n}(\mathbf{g}(\mathbf{x}) - \boldsymbol{\Gamma}(\mathbf{x}) \boldsymbol{\theta}_0) \rightarrow_d \mathcal{N}_m(\mathbf{0}, \boldsymbol{\Sigma}_0),$$

where  $\boldsymbol{\Sigma}_0 \equiv \mathbb{E}_{p_0}[(\boldsymbol{\Gamma}(\mathbf{x}) \boldsymbol{\theta}_0 - \mathbf{g}(\mathbf{x}))(\boldsymbol{\Gamma}(\mathbf{x}) \boldsymbol{\theta}_0 - \mathbf{g}(\mathbf{x}))^\top]$ , as long as  $\boldsymbol{\Sigma}_0$  exists (C2). Thus, by Slutsky's theorem,

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \equiv \sqrt{n}(\boldsymbol{\Gamma}(\mathbf{x})^{-1}(\mathbf{g}(\mathbf{x}) - \boldsymbol{\Gamma}(\mathbf{x}) \boldsymbol{\theta}_0)) \rightarrow_d \mathcal{N}_r(\mathbf{0}, \boldsymbol{\Gamma}_0^{-1} \boldsymbol{\Sigma} \boldsymbol{\Gamma}_0^{-1}),$$

as long as  $\boldsymbol{\Gamma}_0$  is invertible (C2).

For the second half of the theorem, (C2)  $\mathbb{E}_{p_0} \boldsymbol{\Gamma}(\mathbf{x}) < \infty$  and  $\mathbb{E}_{p_0} \mathbf{g}(\mathbf{x}) < \infty$  implies  $\mathbb{E}_{p_0} |\boldsymbol{\Gamma}(\mathbf{x})| < \infty$  and  $\mathbb{E}_{p_0} |\mathbf{g}(\mathbf{x})| < \infty$ , so by strong law of large numbers (and a union bound on at most  $k^2$  null sets)

$$\boldsymbol{\Gamma}(\mathbf{x}) \rightarrow_{\text{a.s.}} \boldsymbol{\Gamma}_0, \quad \mathbf{g}(\mathbf{x}) \rightarrow_{\text{a.s.}} \mathbf{g}_0.$$

Then outside a null set,

$$\hat{\boldsymbol{\theta}} \equiv \boldsymbol{\Gamma}(\mathbf{x})^{-1} \mathbf{g}(\mathbf{x}) \rightarrow_{\text{a.s.}} \boldsymbol{\Gamma}_0^{-1} \mathbf{g}_0 = \boldsymbol{\theta}_0.$$

$\square$

*Proof for Example 3.1.* We choose to estimate  $\theta \equiv \mu/\sigma^2$ . Then by (2.9) and (2.10),

$$\begin{aligned}\hat{\mu}_h &= \sigma^2 \hat{\theta} \equiv \sigma^2 \Gamma(\mathbf{x})^{-1} g(\mathbf{x}) \\ &= -\sigma^2 \left[ \sum_{i=1}^n h(X_i) t'(X_i)^2 \right]^{-1} \left[ \sum_{i=1}^n h(X_i) b'(X_i) t'(X_i) + h(X_i) t''(X_i) + h'(X_i) t'(X_i) \right] \\ &= -\sigma^2 \left[ \sum_{i=1}^n h(X_i) \right]^{-1} \left[ \sum_{i=1}^n -h(X_i) \frac{X_i}{\sigma^2} + h'(X_i) \right].\end{aligned}$$

By Theorem 4,

$$\begin{aligned}\sqrt{n}(\hat{\mu}_h - \mu_0) &\rightarrow_d \mathcal{N} \left( 0, \frac{\sigma^4 \mathbb{E}_0 \left[ -h(X) \frac{X - \mu_0}{\sigma^2} + h'(X) \right]^2}{\mathbb{E}_0^2[h(X)]} \right) \\ &\sim \mathcal{N} \left( 0, \frac{\mathbb{E}_0 \left[ -h(X)(X - \mu_0) + \sigma^2 h'(X) \right]^2}{\mathbb{E}_0^2[h(X)]} \right).\end{aligned}$$

By integration by parts, (suppressing the dependence of  $p_{\mu_0}$  on  $\mu_0$ )

$$\begin{aligned}&\mathbb{E}_0[h(X)h'(X)(X - \mu_0)] \\ &= \int_0^\infty h'(x)h(x)(x - \mu_0)p(x) dx \\ &= \int_0^\infty h(x)(x - \mu_0)p(x) dh(x) \\ &= h^2(x)(x - \mu_0)p(x)|_0^\infty - \int h(x) dh(x)(x - \mu_0)p(x) \\ &= - \int h^2(x)p(x) dx - \int h(x)h'(x)(x - \mu_0)p(x) dx + \int h^2(x) \frac{(x - \mu_0)^2}{\sigma^2} p(x) dx,\end{aligned}$$

where the last step follows from the assumptions  $\lim_{x \searrow 0^+} h(x) = 0$  and  $\lim_{x \nearrow +\infty} h^2(x)(x - \mu_0)p_{\mu_0}(x) = 0$ . So

$$\mathbb{E}_0[h(X)h'(X)(X - \mu_0)] = \frac{\mathbb{E}[h^2(X)((X - \mu_0)^2/\sigma^2 - 1)]}{2}. \quad (\text{A.3})$$

The asymptotic variance is thus

$$\begin{aligned}&\frac{\mathbb{E}_0 \left[ -h(X)(X - \mu_0) + \sigma^2 h'(X) \right]^2}{\mathbb{E}_0^2[h(X)]} \\ &= \frac{\mathbb{E}_0 \left[ h^2(X)(X - \mu_0)^2 - 2\sigma^2 h^2(X) \left( (X - \mu_0)^2/\sigma^2 - 1 \right) / 2 + \sigma^4 h'^2(X) \right]}{\mathbb{E}_0^2[h(X)]}\end{aligned}$$

$$= \frac{\mathbb{E}_0[\sigma^2 h^2(X) + \sigma^4 h'^2(X)]}{\mathbb{E}_0^2[h(X)]}.$$

The Cramér-Rao lower bound follows from taking the second derivative of  $\log p_{\mu_0}$  with respect to  $\mu_0$ .  $\square$

*Proof for Example 3.2.* We estimate  $\theta \equiv 1/\sigma^2$ . By (2.9) and (2.10),

$$\begin{aligned} \hat{\theta} &\equiv \Gamma(\mathbf{x})^{-1} g(\mathbf{x}) \\ &= - \left[ \sum_{i=1}^n h(X_i) t'(X_i)^2 \right]^{-1} \left[ \sum_{i=1}^n h(X_i) b'(X_i) t'(X_i) + h(X_i) t''(X_i) + h'(X_i) t'(X_i) \right] \\ &= \left[ \sum_{i=1}^n h(X_i) (X_i - \mu)^2 \right]^{-1} \left[ \sum_{i=1}^n h(X_i) + h'(X_i) (X_i - \mu) \right]. \end{aligned}$$

By Theorem 4,  $\sqrt{n}(\hat{\theta} - \theta) \rightarrow_d \mathcal{N}(0, \varsigma^2)$ , where

$$\begin{aligned} \varsigma^2 &\equiv \frac{\mathbb{E}_0[h(X)((X - \mu)^2/\sigma_0^2 - 1) - h'(X)(X - \mu)]^2}{\mathbb{E}_0^2[h(X)(X - \mu)^2]} \\ &= \frac{1}{\mathbb{E}_0^2[h(X)(X - \mu)^2]} \left( \mathbb{E}_0[h^2(X)(X - \mu)^4/\sigma_0^4 - 2h^2(X)(X - \mu)^2/\sigma_0^2 + h^2(X) \right. \\ &\quad \left. + h'^2(X)(X - \mu)^2 - 2h(X)h'(X)(X - \mu)^3/\sigma_0^2 + 2h(X)h'(X)(X - \mu) \right). \end{aligned}$$

By integration by parts, (suppressing the dependence of  $p_{\sigma_0^2}$  on  $\sigma_0^2$ )

$$\begin{aligned} &\mathbb{E}_0[h(X)h'(X)(X - \mu)^3] \\ &= \int_0^\infty h'(x)h(x)(x - \mu)^3 p(x) dx \\ &= \int_0^\infty h(x)(x - \mu)^3 p(x) dh(x) \\ &= h^2(x)(x - \mu)^3 p(x) \Big|_0^\infty - \int h(x) dh(x)(x - \mu)^3 p(x) \\ &= - \int h(x)h'(x)(x - \mu)^3 p(x) dx - 3 \int h^2(x)(x - \mu)^2 p(x) dx + \int h^2(x) \frac{(x - \mu)^4}{\sigma_0^2} p(x) dx, \end{aligned}$$

where the last step follows from the assumptions  $\lim_{x \searrow 0^+} h(x) = 0$  and  $\lim_{x \nearrow +\infty} h^2(x)(x - \mu)^3 p_{\sigma_0^2}(x) = 0$ . Combining this with (A.3) we get

$$\sqrt{n}(\hat{\theta} - \theta) \rightarrow_d \mathcal{N}(0, \varsigma^2) \sim \mathcal{N}\left(0, \frac{2\mathbb{E}_0[h^2(X)(X - \mu)^2/\sigma_0^2] + \mathbb{E}_0[h'^2(X - \mu)^2]}{\mathbb{E}_0^2[h(X)(X - \mu)^2]}\right),$$

and so by the delta method, for  $\hat{\sigma}_k^2 \equiv \hat{\theta}^{-1}$ ,

$$\sqrt{n}(\hat{\sigma}_h^2 - \sigma_0^2) \rightarrow_d \mathcal{N}\left(0, \frac{2\sigma_0^6 \mathbb{E}_0[h^2(X)(X - \mu)^2] + \sigma_0^8 \mathbb{E}_0[h'^2(X - \mu)^2]}{\mathbb{E}_0^2[h(X)(X - \mu)^2]}\right).$$

The Cramér-Rao lower bound follows from taking the second derivative of  $\log p_{\sigma_0^2}$  with respect to  $\sigma_0^2$ .  $\square$

### A.1.3 Proof of Theorems in Section 2.5

*Proof of Theorem 5.*

*Case  $b \neq 0$ :* We use a strategy similar to that of Inouye et al. (2016). Let  $\mathcal{V}_1 = \{\mathbf{v} : \|\mathbf{v}\|_1 = 1, \mathbf{v} \in \mathbb{R}_+^m\}$ . Then by Fubini-Tonelli the normalizing constant is,

$$\begin{aligned} & \int_{\mathbb{R}_+^m} \exp\left(\boldsymbol{\eta}^\top \frac{\mathbf{x}^b - \mathbf{1}_m}{b} - \frac{1}{2a} \mathbf{x}^{a\top} \mathbf{K} \mathbf{x}^a\right) d\mathbf{x} \\ &= \int_{\mathcal{V}_1} \int_0^\infty \exp\left(\boldsymbol{\eta}^\top \frac{z^b \mathbf{v}^b - \mathbf{1}_m}{b} - \frac{1}{2a} z^{2a} \mathbf{v}^{a\top} \mathbf{K} \mathbf{v}^a\right) dz d\mathbf{v} \\ &\propto \int_{\mathcal{V}_1} \int_0^\infty \exp\left(z^b (\boldsymbol{\eta}^\top \mathbf{v}^b)/b - z^{2a} (\mathbf{v}^{a\top} \mathbf{K} \mathbf{v}^a)/(2a)\right) dz d\mathbf{v}. \end{aligned}$$

Here  $\mathcal{V}_1$  is compact and the inner integral, if finite, is continuous in  $\mathbf{v}$ . It thus suffices to show that the inner integral is finite at every single  $\mathbf{v} \in \mathcal{V}_1$ .

Fixing  $\mathbf{v} \in \mathcal{V}_1$ , write  $A \equiv A(\mathbf{v}) \equiv \mathbf{v}^{a\top} \mathbf{K} \mathbf{v}^a / (2a)$  and  $B \equiv B(\mathbf{v}) \equiv (\boldsymbol{\eta}^\top \mathbf{v}^b) / b$ . We need to show that

$$N(A, B, a, b) \equiv \int_0^\infty \exp(-Az^{2a} + Bz^b) dz < +\infty.$$

Recall that (CC1)  $\mathbf{v}^\top \mathbf{K} \mathbf{v} > 0$  for all  $\mathbf{v} \in \mathbb{R}_+^m \setminus \{\mathbf{0}\}$ , so  $A > 0$ .

- (i) Suppose  $B \leq 0$ . Then  $N(A, B, a, b) \leq \int_0^\infty \exp(-Az^{2a}) dz = A^{-a/2} \Gamma(1 + 1/(2a))$ , a finite constant since  $A > 0$  and  $a > 0$ .
- (ii) Suppose  $B > 0$ . We first want to bound  $\exp(-Az^{2a} + Bz^b) \leq N_0 \exp(-Az^{2a}/2)$  by some finite constant  $N_0 > 0$ , so that  $N(A, B, a, b) \leq N_0 \int_0^\infty \exp(-Az^{2a}/2) dz$ , a finite constant for  $a > 0$ . Thus, it remains to give conditions so that  $\exp(-Az^{2a}/2 + Bz^b)$

is bounded by some finite constant  $N_0$ , which by continuity only requires a finite limit as  $z \searrow 0$  and as  $z \nearrow +\infty$ . As  $z \nearrow +\infty$ ,  $Bz^b \nearrow +\infty$ , while  $-Az^{2a}/2 \searrow -\infty$ . We thus need  $b < 2a$  so that the sum of the two does not go to positive infinity. On the other hand, as  $z \searrow 0$ ,  $-Az^{2a}/2 \nearrow 0$ , so we need  $b > 0$ , otherwise  $z^b \nearrow +\infty$ . In conclusion, we require that  $2a > b > 0$ .

It thus suffices to require (CC1) and (CC2)  $2a > b > 0$  to eliminate restrictions on  $B$ , and hence on  $\boldsymbol{\eta}$ . That is,  $\boldsymbol{\eta}$  can take value in the entirety of  $\mathbb{R}^m$ .

*Case  $b = 0$ :* Again in (CC1) we assume  $\mathbf{v}^\top \mathbf{K} \mathbf{v} > 0$  for all  $\mathbf{v} \in \mathbb{R}_+^m \setminus \{\mathbf{0}\}$ . Since  $\mathcal{V}_2 \equiv \{\mathbf{v} : \|\mathbf{v}\|_2 = 1, \mathbf{v} \in \mathbb{R}_+^m\}$  is compact and  $\mathbf{v}^\top \mathbf{K} \mathbf{v}$  is continuous in  $\mathbf{v}$  and strictly positive on  $\mathcal{V}_2$ , the image of  $\mathcal{V}_2$  under  $\mathbf{v}^\top \mathbf{K} \mathbf{v}$  is a compact subset of  $(0, \infty)$ , i.e.  $N_{\mathbf{K}} \equiv \inf_{\mathbf{v} \in \mathbb{R}_+^m \setminus \{\mathbf{0}\}} \mathbf{v}^\top \mathbf{K} \mathbf{v} / \mathbf{v}^\top \mathbf{v} \equiv \inf_{\mathbf{v} \in \mathcal{V}_2} \mathbf{v}^\top \mathbf{K} \mathbf{v} > 0$ . We thus have

$$\begin{aligned} & \int_{\mathbb{R}_+^m} \exp\left(\boldsymbol{\eta}^\top \log(\mathbf{x}) - \frac{1}{2a} \mathbf{x}^{a\top} \mathbf{K} \mathbf{x}^a\right) d\mathbf{x} \\ & \leq \int_{\mathbb{R}_+^m} \exp\left(\boldsymbol{\eta}^\top \log(\mathbf{x}) - \frac{N_{\mathbf{K}}}{2a} \mathbf{x}^{a\top} \mathbf{x}^a\right) d\mathbf{x} \\ & = \prod_{j=1}^m \int_0^\infty \exp\left(\eta_j \log(x_j) - \frac{N_{\mathbf{K}}}{2a} x_j^{2a}\right) dx_j \\ & = \prod_{j=1}^m \left[ \Gamma\left(\frac{\eta_j + 1}{2a}\right) \frac{(N_{\mathbf{K}}/2a)^{-\frac{\eta_j + 1}{2a}}}{2a} \right], \end{aligned}$$

where the integration follows by change of variable and requires  $a > 0$ . Assuming  $a > 0$ , the last quantity is finite if and only if  $\boldsymbol{\eta} \succ -\mathbf{1}_m$ , by definition of the gamma function.

In conclusion, given conditions (CC1)  $\inf_{\mathbf{v} \in \mathbb{R}_+^m \setminus \{\mathbf{0}\}} \mathbf{v}^\top \mathbf{K} \mathbf{v} > 0$ , (CC2)  $2a > b > 0$ , and (CC3)  $a > 0$ ,  $b = 0$  and  $\boldsymbol{\eta} \succ -\mathbf{1}_m$ , the unnormalized density (2.14) has a finite normalizing constant when (CC1) and (CC2) both hold, or (CC1) and (CC3) both hold.

The centered settings, where the term involving  $\mathbf{x}^b$  is excluded, can be considered as a special case of both (1) and (2) with  $\boldsymbol{\eta} \equiv \mathbf{0}$ , and thus (CC1) and  $a > 0$  are sufficient.  $\square$

*Proof of Theorem 6.* Recall assumptions (A1) and (A2):

$$(A1) \quad p_0(\mathbf{x})h_j(x_j)\partial_j \log p(\mathbf{x}) \Big|_{x_j \nearrow +\infty}^{x_j \searrow 0^+} = 0, \quad \forall \mathbf{x}_{-j} \in \mathbb{R}_+^{m-1}, \quad \forall p \in \mathcal{P}_+;$$

$$(A2) \quad \mathbb{E}_{p_0} \|\nabla \log p(\mathbf{X}) \circ \mathbf{h}^{1/2}(\mathbf{X})\|_2^2 < +\infty, \quad \mathbb{E}_{p_0} \|(\nabla \log p(\mathbf{X}) \circ \mathbf{h}(\mathbf{X}))'\|_1 < +\infty, \quad \forall p \in \mathcal{P}_+.$$

Let  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  be the true parameters so that  $p_0 \in \mathcal{P}_+$ , with  $\mathcal{P}_+$  corresponding to a parameter space in which all parameters satisfy the conditions for a finite normalizing constant. We now give sufficient conditions for  $h$  to satisfy (A1) and (A2).

*Conditions for (A1):* Fix  $j = 1, \dots, m$  and  $\mathbf{x}_{-j} \in \mathbb{R}_+^{m-1}$ . We show that the conditions on  $h_j$  imply that the limits go to 0 as  $x_j \nearrow +\infty$  and as  $x_j \searrow 0^+$ , which is stronger than (A1); in fact, from (A.4) below, the limits cannot go to a nonzero finite constant assuming an  $h$  with polynomial tail, since  $a > 0$  and  $B_1 \equiv \kappa_{0,jj} > 0$  for all  $j$ . Now,

$$\begin{aligned} & p_0(\mathbf{x})h_j(x_j)\partial_j \log p(\mathbf{x}) \\ & \propto h_j(x_j) \exp\left(-\frac{1}{2a}\mathbf{x}^{a\top}\mathbf{K}_0\mathbf{x}^a + \boldsymbol{\eta}_0^\top \frac{\mathbf{x}^b - \mathbf{1}_m}{b}\right) \partial_j \left(-\frac{1}{2a}\mathbf{x}^{a\top}\mathbf{K}\mathbf{x}^a + \boldsymbol{\eta}^\top \frac{\mathbf{x}^b - \mathbf{1}_m}{b}\right) \\ & \propto h_j(x_j) \exp\left(-\frac{1}{a}(\mathbf{k}_{0,j,-j}^\top \mathbf{x}_{-j}^a)x_j^a - \frac{\kappa_{0,jj}}{2a}x_j^{2a} + \eta_{0,j} \frac{x_j^b - 1}{b}\right) \times \\ & \quad \left(-\mathbf{k}_{j,-j}^\top \mathbf{x}_{-j}^a x_j^{a-1} - \kappa_{jj}x_j^{2a-1} + \eta_j x_j^{b-1}\right) \\ & \equiv h_j(x_j) \exp\left(\frac{A_1 x_j^a}{a} + \frac{B_1 x_j^{2a}}{2a} + C_1 \frac{x_j^b - 1}{b}\right) (A_2 x_j^{a-1} + B_2 x_j^{2a-1} + C_2 x_j^{b-1}), \end{aligned} \quad (A.4)$$

where  $A_1 \equiv -\mathbf{k}_{0,j,-j}^\top \mathbf{x}_{-j}^a$ ,  $A_2 \equiv -\mathbf{k}_{j,-j}^\top \mathbf{x}_{-j}^a$ ,  $B_1 \equiv -\kappa_{0,jj} < 0$  and  $B_2 \equiv -\kappa_{jj} < 0$  by condition (CC1). Finally  $C_1 \equiv \eta_{0,j}$ ,  $C_2 \equiv \eta_j$ .

- (1) Let  $x_j \nearrow +\infty$ . If  $b > 0$ , since  $2a > b > 0$  and  $B_1 < 0$ , the exponential term in (A.4) decreases to 0 exponentially and its reciprocal dominates any polynomial functions. Thus, the entire product goes to 0 if  $h_j(x_j)$  grows no faster than polynomially as  $x_j \nearrow +\infty$ . If  $b = 0$ , the  $C_1 \log x_j$  term is again dominated by  $B_1 x_j^{2a}/(2a)$ , and the same conclusion holds.

(2) Let  $x_j \searrow 0$ .

(i) Let  $b > 0$ . Then the exponential term in (A.4) goes to constant  $\exp(-C_1/b)$ , and we only need

$$\lim_{x_j \searrow 0^+} h_j(x_j)(A_2 x_j^{a-1} + B_2 x_j^{2a-1} + C_2 x_j^{b-1}) = 0. \quad (\text{A.5})$$

- If  $a > 1$  and  $b > 1$ , the second term in (A.5) is a polynomial with three terms having powers  $\geq \min\{a-1, b-1\}$ . The product goes to zero if and only if  $h_j(x_j) = o(x_j^{\max\{1-a, 1-b\}})$  as  $x_j \searrow 0$ . Note that this is satisfied by any  $h_j$  that has a finite right limit at 0.
- If  $a = 1$  and  $b \geq 1$ , or  $a \geq 1$  and  $b = 1$ , then the second term in (A.5) is a polynomial of non-negative power plus a potentially nonzero constant. A sufficient condition for (A.5) is thus  $\lim_{x_j \searrow 0} h_j(x_j) = 0$ .
- If  $a < 1$  or  $b < 1$ , then the second part in (A.5) is a polynomial having terms with negative degree  $\geq \min\{a-1, b-1\}$ . To counteract this a sufficient condition is  $h_j(x_j) = o(x_j^{\max\{1-a, 1-b\}})$ .

In conclusion,  $\lim_{x_j \searrow 0^+} p_0(\mathbf{x})h_j(x_j)\partial_j \log p(\mathbf{x}) = 0$  if and only if

$$\lim_{x_j \searrow 0^+} h_j(x_j)/x_j^{\max\{1-a, 1-b\}} = 0.$$

(ii) Now assume  $b = 0$ . Then, (A.4) now becomes

$$h_j(x_j) \exp\left(A_1 x_j^a/a + B_1 x_j^{2a}/(2a) + C_1 \log x_j\right) \left(A_2 x_j^{a-1} + B_2 x_j^{2a-1} + C_2/x_j\right).$$

With  $C_1 \log x_j$  dominating, the exponential part scales as  $x_j^{C_1}$ . We thus require

$$\lim_{x_j \searrow 0^+} h_j(x_j)(A_2 x_j^{a-1+C_1} + B_2 x_j^{2a-1+C_1} + C_2 x_j^{C_1-1}) = 0,$$

which by the previous discussion on (A.5) holds if and only if

$$\lim_{x_j \searrow 0^+} h_j(x_j)/x_j^{1-C_1} = 0$$

since  $1-a-C_1 < 1-C_1$ .

In summary, (A1) is satisfied if  $h_j(x_j)$  grows at most polynomially as  $x_j \nearrow +\infty$ , and  $\lim_{x_j \searrow 0^+} h_j(x_j)/x_j^{\max\{1-a, 1-b\}} = 0$  if  $b > 0$ , or  $\lim_{x_j \searrow 0^+} h_j(x_j)/x_j^{1-\eta_{0,j}} = 0$  if  $b = 0$ .

*Conditions for (A2):* For (A2), we consider powers of  $x$  as the  $h$  functions for simplicity; conclusions for other functions that have the same tail behavior (big-O scaling) as  $x \searrow 0$  and  $x \nearrow +\infty$  follow similarly. Sufficiency results for piecewise power functions follow by partitioning, and similarly for other functions  $h$  whose function values and derivatives can be bounded by those of some piecewise power function (e.g. truncated powers), since (A2) is on integrability of products involving positive powers of  $h$  and  $h'$ .

Let  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  be the true parameters from the parameter space that satisfies the conditions for finite normalizing constant. By part (2) of the proof of Theorem 5, the assumption that  $\mathbf{K}_0$  satisfies (CC1) implies that  $\min_{\mathbf{v} \in \mathbb{R}_+^m \setminus \{0\}} \mathbf{v}^\top \mathbf{K}_0 \mathbf{v} / \mathbf{v}^\top \mathbf{v} \equiv N_{\mathbf{K}_0} > 0$ . Then we have the following decomposition

$$\begin{aligned} p_{\mathbf{K}_0, \boldsymbol{\eta}_0}(\mathbf{x}) &\equiv \exp\left(-\frac{1}{2a} \mathbf{x}^{a\top} \mathbf{K}_0 \mathbf{x}^a + \boldsymbol{\eta}_0^\top \frac{\mathbf{x}^b - \mathbf{1}_m}{b}\right) \\ &\leq \prod_{j=1}^m \exp\left(-\frac{N_{\mathbf{K}_0}}{2a} x_j^{2a} + \eta_{0,j} \frac{x_j^b - 1}{b}\right). \end{aligned}$$

Then for any other  $\mathbf{K}$  and  $\boldsymbol{\eta}$  in the parameter space, for the first part of (A2) it suffices to show for any  $j = 1, \dots, m$  that  $D < \infty$ , where

$$\begin{aligned} D &\equiv \int_{\mathbb{R}_+^m} \prod_{j=1}^m \exp\left(-\frac{N_{\mathbf{K}_0}}{2a} x_j^{2a} + \eta_{0,j} \frac{x_j^b - 1}{b}\right) h_j(x_j) \times \\ &\quad \left(-\kappa_{jj} x_j^{2a-1} - \sum_{i \neq j} \kappa_{ji} x_i^a x_j^{a-1} + \eta_j x_j^{b-1}\right)^2 d\mathbf{x} \\ &\geq \int_{\mathbb{R}_+^m} p_0(\mathbf{x}) h_j(x_j) (\partial_j \log p(\mathbf{x}))^2 d\mathbf{x}. \end{aligned}$$

Note that

$$\left(-\kappa_{jj} x_j^{2a-1} - \sum_{i \neq j} \kappa_{ji} x_i^a x_j^{a-1} + \eta_j x_j^{b-1}\right)^2$$

$$\begin{aligned}
&= \kappa_{jj}^2 x_j^{4a-2} + \sum_{i \neq j, \ell \neq j} \kappa_{ji} \kappa_{j\ell} x_i^a x_\ell^a x_j^{2a-2} + \eta_j^2 x_j^{2b-2} + 2 \sum_{i \neq j} \kappa_{jj} \kappa_{ji} x_i^a x_j^{3a-2} \\
&\quad - 2 \sum_{i \neq j} \kappa_{ji} \eta_j x_i^a x_j^{a+b-2} - 2 \kappa_{jj} \eta_j x_j^{2a+b-2}.
\end{aligned}$$

Thus, plugging this back in the definition of  $D$ , we can split  $D$  into a sum of six terms  $D_1$  through  $D_6$ , each of which is a sum of terms of the form

$$\begin{aligned}
&\int_{\mathbb{R}_+^m} \prod_{k=1}^m \exp\left(-\frac{N_{\mathbf{K}_0}}{2a} x_k^{2a} + \eta_{0,k} \frac{x_k^b - 1}{b}\right) h_j(x_j) x_i^{\text{pow}_i} x_\ell^{\text{pow}_\ell} x_j^{\text{pow}_j} \, d\mathbf{x} \\
&= \prod_{k \neq j} \int_0^\infty \exp\left(-\frac{N_{\mathbf{K}_0}}{2a} x_k^{2a} + \eta_{0,k} \frac{x_k^b - 1}{b}\right) x_k^{\text{pow}_k} \, dx_k \\
&\quad \times \int_0^\infty \exp\left(-\frac{N_{\mathbf{K}_0}}{2a} x_j^{2a} + \eta_{0,j} \frac{x_j^b - 1}{b}\right) h_j(x_j) x_j^{\text{pow}_j} \, dx_j
\end{aligned}$$

times a constant involving  $\mathbf{K}$  and  $\eta_j$ , where  $\text{pow}_k \geq 0$  for each  $k \neq j$ . We have thus decomposed the integral into a product of univariate integrals. Note that

$$\int_0^\infty \exp\left(-\frac{N_{\mathbf{K}_0}}{2a} x_i^{2a} + \eta_{0,i} \frac{x_i^b - 1}{b}\right) x_i^{\text{pow}_i} \, dx_i$$

is finite for all  $\text{pow}_i \geq 0$  regardless of whether  $b$  is nonzero, since we assumed  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  to lie in the parameter space with a finite normalizing constant. Indeed, if  $b > 0$  then the terms in the exponential is a regular polynomial with positive degree and a negative leading term; if  $b = 0$  then integrability follows from  $\eta_{0,i} + \text{pow}_i \geq \eta_{0,i} > -1$ . Thus, we only need to consider the univariate integral that involve the  $x_j$  terms, namely

$$\int_0^\infty \exp\left(-\frac{N_{\mathbf{K}_0}}{2a} x_j^{2a} + \eta_{0,j} \frac{x_j^b - 1}{b}\right) h_j(x_j) x_j^{\text{pow}_j} \, dx_j,$$

where  $\text{pow}_j$  takes value in  $\{4a-2, 2a-2, 2b-2, 3a-2, a+b-2, 2a+b-2\} \subseteq [2 \min\{a, b\} - 2, 4a-2]$ . We split the integral into two parts over  $[0, 1]$  and  $[1, \infty]$ , respectively.

- If  $b > 0$ , on  $[0, 1]$  the exponential part is bounded above and below by positive constants, and for (A1) we require  $h_j(x) = o(x^{1-\min\{a,b\}})$  as  $x \searrow 0^+$ , so the integrand is  $o(x^{\min\{a,b\}-1}) = o(x^{-1})$  and is thus integrable on  $[0, 1]$ . The integrand on  $[1, \infty)$  is integrable as in (A1) we assume  $h$  to grow at most polynomially.

- If  $b = 0$ ,  $\text{pow}_j \in [-2, 4a - 2]$  and the integrand becomes

$$\exp(-N_{\mathbf{K}_0} x_j^{2a} / (2a)) h_j(x_j) x_j^{\text{pow}_j + \eta_{0,j}}.$$

On  $[0, 1]$ , (A1) requires  $h_j(x) = o(x^{1 - \min_j \eta_{0,j}})$ , so  $h_j(x_j) x_j^{\text{pow}_j + \eta_{0,j}} = o(x^{-1})$  and the integrand is again integrable. Integrability on  $[1, \infty)$  follows similarly to the case with  $b > 0$ .

Now consider the second part of (A2). By definition  $\mathbb{E}_{p_0} \|(\nabla \log p(\mathbf{X}) \circ \mathbf{h}(\mathbf{X}))'\|_1$  equals

$$\begin{aligned} & \int_{\mathbb{R}_+^m} p_0(\mathbf{x}) \sum_{j=1}^m |h'_j(X_j) \partial_j \log p(\mathbf{X}) + h_j(X_j) \partial_j^2 \log p(\mathbf{X})| \, d\mathbf{x} \\ & \leq \sum_{j=1}^m \int_{\mathbb{R}_+^m} \exp\left(-\frac{N_{\mathbf{K}_0}}{2a} x_j^{2a} + \eta_{0,j} \frac{x_j^b - 1}{b}\right) \left| h'_j(x_j) \left(-\kappa_{jj} x_j^{2a-1} - \sum_{i \neq j} \kappa_{ji} x_i^a x_j^{a-1} + \eta_j x_j^{b-1}\right) \right. \\ & \quad \left. + h_j(x_j) \left(-\kappa_{jj}(2a-1) x_j^{2a-2} - \sum_{i \neq j} \kappa_{ji}(a-1) x_i^a x_j^{a-2} + (b-1) \eta_j x_j^{b-2}\right) \right| \, d\mathbf{x}. \end{aligned}$$

By the triangle inequality and the fact that  $h_j \geq 0$  and  $h'_j \geq 0$ , similar to the proof for the first part, for each  $j$  the integral can be bounded by a sum of six integrals, each of the form

$$\text{const} \times \int_{\mathbb{R}_+^m} \prod_{k=1}^m \exp\left(-\frac{N_{\mathbf{K}_0}}{2a} x_k^{2a} + \eta_{0,k} \frac{x_k^b - 1}{b}\right) h_j(x_j) x_i^{\text{pow}_i} x_j^{\text{pow}_j} \, d\mathbf{x},$$

or with  $h_j$  replaced by  $h'_j$ . Finiteness thus follows from the same type of discussion by noting that  $h_j(x) = o(x^{1 - \min\{a,b\}})$  and  $h'_j(x) = o(x^{-\min\{a,b\}})$ .

