

Directed Graphical Models and Causal Discovery for Zero-Inflated Data

Shiqing Yu

Department of Statistics, University of Washington, Seattle, U.S.A.

SYU.PHD@GMAIL.COM

Mathias Drton

Department of Mathematics and Munich Data Science Institute, Technical University of Munich, Germany

MATHIAS.DRTON@TUM.DE

Ali Shojaie

Department of Biostatistics, University of Washington, Seattle, U.S.A.

ASHOJAIE@UW.EDU

Editors:

Mihaela van der Schaar, Dominik Janzing and Cheng Zhang

Abstract

With advances in technology, gene expression measurements from single cells can be used to gain refined insights into regulatory relationships among genes. Directed graphical models are well-suited to explore such (cause-effect) relationships. However, statistical analyses of single cell data are complicated by the fact that the data often show zero-inflated expression patterns. To address this challenge, we propose directed graphical models that are based on Hurdle conditional distributions parametrized in terms of polynomials in parent variables and their 0/1 indicators of being zero or nonzero. While directed graphs for Gaussian models are only identifiable up to an equivalence class in general, we show that, under a natural and weak assumption, the exact directed acyclic graph of our zero-inflated models can be identified. We propose methods for graph recovery, apply our model to real single-cell gene expression data on T helper cells, and show simulated experiments that validate the identifiability and graph estimation methods in practice.

Keywords: Bayesian network, causal discovery, directed acyclic graph, identifiability

1. Introduction

Graphical models specify conditional independence relations among variables in a random vector 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, Part V). 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 Gaussian graphical models, 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

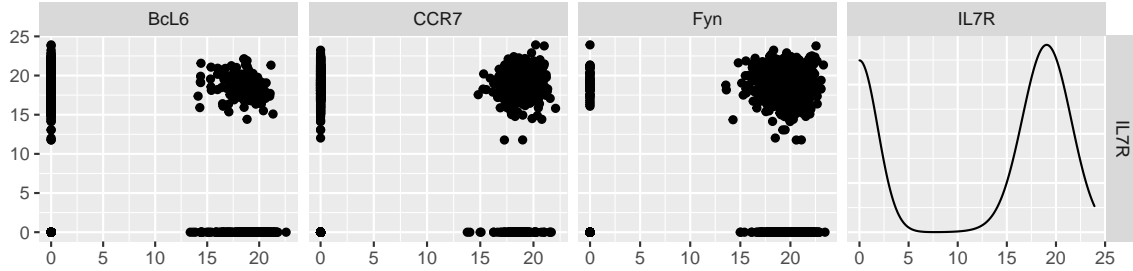


Figure 1: Kernel density of a selected gene (IL7R) and scatter plots of its relationship with three other genes (Bcl6, CCR7, Fyn) from the T helper cell data analyzed in Section 6.

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). This pattern is shown in Figure 1, which is based on a single-cell dataset with 1951 measurements from eight healthy donors, which we analyze in Section 6. The figure clearly shows the large number of zero values in each gene as well as the nonlinear relationships between genes. These effects create challenges for existing causal discovery procedures. Specifically, disentangling the on/off status of the genes from the quantitative mRNA expression levels requires new methods that account for the zero inflation, i.e., the excessive number of zeros in the data. Motivated by the Gaussian-like distribution of the mRNA expression levels, a natural choice is to model the causal network using conditional zero-inflated Gaussian distributions.

In this paper, we propose two versions of directed graphical models for zero-inflated data, and prove that under a weak assumption the exact DAG can be recovered from the joint distribution. Our new graphical models build on the recently-proposed Hurdle graphical model of McDavid et al. (2019), but facilitate estimation of DAGs from observational single-cell sequencing data. In contrast to McDavid et al. (2019), our models are also not limited to zero-inflated Gaussian distributions, as we allow variables that are “on” to be non-linear polynomial functions of other variables and stochastic noise. The proposed model and corresponding identifiability theory differs from the recent proposal of Choi et al. (2020), in which the data are always counts, with additional zero-inflation. Specifically, we model the on/off status of each gene, conditional on its parents, with a Bernoulli random variable. Then, conditional on the event that the gene is on, the expression level is modeled by a Gaussian distribution depending on the parents. After presenting two directed graphical models for zero-inflated data in Section 2, in Section 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 7). 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. In Section 4, we introduce different methods for estimation

of the DAG. Simulation studies supporting the use of these methods are given in Section 5, and they are then applied to the T-follicular helper cell dataset (Section 6). Throughout the paper, we use subscripts to 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., $\mathbf{y}_V = (y_V)_{V \in \mathcal{V}}$.

2. Directed Graphical Models for Zero-Inflated Data

2.1. Hurdle Joint Distributions for Zero-Inflated Continuous Observations

We start by reviewing the undirected Hurdle graphical model. [McDavid et al. \(2019\)](#) proposed a *Hurdle joint distribution* with density

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

with respect to λ^m , where λ 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 $\mathbf{1}_y \equiv (\mathbb{1}_{\{y_1 \neq 0\}}, \dots, \mathbb{1}_{\{y_m \neq 0\}}) \in \{0, 1\}^m$ captures which components of \mathbf{y} are non-zero.

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

$$p(Y_i = y_i | \mathbf{Y}_{-i} = \mathbf{y}_{-i}) = f(y_i; \alpha_{ii} + \alpha_{i,-i}^\top \mathbf{1}_{\mathbf{y}_{-i}} + \beta_{i,-i}^\top \mathbf{y}_{-i}, \beta_{ii} + \beta_{-i,i}^\top \mathbf{1}_{\mathbf{y}_{-i}} - \mathbf{k}_{i,-i}^\top \mathbf{y}_{-i}, k_{ii}), \quad (2)$$

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

2.2. Hurdle Conditionals

The observation in (2) above gives rise to the following definition. Recall that λ is the sum of a point mass at 0 and the Lebesgue measure on \mathbb{R} .

Definition 1 ((α, β, k) -Hurdle conditionals) *Let X be a scalar random variable, and let \mathbf{Z} be an m -dimensional random vector. We say that the conditional distribution of X given \mathbf{Z} is of (α, β, k) -Hurdle type if it admits conditional densities with respect to λ of the form*

$$p(X = x | \mathbf{Z} = \mathbf{z}) = f_{\alpha, \beta, k}^{(m)}(X | \mathbf{Z}) \equiv \frac{\exp(\alpha(\mathbf{z}) \mathbb{1}_x + \beta(\mathbf{z})x - kx^2/2)}{\sqrt{2\pi/k} \exp(\alpha(\mathbf{z}) + \beta^2(\mathbf{z})/(2k)) + 1}. \quad (3)$$

Here, α and β 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 2 ((p, μ, σ^2) -Hurdle conditionals) *Let X be a scalar random variable, and let \mathbf{Z} be an m -dimensional random vector. We say that the conditional distribution of X given \mathbf{Z} is of (p, μ, σ^2) -Hurdle type if it admits conditional densities with respect to λ of the form*

$$p(X = x | \mathbf{Z} = \mathbf{z}) = f_{p, \mu, \sigma^2}^{(m)}(X | \mathbf{Z}) \equiv (1 - p(\mathbf{z}))(1 - \mathbb{1}_x) + p(\mathbf{z}) \mathbb{1}_x \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu(\mathbf{z}))^2}{2\sigma^2}\right). \quad (4)$$

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

It is easy to show that the two parametrizations (3) and (4) are connected through

$$\log \frac{p}{1-p} = \alpha + \frac{\beta^2}{2k} - \frac{1}{2} \log\left(\frac{k}{2\pi}\right), \quad \mu = \frac{\beta}{k}, \quad \sigma^2 = \frac{1}{k}. \quad (5)$$

That is, the conditional log odds of being nonzero is linear in α and quadratic in β , and the conditional Gaussian mean is proportional to β . While the (α, β, 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 σ^2 parameters of the conditional Gaussian distribution. We thus refer to (3) as the *canonical parametrization*, and (4) as the *moment parametrization*.

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

$$f(\mathbf{y}) = \prod_{V \in \mathcal{V}} f_V(y_V | \mathbf{y}_{\text{pa}(V)}), \quad (6)$$

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 2.1, we observed that, for the Hurdle joint distributions (1), the conditional distribution of any Y_i given the others is an (α, β, k) -Hurdle with k constant, and α and β linear functions of those variables (and their indicators) that are conditionally dependent on Y_i ; see (2). Motivated by this fact, we specify directed graphical models for zero-inflated data by assuming the conditional densities in the factorization in (6) to be (α, β, k) - or (p, μ, σ^2) -Hurdle conditionals. We then assume the parameters in these conditionals to be *Hurdle polynomials* in its parents, as defined now.

Definition 3 (Hurdle polynomials) *Let $\mathbf{Y} = (Y_V)_{V \in \mathcal{V}} \in \mathbb{R}^m$ be a 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*

$$\mathcal{H}(\mathbf{Y}; \mathcal{S}) \equiv \left\{ c_0 + \sum_{j=1}^T c_j \prod_{U \in \mathcal{U}_j} Y_U^{d_{j,U}} \prod_{V \in \mathcal{V}_j} \mathbb{1}_{Y_V}, \quad c_0 \in \mathbb{R}, T \in \mathbb{N}, \right. \\ \left. c_j \neq 0, \mathcal{U}_j \subseteq \mathcal{S}, \mathcal{V}_j \subseteq \mathcal{S} \setminus \mathcal{U}_j, d_{j,U} \in \mathbb{N} \quad \forall U \in \mathcal{U}_j \quad \forall j = 1, \dots, T \right\}, \quad (7)$$

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 (7) is $\max_{j=1, \dots, T} \sum_{U \in \mathcal{U}_j} d_{j,U} + |\mathcal{V}_j|$. Here $|\cdot|$ denotes the set cardinality.

In the definition (7), for the j -th term in the polynomial, c_j is its polynomial coefficient, \mathcal{V}_j is the set of nodes that define the term only through their indicators, while \mathcal{U}_j is the set of those whose values directly define the term, with $\{d_{j,U}\}_{U \in \mathcal{U}_j}$ the corresponding exponents.

We are now ready to formally define our models.

Definition 4 (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 λ^m) that factors as in (6) with each conditional density $f_V(y_V | \mathbf{y}_{\text{pa}(V)})$ being a Hurdle conditional

- (1) in the (α, β, k) -parametrization with parameters α_V , β_V and k_V , where k_V is constant, α_V and β_V are Hurdle polynomials in $\mathbf{y}_{\text{pa}(V)}$; or
- (2) in the (p, μ, σ^2) -parametrization with parameters p_V , μ_V and σ_V^2 , where σ_V^2 is constant, $\log(p_V/(1-p_V))$ and μ_V are Hurdle polynomials in $\mathbf{y}_{\text{pa}(V)}$.

It is clear from (5) that if we let 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.

3. Identifiability

3.1. Strong Identifiability

As we show next, the directed graphical models from Definition 4 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 a mild assumption on the Hurdle conditionals. Let $\pi(\mathbf{y}_{\mathcal{S}}) \in \mathcal{H}(\mathbf{Y}; \mathcal{S})$ 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 $\pi(\mathbf{y}_{\mathcal{S}})$ 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 5 (Strong Hurdle polynomials) Let $\pi(\mathbf{y}_{\mathcal{S}}) \in \mathcal{H}(\mathbf{Y}; \mathcal{S})$. We say $\pi(\mathbf{y}_{\mathcal{S}})$ is 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 $d \geq 1$.

Our first theorem gives an identifiability result that invokes a faithfulness assumption; see Section 15.3.2 of [Maathuis et al. \(2019\)](#) for a definition and discussion of faithfulness.

Theorem 6 (DAG identifiability with strong Hurdle polynomials) Let $f(\mathbf{y})$ be a joint density with respect to λ^m that is faithful w.r.t. and factors according to a DAG $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, as in (6). Suppose for each $V \in \mathcal{V}$, the conditional $f_V(y_V | \mathbf{y}_{\text{pa}(V)})$ is of Hurdle type with parameters

(α_V, β_V, k_V) or (p_V, μ_V, σ_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 there does not exist any other DAG $\mathcal{G}' \neq \mathcal{G}$ such that $f(\mathbf{y})$ factors and is faithful w.r.t. \mathcal{G}' .

In the proof in Appendix A, we show that the 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.

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 μ and inverse variance ν .

Definition 7 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 exist a constant ν_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 (a) $\mu_1^V(\mathbf{y}_{\mathcal{P}}) \neq \mu_2^V(\mathbf{y}_{\mathcal{P}})$; or (b) $\nu_1^V \neq \nu_2^V(\mathbf{y}_{\mathcal{P}})$ 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 $\mathbf{y}_{\mathcal{P}}$; the inverse variance parameter is an absolute constant for one of these distributions and polynomial in $\mathbf{y}_{\mathcal{P}}$ for the other. 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 an observation from Proposition 29(ii) in [Peters et al. \(2014\)](#); see Section 1.8 of [Maathuis et al. \(2019\)](#) for background on Markov properties.

Proposition 8 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 9 Proposition 8 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. Similar to Proposition 8, our results also require faithfulness, which is natural in this setting, as limits of parameters recover the Gaussian/binary case for which faithful distributions exist.

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 8.

Theorem 10 (General Identifiability) *Let \mathbf{Y} , \mathcal{G} , \mathcal{G}' , W , U , \mathcal{P} be as in Proposition 8. Let \mathbf{Y} have a λ^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 (3), or equivalently $(p_V^{\mathcal{H}})_{V \in \mathcal{V}}$ and $(\mu_V^{\mathcal{H}})_{V \in \mathcal{V}}$ from (4), that are Hurdle polynomials of the form (7), 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 8 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 10.

The result is proved in Appendix A. 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 A 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 11 (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 satisfying the conditions of Theorem 10 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 6, meaning that the sets of Markov and faithful distributions associated to \mathcal{G} and \mathcal{G}' must be disjoint.*

Theorem 6 and Corollary 11 state that the DAGs are perfectly identifiable if $m = 2, 3$ or if we assume the Hurdle polynomials to be *strong*; Theorem 10 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 Appendix B we show scatter plots of simulated data that give some indication of how Markov equivalent graphs may be differentiated under our models.

