

# Forest Density Estimation

**Han Liu**

*Department of Biostatistics and Computer Science  
Johns Hopkins University  
Baltimore, MD 21210, USA*

HANLIU@CS.JHU.EDU

**Min Xu**

**Haijie Gu**

**Anupam Gupta**

**John Lafferty\***

*School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213, USA*

MINX@CS.CMU.EDU  
HAIJIE@CS.CMU.EDU  
ANUPAMG@CS.CMU.EDU  
LAFFERTY@CS.CMU.EDU

**Larry Wasserman**

*Department of Statistics  
Carnegie Mellon University  
Pittsburgh, PA 15213, USA*

LARRY@STAT.CMU.EDU

**Editor:** Tong Zhang

## Abstract

We study graph estimation and density estimation in high dimensions, using a family of density estimators based on forest structured undirected graphical models. For density estimation, we do not assume the true distribution corresponds to a forest; rather, we form kernel density estimates of the bivariate and univariate marginals, and apply Kruskal's algorithm to estimate the optimal forest on held out data. We prove an oracle inequality on the excess risk of the resulting estimator relative to the risk of the best forest. For graph estimation, we consider the problem of estimating forests with restricted tree sizes. We prove that finding a maximum weight spanning forest with restricted tree size is NP-hard, and develop an approximation algorithm for this problem. Viewing the tree size as a complexity parameter, we then select a forest using data splitting, and prove bounds on excess risk and structure selection consistency of the procedure. Experiments with simulated data and microarray data indicate that the methods are a practical alternative to Gaussian graphical models.

**Keywords:** kernel density estimation, forest structured Markov network, high dimensional inference, risk consistency, structure selection consistency

## 1. Introduction

One way to explore the structure of a high dimensional distribution  $P$  for a random vector  $X = (X_1, \dots, X_d)$  is to estimate its undirected graph. The undirected graph  $G$  associated with  $P$  has  $d$  vertices corresponding to the variables  $X_1, \dots, X_d$ , and omits an edge between two nodes  $X_i$  and  $X_j$  if and only if  $X_i$  and  $X_j$  are conditionally independent given the other variables. Currently, the most popular methods for estimating  $G$  assume that the distribution  $P$  is Gaussian. Finding the

---

\*. John Lafferty is also in the Department of Statistics at Carnegie Mellon University.

graphical structure in this case amounts to estimating the inverse covariance matrix  $\Omega$ ; the edge between  $X_j$  and  $X_k$  is missing if and only if  $\Omega_{jk} = 0$ . Algorithms for optimizing the  $\ell_1$ -regularized log-likelihood have recently been proposed that efficiently produce sparse estimates of the inverse covariance matrix and the underlying graph (Banerjee et al., 2008; Friedman et al., 2007).

In this paper our goal is to relax the Gaussian assumption and to develop nonparametric methods for estimating the graph of a distribution. Of course, estimating a high dimensional distribution is impossible without making any assumptions. The approach we take here is to force the graphical structure to be a forest, where each pair of vertices is connected by at most one path. Thus, we relax the distributional assumption of normality but we restrict the family of undirected graphs that are allowed.

If the graph for  $P$  is a forest, then a simple conditioning argument shows that its density  $p$  can be written as

$$p(x) = \prod_{(i,j) \in E} \frac{p(x_i, x_j)}{p(x_i)p(x_j)} \prod_{k=1}^d p(x_k)$$

where  $E$  is the set of edges in the forest (Lauritzen, 1996). Here  $p(x_i, x_j)$  is the bivariate marginal density of variables  $X_i$  and  $X_j$ , and  $p(x_k)$  is the univariate marginal density of the variable  $X_k$ . With this factorization, we see that it is only necessary to estimate the bivariate and univariate marginals. Given any distribution  $P$  with density  $p$ , there is a tree  $T$  and a density  $p_T$  whose graph is  $T$  and which is closest in Kullback-Leibler divergence to  $p$ . When  $P$  is known, then the best fitting tree distribution can be obtained by Kruskal's algorithm (Kruskal, 1956), or other algorithms for finding a maximum weight spanning tree. In the discrete case, the algorithm can be applied to the estimated probability mass function, resulting in a procedure originally proposed by Chow. and Liu (1968). Here we are concerned with continuous random variables, and we estimate the bivariate marginals with nonparametric kernel density estimators.

In high dimensions, fitting a fully connected spanning tree can be expected to overfit. We regulate the complexity of the forest by selecting the edges to include using a data splitting scheme, a simple form of cross validation. In particular, we consider the family of forest structured densities that use the marginal kernel density estimates constructed on the first partition of the data, and estimate the risk of the resulting densities over a second, held out partition. The optimal forest in terms of the held out risk is then obtained by finding a maximum weight spanning forest for an appropriate set of edge weights.

A closely related approach is proposed by Bach and Jordan (2003), where a tree is estimated for the random vector  $Y = WX$  instead of  $X$ , where  $W$  is a linear transformation, using an algorithm that alternates between estimating  $W$  and estimating the tree  $T$ . Kernel density estimators are used, and a regularization term that is a function of the number of edges in the tree is included to bias the optimization toward smaller trees. We omit the transformation  $W$ , and we use a data splitting method rather than penalization to choose the complexity of the forest.

While tree and forest structured density estimation has been long recognized as a useful tool, there has been little theoretical analysis of the statistical properties of such density estimators. The main contribution of this paper is an analysis of the asymptotic properties of forest density estimation in high dimensions. We allow both the sample size  $n$  and dimension  $d$  to increase, and prove oracle results on the risk of the method. In particular, we assume that the univariate and bivariate

marginal densities lie in a Hölder class with exponent  $\beta$  (see Section 4 for details), and show that

$$R(\widehat{p}_{\widehat{F}}) - \min_F R(\widehat{p}_F) = O_P \left( \sqrt{\log(nd)} \left( \frac{k^* + \widehat{k}}{n^{\beta/(2+2\beta)}} + \frac{d}{n^{\beta/(1+2\beta)}} \right) \right)$$

where  $R$  denotes the risk, the expected negative log-likelihood,  $\widehat{k}$  is the number of edges in the estimated forest  $\widehat{F}$ , and  $k^*$  is the number of edges in the optimal forest  $F^*$  that can be constructed in terms of the kernel density estimates  $\widehat{p}$ .

In addition to the above results on risk consistency, we establish conditions under which

$$\mathbb{P} \left( \widehat{F}_d^{(k)} = F_d^{*(k)} \right) \rightarrow 1$$

as  $n \rightarrow \infty$ , where  $F_d^{*(k)}$  is the *oracle forest*—the best forest with  $k$  edges; this result allows the dimensionality  $d$  to increase as fast as  $o(\exp(n^{\beta/(1+\beta)}))$ , while still having consistency in the selection of the oracle forest.