We conclude that if the true and the proposed parameters give densities with finite normalizing constants, and if  $h$  satisfies assumption (A1), then (A2) is automatically satisfied.

In the centered case where we assume  $\boldsymbol{\eta} \equiv \mathbf{0}$ , we only need  $\lim_{x_j \searrow 0^+} h_j(x_j) / x_j^{1-a} = 0$  as it is a special case with  $b = 2a$ .  $\square$

#### A.1.4 Proof of Theorems in Section 2.6

*Proof of Corollary 8.* By Theorem 7, under assumptions in that theorem, the support of  $\hat{\boldsymbol{\Psi}}$  is a subset of the true support of  $\boldsymbol{\Psi}_0$ , and  $\|\hat{\boldsymbol{\Psi}} - \boldsymbol{\Psi}_0\|_\infty \leq \frac{\epsilon_{\mathbf{r}_0}}{2-\alpha} \lambda$ . Since  $\boldsymbol{\Psi}_0$  has  $|S_0|$  nonzero

entries,

$$\|\hat{\Psi} - \Psi_0\|_F = \left[ \sum_{\Psi_{0,jk} \neq 0} (\hat{\Psi}_{jk} - \Psi_{0,jk})^2 \right]^{1/2} \leq \sqrt{|S_0|} \|\hat{\Psi} - \Psi_0\|_\infty \leq \frac{c_{\Gamma_0}}{2 - \alpha} \lambda \sqrt{|S_0|}.$$

Similarly, by the definition of matrix  $\ell_\infty$ - $\ell_\infty$  norm,

$$\|\hat{\Psi} - \Psi_0\|_2 \leq \|\hat{\Psi} - \Psi_0\|_\infty = \max_{j=1, \dots, m} \sum_{k=1}^m |\hat{\Psi}_{jk} - \Psi_{0,jk}| \leq \frac{c_{\Gamma_0}}{2 - \alpha} \lambda d_{\Psi_0}.$$

The result follows by also noting that  $\|\hat{\Psi} - \Psi_0\|_2 \leq \|\hat{\Psi} - \Psi_0\|_F$ .  $\square$

*Proof of Theorem 9.* The proof is based on Theorem 7 and a probabilistic bound on  $\|\Gamma_\gamma - \Gamma_0\|_\infty$ , where in the case of centered Gaussian  $\Gamma = \text{diag}(\mathbf{xx}^\top, \dots, \mathbf{xx}^\top)$ . Denote  $\Sigma_0 = \mathbf{K}_0^{-1}$ . In particular, given  $\tau > 2$  we wish to show that for  $\epsilon = 80\sqrt{2}c_0 \max_j(\Sigma_{0,jj})$ , assuming  $c_0 \equiv \sqrt{(\tau \log m + \log 4)/n} < 1/\sqrt{2}$ ,

$$\mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_j^{(i)} X_k^{(i)} + \gamma_{1j} \mathbf{1}_{\{j=k\}} - \mathbb{E} X_j X_k \right| > \epsilon \right) \leq m^{2-\tau},$$

and so the results follow from Theorem 7.

By Lemma 1 of Ravikumar et al. (2011), since  $X_j/\sqrt{\Sigma_{0,jj}}$  is Gaussian with mean 0 and standard deviation 1,

$$\mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_j^{(i)} X_k^{(i)} - \mathbb{E} X_j X_k \right| > t \right) \leq 4 \exp \left( -\frac{nt^2}{3200 \max_j(\Sigma_{0,jj})^2} \right)$$

for  $t \in (0, 40 \max_j(\Sigma_{0,jj}))$ . Denote the event as  $\mathcal{E}_{j,k}(t)$ . Note that  $\mathbb{E} X_j^2 \leq \max_j \Sigma_{0,jj} = \epsilon/(80\sqrt{2}c_0)$ . Then letting  $t = \epsilon/2$  and conditioning on the complement of  $\mathcal{E}_{j,j}(\epsilon/2)$ , we have

$$n^{-1} \sum_{i=1}^n X_j^{(i)2} \leq \mathbb{E} X_j^2 + \epsilon/2 \leq \frac{\epsilon}{2} \left( 1 + \frac{1}{40\sqrt{2}c_0} \right).$$

Thus, choosing  $\gamma_{\ell j} = (\delta - 1) \sum_{i=1}^n X_j^{(i)2}/n$  for  $\ell = 1, \dots, m$  ( $\Gamma$  has  $m$  identical blocks) with  $1 < \delta < 1 + (1 + 1/(40\sqrt{2}c_0))^{-1}$ , by the triangle inequality and a union bound we have

$$\mathbb{P} \left( \max_{j,k} \left| n^{-1} \sum_{i=1}^n X_j^{(i)} X_k^{(i)} + \gamma_{1j} \mathbf{1}_{\{j=k\}} - \mathbb{E} X_j X_k \right| > \epsilon \right) \leq \mathbb{P}(\mathcal{E}_{j,k}(\epsilon/2)) = m^{2-\tau}.$$

Since  $\tau > 2$ , it holds that  $1 + (1 + 1/(40\sqrt{2}c_0))^{-1} = 2 - (1 + 40\sqrt{2}c_0)^{-1}$  is larger than  $2 - (1 + 80\sqrt{\log m/n})^{-1} \equiv C(n, m)$ , so it is safe to choose any  $\delta \in (1, C(n, m))$ . Thus by the requirement on  $\epsilon$ , the theorem statement holds when  $n > \max(c^*c_1^2d_{\mathbf{K}}^2, 2)(\tau \log m + \log 4)$  with  $c^* = 12800 \max_j(\Sigma_{0,jj})^2$ .  $\square$

*Proof of Theorem 10.* The proof of Theorem 7 from Lin et al. (2016) does not rely on the fact that the original  $\mathbf{\Gamma}$  is an unbiased estimator for the population  $\mathbf{\Gamma}_0$ , but instead only requires one to bound  $\|\mathbf{\Gamma} - \mathbf{\Gamma}_0\|_\infty$ . Thus, for  $\mathbf{\Gamma}_\gamma = \mathbf{\Gamma} + \text{diag}(\gamma)$ , by Theorem 7 it suffices to prove that for any  $\tau > 3$ , we can bound  $\|\mathbf{\Gamma}(\mathbf{x}) + \text{diag}(\gamma(\mathbf{x})) - \mathbf{\Gamma}_0\|_\infty$  by some  $\epsilon_1$  and  $\|\mathbf{g}(\mathbf{x}) - \mathbf{g}_0\|_\infty$  by some  $\epsilon_2$ , uniformly with probability  $1 - m^{3-\tau}$ . Recall from (2.19) that the  $j^{\text{th}}$  block of  $\mathbf{\Gamma}_\gamma \in \mathbb{R}^{m^2 \times m^2}$  has  $(k, \ell)$ -th entry

$$n^{-1} \sum_{i=1}^n X_k^{(i)} X_\ell^{(i)} h_j \left( X_j^{(i)} \right) + \gamma_{kj} \cdot \mathbf{1}_{\{k=\ell\}}.$$

The entry in  $\mathbf{g} \in \mathbb{R}^{m^2}$  (obtained by linearizing a  $m \times m$  matrix) corresponding to  $(j, k)$  is

$$n^{-1} \sum_{i=1}^n X_k^{(i)} h'_j \left( X_j^{(i)} \right) + n^{-1} \mathbf{1}_{\{j=k\}} \sum_{i=1}^n h_j \left( X_j^{(i)} \right).$$

Denote  $M \equiv \max_j \sup_{x>0} h_j(x)$ ,  $M' \equiv \max_j \sup_{x>0} h'_j(x)$ , and  $c_{\mathbf{X}} \equiv 2 \max_j (2\sqrt{\Sigma_{jj}} + \sqrt{e} \mathbb{E}_0 X_j)$ . Using results for sub-Gaussian random variables from Lemma 36.2 in Appendix A.2, we have for any  $t_1 > 0$ ,

$$\mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_k^{(i)} X_\ell^{(i)} h_j \left( X_j^{(i)} \right) - \mathbb{E}_0 X_k X_\ell h_j(X_j) \right| > t_1 \right) \leq 2 \exp \left( - \min \left( \frac{nt_1^2}{2M^2 c_{\mathbf{X}}^4}, \frac{nt_1}{2M c_{\mathbf{X}}^2} \right) \right).$$

Thus, choosing  $\epsilon_1 \equiv 2M c_{\mathbf{X}}^2 c_{n,m}$ , where  $c_{n,m} \equiv \max \left\{ \frac{2(\log m^\tau + \log 6)}{n}, \sqrt{\frac{2(\log m^\tau + \log 6)}{n}} \right\}$ , for  $\gamma_{kj} \leq \epsilon_1/2$ , we have

$$\begin{aligned} & \mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_k^{(i)} X_\ell^{(i)} h_j \left( X_j^{(i)} \right) + \gamma_{kj} \mathbf{1}_{\{k=\ell\}} - \mathbb{E}_0 X_k X_\ell h_j(X_j) \right| > \epsilon_1 \right) \\ & \leq \mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_k^{(i)} X_\ell^{(i)} h_j \left( X_j^{(i)} \right) - \mathbb{E}_0 X_k X_\ell h_j(X_j) \right| > \epsilon_1/2 \right) \end{aligned} \quad (\text{A.6})$$

$$\leq 2 \exp \left( - \min \left( \frac{n\epsilon_1^2}{8M^2c_{\mathbf{X}}^4}, \frac{n\epsilon_1}{4Mc_{\mathbf{X}}^2} \right) \right) \leq \frac{1}{3m^\tau}. \quad (\text{A.7})$$

Denote the event inside the probability in (A.6) as  $\mathcal{E}_{k,\ell,j}(\epsilon_1/2)$ . By definition,

$$c_{\mathbf{X}}^2 = 4 \max_k \left( 4\Sigma_{kk} + 4\sqrt{e}\sqrt{\Sigma_{kk}} \mathbb{E}_0 X_k + e(\mathbb{E}_0 X_k)^2 \right) \geq 4e \max_k \left( \Sigma_{kk} + (\mathbb{E}_0 X_k)^2 \right).$$

By Lemmas 35.2 and 36.1 from Appendix A.2,  $\text{var}(X_k) \leq \Sigma_{kk}$ , so  $c_{\mathbf{X}}^2 \geq 4e \max_k \mathbb{E}_0 X_k^2 \geq 4e \mathbb{E}_0 X_k^2 h_j(X_j)/M$ . Thus, setting  $\epsilon_1 = 2Mc_{\mathbf{X}}^2 c_{n,m}$ , on the complement of  $\mathcal{E}_{k,k,j}(\epsilon_1/2)$  we have

$$n^{-1} \sum_{i=1}^n X_k^{(i)2} h_j \left( X_j^{(i)} \right) \leq \mathbb{E}_0 X_k^2 h_j(X_j) + \epsilon_1/2 \leq \frac{\epsilon_1}{2} \left( 1 + \frac{1}{4ec_{n,m}} \right).$$

Then

$$\frac{1}{1 + 1/(4ec_{n,m})} \frac{1}{n} \sum_{i=1}^n X_k^{(i)2} h_j \left( X_j^{(i)} \right) \leq \epsilon_1/2 \quad (\text{A.8})$$

on the complement of  $\mathcal{E}_{k,k,j}(\epsilon_1/2)$ , again with  $c_{n,m} \equiv \max \left\{ \frac{2(\log m^\tau + \log 6)}{n}, \sqrt{\frac{2(\log m^\tau + \log 6)}{n}} \right\}$ . Note that the multiplier on the left of (A.8) is increasing in  $c_{n,m}$ , and that  $2(\log m^\tau + \log 6) > 6 \log m$  by the assumption that  $\tau > 3$ . Thus, if we let

$$\gamma_{kj} \equiv \frac{1}{1 + 1/\left(4e \max \left\{ 6 \log m/n, \sqrt{6 \log m/n} \right\}\right)} \frac{1}{n} \sum_{i=1}^n X_k^{(i)2} h_j \left( X_j^{(i)} \right),$$

which is just a constant multiple of the  $(k, k)$ -th entry of  $\mathbf{\Gamma}_j$  itself, with the constant explicitly calculable and a function of  $p$  and  $n$  only, then for  $k = \ell$

$$\mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_k^{(i)} X_\ell^{(i)} h_j \left( X_j^{(i)} \right) + \gamma_{kj} - \mathbb{E}_0 X_k X_\ell h_j(X_j) \right| \geq \epsilon_1 \right) \leq \mathbb{P}(\mathcal{E}_{k,\ell,j}(\epsilon_1/2)) \leq \frac{1}{3m^\tau}.$$

Since this also holds for  $k \neq \ell$  without the  $\gamma_{kj}$  term, by a union bound over  $m^3$  events,

$$\mathbb{P} \left( \max_{j,k,\ell} \left| n^{-1} \sum_{i=1}^n X_k^{(i)} X_\ell^{(i)} h_j \left( X_j^{(i)} \right) + \gamma_{kj} \mathbf{1}_{\{k=\ell\}} - \mathbb{E}_0 X_k X_\ell h_j(X_j) \right| \geq \epsilon_1 \right) \leq \frac{1}{3m^{\tau-3}}. \quad (\text{A.9})$$

Now, on the other hand, Lemma 36.1 and Hoeffding's inequality give for any  $t_{2,1}, t_{2,2} > 0$  that

$$\mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_k^{(i)} h'_j \left( X_j^{(i)} \right) - \mathbb{E}_0 X_k h'_j(X_j) \right| \geq t_{2,1} \right) \leq 2 \exp \left( - \frac{nt_{2,1}^2}{2M^2c_{\mathbf{X}}^2} \right),$$

$$\mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n h_j \left( X_j^{(i)} \right) - \mathbb{E}_0 h_j(X_j) \right| \geq t_{2,2} \right) \leq 2 \exp \left( -2nt_{2,2}^2/M^2 \right).$$

Choosing  $\epsilon_{2,1} \equiv \sqrt{2M'c_{\mathbf{X}} \sqrt{\frac{\log m^{\tau-1} + \log 6}{n}}}$ ,  $\epsilon_{2,2} \equiv M \sqrt{\frac{\log m^{\tau-2} + \log 6}{2n}}$  and taking union bounds over  $m^2$ , and  $m$  events, respectively, we have

$$\mathbb{P} \left( \max_{j,k} \left| n^{-1} \sum_{i=1}^n X_k^{(i)} h'_j \left( X_j^{(i)} \right) - \mathbb{E}_0 X_k h'_j(X_j) \right| \geq \epsilon_{2,1} \right) \leq \frac{1}{3m^{\tau-3}}, \quad (\text{A.10})$$

$$\mathbb{P} \left( \max_j \left| n^{-1} \sum_{i=1}^n h_j \left( X_j^{(i)} \right) - \mathbb{E}_0 h_j(X_j) \right| \geq \epsilon_{2,2} \right) \leq \frac{1}{3m^{\tau-3}}. \quad (\text{A.11})$$

Hence, by (A.9) (A.10) (A.11), with probability at least  $1 - m^{3-\tau}$ ,  $\|\Gamma_\gamma(\mathbf{x}) - \Gamma_0\|_\infty < \epsilon_1$  and  $\|\mathbf{g}(\mathbf{x}) - \mathbf{g}_0\|_\infty < \epsilon_2 \equiv \epsilon_{2,1} + \epsilon_{2,2}$ . Consider any  $\tau > 3$ , and let

$$\begin{aligned} c_2 &\equiv \frac{6}{\alpha} c_{\Gamma_0}, \\ n &> \max \left\{ 2M^2 c_{\mathbf{X}}^4 c_2^2 d_{\mathbf{K}_0}^2 (\tau \log m + \log 6), 2M c_{\mathbf{X}}^2 c_2 d_{\mathbf{K}_0} (\tau \log m + \log 6) \right\}, \\ \lambda &> \frac{3(2-\alpha)}{\alpha} \max \{ c_{\mathbf{K}_0} \epsilon_1, \epsilon_2 \} \\ &\equiv \frac{3(2-\alpha)}{\alpha} \max \left\{ 4M c_{\mathbf{K}_0} c_{\mathbf{X}}^2 \frac{(\log m^\tau + \log 6)}{n}, \right. \\ &\quad \left. 2M c_{\mathbf{K}_0} c_{\mathbf{X}}^2 \sqrt{\frac{2(\log m^\tau + \log 6)}{n}}, \sqrt{2M'c_{\mathbf{X}} \sqrt{\frac{\log m^{\tau-1} + \log 6}{n}}} + M \sqrt{\frac{\log m^{\tau-2} + \log 6}{2n}} \right\}. \end{aligned}$$

Then  $d_{\mathbf{K}_0} \epsilon_1 \leq \alpha/(6c_{\Gamma_0})$  and the results follow from Theorem 7.  $\square$

*Proof of Theorem 11.* Similar to the proof of Theorem 10, by Theorem 7 it suffices to prove that for any  $\tau > 3$ , we can bound  $\|\Gamma_\gamma(\mathbf{x}) - \Gamma_0\|_\infty$  by some  $\epsilon_1$  and  $\|\mathbf{g}(\mathbf{x}) - \mathbf{g}_0\|_\infty$  by some  $\epsilon_2$ , uniformly with probability  $1 - m^{3-\tau}$ . Recall that  $\Gamma \in \mathbb{R}^{(m^2+m) \times (m^2+m)}$  is a rearrangement of  $\Gamma^{(*)}$ , which is in turn formed by  $\Gamma_{11} \in \mathbb{R}^{m^2 \times m^2}$ ,  $\Gamma_{12} \in \mathbb{R}^{m^2 \times m}$ ,  $\Gamma_{12}^\top$  and  $\Gamma_{22} \in \mathbb{R}^{m \times m}$ , all of which are block-diagonal with  $m$  blocks.

The  $j^{\text{th}}$  block of  $\Gamma_{11} \in \mathbb{R}^{m^2 \times m^2}$  has  $(k, \ell)$ -th entry

$$n^{-1} \sum_{i=1}^n X_k^{(i)} X_\ell^{(i)} h_j \left( X_j^{(i)} \right),$$

the  $k^{\text{th}}$  entry in the  $j^{\text{th}}$  block of  $\mathbf{\Gamma}_{12}$  is

$$-n^{-1} \sum_{i=1}^n X_k^{(i)} h_j \left( X_j^{(i)} \right),$$

the  $j^{\text{th}}$  diagonal entry of  $\mathbf{\Gamma}_{22}$  is

$$n^{-1} \sum_{i=1}^n h_j \left( X_j^{(i)} \right).$$

On the other hand,  $\mathbf{g} \in \mathbb{R}^{(m^2+m)}$  is a rearrangement of  $\mathbf{g}^{(*)} \equiv [\mathbf{g}_1^\top, \mathbf{g}_2^\top]^\top$ , where the entry in  $\mathbf{g}_1 \in \mathbb{R}^{m^2}$  (obtained by linearizing a  $m \times m$  matrix) corresponding to  $(j, k)$ , is

$$n^{-1} \sum_{i=1}^n X_k^{(i)} h'_j \left( X_j^{(i)} \right) + n^{-1} \mathbf{1}_{\{j=k\}} \sum_{i=1}^n h_j \left( X_j^{(i)} \right),$$

while the  $j$ -th component of  $\mathbf{g}_2 \in \mathbb{R}^m$  is

$$-n^{-1} \sum_{i=1}^n h'_j \left( X_j^{(i)} \right).$$

Recalling that the bounds in Lemma 36 also hold when  $\boldsymbol{\mu} \neq 0$ , we may then use bounds similar to those in the proof of Theorem 10, and use union bounds to arrive at analogous consistency results, modulus different constants. The amplifiers  $\gamma$  can be incorporated analogously.  $\square$

## A.2 Auxiliary Lemmas and Definitions

In this appendix, to simplify notation, when it is clear from the context, the operator  $\mathbb{E}$  is defined as the expectation under the true distribution, unless otherwise noted.

**Definition 16** (Sub-Gaussian and Sub-Exponential Variables).

The sub-Gaussian ( $r = 2$ ) and sub-exponential ( $r = 1$ ) norms of a random variable are

$$\|X\|_{\psi_r} \equiv \sup_{q \geq 1} q^{-1/r} (\mathbb{E}|X|^{rq})^{1/(rq)} \equiv \sup_{q \geq 1} q^{-1/r} \|X\|_{rq}.$$

If  $\|X\|_{\psi_2} < \infty$  we say  $X$  is sub-Gaussian; if  $\|X\|_{\psi_1} < \infty$  we call  $X$  sub-exponential. For a zero-mean sub-Gaussian random variable  $X$  also define the sub-Gaussian parameter

$$\tau(X) = \inf\{\tau \geq 0 : \mathbb{E} \exp(tX) \leq \exp(\tau^2 t^2/2), \forall t \in \mathbb{R}\}.$$

The definition of sub-Gaussian norm here allows for a non-centered variable and differs from the one in Vershynin (2012), which uses  $\|X\|_q$ . Instead, it coincides with  $\theta_2$  in Buldygin and Kozachenko (2000). The sub-Gaussian parameter is defined as in Buldygin and Kozachenko (2000) and the sub-exponential norm as in Vershynin (2012).

**Lemma 35** (Properties of Sub-Gaussian and Sub-Exponential Variables).

1) For any  $X$  and  $r = 1, 2$ ,  $\|X - \mathbb{E}X\|_{\psi_r} \leq 2\|X\|_{\psi_r}$  and  $\|X\|_{\psi_r} \leq \|X - \mathbb{E}X\|_{\psi_r} + |\mathbb{E}X|$ , as long as the expectation and norms are finite.

2) (Buldygin and Kozachenko, 2000)  $\tau(X)$  is a norm on the space of all zero-mean sub-Gaussian variables; so  $\tau(X + Y) \leq \tau(X) + \tau(Y)$ . If  $X$  is zero-mean sub-Gaussian, then  $\text{var}(X) \leq \tau^2(X)$ ,  $\|X\|_{\psi_2} \leq 2\tau(X)/\sqrt{e}$ ,  $\tau(X) \leq \sqrt{e}\|X\|_{\psi_2}$ . If  $X_1, \dots, X_n$  are i.i.d. zero-mean sub-Gaussian,  $\tau(n^{-1} \sum_{i=1}^n X_i) \leq n^{-1/2}\tau(X_i)$ .

3) If  $X_1$  and  $X_2$  are sub-Gaussian (not necessarily independent) with  $\|X_1\|_{\psi_2} \leq K_1$  and  $\|X_2\|_{\psi_2} \leq K_2$ , then  $X_1X_2$  is sub-exponential with  $\|X_1X_2\|_{\psi_1} \leq K_1K_2$ .

4) (Buldygin and Kozachenko, 2000) If  $X$  is zero-mean sub-Gaussian and  $q > 0$ , then

$$\mathbb{E}|X|^q \leq 2(q/e)^{q/2}\tau^q(X).$$

5) (Buldygin and Kozachenko, 2000) If  $X_1, \dots, X_n$  are independent zero-mean, sub-Gaussian variables, then for any  $\epsilon > 0$ ,

$$\begin{aligned} \mathbb{P}(|X_1| \geq \epsilon) &\leq 2 \exp\left(-\frac{\epsilon^2}{2\tau^2(X_1)}\right), \\ \mathbb{P}\left(\left|n^{-1} \sum_{i=1}^n X_i\right| > \epsilon\right) &\leq 2 \exp\left(-\frac{n\epsilon^2}{2 \max_i \tau^2(X_i)}\right). \end{aligned}$$

6) (Vershynin, 2012) If  $X_1, \dots, X_n$  are independent zero-mean sub-exponential random variables with  $K \geq \max_i \|X_i\|_{\psi_1}$ , then for any  $\epsilon > 0$ ,

$$\mathbb{P}(|X_1| \geq \epsilon) \leq 2 \exp\left(-\min\left(\frac{\epsilon^2}{8e^2K^2}, \frac{\epsilon}{4eK}\right)\right),$$

$$\mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_i \right| \geq \epsilon \right) \leq 2 \exp \left( - \min \left( \frac{n\epsilon^2}{8e^2 K^2}, \frac{n\epsilon}{4eK} \right) \right).$$

7) (Boucheron et al., 2013) If for  $X_i$  i.i.d. there exists some  $B > 0$  such that

$$\sup_{q \geq 2} \left( \frac{\mathbb{E}|X|^q}{q!} \right)^{1/q} \leq B/2$$

then for all  $\epsilon > 0$ ,

$$\mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n (X_i - \mathbb{E}X_i) \right| \geq \epsilon \right) \leq 2 \exp \left( - \min \left( \frac{n\epsilon^2}{2B^2}, \frac{n\epsilon}{2B} \right) \right).$$

*Proof.* 1) For  $r = 1, 2$ , by the triangle inequality,  $\|X - \mathbb{E}X\|_{\psi_r} \leq \|X\|_{\psi_r} + \|\mathbb{E}X\|_{\psi_r} = \|X\|_{\psi_r} + |\mathbb{E}X| \leq \|X\|_{\psi_r} + \mathbb{E}|X| \leq 2\|X\|_{\psi_r}$ , where in the last step we used the definition of  $\|\cdot\|_{\psi_r}$  with  $q = 1$  for  $r = 1$  and  $\mathbb{E}|X| \leq (\mathbb{E}|X|^2)^{1/2}$  with  $q = 2$  for  $r = 2$ . On the other hand,  $\|X\|_{\psi_r} \leq \|X - \mathbb{E}X\|_{\psi_r} + \|\mathbb{E}X\|_{\psi_r} = \|X - \mathbb{E}X\|_{\psi_r} + |\mathbb{E}X|$ .

2) These follow from Theorems 1.2 and 1.3 and Lemmas 1.2 and 1.7 from Buldygin and Kozachenko (2000), and  $\sqrt[4]{3.1}e^{9/16}/\sqrt{2} \approx 1.6467 \leq 1.6487 \approx \sqrt{e}$ .

3) By Hölder's inequality (or Cauchy-Schwarz),

$$\begin{aligned} \|X_1 X_2\|_{\psi_1} &= \sup_{q \geq 1} q^{-1} (\mathbb{E}|X_1 X_2|^q)^{1/q} = \sup_{q \geq 1} q^{-1} (\mathbb{E}|X_1^q X_2^q|)^{1/q} \\ &\leq \sup_{q \geq 1} q^{-1} [(\mathbb{E}|X_1|^{2q})^{1/2} (\mathbb{E}|X_2|^{2q})^{1/2}]^{1/q} \\ &\leq \sup_{q \geq 1} [q^{-1/2} (\mathbb{E}|X_1|^{2q})^{1/2q}] \sup_{q \geq 1} [q^{-1/2} (\mathbb{E}|X_2|^{2q})^{1/2q}] \\ &= \|X_1\|_{\psi_2} \|X_2\|_{\psi_2} \leq K_1 K_2. \end{aligned}$$

4-6) These are Lemma 1.4 and Theorem 1.5 in Buldygin and Kozachenko (2000), and a consequence of Corollary 5.17 in Vershynin (2012).

7) By Theorem 2.10 of Boucheron et al. (2013) wherein we let  $v \equiv nB^2/2$  and  $c \equiv B/2$ , we have

$$\mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n (X_i - \mathbb{E}X_i) \right| \geq \epsilon \right) \leq 2 \exp \left( - \frac{n\epsilon^2}{B^2 + B\epsilon} \right)$$

for all  $\epsilon > 0$ . (Theorem 2.10 gives an one-sided bound; bound for the other side is obtained by taking  $X_i = -X_i$ ). The inequality follows by splitting into cases  $\epsilon \leq B$  and  $\epsilon > B$ .  $\square$

**Lemma 36.** *Suppose  $\mathbf{X}$  follows a truncated normal distribution on  $\mathbb{R}_+^m$  with parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma} = \mathbf{K}^{-1} \succ \mathbf{0}$ . Let  $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(n)}$  be i.i.d. copies of  $\mathbf{X}$ , with  $j$ -th component of the  $i$ -th copy being  $X_j^{(i)}$ . Then*

1. For  $j = 1, \dots, p$ ,  $\tau(X_j - \mathbb{E}X_j) \leq \sqrt{\Sigma_{jj}}$ . That is, the sub-Gaussian parameter of any marginal distribution of  $\mathbf{X}$ , after centering, is bounded by the square root of its corresponding diagonal entry in the covariance parameter  $\boldsymbol{\Sigma}$ . Then, for any  $\epsilon > 0$ ,

$$\mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_j^{(i)} - \mathbb{E}X_j \right| > \epsilon \right) \leq 2 \exp \left( -\frac{n\epsilon^2}{2\Sigma_{jj}} \right).$$

In particular, if  $h_0$  is a function bounded by  $M_0$ , then for any  $\epsilon > 0$ ,

$$\begin{aligned} \mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_j^{(i)} h_0 \left( X_k^{(i)} \right) - \mathbb{E}X_j h_0(X_k) \right| \geq \epsilon \right) &\leq 2 \exp \left( -\frac{n\epsilon^2}{8M_0^2(2\sqrt{\Sigma_{jj}} + \sqrt{e} \mathbb{E}X_j)^2} \right), \\ \tau \left( n^{-1} \sum_{i=1}^n X_j^{(i)} h_0 \left( X_k^{(i)} \right) - \mathbb{E}X_j h_0(X_k) \right) &\leq \frac{2M_0}{\sqrt{n}} \left( 2\sqrt{\Sigma_{jj}} + \sqrt{e} \mathbb{E}X_j \right), \\ \left\| n^{-1} \sum_{i=1}^n X_j^{(i)} h_0 \left( X_k^{(i)} \right) - \mathbb{E}X_j h_0(X_k) \right\|_{\psi_2} &\leq \frac{4M_0}{\sqrt{en}} \left( 2\sqrt{\Sigma_{jj}} + \sqrt{e} \mathbb{E}X_j \right). \end{aligned}$$

2. For  $j, k, \ell \in \{1, \dots, p\}$ , if  $h_0$  is a function bounded by  $M_0$ , then

$$\|X_j X_k h_0(X_\ell) - \mathbb{E}X_j X_k h_0(X_\ell)\|_{\psi_1} \leq \frac{M_0}{2e} c_{\mathbf{X}}^2, \quad (\text{A.12})$$

where  $c_{\mathbf{X}} \equiv 2 \max_j (2\sqrt{\Sigma_{jj}} + \sqrt{e} \mathbb{E}X_j)$ . In particular, for any  $\epsilon > 0$ ,

$$\begin{aligned} \mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_j^{(i)} X_k^{(i)} h_0 \left( X_\ell^{(i)} \right) - \mathbb{E}X_j X_k h_0(X_\ell) \right| > \epsilon \right) \\ \leq 2 \exp \left( -\min \left( \frac{n\epsilon^2}{2M_0^2 c_{\mathbf{X}}^4}, \frac{n\epsilon}{2M_0 c_{\mathbf{X}}^2} \right) \right). \end{aligned}$$

*Proof of Lemma 36.* 1. Without loss of generality choose  $j = 1$ . By the definition of sub-Gaussian parameters, we need to show that for all  $t \in \mathbb{R}$ ,

$$\mathbb{E} \exp(tX_1) \leq \exp(t^2 \Sigma_{11}/2 + t \mathbb{E} X_1),$$

which is equivalent to

$$t^2 \Sigma_{11}/2 + t \mathbb{E} X_1 - \log \mathbb{E} \exp(tX_1) \geq 0 \quad \forall t \in \mathbb{R}. \quad (\text{A.13})$$

Since the left-hand side of (A.13) equals 0 at  $t = 0$ , it suffices to show that its derivative,

$$t \Sigma_{11} + \mathbb{E} X_1 - \frac{d \log \mathbb{E} \exp(tX_1)}{dt} = t \Sigma_{11} + \mathbb{E} X_1 - \frac{\frac{d \mathbb{E} \exp(tX_1)}{dt}}{\mathbb{E} \exp(tX_1)}, \quad (\text{A.14})$$

is non-negative on  $(0, \infty)$  and non-positive on  $(-\infty, 0)$ . By properties of moment-generating functions,  $\frac{d}{dt} \mathbb{E} \exp(tX_1)$  evaluated at  $t = 0$  equals  $\mathbb{E} X_1$ , so (A.14) equals 0 at  $t = 0$ . It in turn suffices to show the derivative of (A.14), namely

$$\Sigma_{11} - \frac{d^2 \log \mathbb{E} \exp(tX_1)}{dt^2} \quad (\text{A.15})$$

is non-negative in  $t \in \mathbb{R}$ .

Given any vector  $\mathbf{v} \in \mathbb{R}^p$ , define  $\mathbb{R}_+^p - \mathbf{v} \equiv \{\mathbf{u} - \mathbf{v} : \mathbf{u} \in \mathbb{R}_+^p\}$ . By Tallis (1961), denoting the first column of  $\Sigma$  as  $\Sigma_1$ , the moment-generating function of the marginal distribution of  $X_1$  is

$$\frac{\int_{\mathbb{R}_+^p - (\boldsymbol{\mu} + t \Sigma_1)} \exp\left(-\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x}\right) d\mathbf{x}}{\int_{\mathbb{R}_+^p - \boldsymbol{\mu}} \exp\left(-\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x}\right) d\mathbf{x}} \exp\left(t \mu_1 + \frac{1}{2} t^2 \Sigma_{11}\right).$$

(A.15) thus becomes

$$-\frac{d^2}{dt^2} \log \int_{\mathbb{R}_+^p - (\boldsymbol{\mu} + t \Sigma_1)} \exp\left(-\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x}\right) d\mathbf{x}.$$

Showing this is non-negative in  $t \in \mathbb{R}$  is equivalent to showing that the integral itself is log-concave in  $t$ . But

$$\int_{\mathbb{R}_+^p - (\boldsymbol{\mu} + t \Sigma_1)} \exp\left(-\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x}\right) d\mathbf{x} = \int_{\mathbb{R}^p} \exp\left(-\frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x}\right) \mathbf{1}_{\mathbb{R}_+^p - \boldsymbol{\mu}}(\mathbf{x} + t \Sigma_1) d\mathbf{x}$$

with  $\exp\left(-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1}\mathbf{x}\right)$  log-concave in  $\mathbf{x}$  and  $\mathbf{1}_{\mathbb{R}_+^p - \boldsymbol{\mu}}(\mathbf{x} + t\Sigma_1)$  log-concave in  $(\mathbf{x}, t)$  since  $\mathbb{R}_+^p - \boldsymbol{\mu}$  is a convex set. Since log-concavity is closed under multiplication and integration over  $\mathbb{R}^p$ , the integral is indeed log-concave, and our proof of the bound on the sub-Gaussian parameter of  $X_j - \mathbb{E}X_j$  is complete. The tail bound follows from 5) of Lemma 35.