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 ℓ of any DAG model can be decomposed into the sum of conditional (or nodewise) log-likelihood functions ℓ^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 (α, β, k) -parametrization from (3), the nodewise log-likelihood function is

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

for the moment (p, μ, σ^2) -parametrization from (4) it is

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

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 μ_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 ℓ^V as a function of the parameters (i.e., coefficients) in those polynomials.

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 (α, β, k) -parametrization, the log-likelihood function is convex in α_V, β_V and k_V . Moreover, α_V and β_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, μ, σ^2) -parametrization (4), on the other hand, can be easily solved by separately fitting a logistic regression to p_V and a linear regression to μ_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 (α, β, k) -parametrization with linear Hurdle polynomials (i.e., degree 1) naturally comes from conditional distributions of the joint distribution defined for undirected graphical models in [McDavid et al. \(2019\)](#). However, at least for higher degrees, the (p, μ, σ^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 us to use optimized standard regression solvers for model fitting. The (p, μ, σ^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. In our implementation, we use an ℓ_2 regularization and select its tuning parameter using the Bayesian information criterion (BIC). We also assume the highest degree of Hurdle polynomials and select the degree by optimizing BIC simultaneously over the degree and the ℓ_2 penalty, so the degree is separately optimized for each regression (combination of node and its candidate parent set).

4.2. Graph Search

To estimate the underlying DAG, we 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 ν is the total number of parameters in the model, n is the sample size, and ℓ is the log-likelihood as introduced in Section 4.

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 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 $m < 16$.

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).

There are various approaches that may help accelerating the estimation process. As an example, one can use caching ([Ramsey et al. 2017](#)) and dynamic updating ([Goudie and Mukherjee 2016](#)) to save time on computing the likelihoods and checking acyclicity in the current estimated graph. To speed up the estimation, we cache the BICs of all the nodewise regressions that have been fit so far, which requires little memory overhead. As the greedy search may be stuck in a local minimum, the most obvious way to circumvent this is to run the greedy algorithm initialized with multiple random DAGs with the same number of nodes and different sparsity levels, and choose the output that has the lowest BIC. Alternatively, one can first estimate an undirected graph using the method of [McDavid et al. \(2019\)](#), and initialize the search with multiple directed graphs whose moral graph is the estimated undirected graph. Moreover, to scale to larger m , we can first use the procedure of [McDavid et al. \(2019\)](#) to identify the connected components of the estimated undirected graph and then estimate the directed edges in each connected component. This procedure is justified by the fact that the connected components for the underlying true undirected and directed graphs coincide.

5. Numerical Experiments

Our numerical studies in this section aim to verify identifiability and exact DAG recovery. Due to space limitation, we present the results for the exhaustive search in the main paper. Following the discussion in Section 4.2, we use our self-implemented greedy search (GDS) ([Chickering 2003](#)) with BIC score, as well as an exhaustive search with dynamic programming ([Silander and Myllymäki](#)

2006). The results for the greedy search—which facilitates estimation of much larger DAGs—show similar trends and are presented in Appendix B.

To illustrate the performance of exhaustive search, we consider three DAG structures: (i) chain graph with $m = 10$, (ii) complete graph with $m = 10$, and (iii) lattice graph with $m = 9$. For each structure, we consider true generating conditional distributions using the following parametrizations: a) (α, β, k) -(canonical) parametrization with *linear* Hurdle polynomials, b) (p, μ, σ^2) -(moment) parametrization with *linear* Hurdle polynomials, and c) (p, μ, σ^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 (5), distributions represented by a) and b) are disjoint because $\log(p/(1-p))$ is a weighted sum of α and β^2 .

Recall the definition of Hurdle conditionals in (3) and (4) in Section 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 ϕ is the standard normal density. Otherwise, for parametrization a), we use Hurdle conditionals with parameters $k_V = 1$, $\alpha_V(\mathbf{y}_{\text{pa}(V)}) = \beta_V(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} (\mathbf{1}_{y_U} - y_U)$; similarly for parametrization b) we take $\sigma_V^2 = 1$, $\log \frac{p_V}{1-p_V}(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} (\mathbf{1}_{y_U} + y_U)$ and $\mu_V(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} (\mathbf{1}_{y_U} - y_U)$; finally, for parametrization c) we take $\sigma_V^2 = 1$ and $\log \frac{p_V}{1-p_V}(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} \left(\mathbf{1}_{y_U} + y_U + \frac{y_U^2}{10} \right) + \frac{1}{10} \sum_{\substack{U, W \in \text{pa}(V) \\ U \neq W}} (\mathbf{1}_{y_U} + y_U)(\mathbf{1}_{y_W} + y_W)$, and $\mu_V(\mathbf{y}_{\text{pa}(V)}) = \sum_{U \in \text{pa}(V)} \left(\mathbf{1}_{y_U} - y_U - \frac{y_U^2}{10} \right) + \frac{1}{10} \sum_{\substack{U, W \in \text{pa}(V) \\ U \neq W}} (\mathbf{1}_{y_U} \mathbf{1}_{y_W} - y_U \mathbf{1}_{y_W} - y_V \mathbf{1}_{y_U} - y_V y_U)$. 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 μ_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 (see Figure S2 in the Appendix for some scatter plots of the data generated). To assess the effect of misspecified parametrizations, for each combination of true DAG and true data generating parametrization— (α, β, k) -linear and (p, μ, σ^2) -linear and quadratic—we estimate the DAG using all three parametrizations for generating data.

The results are shown in Figures 2. Due to space limitation, only results for correctly specified models are presented in the main paper, and the expanded results with misspecified models are given in the Appendix. Each row of the figure corresponds to one of the three graphs (chain, complete, lattice) and each column corresponds to results using one estimating parametrization. The plots show the average true positive rate (TPR) and false discovery rate (FDR) over $B = 100$ iterations, defined as $\text{TPR} = |\hat{S} \cap S_0|/|S_0|$ and $\text{FDR} = |\hat{S} \setminus S_0|/\max\{|\hat{S}|, 1\}$, where \hat{S} denotes the estimated set of (directed) edges, and S_0 the set of true edges.

The results indeed indicate that in all settings, exhaustive search with correct parametrization almost always identifies the exact DAG for large n . The results in the Appendix show that model misspecification does not seem to negatively impact the performance. Overall, our simulation studies confirm the identifiability theory (Theorem 6).

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⁺PD1⁺) cells that produce high levels of proteins CXCR5 and PD1, while we do not make this restriction. Instead, we add the indicators of CXCR5^{+/-} and PD1^{+/-} 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. Following [Section 4.2](#) we first estimate the connected components using the method from [McDavid et al. \(2019\)](#) for undirected graphs, and proceed with estimation of DAGs for each component. We use the (p, μ, σ^2) -parametrization as it is more flexible than the (α, β, 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 5](#), the (p, μ, σ^2) is also more robust than (α, β, 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 shows promising performance in the simulations in [Appendix B](#). We also use the stability selection procedure of [Shah and Samworth \(2013\)](#) to control 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.

[Figure 3](#) shows the estimated directed network along with the undirected network obtained using the method of [McDavid et al. \(2019\)](#). Overall, the estimated DAG structure is very similar to the undirected graph, with few differences including isolated nodes that only have a single edge with weak association in the undirected network.

7. Discussion

Motivated by the recent advent of single-cell sequencing technologies, we propose new methods for learning DAGs from zero-inflated data. Our procedures take advantage of two key features of single-cell transcriptomics data, namely, the zero-inflation, and the large number of observations from individual samples. 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.

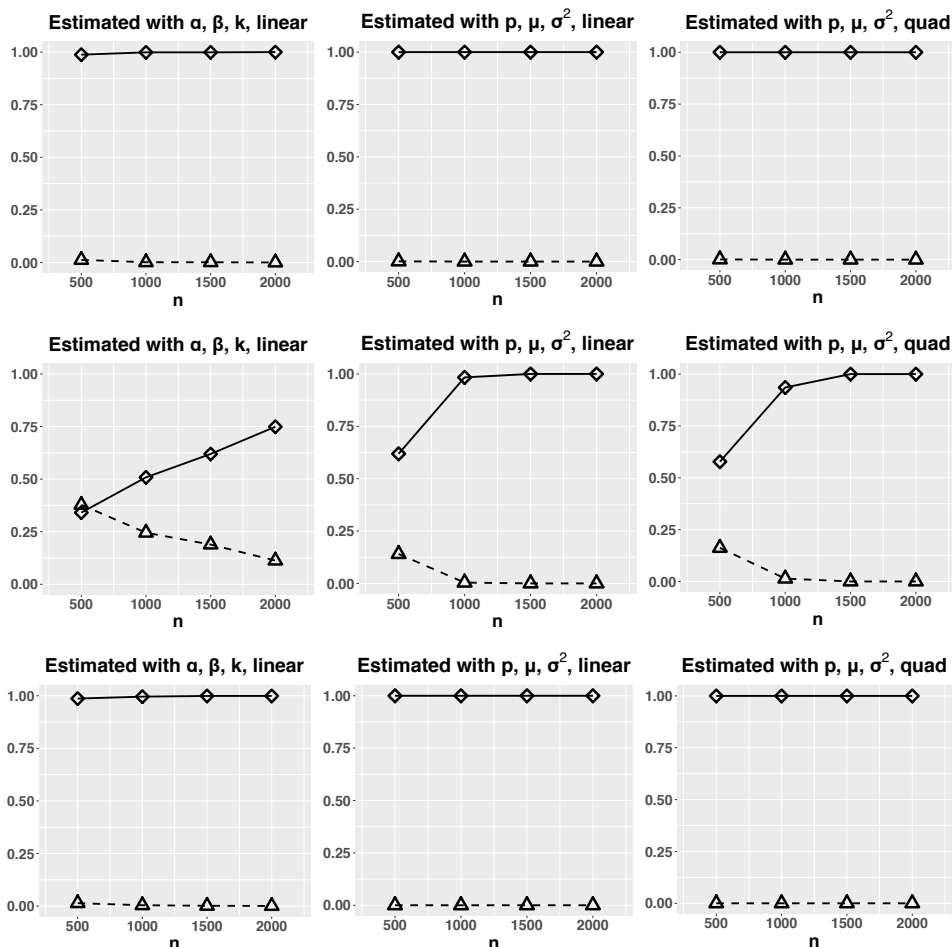


Figure 2: Simulation results for exhaustive search. Each row corresponds to a different graph (chain, complete, lattice). In all cases, estimation methods match true parametrizations. ‘ \diamond ’ with solid lines: true positive rate; ‘ \triangle ’ with dashed lines: false discovery rate.

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 (α, β, k) that is linked to the undirected graphs studied in [McDavid et al. \(2019\)](#), and the other called (p, μ, σ^2) that directly models the conditional moments. Both approaches have computational appeal. In particular, the (α, β, k) -parametrization leads to convex loss functions in the parameters to be estimated, while the (p, μ, σ^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.

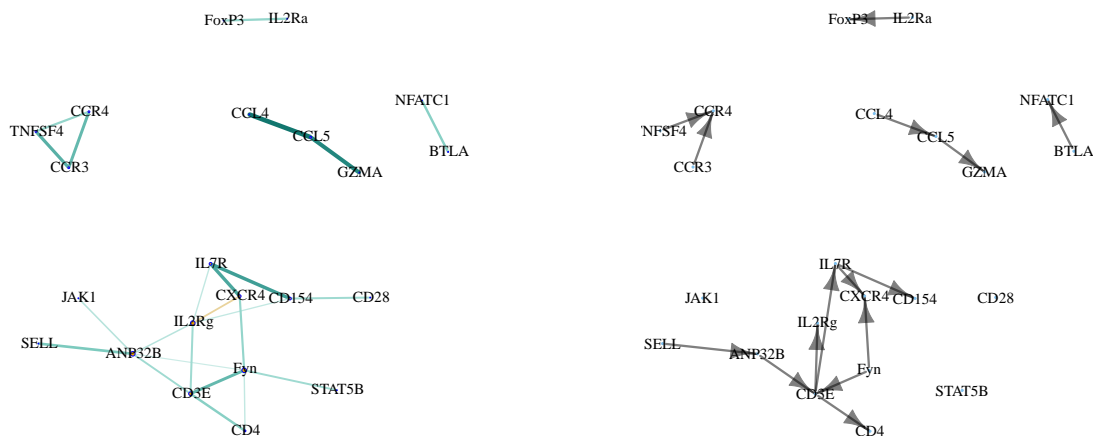


Figure 3: Estimated graph for T helper single cell data. Left: Undirected graph using the method of [McDavid et al. \(2019\)](#), with edge width and saturation representing the edge strength. Right: Directed graph using our method with stability selection [Shah and Samworth \(2013\)](#) to control FDR.

Several extension of our work would be of interest. The first is to prove our conjecture that the sets of distributions from which the DAG is not identifiable are also empty for graphs with more than 3 nodes. The second is proving the consistency and investigating finite sample properties of the proposed estimation procedures. Finally, it would be interesting to extend our model to zero-inflated distributions under a truncation to the nonnegative orthant \mathbb{R}_+^m , which would be of interest for nonnegative *omics* data by generalizing the *score matching* loss ([Hyvärinen 2005, 2007](#); [Lyu 2009](#); [Yu et al. 2019](#)) to data of mixed type.

Acknowledgments

The authors gratefully acknowledge grant DMS/NIGMS-1561814 from the US National Science Foundation (NSF) and grant R01-GM114029 from the US National Institutes of Health (NIH). This project has also received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement No 883818). The authors also thank Andrew McDavid, Jonas Peters, and Steffen Lauritzen for helpful discussions.