Among the only other previous work analyzing tree structured graphical models is Tan et al. (2011) and Chechetka and Guestrin (2007). Tan et al. (2011) analyze the error exponent in the rate of decay of the error probability for estimating the tree, in the fixed dimension setting, and Chechetka and Guestrin (2007) give a PAC analysis. An extension to the Gaussian case is given by Tan et al. (2010).

We also study the problem of estimating forests with restricted tree sizes. In many applications, one is interested in obtaining a graphical representation of a high dimensional distribution to aid in interpretation. For instance, a biologist studying gene interaction networks might be interested in a visualization that groups together genes in small sets. Such a clustering approach through density estimation is problematic if the graph is allowed to have cycles, as this can require marginal densities to be estimated with many interacting variables. Restricting the graph to be a forest circumvents the curse of dimensionality by requiring only univariate and bivariate marginal densities. The problem of clustering the variables into small interacting sets, each supported by a tree-structured density, becomes the problem of estimating a maximum weight spanning forest with a restriction on the size of each component tree. As we demonstrate, estimating restricted tree size forests can also be useful in model selection for the purpose of risk minimization. Limiting the tree size gives another way of regulating tree complexity that provides larger family of forest to select from in the data splitting procedure.

While the problem of finding a maximum weight forest with restricted tree size may be natural, it appears not to have been studied previously. We prove that the problem is NP-hard through a reduction from the problem of Exact 3-Cover (Garey and Johnson, 1979), where we are given a set  $X$  and a family  $\mathcal{S}$  of 3-element subsets of  $X$ , and must choose a subfamily of disjoint 3-element subsets to cover  $X$ . While finding the exact optimum is hard, we give a practical 4-approximation algorithm for finding the optimal tree restricted forest; that is, our algorithm outputs a forest whose weight is guaranteed to be at least  $\frac{1}{4}w(F^*)$ , where  $w(F^*)$  is the weight of the optimal forest. This approximation guarantee translates into excess risk bounds on the constructed forest using our previous analysis. Our experimental results with this approximation algorithm show that it can be effective in practice for forest density estimation.

In Section 2 we review some background and notation. In Section 3 we present a two-stage algorithm for estimating high dimensional densities supported by forests, and we provide a theoretical

analysis of the algorithm in Section 4, with the detailed proofs collected in an appendix. In Section 5, we explain how to estimate maximum weight forests with restricted tree size. In Section 6 we present experiments with both simulated data and gene microarray data sets, where the problem is to estimate the gene-gene association graphs.

## 2. Preliminaries and Notation

Let  $p^*(x)$  be a probability density with respect to Lebesgue measure  $\mu(\cdot)$  on  $\mathbb{R}^d$  and let  $X^{(1)}, \dots, X^{(n)}$  be  $n$  independent identically distributed  $\mathbb{R}^d$ -valued data vectors sampled from  $p^*(x)$  where  $X^{(i)} = (X_1^{(i)}, \dots, X_d^{(i)})$ . Let  $\mathcal{X}_j$  denote the range of  $X_j^{(i)}$  and let  $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_d$ . For simplicity we assume that  $\mathcal{X}_j = [0, 1]$ .

A graph is a forest if it is acyclic. If  $F$  is a  $d$ -node undirected forest with vertex set  $V_F = \{1, \dots, d\}$  and edge set  $E(F) \subset \{1, \dots, d\} \times \{1, \dots, d\}$ , the number of edges satisfies  $|E(F)| < d$ , noting that we do not restrict the graph to be connected. We say that a probability density function  $p(x)$  is *supported by a forest  $F$*  if the density can be written as

$$p_F(x) = \prod_{(i,j) \in E(F)} \frac{p(x_i, x_j)}{p(x_i)p(x_j)} \prod_{k \in V_F} p(x_k), \quad (1)$$

where each  $p(x_i, x_j)$  is a bivariate density on  $\mathcal{X}_i \times \mathcal{X}_j$ , and each  $p(x_k)$  is a univariate density on  $\mathcal{X}_k$ . More details can be found in Lauritzen (1996).

Let  $\mathcal{F}_d$  be the family of forests with  $d$  nodes, and let  $\mathcal{P}_d$  be the corresponding family of densities:

$$\mathcal{P}_d = \left\{ p \geq 0 : \int_{\mathcal{X}} p(x) d\mu(x) = 1, \text{ and } p(x) \text{ satisfies (1) for some } F \in \mathcal{F}_d \right\}. \quad (2)$$

To bound the number of labeled spanning forests on  $d$  nodes, note that each such forest can be obtained by forming a labeled tree on  $d + 1$  nodes, and then removing node  $d + 1$ . From Cayley's formula (Cayley, 1889; Aigner and Ziegler, 1998), we then obtain the following.

**Proposition 1** *The size of the collection  $\mathcal{F}_d$  of labeled forests on  $d$  nodes satisfies*

$$|\mathcal{F}_d| < (d + 1)^{d-1}.$$

Define the oracle forest density

$$q^* = \arg \min_{q \in \mathcal{P}_d} D(p^* \| q) \quad (3)$$

where the Kullback-Leibler divergence  $D(p \| q)$  between two densities  $p$  and  $q$  is

$$D(p \| q) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx,$$

under the convention that  $0 \log(0/q) = 0$ , and  $p \log(p/0) = \infty$  for  $p \neq 0$ . The following is proved by Bach and Jordan (2003).

**Proposition 2** Let  $q^*$  be defined as in (3). There exists a forest  $F^* \in \mathcal{F}_d$ , such that

$$q^* = p_{F^*}^* = \prod_{(i,j) \in E(F^*)} \frac{p^*(x_i, x_j)}{p^*(x_i) p^*(x_j)} \prod_{k \in V_{F^*}} p^*(x_k) \quad (4)$$

where  $p^*(x_i, x_j)$  and  $p^*(x_i)$  are the bivariate and univariate marginal densities of  $p^*$ .

For any density  $q(x)$ , the negative log-likelihood risk  $R(q)$  is defined as

$$R(q) = -\mathbb{E} \log q(X) = - \int_{\mathcal{X}} p^*(x) \log q(x) dx$$

where the expectation is defined with respect to the distribution of  $X$ .