Now by 1) and 2) of Lemma 35,

$$\|X_j\|_{\psi_2} \leq 2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j.$$

If  $h_0$  is a function bounded by  $M_0$ , then by definition

$$\|X_j h_0(X_k)\|_{\psi_2} \leq M_0 \left(2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j\right).$$

By 1) and 2) of Lemma 35 again,

$$\begin{aligned} \tau(X_j h_0(X_k) - \mathbb{E}X_j h_0(X_k)) &\leq \sqrt{e} \|X_j h_0(X_k) - \mathbb{E}X_j h_0(X_k)\|_{\psi_2} \\ &\leq 2\sqrt{e} \|X_j h_0(X_k)\|_{\psi_2} \\ &\leq 2M_0(2\sqrt{\Sigma_{jj}} + \sqrt{e} \mathbb{E}X_j). \end{aligned}$$

The tail bound thus follows from the first inequality using 5) of Lemma 35. By 2) of the Lemma 35,

$$\begin{aligned} \tau\left(n^{-1} \sum_{i=1}^n X_j^{(i)} h_0\left(X_k^{(i)}\right) - \mathbb{E}X_j h_0(X_k)\right) &\leq \frac{2M_0}{\sqrt{n}} \left(2\sqrt{\Sigma_{jj}} + \sqrt{e} \mathbb{E}X_j\right), \\ \left\|n^{-1} \sum_{i=1}^n X_j^{(i)} h_0\left(X_k^{(i)}\right) - \mathbb{E}X_j h_0(X_k)\right\|_{\psi_2} &\leq \frac{4M_0}{\sqrt{en}} \left(2\sqrt{\Sigma_{jj}} + \sqrt{e} \mathbb{E}X_j\right). \end{aligned}$$

2. By the proof of 1) of this lemma,  $\|X_j\|_{\psi_2} \leq 2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j$ , and by 3) of Lemma 35,

$$\|X_j X_k\|_{\psi_1} \leq \left(2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j\right) \left(2\sqrt{\Sigma_{kk}/e} + \mathbb{E}X_k\right) \leq \max_j \left(2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j\right)^2.$$

Since  $h_0$  is a function bounded by  $M_0$ , by definition

$$\|X_j X_k h_0(X_\ell)\|_{\psi_1} \leq M_0 \max_j \left(2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j\right)^2.$$

Then by 1) of Lemma 35 again,

$$\|X_j X_k h_0(X_\ell) - \mathbb{E}X_j X_k h_0(X_\ell)\|_{\psi_1} \leq 2M_0 \max_j \left( 2\sqrt{\Sigma_{jj}/e} + \mathbb{E}X_j \right)^2.$$

The tail bound then follows from 6) of Lemma 35.  $\square$

Although not used for our consistency results, in the special case of  $h_0 \equiv 1$ , we also have the following lemma. The notable difference between bounds (A.16) below and (A.12) from Lemma 36.2 is in the constants and dependency on  $\mathbb{E}X_j$ : The constants in the denominator in the right-hand side of (A.12) is smaller and thus gives a tighter bound, but (A.16) is preferred when  $\mathbb{E}X_j$  is notably large compared to  $\sqrt{\Sigma_{jj}}$ , since the constant is only linear in  $\mathbb{E}X_j$ .

**Lemma 37.** *Consider the setting in Lemma 36. Then for  $j, k \in \{1, \dots, p\}$ , for any  $\epsilon > 0$ ,*

$$\mathbb{P} \left( n^{-1} \left| \sum_{i=1}^n X_j^{(i)} X_k^{(i)} - \mathbb{E}X_j X_k \right| \geq \epsilon \right) \leq 4 \exp \left( - \min \left( \frac{2n\epsilon^2}{C_1^2}, \frac{n\epsilon}{C_1} \right) \right), \quad (\text{A.16})$$

where  $C_1 \equiv 91 \max_j \Sigma_{jj} + 72 \max_j \mathbb{E}X_j \max_j \sqrt{\Sigma_{jj}}$ .

*Proof of Lemma 37.* We use a proof similar to Lemma 1 in Ravikumar et al. (2011) (note that  $\mathbb{E}X_j$  may be nonzero in our case). Define

$$U_{jk}^{(i)} \equiv X_j^{(i)} + X_k^{(i)}, \quad U_{jk} \equiv X_j + X_k, \quad V_{jk}^{(i)} \equiv X_j^{(i)} - X_k^{(i)}, \quad V_{jk} \equiv X_j - X_k.$$

Since  $X_j^{(i)} X_k^{(i)} = \frac{1}{4} \left( U_{jk}^{(i)2} - V_{jk}^{(i)2} \right)$ , by union bound we have

$$\begin{aligned} & \mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n X_j^{(i)} X_k^{(i)} - \mathbb{E}X_j X_k \right| \geq \epsilon \right) \\ & \leq \mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n U_{jk}^{(i)2} - \mathbb{E}U_{jk}^2 \right| \geq 2\epsilon \right) + \mathbb{P} \left( \left| n^{-1} \sum_{i=1}^n V_{jk}^{(i)2} - \mathbb{E}V_{jk}^2 \right| \geq 2\epsilon \right). \end{aligned}$$

We next define

$$Z_{jk}^{(i)} \equiv U_{jk}^{(i)2} - \mathbb{E}U_{jk}^2 = A_{jk}^{(i)2} + B_{jk}^{(i)} + C_{jk}, \quad \bar{X}_j^{(i)} \equiv X_j^{(i)} - \mathbb{E}X_j,$$

$$A_{jk}^{(i)} \equiv \bar{X}_j^{(i)} + \bar{X}_k^{(i)}, \quad B_{jk}^{(i)} \equiv 2(\mathbb{E}X_j + \mathbb{E}X_k)(\bar{X}_j^{(i)} + \bar{X}_k^{(i)}), \quad C_{jk} \equiv -\mathbb{E}(\bar{X}_j^{(i)} + \bar{X}_k^{(i)})^2.$$

Then since  $\tau$  is a norm by 2) of Lemma 35,  $A_{jk}$  is sub-Gaussian with parameter  $\leq \sqrt{\Sigma_{jj}} + \sqrt{\Sigma_{kk}}$ , and  $B_{jk}$  is sub-Gaussian with parameter  $\leq 2(\mathbb{E}X_j + \mathbb{E}X_k)(\sqrt{\Sigma_{jj}} + \sqrt{\Sigma_{kk}})$ . Using 4) of Lemma 35 together with the inequality  $(a + b + c)^q \leq (3 \max\{a, b, c\})^q \leq 3^q(a^q + b^q + c^q)$  for all  $a, b, c \geq 0$  and  $q > 0$ , we have for any  $q \geq 2$

$$\begin{aligned} (\mathbb{E}|Z_{jk}|^q)^{1/q} &\leq (3^q (\mathbb{E}|A_{jk}|^{2q} + \mathbb{E}|B_{jk}|^q + |C_{jk}|^q))^{1/q} \\ &\leq 3^{1+1/q} ((\mathbb{E}|A_{jk}|^{2q})^{1/q} + (\mathbb{E}|B_{jk}|^q)^{1/q} + |C_{jk}|) \\ &\leq 3^{1+1/q} \left( 2^{1/q}(2q/e) \left( \sqrt{\Sigma_{jj}} + \sqrt{\Sigma_{kk}} \right)^2 \right. \\ &\quad \left. + 2^{1/q} \sqrt{q/e} 2(\mathbb{E}X_j + \mathbb{E}X_k)(\sqrt{\Sigma_{jj}} + \sqrt{\Sigma_{kk}}) + \text{var}(X_j + X_k) \right). \end{aligned}$$

Using  $\text{var}(X + Y) \leq 2(\text{var}(X) + \text{var}(Y))$  and the fact that  $\text{var}(X_j) = \text{var}(X_j - \mathbb{E}X_j) \leq \tau^2(X_j - \mathbb{E}X_j) \leq \Sigma_{jj}$  (by 2) of Lemma 35 and 1) of Lemma 36, we then have

$$\begin{aligned} \left( \frac{\mathbb{E}|Z_{jk}|^q}{q!} \right)^{1/q} &\leq 3^{1+1/q} \frac{2^{3+1/q}(q/e) \max_j \Sigma_{jj} + 2^{3+1/q} \sqrt{q/e} \max_j \mathbb{E}X_j \cdot \max_j \sqrt{\Sigma_{jj}} + 4 \max_j \Sigma_{jj}}{(q!)^{1/q}}. \end{aligned}$$

Since all three coefficients involving  $q$  are decreasing in  $q \geq 2$ , we have

$$\sup_{q \geq 2} \left( \frac{\mathbb{E}|Z_{jk}|^q}{q!} \right)^{1/q} \leq (48\sqrt{3}/e + 6\sqrt{6}) \max_j \Sigma_{jj} + 24\sqrt{6}/e \max_j \mathbb{E}X_j \max_j \sqrt{\Sigma_{jj}}.$$

Thus by 7) of Lemma 35, letting  $B \equiv (91 \max_j \Sigma_{jj} + 72 \max_j \mathbb{E}X_j \max_j \sqrt{\Sigma_{jj}})$ , we have for all  $\epsilon > 0$ :

$$\mathbb{P} \left( n^{-1} \left| \sum_{i=1}^n Z_{jk}^{(i)} \right| \geq 2\epsilon \right) \leq 2 \exp \left( - \min \left( \frac{2n\epsilon^2}{B^2}, \frac{n\epsilon}{B} \right) \right).$$

A tail bound for the sample average of  $V_{jk}^2$  can be similarly derived, and the result follows.  $\square$

### A.3 Simulation Results for Erdős-Rényi Graphs

We revisit the simulations from Section 2.7 but use Erdős-Rényi (ER) graphs in which each possible edge is independently included with probability  $\pi$ . Independent uniform draws from

$[0.5, 1]$  are used to fill the non-zero off-diagonal entries of the symmetric matrix  $\mathbf{K}_0$ . The diagonal elements are set such that  $\mathbf{K}_0$  has minimum eigenvalue 0.1. We choose  $\pi = 0.08$  for  $n = 1000$ , and  $\pi = 0.02$  for  $n = 80$ .

### A.3.1 Truncated GGMs

In this section we present the results for truncated GGMs.

#### Choice of $h$

The results for truncated centered GGMs are reported in Table A.1 and Figure A.1. Those for truncated non-centered GGMs using the profiled estimator are in Table A.2 and Figure A.2.

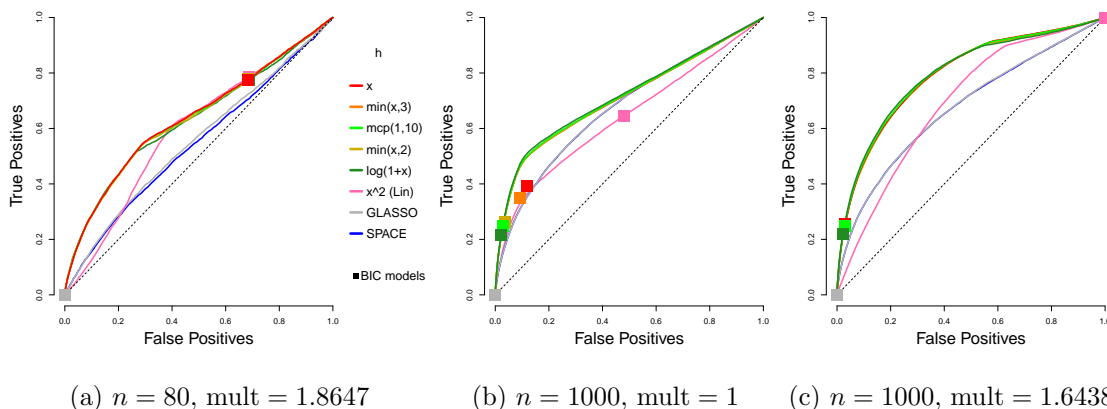


Figure A.1: Average ROC curves of our *centered* estimator for  $m = 100$  variables and two sample sizes  $n$  under various choices of  $h$ , compared to SPACE and GLASSO, for the *truncated centered GGM* case. Squares indicate average true positive rate (TPR) and false positive rate (FPR) of models picked by eBIC with refitting for the estimator in the same color.

Centered, $n = 80$ , multiplier 1.8647					
min(log(1 + $x$ ), $c$ )			min( $x$ , $c$ )		
$c$	Mean	sd	$c$	Mean	sd
$\infty$	0.632	0.036	$\infty$	0.638	0.035
2	0.632	0.036	3	0.638	0.035
1	0.630	0.035	2	0.635	0.035
0.5	0.613	0.033	1	0.623	0.033
MCP(1, $c$ )			SCAD(1, $c$ )		
$c$	Mean	sd	$c$	Mean	sd
10	0.637	0.035	10	0.638	0.035
5	0.636	0.036	5	0.637	0.035
1	0.617	0.033	2	0.632	0.035
$x^{1.5}$ : (0.627, 0.032)			$x^2$ : (0.595, 0.028)		
GLASSO (0.553, 0.029)			SPACE: (0.544, 0.026)		
NS: (0.543, 0.028)			SJ: (0.519, 0.028)		

Centered, $n = 1000$ , multiplier 1						Centered, $n = 1000$ , multiplier 1.6438					
min(log(1 + $x$ ), $c$ )			min( $x$ , $c$ )			min(log(1 + $x$ ), $c$ )			min( $x$ , $c$ )		
$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd
$\infty$	0.716	0.016	2	0.710	0.016	$\infty$	0.796	0.014	$\infty$	0.795	0.014
2	0.716	0.016	3	0.710	0.016	2	0.796	0.014	3	0.794	0.014
1	0.715	0.016	1	0.710	0.017	1	0.794	0.014	2	0.792	0.014
0.5	0.694	0.017	$\infty$	0.709	0.016	0.5	0.772	0.015	1	0.784	0.015
MCP(1, $c$ )			SCAD(1, $c$ )			MCP(1, $c$ )			SCAD(1, $c$ )		
$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd
5	0.714	0.016	2	0.713	0.016	5	0.796	0.014	5	0.795	0.014
10	0.711	0.016	5	0.711	0.016	10	0.796	0.014	10	0.795	0.014
1	0.707	0.017	10	0.710	0.016	1	0.778	0.015	2	0.793	0.014
$x^{1.5}$ : (0.678, 0.016)			$x^2$ : (0.64, 0.017)			$x^{1.5}$ : (0.757, 0.015)			$x^2$ : (0.693, 0.016)		
GLASSO: (0.675, 0.016)			SPACE: (0.675, 0.016)			GLASSO: (0.675, 0.016)			SPACE: (0.675, 0.016)		
NS: (0.675, 0.016)			SJ: (0.624, 0.017)			NS: (0.675, 0.016)			SJ: (0.624, 0.017)		

Table A.1: Mean and standard deviation of areas under the ROC curves (AUC) using different estimators in the centered setting, with  $n = 80$  and multiplier 1.8647, or  $n = 1000$  and multipliers 1 and 1.6438. Methods include our estimator with different choices of  $h$ , GLASSO, SPACE, neighborhood selection (NS), and Space JAM (SJ).

Non-centered profiled, $n = 80$ , multiplier 1.8647					
min(log(1 + x), c)			min(x, c)		
$c$	Mean	sd	$c$	Mean	sd
1	0.588	0.034	3	0.588	0.033
$\infty$	0.588	0.034	$\infty$	0.588	0.033
2	0.588	0.034	2	0.588	0.033
0.5	0.576	0.033	1	0.583	0.033
MCP(1, c)			SCAD(1, c)		
$c$	Mean	sd	$c$	Mean	sd
5	0.588	0.033	5	0.588	0.033
10	0.588	0.033	10	0.588	0.033
1	0.581	0.033	2	0.587	0.033
$x^{1.5}$ : (0.582,0.028)			$x^2$ : (0.576,0.028)		
GLASSO: (0.572,0.033)			SPACE: (0.562,0.031)		
NS: (0.560,0.032)			SJ: (0.535,0.027)		

Non-centered profiled, $n = 1000$ , multiplier 1						Non-centered profiled, $n = 1000$ , multiplier 1.6438					
min(log(1 + x), c)			min(x, c)			min(log(1 + x), c)			min(x, c)		
$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd
2	0.692	0.022	1	0.687	0.022	2	0.705	0.021	$\infty$	0.705	0.022
$\infty$	0.692	0.022	$\infty$	0.686	0.022	$\infty$	0.705	0.021	3	0.705	0.021
1	0.691	0.022	3	0.685	0.022	1	0.703	0.021	2	0.702	0.022
0.5	0.684	0.02	2	0.685	0.022	0.5	0.683	0.019	1	0.695	0.021
MCP(1, c)			SCAD(1, c)			MCP(1, c)			SCAD(1, c)		
$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd	$c$	Mean	sd
5	0.689	0.022	2	0.687	0.022	5	0.706	0.021	10	0.705	0.022
1	0.689	0.020	5	0.687	0.022	10	0.706	0.022	5	0.705	0.022
10	0.687	0.022	10	0.686	0.022	1	0.690	0.019	2	0.703	0.022
$x^{1.5}$ : (0.663,0.020)			$x^2$ : (0.638,0.019)			$x^{1.5}$ : (0.689,0.021)			$x^2$ : (0.664,0.019)		
GLASSO (0.700,0.022)			SPACE: (0.699,0.022)			GLASSO (0.700,0.022)			SPACE: (0.699,0.022)		
NS: (0.699,0.022)			SJ: (0.655,0.021)			NS: (0.699,0.022)			SJ: (0.655,0.021)		

Table A.2: Mean and standard deviation of AUC using different profiled estimators in the non-centered setting, with  $n = 80$  and multiplier 1.8647, or  $n = 1000$  and multipliers 1 and 1.6438. Methods as for Table A.1.

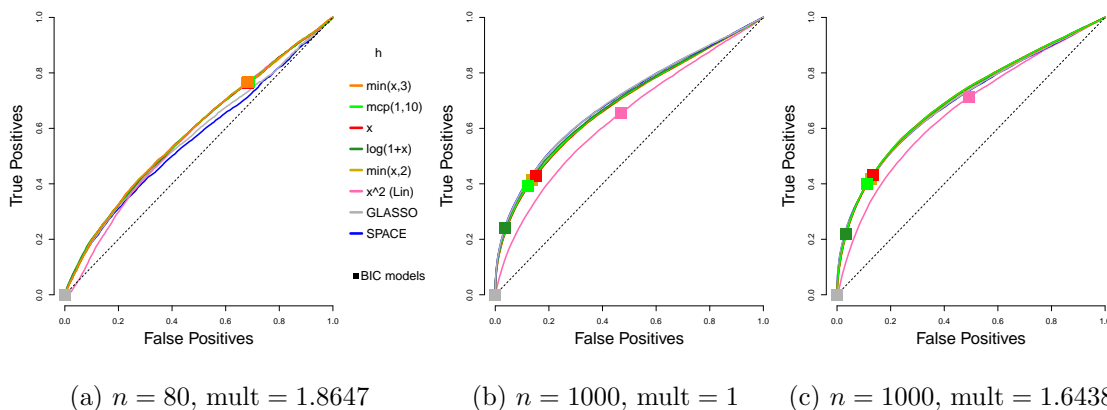


Figure A.2: Average for the truncated non-centered GGM case.  $n = 80$  or  $1000$ ,  $m = 100$ .

### Choice of multiplier

The results for truncated centered GGMs where each curve represents a different multiplier are shown in Figure A.3, and those for truncated non-centered GGMs are in Figure A.4, where each curve corresponds to a different ratio  $\lambda_{\mathbf{K}}/\lambda_{\eta}$ .

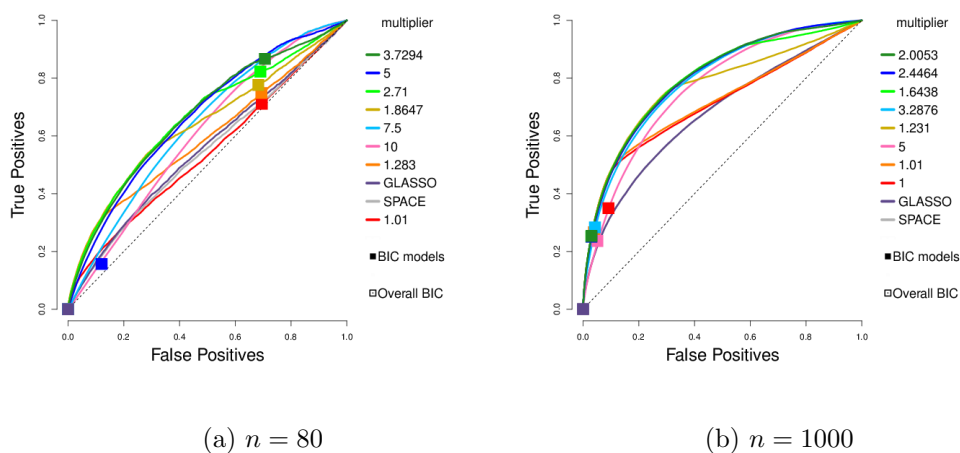
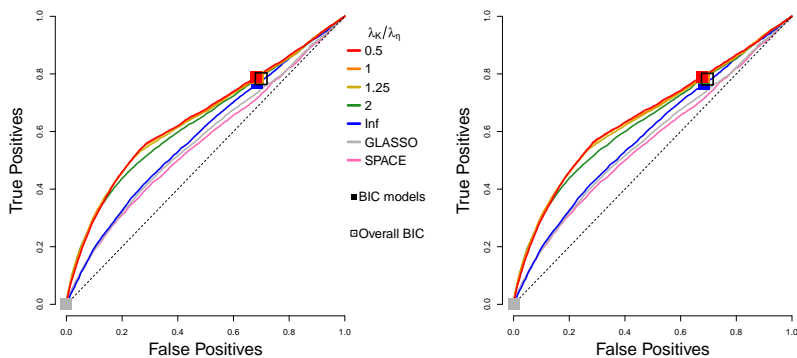
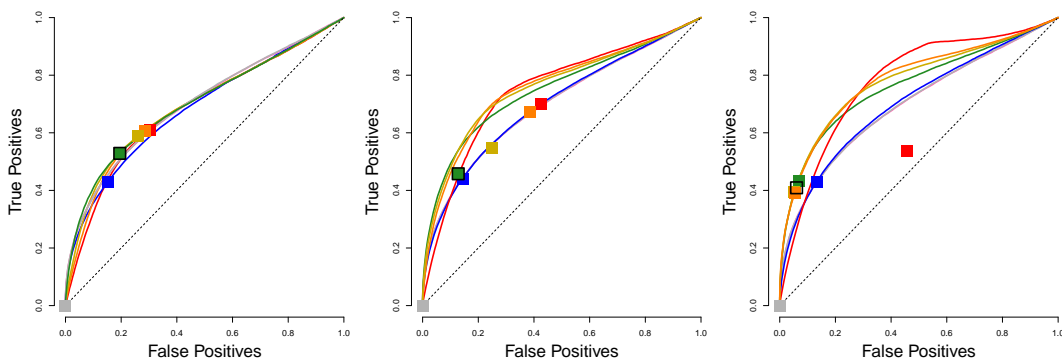


Figure A.3: Performance of  $\min(x, 3)$  for truncated centered GGMs with different multipliers, compared to GLASSO and SPACE, in the centered setting,  $n = 80$  or  $1000$ .



(a)  $n = 80$ , mult = 1.7897

(b)  $n = 80$ , mult = 1.8647



(c)  $n = 1000$ , mult = 1

(d)  $n = 1000$ , mult = 1.2310

(e)  $n = 1000$ , mult = 1.6438

Figure A.4: Performance of the non-centered estimator with  $h(x) = \min(x, 3)$ . Each curve corresponds to a different choice of  $\lambda_{\mathbf{K}}/\lambda_{\boldsymbol{\eta}}$ . Squares indicate models picked by eBIC with refit. The square with black outline has the highest eBIC among all models (combinations of  $\lambda_{\mathbf{K}}$ ,  $\lambda_{\boldsymbol{\eta}}$ ). The multipliers correspond to medium or high for  $n = 80$ , and low, medium and high for  $n = 1000$ , respectively.

A.3.2 Other a/b Models

Figure A.5 exhibits the results for the exponential models.

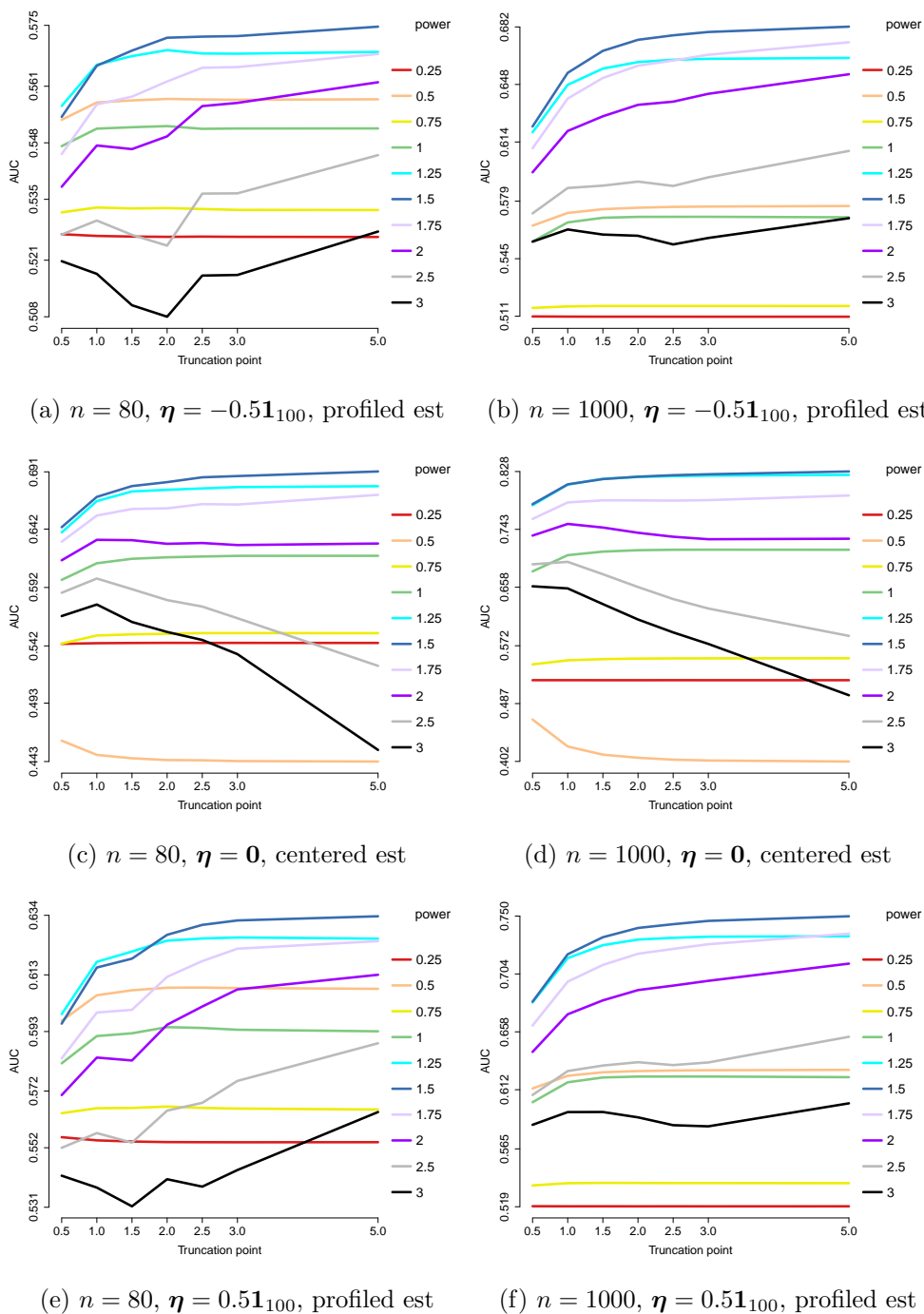
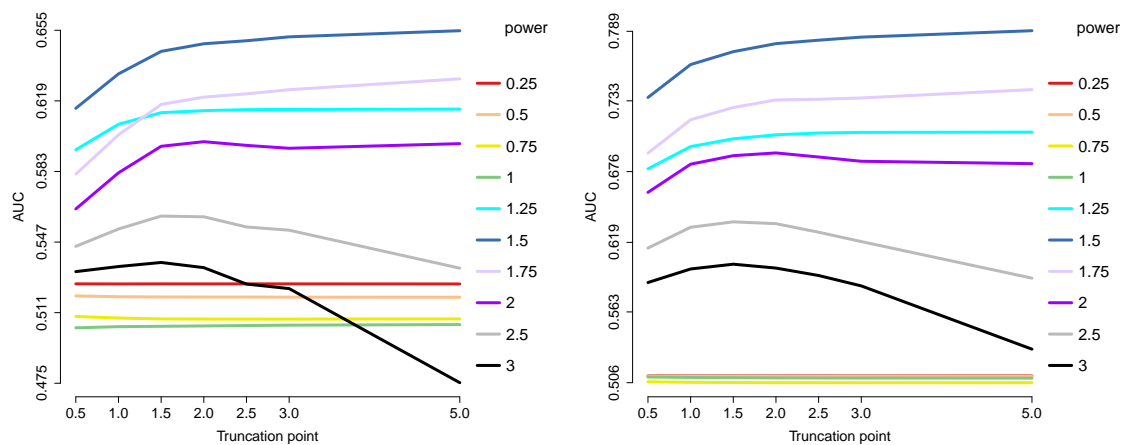


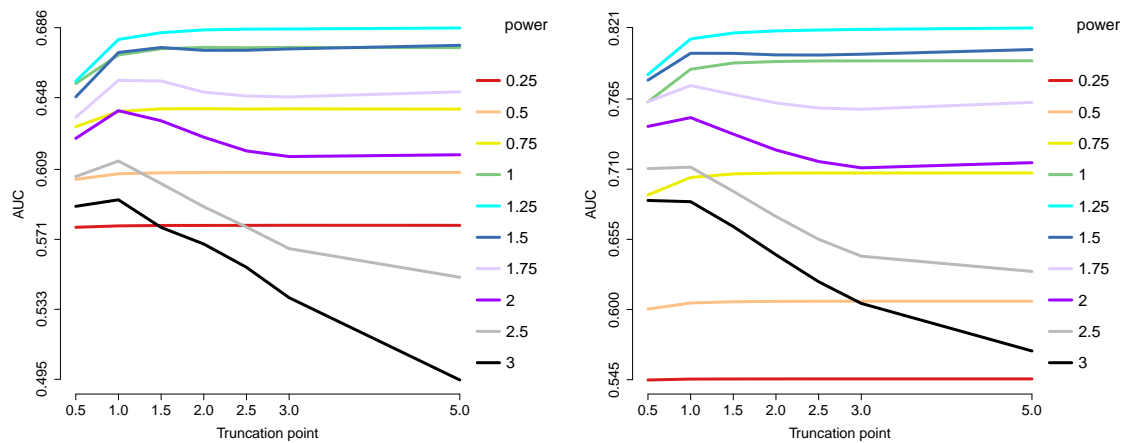
Figure A.5: AUCs for edge recovery using generalized score matching for exponential models.

Figure A.6 displays the results for the gamma models.



(a)  $n = 80$ ,  $\boldsymbol{\eta} = -0.51_{100}$ , profiled estimator

(b)  $n = 1000$ ,  $\boldsymbol{\eta} = -0.51_{100}$ , profiled estimator



(c)  $n = 80$ ,  $\boldsymbol{\eta} = 0.51_{100}$ , profiled estimator

(d)  $n = 1000$ ,  $\boldsymbol{\eta} = 0.51_{100}$ , profiled estimator

Figure A.6: AUCs for edge recovery using generalized score matching for the gamma models.

Figures A.7 and A.8 demonstrate the results for  $a = 3/2$  and  $b = 1/2$  or  $b = 0$ , respectively.