References

- Reka Albert. Scale-free networks in cell biology. *Journal of Cell Science*, 118(21):4947–4957, 2005.
- 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.
- Junsouk Choi, Robert Chapkin, and Yang Ni. Bayesian causal structural learning with zero-inflated poisson bayesian networks. *Advances in Neural Information Processing Systems*, 33:5887–5897, 2020.
- Robert JB Goudie and Sach Mukherjee. A Gibbs sampler for learning DAGs. *The Journal of Machine Learning Research*, 17(1):1032–1070, 2016.
- Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(Apr):695–709, 2005.
- Aapo Hyvärinen. Some extensions of score matching. *Computational Statistics & Data Analysis*, 51(5):2499–2512, 2007.
- 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.
- Masashi Okamoto. Distinctness of the eigenvalues of a quadratic form in a multivariate sample. *The Annals of Statistics*, 1:763–765, 1973.
- 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.
- Joseph Ramsey, Madelyn Glymour, Ruben Sanchez-Romero, and Clark Glymour. A million variables and more: the fast greedy equivalence search algorithm for learning high-dimensional graphical causal models, with an application to functional magnetic resonance images. *International journal of data science and analytics*, 3(2):121–129, 2017.

- 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.
- Y. Samuel Wang and Mathias Drton. High-dimensional causal discovery under non-Gaussianity. *Biometrika*, 107(1):41–59, 2020.
- Shiqing Yu, Mathias Drton, and Ali Shojaie. Generalized score matching for non-negative data. *Journal of Machine Learning Research*, 20(76):1–70, 2019.

A. Proofs

In this appendix we present proofs for the theorems and corollaries in the paper.

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 12 *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 ℓ_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

$$\sum_{i=1}^{N_a} a_0^i f^{(m)}(\mathbf{y}; \mathbf{a}^i) = \sum_{i=1}^{N_b} b_0^i f^{(m)}(\mathbf{y}; \mathbf{b}^i) \quad (8)$$

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 (8) are a sum of distinct exponentials of polynomials.

Then we must have $N_a = N_b$ and there is a permutation π 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 (8).

Proof [Proof of Lemma 12] First note that both sides of (8) 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. (8) 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 $a_0 \exp(a_1 x + \dots + a_p x^p)$ goes to

- (i) $a_0 \neq 0$ if $a_1 = \dots = a_p = 0$, 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 (8) 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, $\mathbf{a}^i - \mathbf{a}^j \neq 0$, so $d_{\max \neq 0}(\mathbf{a}^i - \mathbf{a}^j)$ exists and

this rearrangement is possible. Now dividing both sides of (8) by $f^{(1)}(y_1; \mathbf{a}^1) = \exp(a_1^1 y_1 + \dots + a_p^1 y_1^p)$ we have

$$a_0^1 + \sum_{i=2}^{N_a} a_0^i f^{(1)}(y_1; \mathbf{a}^i - \mathbf{a}^1) = \sum_{i=1}^{N_b} b_0^i f^{(1)}(y_1; \mathbf{b}^i - \mathbf{a}^1). \quad (9)$$

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 (9) 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 (9) 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

$$\sum_{i=N_b+1}^{N_a} a_0^i f^{(1)}(y_1; \mathbf{a}^i - \mathbf{a}^{N_b}) = \text{const.}$$

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 $a_0^{N_b+1} + \sum_{i=N_b+2}^{N_a} a_0^i f^{(1)}(y_1; \mathbf{a}^i - \mathbf{a}^{N_b+1}) = f^{(1)}(y_1; \mathbf{a}^{N_b} - \mathbf{a}^{N_b+1})$. 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 $\sum_{d \in \mathcal{D}, d_1=d} a_d^i \prod_{j=2}^m y_j^{d_j}$, 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 (8) 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)),

$$\sum_{\ell=1}^C f^{(1)}(y_1; \mathbf{c}_{1*}^\ell(\mathbf{y}_{-1})) \left\{ \sum_{j=1}^{n_\ell^a} a_0^{k_{\ell j}^a} f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{a}_{-1}^{k_{\ell j}^a}) + \sum_{j=1}^{n_\ell^b} b_0^{k_{\ell j}^b} f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{b}_{-1}^{k_{\ell j}^b}) \right\} = 0, \quad (10)$$

where $C > 0$, each \mathbf{c}_1^ℓ (coefficients for \mathbf{c}_{1*}^ℓ) is some \mathbf{a}_1^i or \mathbf{b}_1^i , and $\{\mathbf{c}_1^\ell\}_{\ell=1}^C$ are distinct. Here, let $\{k_{11}^a, \dots, k_{1,n_1^a}^a, \dots, k_{C1}^a, \dots, k_{C,n_C^a}^a\}$ be a permutation of $\{1, \dots, N_a\}$, and $\{k_{11}^b, \dots, k_{1,n_1^b}^b, \dots, k_{C1}^b, \dots, k_{C,n_C^b}^b\}$ a permutation of $\{1, \dots, N_b\}$.

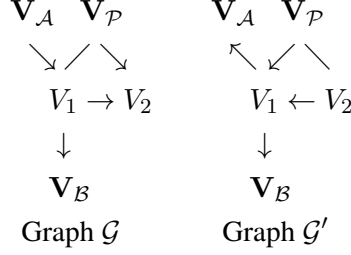
Since $\{\mathbf{c}_1^\ell\}_{\ell=1}^C$ are distinct, $\{\mathbf{c}_{1*}^\ell\}_{\ell=1}^C$ are distinct finite polynomials in $\mathbf{y}_{-1} \in \mathbb{R}^{m-1}$. For each pair of such distinct polynomials, the lemma of Okamoto (1973) implies that they only agree at a Lebesgue-null subset of \mathbb{R}^{n-1} , 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 (10) 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

$$\sum_{j=1}^{n_\ell^a} a_0^{k_{\ell j}^a} f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{a}_{-1}^{k_{\ell j}^a}) = \sum_{j=1}^{n_\ell^b} -b_0^{k_{\ell j}^b} f^{(m-1)}(\mathbf{y}_{-1}; \mathbf{b}_{-1}^{k_{\ell j}^b}) \quad (11)$$

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 (11) has its unique coefficients, and similarly for the right-hand side. Since (11) 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 π 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 ℓ , $N_a = \sum_{\ell=1}^C n_\ell^a = \sum_{\ell=1}^C n_\ell^b = N_b$, and we have thus again matched each \mathbf{a}^ℓ with a \mathbf{b}^ℓ as well as the corresponding a_0 's with b_0 's. This ends the proof for m , and the entire proof. \blacksquare

Proof [Proof of Theorem 6] 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 8, 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 8 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 (α, β, k) -parametrization only, since the result for the (p, μ, σ^2) naturally follows from their relationship (5). 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}_A , \mathbf{V}_B and \mathbf{V}_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 α'_2 and β'_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 α'_U and β'_U to 0. Since the two joint distributions derived from both graphs must be proportional to each other, we get for $y_2 \neq 0$

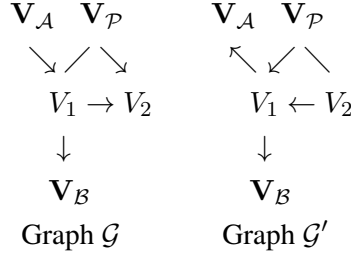
$$\exp[y_2\{\beta'_2(\mathbf{0}) - \beta_2(\mathbf{0})\} - (k'_2 - k_2)y_2^2/2]$$

$$\propto \prod_{U \in \mathcal{P} \cup \{1, 2\} \in \text{pa}_{\mathcal{G}'}(U)} \left[\sqrt{2\pi/k'_U} \exp\{\alpha'_U(y_2, \mathbf{0}) + \beta'_U(y_2, \mathbf{0})^2 / (2k'_U)\} + 1 \right]. \quad (12)$$

Note that $2 \in \text{pa}_{\mathcal{G}'}(1)$ and thus the product on the right of (12) 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 (12) 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 12, 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. \blacksquare

Proof [Proof of Theorem 10] As in the proof of Theorem 6 using Proposition 8, 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 8 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:



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. λ . Here f_1 and f'_2 are continuous functions in \mathbb{R} with no additive constant term, and δ_1 and δ'_2 are constants.

We prove the results in the (α, β, k) -parameterization only, since results for the (p, μ, σ^2) -parameterization would follow from their relationship (5). In our model, we assumed the α and β parameters for each node to be polynomial in the parents and their indicators. We also assumed that for each node, either the β function is nonconstant in any of the parents, or α depends on the value of all of its parents.

Consider a generic β 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 β is nonconstant in any of \mathcal{P} , and write $\beta(\mathbf{y}_{\mathcal{P}})$ equivalently as $\beta(\mathbf{y}_{p_0}, \mathbf{y}_{\mathcal{P}_1})$. 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 β_{-1} through β_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 β_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

$$\mathcal{Y}_{\beta, p_0, \mathcal{P}_1} \equiv \left\{ \mathbf{y}_{\mathcal{P}_1} \in \mathbb{R}^{|\mathcal{P}_1|} : \beta(y_{p_0}, \mathbf{y}_{\mathcal{P}_1}) \text{ nonconstant function in } y_{p_0} \right\}.$$

Thus $\mathbb{R}^{|\mathcal{P}_1|} \setminus \mathcal{Y}_{\beta, p_0, \mathcal{P}_1}$ has zero Lebesgue measure assuming β is nonconstant in its any of \mathcal{P} . Hence, by a similar argument, under the assumptions of the theorem, letting

$$\mathcal{Y}_{\alpha, \beta, p_0, \mathcal{P}_1} \equiv \left\{ \mathbf{y}_{\mathcal{P}_1} \in \mathbb{R}^{|\mathcal{P}_1|} : \beta(y_{p_0}, \mathbf{y}_{\mathcal{P}_1}) \text{ nonconstant function in } y_{p_0} \text{ or } \right. \\ \left. \alpha(y_{p_0}, \mathbf{y}_{\mathcal{P}_1}) \text{ depends on the value of } y_{p_0} \right\},$$

the set $\mathbb{R}^{|\mathcal{P}_1|} \setminus \mathcal{Y}_{\alpha, \beta, p_0, \mathcal{P}_1}$ 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 α and β 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 α and β functions on \mathcal{P} , and write them as scalar functions instead notation-wise. By discussion above, β_2 becomes a nonconstant function in y_1 and β'_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 α'_1 , α_2 , β'_1 and β_2 are automatically univariate functions, with β'_1 and β_2 nonconstant by assumption.

The joint density of $P(y_1, y_2 | \mathbf{Y}_{\mathcal{P}} = \mathbf{y}_{\mathcal{P}})$ w.r.t. λ thus has two characterizations (up to normalizing constants)

$$\frac{\exp\{\mathbb{1}_{y_1} \delta_1 + f_1(y_1) + \mathbb{1}_{y_2} \alpha_2(y_1) + y_2 \beta_2(y_1) - y_2^2 k_2 / 2\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2 / (2k_2)\} + 1} \\ \propto \frac{\exp\{\mathbb{1}_{y_2} \delta'_2 + f'_2(y_2) + \mathbb{1}_{y_1} \alpha'_1(y_2) + y_1 \beta'_1(y_2) - y_1^2 k'_1 / 2\}}{\sqrt{2\pi/k'_1} \exp\{\alpha'_1(y_2) + \{\beta'_1(y_2)\}^2 / (2k'_1)\} + 1}, \quad (13)$$

where $\alpha_2(y_1)$ has the form $c_{\alpha_2, -1} + c_{\alpha_2, 0} \mathbb{1}_{y_1} + c_{\alpha_2, 1} y_1 + \dots + c_{\alpha_2, k} y_1^k$ 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 (13) becomes

$$\frac{\exp\{f_1(y_1) + \alpha_2(y_1) + y_2 \beta_2(y_1) - y_2^2 k_2 / 2\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2 / (2k_2)\} + 1} \mathbb{1}_{y_1} \mathbb{1}_{y_2}, \\ \propto \frac{\exp\{f'_2(y_2) + \alpha'_1(y_2) + y_1 \beta'_1(y_2) - y_1^2 k'_1 / 2\}}{\sqrt{2\pi/k'_1} \exp\{\alpha'_1(y_2) + (\beta'_1(y_2))^2 / (2k'_1)\} + 1} \mathbb{1}_{y_1} \mathbb{1}_{y_2}, \quad (14)$$

for all $(y_1, y_2) \in (\mathbb{R} \setminus \{0\})^2$. (14) 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,

$$\begin{aligned} & \exp\{f_1(y_1) + P_1(y_1, y_2) + P_4(y_2)\} + \exp\{f_1(y_1) + P_1(y_1, y_2)\} \\ & = \exp\{f'_2(y_2) + P_3(y_1, y_2) + P_2(y_1)\} + \exp\{f'_2(y_2) + P_3(y_1, y_2)\}. \end{aligned} \quad (15)$$

Differentiating both sides of (15) with respect to y_1 ,

$$\begin{aligned} & \left[\frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2)\} \right] \exp \{f_1(y_1) + P_1(y_1, y_2) + P_4(y_2)\} \\ & \quad + \exp\{f_1(y_1) + P_1(y_1, y_2)\} \\ & = \left[\frac{\partial}{\partial y_1} \{P_3(y_1, y_2) + P_2(y_1)\} \right] \exp \{f'_2(y_2) + P_3(y_1, y_2) + P_2(y_1)\} \\ & \quad + \left\{ \frac{\partial}{\partial y_1} P_3(y_1, y_2) \right\} \exp \{f'_2(y_2) + P_3(y_1, y_2)\}. \end{aligned} \quad (16)$$

Plugging (15) into the left-hand side of (16),

$$\begin{aligned} & \left[\frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2)\} \right] \left[\exp \{f'_2(y_2) + P_3(y_1, y_2) + P_2(y_1)\} \right. \\ & \quad \left. + \exp \{f'_2(y_2) + P_3(y_1, y_2)\} \right] \\ & = \left[\frac{\partial}{\partial y_1} \{P_3(y_1, y_2) + P_2(y_1)\} \right] \exp \{f'_2(y_2) + P_3(y_1, y_2) + P_2(y_1)\} \\ & \quad + \left\{ \frac{\partial}{\partial y_1} P_3(y_1, y_2) \right\} \exp \{f'_2(y_2) + P_3(y_1, y_2)\}, \end{aligned}$$

which simplifies to

$$\begin{aligned} & \left[\frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2) - P_3(y_1, y_2) - P_2(y_1)\} \right] \\ & \quad \times \exp \{f'_2(y_2) + P_3(y_1, y_2) + P_2(y_1)\} \\ & \quad + \left[\frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2) - P_3(y_1, y_2)\} \right] \\ & \quad \times \exp \{f'_2(y_2) + P_3(y_1, y_2)\} = 0. \end{aligned}$$