It is straightforward to see that the density  $q^*$  defined in (3) also minimizes the negative log-likelihood loss:

$$q^* = \arg \min_{q \in \mathcal{P}_d} D(p^* \| q) = \arg \min_{q \in \mathcal{P}_d} R(q).$$

Let  $\hat{p}(x)$  be the kernel density estimate, we also define

$$\hat{R}(q) = - \int_{\mathcal{X}} \hat{p}(x) \log q(x) dx.$$

We thus define the oracle risk as  $R^* = R(q^*)$ . Using Proposition 2 and Equation (1), we have

$$\begin{aligned} R^* &= R(q^*) = R(p_{F^*}^*) \\ &= - \int_{\mathcal{X}} p^*(x) \left( \sum_{(i,j) \in E(F^*)} \log \frac{p^*(x_i, x_j)}{p^*(x_i) p^*(x_j)} + \sum_{k \in V_{F^*}} \log(p^*(x_k)) \right) dx \\ &= - \sum_{(i,j) \in E(F^*)} \int_{\mathcal{X}_i \times \mathcal{X}_j} p^*(x_i, x_j) \log \frac{p^*(x_i, x_j)}{p^*(x_i) p^*(x_j)} dx_i dx_j - \sum_{k \in V_{F^*}} \int_{\mathcal{X}_k} p^*(x_k) \log p^*(x_k) dx_k \\ &= - \sum_{(i,j) \in E(F^*)} I(X_i; X_j) + \sum_{k \in V_{F^*}} H(X_k), \end{aligned} \quad (5)$$

where

$$I(X_i; X_j) = \int_{\mathcal{X}_i \times \mathcal{X}_j} p^*(x_i, x_j) \log \frac{p^*(x_i, x_j)}{p^*(x_i) p^*(x_j)} dx_i dx_j$$

is the mutual information between the pair of variables  $X_i, X_j$  and

$$H(X_k) = - \int_{\mathcal{X}_k} p^*(x_k) \log p^*(x_k) dx_k$$

is the entropy. While the best forest will in fact be a spanning tree, the densities  $p^*(x_i, x_j)$  are in practice not known. We estimate the marginals using finite data, in terms of a kernel density estimates  $\hat{p}_{n_1}(x_i, x_j)$  over a training set of size  $n_1$ . With these estimated marginals, we consider all forest density estimates of the form

$$\hat{p}_F(x) = \prod_{(i,j) \in E(F)} \frac{\hat{p}_{n_1}(x_i, x_j)}{\hat{p}_{n_1}(x_i) \hat{p}_{n_1}(x_j)} \prod_{k \in V_F} \hat{p}_{n_1}(x_k).$$

Within this family, the best density estimate may not be supported on a full spanning tree, since a full tree will in general be subject to overfitting. Analogously, in high dimensional linear regression, the optimal regression model will generally be a full  $d$ -dimensional fit, with a nonzero parameter for each variable. However, when estimated on finite data the variance of a full model will dominate the squared bias, resulting in overfitting. In our setting of density estimation we will regulate the complexity of the forest by cross validating over a held out set.

There are several different ways to judge the quality of a forest structured density estimator. In this paper we concern ourselves with prediction and structure estimation.

**Definition 3 ((Risk consistency))** For an estimator  $\hat{q}_n \in \mathcal{P}_d$ , the excess risk is defined as  $R(\hat{q}_n) - R^*$ . The estimator  $\hat{q}_n$  is risk consistent with convergence rate  $\delta_n$  if

$$\lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} \mathbb{P}(R(\hat{q}_n) - R^* \geq M\delta_n) = 0.$$

In this case we write  $R(\hat{q}_n) - R^* = O_P(\delta_n)$ .

**Definition 4 ((Estimation consistency))** An estimator  $\hat{q}_n \in \mathcal{P}_d$  is estimation consistent with convergence rate  $\delta_n$ , with respect to the Kullback-Leibler divergence, if

$$\lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} \mathbb{P}(D(p_{F^*}^* \parallel \hat{q}_n) \geq M\delta_n) = 0.$$

**Definition 5 ((Structure selection consistency))** An estimator  $\hat{q}_n \in \mathcal{P}_d$  supported by a forest  $\hat{F}_n$  is structure selection consistent if

$$\mathbb{P}\left(E(\hat{F}_n) \neq E(F^*)\right) \rightarrow 0,$$

as  $n$  goes to infinity, where  $F^*$  is defined in (4).

Later we will show that estimation consistency is almost equivalent to risk consistency. If the true density is given, these two criteria are exactly the same; otherwise, estimation consistency requires stronger conditions than risk consistency.

It is important to note that risk consistency is an oracle property, in the sense that the true density  $p^*(x)$  is not restricted to be supported by a forest; rather, the property assesses how well a given estimator  $\hat{q}$  approximates the best forest density (the oracle) within a class.

### 3. Kernel Density Estimation For Forests

If the true density  $p^*(x)$  were known, by Proposition 2, the density estimation problem would be reduced to finding the best forest structure  $F_d^*$ , satisfying

$$F_d^* = \arg \min_{F \in \mathcal{F}_d} R(p_F^*) = \arg \min_{F \in \mathcal{F}_d} D(p^* \parallel p_F^*).$$

The optimal forest  $F_d^*$  can be found by minimizing the right hand side of (5). Since the entropy term  $H(X) = \sum_k H(X_k)$  is constant across all forests, this can be recast as the problem of finding the maximum weight spanning forest for a weighted graph, where the weight  $w(i, j)$  of the edge connecting nodes  $i$  and  $j$  is  $I(X_i; X_j)$ . Kruskal's algorithm (Kruskal, 1956) is a greedy algorithm

that is guaranteed to find a maximum weight spanning tree of a weighted graph. In the setting of density estimation, this procedure was proposed by Chow. and Liu (1968) as a way of constructing a tree approximation to a distribution. At each stage the algorithm adds an edge connecting that pair of variables with maximum mutual information among all pairs not yet visited by the algorithm, if doing so does not form a cycle. When stopped early, after  $k < d - 1$  edges have been added, it yields the best  $k$ -edge weighted forest.

Of course, the above procedure is not practical since the true density  $p^*(x)$  is unknown. We replace the population mutual information  $I(X_i; X_j)$  in (5) by the plug-in estimate  $\widehat{I}_n(X_i, X_j)$ , defined as

$$\widehat{I}_n(X_i, X_j) = \int_{\mathcal{X}_i \times \mathcal{X}_j} \widehat{p}_n(x_i, x_j) \log \frac{\widehat{p}_n(x_i, x_j)}{\widehat{p}_n(x_i) \widehat{p}_n(x_j)} dx_i dx_j$$

where  $\widehat{p}_n(x_i, x_j)$  and  $\widehat{p}_n(x_i)$  are bivariate and univariate kernel density estimates. Given this estimated mutual information matrix  $\widehat{M}_n = [\widehat{I}_n(X_i, X_j)]$ , we can then apply Kruskal's algorithm (equivalently, the Chow-Liu algorithm) to find the best forest structure  $\widehat{F}_n$ .

Since the number of edges of  $\widehat{F}_n$  controls the number of degrees of freedom in the final density estimator, we need an automatic data-dependent way to choose it. We adopt the following two-stage procedure. First, randomly partition the data into two sets  $\mathcal{D}_1$  and  $\mathcal{D}_2$  of sizes  $n_1$  and  $n_2$ ; then, apply the following steps:

1. Using  $\mathcal{D}_1$ , construct kernel density estimates of the univariate and bivariate marginals and calculate  $\widehat{I}_{n_1}(X_i, X_j)$  for  $i, j \in \{1, \dots, d\}$  with  $i \neq j$ . Construct a full tree  $\widehat{F}_{n_1}^{(d-1)}$  with  $d - 1$  edges, using the Chow-Liu algorithm.
2. Using  $\mathcal{D}_2$ , prune the tree  $\widehat{F}_{n_1}^{(d-1)}$  to find a forest  $\widehat{F}_{n_1}^{(k)}$  with  $\widehat{k}$  edges, for  $0 \leq \widehat{k} \leq d - 1$ .