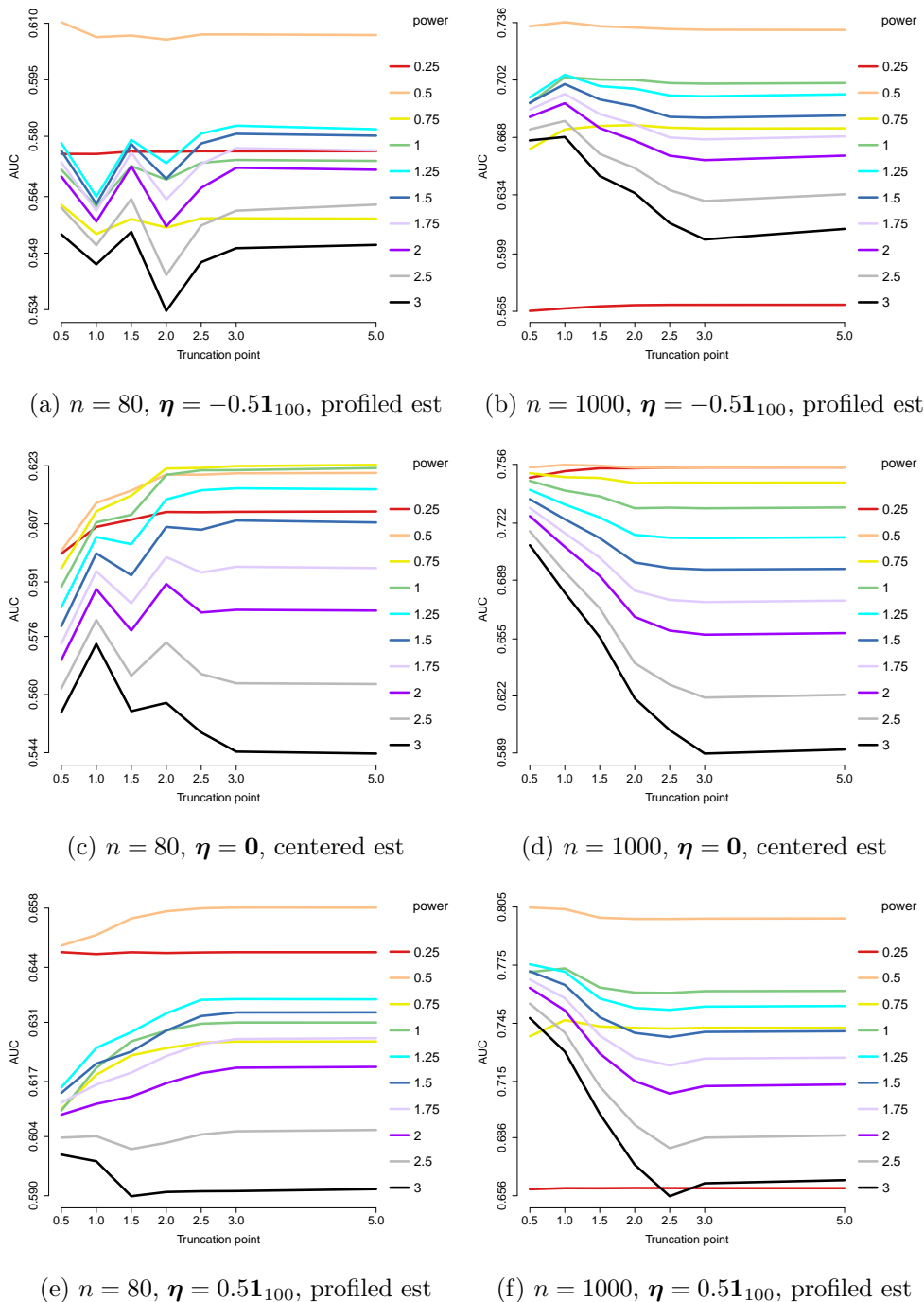
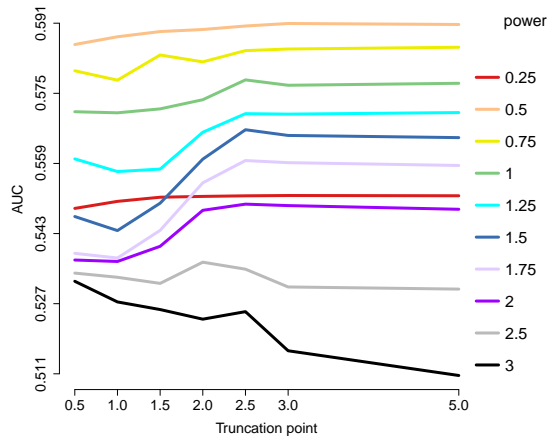
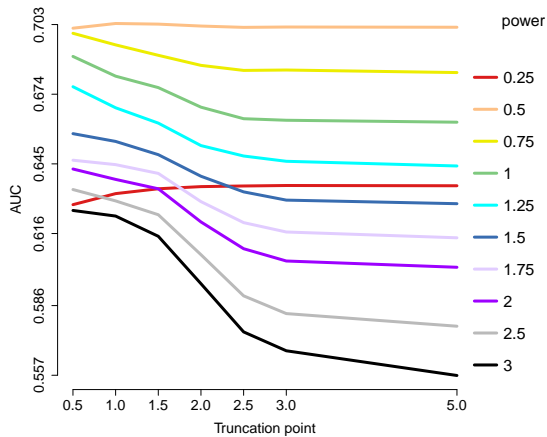


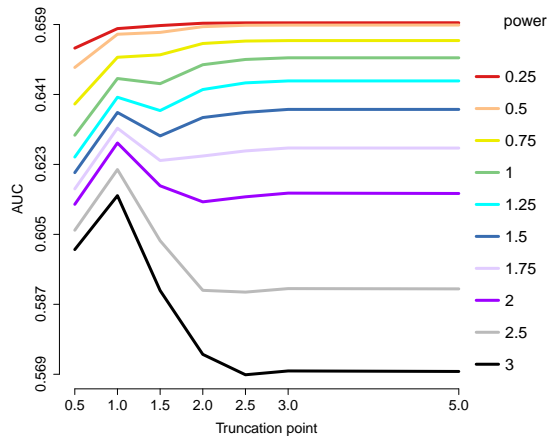
Figure A.7: AUCs for edge recovery using generalized score matching for  $a = 3/2, b = 1/2$ .



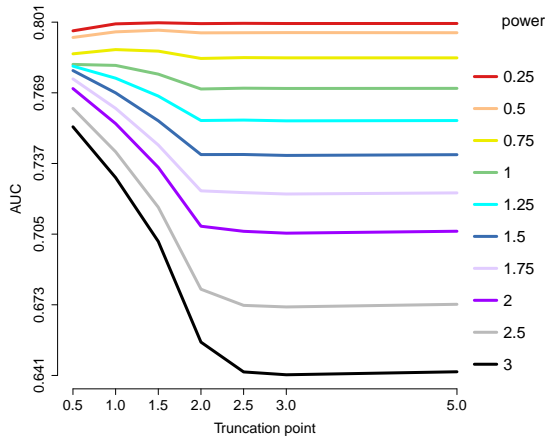
(a)  $n = 80, \eta = -0.51_{100}$ , profiled estimator



(b)  $n = 1000, \eta = -0.51_{100}$ , profiled estimator



(c)  $n = 80, \eta = 0.51_{100}$ , profiled estimator



(d)  $n = 1000, \eta = 0.51_{100}$ , profiled estimator

Figure A.8: AUCs for edge recovery using generalized score matching for  $a = 3/2, b = 0$ .

Appendix B  
**APPENDIX TO CHAPTER 3**

Before proving Lemma 13 we first prove the following lemma.

**Lemma 38.** *Suppose  $h_1, \dots, h_m$  are absolutely continuous in any bounded sub-interval of  $\mathbb{R}_+$ . Then for any  $j = 1, \dots, m$  and any  $\mathbf{x}_{-j} \in \mathcal{S}_{-j, \mathcal{D}}$ ,  $(h_j \circ \varphi_j)$  is absolutely continuous in  $x_j$  in any bounded sub-interval of  $\mathcal{C}_{j, \mathcal{D}}(\mathbf{x}_{-j})$ .*

*Proof of Lemma 38.* In the proof we drop the dependency on  $\mathbf{x}_{-j}$  in notation. By assumption, under Equation 3.5 any bounded sub-interval  $[a, b]$  of  $\mathcal{C}_{j, \mathcal{D}}(\mathbf{x}_{-j})$  must be a sub-interval of  $[a_{k,j}, b_{k,j}]$  for some  $k$  (for simplicity we do not differentiate among  $[a, b]$ ,  $(a, b]$ ,  $[a, b)$  and  $(a, b)$  here).

- (1) If  $a_{k,j} > -\infty$  and  $b_{k,j} < +\infty$ , denote  $C_0 \equiv \min\{C_j, (b_{k,j} - a_{k,j})/2\}$  and rewrite

$$\begin{aligned} & (h_j \circ \varphi_j)(\mathbf{x}) \\ &= h_j(\min(C_j, x_j - a_{k,j}, b_{k,j} - x_j)) \\ &= h_j(x_j - a_{k,j})\mathbb{1}_{x_j \in [a_{k,j}, a_{k,j} + C_0]} + h_j(C_j)\mathbb{1}_{x_j \in [a_{k,j} + C_0, b_{k,j} - C_0]} + h_j(b_{k,j} - x_j)\mathbb{1}_{x_j \in [b_{k,j} - C_0, b_{k,j}]}. \end{aligned}$$

Then by absolute continuity of  $h_j$  in  $[a_{k,j}, b_{k,j}]$  it is apparent that  $(h_j \circ \varphi_j)$  is differentiable in  $x_j$  a.e. with derivative

$$\partial_j(h_j \circ \varphi_j)(\mathbf{x}) = h'_j(x_j - a_{k,j})\mathbb{1}_{x_j \in [a_{k,j}, a_{k,j} + C_0]} - h'_j(b_{k,j} - x_j)\mathbb{1}_{x_j \in [b_{k,j} - C_0, b_{k,j}]}$$

Then by the absolute continuity of  $h_j$  again, for  $x_j \in [a_{k,j}, b_{k,j}]$ ,

$$\begin{aligned} & \int_{a_{k,j}}^{x_j} \partial_j(h_j \circ \varphi_j)(t_j; \mathbf{x}_{-j}) dt_j \\ &= h_j(x_j - a_{k,j})\mathbb{1}_{x_j \in [a_{k,j}, a_{k,j} + C_0]} + h_j(C_j)\mathbb{1}_{x_j \in [a_{k,j} + C_0, b_{k,j} - C_0]} + h_j(b_{k,j} - x_j)\mathbb{1}_{x_j \in [b_{k,j} - C_0, b_{k,j}]} \\ &= (h_j \circ \varphi_j)(\mathbf{x}), \end{aligned}$$

which proves that  $(h_j \circ \varphi_j)(\mathbf{x})$  is absolutely continuous in  $x_j$  in  $[a_{k,j}, b_{k,j}]$ , and hence in  $[a, b] \subset [a_{k,j}, b_{k,j}]$ .

- (2) If  $a_{k,j} > -\infty$  and  $b_{k,j} = +\infty$ , on  $[a, b]$   $(h_j \circ \varphi_j)(\mathbf{x}) = h_j(\min(C_j, x_j - a_{k,j}))$  is an absolutely continuous function in a linear function of  $x_j$  truncated above by  $C_j$ , and is thus trivially absolutely continuous in  $[a, b]$ .

(3) If  $a_{k,j} = -\infty$  and  $b_{k,j} < +\infty$ , on  $[a, b]$   $(h_j \circ \varphi_j)(\mathbf{x}) = h_j(\min(C_j, b_{k,j} - x_j))$  is an absolutely continuous function in a linear function of  $x_j$  truncated above by  $C_j$ , and is thus trivially absolutely continuous in  $[a, b]$ .

(4) If  $a_{k,j} = -\infty$  and  $b_{k,j} = +\infty$ ,  $(h_j \circ \varphi_j)(\mathbf{x}) = h_j(C_j)$  is constant and hence trivially absolutely continuous in  $[a, b]$ .

□

*Proof of Lemma 13.* By simple manipulation

$$\begin{aligned}
J_{\mathbf{h}, \mathcal{C}, \mathcal{D}}(p) &\equiv \frac{1}{2} \int_{\mathcal{D}} p_0(\mathbf{x}) \left\| \nabla \log p(\mathbf{x}) \odot (\mathbf{h} \circ \boldsymbol{\varphi})^{1/2}(\mathbf{x}) - \nabla \log p_0(\mathbf{x}) \odot (\mathbf{h} \circ \boldsymbol{\varphi})^{1/2}(\mathbf{x}) \right\|_2^2 d\mathbf{x} \\
&= \frac{1}{2} \sum_{j=1}^m \int_{\mathcal{D}} p_0(\mathbf{x}) (h_j \circ \varphi_j)(\mathbf{x}) (\partial_j \log p_0(\mathbf{x}) - \partial_j \log p(\mathbf{x}))^2 d\mathbf{x} \\
&= \frac{1}{2} \sum_{j=1}^m \int_{\mathcal{D}} p_0(\mathbf{x}) (h_j \circ \varphi_j)(\mathbf{x}) (\partial_j \log p(\mathbf{x}))^2 d\mathbf{x} \\
&\quad - \sum_{j=1}^m \int_{\mathcal{D}} p_0(\mathbf{x}) (h_j \circ \varphi_j)(\mathbf{x}) \partial_j \log p_0(\mathbf{x}) \partial_j \log p(\mathbf{x}) d\mathbf{x} + \text{const.} \tag{B.1}
\end{aligned}$$

By (B.1) it suffices to prove for all  $j = 1, \dots, m$  that

$$\int_{\mathcal{D}} p_0(\mathbf{x}) (h_j \circ \varphi_j)(\mathbf{x}) \partial_j \log p_0(\mathbf{x}) \partial_j \log p(\mathbf{x}) d\mathbf{x} = - \int_{\mathcal{D}} p_0(\mathbf{x}) \partial_j [(h_j \circ \varphi_j)(\mathbf{x}) \partial_j \log p(\mathbf{x})] d\mathbf{x}. \tag{B.2}$$

$\int_{\mathcal{D}} p_0(\mathbf{x}) \left\| \nabla \log p(\mathbf{x}) \odot (\mathbf{h} \circ \boldsymbol{\varphi})^{1/2}(\mathbf{x}) \right\|_2^2 d\mathbf{x}$  and  $\int_{\mathcal{D}} p_0(\mathbf{x}) \left\| \nabla \log p_0(\mathbf{x}) \odot (\mathbf{h} \circ \boldsymbol{\varphi})^{1/2}(\mathbf{x}) \right\|_2^2 d\mathbf{x}$  are both finite by assumption, so by  $|2ab| \leq a^2 + b^2$  the integrand in the left-hand side of (B.2) is integrable. Then by Fubini-Tonelli

$$\begin{aligned}
&\int_{\mathcal{D}} p_0(\mathbf{x}) (h_j \circ \varphi_j)(\mathbf{x}) \partial_j \log p_0(\mathbf{x}) \partial_j \log p(\mathbf{x}) d\mathbf{x} \\
&= \int_{\mathcal{S}_{-j}} \int_{\mathcal{C}_j(\mathbf{x}_{-j})} \underbrace{(h_j \circ \varphi_j)(\mathbf{x}) \partial_j p_0(\mathbf{x}) \partial_j \log p(\mathbf{x})}_{\equiv f(\mathbf{x})} dx_j d\mathbf{x}_{-j} \\
&= \int_{\mathcal{S}_{-j}} \int_{\mathbb{R}} \mathbb{1}_{\mathcal{C}_j(\mathbf{x}_{-j})}(x_j) f(\mathbf{x}) dx_j d\mathbf{x}_{-j}
\end{aligned}$$

$$\begin{aligned}
&= \int_{\mathcal{S}_{-j}} \int_{\mathbb{R}} \left[ \sum_{k=1}^{K_j(\mathbf{x}_{-j})} \mathbb{1}_{[a_{k,j}(\mathbf{x}_{-j}), b_{k,j}(\mathbf{x}_{-j})]}(x_j) \right] f(x_j; \mathbf{x}_{-j}) dx_j d\mathbf{x}_{-j} \\
&= \int_{\mathcal{S}_{-j}} \left[ \sum_{k=1}^{K_j(\mathbf{x}_{-j})} \int_{a_{k,j}(\mathbf{x}_{-j})}^{b_{k,j}(\mathbf{x}_{-j})} f(x_j; \mathbf{x}_{-j}) dx_j \right] d\mathbf{x}_{-j} \tag{B.3}
\end{aligned}$$

where the interchangeability of integration and (potentially infinite) summation is justified by Fubini-Tonelli again. Then using the decomposition of the domain in (3.5) while omitting the dependency of  $a_{k,j}$  and  $b_{k,j}$  on  $\mathbf{x}_{-j}$  in notation, for a.e.  $\mathbf{x}_{-j} \in \mathcal{S}_{-j}$  and any  $k = 1, \dots, K_j(\mathbf{x}_{-j})$  we have

$$\begin{aligned}
\int_{a_{k,j}}^{b_{k,j}} f(\mathbf{x}) dx_j &= \int_{a_{k,j}}^{b_{k,j}} (h_j \circ \varphi_j)(\mathbf{x}) \partial_j p_0(\mathbf{x}) \partial_j \log p(\mathbf{x}) dx_j \\
&= \lim_{x_j \nearrow b_{k,j}} (h_j \circ \varphi_j)(\mathbf{x}) p_0(\mathbf{x}) \partial_j \log p(\mathbf{x}) - \lim_{x_j \searrow a_{k,j}^+} (h_j \circ \varphi_j)(\mathbf{x}) p_0(\mathbf{x}) \partial_j \log p(\mathbf{x}) \\
&\quad - \int_{a_{k,j}}^{b_{k,j}} p_0(\mathbf{x}) \partial_j [(h_j \circ \varphi_j)(\mathbf{x}) \partial_j \log p(\mathbf{x})] dx_j \\
&= - \int_{a_{k,j}}^{b_{k,j}} p_0(\mathbf{x}) \partial_j [(h_j \circ \varphi_j)(\mathbf{x}) \partial_j \log p(\mathbf{x})] dx_j,
\end{aligned}$$

by integration by parts and by Assumption (A1) on the limits going to 0. The integration by parts is justified by the fundamental theorem of calculus for absolutely continuous functions (Lemma 38) as well as the product rule (cf. proof of Lemma 34). Thus, by going backwards using Fubini-Tonelli twice again, (B.3) becomes

$$\begin{aligned}
&\int_{\mathcal{S}_{-j}} \left\{ - \sum_{k=1}^{K_j(\mathbf{x}_{-j})} \int_{a_{k,j}(\mathbf{x}_{-j})}^{b_{k,j}(\mathbf{x}_{-j})} p_0(\mathbf{x}) \partial_j [(h_j \circ \varphi_j)(\mathbf{x}) \partial_j \log p(\mathbf{x})] dx_j \right\} d\mathbf{x}_{-j} \\
&= - \int_{\mathcal{S}_{-j}} \int_{\mathcal{C}_j(\mathbf{x}_{-j})} p_0(\mathbf{x}) \partial_j [(h_j \circ \varphi_j)(\mathbf{x}) \partial_j \log p(\mathbf{x})] dx_j d\mathbf{x}_{-j} \\
&= - \int_{\mathcal{D}} p_0(\mathbf{x}) \partial_j [(h_j \circ \varphi_j)(\mathbf{x}) \partial_j \log p(\mathbf{x})] d\mathbf{x},
\end{aligned}$$

proving (B.2).  $\square$

*Proof of Theorem 16.* Note that the condition that  $\mathbf{v}^{\top} \mathbf{K} \mathbf{v}^a > 0 \forall \mathbf{v} \in \mathcal{D} \setminus \{\mathbf{0}\}$  implies that  $\mathbf{v}^{\top} \mathbf{K} \mathbf{v}^a > 0 \forall \mathbf{v} \in \mathcal{D}_+ \equiv \{\mathbf{v} / \|\mathbf{v}\|_2 : \mathbf{v} \in \mathcal{D} \setminus \{\mathbf{0}\}\} \subseteq \{\mathbf{v} \in \mathbb{R}^m : \|\mathbf{v}\|_2 = 1\} \equiv \mathbb{S}^{m-1}$  with  $\mathbb{S}^{m-1}$

compact, so

$$N_{\mathbf{K}} \equiv \inf_{\mathbf{v} \in \mathcal{D} \setminus \{\mathbf{0}\}} \mathbf{v}^{a\top} \mathbf{K} \mathbf{v}^a / \mathbf{v}^{a\top} \mathbf{v}^a = \inf_{\mathbf{v} \in \mathcal{D}_+} \mathbf{v}^{a\top} \mathbf{K} \mathbf{v}^a / \mathbf{v}^{a\top} \mathbf{v}^a \geq \inf_{\mathbf{v} \in \mathbb{S}^{m-1}} \mathbf{v}^{a\top} \mathbf{K} \mathbf{v}^a / \mathbf{v}^{a\top} \mathbf{v}^a > 0.$$

(1) *Case  $a > 0$  and  $b > 0$  (CC1, CC2):* Since  $p$  is bounded everywhere, it is integrable over a bounded  $\mathcal{D}$  (proving (CC1)). Otherwise, assume  $\mathcal{D}$  is unbounded. If either  $a$  or  $b$  is non-integer, then  $\mathcal{D} \subset \mathbb{R}_+^m$  and a sufficient condition is  $\mathbf{v}^a \mathbf{K} \mathbf{v}^a > 0 \forall \mathbf{v} \in \mathcal{D} \setminus \{\mathbf{0}\}$  and either  $\boldsymbol{\eta}^\top \mathbf{v}^b \leq 0 \forall \mathbf{v} \in \mathcal{D}$  or  $2a > b > 0$ , corresponding to (i) and (ii) in the Proof of Theorem 5 in Section A.1.3, respectively. If  $a$  and  $b$  are both integers,  $\mathcal{D}$  may be any subset of  $\mathbb{R}^m$  and the same sufficient condition can be implied following the same proof in Section A.1.3, with integration over  $(-\infty, +\infty)$  instead of  $(0, +\infty)$ . This proves (CC2).

(2) *Case  $a > 0$  and  $b = 0$  (CC3):* By definition  $\mathcal{D} \subseteq \mathbb{R}_+^m$ . If  $\mathcal{D}$  is bounded,  $-\frac{1}{2a} \mathbf{x}^{a\top} \mathbf{K} \mathbf{x}^a$  as a continuous function is bounded, and it thus suffices to bound  $\int_{\mathcal{D}} \exp(\boldsymbol{\eta}^\top \log(\mathbf{x})) \, d\mathbf{x} = \int_{\mathcal{D}} \prod_{j=1}^m x_j^{\eta_j} \, d\mathbf{x} \leq \prod_{j=1}^m \int_{\rho_j(\mathcal{D})} x_j^{\eta_j} \, dx_j < +\infty$  if  $\eta_j > -1$  for all  $j$  such that  $0 \in \rho_j(\mathcal{D})$ , where for the  $\leq$  step we used the fact that  $x_j > 0$ . This proves (CC3) (i).

If  $\mathcal{D}$  is unbounded and  $\mathbf{v}^{a\top} \mathbf{K} \mathbf{v}^a > 0$  for all  $\mathbf{v} \in \mathcal{D} \setminus \{\mathbf{0}\}$ , using the fact that  $\exp(\dots) > 0$ ,

$$\begin{aligned} \int_{\mathcal{D}} p_{\boldsymbol{\eta}, \mathbf{K}}(\mathbf{x}) \, d\mathbf{x} &= \int_{\mathcal{D}} \exp(-\mathbf{x}^{a\top} \mathbf{K} \mathbf{x}^a / (2a) + \boldsymbol{\eta}^\top \log(\mathbf{x})) \, d\mathbf{x} \\ &\leq \prod_{j=1}^m \int_{\rho_j(\mathcal{D})} \exp(-N_{\mathbf{K}} x_j^{2a} / (2a) + \eta_j \log(x_j)) \, dx_j. \end{aligned}$$

Note that the indefinite integral of the last display is

$$-\frac{1}{2a} x^{1+\eta_j} \left( \frac{N_{\mathbf{K}}}{2a} x^{2a} \right)^{-(1+\eta_j)/(2a)} \Gamma \left[ \frac{1+\eta_j}{2a}, \frac{N_{\mathbf{K}} x^{2a}}{2a} \right]$$

so the definite integral is finite if and only if  $\eta_j > -1$  for all  $j$  s.t.  $0 \in \rho_j(\mathcal{D})$ . This proves (CC3) (ii).

If  $\mathcal{D}$  is unbounded and  $\mathbf{v}^{a\top} \mathbf{K} \mathbf{v}^a \geq 0$  for all  $\mathbf{v} \in \mathcal{D}$ ,  $\int_{\mathcal{D}} p_{\boldsymbol{\eta}, \mathbf{K}}(\mathbf{x}) \, d\mathbf{x} \leq \prod_{j=1}^m \int_{\rho_j(\mathcal{D})} x_j^{\eta_j} \, dx_j < \infty$  if  $\eta_j > -1$  for all  $j$  s.t.  $0 \in \rho_j(\mathcal{D})$  and  $\eta_j < -1$  for all  $j$  s.t.  $\rho_j(\mathcal{D})$  is unbounded. This proves (CC3) (iii).

(3) *Case  $a = 0$ ,  $\mathcal{D}$  is bounded and  $0 \notin \rho_j(\mathcal{D})$  for all  $j$  (CC4):* If  $\mathcal{D}$  is bounded and  $0 \notin \rho_j(\mathcal{D})$  for all  $j$ , then  $\log(\mathcal{D})$  is bounded, and since the integrand is continuous and bounded, the integral is finite without any further requirements.

(4) *Case  $a = 0$  and  $b = 0$  (CC5):* Assume  $\log(\mathbf{x})^\top \mathbf{K} \log(\mathbf{x}) > 0$  for all  $\mathbf{x} \in \mathcal{D}$ , then

$$\begin{aligned} \int_{\mathcal{D}} p_{\boldsymbol{\eta}, \mathbf{K}}(\mathbf{x}) \, d\mathbf{x} &= \int_{\mathcal{D}} \exp\left(-\frac{1}{2} \log(\mathbf{x})^\top \mathbf{K} \log(\mathbf{x}) + \boldsymbol{\eta}^\top \log(\mathbf{x})\right) \, d\mathbf{x} \\ &= \int_{\log(\mathcal{D})} \exp\left(-\frac{1}{2} \mathbf{x}^\top \mathbf{K} \mathbf{x} + (\boldsymbol{\eta} + \mathbf{1}_m)^\top \mathbf{x}\right) \, d\mathbf{x} \\ &< \prod_{j=1}^m \int_{\log(\rho_j(\mathcal{D}))} \exp\left(-N_{\mathbf{K}} x_j^2 / 2 + (\eta_j + 1)x_j\right) \, dx_j \\ &< \prod_{j=1}^m \int_{-\infty}^{\infty} \exp\left(-N_{\mathbf{K}} x_j^2 / 2 + (\eta_j + 1)x_j\right) \, dx_j < +\infty \end{aligned}$$

since the integrand is proportional to a univariate Gaussian density.

(5) *Case  $a = 0$  and  $b > 0$  (CC6, CC7):* Assume  $\log(\mathbf{x})^\top \mathbf{K} \log(\mathbf{x}) > 0$  for all  $\mathbf{x} \in \mathcal{D}$  and  $\eta_j \leq 0$  for all  $j$  s.t.  $\rho_j(\mathcal{D})$  is unbounded (from above). Then

$$\begin{aligned} \int_{\mathcal{D}} p_{\boldsymbol{\eta}, \mathbf{K}}(\mathbf{x}) \, d\mathbf{x} &= \int_{\mathcal{D}} \exp\left(-\frac{1}{2} \log(\mathbf{x})^\top \mathbf{K} \log(\mathbf{x}) + \boldsymbol{\eta}^\top \mathbf{x}^b\right) \, d\mathbf{x} \\ &= \int_{\log(\mathcal{D})} \exp\left(-\frac{1}{2} \mathbf{x}^\top \mathbf{K} \mathbf{x} + \mathbf{1}_m^\top \mathbf{x} + \boldsymbol{\eta}^\top \exp(b\mathbf{x})\right) \, d\mathbf{x} \\ &< \prod_{j=1}^m \int_{\log(\rho_j(\mathcal{D}))} \exp\left(-N_{\mathbf{K}} x_j^2 / 2 + x_j + \eta_j \exp(bx_j)\right) \, dx_j \\ &\leq \prod_{j=1}^m \int_{-\infty}^{\infty} c_j \exp\left(-N_{\mathbf{K}} x_j^2 / 2 + x_j\right) \, dx_j < +\infty, \end{aligned}$$

where  $c_j \equiv 1$  if  $\eta_j \leq 0$  or  $c_j \equiv \exp\left(\eta_j (\sup \rho_j(\mathcal{D}))^b\right) > +\infty$  otherwise. This proves (CC6).

Finally, if  $\log(\mathcal{D})$  is unbounded and  $\log(\mathbf{x})^\top \mathbf{K} \log(\mathbf{x}) \geq 0$  for all  $\mathbf{x} \in \mathcal{D}$ , the integral is bounded by

$$\prod_{j=1}^m \int_{\log(\rho_j(\mathcal{D}))} \exp\left(x_j + \eta_j \exp(bx_j)\right) \, dx_j$$

which is finite if and only if  $\eta_j < 0$  for all  $j$  s.t.  $\rho_j(\mathcal{D})$  is unbounded (from above). This proves (CC7).

□

*Proof of Theorem 18.* It suffices to consider the case  $\mathcal{D} = \mathbb{R}_+^m$  for general  $a$  and  $b$  as well as  $\mathcal{D} = \mathbb{R}^m$  for integer  $a > 0$  and  $b > 0$  (so that (3.17) is well defined on  $\mathbb{R}^m$ ): For (A.1), the irregularities only occur at the boundary points, but with the composition  $(h_j \circ \varphi_j)(\mathbf{x})$  with  $x_j$  approaching any finite boundary point behaves like  $h_j(x_j)$  with  $x_j \searrow 0^+$  in  $\mathcal{D} = \mathbb{R}_+^m$ , and  $(h_j \circ \varphi_j)(\mathbf{x})$  with  $x_j \rightarrow \infty$  behaves like  $h_j(x_j)$  with  $x_j \rightarrow \infty$  in  $\mathcal{D} = \mathbb{R}_+^m$  (or  $\mathbb{R}^m$  if applicable). For (A.2), obviously integrability over  $\mathcal{D}$  follows from that over  $\mathcal{D} = \mathbb{R}_+^m$  or  $\mathbb{R}^m$ . (A.3) is trivially satisfied by a power function  $h_j$ .

As in the proof of Theorem 16,  $N_{\mathbf{K}} \equiv \inf_{\mathbf{v} \in \mathcal{D}} \mathbf{v}^a \top \mathbf{K} \mathbf{v}^a / \mathbf{v}^a \top \mathbf{v}^a > 0$ .

(1) The case for  $a > 0$  and  $b \geq 0$  and  $\mathcal{D} = \mathbb{R}_+^m$  is covered in Chapter 2. The proof for the case for  $a > 0$  and  $b > 0$  and  $\mathcal{D} = \mathbb{R}^m$  is analogous and omitted.

(2) *Case  $a = 0$  and  $b = 0$ :*

$$\begin{aligned}
& |p_0(\mathbf{x}) \partial_j \log p(\mathbf{x})| \\
& \propto \exp \left( -\frac{1}{2} \log(\mathbf{x}) \top \mathbf{K}_0 \log(\mathbf{x}) + \boldsymbol{\eta}_0 \top \log(\mathbf{x}) \right) |x_j^{-1} (\eta_j - \boldsymbol{\kappa}_{j,-j} \top \log(\mathbf{x}_{-j})) - \kappa_{jj} x_j^{-1} \log x_j| \\
& \leq |(\eta_j - \boldsymbol{\kappa}_{j,-j} \top \log \mathbf{x}_{-j}) \exp [-N_{\mathbf{K}_0} (\log x_j)^2 / 2 + (\eta_j - 1) \log x_j]| \\
& - \kappa_{jj} \exp [-N_{\mathbf{K}_0} (\log x_j)^2 / 2 + (\eta_j - 1) \log x_j] \log x_j| \times \prod_{k \neq m} \exp (-N_{\mathbf{K}_0} (\log x_k)^2 / 2 + \eta_j \log x_k) \\
& \propto \mathcal{O} [\exp (-N_{\mathbf{K}_0} y_j^2 / 2 + (\eta_j - 1) y_j)] + \mathcal{O} [\exp (-N_{\mathbf{K}_0} y_j^2 / 2 + (\eta_j - 1) y_j) y_j]
\end{aligned}$$

which apparently vanishes as  $x_j \searrow 0^+$  and  $x_j \nearrow +\infty$  with  $y_j \equiv \log(x_j)$  since it is dominated by a constant times a Gaussian density in  $y_j$ . Thus, by Proposition 14, (A.1) is satisfied with any  $\alpha_j \geq 0$ . Likewise, for (A.2),

$$\begin{aligned}
& \int_{\mathbb{R}_+^m} p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) \odot (\mathbf{h} \circ \boldsymbol{\varphi})^{1/2}(\mathbf{x})\|_2^2 d\mathbf{x} \\
& \leq \text{const} \cdot \sum_{j=1}^m \int_{\mathbb{R}_+^m} \prod_{k=1}^m \exp [-N_{\mathbf{K}_0} (\log x_k)^2 / 2 + \eta_k \log(x_k)] \times \\
& \quad h_j(x_j) [x_j^{-1} (\eta_j - \boldsymbol{\kappa}_{j,-j} \top \log(\mathbf{x}_{-j})) - \kappa_{jj} x_j^{-1} \log x_j]^2 d\mathbf{x},
\end{aligned}$$

which can be decomposed into a sum of products of univariate integrals of the form

$$\text{const} \cdot \exp(-N_{\mathbf{K}_0}(\log x_j)^2/2 + A \log(x_j)) (\log x_j)^B (h_j(x_j))^C$$

with  $B = 0, 1, 2$ ,  $C = 0, 1$ , and constants  $A$ . With  $h_j(x_j) = x_j^{\alpha_j}$  for any  $\alpha_j \geq 0$  this is bounded by some Gaussian density in  $\log x_j$ , and thus  $\int_{\mathbb{R}_+^m} p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) \odot (\mathbf{h} \circ \boldsymbol{\varphi})^{1/2}(\mathbf{x})\|_2^2 d\mathbf{x} < +\infty$ . Proof for  $\int_{\mathbb{R}_+^m} p_0(\mathbf{x}) \|\nabla \log p(\mathbf{x}) \odot (\mathbf{h} \circ \boldsymbol{\varphi})(\mathbf{x})\|_1 d\mathbf{x} < +\infty$  is similar and is omitted.