Since $\exp \{f'_2(y_2) + P_3(y_1, y_2)\} \neq 0$, this becomes

$$\begin{aligned} & \left[\frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2) - P_3(y_1, y_2) - P_2(y_1)\} \right] \exp \{P_2(y_1)\} \\ & \quad + \left[\frac{\partial}{\partial y_1} \{f_1(y_1) + P_1(y_1, y_2) - P_3(y_1, y_2)\} \right] = 0. \end{aligned} \quad (17)$$

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

$$P_1(y_1, y_2) - P_3(y_1, y_2) = \alpha_2(y_1) + y_2\beta_2(y_1) - y_2^2k_2/2 - \alpha'_1(y_2) - y_1\beta'_1(y_2) + y_1^2k'_1/2. \quad (18)$$

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 β_2 and β'_1 are both linear with the same coefficient on the linear term. Now that β_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 β_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 (18)

$$\begin{aligned} &P_1(y_1, y_2) - P_3(y_1, y_2) \\ &= (\alpha_2(y_1) - y_1\beta'_{1;-1,0} + y_1^2k'_1/2) - (\alpha'_1(y_2) - y_2\beta_{2;-1,0} + y_2^2k_2/2) \\ &\equiv P_{1,3}(y_1) - (\text{function in } y_2 \text{ only}). \end{aligned}$$

Plugging this into (17), 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 $y_1 \neq 0$,

$$\begin{aligned} \exp(f_1(y_1)) &\propto \frac{1 + \exp\{P_2(y_1)\}}{\exp\{P_{1,3}(y_1)\}} \\ &= \frac{1 + \sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\}}{\exp\{\alpha_2(y_1) - \beta'_{1;-1,0}y_1 + y_1^2k'_1/2\}} \\ &= \exp\{-\alpha_2(y_1) + y_1\beta'_{1;-1,0} - y_1^2k'_1/2\} \\ &\quad + \sqrt{2\pi/k_2} \exp\{y_1\beta'_{1;-1,0} + \beta_2(y_1)^2/(2k_2) - y_1^2k'_1/2\}. \quad (19) \end{aligned}$$

Now condition on the event $\mathbb{1}_{y_1} = 1$ and $\mathbb{1}_{y_2} = 0$. Then (13) becomes

$$\frac{\exp\{f_1(y_1)\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \mathbb{1}_{y_1} \propto \exp\{y_1\beta'_1(0) - y_1^2 k'_1/2\} \mathbb{1}_{y_1},$$

which implies that for $y_1 \neq 0$,

$$\begin{aligned} \exp\{f_1(y_1)\} &\propto \exp\{y_1\beta'_1(0) - y_1^2 k'_1/2\} \\ &\quad + \sqrt{2\pi/k_2} \exp\{y_1\beta'_1(0) - y_1^2 k'_1/2 + \alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\}. \end{aligned} \quad (20)$$

Applying Lemma 12 to (19) and (20), by matching the terms we have (conditional on $y_1 \neq 0$) either

$$-\alpha_2(y_1) + y_1\beta'_{1,-1,0} = y_1\beta'_1(0) + \text{const}; \quad \text{or} \quad (21)$$

$$-\alpha_2(y_1) + y_1\beta'_{1,-1,0} = y_1\beta'_1(0) + \alpha_2(y_1) + \beta_2(y_1)^2/(2k_2) + \text{const} \quad \text{and}$$

$$y_1\beta'_{1,-1,0} + \beta_2(y_1)^2/(2k_2) = y_1\beta'_1(0) + \text{const}. \quad (22)$$

Conditional on $y_1 \neq 0$, in the first case (21), $\alpha_2(y_1) = y_1 c_{\beta'_1,0} + \text{const}$; in the second case (22), $\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2) = \text{const}$ and $\beta_2(y_1)^2/(2k_2) = -y_1 c_{\beta'_1,0} + \text{const}$, 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 (21). Thus, in either case, $\alpha_2(y_1) = c_{\alpha_2,0} \mathbb{1}_{y_1} + y_1 c_{\beta'_1,0} + \text{const}$, i.e. α_2 is linear (or constant) in $y_1 \neq 0$ with coefficient on y_1 equal to $c_{\beta'_1,0}$. By (21) for $y_1 \neq 0$,

$$\begin{aligned} \exp\{f_1(y_1)\} &\propto \exp\{y_1\beta'_1(0) - y_1^2 k'_1/2\} \\ &\quad + \sqrt{2\pi/k_2} \exp\{y_1\beta'_{1,-1,0} + \beta_2(y_1)^2/(2k_2) - y_1^2 k'_1/2\}, \end{aligned} \quad (23)$$

clearly a single univariate Gaussian or a mixture of two univariate Gaussian distributions (since β_2 is at most linear in y_1). Similarly, we must have $\alpha'_1(y_2) = y_2\beta_{2,-1,0} - y_2\beta_2(0) + \text{const} = y_2 c_{\beta_{2,0}} + \text{const}$ for $y_2 \neq 0$, and for $y_2 \neq 0$

$$\begin{aligned} \exp\{f'_2(y_2)\} &\propto \exp\{y_2\beta_2(0) - y_2^2 k_2/2\} \\ &\quad + \sqrt{2\pi/k'_1} \exp\{y_2\beta_{2,-1,0} + \beta'_1(y_2)^2/(2k'_1) - y_2^2 k_2/2\}. \end{aligned} \quad (24)$$

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 (23), $\beta'_1(0) = \beta'_{1,-1,0}$ and $\beta_2(y_1)$ is a constant given $\mathbb{1}_{y_1}$, 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 (13) 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

$$\begin{aligned} \exp\{f'_2(y_2)\} &\propto \frac{\exp\{f_1(0) + \alpha_2(0) + y_2\beta_2(0) - y_2^2 k_2/2\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(0) + \beta_2(0)^2/(2k_2)\} + 1} \\ &\quad + \exp\{y_2\beta_{2,-1,0} - y_2^2 k_2/2\} \int_{\mathbb{R}} \frac{\exp\{\delta_1 + f_1(y_1) + \alpha_2(y_1)\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} dy_1, \end{aligned}$$

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 (24) 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 β'_1 is an absolute constant in y_2 and $\mathbb{1}_{y_2}$, and also that α_2 may depend on y_1 only through $\mathbb{1}_{y_1}$, a contradiction to the assumption of the theorem.

Thus, (23) and (24) 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}' . ■

Proof [Proof of Corollary 11] When $|\mathcal{V}| = 2$, in Proposition 8 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 8 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 13 *In the proof of Theorem 10, 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 (24) in the proof must hold. Specifically, conditioning on almost every $\mathbf{y}_{\mathcal{P}}$, α_2 and β_2 in \mathcal{G} as well as α'_1 and β'_1 in \mathcal{G}' can be at most linear in y_1 and y_2 , respectively, namely*

$$\begin{aligned} \beta'_1(y_2) &= c_{\beta'_1,-1} + c_{\beta'_1,0}\mathbb{1}_{y_2} + c_{\beta'_1,1}y_2, & \beta_2(y_1) &= c_{\beta_2,-1} + c_{\beta_2,0}\mathbb{1}_{y_1} + c_{\beta_2,1}y_1, \\ \alpha'_1(y_2) &= c_{\alpha'_1,-1} + c_{\alpha'_1,0}\mathbb{1}_{y_2} + c_{\alpha'_1,1}y_2, & \alpha_2(y_1) &= c_{\alpha_2,-1} + c_{\alpha_2,0}\mathbb{1}_{y_1} + c_{\alpha_2,1}y_1, \end{aligned}$$

with coefficients depending on $\mathbf{y}_{\mathcal{P}}$ where

$$c_{\alpha'_1,1} = c_{\beta_2,0}, \quad c_{\alpha_2,1} = c_{\beta'_1,0}, \quad c_{\beta'_1,1} = c_{\beta_2,1}. \quad (25)$$

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 [Proof of Remark 13]

By (13), (23), (24), the joint distribution of Y_1 and Y_2 conditional on $\mathbf{Y}_{\mathcal{P}}$ has two characterizations (up to normalizing constants)

$$\begin{aligned}
 & \frac{\exp\{\mathbb{1}_{y_1}\delta_1 + y_1\beta_1'(0) - y_1^2k_1'/2 + \mathbb{1}_{y_2}\alpha_2(y_1) + y_2\beta_2(y_1) - y_2^2k_2/2\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
 & + \frac{\sqrt{2\pi/k_2} \exp\{\mathbb{1}_{y_1}\delta_1 + y_1\beta_{1;-1,0}' + \beta_2(y_1)^2/(2k_2) - y_1^2k_1'/2 + \mathbb{1}_{y_2}\alpha_2(y_1) + y_2\beta_2(y_1) - y_2^2k_2/2\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
 & \propto \frac{\exp\{\mathbb{1}_{y_2}\delta_2' + y_2\beta_2(0) - y_2^2k_2/2 + \mathbb{1}_{y_1}\alpha_1'(y_2) + y_1\beta_1'(y_2) - y_1^2k_1'/2\}}{\sqrt{2\pi/k_1'} \exp\{\alpha_1'(y_2) + \beta_1'(y_2)^2/(2k_1')\} + 1} \\
 & + \frac{\sqrt{2\pi/k_1'} \exp\{\mathbb{1}_{y_2}\delta_2' + y_2\beta_{2;-1,0}' + \beta_1'(y_2)^2/(2k_1') - y_2^2k_2/2 + \mathbb{1}_{y_1}\alpha_1'(y_2) + y_1\beta_1'(y_2) - y_1^2k_1'/2\}}{\sqrt{2\pi/k_1'} \exp\{\alpha_1'(y_2) + \beta_1'(y_2)^2/(2k_1')\} + 1}.
 \end{aligned} \tag{26}$$

Divide both sides by $\exp(y_1\beta_1'(0) + y_2\beta_2(0) - y_1^2k_1'/2 - y_2^2k_2/2)$ and expanding $\beta_1'(y_2)$ and $\beta_2(y_1)$, this becomes

$$\begin{aligned}
 & \frac{\exp\{\mathbb{1}_{y_1}\delta_1 + \mathbb{1}_{y_2}\alpha_2(y_1) + y_2c_{\beta_2,0}\mathbb{1}_{y_1} + y_1y_2c_{\beta_2,1}\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
 & + \frac{\sqrt{2\pi/k_2} \exp\{\mathbb{1}_{y_1}\delta_1 + y_1c_{\beta_1',0} + \beta_2(y_1)^2/(2k_2) + \mathbb{1}_{y_2}\alpha_2(y_1) + y_2c_{\beta_2,0}\mathbb{1}_{y_1} + y_1y_2c_{\beta_2,1}\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
 & \propto \frac{\exp\{\mathbb{1}_{y_2}\delta_2' + \mathbb{1}_{y_1}\alpha_1'(y_2) + y_1c_{\beta_1',0}\mathbb{1}_{y_2} + y_1y_2c_{\beta_1',1}\}}{\sqrt{2\pi/k_1'} \exp\{\alpha_1'(y_2) + \beta_1'(y_2)^2/(2k_1')\} + 1} \\
 & + \frac{\sqrt{2\pi/k_1'} \exp\{\mathbb{1}_{y_2}\delta_2' + y_2c_{\beta_2,0} + \beta_1'(y_2)^2/(2k_1') + \mathbb{1}_{y_1}\alpha_1'(y_2) + y_1c_{\beta_1',0}\mathbb{1}_{y_2} + y_1y_2c_{\beta_1',1}\}}{\sqrt{2\pi/k_1'} \exp\{\alpha_1'(y_2) + \beta_1'(y_2)^2/(2k_1')\} + 1}.
 \end{aligned}$$

Now expanding $\alpha_1'(y_2)$ and $\alpha_2(y_1)$ and using the relationships in (25), we divide both sides by $\exp(y_1c_{\alpha_2,1}\mathbb{1}_{y_2} + y_2c_{\beta_2,0}\mathbb{1}_{y_1} + y_1y_2c_{\beta_2,1}) = \exp(y_1c_{\beta_1',0}\mathbb{1}_{y_2} + y_2c_{\alpha_1',1}\mathbb{1}_{y_1} + y_1y_2c_{\beta_2,1})$ and get

$$\begin{aligned}
 & \frac{\exp\{\mathbb{1}_{y_1}\delta_1 + \mathbb{1}_{y_2}(c_{\alpha_2,-1} + c_{\alpha_2,0}\mathbb{1}_{y_1})\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
 & + \frac{\sqrt{2\pi/k_2} \exp\left\{\mathbb{1}_{y_1}\delta_1 + y_1c_{\beta_1',0} + \beta_2(y_1)^2/(2k_2) + \mathbb{1}_{y_2}(c_{\alpha_2,-1} + c_{\alpha_2,0}\mathbb{1}_{y_1})\right\}}{\sqrt{2\pi/k_2} \exp\{\alpha_2(y_1) + \beta_2(y_1)^2/(2k_2)\} + 1} \\
 & = C_0 \frac{\exp\{\mathbb{1}_{y_2}\delta_2' + \mathbb{1}_{y_1}(c_{\alpha_1',-1} + c_{\alpha_1',0}\mathbb{1}_{y_2})\}}{\sqrt{2\pi/k_1'} \exp\{\alpha_1'(y_2) + \beta_1'(y_2)^2/(2k_1')\} + 1} \\
 & + C_0 \frac{\sqrt{2\pi/k_1'} \exp\{\mathbb{1}_{y_2}\delta_2' + y_2c_{\beta_2,0} + \beta_1'(y_2)^2/(2k_1') + \mathbb{1}_{y_1}(c_{\alpha_1',-1} + c_{\alpha_1',0}\mathbb{1}_{y_2})\}}{\sqrt{2\pi/k_1'} \exp\{\alpha_1'(y_2) + \beta_1'(y_2)^2/(2k_1')\} + 1}
 \end{aligned} \tag{27}$$

for some C_0 . Setting $\mathbb{1}_{y_1} = \mathbb{1}_{y_2} = 0$ (27) becomes

$$\frac{1 + \sqrt{2\pi/k_2} \exp\{c_{\beta_2,-1}^2/(2k_2)\}}{\sqrt{2\pi/k_2} \exp\{c_{\alpha_2,-1} + c_{\beta_2,-1}^2/(2k_2)\} + 1} = C_0 \frac{1 + \sqrt{2\pi/k'_1} \exp\{c_{\beta'_{1,-1}}^2/(2k'_1)\}}{\sqrt{2\pi/k'_1} \exp\{c_{\alpha'_{1,-1}} + c_{\beta'_{1,-1}}^2/(2k'_1)\} + 1}, \quad (28)$$

and with $\mathbb{1}_{y_1} \neq 0, \mathbb{1}_{y_2} = 0$ (27) becomes

$$\begin{aligned} \exp(\delta_1) \frac{1 + \sqrt{2\pi/k_2} \exp\{y_1 c_{\beta'_{1,0}} + \beta_2(y_1)^2/(2k_2)\}}{\sqrt{2\pi/k_2} \exp\{c_{\alpha_2,-1} + c_{\alpha_2,0} + c_{\alpha_2,1}y_1 + \beta_2(y_1)^2/(2k_2)\} + 1} \\ = C_0 \exp(c_{\alpha'_{1,-1}}) \frac{1 + \sqrt{2\pi/k'_1} \exp\{c_{\beta'_{1,-1}}^2/(2k'_1)\}}{\sqrt{2\pi/k'_1} \exp\{c_{\alpha'_{1,-1}} + c_{\beta'_{1,-1}}^2/(2k'_1)\} + 1}. \end{aligned} \quad (29)$$