Once  $\widehat{F}_{n_1}^{(k)}$  is obtained in Step 2, we can calculate  $\widehat{p}_{\widehat{F}_{n_1}^{(k)}}$  according to (1), using the kernel density estimates constructed in Step 1.

### 3.1 Step 1: Estimating the Marginals

Step 1 is carried out on the data set  $\mathcal{D}_1$ . Let  $K(\cdot)$  be a univariate kernel function. Given an evaluation point  $(x_i, x_j)$ , the bivariate kernel density estimate for  $(X_i, X_j)$  based on the observations  $\{X_i^{(s)}, X_j^{(s)}\}_{s \in \mathcal{D}_1}$  is defined as

$$\widehat{p}_{n_1}(x_i, x_j) = \frac{1}{n_1} \sum_{s \in \mathcal{D}_1} \frac{1}{h_2^2} K\left(\frac{X_i^{(s)} - x_i}{h_2}\right) K\left(\frac{X_j^{(s)} - x_j}{h_2}\right), \quad (6)$$

where we use a product kernel with  $h_2 > 0$  be the bandwidth parameter. The univariate kernel density estimate  $\widehat{p}_{n_1}(x_k)$  for  $X_k$  is

$$\widehat{p}_{n_1}(x_k) = \frac{1}{n_1} \sum_{s \in \mathcal{D}_1} \frac{1}{h_1} K\left(\frac{X_k^{(s)} - x_k}{h_1}\right), \quad (7)$$

---

**Algorithm 1** Chow-Liu (Kruskal)
 

---

- 1: **Input** data  $\mathcal{D}_1 = \{X^{(1)}, \dots, X^{(n_1)}\}$ .
  - 2: Calculate  $\widehat{M}_{n_1}$ , according to (6), (7), and (8).
  - 3: Initialize  $E^{(0)} = \emptyset$
  - 4: **for**  $k = 1, \dots, d - 1$  **do**
  - 5:    $(i^{(k)}, j^{(k)}) \leftarrow \arg \max_{(i,j)} \widehat{M}_{n_1}(i, j)$  such that  $E^{(k-1)} \cup \{(i^{(k)}, j^{(k)})\}$  does not contain a cycle
  - 6:    $E^{(k)} \leftarrow E^{(k-1)} \cup \{(i^{(k)}, j^{(k)})\}$
  - 7: **Output** tree  $\widehat{F}_{n_1}^{(d-1)}$  with edge set  $E^{(d-1)}$ .
- 

where  $h_1 > 0$  is the univariate bandwidth. Detailed specifications for  $K(\cdot)$  and  $h_1, h_2$  will be discussed in the next section.

We assume that the data lie in a  $d$ -dimensional unit cube  $\mathcal{X} = [0, 1]^d$ . To calculate the empirical mutual information  $\widehat{I}_{n_1}(X_i, X_j)$ , we need to numerically evaluate a two-dimensional integral. To do so, we calculate the kernel density estimates on a grid of points. We choose  $m$  evaluation points on each dimension,  $x_{1i} < x_{2i} < \dots < x_{mi}$  for the  $i$ th variable. The mutual information  $\widehat{I}_{n_1}(X_i, X_j)$  is then approximated as

$$\widehat{I}_{n_1}(X_i, X_j) = \frac{1}{m^2} \sum_{k=1}^m \sum_{\ell=1}^m \widehat{p}_{n_1}(x_{ki}, x_{\ell j}) \log \frac{\widehat{p}_{n_1}(x_{ki}, x_{\ell j})}{\widehat{p}_{n_1}(x_{ki}) \widehat{p}_{n_1}(x_{\ell j})}. \quad (8)$$

The approximation error can be made arbitrarily small by choosing  $m$  sufficiently large. As a practical concern, care needs to be taken that the factors  $\widehat{p}_{n_1}(x_{ki})$  and  $\widehat{p}_{n_1}(x_{\ell j})$  in the denominator are not too small; a truncation procedure can be used to ensure this. Once the  $d \times d$  mutual information matrix  $\widehat{M}_{n_1} = [\widehat{I}_{n_1}(X_i, X_j)]$  is obtained, we can apply the Chow-Liu (Kruskal) algorithm to find a maximum weight spanning tree.

### 3.2 Step 2: Optimizing the Forest

The full tree  $\widehat{F}_{n_1}^{(d-1)}$  obtained in Step 1 might have high variance when the dimension  $d$  is large, leading to overfitting in the density estimate. In order to reduce the variance, we prune the tree; that is, we choose forest with  $k \leq d - 1$  edges. The number of edges  $k$  is a tuning parameter that induces a bias-variance tradeoff.

In order to choose  $k$ , note that in stage  $k$  of the Chow-Liu algorithm we have an edge set  $E^{(k)}$  (in the notation of the Algorithm 1) which corresponds to a forest  $\widehat{F}_{n_1}^{(k)}$  with  $k$  edges, where  $\widehat{F}_{n_1}^{(0)}$  is the union of  $d$  disconnected nodes. To select  $k$ , we choose among the  $d$  trees  $\widehat{F}_{n_1}^{(0)}, \widehat{F}_{n_1}^{(1)}, \dots, \widehat{F}_{n_1}^{(d-1)}$ .

Let  $\widehat{p}_{n_2}(x_i, x_j)$  and  $\widehat{p}_{n_2}(x_k)$  be defined as in (6) and (7), but now evaluated solely based on the held-out data in  $\mathcal{D}_2$ . For a density  $p_F$  that is supported by a forest  $F$ , we define the held-out negative log-likelihood risk as

$$\begin{aligned} \widehat{R}_{n_2}(p_F) &= - \sum_{(i,j) \in E_F} \int_{\mathcal{X}_i \times \mathcal{X}_j} \widehat{p}_{n_2}(x_i, x_j) \log \frac{p(x_i, x_j)}{p(x_i) p(x_j)} dx_i dx_j - \sum_{k \in V_F} \int_{\mathcal{X}_k} \widehat{p}_{n_2}(x_k) \log p(x_k) dx_k. \end{aligned} \quad (9)$$

The selected forest is then  $\widehat{F}_{n_1}^{(\widehat{k})}$  where

$$\widehat{k} = \arg \min_{k \in \{0, \dots, d-1\}} \widehat{R}_{n_2} \left( \widehat{P}_{\widehat{F}_{n_1}^{(k)}} \right)$$

and where  $\widehat{p}_{\widehat{F}_{n_1}^{(k)}}$  is computed using the density estimate  $\widehat{p}_{n_1}$  constructed on  $\mathcal{D}_1$ .