(3) *Case  $a = 0$  and  $b > 0$* : Recall  $\rho_j(\mathcal{D}) \equiv \overline{\{x_j : \mathbf{x} \in \mathcal{D}\}}$ . Let  $\rho_j^*(\mathcal{D}) \equiv \sup \rho_j(\mathcal{D})$ . Since we assume that  $\eta_j \leq 0$  for any  $j$  such that  $\rho_j^*(\mathcal{D}) < +\infty$ ,

$$\begin{aligned} & p_0(\mathbf{x}) \partial_j \log p(\mathbf{x}) \\ & \propto \exp\left(-\frac{1}{2} \log(\mathbf{x})^\top \mathbf{K}_0 \log(\mathbf{x}) + \frac{1}{b} \boldsymbol{\eta}_0^\top \mathbf{x}^b\right) [\eta_j x_j^{b-1} - x_j^{-1} \boldsymbol{\kappa}_{j,-j}^\top \log(\mathbf{x}_{-j}) - \kappa_{jj} x_j^{-1} \log x_j] \\ & \leq \exp\left(-\frac{1}{2} \log(\mathbf{x})^\top \mathbf{K}_0 \log(\mathbf{x}) + \frac{1}{b} \sum_{j: \rho_j^*(\mathcal{D}) < +\infty} \eta_{0j} (\rho_j^*(\mathcal{D}))^b\right) \times \\ & \quad [-x_j^{-1} \boldsymbol{\kappa}_{j,-j}^\top \log(\mathbf{x}_{-j}) - \kappa_{jj} x_j^{-1} \log x_j] \\ & \propto \exp\left(-\frac{1}{2} \log(\mathbf{x})^\top \mathbf{K}_0 \log(\mathbf{x})\right) [-x_j^{-1} \boldsymbol{\kappa}_{j,-j}^\top \log(\mathbf{x}_{-j}) - \kappa_{jj} x_j^{-1} \log x_j] \end{aligned}$$

is bounded by the corresponding quantity in the  $a = b = 0$  case with  $\boldsymbol{\eta} = \mathbf{0}_m$ , and (A.1) is thus satisfied. Similarly, the two quantities for (A.2) are bounded by a constant times those in the  $a = b = 0$  case with  $\boldsymbol{\eta} = \mathbf{0}_m$  and (A.2) is thus also satisfied.  $\square$

*Proof of Theorem 19.* Fix  $j = 1, \dots, m-1$  and  $\mathbf{x}_{-j,-m} \in \mathcal{S}_{-j, \mathcal{D}_{-m}}$ , i.e.  $\mathbf{x}_{-j,-m} \in \mathbb{R}_+^{m-2}$  such that  $\mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m} < 1$ . In our discussion, for ease of notation, given  $x_j$  and  $\mathbf{x}_{-j,-m}$ , we may still write  $x_m \equiv 1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m}$ , a function in  $x_j$  and  $\mathbf{x}_{-j,-m}$ , and for simplicity we may drop its dependence on  $x_j$  and  $\mathbf{x}_{-j,-m}$ . Note that  $\mathcal{C}_j(\mathbf{x}_{-j,-m}) = (0, 1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m})$ .

(1) *Case  $a > 0$  and  $b \geq 0$* : For (A.1), we need  $p_0(\mathbf{x}_{-m}) \partial_j \log p(\mathbf{x}_{-m})(h_j \circ \varphi_j)(\mathbf{x}_{-m}) \rightarrow 0$  as  $x_j \searrow 0^+$  and  $x_j \nearrow 1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m}$ . As  $x_j$  goes to any finite constant, by (3.25)  $p_0(\mathbf{x}_{-m})$  converges to a non-zero constant when  $b > 0$ , or a finite constant times the limit of  $x_j^{\eta_j} x_m(x_j)^{\eta_m}$  when  $b = 0$ . Note that

$$\begin{aligned} \partial_j \log p(\mathbf{x}_{-m}) = & -(\boldsymbol{\kappa}_{-m,j}^\top \mathbf{x}_{-m}^a) x_j^{a-1} + (\boldsymbol{\kappa}_{-m,m}^\top \mathbf{x}_{-m}^a) x_m^{a-1}(x_j) - x_m^a(x_j) \kappa_{jm} x_j^{a-1} \\ & + x_m^{2a-1}(x_j) \kappa_{mm} + \eta_j x_j^{b-1} - \eta_m x_m^{b-1}(x_j). \end{aligned} \quad (\text{B.4})$$

i) If  $b > 0$ , by arguments above we only consider  $(h_j \circ \varphi_j)(\mathbf{x}_{-m}) \partial_j \log p(\mathbf{x}_{-m})$ .

a) As  $x_j \searrow 0^+$ ,  $x_m \nearrow 1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m} > 0$ , so  $\partial_j \log p(\mathbf{x}_{-m}) = \mathcal{O}(x_j^{a-1}) + \mathcal{O}(x_j^{b-1}) + \mathcal{O}(1)$ . Thus we need  $\alpha_j > \max\{0, 1-a, 1-b\}$  so that  $(h_j \circ \varphi_j)(\mathbf{x}_{-m}) \partial_j \log p(\mathbf{x}_{-m}) \rightarrow 0$ .

b) The case where  $x_j \nearrow 1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m}$  and  $x_m \searrow 0^+$  is an analog of a) by noting that  $\varphi_j$  is symmetric in  $x_j$  about the midpoint of its domain  $(1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m})/2$ .

ii) If  $b = 0$ , we need  $x_j^{\eta_j} x_m^{\eta_m}(x_j)(h_j \circ \varphi_j)(\mathbf{x}_{-m}) \partial_j \log p(\mathbf{x}_{-m}) \rightarrow 0$ . Note that this quantity has the same form as in a) just with  $\eta_j$  or  $\eta_m$  added to the  $a$  and  $b (= 0)$  in the exponents, we thus require  $\alpha_j > \max\{0, 1-a-\eta_j, 1-a-\eta_m, 1-\eta_j, 1-\eta_m\} = \max\{0, 1-\eta_j, 1-\eta_m\}$ .

In conclusion, (A.1) requires  $\alpha_j \geq \max\{0, 1-a, 1-b\}$  for  $b > 0$  or  $\alpha_j > \max\{0, 1-\eta_j\}$ . For (A.2), we only prove the first integrability condition, since the second integrability condition is similar. For the first, we need to show that

$$\int_{\mathbf{x}_{-m} > \mathbf{0}, \mathbf{1}_{m-1}^\top \mathbf{x}_{-m} < 1} p_0(\mathbf{x}_{-m})(h_j \circ \varphi_j)(\mathbf{x}_{-m}) (\partial_j \log p(\mathbf{x}_{-m}))^2 d\mathbf{x}_{-m} < +\infty.$$

Using the fact that  $\mathbf{0} \prec \mathbf{x}_{-m} \prec \mathbf{1}$  and  $0 < x_j^a < 1$ ,  $0 < x_m^a < 1$  with the triangle inequality multiple times, we have

$$\begin{aligned} |\partial_j \log p(\mathbf{x}_{-m})| & \leq \sum_{i=1}^{m-1} |\kappa_{ij}| x_j^{a-1} + \sum_{i=1}^{m-1} |\kappa_{im}| x_m^{a-1} + |\kappa_{jm}| x_j^{a-1} + x_m^{a-1} |\kappa_{mm}| + |\eta_j| x_j^{b-1} + |\eta_m| x_m^{b-1} \\ & \leq \|\mathbf{K}\|_1 x_j^{a-1} + \|\mathbf{K}\|_1 x_m^{a-1} + |\eta_j| x_j^{b-1} + |\eta_m| x_m^{b-1}, \end{aligned}$$

where  $\|\mathbf{K}\|_1 \equiv \max_{j=1,\dots,m} \sum_{i=1}^m |\kappa_{ij}|$ . We again consider the following two cases.

i) If  $b > 0$ ,  $p_0(\mathbf{x}_{-m})$  is bounded by an absolute constant, which we therefore ignore. We first fix  $\mathbf{x}_{-j,-m}$  and denote  $y_j(\mathbf{x}_{-j,-m}) \equiv 1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m} = x_j + x_m$ . Then

$$\int_0^{y_j} (h_j \circ \varphi_j)(\mathbf{x}_{-m}) (\partial_j \log p(\mathbf{x}_{-m}))^2 dx_j$$

$$\begin{aligned}
&\leq \int_0^{y_j} (h_j \circ \varphi_j)(\mathbf{x}_{-m}) \left( \|\mathbf{K}\|_1 x_j^{a-1} + \|\mathbf{K}\|_1 x_m^{a-1}(x_j) + |\eta_j| x_j^{b-1} + |\eta_m| x_m^{b-1}(x_j) \right)^2 dx_j \\
&\leq \int_0^{y_j/2} (h_j \circ \varphi_j)(\mathbf{x}_{-m}) \left( \|\mathbf{K}\|_1 x_j^{a-1} + |\eta_j| x_j^{b-1} + c_{1,m}(\mathbf{x}_{-j,-m}) \right)^2 dx_j \\
&\quad + \int_{y_j/2}^{y_j} (h_j \circ \varphi_j)(\mathbf{x}_{-m}) \left( \|\mathbf{K}\|_1 x_m^{a-1}(x_j) + |\eta_m| x_m^{b-1}(x_j) + c_{1,j}(\mathbf{x}_{-j,-m}) \right)^2 dx_j \\
&= \int_0^{y_j/2} h_j(x_j) \left( \|\mathbf{K}\|_1 x_j^{a-1} + |\eta_j| x_j^{b-1} + c_{1,m}(\mathbf{x}_{-j,-m}) \right)^2 dx_j \\
&\quad + \int_0^{y_j/2} h_j(x_j) \left( \|\mathbf{K}\|_1 x_j^{a-1} + |\eta_m| x_j^{b-1} + c_{1,j}(\mathbf{x}_{-j,-m}) \right)^2 dx_j
\end{aligned}$$

where in the last step we used change of variable  $x_j \leftarrow x_m(x_j) = y_j - x_j$  for the second term, and where

$$0 < c_{1,j} \equiv \max_{y_j/2 \leq x_j \leq y_j} \left( \|\mathbf{K}\|_1 x_j^{a-1} + |\eta_j| x_j^{b-1} \right) = \mathcal{O}(y_j^{a-1}) + \mathcal{O}(y_j^{b-1}) + \mathcal{O}(1) < +\infty,$$

where the  $\mathcal{O}$  depends on  $\mathbf{K}$  and  $\boldsymbol{\eta}$ . We thus have

$$\begin{aligned}
&\int_{\substack{\mathbf{x}_{-m} \succ \mathbf{0}, \\ \mathbf{1}_{m-1}^\top \mathbf{x}_{-m} < 1}} (h_j \circ \varphi_j)(\mathbf{x}_{-m}) (\partial_j \log p(\mathbf{x}_{-m}))^2 d\mathbf{x}_{-m} \\
&= \int_{\substack{\mathbf{x}_{-j,-m} \succ \mathbf{0}, \\ y_j(\mathbf{x}_{-j,-m}) > 0}} \int_0^{y_j(\mathbf{x}_{-j,-m})} (h_j \circ \varphi_j)(\mathbf{x}_{-m}) (\partial_j \log p(\mathbf{x}_{-m}))^2 dx_j d\mathbf{x}_{-j,-m} \\
&\leq \int_{\substack{\mathbf{x}_{-j,-m} \succ \mathbf{0}, \\ y_j(\mathbf{x}_{-j,-m}) > 0}} \int_0^{y_j(\mathbf{x}_{-j,-m})/2} h_j(x_j) \left( \mathcal{O}(x_j^{a-1}) + \mathcal{O}(x_j^{b-1}) + \mathcal{O}(y_j^{a-1}) + \mathcal{O}(y_j^{b-1}) + \mathcal{O}(1) \right)^2 dx_j d\mathbf{x}_{-j,-m} \\
&\leq \int_{\substack{\mathbf{x}_{-j,-m} \succ \mathbf{0}, \\ y_j(\mathbf{x}_{-j,-m}) > 0}} \int_0^{y_j(\mathbf{x}_{-j,-m})/2} h_j(x_j) \left( \mathcal{O}(x_j^{a-1}) + \mathcal{O}(x_j^{b-1}) + \mathcal{O}(1) \right)^2 dx_j d\mathbf{x}_{-j,-m} \\
&\quad + \int_{\substack{\mathbf{x}_{-j,-m} \succ \mathbf{0}, \\ y_j(\mathbf{x}_{-j,-m}) > 0}} \int_0^{y_j(\mathbf{x}_{-j,-m})/2} h_j(x_j) \left( \mathcal{O}(y_j^{a-1}) + \mathcal{O}(y_j^{b-1}) + \mathcal{O}(1) \right)^2 dx_j d\mathbf{x}_{-j,-m} \\
&\leq \int_{\mathbf{1} \succ \mathbf{x}_{-j,-m} \succ \mathbf{0}} \int_0^1 x_j^{\alpha_j} \left( \mathcal{O}(x_j^{a-1}) + \mathcal{O}(x_j^{b-1}) + \mathcal{O}(1) \right)^2 dx_j d\mathbf{x}_{-j,-m} \\
&\quad + \int_{\substack{\mathbf{x}_{-j,-m} \succ \mathbf{0}, \\ y_j(\mathbf{x}_{-j,-m}) > 0}} \frac{(y_j/2)^{\alpha_j+1}}{\alpha_j+1} \left( \mathcal{O}(y_j^{a-1}) + \mathcal{O}(y_j^{b-1}) + \mathcal{O}(1) \right)^2 d\mathbf{x}_{-j,-m} \\
&\leq \int_0^1 \mathcal{O}\left(x_j^{2a-2+\alpha_j}\right) + \mathcal{O}\left(x_j^{2b-2+\alpha_j}\right) + \mathcal{O}\left(x_j^{\alpha_j}\right) dx_j
\end{aligned}$$

$$\begin{aligned}
& + \sum_{p \in \{2a-1+\alpha_j, 2b-1+\alpha_j, \alpha_j+1\}} \int_{\substack{\mathbf{x}_{-j,-m} \succ \mathbf{0}, \\ \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m} < 1}} \mathcal{O}((1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m})^p) d\mathbf{x}_{-j,-m} \\
& = \int_0^1 o(x_j^{a-1}) + o(x_j^{b-1}) + o(x^0) dx_j + \sum_{p \in \{a,b,1\}} \mathcal{O}(\Gamma(p+1)/\Gamma(p+m-1)) < +\infty
\end{aligned}$$

for  $\alpha_j > \max\{0, 1-a, 1-b\}$ , where the second term of the last quantity follows from the inverse normalizing constant of the Dirichlet distribution with parameters  $(\mathbf{1}_{m-2}, p+1)$ .

- ii) If  $b = 0$ ,  $p_0(\mathbf{x}_{-m})$  is now bounded by  $C_2 \prod_{j=1}^m x_j^{\eta_{0j}}$  where  $C_2$  is the inverse normalizing constant of  $p_0(\mathbf{x}_{-m})$  times  $\sup_{\mathbf{x} \succ \mathbf{0}, \mathbf{1}^\top \mathbf{x} = 1} \exp(-\mathbf{x}^{a^\top} \mathbf{K}_0 \mathbf{x}^a / (2a))$ , a positive and finite constant. Then by the same reasoning as in i), with  $y_j(\mathbf{x}_{-j,-m}) \equiv 1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m}$  and noting that  $\boldsymbol{\eta} \succ -\mathbf{1}_m$ ,

$$\begin{aligned}
& \int_{\substack{\mathbf{x}_{-m} \succ \mathbf{0}, \mathbf{1}_{m-1}^\top \mathbf{x}_{-m} < 1}} p_0(\mathbf{x}_{-m})(h_j \circ \varphi_j)(\mathbf{x}_{-m}) (\partial_j \log p(\mathbf{x}_{-m}))^2 d\mathbf{x}_{-m} \\
& \leq \int_{\substack{\mathbf{x}_{-j,-m} \succ \mathbf{0}, y_j > 0}} \int_0^{y_j/2} C_2 \prod_{k=1}^m x_k^{\eta_{0k}} h_j(x_j) (\mathcal{O}(x_j^{a-1}) + \mathcal{O}(x_j^{-1}) + \mathcal{O}(1))^2 dx_j d\mathbf{x}_{-j,-m} \\
& \quad + \int_{\substack{\mathbf{x}_{-j,-m} \succ \mathbf{0}, y_j > 0}} \int_0^{y_j/2} C_2 \prod_{k=1}^m x_k^{\eta_{0k}} h_j(x_j) (\mathcal{O}(y_j^{a-1}) + \mathcal{O}(y_j^{-1}) + \mathcal{O}(1))^2 dx_j d\mathbf{x}_{-j,-m} \\
& \leq C_2 \prod_{k \neq j, m} \int_0^1 x_k^{\eta_{0k}} dx_k \int_0^1 x_j^{\eta_{0k} + \alpha_j} (\mathcal{O}(x_j^{-1}) + \mathcal{O}(1))^2 dx_j \\
& \quad + C_2 \int_{\substack{\mathbf{x}_{-j,-m} \succ \mathbf{0}, y_j > 0}} \frac{(y_j/2)^{\alpha_j + 1 + \eta_{0j}}}{\alpha_j + 1 + \eta_{0j}} \prod_{k \neq j, m} x_k^{\eta_{0k}} (\mathcal{O}(y_j^{-1}) + \mathcal{O}(1))^2 d\mathbf{x}_{-j,-m} \\
& \leq C_2 \prod_{k \neq j, m} \frac{1}{\eta_{0k} + 1} \int_0^1 \mathcal{O}(x_j^{-1}) + \mathcal{O}(x_j) dx_j \\
& \quad + \sum_{p \in \{0,2\}} C_2 \int_{\substack{\mathbf{x}_{-j,-m} \succ \mathbf{0}, \\ y_j > 0}} \prod_{k \neq j, m} x_k^{\eta_{0k}} \mathcal{O}(y_j^p) d\mathbf{x}_{-j,-m} \\
& < +\infty
\end{aligned}$$

since the integral in the second term is the inverse normalizing constant of the Dirichlet distribution with parameters  $(\boldsymbol{\eta}_{0,-j,-m} + \mathbf{1}_{m-2}, p+1)$ , i.e.  $\frac{\Gamma(p+1) \prod_{k \neq j, m} \Gamma(\eta_{0k} + 1)}{\Gamma(\mathbf{1}_{m-2}^\top \boldsymbol{\eta}_{0,-j,-m} + p + m - 1)} < +\infty$ .

This ends the proof for the first integrability condition for (A.1) for  $a > 0$ . For the second half, the integrand we consider is  $p_0(\mathbf{x}_{-m}) |\partial_j (\partial_j \log p(\mathbf{x}_{-m})(h_j \circ \varphi)(\mathbf{x}_{-m}))|$ . The arguments

are similar to those for the first condition, where we first bound  $\partial_{jj} \log p(\mathbf{x})$  using sums of products of powers of  $\mathbf{x}$ . Then for each fixed  $\mathbf{x}_{-j,-m}$  we split the domain of  $x_j$  into two halves and deal with the potential singularity at  $x_j \searrow 0^+$ , where one can show that the requirement on  $\alpha_j$  is just enough for the integrand to be  $o(x_j^{-1})$  and thus the integral is finite. The detailed proof is tedious and is omitted.

(2) *Case a = 0 and b ≥ 0:* First consider  $b = 0$ . We again write  $y_j \equiv 1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m}$ . Fixing  $\mathbf{x}_{-j,-m} \in \mathcal{S}_{-j, \mathcal{D}_{-m}}$ ,

$$\begin{aligned} & p_0(\mathbf{x}_{-m}) |\partial_j \log p(\mathbf{x}_{-m})| \tag{B.5} \\ \propto & \exp \left[ -\frac{1}{2} \log(\mathbf{x}_{-m})^\top \mathbf{K}_{0,-m,-m} \log(\mathbf{x}_{-m}) - \log(\mathbf{x}_{-m})^\top \boldsymbol{\kappa}_{0,-m,m} \log x_m \right. \\ & \left. - \frac{1}{2} \kappa_{0,m,m} (\log x_m)^2 + \boldsymbol{\eta}_{0,-m}^\top \log(\mathbf{x}_{-m}) + \eta_{0,m} \log x_m \right] \times \\ & \left| -\boldsymbol{\kappa}_{-m,j}^\top \log(\mathbf{x}_{-m})/x_j + \boldsymbol{\kappa}_{-m,m}^\top \log(\mathbf{x}_{-m})/x_m - \kappa_{jm} \log x_m/x_j \right. \\ & \left. + \kappa_{mm} \log x_m/x_m + \eta_j/x_j - \eta_m/x_m \right|. \tag{B.6} \end{aligned}$$

$$\begin{aligned} \leq & \prod_{k \neq j,m} \exp \left[ -\frac{N_{\mathbf{K}_{0,-m,-m}}}{2} (\log x_k)^2 + \eta_{0,k} \log x_k \right] \exp \left[ -\frac{N_{\mathbf{K}_{0,-m,-m}}}{2} (\log x_j)^2 + \eta_{0,j} \log x_j \right. \\ & \left. - \log(\mathbf{x}_{-m})^\top \boldsymbol{\kappa}_{0,-m,m} \log x_m - \frac{1}{2} \kappa_{0,m,m} (\log x_m)^2 + \eta_{0,m} \log x_m \right] \times \\ & \left| -\boldsymbol{\kappa}_{-m,j}^\top \log(\mathbf{x}_{-m})/x_j + \boldsymbol{\kappa}_{-m,m}^\top \log(\mathbf{x}_{-m})/x_m - \kappa_{jm} \log x_m/x_j \right. \\ & \left. + \kappa_{mm} \log x_m/x_m + \eta_j/x_j - \eta_m/x_m \right|. \tag{B.7} \end{aligned}$$

which is  $\mathcal{O}(\exp(\mathcal{O}((\log x_j)^2) + \mathcal{O}(\log x_j) + \mathcal{O}(\log \log x_j)))$  as  $x_j \searrow 0^+$ . Since the coefficient on the leading term is negative the entire term goes to 0. By symmetry the quantity goes to zero also when  $x_j \nearrow y_j^-$ . Thus, (A.1) holds for any  $\alpha_j \geq 0$ .

Similarly, both  $p_0(\mathbf{x}_{-m})(h_j \circ \varphi_j)(\mathbf{x}_{-m})(\partial_j \log p(\mathbf{x}_{-m}))^2$  and  $p_0(\mathbf{x}_{-m})|\partial(\partial_j \log p(\mathbf{x}_{-m}))(h_j \circ \varphi)(\mathbf{x}_{-m})|$  are products of  $\prod_{j=1}^m \mathcal{O}(\exp(-(\log x_j)^2))$  times a polynomial, and are thus bounded and go to 0 at the boundaries of  $\mathcal{D}_{-m}$ . Thus extending the

integrands to 0 at the boundaries, they are continuous and bounded in the compact  $\overline{\mathcal{D}_{-m}}$ , so integrals  $\int_{\mathcal{D}_{-m}} p_0(\mathbf{x}_{-m})(h_j \circ \varphi_j)(\mathbf{x}_{-m})(\partial_j \log p(\mathbf{x}_{-m}))^2 d\mathbf{x}_{-m}$  and  $\int_{\mathcal{D}_{-m}} p_0(\mathbf{x}_{-m}) |\partial(\partial_j \log p(\mathbf{x}_{-m})(h_j \circ \varphi_j)(\mathbf{x}_{-m}))| d\mathbf{x}_{-m}$  are finite, thus proving (A.2).

For  $a = 0$  and  $b > 0$ ,  $\boldsymbol{\eta} \preceq \mathbf{0}$ , and the proof is similar and is omitted. In particular,  $p_0(\mathbf{x}_{-m})$  is bounded by that with  $a = 0$ ,  $b = 0$ ,  $\boldsymbol{\eta} \equiv \mathbf{0}_m$ , and thus its product with any polynomial is bounded and goes to 0 at the boundary of  $\overline{\mathcal{D}_{-m}}$ .  $\square$

*Proof of Theorem 20.* For notational simplicity, denote  $\tilde{\mathbf{K}} \equiv \mathbf{K}_1 - \mathbf{K}_2$  with columns  $\tilde{\boldsymbol{\kappa}}_1, \dots, \tilde{\boldsymbol{\kappa}}_m$ , and denote  $\tilde{\boldsymbol{\eta}} \equiv \boldsymbol{\eta}_1 - \boldsymbol{\eta}_2$ . Assume that either  $\tilde{\mathbf{K}} \neq \mathbf{0}_{m \times m}$  or  $\tilde{\boldsymbol{\eta}} \neq \mathbf{0}_m$ , otherwise there is nothing to prove. By (3.26), writing  $x_m \equiv 1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m}$  and  $\mathbf{x} = (\mathbf{x}_{-m}; x_m)$  and taking the gradient of the log of both sides of the equation with respect to  $x_j$ ,  $j = 1, \dots, m-1$ , we have

$$\left( x_j^{a-1} \tilde{\boldsymbol{\kappa}}_{j,j} - x_m^{a-1} \tilde{\boldsymbol{\kappa}}_{j,m} \right)^\top \mathbf{x}^a = \tilde{\eta}_j x_j^{b-1} - \tilde{\eta}_m x_m^{b-1} \quad (\text{B.8})$$

for all  $\mathbf{x}_{-m} \in \mathcal{D}_{-m} \equiv \{ \mathbf{x}_{-m} \in \mathbb{R}_+^{m-1} \mid \mathbf{x}_{-m} \succ \mathbf{0}, \mathbf{1}_{m-1}^\top \mathbf{x}_{-m} < 1 \}$ . In the following when  $a = 0$  by  $x^a$  we mean  $\log(x)$ , and by  $x^{a-1}$  we mean  $1/x$  and we do not treat this case differently as the same expressions still hold.

(1) Suppose  $\left( x_j^{a-1} \tilde{\boldsymbol{\kappa}}_{j,j} - x_m^{a-1} \tilde{\boldsymbol{\kappa}}_{j,m} \right)_{-j,-m} = \mathbf{0}_{m-2}$  for all  $\mathbf{x}_{-m} \in \mathcal{D}_{-m}$  and  $x_m = 1 - \mathbf{1}_{-m}^\top \mathbf{x}_{-m}$  or  $\tilde{\eta}_j$ .

(i) Suppose  $a = 1$ , then (B.8) becomes  $(\tilde{\boldsymbol{\kappa}}_{j,j} - \tilde{\boldsymbol{\kappa}}_{j,m}) x_j + (\tilde{\boldsymbol{\kappa}}_{m,j} - \tilde{\boldsymbol{\kappa}}_{m,m}) x_m = \tilde{\eta}_j x_j^{b-1} - \tilde{\eta}_m x_m^{b-1}$ , and we must have  $b = 2$  or  $b = 1$  or  $\tilde{\eta}_j = \tilde{\eta}_m = 0$ , i.e.  $\boldsymbol{\eta}_1 = \boldsymbol{\eta}_2$ .

(ii) Suppose  $a \neq 1$ . The assumption implies that  $(\tilde{\boldsymbol{\kappa}}_j)_{-j,-m} = (\tilde{\boldsymbol{\kappa}}_m)_{-j,-m} = \mathbf{0}_{m-2}$ . Then (B.8) becomes  $(x_j^{a-1} \tilde{\boldsymbol{\kappa}}_{j,j} - x_m^{a-1} \tilde{\boldsymbol{\kappa}}_{j,m}) x_j^a + (x_j^{a-1} \tilde{\boldsymbol{\kappa}}_{m,j} - x_m^{a-1} \tilde{\boldsymbol{\kappa}}_{m,m}) x_m^a = \tilde{\eta}_j x_j^{b-1} - \tilde{\eta}_m x_m^{b-1}$ . Since  $x_j > 0$  and  $x_m > 0$  are arbitrary (as  $\mathbf{1}_{-m}^\top \mathbf{x}_{-m}$  can vary) as long as  $x_j + x_m < 1$ , the cross terms must not exist, and so  $\tilde{\boldsymbol{\kappa}}_{j,m} = \tilde{\boldsymbol{\kappa}}_{m,j} = 0$ . It thus follows that  $\tilde{\boldsymbol{\kappa}}_{-j,j} = \tilde{\boldsymbol{\kappa}}_{-m,m} = \mathbf{0}_{m-1}$  and hence  $\tilde{\mathbf{K}}$  is diagonal, and the original equality becomes  $-\frac{1}{2} \text{diag}(\tilde{\mathbf{K}})^\top (\mathbf{x}^a)^2 + \tilde{\boldsymbol{\eta}}^\top \mathbf{x}^b = 0$ , in which by  $\mathbf{x}^0$  we mean  $\log(\mathbf{x})$ . Thus we must have  $2a = b \neq 0$  and  $\mathbf{K}_1 - \mathbf{K}_2 = 2\boldsymbol{\eta}_1 - 2\boldsymbol{\eta}_2$ .

(2) Now fix  $x_j$  and  $x_m$  such that  $\left(x_j^{a-1}\tilde{\boldsymbol{\kappa}}_{j,j} - x_m^{a-1}\tilde{\boldsymbol{\kappa}}_{m,m}\right)_{-j,-m} \neq \mathbf{0}_{m-2}$ . Note that  $\mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m} = 1 - x_j - x_m$  is also fixed. Now the right-hand side and the first vector on the left-hand side of (B.8) are both constant, while  $\mathbf{x}_{-j,-m}$  is allowed to vary freely as long as their sum is fixed. A necessary condition of (B.8) is thus

$$\left(x_j^{a-1}\tilde{\boldsymbol{\kappa}}_{j,j} - x_m^{a-1}\tilde{\boldsymbol{\kappa}}_{m,m}\right)_{-j,-m}^\top \mathbf{x}_{-j,-m}^a = \text{const depending on } x_j \text{ and } x_m \text{ only} \quad (\text{B.9})$$

for all  $\mathbf{x}_{-j,-m}^a \in \mathcal{U}_{x_j, x_m} \equiv \{\mathbf{y}^a : \mathbf{y} \succ \mathbf{0}_{m-2}, \mathbf{1}_{m-2}^\top \mathbf{y} = \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m} = 1 - x_j - x_m\}$ .

Suppose  $a \neq 1$ . Then  $\mathcal{U}_{x_j, x_m}$  is not entirely on a hyperplane, and by assumption  $\left(x_j^{a-1}\tilde{\boldsymbol{\kappa}}_{j,j} - x_m^{a-1}\tilde{\boldsymbol{\kappa}}_{m,m}\right)_{-j,-m}$  is not a zero vector, so the equality cannot hold. We thus have  $a = 1$ , so that  $\mathcal{U}_{x_j, x_m}$  lies on the hyperplane  $\mathcal{H}_{x_j, x_m} \equiv \{\mathbf{y} : \mathbf{1}_{m-2}^\top \mathbf{y} = 1 - x_j - x_m\}$ . Since  $\mathcal{H}_{x_j, x_m} \equiv \{c\mathbf{x}_{-j,-m} : c \in \mathbb{R}, \mathbf{x} \in \mathcal{U}_{x_j, x_m}\}$ , (B.9) must hold for all  $\mathbf{x}_{-j,-m}$  in the hyperplane  $\mathcal{H}_{x_j, x_m}$ , and by the assumption that  $\left(\tilde{\boldsymbol{\kappa}}_{j,j} - \tilde{\boldsymbol{\kappa}}_{m,m}\right)_{-j,-m}$  is nonzero it must be a constant multiple of  $\mathbf{1}_m$ , and the right-hand side of (B.9) is hence  $c_0(1 - x_j - x_m)$  for some absolute constant  $c_0 \neq 0$  assuming  $\left(\tilde{\boldsymbol{\kappa}}_{j,j} - \tilde{\boldsymbol{\kappa}}_{m,m}\right)_{-j,-m} = c_0 \mathbf{1}_m$ . Plugging this back in (B.8) we get

$$c_0(1 - x_j - x_m) + \left(\tilde{\kappa}_{j,j} - \tilde{\kappa}_{j,m}\right) x_j + \left(\tilde{\kappa}_{m,j} - \tilde{\kappa}_{m,m}\right) x_m = \tilde{\eta}_j x_j^{b-1} - \tilde{\eta}_m x_m^{b-1},$$

and hence as in (1) (i) we have  $b = 2$  or  $b = 1$  or  $\boldsymbol{\eta}_1 = \boldsymbol{\eta}_2$ .