Since the right-hand side of (29) is a constant, by matching the numerator and the denominator of the left-hand side using Lemma 12, 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 β'_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

$$\alpha_{2,-1,0} \equiv c_{\alpha_2,-1} + c_{\alpha_2,0} = 0, \quad \text{and by symmetry} \quad \alpha'_{1,-1,0} \equiv c_{\alpha'_{1,-1}} + c_{\alpha'_{1,0}} = 0. \quad (30)$$

Thus the left-hand side of (29) is just $\exp(\delta_1)$. Note that the right-hand side of (29) is $\exp(c_{\alpha'_{1,-1}})$ times the right-hand side of (28). So by equating the left-hand side of (29) with $\exp(c_{\alpha'_{1,-1}})$ times the left-hand side of (28) we have

$$\exp(\delta_1) = \exp(c'_{\alpha_1,-1}) \frac{1 + \sqrt{2\pi/k_2} \exp\{c_{\beta_2,-1}^2/(2k_2)\}}{\sqrt{2\pi/k_2} \exp\{c_{\alpha_2,-1} + c_{\beta_2,-1}^2/(2k_2)\} + 1} \quad (31)$$

and similarly

$$\exp(\delta'_2) = \exp(c_{\alpha_2,-1}) \frac{1 + \sqrt{2\pi/k'_1} \exp\{c_{\beta'_{1,-1}}^2/(2k'_1)\}}{\sqrt{2\pi/k'_1} \exp\{c_{\alpha'_{1,-1}} + c_{\beta'_{1,-1}}^2/(2k'_1)\} + 1}. \quad (32)$$

Now by (30), with $\mathbb{1}_{y_1} = \mathbb{1}_{y_2} = 1$, (27) simplifies to $\exp(\delta_1) = C_0 \cdot \exp(\delta'_2)$. Thus by (28), (31) and (32), 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 (30), we get

$$c_{\alpha'_{1,-1}} = c_{\alpha_2,-1}, \quad c_{\alpha'_{1,0}} = c_{\alpha_2,0}, \quad c_{\alpha'_{1,-1}} + c_{\alpha'_{1,0}} = c_{\alpha_2,-1} + c_{\alpha_2,0} = 0. \quad (33)$$

Note that this result holds as long as we assume identifiability does not hold. \blacksquare

B. Results of Additional Numerical Experiments

Here we provide additional details and results of numerical experiments.

Figure 4 shows the true DAG structures used in the simulation studies in Section 5.

In Figure S5, 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, μ, 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. It is worth noting 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.

B.1. Additional Results for Exhaustive Search

Recall that we consider the following DAG structures: (i) chain graph with $m = 10$, (ii) complete graph with $m = 10$, (iii) lattice graph with $m = 9$; see Figure 4 in B for an illustration of the DAG structures.

The results for correctly specified models are shown in Figures S6–S8. Each figure has one true underlying DAG from those mentioned above. In all figures, each row indicates one choice of true data generating parametrization— (α, β, k) -linear, and (p, μ, σ^2) -linear and quadratic—and each column shows the results using each estimating parametrization. Thus, plots on the diagonal (with

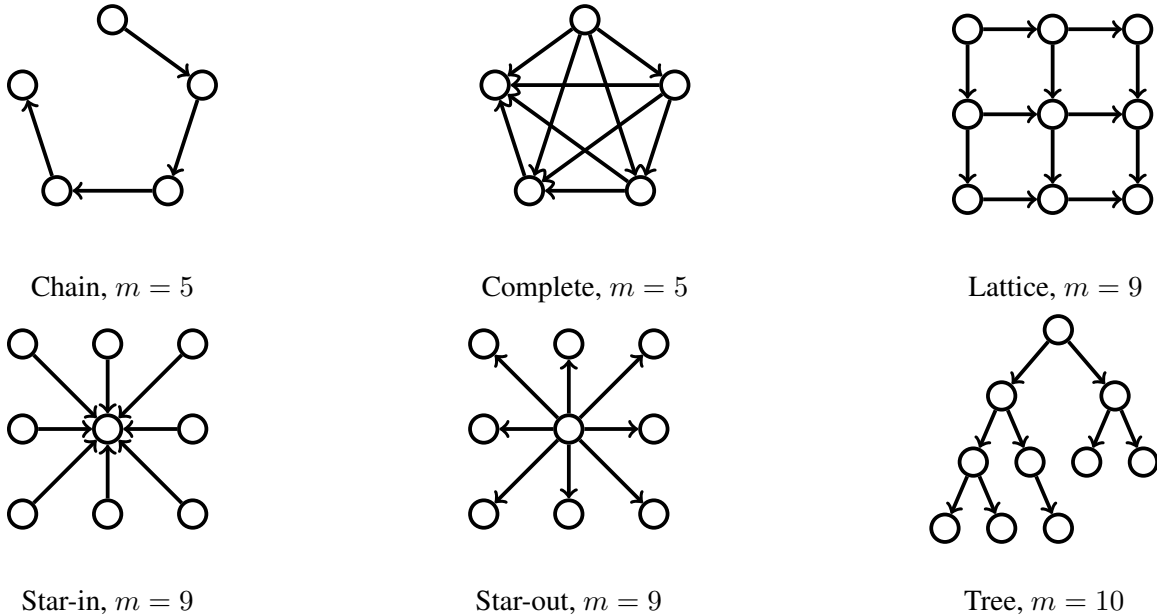


Figure 4: Example graph structures used in our experiments.

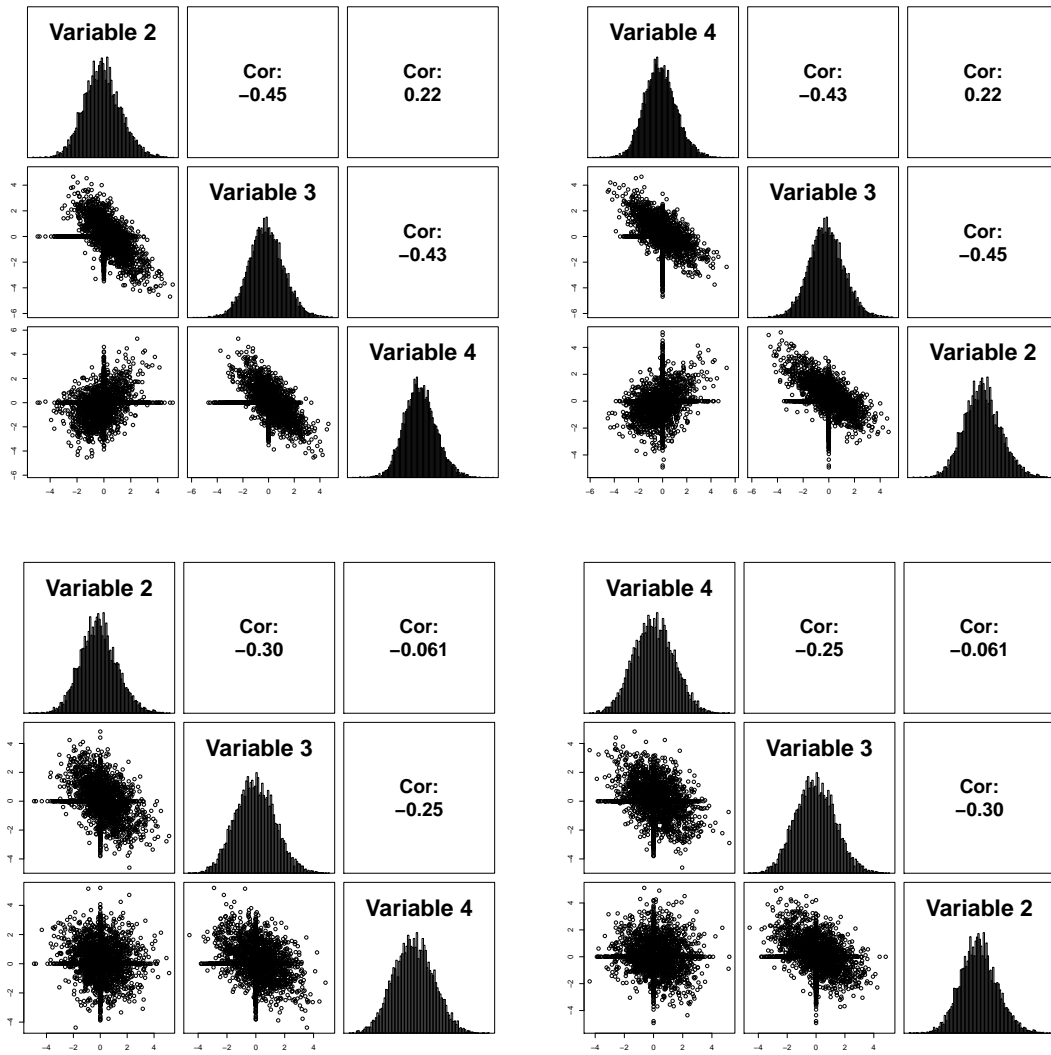


Figure S5: Pairwise scatterplots of zero-inflated data generated using chain graphs (upper row) and complete graphs (lower row), both with topological ordering $1 \rightarrow 2 \rightarrow 3 \rightarrow 4 \rightarrow 5$; only nodes 2, 3 and 4 are plotted. Plots on the left are plotted in the order 2, 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 methods can still determine the correct topological ordering.

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 results indeed indicate that in all settings, exhaustive search with correct parametrization almost always identifies the exact DAG for large n . Surprisingly, model misspecification does not seem to negatively impact the results by a significant amount. Overall, our simulation studies confirm the identifiability theory (Theorem 6).

Figure S9–S11 show exact recovery rates. In the plots, exact success rates are measured by the percentage of times (out of $B = 100$ iterations for each setting) that the exact DAG is recovered, whereas the equivalent success rates stand for the percentage of times the equivalence class of DAG is correctly identified, and are thus no less than the exact success rates.

B.2. Results for Greedy Search (GDS)

To evaluate the performance of greedy search we consider the following graphs: (i) chain graph with $m = 100$, (ii) complete graph with $m = 10$, (iii) lattice graph with $m = 100$, (iv) star-in graph with $m = 20$ ($j \rightarrow 1$ for $j = 2, \dots, m$), (v) star-out graph with $m = 100$ ($1 \rightarrow j$ for $j = 2, \dots, m$), (vi) tree graph with $m = 100$ ($j \rightarrow 2j$ for $j \leq \lfloor m/2 \rfloor$ and $j \rightarrow 2j + 1$ for $j \leq \lfloor (m - 1)/2 \rfloor$).

Results for the greedy search algorithm are shown in Figures S12 and S13, where each row corresponds to a different true graph, and each column corresponds to one of the three aforementioned parametrizations, where for simplicity we only present the results with correctly specified parametrization.

The results indicate that GDS works reasonably well in all settings but may require larger samples for recovering the structure of complete/ very dense graphs, or graphs with high in-degrees. 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 [Shah and Samworth \(2013\)](#) can further improve the GDS results.