For computational simplicity, we can also estimate  $\widehat{k}$  as

$$\begin{aligned} \widehat{k} &= \arg \max_{k \in \{0, \dots, d-1\}} \frac{1}{n_2} \sum_{s \in \mathcal{D}_2} \log \left( \prod_{(i,j) \in E^{(k)}} \frac{\widehat{p}_{n_1}(X_i^{(s)}, X_j^{(s)})}{\widehat{p}_{n_1}(X_i^{(s)}) \widehat{p}_{n_1}(X_j^{(s)})} \prod_{k \in V_{\widehat{F}_{n_1}^{(k)}}} \widehat{p}_{n_1}(X_k^{(s)}) \right) \\ &= \arg \max_{k \in \{0, \dots, d-1\}} \frac{1}{n_2} \sum_{s \in \mathcal{D}_2} \log \left( \prod_{(i,j) \in E^{(k)}} \frac{\widehat{p}_{n_1}(X_i^{(s)}, X_j^{(s)})}{\widehat{p}_{n_1}(X_i^{(s)}) \widehat{p}_{n_1}(X_j^{(s)})} \right). \end{aligned}$$

This minimization can be efficiently carried out by iterating over the  $d-1$  edges in  $\widehat{F}_{n_1}^{(d-1)}$ .

Once  $\widehat{k}$  is obtained, the final forest density estimate is given by

$$\widehat{p}_n(x) = \prod_{(i,j) \in E^{(\widehat{k})}} \frac{\widehat{p}_{n_1}(x_i, x_j)}{\widehat{p}_{n_1}(x_i) \widehat{p}_{n_1}(x_j)} \prod_k \widehat{p}_{n_1}(x_k).$$

**Remark 6** For computational efficiency, Step 1 can be carried out simultaneously with Step 2. In particular, during the Chow-Liu iteration, whenever an edge is added to  $E^{(k)}$ , the log-likelihood of the resulting density estimator on  $\mathcal{D}_2$  can be immediately computed. A more efficient algorithm to speed up the computation of the mutual information matrix is discussed in Appendix B.

### 3.3 Building a Forest on Held-out Data

Another approach to estimating the forest structure is to estimate the marginal densities on the training set, but only build graphs on the held-out data. To do so, we first estimate the univariate and bivariate kernel density estimates using  $\mathcal{D}_1$ , denoted by  $\widehat{p}_{n_1}(x_i)$  and  $\widehat{p}_{n_1}(x_i, x_j)$ . We also construct a new set of univariate and bivariate kernel density estimates using  $\mathcal{D}_2$ ,  $\widehat{p}_{n_2}(x_i)$  and  $\widehat{p}_{n_2}(x_i, x_j)$ . We then estimate the ‘‘cross-entropies’’ of the kernel density estimates  $\widehat{p}_{n_1}$  for each pair of variables by computing

$$\begin{aligned} \widehat{I}_{n_2, n_1}(X_i, X_j) &= \int \widehat{p}_{n_2}(x_i, x_j) \log \frac{\widehat{p}_{n_1}(x_i, x_j)}{\widehat{p}_{n_1}(x_i) \widehat{p}_{n_1}(x_j)} dx_i dx_j \\ &\approx \frac{1}{m^2} \sum_{k=1}^m \sum_{\ell=1}^m \widehat{p}_{n_2}(x_{ki}, x_{\ell j}) \log \frac{\widehat{p}_{n_1}(x_{ki}, x_{\ell j})}{\widehat{p}_{n_1}(x_{ki}) \widehat{p}_{n_1}(x_{\ell j})}. \end{aligned} \quad (10)$$

Our method is to use  $\widehat{I}_{n_2, n_1}(X_i, X_j)$  as edge weights on a full graph and run Kruskal’s algorithm until we encounter edges with negative weight. Let  $\mathcal{F}$  be the set of all forests and  $\widehat{w}_{n_2}(i, j) = \widehat{I}_{n_2, n_1}(X_i, X_j)$ . The final forest is then

$$\widehat{F}_{n_2} = \arg \max_{F \in \mathcal{F}} \widehat{w}_{n_2}(F) = \arg \min_{F \in \mathcal{F}} \widehat{R}_{n_2}(\widehat{P}_F)$$

By building a forest on held-out data, we directly cross-validate over *all* forests.

## 4. Statistical Properties

In this section we present our theoretical results on risk consistency, structure selection consistency, and estimation consistency of the forest density estimate  $\hat{p}_n = \hat{p}_{\hat{F}_d^{(k)}}$ .

To establish some notation, we write  $a_n = \Omega(b_n)$  if there exists a constant  $c$  such that  $a_n \geq cb_n$  for sufficiently large  $n$ . We also write  $a_n \asymp b_n$  if there exists a constant  $c$  such that  $a_n \leq cb_n$  and  $b_n \leq ca_n$  for sufficiently large  $n$ . Given a  $d$ -dimensional function  $f$  on the domain  $\mathcal{X}$ , we denote its  $L_2(P)$ -norm and sup-norm as

$$\|f\|_{L_2(P)} = \sqrt{\int_{\mathcal{X}} f^2(x) dP_X(x)}, \quad \|f\|_{\infty} = \sup_{x \in \mathcal{X}} |f(x)|$$

where  $P_X$  is the probability measure induced by  $X$ . Throughout this section, all constants are treated as generic values, and as a result they can change from line to line.

In our use of a data splitting scheme, we always adopt equally sized splits for simplicity, so that  $n_1 = n_2 = n/2$ , noting that this does not affect the final rate of convergence.

### 4.1 Assumptions on the Density

Fix  $\beta > 0$ . For any  $d$ -tuple  $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{N}^d$  and  $x = (x_1, \dots, x_d) \in \mathcal{X}$ , we define  $x^\alpha = \prod_{j=1}^d x_j^{\alpha_j}$ . Let  $D^\alpha$  denote the differential operator

$$D^\alpha = \frac{\partial^{\alpha_1 + \dots + \alpha_d}}{\partial x_1^{\alpha_1} \dots \partial x_d^{\alpha_d}}.$$

For any real-valued  $d$ -dimensional function  $f$  on  $\mathcal{X}$  that is  $\lfloor \beta \rfloor$ -times continuously differentiable at point  $x_0 \in \mathcal{X}$ , let  $P_{f, x_0}^{(\beta)}(x)$  be its Taylor polynomial of degree  $\lfloor \beta \rfloor$  at point  $x_0$ :

$$P_{f, x_0}^{(\beta)}(x) = \sum_{\alpha_1 + \dots + \alpha_d \leq \lfloor \beta \rfloor} \frac{(x - x_0)^\alpha}{\alpha_1! \dots \alpha_d!} D^\alpha f(x_0).$$

Fix  $L > 0$ , and denote by  $\Sigma(\beta, L, r, x_0)$  the set of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  that are  $\lfloor \beta \rfloor$ -times continuously differentiable at  $x_0$  and satisfy

$$\left| f(x) - P_{f, x_0}^{(\beta)}(x) \right| \leq L \|x - x_0\|_2^\beta, \quad \forall x \in \mathcal{B}(x_0, r)$$

where  $\mathcal{B}(x_0, r) = \{x : \|x - x_0\|_2 \leq r\}$  is the  $L_2$ -ball of radius  $r$  centered at  $x_0$ . The set  $\Sigma(\beta, L, r, x_0)$  is called the  $(\beta, L, r, x_0)$ -locally Hölder class of functions. Given a set  $A$ , we define

$$\Sigma(\beta, L, r, A) = \bigcap_{x_0 \in A} \Sigma(\beta, L, r, x_0).$$

The following are the regularity assumptions we make on the true density function  $p^*(x)$ .