□

*Proof of Theorem 21.* For any  $\mathbf{x}_{-m} \in \mathcal{D}_{-m}$ , write  $\iota_{+m}(\mathbf{x}_{-m}) = (\mathbf{x}_{-m}, 1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m})$ . We first prove the finiteness of the normalizing constant. If  $\mathbf{K}$  is positive definite, the inverse normalizing constant is

$$\begin{aligned} & \int_{\mathcal{D}_{-m}} \exp\left(-\frac{1}{2} \log \iota_{+m}(\mathbf{x}_{-m})^\top \mathbf{K} \log \iota_{+m}(\mathbf{x}_{-m}) + \boldsymbol{\eta}^\top \log \iota_{+m}(\mathbf{x}_{-m})\right) d\mathbf{x}_{-m} \\ & \leq \int_{\mathcal{D}_{-m}} \exp\left(\sum_{j=1}^m \left(-\lambda_{\min}(\mathbf{K}) (\log \iota_{+m}(\mathbf{x}_{-m}))_j^2 + \eta_j (\log \iota_{+m}(\mathbf{x}_{-m}))_j\right)\right) d\mathbf{x}_{-m} \\ & \leq \int_{\mathcal{D}_{-m}} \exp\left(\sum_{j=1}^m \frac{\eta_j^2}{4\lambda_{\min}(\mathbf{K})}\right) d\mathbf{x}_{-m} < +\infty, \end{aligned}$$

proving (1). Now assume  $\mathbf{K}$  is no longer positive definite. If  $\mathbf{K}\mathbf{1}_m = \mathbf{0}_m$ , for any  $\mathbf{a} \in \mathbb{R}^m$ ,

$$\begin{aligned} \mathbf{a}^\top \mathbf{K} \mathbf{a} &= (\mathbf{a}_{-j}, a_j)^\top \begin{bmatrix} \mathbf{K}_{-j,-j} & \boldsymbol{\kappa}_{-j,j} \\ \boldsymbol{\kappa}_{-j,j}^\top & \kappa_{jj} \end{bmatrix} (\mathbf{a}_{-j}, a_j) \\ &= (\mathbf{a}_{-j}, a_j)^\top \begin{bmatrix} \mathbf{K}_{-j,-j} & -\mathbf{K}_{-j,-j} \mathbf{1}_{m-1} \\ -\mathbf{1}_{m-1}^\top \mathbf{K}_{-j,-j} & \mathbf{1}_{m-1}^\top \mathbf{K}_{-j,-j} \mathbf{1}_{m-1} \end{bmatrix} (\mathbf{a}_{-j}, a_j) \\ &= (\mathbf{a}_{-j} - a_j \mathbf{1}_m)^\top \mathbf{K}_{-j,-j} (\mathbf{a}_{-j} - a_j \mathbf{1}_{m-1}), \end{aligned}$$

which is zero if and only if  $\mathbf{a}_{-j} = a_j \mathbf{1}_{m-1}$ , i.e.  $a_1 = \dots = a_m$ , and is positive otherwise. Thus, if  $\mathbf{K}\mathbf{1}_m = \mathbf{0}_m$ , the condition that  $\mathbf{K}_{-j,-j}$  is positive definite for some  $j = 1, \dots, m$  is equivalent to that  $\mathbf{K}_{-j,-j}$  is positive definite for all  $j = 1, \dots, m$ , and implies that  $\mathbf{K}$  is positive semi-definite.

If  $\mathbf{K}$  is positive semi-definite and  $\boldsymbol{\eta} \succ -\mathbf{1}_m$ , the inverse normalizing constant is

$$\begin{aligned} &\int_{\mathcal{D}_{-m}} \exp\left(-\frac{1}{2} \log \iota_{+m}(\mathbf{x}_{-m})^\top \mathbf{K} \log \iota_{+m}(\mathbf{x}_{-m}) + \boldsymbol{\eta}^\top \log \iota_{+m}(\mathbf{x}_{-m})\right) d\mathbf{x}_{-m} \\ &\leq \int_{\mathcal{D}_{-m}} \exp(\boldsymbol{\eta}^\top \log \iota_{+m}(\mathbf{x}_{-m})) d\mathbf{x}_{-m} = \frac{\prod_{j=1}^m \Gamma(\eta_j + 1)}{\Gamma(\mathbf{1}_m^\top \boldsymbol{\eta} + m)} < +\infty \end{aligned}$$

since the last quantity is the inverse normalizing constant of the Dirichlet distribution with parameters  $(\boldsymbol{\eta} + \mathbf{1}_m)$ , proving (3).

On the other hand, suppose  $\mathbf{K}\mathbf{1}_m = \mathbf{0}_m$  and  $\mathbf{K}_{-j,-j}$  is positive definite for some/all  $j = 1, \dots, m$ . Again letting  $\mathbf{x}_{-m}$  be the free variables and letting  $x_m = 1 - \mathbf{1}_{m-1}^\top \mathbf{x}_{-m}$ , define the *additive log-ratio transformation* applied to  $\mathbf{x}$ :  $\mathbf{y}_{-m} \equiv \text{alt}(\mathbf{x}) \equiv \log \mathbf{x}_{-m} - (\log x_m) \mathbf{1}_{m-1}$ , a random vector supported on  $\mathbb{R}^{m-1}$ . Append an extra  $y_m = 0$  for ease of notation. The transformation is thus bijective and the inverse transformation, the *additive logistic transformation*  $\mathbf{x} = \exp(\mathbf{y}) / \mathbf{1}_m^\top \exp(\mathbf{y})$ . Since

$$\partial x_k / \partial y_j = -x_k x_j, \quad \partial x_j / \partial y_j = x_j (1 - x_j)$$

for  $j \neq k$ ,  $j, k = l, \dots, m-1$ , we have

$$\left| \frac{\partial \mathbf{x}_{-m}}{\partial \mathbf{y}_{-m}} \right| = \begin{vmatrix} x_1(1-x_1) & -x_1x_2 & \cdots & -x_1x_{m-1} \\ -x_1x_2 & x_2(1-x_2) & \cdots & -x_2x_{m-1} \\ \vdots & \vdots & \ddots & \vdots \\ -x_1x_{m-1} & -x_2x_{m-1} & \cdots & x_{m-1}(1-x_{m-1}) \end{vmatrix} = \prod_{j=1}^m x_j = \exp(\mathbf{1}_m^\top \log \mathbf{x}).$$

Then by  $\mathbf{1}_m^\top \mathbf{K} = \mathbf{K} \mathbf{1}_m = \mathbf{0}_m$ ,  $\mathbf{y}_{-m}$  has density proportional to

$$\begin{aligned} & p(\mathbf{x}_{-m}) |\partial \mathbf{x}_{-m} / \partial \mathbf{y}_{-m}| \\ \propto & \exp \left( -\frac{1}{2} (\log \mathbf{x} - (\log x_m) \mathbf{1}_m)^\top \mathbf{K} (\log \mathbf{x} - (\log x_m) \mathbf{1}_m) \right. \\ & \left. + (\boldsymbol{\eta} + \mathbf{1}_m)^\top (\log \mathbf{x} - (\log x_m) \mathbf{1}_m) + (\log x_m) \mathbf{1}_m^\top (\boldsymbol{\eta} + \mathbf{1}_m) \right) \\ = & \exp \left( -\frac{1}{2} \mathbf{y}^\top \mathbf{K} \mathbf{y} + (\boldsymbol{\eta} + \mathbf{1}_m)^\top \mathbf{y} + \mathbf{1}_m^\top (\boldsymbol{\eta} + \mathbf{1}_m) \log x_m \right) \\ = & \exp \left( -\frac{1}{2} \mathbf{y}_{-m}^\top \mathbf{K}_{-m,-m} \mathbf{y}_{-m} + (\boldsymbol{\eta}_{-m} + \mathbf{1}_{m-1})^\top \mathbf{y}_{-m} - (\mathbf{1}_m^\top \boldsymbol{\eta} + m) \log (1 + \mathbf{1}_{m-1}^\top \exp(\mathbf{y}_{-m})) \right). \end{aligned}$$

Note that the  $\log x_m = -\log (1 + \mathbf{1}_{m-1}^\top \exp(\mathbf{y}_{-m})) < 0$ , so for  $\mathbf{1}_m^\top \boldsymbol{\eta} + m \geq 0$  the last display is always upper-bounded by a constant times a normal density with a positive definite inverse covariance matrix  $\mathbf{K}_{-m,-m}$ , and thus the normalizing constant is finite, thus proving (2).

As for (A.1), fix  $j = 1, \dots, m-1$  and any  $\ell \in \{1, \dots, m-1\} \setminus \{j\}$ , and write  $\mathbf{z} \equiv \log \mathbf{x} - (\log x_\ell) \mathbf{1}_m$ . Fix any  $\mathbf{x}_{-j,-m} \in \mathbb{R}^{m-2}$  with  $\mathbf{x}_{-j,-m} \succ \mathbf{0}$ ,  $\mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m} < 1$ . Then if (1)  $\mathbf{K}_0$  is positive definite or (2)  $\mathbf{K}_0 \mathbf{1}_m = \mathbf{0}_m$  and  $\mathbf{K}_{0,-\ell,-\ell}$  is positive definite, by the proof above,  $p_0(\mathbf{x}_{-m}) x_i^t$  is upper bounded by a finite constant depending on  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0$  for any  $t \in \mathbb{R}$  and  $i = j, m$ , since it is a constant times the density with parameters  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0 + t \mathbf{e}_i$ , and since we did not impose any restriction on the  $\boldsymbol{\eta}$  parameter. On the other hand, for (3)  $\mathbf{K}_0 \mathbf{1}_m = \mathbf{0}_m$ ,  $\mathbf{K}_{0,-\ell,-\ell}$  is positive semi-definite and  $\boldsymbol{\eta}_0 \succ -\mathbf{1}_m$ ,

$$\begin{aligned} p_0(\mathbf{x}_{-m}) & \propto \exp \left( -\frac{1}{2} \log \mathbf{x}^\top \mathbf{K}_0 \log \mathbf{x} + \boldsymbol{\eta}_0^\top \log \mathbf{x} \right) \\ & = \exp \left( -\frac{1}{2} \mathbf{z}_{-\ell}^\top \mathbf{K}_{0,-\ell,-\ell} \mathbf{z}_{-\ell} + \boldsymbol{\eta}_{0,-\ell}^\top \mathbf{z}_{-\ell} + (\mathbf{1}_m^\top \boldsymbol{\eta}_0) \log x_\ell \right) \end{aligned}$$

$$\begin{aligned}
&\leq \exp(\boldsymbol{\eta}_{0,-\ell}^\top \mathbf{z}_{-\ell} + (\mathbf{1}_m^\top \boldsymbol{\eta}_0) \log x_\ell) \\
&\propto \exp(\eta_{0,j} z_j + \eta_{0,m} z_m).
\end{aligned} \tag{B.10}$$

On the other hand,

$$\begin{aligned}
&|\partial_j \log p(\mathbf{x}_{-m})| \min\{x_j, x_m\}^{\alpha_j} \\
&= |-\boldsymbol{\kappa}_{\cdot,j}^\top \log \mathbf{x}/x_j + \boldsymbol{\kappa}_{\cdot,m}^\top \log \mathbf{x}/x_m + \eta_j/x_j - \eta_m/x_m| \min\{x_j, x_m\}^{\alpha_j} \\
&= (|\boldsymbol{\kappa}_{\cdot,j}^\top \log \mathbf{x}/x_j| + |\boldsymbol{\kappa}_{\cdot,m}^\top \log \mathbf{x}/x_m| + |\eta_j/x_j| + |\eta_m/x_m|) \min\{x_j, x_m\}^{\alpha_j} \\
&\leq (|\boldsymbol{\kappa}_{\cdot,j}^\top \log \mathbf{x}| + |\boldsymbol{\kappa}_{\cdot,m}^\top \log \mathbf{x}| + |\eta_j| + |\eta_m|) \min\{x_j, x_m\}^{\alpha_j-1}.
\end{aligned}$$

Thus, as  $x_j \searrow 0^+$  or  $x_m \searrow 0^+$  (i.e.  $x_j \nearrow 1 - \mathbf{1}_{m-2}^\top \mathbf{x}_{-j,-m}$ ), by multiplying the two bounds we have  $p_0(\mathbf{x}_{-m}) |\partial_j \log p(\mathbf{x}_{-m})| (h_j \circ \varphi_j)(\mathbf{x}) \searrow 0^+$  for any  $\alpha_j$  for (1) and (2) (by letting  $t$  to e.g.  $\alpha_j - 2$  in the discussion above), or for (3) by a constant times

$$(|\boldsymbol{\kappa}_{\cdot,j}^\top \log \mathbf{x}| + |\boldsymbol{\kappa}_{\cdot,m}^\top \log \mathbf{x}| + |\eta_j| + |\eta_m|) \min\{x_j, x_m\}^{\alpha_j-1} x_j^{\eta_{0,j}} x_m^{\eta_{0,m}} \searrow 0^+$$

if  $\alpha_j > \max\{1 - \eta_{0,j}, 1 - \eta_{0,m}\}$ .

As for (A.2), the results follow by a similar discussion for the Gamma model ( $a$ - $b$  model with  $b = 0$ ) on the standard simplex in Section 3.6.  $\square$

*Proof of Theorem 24.* It suffices to bound  $\boldsymbol{\Gamma}$  and  $\mathbf{g}$  using their forms in Section 3.5.2 and apply Theorem 1 in Lin et al. (2016). Thus, we first find the bounds of  $(h_j \circ \varphi_j)(\mathbf{x}) x_j^{p_j} x_k^{p_k} x_\ell^{p_\ell}$  with  $h_j(x) = x^{\alpha_j}$ ,  $\alpha_j \geq 0$ ,  $\alpha_j \geq -p_j$ ,  $p_j \in \mathbb{R}$ ,  $p_k \geq 0$ ,  $p_\ell \geq 0$  and  $x_i \in [u_i, v_i]$  for  $i = 1, \dots, m$ . Suppose without loss of generality that  $j, k, \ell$  are all different, as

$$\max_{x_j} f_{j,1}(x_j) \max_{x_j} f_{j,2}(x_j) \max_{x_j} f_{j,3}(x_j) \geq \max_{x_j} (f_{j,1}(x_j) f_{j,2}(x_j) f_{j,3}(x_j)) \geq 0$$

for any nonnegative functions  $f_{j,1}, f_{j,2}, f_{j,3}$ .

As  $x_j$  approaches its boundary,  $\varphi_j(\mathbf{x}) \searrow 0^+$  and hence  $(h_j \circ \varphi_j)(\mathbf{x}) x_j^{p_j} \searrow 0^+$  if  $\alpha_j > -p_j$ . The lower bound 0 for  $(h_j \circ \varphi_j)(\mathbf{x}) x_j^{p_j} x_k^{p_k} x_\ell^{p_\ell}$  is thus tight enough.

As for the upper bound, the only way for the quantity to be unbounded from above is when  $x_j \searrow 0^+$  and  $p_j < 0$ , but as  $x_j \searrow 0^+$ ,  $(h_j \circ \varphi_j)(\mathbf{x}) = x_j^{\alpha_j}$  so this cannot happen with

the choice of  $\alpha_j \geq -p_j$ . Noting that  $h_j$  is monotonically increasing, we consider the following cases:

- (1) Suppose  $x_j \geq (u_j + v_j)/2$ . Then  $(h_j \circ \varphi_j)(\mathbf{x}) \leq h_j(\min\{C_j, v_j - x_j\}) \leq h_j(\min\{C_j, (v_j - u_j)/2\}) \leq \min\{C_j^{\alpha_j}, (v_j - u_j)^{\alpha_j}/2^{\alpha_j}\}$ ,  $x_j^{p_j} \leq (u_j + v_j)^{p_j}/2^{p_j}$  if  $p_j < 0$  or  $x_j^{p_j} \leq v_j^{p_j}$  if  $p_j \geq 0$ .
- (2) Suppose  $x_j \leq (u_j + v_j)/2$ . Then  $(h_j \circ \varphi_j)(\mathbf{x})x_j^{p_j} \leq h_j(\min\{C_j, x_j - u_j\})x_j^{p_j} = \min\{C_j^{\alpha_j}, (x_j - u_j)^{\alpha_j}\}x_j^{p_j}$ . Now let  $f(x) = (\min\{C_j, x - u_j\})^{\alpha_j}x^{p_j}$ . Then  $(\log f(x))' = \alpha_j/(x - u_j)\mathbf{1}_{x < u_j + C_j} + p_j/x$ . For  $x \geq u_j + C_j$  this has the same sign as  $p_j$ , otherwise it is equal to  $((\alpha_j + p_j)x - u_j p_j)/(x(x - u_j)) \geq 0$  on  $(u_j, v_j)$  since  $x > u_j$ ,  $\alpha_j \geq -p_j$  and  $\alpha_j \geq 0$ . This implies that if  $p_j \geq 0$  or  $v_j - u_j \leq 2C_j$ ,  $f$  is increasing on  $(u_j, (u_j + v_j)/2)$ , and so  $(h_j \circ \varphi_j)(\mathbf{x})x_j^{p_j} \leq \min\{C_j, (v_j - u_j)/2\}^{\alpha_j} (u_j + v_j)^{p_j}/2^{p_j}$ ; otherwise,  $f$  is increasing on  $(u_j, u_j + C_j)$  and decreasing on  $(u_j + C_j, (u_j + v_j)/2)$ , so  $(h_j \circ \varphi_j)(\mathbf{x})x_j^{p_j} \leq C_j^{\alpha_j} (u_j + C_j)^{p_j}$ .

Thus, defining

$$\zeta(C_j, u_j, v_j, \alpha_j, p_j) \equiv \begin{cases} \min\{C_j, (v_j - u_j)/2\}^{\alpha_j} (u_j + v_j)^{p_j}/2^{p_j}, & p_j < 0, v_j - u_j \leq 2C_j, \\ \min\{C_j, (v_j - u_j)/2\}^{\alpha_j} (u_j + C_j)^{p_j}, & p_j < 0, v_j - u_j > 2C_j, \\ \min\{C_j, (v_j - u_j)/2\}^{\alpha_j} v_j^{p_j}, & p_j \geq 0, \end{cases}$$

we have  $0 \leq (h_j \circ \varphi_j)(\mathbf{x})x_j^{p_j}x_k^{p_k}x_\ell^{p_\ell} \leq \zeta(C_j, u_j, v_j, \alpha_j, p_j)v_k^{p_k}v_\ell^{p_\ell}$ . Now assume additionally that  $\alpha_j \geq \max\{1, 1 - p_j\}$ , then by  $h_j'(x) = \alpha_j x_j^{\alpha_j - 1}$ ,  $0 \leq \partial_j(h_j \circ \varphi_j)(\mathbf{x})x_j^{p_j}x_k^{p_k} \leq \alpha_j \zeta(C_j, u_j, v_j, \alpha_j - 1, p_j)v_k^{p_k}$ .

First assume  $a > 0$ . Then assuming  $\alpha_1, \dots, \alpha_m \geq \max\{1, 2 - 2a, 2 - 2b, 1 - a, 2 - a, 2 - b\} = \max\{1, 2 - a, 2 - b\}$ , using the form of  $\mathbf{\Gamma}$  and  $\mathbf{g}$  in Section 3.5.2, for all  $j, k, \ell$  we have  $0 \leq \gamma_{j,k,\ell}(\mathbf{x}) \leq \mathfrak{r} \equiv \max_{j,k=1,\dots,m} \max\{\zeta(C_j, u_j, v_j, \alpha_j, 2a - 2)v_k^{2a}, \zeta(C_j, u_j, v_j, \alpha_j, 2b - 2)\}$  and  $0 \leq g_{j,k}(\mathbf{x}) \leq \mathfrak{g} \equiv \max_{j,k=1,\dots,m} \max\{\alpha_j \zeta(C_j, u_j, v_j, \alpha_j - 1, a - 1)v_k^a + |a - 1|\zeta(C_j, u_j, v_j, \alpha_j, a - 2)v_k^a + a\zeta(C_j, u_j, v_j, \alpha_j, 2a - 2), \alpha_j \zeta(C_j, u_j, v_j, \alpha_j - 1, b - 1) + |b - 1|\zeta(C_j, u_j, v_j, \alpha_j, b - 2)\}$ .

Then by Hoeffding's inequality,

$$\mathbb{P} \left( \max_{j,k,\ell} |\gamma_{j,k,\ell} - \mathbb{E}_0 \gamma_{j,k,\ell}| \geq \epsilon_1/2 \right) \leq 2 \exp \left( -n\epsilon_1^2 / (2\varsigma_{\mathbf{r}}^2) \right), \quad (\text{B.11})$$

$$\mathbb{P} \left( \max_{j,k} |g_{j,k} - \mathbb{E}_0 g_{j,k}| \geq \epsilon_2 \right) \leq 2 \exp \left( -2n\epsilon_2^2 / \varsigma_{\mathbf{g}}^2 \right). \quad (\text{B.12})$$

Let  $\epsilon_1 \equiv \varsigma_{\mathbf{r}} \sqrt{2(\log m^\tau + \log 4)/n}$  and  $\epsilon_2 \equiv \varsigma_{\mathbf{g}} \sqrt{(\log m^\tau + \log 4)/(2n)}$ . With the choice of  $\delta \leq 1 + \sqrt{(\log m^\tau + \log 4)/(2n)}$  and by the fact that  $0 \leq \max_{j,k,\ell} \gamma_{j,k,\ell} \leq \varsigma_{\mathbf{r}} = \epsilon_1 / (2\delta - 2)$ , (B.11) and (B.12) imply that

$$\begin{aligned} \mathbb{P} \left( \max_{j,k,\ell} |\delta \gamma_{j,k,\ell} - \mathbb{E}_0 \gamma_{j,k,\ell}| \geq \epsilon_1 \right) &\leq \mathbb{P} \left( \max_{j,k,\ell} |\gamma_{j,k,\ell} - \mathbb{E}_0 \gamma_{j,k,\ell}| + (\delta - 1) \max_{j,k,\ell} \gamma_{j,k,\ell} \geq \epsilon_1 \right) \\ &\leq \mathbb{P} \left( \max_{j,k,\ell} |\gamma_{j,k,\ell} - \mathbb{E}_0 \gamma_{j,k,\ell}| \geq \epsilon_1/2 \right) \leq m^{-\tau}/2, \end{aligned} \quad (\text{B.13})$$

$$\mathbb{P} \left( \max_{j,k} |g_{j,k} - \mathbb{E}_0 g_{j,k}| \geq \epsilon_2 \right) \leq m^{-\tau}/2. \quad (\text{B.14})$$

The results then follow by applying Theorem 1 in Lin et al. (2016).

In the case where  $a = 0$ , and  $u_k > 0$  for all  $k$ ,

$$|(h_j \circ \varphi_j)(\mathbf{x}) x_j^{p_j} \log(x_k) \log(x_\ell)| \leq \zeta(C_j, u_j, v_j, \alpha_j, p_j) \cdot \max\{|\log(u_k) \log(u_\ell)|, |\log(v_k) \log(v_\ell)|\}$$

and everything else follows similarly as for  $a > 0$ .  $\square$

*Proof of Lemma 25.* We show that  $X_j^{2a}$  for  $a > 0$  or  $\log X_j$  for  $a = 0$  is sub-exponential by showing its moment-generating function is finite. Then the sub-exponentiality follows from Theorem 2.13 of Wainwright (2019).

First consider the case where  $a = 0$ . In Corollary 17, we only require  $\mathbf{K}$  to be positive definite without any restrictions on  $\boldsymbol{\eta}$ , and thus for any  $t \in \mathbb{R}$ ,  $\mathbb{E}_0 \exp(t \log X_j)$  is the inverse normalizing constant for the model with parameters  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0 + t\mathbf{e}_j$ , where  $\mathbf{e}_j$  is the vector with the  $j$ -th coordinate equal to 1 and the rest equal to 0, and is thus finite.

Next, consider  $a > 0$ . Corollary 17 requires  $\mathbf{K}_0$  to be positive definite, and in addition  $\boldsymbol{\eta}_0 \succ -\mathbf{1}_m$  if  $b = 0$ . Then, again writing  $\mathbf{x}^0/0 = \log \mathbf{x}$  for the  $b$  part,

$$p_0(\mathbf{x}) \exp(tx_j^{2a}) \propto \exp \left( -\frac{1}{2a} \mathbf{x}^{a\top} \mathbf{K}_0 \mathbf{x}^a + \frac{1}{b} \boldsymbol{\eta}_0^\top \mathbf{x}^b + tx_j^{2a} \right)$$

$$\leq \exp \left( \sum_{k=1}^m \left( (-\lambda_{\min}(\mathbf{K}_0) + 2at\mathbf{1}_{k=j}) x_k^{2a} / (2a) + \eta_{0,k} x_k^b / b \right) \right),$$

a constant times the density for the model with parameters  $\text{diag}(\lambda_{\min}(\mathbf{K}_0) \mathbf{1}_m - 2ate_j)$  and  $\boldsymbol{\eta}_0$ . Thus,  $\mathbb{E}_0 \exp(tX_j^{2a})$  is finite for  $t \in (-\infty, \lambda_{\min}(\mathbf{K}_0)/(2a)) \ni 0$ .  $\square$

*Proof of Corollary 26.* Let the sub-exponential norm of  $X_j^{2a}$  be  $\|X_j^{2a}\|_{\psi_1} \equiv \sup_{q \geq 1} (\mathbb{E}_0 |X_j|^{2aq})^{1/q} / q$ , then by Lemma 35.6)

$$\mathbb{P}(|X_j^{2a} - \mathbb{E}_0 X_j^{2a}| \geq \epsilon_{3,j}) \leq \exp \left( - \min \left( \frac{\epsilon_3^2}{8e^2 \|X_j^{2a}\|_{\psi_1}^2}, \frac{\epsilon_3}{4e \|X_j^{2a}\|_{\psi_1}} \right) \right).$$

Letting

$$\epsilon_{3,j} \equiv \max \left\{ 2\sqrt{2}e \|X_j^{2a}\|_{\psi_1} \sqrt{\log 3 + \log n + \tau \log m + \log(m - |\rho_{\mathcal{D}}^*|)}, \right. \\ \left. 4e \|X_j^{2a}\|_{\psi_1} (\log 3 + \log n + \tau \log m + \log(m - |\rho_{\mathcal{D}}^*|)) \right\},$$

then  $\max\{\mathbb{E}_0 X_j^{2a} - \epsilon_{3,j}, 0\}^{1/(2a)} \leq X_j^{(i)} \leq (\mathbb{E}_0 X_j^{2a} + \epsilon_{3,j})^{1/(2a)}$  for all  $j \notin \rho_{\mathcal{D}}^*$  and  $i = 1, \dots, n$  with probability at least  $1 - 1/(3m^\tau)$ . The rest follows as in the proof of Theorem 24.  $\square$

*Proof of Corollary 27.* It suffices to bound  $\boldsymbol{\Gamma}$  and  $\mathbf{g}$  using their forms in Section 3.6.1 and apply Theorem 1 in Lin et al. (2016). We first bound  $(h_j \circ \varphi_j)(\mathbf{x}) x_j^{p_j} x_m^{p_m}$  with  $h_j(x) = x^{\alpha_j}$ ,  $\alpha_j \geq \max\{0, -p_j, -p_m, -p_j - p_m\}$ ,  $p_j, p_m \in \mathbb{R}$ , and  $0 < x_j + x_m < 1$ . By the definition of  $\varphi$  on simplices,  $\varphi_j(\mathbf{x}) = \min\{C_j, x_j, x_m\}$ , so  $(h_j \circ \varphi_j)(\mathbf{x}) x_j^{p_j} x_m^{p_m} = \min\{C_j, x_j, x_m\}^{\alpha_j} x_j^{p_j} x_m^{p_m}$  and is tightly lower bounded by 0. Noting that  $\min\{x_j, x_m\} < 1/2$ , we consider the following cases.

- (1) If  $C_j < \min\{x_j, x_m\}$ , then  $C_j < 1/2$  and the quantity is  $C_j^{\alpha_j} x_j^{p_j} x_m^{p_m} \leq C_j^{\alpha_j + (p_j)_- + (p_m)_-} < 2^{-\alpha_j - (p_j)_- - (p_m)_-}$  where  $(y)_- = y$  if  $y < 0$  and 0 otherwise.
- (2) Otherwise suppose  $x_j \leq x_m$  and  $x_j \leq C_j$ , then the quantity is equal to  $x_j^{\alpha_j + p_j} x_m^{p_m}$ , which if  $p_m \leq 0$  is upper bounded by  $x_j^{\alpha_j + p_j + p_m} < 2^{-\alpha_j - p_j - p_m} < 1$ ; if  $p_m > 0$  it is upper

bounded by  $((\alpha_j + p_j)/(\alpha_j + p_j + p_m))^{\alpha_j + p_j} (p_m/(\alpha_j + p_j + p_m))^{p_m}$  if  $(\alpha_j + p_j)/(\alpha_j + p_j + p_m) \leq 1/2$  or by  $2^{-\alpha_j - p_j - p_m}$  otherwise. Note that the statement for  $p_m > 0$  covers the one for  $p_m \leq 0$ . The conclusion for  $x_m \leq x_j$  and  $x_m \leq C_j$  follows by symmetry, and note that at most one of  $(\alpha_j + p_j)/(\alpha_j + p_j + p_m) \leq 1/2$  and  $(\alpha_j + p_m)/(\alpha_j + p_j + p_m) \leq 1/2$  can hold.

$$\text{In conclusion, defining } \zeta_2(\alpha_j, p_j, p_m) = \begin{cases} \left(\frac{\alpha_j + p_j}{\alpha_j + p_j + p_m}\right)^{\alpha_j + p_j} \left(\frac{p_m}{\alpha_j + p_j + p_m}\right)^{p_m}, & \text{if } p_m \geq \alpha_j + p_j, \\ \left(\frac{\alpha_j + p_m}{\alpha_j + p_j + p_m}\right)^{\alpha_j + p_m} \left(\frac{p_j}{\alpha_j + p_j + p_m}\right)^{p_j}, & \text{if } p_j \geq \alpha_j + p_m, \\ 2^{-\alpha_j - p_j - p_m}, & \text{otherwise,} \end{cases}$$

we have  $(h_j \circ \varphi_j)(\mathbf{x}) x_j^{p_j} x_m^{p_m} \leq \zeta_2(\alpha_j, p_j, p_m) < 1$ . Similarly  $\partial_j(h_j \circ \varphi_j)(\mathbf{x}) x_j^{p_j} x_m^{p_m} \leq \alpha_j \zeta_2(\alpha_j - 1, p_j, p_m) < \alpha_j$  if  $\alpha_j - 1 \geq \max\{0, -p_j, -p_m, -p_j - p_m\}$ .

Then for all  $j, k, \ell$ , as long as  $\alpha_j \geq \max\{1, 2 - a, 2 - b\}$ , we have  $0 \leq \gamma_{j,k,\ell} < 1$ , and similarly  $0 \leq g_{j,k} < \max_{j=1,\dots,m} \alpha_j + \max\{|a - 1| + 2a, |b - 1|\}$ . The rest follows from the same proof as Theorem 24.

Note that using the form of  $\mathbf{\Gamma}$  in Section 3.6.1, a tighter bound for  $\gamma_{j,k,\ell}$  is  $\max_{j,k=1,\dots,m} \max\{\zeta_2(\alpha_j, 2a - 2, 0), \zeta_2(\alpha_j, 4a - 2, 0), \zeta_2(\alpha_j, 2a - 2, 2a), \zeta_2(\alpha_j, 2b - 2, 0), \zeta_2(\alpha_j, 0, 2a - 2), \zeta_2(\alpha_j, 2a, 2a - 2), \zeta_2(\alpha_j, 0, 4a - 2), \zeta_2(\alpha_j, 0, 2b - 2), \zeta_2(\alpha_j, a - 1, a - 1), \zeta_2(\alpha_j, 2a - 1, 2a - 1), \zeta_2(\alpha_j, a - 1, 3a - 1), \zeta_2(\alpha_j, 3a - 1, a - 1), \zeta_2(\alpha_j, b - 1, b - 1)\}$ , and the one for  $g_{j,k}$  can be similarly written in terms of  $\zeta_2(\alpha_j, \cdot, \cdot)$  and  $\alpha_j \zeta_2(\alpha_j - 1, \cdot, \cdot)$ .  $\square$

*Proof of Lemma 28.* For  $a > 0$  or  $b > 0$ , the proof of Lemma 25 works even for the simplex domain. We thus only consider the case where  $a = b = 0$ . Similar to the proof there, we prove by showing that the moment-generating function of  $\log X_j$  is finite and invoking Theorem 2.13 of Wainwright (2019). According to Theorem 21, assume

- (1)  $\mathbf{K}_0$  is positive definite, or
- (2)  $\mathbf{K}_0 \mathbf{1}_m = \mathbf{0}$ ,  $\mathbf{K}_{0,-k,-k}$  is positive definite for some  $k = 1, \dots, m$ , and  $\mathbf{1}_m^\top \boldsymbol{\eta} + m \geq 0$ , or
- (3)  $\mathbf{K}_0 \mathbf{1}_m = \mathbf{0}$ ,  $\mathbf{K}_0$  is positive semi-definite, and  $\boldsymbol{\eta} \succ -\mathbf{1}_m$ .