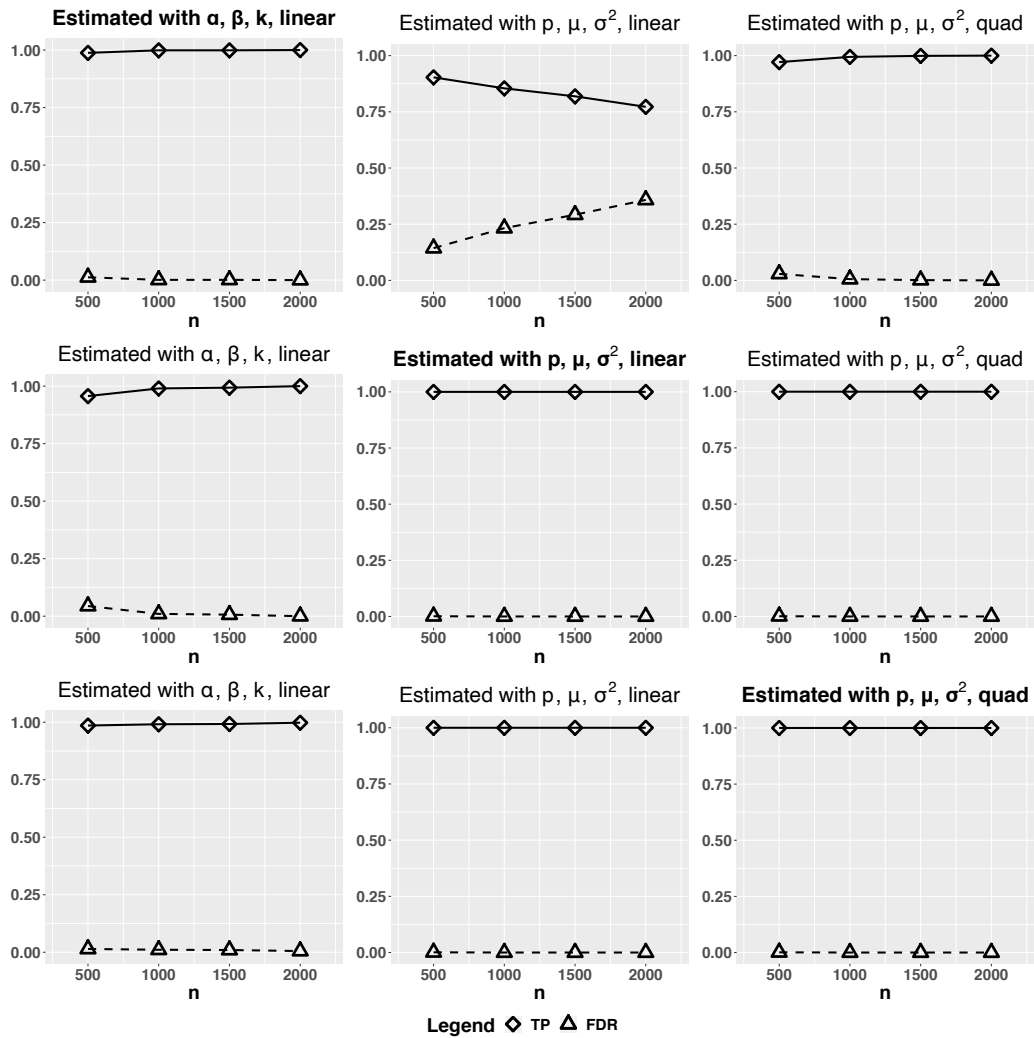


Figure S6: Chain graph, $m = 10$, exhaustive search. Each row corresponds to a different generating parametrization, and each column a different estimating parametrization. Generating and estimating parametrizations agree on the diagonal. ‘ \diamond ’ with solid lines: true positive rate; ‘ \triangle ’ with dashed lines: false discovery rate.

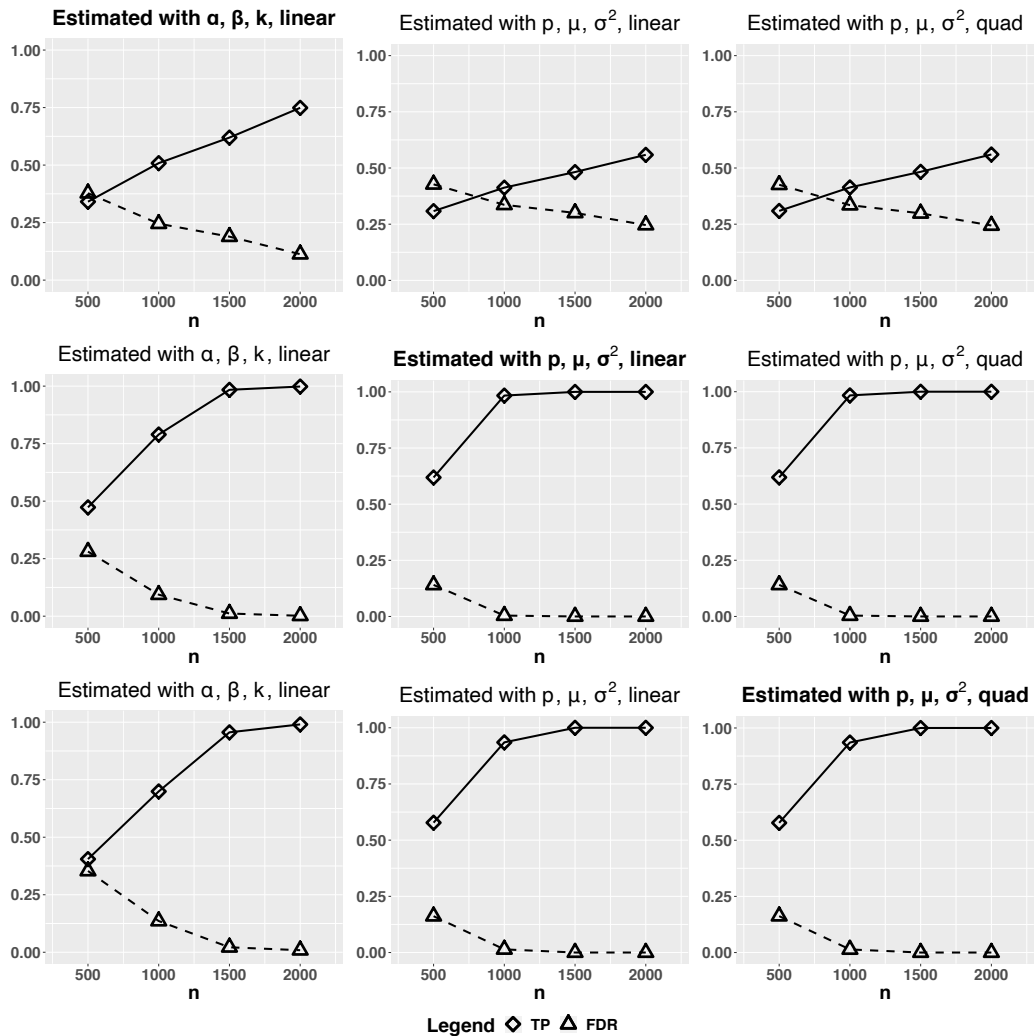


Figure S7: Complete graph, $m = 10$, exhaustive search. Each row corresponds to a different generating parametrization, and each column a different estimating parametrization. Generating and estimating parametrizations agree on the diagonal. ‘ \diamond ’ with solid lines: true positive rate; ‘ \triangle ’ with dashed lines: false discovery rate.

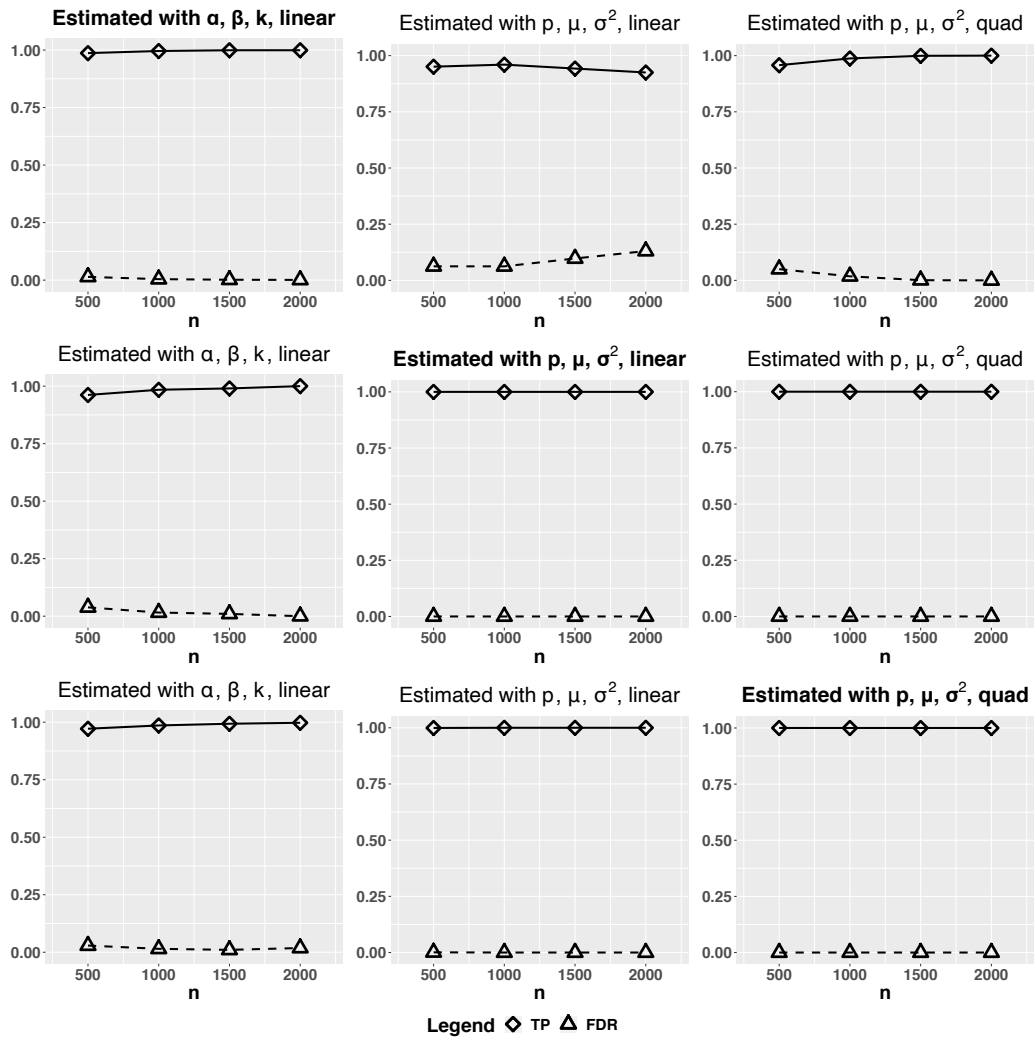


Figure S8: Lattice graph, $m = 9$, exhaustive search. Each row corresponds to a different generating parametrization, and each column a different estimating parametrization. Generating and estimating parametrizations agree on the diagonal. ‘ \diamond ’ with solid lines: true positive rate; ‘ Δ ’ with dashed lines: false discovery rate.

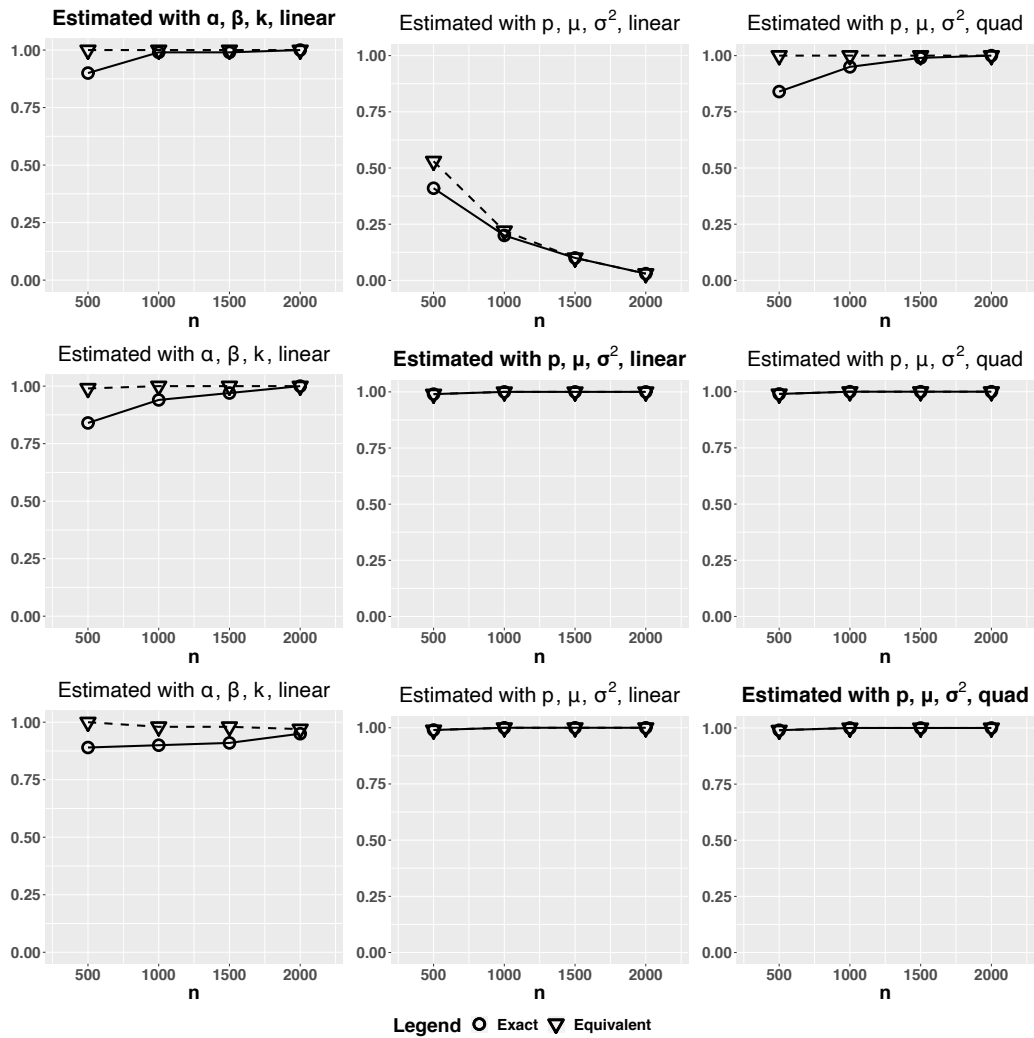


Figure S9: Chain graph, $m = 10$, exhaustive search. Each row corresponds to a different generating parametrization, and each column a different estimating parametrization. Generating and estimating parametrizations agree on the diagonal. ‘o’ with solid lines: success rates of exact DAG recovery; ‘∇’ with dashed lines: success rates for recovery of equivalence class.