**Assumption 1** For any  $1 \leq i < j \leq d$ , we assume

(D1) there exist  $L_1 > 0$  and  $L_2 > 0$  such that for any  $c > 0$  the true bivariate and univariate densities satisfy

$$p^*(x_i, x_j) \in \Sigma\left(\beta, L_2, c(\log n/n)^{\frac{1}{2\beta+2}}, \mathcal{X}_i \times \mathcal{X}_j\right)$$

and

$$p^*(x_i) \in \Sigma\left(\beta, L_1, c(\log n/n)^{\frac{1}{2\beta+1}}, \mathcal{X}_i\right);$$

(D2) there exists two constants  $c_1$  and  $c_2$  such that

$$c_1 \leq \inf_{x_i, x_j \in \mathcal{X}_i \times \mathcal{X}_j} p^*(x_i, x_j) \leq \sup_{x_i, x_j \in \mathcal{X}_i \times \mathcal{X}_j} p^*(x_i, x_j) \leq c_2$$

$\mu$ -almost surely.

These assumptions are mild, in the sense that instead of adding constraints on the joint density  $p^*(x)$ , we only add regularity conditions on the bivariate and univariate marginals.

#### 4.2 Assumptions on the Kernel

An important ingredient in our analysis is an exponential concentration result for the kernel density estimate, due to Giné and Guillou (2002). We first specify the requirements on the kernel function  $K(\cdot)$ .

Let  $(\Omega, \mathcal{A})$  be a measurable space and let  $\mathcal{F}$  be a uniformly bounded collection of measurable functions.

**Definition 7**  $\mathcal{F}$  is a bounded measurable VC class of functions with characteristics  $A$  and  $v$  if it is separable and for every probability measure  $P$  on  $(\Omega, \mathcal{A})$  and any  $0 < \varepsilon < 1$ ,

$$N(\varepsilon \|F\|_{L_2(P)}, \mathcal{F}, \|\cdot\|_{L_2(P)}) \leq \left(\frac{A}{\varepsilon}\right)^v,$$

where  $F(x) = \sup_{f \in \mathcal{F}} |f(x)|$  and  $N(\varepsilon, \mathcal{F}, \|\cdot\|_{L_2(P)})$  denotes the  $\varepsilon$ -covering number of the metric space  $(\Omega, \|\cdot\|_{L_2(P)})$ ; that is, the smallest number of balls of radius no larger than  $\varepsilon$  (in the norm  $\|\cdot\|_{L_2(P)}$ ) needed to cover  $\mathcal{F}$ .

The one-dimensional density estimates are constructed using a kernel  $K$ , and the two-dimensional estimates are constructed using the product kernel

$$K_2(x, y) = K(x) \cdot K(y).$$

**Assumption 2** The kernel  $K$  satisfies the following properties.

(K1)  $\int K(u) du = 1$ ,  $\int_{-\infty}^{\infty} K^2(u) du < \infty$  and  $\sup_{u \in \mathbb{R}} K(u) \leq c$  for some constant  $c$ .

(K2)  $K$  is a finite linear combination of functions  $g$  whose epigraphs  $\text{epi}(g) = \{(s, u) : g(s) \geq u\}$ , can be represented as a finite number of Boolean operations (union and intersection) among sets of the form  $\{(s, u) : Q(s, u) \geq \phi(u)\}$ , where  $Q$  is a polynomial on  $\mathbb{R} \times \mathbb{R}$  and  $\phi$  is an arbitrary real function.

(K3)  $K$  has a compact support and for any  $\ell \geq 1$  and  $1 \leq \ell' \leq \lfloor \beta \rfloor$

$$\int |t|^\beta |K(t)| dt < \infty, \text{ and } \int |K(t)|^\ell dt < \infty, \int t^{\ell'} K(t) dt = 0.$$

Assumptions (K1), (K2) and (K3) are mild. As pointed out by Nolan and Pollard (1987), both the pyramid (truncated or not) kernel and the boxcar kernel satisfy them. It follows from (K2) that the classes of functions

$$\begin{aligned}\mathcal{F}_1 &= \left\{ \frac{1}{h_1} K\left(\frac{u-\cdot}{h_1}\right) : u \in \mathbb{R}, h_1 > 0 \right\} \\ \mathcal{F}_2 &= \left\{ \frac{1}{h_2^2} K\left(\frac{u-\cdot}{h_2}\right) K\left(\frac{t-\cdot}{h_2}\right) : u, t \in \mathbb{R}, h_2 > 0 \right\}\end{aligned}\quad (11)$$

are bounded VC classes, in the sense of Definition 7. Assumption (K3) essentially says that the kernel  $K(\cdot)$  should be  $\beta$ -valid; see Tsybakov (2008) and Definition 6.1 in Rigollet and Vert (2009) for further details about this assumption. Kernels satisfying (K2) include finite linear combinations of functions of the form  $\phi(p(x))$  where  $p$  is a polynomial and  $\phi$  is a bounded function of bounded variation (Giné and Guillou, 2002; Nolan and Pollard, 1987). Therefore, the kernels constructed in terms of Legendre polynomials as in Rigollet and Vert (2009) and Tsybakov (2008), satisfy (K2) and (K3).

We choose the bandwidths  $h_1$  and  $h_2$  used in the one-dimensional and two-dimensional kernel density estimates to satisfy

$$h_1 \asymp \left(\frac{\log n}{n}\right)^{\frac{1}{1+2\beta}} \quad (12)$$

$$h_2 \asymp \left(\frac{\log n}{n}\right)^{\frac{1}{2+2\beta}}. \quad (13)$$

This choice of bandwidths ensures the optimal rate of convergence.

### 4.3 Risk Consistency

Given the above assumptions, we first present a key lemma that establishes the rates of convergence of bivariate and univariate kernel density estimates in the sup norm. The proof of this and our other technical results are provided in Appendix A.