For any suitable  $t$ ,  $\mathbb{E}_0 \exp(t \log X_j)$  is the inverse normalizing constant for the model with parameters  $\mathbf{K}_0$  and  $\boldsymbol{\eta}_0 + t\mathbf{e}_j$ , and is thus finite for (1) with  $t \in \mathbb{R}$  and (3) with  $t \in (-1 - \eta_{0,j}, +\infty) \ni 0$ . For (2), recall that in the proof of Theorem 21 we have shown that for any  $k = 1, \dots, m$ , the density of  $\log \mathbf{X}_{-k} - (\log X_k) \mathbf{1}_{m-1}$  is bounded by a constant times a Gaussian density, and thus  $\mathbb{E}_0 [X_j^t / X_k^t] = \mathbb{E}_0 \exp(t(\log X_j - \log X_k)) < +\infty$  for any  $k = 1, \dots, m$  and  $t \in \mathbb{R}$ . So for any  $t < 0$

$$\begin{aligned} \mathbb{E}_0 X_j^t &\leq \mathbb{E}_0 [X_j^t | X_j \geq 1/m] \mathbb{P}(X_j \geq 1/m) + \sum_{k \neq j} \mathbb{E}_0 [X_j^t | X_k \geq 1/m] \mathbb{P}(X_k \geq 1/m) \\ &\leq m^{-t} \mathbb{P}(X_j \geq 1/m) + \sum_{k \neq j} m^{-t} \mathbb{E}_0 [X_j^t / X_k^t | X_k \geq 1/m] \mathbb{P}(X_k \geq 1/m) \\ &\leq m^{-t} + \sum_{k \neq j} m^{-t} \mathbb{E}_0 [X_j^t / X_k^t] < +\infty. \end{aligned}$$

On the other hand,  $\mathbb{E}_0 X_j^t \leq 1$  for  $t \geq 0$ . Thus,  $\mathbb{E}_0 \exp(t \log X_j) < +\infty$  for any  $t \in \mathbb{R}$  for (2). Hence, for all of (1)–(3) we have  $\mathbb{E}_0 \exp(t \log X_j) < +\infty$  for  $t$  in a neighborhood around 0.  $\square$

*Proof of Corollary 29.* Let the sub-exponential norm of  $\log X_j$  be  $\|\log X_j\|_{\psi_1} \equiv \sup_{q \geq 1} (\mathbb{E}_0 |\log X_j|^q)^{1/q} / q$ , then by Lemma 35.6),

$$\mathbb{P}(-\log X_j + \mathbb{E}_0 \log X_j \geq \epsilon_3) \leq \exp \left( - \min \left( \frac{\epsilon_3^2}{8e^2 \|\log X_j\|_{\psi_1}^2}, \frac{\epsilon_3}{4e \|\log X_j\|_{\psi_1}} \right) \right).$$

Letting

$$\epsilon_3 \equiv \max \left\{ 2\sqrt{2}e \max_j \|\log X_j\|_{\psi_1} \sqrt{\log 3 + \log n + (\tau + 1) \log m}, \right. \\ \left. 4e \max_j \|\log X_j\|_{\psi_1} (\log 3 + \log n + (\tau + 1) \log m) \right\},$$

then  $0 \leq -\log X_j^{(i)} \leq \max_k \mathbb{E}_0 \log X_k + \epsilon_3$  for all  $j = 1, \dots, m$  and  $i = 1, \dots, n$  with probability at least  $1 - 1/(3m^\tau)$ . The rest follows as in the proof of Theorem 24 and Corollary 27.  $\square$

Appendix C

**APPENDIX TO CHAPTER 4**

In this appendix we present proofs for the theorems and corollaries in Chapter 4.

We first prove the following lemma that states that if two sums of distinct (ignoring the multiplicative constant) exponentials of polynomials in  $\mathbf{y} \in \mathbb{R}^m$  agree almost everywhere in  $\mathbb{R}^m$ , then they must have the same number of terms and there must be a 1-1 correspondence between the terms.

**Lemma 39.** *Let the number of variable be  $m \geq 1$  and the degree be  $p \geq 1$ . Let  $\mathcal{D} \equiv \{\mathbf{d} \in \mathbb{Z}_{\geq 0}^m : 1 \leq \sum_{j=1}^m d_j \leq p\}$  be the set of nonnegative integer-valued  $m$ -vectors with  $\ell_1$  norm  $\in [1, p]$ . Given a vector  $\mathbf{a} \in \mathbb{R}^{|\mathcal{D}|}$  indexed by  $\mathbf{d} \in \mathcal{D}$  (i.e.  $a_{\mathbf{d}} \in \mathbb{R}$  for all  $\mathbf{d} \in \mathcal{D}$ ), define*

$$f^{(m)}(\mathbf{y}; \mathbf{a}) \equiv \exp\left(\sum_{\mathbf{d} \in \mathcal{D}} a_{\mathbf{d}} \prod_{j=1}^m y_j^{d_j}\right),$$

*the exponential of the corresponding polynomial of degree  $\leq p$  in  $\mathbf{y} \in \mathbb{R}^m$ . Note that  $f^{(m)}$  does not have a constant term, and has degrees  $\mathbf{d} \in \mathcal{D}$  and coefficients  $\mathbf{a}$ .*

*Suppose we have*

$$\sum_{i=1}^{N_a} a_0^i f^{(m)}(\mathbf{y}; \mathbf{a}^i) = \sum_{i=1}^{N_b} b_0^i f^{(m)}(\mathbf{y}; \mathbf{b}^i) \quad (\text{C.1})$$

*for almost every  $\mathbf{y} \equiv (y_1, \dots, y_m) \in \mathbb{R}^m$  with respect to the Lebesgue measure, where  $N_a \geq 0$ ,  $N_b \geq 0$ ,  $\{\mathbf{a}^i\}_{i=1}^{N_a}$  are  $N_a$  distinct vectors in  $\mathbb{R}^{|\mathcal{D}|}$ ,  $\{\mathbf{b}^i\}_{i=1}^{N_b}$  are  $N_b$  distinct vectors in  $\mathbb{R}^{|\mathcal{D}|}$  (otherwise just combine the coefficients), and  $a_0^i, b_0^i \in \mathbb{R} \setminus \{0\}$  for all  $i$ . In other words, both sides of (C.1) are a sum of distinct exponentials of polynomials.*

*Then we must have  $N_a = N_b$  and there is a permutation  $\pi$  of  $\{1, \dots, N_a\}$  such that  $\mathbf{a}^i = \mathbf{b}^{\pi(i)}$  and  $a_0^i = b_0^{\pi(i)}$ , i.e. there is a 1-1 correspondence between the summands on both sides of (C.1).*

*Proof of Lemma 39.* First note that both sides of (C.1) are continuous functions, and so is their difference, which is 0 almost everywhere by assumption. Thus, the inverse image of the open set  $\mathbb{R} \setminus \{0\}$  under the difference is also open, and must be the empty set since it has measure 0. (C.1) thus holds for all  $\mathbf{y} \in \mathbb{R}^m$ .

We prove by induction on  $m$ , and first show the result for  $m = 1$ . In this case,  $f^{(1)}(y_1; \mathbf{a}) \equiv \exp(a_1 y_1 + \dots + a_p y_1^p)$ , and  $\mathbf{a}$  is just a  $p$ -vector.

First suppose  $N_a \neq 0$  and  $N_b \neq 0$ . Observe that as  $x \nearrow +\infty$ , if  $a_0 \neq 0$ , the function  $a_0 \exp(a_1 x + \cdots + a_p x^p)$  goes to

- (i)  $a_0 \neq 0$  if  $a_1 = \cdots = a_p = 0$ , or
- (ii) 0 if  $a_{d_{\max \neq 0}(\mathbf{a})} < 0$  where  $d_{\max \neq 0}(\mathbf{a})$  is the largest  $d \in \{1, \dots, p\}$  such that  $a_d \neq 0$ , or
- (iii)  $+\infty$  if  $a_{d_{\max \neq 0}(\mathbf{a})} > 0$ .

Rearrange the terms on the left of (C.1) so that for each  $1 \leq i < j \leq N_a$  we have  $(\mathbf{a}^i - \mathbf{a}^j)_{d_{\max \neq 0}(\mathbf{a}^i - \mathbf{a}^j)} > 0$ , and denote this total order as  $\mathbf{a}^i > \mathbf{a}^j$ . Rearrange the right-hand side similarly. By the assumption that  $\{\mathbf{a}^i\}_{i=1}^{N_a}$  are distinct,  $\mathbf{a}^i - \mathbf{a}^j \neq \mathbf{0}$ , so  $d_{\max \neq 0}(\mathbf{a}^i - \mathbf{a}^j)$  exists and this rearrangement is possible. Now dividing both sides of (C.1) by  $f^{(1)}(y_1; \mathbf{a}^1) = \exp(a_1^1 y_1 + \cdots + a_p^1 y_1^p)$  we have

$$a_0^1 + \sum_{i=2}^{N_a} a_0^i f^{(1)}(y_1; \mathbf{a}^i - \mathbf{a}^1) = \sum_{i=1}^{N_b} b_0^i f^{(1)}(y_1; \mathbf{b}^i - \mathbf{a}^1). \quad (\text{C.2})$$

Since  $a_0^1 \neq 0$ , and by the unique maximality of  $\mathbf{a}^1$ , as  $y_1 \nearrow +\infty$ , all terms in the summation on the left go to 0 (case (ii)). Thus, the right-hand side necessarily also goes to  $a_0^1 \neq 0$ , landing us in case (i) for at least one (and only one because  $\mathbf{b}^i$  are unique) term on the right, i.e.  $\mathbf{b}^i - \mathbf{a}^1 = \mathbf{0}$ . (A nonzero finite limit cannot come from a sum of terms that go to  $+\infty$  with positive and negative weights, since they must grow at different rates by uniqueness of  $\mathbf{b}^i - \mathbf{a}^1$ .) Since summands on both sides are sorted, we must have  $\mathbf{b}^1 = \mathbf{a}^1$ .

Then (C.2) becomes  $a_0^1 - b_0^1 + \sum_{i=2}^{N_a} a_0^i f^{(1)}(y_1; \mathbf{a}^i - \mathbf{a}^1) = \sum_{i=2}^{N_b} b_0^i f^{(1)}(y_1; \mathbf{b}^i - \mathbf{a}^1)$ . If  $a_0^1 \neq b_0^1$ , by the same reasoning there exists another  $i \in \{2, \dots, N_b\}$  such that  $\mathbf{b}^i - \mathbf{a}^1 = \mathbf{0}$ , violating uniqueness of  $\{\mathbf{b}^i\}_{i=1}^{N_b}$ . Thus,  $a_0^1 = b_0^1$  and  $\mathbf{a}^1 = \mathbf{b}^1$ , and we have reduced the number of summands on both sides of (C.2) by 1 to

$$\sum_{i=2}^{N_a} a_0^i f^{(1)}(y_1; \mathbf{a}^i - \mathbf{a}^1) = \sum_{i=2}^{N_b} b_0^i f^{(1)}(y_1; \mathbf{b}^i - \mathbf{a}^1).$$

Continuing this process by each time dividing both sides by  $f^{(1)}(y_1; \mathbf{a}^j - \mathbf{a}^{j-1})$ , we would have matched  $\min\{N_a, N_b\}$  pairs of coefficients between the  $a$  and the  $b$  groups. If  $N_a \neq N_b$ , assume  $N_a > N_b$  without loss of generality, then

$$\sum_{i=N_b+1}^{N_1} a_0^i f^{(1)}(y_1; \mathbf{a}^i - \mathbf{a}^{N_b}) = \text{const.}$$

Here the right-hand side is a constant that could be nonzero, because the argument for  $a_0^1 = b_0^1$  in our first elimination step does not apply here. Dividing both sides by  $f^{(1)}(y_1; \mathbf{a}^{N_b+1} - \mathbf{a}^{N_b})$ , we have  $a_0^{N_b+1} + \sum_{i=N_b+2}^{N_1} a_0^i f^{(1)}(y_1; \mathbf{a}^i - \mathbf{a}^{N_b+1}) = f^{(1)}(y_1; \mathbf{a}^{N_b} - \mathbf{a}^{N_b+1})$ . By maximality of  $\mathbf{a}^{N_b+1}$  among  $\mathbf{a}^{N_b+1}, \dots, \mathbf{a}^{N_a}$ , the left-hand side goes to  $a_0^{N_b+1} \neq 0$  as  $y_1 \nearrow +\infty$ , while since  $\mathbf{a}^{N_b} > \mathbf{a}^{N_b+1}$ , the right-hand side goes to  $+\infty$ , a contradiction. Thus,  $N_a = N_b$ ,  $a_0^i = b_0^i$  and  $\mathbf{a}^i = \mathbf{b}^i$  for  $i = 1, \dots, N_a$ , proving the  $m = 1$  case when  $N_a \neq 0$  and  $N_b \neq 0$ .

Now consider the case where one of  $N_a$  and  $N_b$  is 0; assume without loss of generality that  $N_b = 0$ , then by division by  $f^{(1)}(y_1; \mathbf{a}^1)$ , the right-hand side is constant 0, while the left-hand side goes to  $a_0^1 \neq 0$  unless  $N_a = 0$ , so  $N_a = N_b = 0$ .

Now suppose the result holds for some  $m - 1 \geq 1$ , and suppose either  $N_a \neq 0$  or  $N_b \neq 0$ , otherwise there is nothing to prove. We denote  $\mathbf{a}_1$  as the subvector of  $\mathbf{a}$  corresponding to  $\mathbf{d}$  with  $d_1 \geq 1$ , i.e.  $\{a_{\mathbf{d}}\}_{\mathbf{d} \in \mathcal{D}, d_1 \geq 1}$ , and  $\mathbf{a}_{-1}$  as that of  $\mathbf{a}$  with  $d_1 = 0$ . Separating out the terms involving  $y_1$ ,

$$\begin{aligned} f^{(m)}(\mathbf{y}; \mathbf{a}^i) &= \exp \left\{ \sum_{d=1}^p \left( \sum_{\mathbf{d} \in \mathcal{D}, d_1=d} a_{\mathbf{d}}^i \prod_{j=2}^m y_j^{d_j} \right) y_1^d \right\} \exp \left( \sum_{\mathbf{d} \in \mathcal{D}, d_1=0} a_{\mathbf{d}}^i \prod_{j=2}^m y_j^{d_j} \right) \\ &= f^{(1)}(y_1; \mathbf{a}_{1*}^i(\mathbf{y}_{-1})) f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{a}_{-1}^i), \end{aligned}$$

where  $\mathbf{a}_{1*}^i(\mathbf{y}_{-1}) : \mathbb{R}^{m-1} \rightarrow \mathbb{R}^p$  is a vector-valued function in  $\mathbf{y}_{-1}$ , with  $d$ -th coordinate a polynomial  $\sum_{\mathbf{d} \in \mathcal{D}, d_1=d} a_{\mathbf{d}}^i \prod_{j=2}^m y_j^{d_j}$ , and coefficients corresponding to  $\mathbf{a}_{1*}^i$ . Note that there is a one-to-one correspondence between such a function  $\mathbf{a}_{1*}^i$  and vector  $\mathbf{a}_{1*}^i$ . So we can rewrite (C.1) as

$$\sum_{i=1}^{N_a} a_0^i f^{(1)}(y_1; \mathbf{a}_{1*}^i(\mathbf{y}_{-1})) f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{a}_{-1}^i) = \sum_{i=1}^{N_b} b_0^i f^{(1)}(y_1; \mathbf{b}_{1*}^i(\mathbf{y}_{-1})) f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{b}_{-1}^i)$$

for all  $\mathbf{y} \in \mathbb{R}^m$ . Then collecting terms with the same  $f^{(1)}$  (same  $\mathbf{a}_1^i$  ( $\mathbf{a}_{1*}^i$ ) or  $\mathbf{b}_1^i$  ( $\mathbf{b}_{1*}^i$ )),

$$\sum_{\ell=1}^C f^{(1)}(y_1; \mathbf{c}_{1*}^\ell(\mathbf{y}_{-1})) \left\{ \sum_{j=1}^{n_\ell^a} a_0^{k_{\ell j}^a} f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{a}_{-1}^{k_{\ell j}^a}) + \sum_{j=1}^{n_\ell^b} b_0^{k_{\ell j}^b} f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{b}_{-1}^{k_{\ell j}^b}) \right\} = 0, \quad (\text{C.3})$$

where  $C > 0$ , each  $\mathbf{c}_1^\ell$  (coefficients for  $\mathbf{c}_{1*}^\ell$ ) is some  $\mathbf{a}_1^i$  or  $\mathbf{b}_1^i$ , and  $\{\mathbf{c}_1^\ell\}_{\ell=1}^C$  are distinct. Here, let  $\{k_{11}^a, \dots, k_{1,n_1}^a, \dots, k_{C1}^a, \dots, k_{C,n_C}^a\}$  be a permutation of  $\{1, \dots, N_a\}$ , and  $\{k_{11}^b, \dots, k_{1,n_1}^b, \dots, k_{C1}^b, \dots, k_{C,n_C}^b\}$  a permutation of  $\{1, \dots, N_b\}$ .

Since  $\{\mathbf{c}_1^\ell\}_{\ell=1}^C$  are distinct,  $\{\mathbf{c}_{1*}^\ell\}_{\ell=1}^C$  are distinct finite polynomials in  $\mathbf{y}_{-1} \in \mathbb{R}^{m-1}$ . For each pair of such distinct polynomials, the lemma of Okamoto (1973) implies that they only agree at a Lebesgue-null subset of  $\mathbb{R}^{n-1}$ , so all polynomials are distinct except on a null set. Thus, for almost every fixed  $\mathbf{y}_{-1} \in \mathbb{R}^{m-1}$ , the left-hand side of (C.3) is a sum of  $C > 0$  distinct  $f^{(1)}$ 's in  $y_1$  multiplied by constant weights depending on  $\mathbf{y}_{-1}$ . But the right-hand side is a sum of 0 terms, so by the result for  $m = 1$  we necessarily have

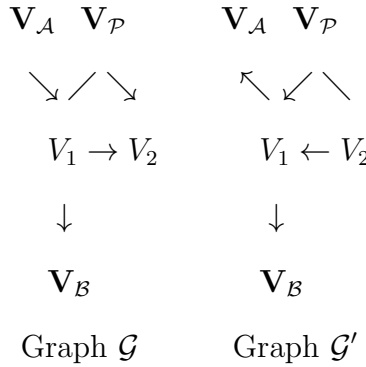
$$\sum_{j=1}^{n_\ell^a} a_0^{k_{\ell j}^a} f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{a}_{-1}^{k_{\ell j}^a}) = \sum_{j=1}^{n_\ell^b} -b_0^{k_{\ell j}^b} f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{b}_{-1}^{k_{\ell j}^b}) \quad (\text{C.4})$$

for all  $\ell = 1, \dots, C$  for almost every  $\mathbf{y}_{-1}$ . Fixing  $\ell \in \{1, \dots, C\}$ , for any  $1 \leq j_1 < j_2 \leq n_\ell^a$ ,  $\mathbf{a}_{-1}^{k_{\ell j_1}^a} \neq \mathbf{a}_{-1}^{k_{\ell j_2}^a}$  and  $\mathbf{a}_{-1}^{k_{\ell j_1}^a} = \mathbf{a}_{-1}^{k_{\ell j_2}^a}$  implies  $\mathbf{a}_{-1}^{k_{\ell j_1}^a} \neq \mathbf{a}_{-1}^{k_{\ell j_2}^a}$ , and similarly for  $\mathbf{b}$ . Thus, each term on the left-hand side of (C.4) has its unique coefficients, and similarly for the right-hand side. Since (C.4) holds for almost every  $\mathbf{y}_{-1}$ , by the result for  $m - 1$  variables, we must have  $n_\ell^a = n_\ell^b$  and each  $a_0^{k_{\ell j}^a} = b_0^{k_{\ell \pi(j)}^b}$  and  $\mathbf{a}_{-1}^{k_{\ell j}^a} = \mathbf{b}_{-1}^{k_{\ell \pi(j)}^b}$  for some permutation  $\pi$  of  $\{1, \dots, n_\ell^a\}$ , which in turn implies  $\mathbf{a}_{-1}^{k_{\ell j}^a} = \mathbf{b}_{-1}^{k_{\ell \pi(j)}^b}$  for all  $j = 1, \dots, n_\ell^a$  by construction of the groups  $\ell = 1, \dots, C$ . Since this holds for all  $\ell$ ,  $N_a = \sum_{\ell=1}^C n_\ell^a = \sum_{\ell=1}^C n_\ell^b = N_b$ , and we have thus again matched each  $\mathbf{a}^\ell$  with a  $\mathbf{b}^\ell$  as well as the corresponding  $a_0$ 's with  $b_0$ 's. This ends the proof for  $m$ , and the entire proof.  $\square$

*Proof of Theorem 30.* Suppose  $\mathcal{G}$  and  $\mathcal{G}'$  have the same node set  $\mathcal{V}$  and are Markov equivalent, otherwise the distributions represented by them are trivially not identical.

Now suppose  $p(\mathbf{Y})$  is Markov and faithful with respect to  $\mathcal{G}$  and  $\mathcal{G}'$ , and factorize w.r.t. both graphs with *strong Hurdle polynomial* parameters. Then by Proposition 31,

there exist  $V_1$  and  $V_2$  such that  $V_1 \rightarrow V_2$  in  $\mathcal{G}$ ,  $V_2 \rightarrow V_1$  in  $\mathcal{G}'$  and  $\mathcal{P} \equiv \text{pa}_{\mathcal{G}}(V_2) \setminus \{V_1\} = \text{pa}_{\mathcal{G}'}(V_1) \setminus \{V_2\}$ . Following the arguments in the proof of Proposition 31 in Peters et al. (2014), recursively marginalizing out nodes without children but having the same parents in both graphs, we eventually obtain structures as follows, where  $\mathcal{A}$  and  $\mathcal{B}$  are some unknown node sets and  $V_2$  does not have any children in Graph  $\mathcal{G}$ .



We consider the  $(\alpha, \beta, k)$ -parametrization only, since the result for the  $(p, \mu, \sigma^2)$  naturally follows from their relationship (4.6). For notational simplicity write  $V_1$  and  $V_2$  as nodes 1 and 2. Suppose after marginalization above we are left with nodes  $\mathcal{V}_0 \subseteq \mathcal{V}$  which include 1, 2,  $\mathbf{V}_{\mathcal{A}}$ ,  $\mathbf{V}_{\mathcal{B}}$  and  $\mathbf{V}_{\mathcal{P}}$  illustrated above. Now let  $Y_U = 0$  for all  $U \in \mathcal{V}_0 \setminus \{2\}$ , and let  $Y_2 \neq 0$ . Then the joint distribution  $p(Y_2 = y_2 \neq 0, \mathbf{y}_{\mathcal{V}_0} = \mathbf{0})$  using  $\mathcal{G}$  is proportional to

$$\prod_{V \in \mathcal{V}_0} \frac{\exp\{\alpha_V(\mathbf{y}_{\text{pa}_{\mathcal{G}}(V)})\mathbb{1}_{y_V} + \beta_V(\mathbf{y}_{\text{pa}_{\mathcal{G}}(V)})y_V - k_V y_V^2/2\}}{\sqrt{2\pi/k_V} \exp\{\alpha_V(\mathbf{y}_{\text{pa}_{\mathcal{G}}(V)}) + \beta_V(\mathbf{y}_{\text{pa}_{\mathcal{G}}(V)})^2/(2k_V)\} + 1} \Bigg|_{y_2 \neq 0, \mathbf{y}_{\mathcal{V}_0 \setminus \{2\}} = \mathbf{0}}$$

$$\propto \exp\{\beta_2(\mathbf{0})y_2 - k_2 y_2^2/2\}$$

since 2 does not have any child in  $\mathcal{G}$ . But using  $\mathcal{G}'$ , the same joint distribution is proportional to

$$\prod_{V \in \mathcal{V}_0} \frac{\exp\{\alpha'_V(\mathbf{y}_{\text{pa}_{\mathcal{G}'}(V)})\mathbb{1}_{y_V} + \beta'_V(\mathbf{y}_{\text{pa}_{\mathcal{G}'}(V)})y_V - k'_V y_V^2/2\}}{\sqrt{2\pi/k'_V} \exp\{\alpha'_V(\mathbf{y}_{\text{pa}_{\mathcal{G}'}(V)}) + \beta'_V(\mathbf{y}_{\text{pa}_{\mathcal{G}'}(V)})^2/(2k'_V)\} + 1} \Bigg|_{y_2 \neq 0, \mathbf{y}_{\mathcal{V}_0 \setminus \{2\}} = \mathbf{0}}$$

$$\propto \exp\{\beta'_2(\mathbf{0})y_2 - k'_2 y_2^2/2\} \times$$

$$\prod_{U \in \mathcal{P} \cup \{1\}, 2 \in \text{pa}_{\mathcal{G}'}(U)} \frac{1}{\sqrt{2\pi/k'_U} \exp\{\alpha'_U(y_2, \mathbf{0}) + \beta'_U(y_2, \mathbf{0})^2/(2k'_U)\} + 1}$$

where in the case where  $\text{pa}_{\mathcal{G}'}(2) = \emptyset$  replace  $\alpha'_2(\mathbf{0})$  and  $\beta'_2(\mathbf{0})$  by constants  $\alpha'_2$  and  $\beta'_2$ , and  $\alpha'_U(y_2, \mathbf{0})$  and  $\beta'_U(y_2, \mathbf{0})$  denote setting all parents other than 2 in the Hurdle polynomials  $\alpha'_U$  and  $\beta'_U$  to  $\mathbf{0}$ . Since the two joint distributions derived from both graphs must be proportional to each other, we get for  $y_2 \neq 0$

$$\begin{aligned} & \exp [y_2\{\beta'_2(\mathbf{0}) - \beta_2(\mathbf{0})\} - (k'_2 - k_2)y_2^2/2] \\ & \propto \prod_{U \in \mathcal{P} \cup \{1\}, 2 \in \text{pa}_{\mathcal{G}'}(U)} \left[ \sqrt{2\pi/k'_U} \exp \{ \alpha'_U(y_2, \mathbf{0}) + \beta'_U(y_2, \mathbf{0})^2/(2k'_U) \} + 1 \right]. \end{aligned} \tag{C.5}$$

Note that  $2 \in \text{pa}_{\mathcal{G}'}(1)$  and thus the product on the right of (C.5) has at least one term. Thus, supposing that for at least one of  $U \in \mathcal{P} \cup \{1\}$  such that  $2 \in \text{pa}_{\mathcal{G}'}(U)$ ,  $\alpha'_U(Y_2, \mathbf{0}) + \beta'_U(Y_2, \mathbf{0})^2/(2k'_U)$  is nonconstant in  $Y_2 \neq 0$ , then the right-hand side of (C.5) can be expanded into a sum of at least two exponentials of polynomials in  $y_2$  (including the constant 1 as a degenerated exponential polynomial), while the left-hand side is a single polynomial in  $y_2$ . This is a contradiction according to Lemma 39, and thus the assumption of having *strong Hurdle polynomials* as the parameters in the Hurdle conditionals implies that  $p(\mathbf{Y})$  cannot be represented by both  $\mathcal{G}$  and  $\mathcal{G}'$ , which ends the proof. □

*Proof of Theorem 32.* As in the proof of Theorem 30 using Proposition 31, under the assumptions there exist  $V_1$  and  $V_2$  such that  $\mathcal{P} \equiv \text{pa}_{\mathcal{G}}(V_2) \setminus \{V_1\} = \text{pa}_{\mathcal{G}'}(V_1) \setminus \{V_2\}$  with  $V_1 \rightarrow V_2$  in  $\mathcal{G}$  and  $V_2 \rightarrow V_1$  in  $\mathcal{G}'$ . Following the arguments in the proof of Proposition 31 in Peters et al. (2014), recursively marginalizing out nodes without children but having the same parents in both graphs, we again obtain structures as follows.



$$\begin{array}{ccc}
V_1 \rightarrow V_2 & & V_1 \leftarrow V_2 \\
\downarrow & & \downarrow \\
\mathbf{V}_{\mathcal{B}} & & \mathbf{V}_{\mathcal{B}} \\
\text{Graph } \mathcal{G} & & \text{Graph } \mathcal{G}'
\end{array}$$

To ease the notation assume we again write  $V_1 = 1$  and  $V_2 = 2$ . Note that the distribution of each node conditional on some other nodes is the sum of a point mass at 0 and a continuous distribution over  $\mathbb{R}$ , which follows by induction and the fact that the indefinite integral of a continuous density is continuous and that the sum of continuous densities is continuous. We focus on the continuous components, and wish to reach the conclusion using the factorization

$$\begin{aligned}
P(y_1, y_2 | \mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}}) &= P(y_1 | \mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}}) P(y_2 | y_1, \mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}}) \\
&= P(y_2 | \mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}}) P(y_1 | y_2, \mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}}),
\end{aligned}$$

where the second terms in both decompositions are a regular Hurdle conditional w.r.t.  $\mathcal{G}$  and  $\mathcal{G}'$ , respectively, and we write the first terms as

$$P(y_1 | \mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}}) \propto \exp\{\mathbb{1}_{y_1} \delta_1 + f_1(y_1)\}$$

and

$$P(y_2 | \mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}}) \propto \exp\{\mathbb{1}_{y_2} \delta'_2 + f'_2(y_2)\}$$

in terms of the conditional densities w.r.t.  $\lambda$ . Here  $f_1$  and  $f'_2$  are continuous functions in  $\mathbb{R}$  with no additive constant term, and  $\delta_1$  and  $\delta'_2$  are constants.

We prove the results in the  $(\alpha, \beta, k)$ -parameterization only, since results for the  $(p, \mu, \sigma^2)$ -parameterization would follow from their relationship (4.6). In our model, we assumed the  $\alpha$  and  $\beta$  parameters for each node to be polynomial in the parents and their indicators. We also assumed that for each node, either the  $\beta$  function is nonconstant in any of the parents, or  $\alpha$  depends on the value of all of its parents.

Consider a generic  $\beta$  function associated with some generic parent set  $\mathcal{P} \equiv \mathcal{P}_1 \sqcup \{p_0\}$  with  $p_0 \notin \mathcal{P}_1 \neq \emptyset$  and suppose that  $\beta$  is nonconstant in any of  $\mathcal{P}$ , and write  $\beta(\mathbf{y}_{\mathcal{P}})$  equivalently

as  $\beta(y_{p_0}, \mathbf{y}_{\mathcal{P}_1})$ . Then  $\beta(\mathbf{y}_{\mathcal{P}})$  has the form  $\beta_{-1}(\mathbf{y}_{\mathcal{P}_1}) + \beta_0(\mathbf{y}_{\mathcal{P}_1})\mathbb{1}_{y_1} + \sum_{i=1}^k \beta_i(\mathbf{y}_{\mathcal{P}_1})y_1^i$ , where by construction  $\beta_{-1}$  through  $\beta_k$  are (potentially constant or even zero) Hurdle polynomials in  $\mathbf{y}_{\mathcal{P}_1}$ , but there must exist some  $j = 0, \dots, k$  such that  $\beta_j$  is nonzero. By the lemma of Okamoto (1973),  $\beta_j(\mathbf{y}_{\mathcal{P}_1}) \neq 0$  for (Lebesgue) almost every  $\mathbf{y}_{\mathcal{P}_1} \in \mathbb{R}^{|\mathcal{P}_1|}$ . Thus,  $\beta(y_{p_0}, \mathbf{y}_{\mathcal{P}_1})$  is nonconstant in  $y_{p_0}$  for almost every  $\mathbf{y}_{\mathcal{P}_1} \in \mathbb{R}^{|\mathcal{P}_1|}$ . Formally, define

$$\mathcal{Y}_{\beta, p_0, \mathcal{P}_1} \equiv \{\mathbf{y}_{\mathcal{P}_1} \in \mathbb{R}^{|\mathcal{P}_1|} : \beta(y_{p_0}, \mathbf{y}_{\mathcal{P}_1}) \text{ nonconstant function in } y_{p_0}\}.$$

Thus  $\mathbb{R}^{|\mathcal{P}_1|} \setminus \mathcal{Y}_{\beta, p_0, \mathcal{P}_1}$  has zero Lebesgue measure assuming  $\beta$  is nonconstant in its any of  $\mathcal{P}$ . Hence, by a similar argument, under the assumptions of the theorem, letting

$$\begin{aligned} \mathcal{Y}_{\alpha, \beta, p_0, \mathcal{P}_1} \equiv \{\mathbf{y}_{\mathcal{P}_1} \in \mathbb{R}^{|\mathcal{P}_1|} : \beta(y_{p_0}, \mathbf{y}_{\mathcal{P}_1}) \text{ nonconstant function in } y_{p_0} \text{ or} \\ \alpha(y_{p_0}, \mathbf{y}_{\mathcal{P}_1}) \text{ depends on the value of } y_{p_0}\}, \end{aligned}$$

$\mathbb{R}^{|\mathcal{P}_1|} \setminus \mathcal{Y}_{\alpha, \beta, p_0, \mathcal{P}_1}$  has zero Lebesgue measure.

Now we go back to  $\mathcal{G}$  and  $\mathcal{G}'$ . Suppose  $\mathcal{P} \neq \emptyset$  and that the Hurdle density of node 2 conditional on  $\{1\} \sqcup \mathcal{P}$  in  $\mathcal{G}$  have  $\alpha$  and  $\beta$  parameters  $\alpha_2(y_1, \mathbf{y}_{\mathcal{P}})$  and  $\beta_2(y_1, \mathbf{y}_{\mathcal{P}})$ , and let those for 1 conditional on  $\{2\} \sqcup \mathcal{P}$  in  $\mathcal{G}'$  be  $\alpha'_1(y_2, \mathbf{y}_{\mathcal{P}})$  and  $\beta'_1(y_2, \mathbf{y}_{\mathcal{P}})$ . We also denote  $\mathcal{Y}_* \equiv \mathcal{Y}_{\alpha_2, \beta_2, 1, \mathcal{P}} \cap \mathcal{Y}_{\alpha'_1, \beta'_1, 2, \mathcal{P}}$ , which by discussion above contains almost every  $\mathbf{y}_{\mathcal{P}} \subset \mathbb{R}^{|\mathcal{P}|}$ .