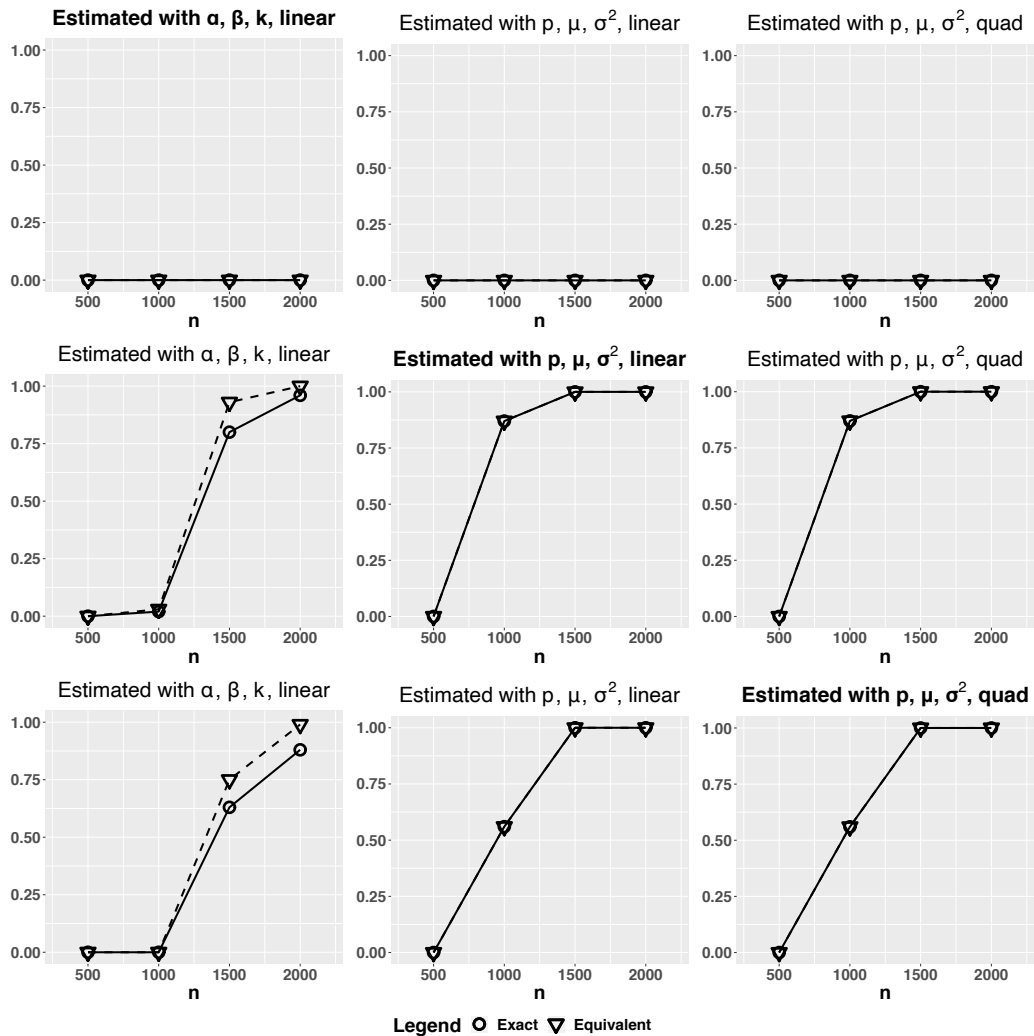


Figure S10: Complete graph, $m = 10$, exhaustive search. Each row corresponds to a different generating parametrization, and each column a different estimating parametrization. Generating and estimating parametrizations agree on the diagonal. ‘ \circ ’ with solid lines: success rates of exact DAG recovery; ‘ ∇ ’ with dashed lines: success rates for recovery of equivalence class.

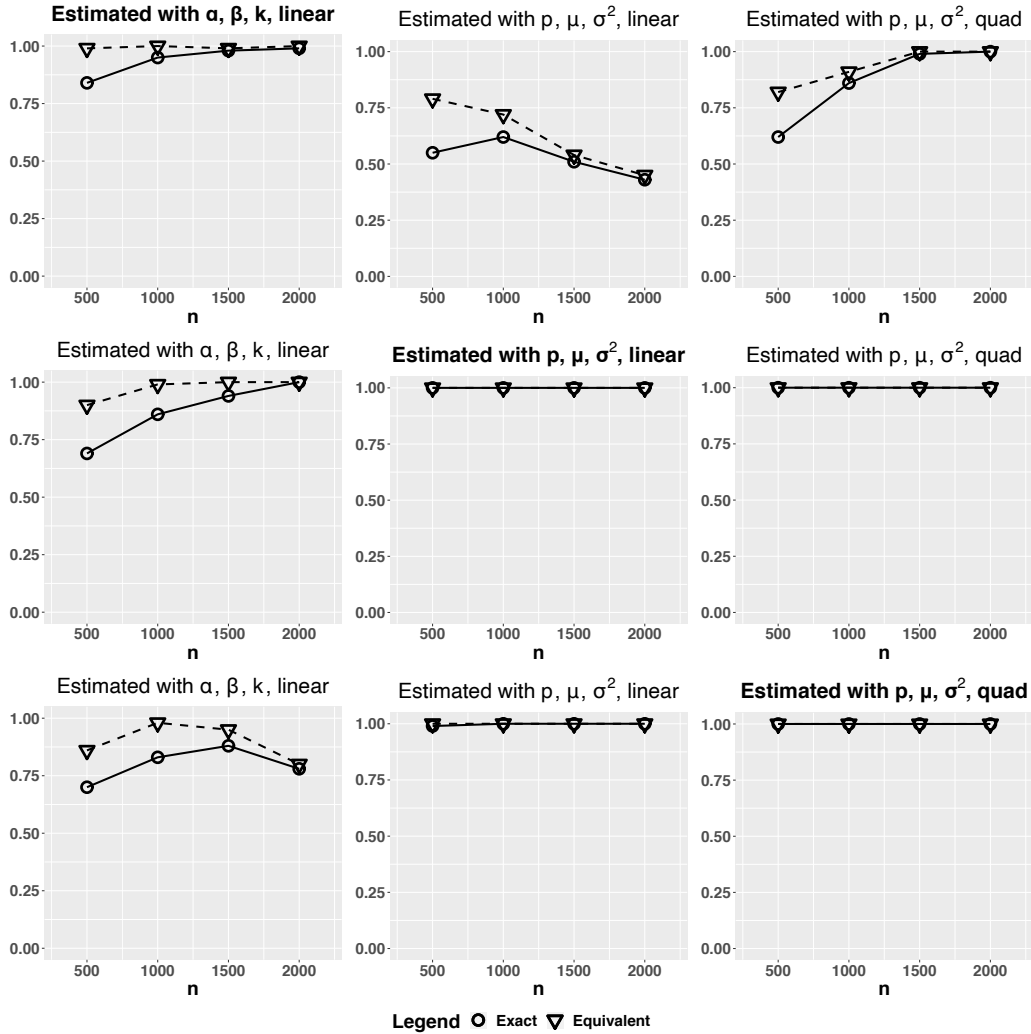


Figure S11: Lattice graph, $m = 9$, exhaustive search. Each row corresponds to a different generating parametrization, and each column a different estimating parametrization. Generating and estimating parametrizations agree on the diagonal. ‘o’ with solid lines: success rates of exact DAG recovery; ‘∇’ with dashed lines: success rates for recovery of equivalence class.

B.3. Details on Estimation of Connected Components

In this section we present simulation results validating our the strategy in Section 4.2, namely first applying the procedure of [McDavid et al. \(2019\)](#), then estimating the directed graphs inside each connected component of its estimated undirected graph. We measure the quality of the connected components (CC) estimated compared to the truth. The underlying true graph is a block-diagonal graph with $m = 100$ nodes, evenly divided into 10 connected components, where each connected component has the exact same setting as the complete graph with $m = 10$ nodes in previous sections. Specifically in each trial, we apply the method of [McDavid et al. \(2019\)](#), pick the estimate

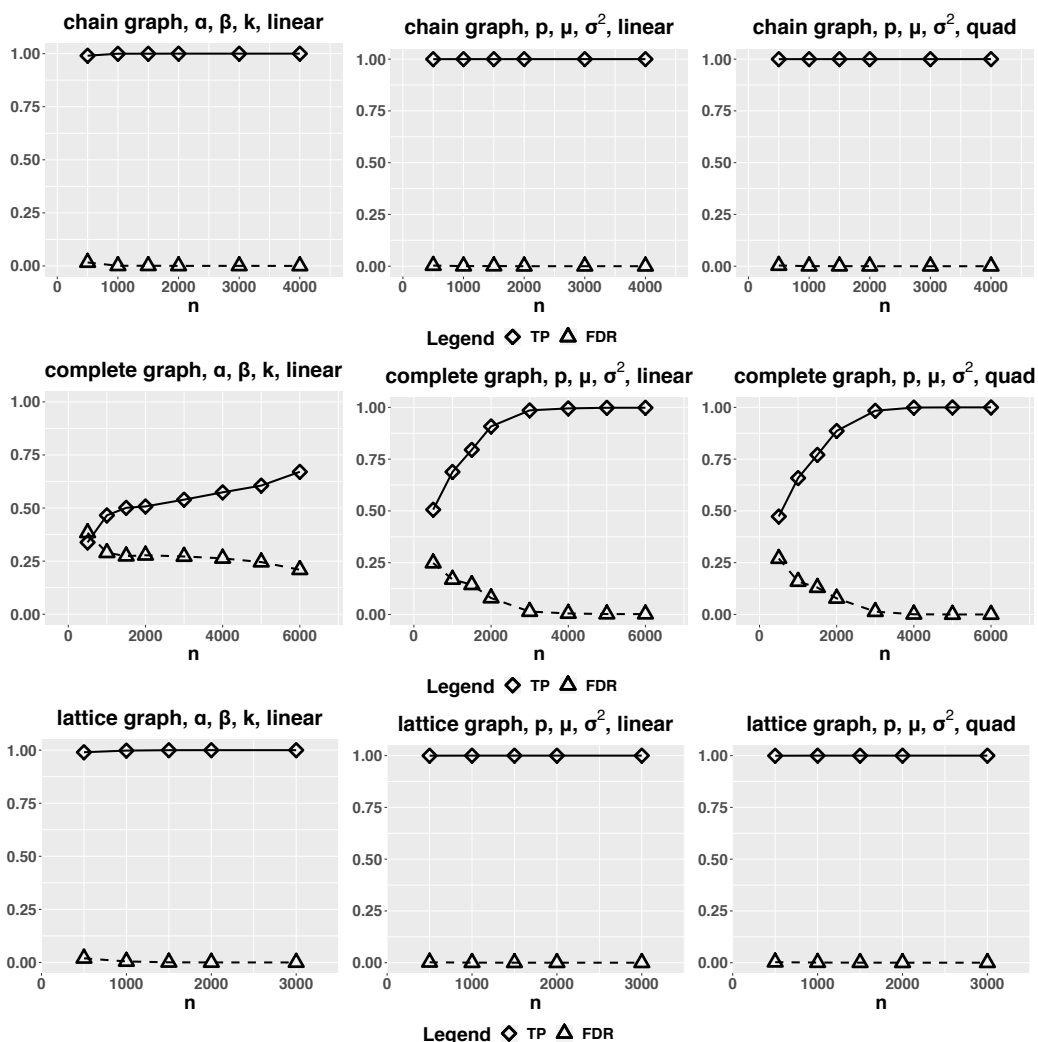


Figure S12: Results for GDS for chain graph with $m = 100$, complete graph with $m = 10$ and lattice graph with $m = 100$. Each row corresponds to a different graph structure, and each column corresponds to a different parametrization; the generating and estimating parametrizations are the same in the results. ‘ \diamond ’ with solid lines: true positive rate; ‘ \triangle ’ with dashed lines: false discovery rate.

that minimizes the BIC (with the exception of metric (b) below) and use the “and” rule to find the undirected graph (since the estimate is asymmetric due to their neighborhood selection method).

As in Section B.2, each column corresponds to a different parametrization (estimating parametrization = truth). As before all results shown are averaged over 100 trials. Each row contains two different metrics with value in $[0, 1]$, which we explain below.

- (a) “Subset/eq CCs”: A trial is count as successful if each true CC is a subset of some estimated CC, i.e. if any two truly connected nodes are estimated to belong to the same CC; the proportion of successful trials is reported.

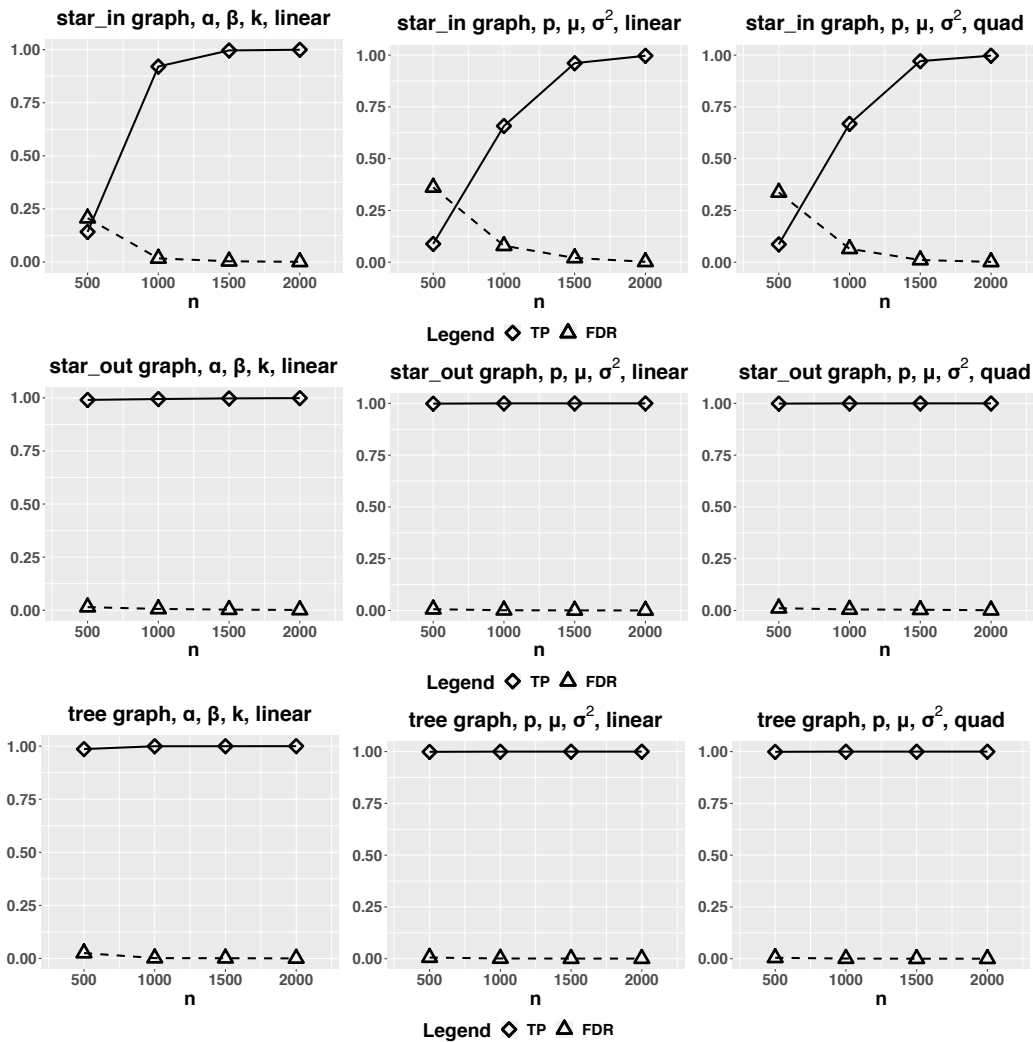


Figure S13: Results for GDS for star.in graph with $m = 20$, star.out graph with $m = 100$ and tree graph with $m = 100$. Each row corresponds to a different graph structure, and each column corresponds to a different parametrization; the generating and estimating parametrizations are the same in the results. ‘ \diamond ’ with solid lines: true positive rate; ‘ Δ ’ with dashed lines: false discovery rate.