**Lemma 8** *Under Assumptions 1 and 2, and choosing bandwidths satisfying (12) and (13), the bivariate and univariate kernel density estimates  $\hat{p}(x_i, x_j)$  and  $\hat{p}(x_k)$  in (6) and (7) satisfy*

$$\begin{aligned}\mathbb{P}\left(\max_{(i,j) \in \{1, \dots, d\} \times \{1, \dots, d\}} \sup_{(x_i, x_j) \in \mathcal{X}_i \times \mathcal{X}_j} |\hat{p}(x_i, x_j) - p^*(x_i, x_j)| \geq \varepsilon\right) \\ \leq c_2 d^2 \exp\left(-c_3 n^{\frac{\beta}{1+\beta}} (\log n)^{\frac{1}{1+\beta}} \varepsilon^2\right)\end{aligned}$$

for  $\varepsilon \geq 4c_4 h_2^\beta$ . Hence, choosing

$$\varepsilon = \Omega\left(4c_4 \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}}\right)$$

we have that

$$\max_{(i,j) \in \{1, \dots, d\} \times \{1, \dots, d\}} \sup_{(x_i, x_j) \in \mathcal{X}_i \times \mathcal{X}_j} |\widehat{p}(x_i, x_j) - p^*(x_i, x_j)| = O_P \left( \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}} \right). \quad (14)$$

Similarly,

$$\mathbb{P} \left( \max_{i \in \{1, \dots, d\}} \sup_{x_i \in \mathcal{X}_i} |\widehat{p}(x_i) - p^*(x_i)| \geq \varepsilon \right) \leq c_5 d \exp \left( -c_6 n^{\frac{2\beta}{1+2\beta}} (\log n)^{\frac{1}{1+2\beta}} \varepsilon^2 \right)$$

and

$$\max_{k \in \{1, \dots, d\}} \sup_{x_k \in \mathcal{X}_k} |\widehat{p}(x_k) - p^*(x_k)| = O_P \left( \sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}} \right). \quad (15)$$

To describe the risk consistency result, let  $\mathcal{P}_d^{(d-1)} = \mathcal{P}_d$  be the family of densities that are supported by forests with at most  $d-1$  edges, as already defined in (2). For  $0 \leq k \leq d-1$ , we define  $\mathcal{P}_d^{(k)}$  as the family of  $d$ -dimensional densities that are supported by forests with at most  $k$  edges. Then

$$\mathcal{P}_d^{(0)} \subset \mathcal{P}_d^{(1)} \subset \dots \subset \mathcal{P}_d^{(d-1)}. \quad (16)$$

Now, due to the nesting property (16), we have

$$\inf_{q_F \in \mathcal{P}_d^{(0)}} R(q_F) \geq \inf_{q_F \in \mathcal{P}_d^{(1)}} R(q_F) \geq \dots \geq \inf_{q_F \in \mathcal{P}_d^{(d-1)}} R(q_F).$$

We first analyze the forest density estimator obtained using a fixed number of edges  $k < d$ ; specifically, consider stopping the Chow-Liu algorithm in Stage 1 after  $k$  iterations. This is in contrast to the algorithm described in 3.2, where the pruned tree size is automatically determined on the held out data. While this is not very realistic in applications, since the tuning parameter  $k$  is generally hard to choose, the analysis in this case is simpler, and can be directly exploited to analyze the more complicated data-dependent method.

**Theorem 9 (Risk consistency)** *Let  $\widehat{p}_{\widehat{F}_d^{(k)}}$  be the forest density estimate with  $|E(\widehat{F}_d^{(k)})| = k$ , obtained after the first  $k$  iterations of the Chow-Liu algorithm, for some  $k \in \{0, \dots, d-1\}$ . Under Assumptions 1 and 2, we have*

$$R(\widehat{p}_{\widehat{F}_d^{(k)}}) - \inf_{q_F \in \mathcal{P}_d^{(k)}} R(q_F) = O_P \left( k \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}} + d \sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}} \right).$$

Note that this result allows the dimension  $d$  to increase at a rate  $o\left(\sqrt{n^{2\beta/(1+2\beta)}/\log n}\right)$  and the number of edges  $k$  to increase at a rate  $o\left(\sqrt{n^{\beta/(1+\beta)}/\log n}\right)$ , with the excess risk still decreasing to zero asymptotically.

The above results can be used to prove a risk consistency result for the data-dependent pruning method using the data-splitting scheme described in Section 3.2.

**Theorem 10** Let  $\widehat{p}_{\widehat{F}_d^{(\widehat{k})}}$  be the forest density estimate using the data-dependent pruning method in Section 3.2, and let  $\widehat{p}_{\widehat{F}_d^{(k)}}$  be the estimate with  $|E(\widehat{F}_d^{(k)})| = k$  obtained after the first  $k$  iterations of the Chow-Liu algorithm. Under Assumptions 1 and 2, we have

$$R(\widehat{p}_{\widehat{F}_d^{(\widehat{k})}}) - \min_{0 \leq k \leq d-1} R(\widehat{p}_{\widehat{F}_d^{(k)}}) = O_P \left( (k^* + \widehat{k}) \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}} + d \sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}} \right)$$

where  $k^* = \arg \min_{0 \leq k \leq d-1} R(\widehat{p}_{\widehat{F}_d^{(k)}})$ .

The proof of this theorem is given in the appendix. A parallel result can be obtained for the method described in Section 3.3, which builds the forest by running Kruskal's algorithm on the heldout data.

**Theorem 11** Let  $\widehat{F}_{n_2}$  be the forest obtained using Kruskal's algorithm on held-out data, and let  $\widehat{k} = |\widehat{F}_{n_2}|$  be the number of edges in  $\widehat{F}_{n_2}$ . Then

$$R(\widehat{p}_{\widehat{F}_{n_2}}) - \min_{F \in \mathcal{F}} R(\widehat{p}_F) = O_P \left( (k^* + \widehat{k}) \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}} + d \sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}} \right)$$

where  $k^* = |F^*|$  is the number of edges in the optimal forest  $F^* = \arg \min_{F \in \mathcal{F}} R(\widehat{p}_F)$ .

#### 4.4 Structure Selection Consistency

In this section, we provide conditions guaranteeing that the procedure is structure selection consistent. Again, we do not assume the true density  $p^*(x)$  is consistent with a forest; rather, we are interested in comparing the estimated forest structure to the oracle forest which minimizes the risk. In this way our result differs from that in Tan et al. (2011), although there are similarities in the analysis.