From now on we thus fix  $\mathbf{y}_{\mathcal{P}} \in \mathcal{Y}_*$  and condition on  $\mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}}$ , and omit the dependency of the  $\alpha$  and  $\beta$  functions on  $\mathcal{P}$ , and write them as scalar functions instead notation-wise. By discussion above,  $\beta_2$  becomes a nonconstant function in  $y_1$  and  $\beta'_1$  becomes a nonconstant function in  $y_2$ . Note that for  $\mathcal{P} = \emptyset$ , we do not fix or condition on any parent variables and  $\alpha'_1$ ,  $\alpha_2$ ,  $\beta'_1$  and  $\beta_2$  are automatically univariate functions, with  $\beta'_1$  and  $\beta_2$  nonconstant by assumption.

The joint density of  $P(y_1, y_2 | \mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}})$  w.r.t.  $\lambda$  thus has two characterizations (up to normalizing constants)

$$\frac{\exp\{\mathbb{1}_{y_1} \delta_1 + f_1(y_1) + \mathbb{1}_{y_2} \alpha_2(y_1) + y_2 \beta_2(y_1) - y_2^2 k_2 / 2\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2 / (2k_2)\} + 1}$$

$$\propto \frac{\exp\{\mathbb{1}_{y_2}\delta'_2 + f'_2(y_2) + \mathbb{1}_{y_1}\alpha'_1(y_2) + y_1\beta'_1(y_2) - y_1^2k'_1/2\}}{\sqrt{2\pi/k'_1} \exp[\alpha'_1(y_2) + \{\beta'_1(y_2)\}^2/(2k'_1)] + 1}, \quad (\text{C.6})$$

where  $\alpha_2(y_1)$  has the form  $c_{\alpha_2,-1} + c_{\alpha_2,0}\mathbb{1}_{y_1} + c_{\alpha_2,1}y_1 + \cdots + c_{\alpha_2,k}y_1^k$  with coefficients being polynomials in  $\mathbf{y}_{\mathcal{P}}$  and their indicators (or constants if  $\mathcal{P} = \emptyset$ ), and similarly for  $\beta_2(y_1)$ ,  $\alpha'_1(y_2)$  and  $\beta'_1(y_2)$ . Note that if the values of  $\mathbb{1}_{y_1}$  and  $\mathbb{1}_{y_2}$  are given, these four functions are just polynomials in  $y_1$  and  $y_2$ , respectively.

First condition on the event  $\mathbb{1}_{y_1} = \mathbb{1}_{y_2} = 1$  that has a positive probability. Then (C.6) becomes

$$\begin{aligned} & \frac{\exp\{f_1(y_1) + \alpha_2(y_1) + y_2\beta_2(y_1) - y_2^2k_2/2\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \mathbb{1}_{y_1} \mathbb{1}_{y_2}, \\ & \propto \frac{\exp\{f'_2(y_2) + \alpha'_1(y_2) + y_1\beta'_1(y_2) - y_1^2k'_1/2\}}{\sqrt{2\pi/k'_1} \exp\{\alpha'_1(y_2) + (\beta'_1(y_2))^2/(2k'_1)\} + 1} \mathbb{1}_{y_1} \mathbb{1}_{y_2}, \end{aligned} \quad (\text{C.7})$$

for all  $(y_1, y_2) \in (\mathbb{R} \setminus \{0\})^2$ . (C.7) has the form

$$\frac{\exp\{f_1(y_1) + P_1(y_1, y_2)\}}{\exp\{P_2(y_1)\} + 1} = \frac{\exp\{f'_2(y_2) + P_3(y_1, y_2)\}}{\exp\{P_4(y_2)\} + 1},$$

where  $P_1$  and  $P_3$  are polynomials in  $y_1$  and  $y_2$  simultaneously, possibly with interactions from the  $y_2\beta_2(y_1)$  and  $y_1\beta'_1(y_2)$  terms, and  $P_2$  and  $P_4$  are univariate polynomials in  $y_1$ ,  $y_2$ , respectively. By cross-multiplication,

$$\begin{aligned} & \exp\{f_1(y_1) + P_1(y_1, y_2) + P_4(y_2)\} + \exp\{f_1(y_1) + P_1(y_1, y_2)\} \\ & = \exp\{f'_2(y_2) + P_3(y_1, y_2) + P_2(y_1)\} + \exp\{f'_2(y_2) + P_3(y_1, y_2)\}. \end{aligned} \quad (\text{C.8})$$

Differentiating both sides of (C.8) with respect to  $y_1$ ,

$$\begin{aligned} & \left[ \frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2)\} \right] \exp \{f_1(y_1) + P_1(y_1, y_2) + P_4(y_2)\} \\ & \quad + \exp\{f_1(y_1) + P_1(y_1, y_2)\} \\ & = \left[ \frac{\partial}{\partial y_1} \{P_3(y_1, y_2) + P_2(y_1)\} \right] \exp \{f'_2(y_2) + P_3(y_1, y_2) + P_2(y_1)\} \\ & \quad + \left\{ \frac{\partial}{\partial y_1} P_3(y_1, y_2) \right\} \exp \{f'_2(y_2) + P_3(y_1, y_2)\}. \end{aligned} \quad (\text{C.9})$$

Plugging (C.8) into the left-hand side of (C.9),

$$\begin{aligned}
& \left[ \frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2)\} \right] [\exp \{f'_2(y_2) + P_3(y_1, y_2) + P_2(y_1)\} \\
& \qquad \qquad \qquad + \exp \{f'_2(y_2) + P_3(y_1, y_2)\}] \\
& = \left[ \frac{\partial}{\partial y_1} \{P_3(y_1, y_2) + P_2(y_1)\} \right] \exp \{f'_2(y_2) + P_3(y_1, y_2) + P_2(y_1)\} \\
& \qquad \qquad \qquad + \left\{ \frac{\partial}{\partial y_1} P_3(y_1, y_2) \right\} \exp \{f'_2(y_2) + P_3(y_1, y_2)\},
\end{aligned}$$

which simplifies to

$$\begin{aligned}
& \left[ \frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2) - P_3(y_1, y_2) - P_2(y_1)\} \right] \\
& \qquad \qquad \qquad \times \exp \{f'_2(y_2) + P_3(y_1, y_2) + P_2(y_1)\} \\
& \qquad \qquad \qquad + \left[ \frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2) - P_3(y_1, y_2)\} \right] \\
& \qquad \qquad \qquad \qquad \qquad \qquad \times \exp \{f'_2(y_2) + P_3(y_1, y_2)\} = 0.
\end{aligned}$$

Since  $\exp \{f'_2(y_2) + P_3(y_1, y_2)\} \neq 0$ , this becomes

$$\begin{aligned}
& \left[ \frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2) - P_3(y_1, y_2) - P_2(y_1)\} \right] \exp \{P_2(y_1)\} \\
& \qquad \qquad \qquad + \left[ \frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2) - P_3(y_1, y_2)\} \right] = 0. \quad (\text{C.10})
\end{aligned}$$

Focusing on the components that involve  $y_2$ , we see that

$$\left[ \frac{\partial}{\partial y_1} \{P_1(y_1, y_2) - P_3(y_1, y_2)\} \right] [\exp \{P_2(y_1)\} + 1]$$

does not depend on  $y_2$ . Since  $(\exp(P_2(y_1)) + 1) > 0$ , we have

$$\frac{\partial^2}{\partial y_1 \partial y_2} \{P_1(y_1, y_2) - P_3(y_1, y_2)\} = 0.$$

Recall that

$$P_1(y_1, y_2) - P_3(y_1, y_2) = \alpha_2(y_1) + y_2 \beta_2(y_1) - y_2^2 k_2 / 2 - \alpha'_1(y_2) - y_1 \beta'_1(y_2) + y_1^2 k'_1 / 2. \quad (\text{C.11})$$

So  $0 = \frac{\partial^2}{\partial y_1 \partial y_2} \{P_1(y_1, y_2) - P_3(y_1, y_2)\} = \frac{d\beta_2(y_1)}{dy_1} - \frac{d\beta'_1(y_2)}{dy_2}$  implies that  $\beta_2$  and  $\beta'_1$  are both linear with the same coefficient on the linear term. Now that  $\beta_2$  has the form  $\beta_2(y_1) = c_{\beta_2, -1} + c_{\beta_2, 0}\mathbb{1}_{y_1} + c_{\beta_2, 1}y_1$ , write  $\beta_{2; -1, 0} \equiv c_{\beta_2, -1} + c_{\beta_2, 0} = \beta_2(0) + c_{\beta_2, 0}$  as a shorthand notation for  $\beta_2$  with indicator set to 1 while  $y_1$  set to 0. Similarly define  $\beta'_{1; -1, 0} \equiv c_{\beta'_1, -1} + c_{\beta'_1, 0} = \beta'_1(0) + c_{\beta'_1, 0}$ . Then for  $y_1, y_2 \neq 0$  since  $c_{\beta_2, 1} = c_{\beta'_1, 1}$ , we necessarily have

$$\begin{aligned} y_2\beta_2(y_1) - y_1\beta'_1(y_2) &= y_2(c_{\beta_2, -1} + c_{\beta_2, 0} + c_{\beta_2, 1}y_1) - y_1(c_{\beta'_1, -1} + c_{\beta'_1, 0} + c_{\beta'_1, 1}y_2) \\ &= y_2\beta_{2; -1, 0} - y_1\beta'_{1; -1, 0}, \end{aligned}$$

and so by (C.11)

$$\begin{aligned} &P_1(y_1, y_2) - P_3(y_1, y_2) \\ &= (\alpha_2(y_1) - y_1\beta'_{1; -1, 0} + y_1^2k'_1/2) - (\alpha'_1(y_2) - y_2\beta_{2; -1, 0} + y_2^2k_2/2) \\ &\equiv P_{1,3}(y_1) - (\text{function in } y_2 \text{ only}). \end{aligned}$$

Plugging this into (C.10), we get

$$\left[ \frac{d}{dy_1} \{f_1(y_1) + P_{1,3}(y_1) - P_2(y_1)\} \right] \exp\{P_2(y_1)\} + \left[ \frac{d}{dy_1} \{f_1(y_1) + P_{1,3}(y_1)\} \right]$$

equals 0, or equivalently

$$\left[ \frac{d}{dy_1} \{f_1(y_1) + P_{1,3}(y_1)\} \right] [\exp\{P_2(y_1)\} + 1] = \left\{ \frac{d}{dy_1} P_2(y_1) \right\} \exp\{P_2(y_1)\}.$$

Then

$$\begin{aligned} f_1(y_1) &= \int \frac{\exp\{P_2(y_1)\} \{dP_2(y_1)/dy_1\}}{\exp\{P_2(y_1)\} + 1} dy_1 - P_{1,3}(y_1) \\ &= \log[1 + \exp\{P_2(y_1)\}] - P_{1,3}(y_1) + \text{const.} \end{aligned}$$

So for  $y_1 \neq 0$ ,

$$\begin{aligned} \exp(f_1(y_1)) &\propto \frac{1 + \exp\{P_2(y_1)\}}{\exp\{P_{1,3}(y_1)\}} \\ &= \frac{1 + \sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\}}{\exp\{\alpha_2(y_1) - \beta'_{1; -1, 0}y_1 + y_1^2k'_1/2\}} \end{aligned}$$

$$\begin{aligned}
&= \exp\{-\alpha_2(y_1) + y_1\beta'_{1,-1,0} - y_1^2k'_1/2\} \\
&\quad + \sqrt{2\pi/k_2} \exp\{y_1\beta'_{1,-1,0} + \beta_2(y_1)^2/(2k_2) - y_1^2k'_1/2\}. \tag{C.12}
\end{aligned}$$

Now condition on the event  $\mathbf{1}_{y_1} = 1$  and  $\mathbf{1}_{y_2} = 0$ . Then (C.6) becomes

$$\frac{\exp\{f_1(y_1)\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \mathbf{1}_{y_1} \propto \exp\{y_1\beta'_1(0) - y_1^2k'_1/2\} \mathbf{1}_{y_1},$$

which implies that for  $y_1 \neq 0$ ,

$$\begin{aligned}
\exp\{f_1(y_1)\} &\propto \exp\{y_1\beta'_1(0) - y_1^2k'_1/2\} \\
&\quad + \sqrt{2\pi/k_2} \exp\{y_1\beta'_1(0) - y_1^2k'_1/2 + \alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\}. \tag{C.13}
\end{aligned}$$

Applying Lemma 39 to (C.12) and (C.13), by matching the terms we have (conditional on  $y_1 \neq 0$ ) either

$$-\alpha_2(y_1) + y_1\beta'_{1,-1,0} = y_1\beta'_1(0) + \text{const}; \quad \text{or} \tag{C.14}$$

$$-\alpha_2(y_1) + y_1\beta'_{1,-1,0} = y_1\beta'_1(0) + \alpha_2(y_1) + \beta_2(y_1)^2/(2k_2) + \text{const} \quad \text{and}$$

$$y_1\beta'_{1,-1,0} + \beta_2(y_1)^2/(2k_2) = y_1\beta'_1(0) + \text{const}. \tag{C.15}$$

Conditional on  $y_1 \neq 0$ , in the first case (C.14),  $\alpha_2(y_1) = y_1c_{\beta'_1,0} + \text{const}$ ; in the second case (C.15),  $\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2) = \text{const}$  and  $\beta_2(y_1)^2/(2k_2) = -y_1c_{\beta'_1,0} + \text{const}$ , which implies  $\beta_2(y_1) = \text{const}$  and  $\alpha_2(y_1) = \text{const}$  for  $y_1 \neq 0$ , and  $c_{\beta'_1,0} = 0$ , which in turn implies (C.14). Thus, in either case,  $\alpha_2(y_1) = c_{\alpha_2,0}\mathbf{1}_{y_1} + y_1c_{\beta'_1,0} + \text{const}$ , i.e.  $\alpha_2$  is linear (or constant) in  $y_1 \neq 0$  with coefficient on  $y_1$  equal to  $c_{\beta'_1,0}$ . By (C.14) for  $y_1 \neq 0$ ,

$$\begin{aligned}
\exp\{f_1(y_1)\} &\propto \exp\{y_1\beta'_1(0) - y_1^2k'_1/2\} \\
&\quad + \sqrt{2\pi/k_2} \exp\{y_1\beta'_{1,-1,0} + \beta_2(y_1)^2/(2k_2) - y_1^2k'_1/2\}, \tag{C.16}
\end{aligned}$$

clearly a single univariate Gaussian or a mixture of two univariate Gaussian distributions (since  $\beta_2$  is at most linear in  $y_1$ ). Similarly, we must have  $\alpha'_1(y_2) = y_2\beta_{2,-1,0} - y_2\beta_2(0) + \text{const} = y_2c_{\beta_2,0} + \text{const}$  for  $y_2 \neq 0$ , and for  $y_2 \neq 0$

$$\begin{aligned} \exp\{f'_2(y_2)\} &\propto \exp\{y_2\beta_2(0) - y_2^2k_2/2\} \\ &\quad + \sqrt{2\pi/k'_1} \exp\{y_2\beta_{2;-1,0} + \beta'_1(y_2)^2/(2k'_1) - y_2^2k_2/2\}. \end{aligned} \quad (\text{C.17})$$

Now suppose by contradiction that  $\exp\{f_1(y_1)\}$  given  $y_1 \neq 0$  has only one Gaussian component, instead of being a sum of two Gaussian densities. Then by (C.16),  $\beta'_1(0) = \beta'_{1;-1,0}$  and  $\beta_2(y_1)$  is a constant given  $\mathbf{1}_{y_1}$ , i.e.  $\beta_2(y_1) = c_{\beta_2,-1} + c_{\beta_2,0} = \beta_{2;-1,0}$  for  $y_1 \neq 0$ . Plugging this into the left-hand side of (C.6) and integrating w.r.t.  $\lambda(y_1)$ , the continuous part ( $y_2 \neq 0$ ) of the marginal distribution of  $y_2$  given  $\mathbf{Y}_{\mathcal{P}} \equiv \mathbf{y}_{\mathcal{P}}$  is

$$\begin{aligned} \exp\{f'_2(y_2)\} &\propto \frac{\exp\{f_1(0) + \alpha_2(0) + y_2\beta_2(0) - y_2^2k_2/2\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(0) + \beta_2(0)^2/(2k_2)\} + 1} \\ &\quad + \exp\{y_2\beta_{2;-1,0} - y_2^2k_2/2\} \int_{\mathbb{R}} \frac{\exp\{\delta_1 + f_1(y_1) + \alpha_2(y_1)\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} dy_1, \end{aligned}$$

which is a mixture between  $\mathcal{N}(\beta_2(0)/k_2, 1/k_2)$  and  $\mathcal{N}(\beta_{2;-1,0}/k_2, 1/k_2)$ , i.e. the variance in both components are equal. Note that the integral in the second term is a Lebesgue integral. This together with (C.17) implies that  $\beta'_1(y_2)$  cannot depend on the value of  $y_2$  given  $y_2 \neq 0$ , i.e.  $\beta'_1(y_2) = c_{\beta'_1,-1} + c_{\beta'_1,0} = \beta'_{1;-1,0}$ . Since we already know that  $\beta'_1(0) = \beta'_{1;-1,0}$  by discussion above, this implies that  $\beta'_1$  is an absolute constant in  $y_2$  and  $\mathbf{1}_{y_2}$ , and also that  $\alpha_2$  may depend on  $y_1$  only through  $\mathbf{1}_{y_1}$ , a contradiction to the assumption of the theorem.

Thus, (C.16) and (C.17) will both have to be mixtures of precisely two Gaussians, and so by definition the joint distribution  $p(\mathbf{Y})$  of  $\mathbf{Y}$  must be of 2-Gaussian type with respect to  $\mathcal{G}$  and  $\mathcal{G}'$ .  $\square$

*Proof of Corollary 33.* When  $|\mathcal{V}| = 2$ , in Proposition 31 we always have  $\mathcal{P} = \emptyset$  and  $V_1$  does not have a parent in  $\mathcal{G}$ , so  $P(Y_{V_1} = y | Y_{V_1} \neq 0)$  by definition is just a Gaussian, not a mixture two Gaussians, and hence  $p(\mathbf{Y})$  cannot be of 2-Gaussian type with respect to any pairs of distinct Markov equivalent graphs.

Now consider  $|\mathcal{V}| = 3$ , and assume the two vertices with reversible edges in Proposition 31 are  $V_1$  and  $V_2$ , and that  $V_1 \rightarrow V_2$  in  $\mathcal{G}$  and  $V_1 \leftarrow V_2$  in  $\mathcal{G}'$ . If neither  $V_1$  or  $V_2$  has  $V_3$  as its parent in both graphs, then we can marginalize  $V_3$  out and it reduces to the 2-d case.

Suppose otherwise. Then we must have (1)  $V_1 \rightarrow V_2 \leftarrow V_3$  in  $\mathcal{G}$ , or (2)  $V_2 \rightarrow V_1 \leftarrow V_3$  in  $\mathcal{G}'$ , or (3) an additional edge between  $V_1$  and  $V_3$  added to (1), or (4) an additional edge between  $V_2$  and  $V_3$  added to (2).

For (1) and (2) both graphs are the only graph in their Markov equivalence class; for (3) the reversible edge becomes  $V_1 - V_3$  violating the assumption (and in fact one can marginalize out the common child  $V_2$  and get back to the 2-d case), and similarly for (4). Thus, we have again ruled out the possibility of any pair of distinct Markov equivalent graphs with respect to which  $p(\mathbf{Y})$  can be of 2-Gaussian type. □

**Remark 6.** *In the proof of Theorem 32, we proved that whenever  $p(\mathbf{Y})$  factorizes with respect to two distinct graphs  $\mathcal{G}$  and  $\mathcal{G}'$  (whenever identifiability does not hold), everything up to (C.17) in the proof must hold. Specifically, conditioning on almost every  $\mathbf{y}_{\mathcal{P}}$ ,  $\alpha_2$  and  $\beta_2$  in  $\mathcal{G}$  as well as  $\alpha'_1$  and  $\beta'_1$  in  $\mathcal{G}'$  can be at most linear in  $y_1$  and  $y_2$ , respectively, namely*

$$\begin{aligned} \beta'_1(y_2) &= c_{\beta'_1,-1} + c_{\beta'_1,0}\mathbf{1}_{y_2} + c_{\beta'_1,1}y_2, & \beta_2(y_1) &= c_{\beta_2,-1} + c_{\beta_2,0}\mathbf{1}_{y_1} + c_{\beta_2,1}y_1, \\ \alpha'_1(y_2) &= c_{\alpha'_1,-1} + c_{\alpha'_1,0}\mathbf{1}_{y_2} + c_{\alpha'_1,1}y_2, & \alpha_2(y_1) &= c_{\alpha_2,-1} + c_{\alpha_2,0}\mathbf{1}_{y_1} + c_{\alpha_2,1}y_1, \end{aligned}$$

with coefficients depending on  $\mathbf{y}_{\mathcal{P}}$  where

$$c_{\alpha'_1,1} = c_{\beta_2,0}, \quad c_{\alpha_2,1} = c_{\beta'_1,0}, \quad c_{\beta'_1,1} = c_{\beta_2,1}. \tag{C.18}$$

It is noted that, although not used in deriving our conclusion involving 2-Gaussian type distributions, we in addition also have the following results.

$$c_{\alpha'_1,-1} = c_{\alpha_2,-1}, \quad c_{\alpha'_1,0} = c_{\alpha_2,0}, \quad c_{\alpha'_1,-1} + c_{\alpha'_1,0} = c_{\alpha_2,-1} + c_{\alpha_2,0} = 0.$$

These might shed some light on how to show that distributions of 2-Gaussian type do not exist for a general  $m \geq 4$ .

*Proof of Remark 6.* By (C.6), (C.16), (C.17), the joint distribution of  $Y_1$  and  $Y_2$  conditional on  $\mathbf{Y}_{\mathcal{P}}$  has two characterizations (up to normalizing constants)

$$\frac{\exp\{\mathbf{1}_{y_1}\delta_1 + y_1\beta'_1(0) - y_1^2k'_1/2 + \mathbf{1}_{y_2}\alpha_2(y_1) + y_2\beta_2(y_1) - y_2^2k_2/2\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\}} + 1$$

$$\begin{aligned}
& + \frac{\sqrt{2\pi/k_2} \exp\{\mathbf{1}_{y_1} \delta_1 + y_1 \beta'_{1,-1,0} + \beta_2(y_1)^2/(2k_2) - y_1^2 k'_1/2 + \mathbf{1}_{y_2} \alpha_2(y_1) + y_2 \beta_2(y_1) - y_2^2 k_2/2\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
& \propto \frac{\exp\{\mathbf{1}_{y_2} \delta'_2 + y_2 \beta_2(0) - y_2^2 k_2/2 + \mathbf{1}_{y_1} \alpha'_1(y_2) + y_1 \beta'_1(y_2) - y_1^2 k'_1/2\}}{\sqrt{2\pi/k'_1} \exp\{\alpha'_1(y_2) + \beta'_1(y_2)^2/(2k'_1)\} + 1} \\
& + \frac{\sqrt{2\pi/k'_1} \exp\{\mathbf{1}_{y_2} \delta'_2 + y_2 \beta_{2,-1,0} + \beta'_1(y_2)^2/(2k'_1) - y_2^2 k_2/2 + \mathbf{1}_{y_1} \alpha'_1(y_2) + y_1 \beta'_1(y_2) - y_1^2 k'_1/2\}}{\sqrt{2\pi/k'_1} \exp\{\alpha'_1(y_2) + \beta'_1(y_2)^2/(2k'_1)\} + 1}.
\end{aligned} \tag{C.19}$$

Divide both sides by  $\exp(y_1 \beta'_1(0) + y_2 \beta_2(0) - y_1^2 k'_1/2 - y_2^2 k_2/2)$  and expanding  $\beta'_1(y_2)$  and  $\beta_2(y_1)$ , this becomes

$$\begin{aligned}
& \frac{\exp\{\mathbf{1}_{y_1} \delta_1 + \mathbf{1}_{y_2} \alpha_2(y_1) + y_2 c_{\beta_2,0} \mathbf{1}_{y_1} + y_1 y_2 c_{\beta_2,1}\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
& + \frac{\sqrt{2\pi/k_2} \exp\{\mathbf{1}_{y_1} \delta_1 + y_1 c_{\beta'_1,0} + \beta_2(y_1)^2/(2k_2) + \mathbf{1}_{y_2} \alpha_2(y_1) + y_2 c_{\beta_2,0} \mathbf{1}_{y_1} + y_1 y_2 c_{\beta_2,1}\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
& \propto \frac{\exp\{\mathbf{1}_{y_2} \delta'_2 + \mathbf{1}_{y_1} \alpha'_1(y_2) + y_1 c_{\beta'_1,0} \mathbf{1}_{y_2} + y_1 y_2 c_{\beta'_1,1}\}}{\sqrt{2\pi/k'_1} \exp\{\alpha'_1(y_2) + \beta'_1(y_2)^2/(2k'_1)\} + 1} \\
& + \frac{\sqrt{2\pi/k'_1} \exp\{\mathbf{1}_{y_2} \delta'_2 + y_2 c_{\beta_2,0} + \beta'_1(y_2)^2/(2k'_1) + \mathbf{1}_{y_1} \alpha'_1(y_2) + y_1 c_{\beta'_1,0} \mathbf{1}_{y_2} + y_1 y_2 c_{\beta'_1,1}\}}{\sqrt{2\pi/k'_1} \exp\{\alpha'_1(y_2) + \beta'_1(y_2)^2/(2k'_1)\} + 1}.
\end{aligned}$$

Now expanding  $\alpha'_1(y_2)$  and  $\alpha_2(y_1)$  and using the relationships in (C.18), we divide both sides by  $\exp(y_1 c_{\alpha_2,1} \mathbf{1}_{y_2} + y_2 c_{\beta_2,0} \mathbf{1}_{y_1} + y_1 y_2 c_{\beta_2,1}) = \exp(y_1 c_{\beta'_1,0} \mathbf{1}_{y_2} + y_2 c_{\alpha'_1,1} \mathbf{1}_{y_1} + y_1 y_2 c_{\beta_2,1})$  and get

$$\begin{aligned}
& \frac{\exp\{\mathbf{1}_{y_1} \delta_1 + \mathbf{1}_{y_2} (c_{\alpha_2,-1} + c_{\alpha_2,0} \mathbf{1}_{y_1})\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
& + \frac{\sqrt{2\pi/k_2} \exp\{\mathbf{1}_{y_1} \delta_1 + y_1 c_{\beta'_1,0} + \beta_2(y_1)^2/(2k_2) + \mathbf{1}_{y_2} (c_{\alpha_2,-1} + c_{\alpha_2,0} \mathbf{1}_{y_1})\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
& = C_0 \frac{\exp\{\mathbf{1}_{y_2} \delta'_2 + \mathbf{1}_{y_1} (c_{\alpha'_1,-1} + c_{\alpha'_1,0} \mathbf{1}_{y_2})\}}{\sqrt{2\pi/k'_1} \exp\{\alpha'_1(y_2) + \beta'_1(y_2)^2/(2k'_1)\} + 1} \\
& + C_0 \frac{\sqrt{2\pi/k'_1} \exp\{\mathbf{1}_{y_2} \delta'_2 + y_2 c_{\beta_2,0} + \beta'_1(y_2)^2/(2k'_1) + \mathbf{1}_{y_1} (c_{\alpha'_1,-1} + c_{\alpha'_1,0} \mathbf{1}_{y_2})\}}{\sqrt{2\pi/k'_1} \exp\{\alpha'_1(y_2) + \beta'_1(y_2)^2/(2k'_1)\} + 1}
\end{aligned} \tag{C.20}$$

for some  $C_0$ . Setting  $\mathbf{1}_{y_1} = \mathbf{1}_{y_2} = 0$  (C.20) becomes

$$\frac{1 + \sqrt{2\pi/k_2} \exp\{c_{\beta_2,-1}^2/(2k_2)\}}{\sqrt{2\pi/k_2} \exp\{c_{\alpha_2,-1} + c_{\beta_2,-1}^2/(2k_2)\} + 1} = C_0 \frac{1 + \sqrt{2\pi/k'_1} \exp\{c_{\beta'_{11},-1}^2/(2k'_1)\}}{\sqrt{2\pi/k'_1} \exp\{c_{\alpha'_{1,-1}} + c_{\beta'_{1,-1}}^2/(2k'_1)\} + 1}, \quad (\text{C.21})$$

and with  $\mathbb{1}_{y_1} \neq 0$ ,  $\mathbb{1}_{y_2} = 0$  (C.20) becomes

$$\begin{aligned} \exp(\delta_1) \frac{1 + \sqrt{2\pi/k_2} \exp\{y_1 c_{\beta'_{1,0}} + \beta_2(y_1)^2/(2k_2)\}}{\sqrt{2\pi/k_2} \exp\{c_{\alpha_2,-1} + c_{\alpha_2,0} + c_{\alpha_2,1}y_1 + \beta_2(y_1)^2/(2k_2)\} + 1} \\ = C_0 \exp(c_{\alpha'_{1,-1}}) \frac{1 + \sqrt{2\pi/k'_1} \exp\{c_{\beta'_{11},-1}^2/(2k'_1)\}}{\sqrt{2\pi/k'_1} \exp\{c_{\alpha'_{1,-1}} + c_{\beta'_{1,-1}}^2/(2k'_1)\} + 1}. \end{aligned} \quad (\text{C.22})$$

Since the right-hand side of (C.22) is a constant, by matching the numerator and the denominator of the left-hand side using Lemma 39, we must have either (i)  $y_1 c_{\beta'_{1,0}} + \beta_2(y_1)^2/(2k_2) = c_{\alpha_2,-1} + c_{\alpha_2,0} + c_{\alpha_2,1}y_1 + \beta_2(y_1)^2/(2k_2)$ , or (ii)  $y_1 c_{\beta'_{1,0}} + \beta_2(y_1)^2/(2k_2) = \text{const}$  for  $y_1 \neq 0$ . But (ii) implies that  $c_{\beta_2,1} = c_{\beta'_{1,0}} = 0$ , which by  $c_{\beta'_{1,1}} = c_{\beta_2,1}$  implies that  $\beta'_1$  is an absolute constant in  $y_2 \in \mathbb{R}$ , a violation to the assumption. Thus (i) holds, and by  $c_{\beta'_{1,0}} = c_{\alpha_2,1}$  this implies that

$$\alpha_{2,-1,0} \equiv c_{\alpha_2,-1} + c_{\alpha_2,0} = 0, \quad \text{and by symmetry} \quad \alpha'_{1,-1,0} \equiv c_{\alpha'_{1,-1}} + c_{\alpha'_{1,0}} = 0. \quad (\text{C.23})$$

Thus the left-hand side of (C.22) is just  $\exp(\delta_1)$ . Note that the right-hand side of (C.22) is  $\exp(c_{\alpha'_{1,-1}})$  times the right-hand side of (C.21). So by equating the left-hand side of (C.22) with  $\exp(c_{\alpha'_{1,-1}})$  times the left-hand side of (C.21) we have

$$\exp(\delta_1) = \exp(c'_{\alpha_{1,-1}}) \frac{1 + \sqrt{2\pi/k_2} \exp\{c_{\beta_2,-1}^2/(2k_2)\}}{\sqrt{2\pi/k_2} \exp\{c_{\alpha_2,-1} + c_{\beta_2,-1}^2/(2k_2)\} + 1} \quad (\text{C.24})$$

and similarly

$$\exp(\delta'_2) = \exp(c_{\alpha_2,-1}) \frac{1 + \sqrt{2\pi/k'_1} \exp\{c_{\beta'_{11},-1}^2/(2k'_1)\}}{\sqrt{2\pi/k'_1} \exp\{c_{\alpha'_{1,-1}} + c_{\beta'_{1,-1}}^2/(2k'_1)\} + 1}. \quad (\text{C.25})$$

Now by (C.23), with  $\mathbb{1}_{y_1} = \mathbb{1}_{y_2} = 1$ , (C.20) simplifies to  $\exp(\delta_1) = C_0 \cdot \exp(\delta'_2)$ . Thus by

(C.21), (C.24) and (C.25), one get

$$C_0 = \frac{\exp(\delta_1)}{\exp(\delta'_2)} = \frac{\exp(c'_{\alpha_1,-1}) \frac{1 + \sqrt{2\pi/k_2} \exp\{c_{\beta_2,-1}^2/(2k_2)\}}{\sqrt{2\pi/k_2} \exp\{c_{\alpha_2,-1} + c_{\beta_2,-1}^2/(2k_2) + 1\}}}{\exp(c_{\alpha_2,-1}) \frac{1 + \sqrt{2\pi/k'_1} \exp\{c_{\beta'_1,-1}^2/(2k'_1)\}}{\sqrt{2\pi/k'_1} \exp\{c_{\alpha'_1,-1} + c_{\beta'_1,-1}^2/(2k'_1) + 1\}}} = \frac{\exp(c'_{\alpha_1,-1})}{\exp(c_{\alpha_2,-1})} C_0$$

and thus  $c'_{\alpha_1,-1} = c_{\alpha_2,-1}$ . Combining with (C.23), we get

$$c_{\alpha'_1,-1} = c_{\alpha_2,-1}, \quad c_{\alpha'_1,0} = c_{\alpha_2,0}, \quad c_{\alpha'_1,-1} + c_{\alpha'_1,0} = c_{\alpha_2,-1} + c_{\alpha_2,0} = 0. \quad (\text{C.26})$$

Note that this result holds as long as we assume identifiability does not hold.

□