- (b) “Correct CCs-Oracle”: From the solution path of [McDavid et al. \(2019\)](#) we take the graph assuming we know the true number of CCs. Then a trial counts as successful if the estimated CCs are exactly equal to the truth, and the proportion of successful trials is reported.
- (c) “Correct #CCs”: A trial is treated as successful if the estimated number of CCs is equal to the truth (10), and the proportion of successful trials is reported.
- (d) “Avg #CCs/10”: Average of the estimated numbers of CCs over 100 trials, divided by the truth (10).

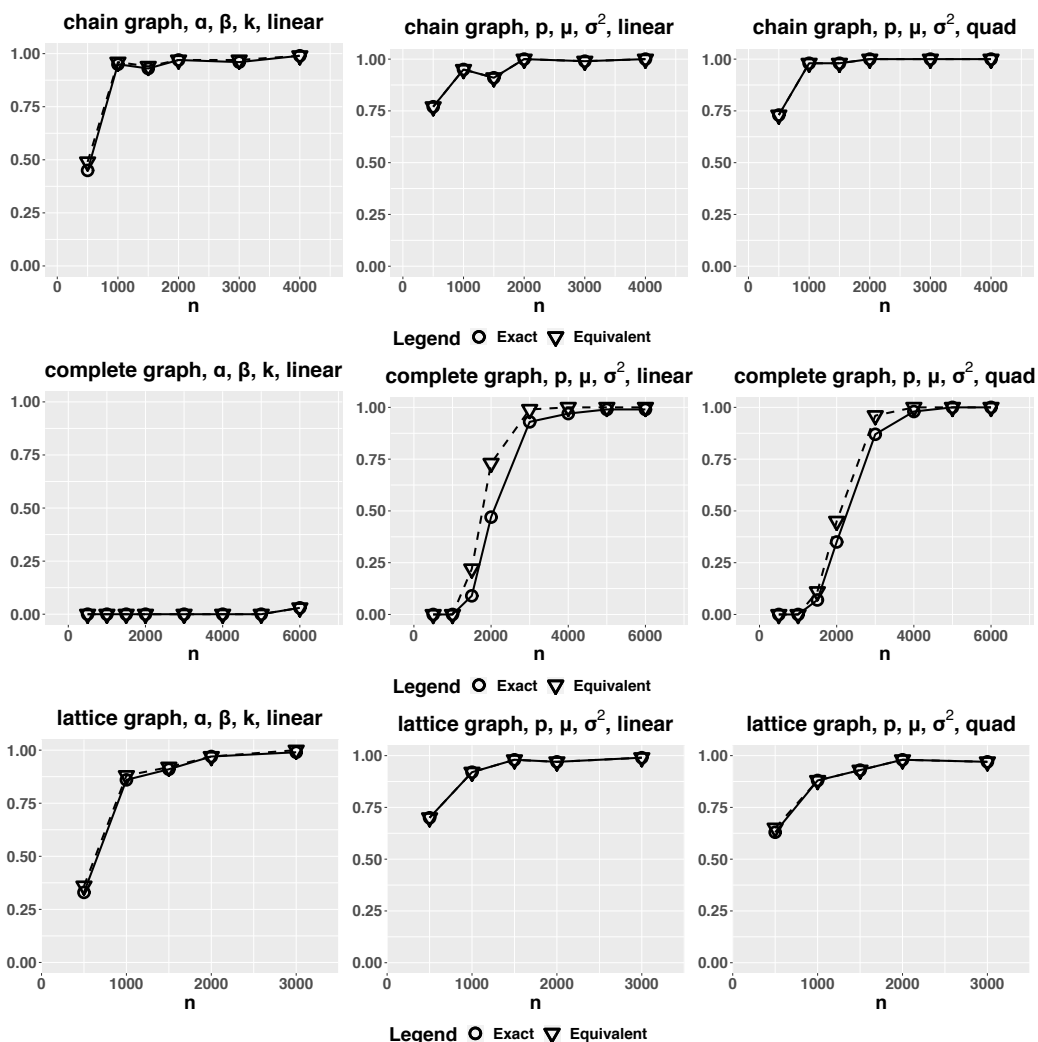


Figure S14: Results for GDS for chain graph with $m = 100$, complete graph with $m = 10$ and lattice graph with $m = 100$. Each row corresponds to a different graph structure, and each column corresponds to a different parametrization; the generating and estimating parametrizations are the same in the results. ‘o’ with solid lines: success rates of exact DAG recovery; ‘∇’ with dashed lines: success rates for recovery of equivalence class.

(e)&(f) “TP/FDR”: The true positive rate and false discovery rate for undirected graph recovery.

A high (a) metric guards against the mistake of failing to keep two truly connected nodes in the same CC, while (c) and (d) measure how many CCs the procedure actually generates, since the trivial case where all nodes form a single CC is undesirable. Note this characterizes the statistical versus computational trade-off discussed in Section 4.2. Metric (b), on the other hand, attempts to test if the procedure can become perfect had it known the true number of CCs. Metrics (e) and (f) provide additional information from the edge recovery perspective in terms of undirected graphs.

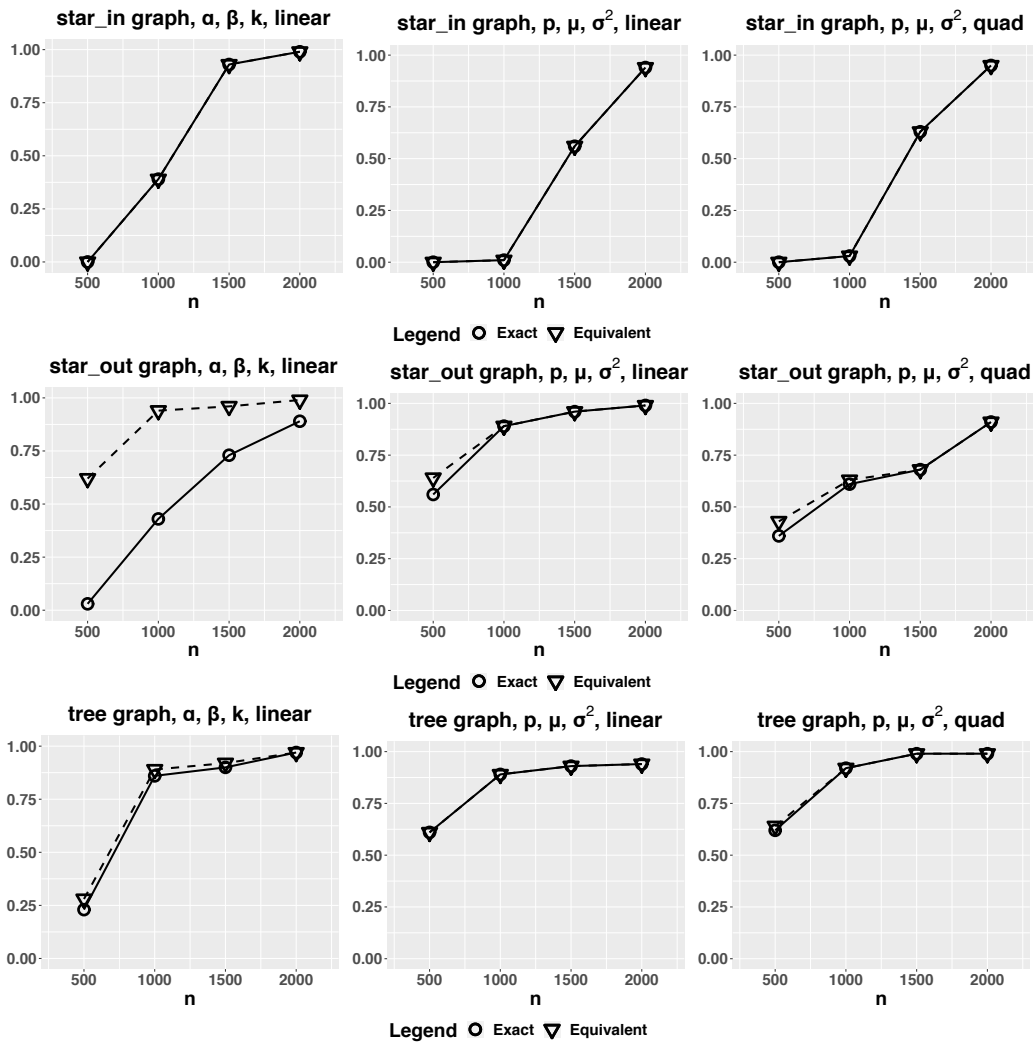


Figure S15: Results for GDS for star_in graph with $m = 20$, star_out graph with $m = 100$ and tree graph with $m = 100$. Each row corresponds to a different graph structure, and each column corresponds to a different parametrization; the generating and estimating parametrizations are the same in the results. ‘o’ with solid lines: success rates of exact DAG recovery; ‘∇’ with dashed lines: success rates for recovery of equivalence class.

Except for (f) which should be close to 0, one would hope for (a)–(e) to be close to 1, which indeed is the case, except for the number of CCs for the α, β, k parametrization.

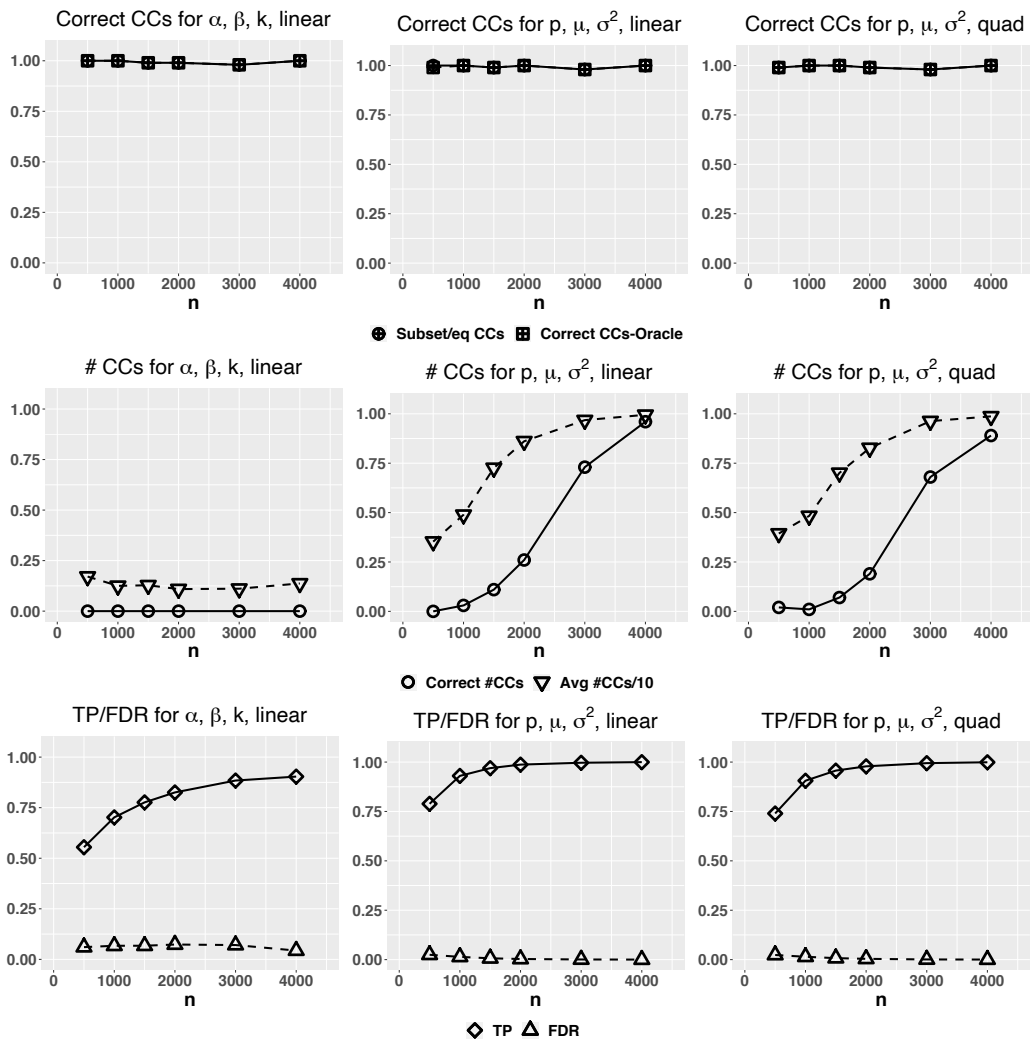


Figure S16: Results for the connected components (CC) estimated using the procedure in [McDavid et al. \(2019\)](#), compared against the truth (block-diagonal graph with 10 complete graphs each with $m = 10$). See discussion in Sections 4.2 and B.3.