By Proposition 2, we can define

$$p_{F_d^{(k)}}^* = \arg \min_{q_F \in \mathcal{P}_d^{(k)}} R(q_F).$$

Thus  $F_d^{(k)}$  is the optimal forest within  $\mathcal{P}_d^{(k)}$  that minimizes the negative log-likelihood loss. Let  $\widehat{F}_d^{(k)}$  be the estimated forest structure, fixing the number of edges at  $k$ ; we want to study conditions under which

$$\mathbb{P} \left( \widehat{F}_d^{(k)} = F_d^{(k)} \right) \rightarrow 1.$$

Let us first consider the population version of the algorithm—if the algorithm cannot recover the best forest  $F_d^{(k)}$  in this ideal case, there is no hope for stable recovery in the data version. The key observation is that the graph selected by the Chow-Liu algorithm only depends on the relative order of the edges with respect to mutual information, not on the specific mutual information values. Let

$$\mathcal{E} = \left\{ \{(i, j), (k, \ell)\} : i < j \text{ and } k < \ell, j \neq \ell \text{ and } i, j, k, \ell \in \{1, \dots, d\} \right\}.$$

The cardinality of  $\mathcal{E}$  is

$$|\mathcal{E}| = O(d^4).$$

Let  $e = (i, j)$  be an edge; the corresponding mutual information associated with  $e$  is denoted as  $I_e$ . If for all  $(e, e') \in \mathcal{E}$ , we have  $I_e \neq I_{e'}$ , the population version of the Chow-Liu algorithm will always obtain the unique solution  $F_d^{(k)}$ . However, this condition is, in a sense, both too weak and too strong. It is too weak because the sample estimates of the mutual information values will only approximate the population values, and could change the relative ordering of some edges. However, the assumption is too strong because, in fact, the relative order of many edge pairs might be changed without affecting the graph selected by the algorithm. For instance, when  $k \geq 2$  and  $I_e$  and  $I_{e'}$  are the largest two mutual information values, it is guaranteed that  $e$  and  $e'$  will both be included in the learned forest  $F_d^{(k)}$  whether  $I_e > I_{e'}$  or  $I_e < I_{e'}$ .

Define the *crucial set*  $\mathcal{J} \subset \mathcal{E}$  to be a set of pairs of edges  $(e, e')$  such that  $I_e \neq I_{e'}$  and flipping the relative order of  $I_e$  and  $I_{e'}$  changes the learned forest structure in the population Chow-Liu algorithm, with positive probability. Here, we assume that the Chow-Liu algorithm randomly selects an edge when a tie occurs.

The cardinality  $|\mathcal{J}|$  of the crucial set is a function of the true density  $p^*(x)$ , and we can expect  $|\mathcal{J}| \ll |\mathcal{E}|$ . The next assumption provides a sufficient condition for the two-stage procedure to be structure selection consistent.

**Assumption 3** *Let the crucial set  $\mathcal{J}$  be defined as before. Suppose that*

$$\min_{((i,j),(k,\ell)) \in \mathcal{J}} |I(X_i; X_j) - I(X_k; X_\ell)| \geq 2L_n$$

where  $L_n = \Omega \left( \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}} \right)$ .

This assumption is strong, but is satisfied in many cases. For example, in a graph with population mutual informations differing by a constant, the assumption holds. Assumption 3 is trivially satisfied if  $\frac{n^{\beta/(1+\beta)}}{\log n + \log d} \rightarrow \infty$ .

**Theorem 12 (Structure selection consistency)** *Let  $F_d^{(k)}$  be the optimal forest within  $\mathcal{P}_d^{(k)}$  that minimizes the negative log-likelihood loss. Let  $\hat{F}_d^{(k)}$  be the estimated forest with  $|E_{\hat{F}_d^{(k)}}| = k$ . Under Assumptions 1, 2, and 3, we have*

$$\mathbb{P} \left( \hat{F}_d^{(k)} = F_d^{(k)} \right) \rightarrow 1$$

as  $n \rightarrow \infty$ .

The proof shows that our method is structure selection consistent as long as the dimension increases as  $d = o(\exp(n^{\beta/(1+\beta)}))$ ; in this case the error decreases at the rate

$$o \left( \exp \left( 4 \log d - c (\log n)^{\frac{1}{1+\beta}} \log d \right) \right).$$

### 4.5 Estimation Consistency

Estimation consistency can be easily established using the structure selection consistency result above. Define the event  $\mathcal{M}_k = \{\hat{F}_d^{(k)} = F_d^{(k)}\}$ . Theorem 12 shows that  $\mathbb{P}(\mathcal{M}_k^c) \rightarrow 0$  as  $n$  goes to infinity.

**Lemma 13** Let  $\widehat{p}_{\widehat{F}_d^{(k)}}$  be the forest-based kernel density estimate for some fixed  $k \in \{0, \dots, d-1\}$ , and let

$$p_{F_d}^* = \arg \min_{q_F \in \mathcal{P}_d^{(k)}} R(q_F).$$

Under the assumptions of Theorem 12,

$$D(p_{F_d}^* \| \widehat{p}_{\widehat{F}_d^{(k)}}) = R(\widehat{p}_{\widehat{F}_d^{(k)}}) - R(p_{F_d}^*)$$

on the event  $\mathcal{M}_k$ .

**Proof** According to Bach and Jordan (2003), for a given forest  $F$  and a target distribution  $p^*(x)$ ,

$$D(p^* \| q_F) = D(p^* \| p_F^*) + D(p_F^* \| q_F) \quad (17)$$

for all distributions  $q_F$  that are supported by  $F$ . We further have

$$D(p^* \| q) = \int_{\mathcal{X}} p^*(x) \log p^*(x) - \int_{\mathcal{X}} p^*(x) \log q(x) dx = \int_{\mathcal{X}} p^*(x) \log p^*(x) dx + R(q) \quad (18)$$

for any distribution  $q$ . Using (17) and (18), and conditioning on the event  $\mathcal{M}_k$ , we have

$$\begin{aligned} D(p_{F_d}^* \| \widehat{p}_{\widehat{F}_d^{(k)}}) &= D(p^* \| \widehat{p}_{\widehat{F}_d^{(k)}}) - D(p^* \| p_{F_d}^*) \\ &= \int_{\mathcal{X}} p^*(x) \log p^*(x) dx + R(\widehat{p}_{\widehat{F}_d^{(k)}}) - \int_{\mathcal{X}} p^*(x) \log p^*(x) dx - R(p_{F_d}^*) \\ &= R(\widehat{p}_{\widehat{F}_d^{(k)}}) - R(p_{F_d}^*), \end{aligned}$$

which gives the desired result. ■

The above lemma combined with Theorem 9 allows us to obtain the following estimation consistency result, the proof of which is omitted.

**Corollary 14 (Estimation consistency)** Under Assumptions 1, 2, and 3, we have

$$D(p_{F_d}^* \| \widehat{p}_{\widehat{F}_d^{(k)}}) = O_P \left( k \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}} + d \sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}} \right)$$

conditioned on the event  $\mathcal{M}_k$ .

## 5. Tree Restricted Forests

We now turn to the problem of estimating forests with restricted tree sizes. As discussed in the introduction, clustering problems motivate the goal of constructing forest structured density estimators where each connected component has a restricted number of edges. But estimating restricted tree size forests can also be useful in model selection for the purpose of risk minimization, since the maximum subtree size can be viewed as an additional complexity parameter.

---

**Algorithm 2** Approximate Max Weight  $t$ -Restricted Forest

---

- 1: **Input** graph  $G$  with positive edge weights, and positive integer  $t \geq 2$ .
- 2: Sort edges in decreasing order of weight.
- 3: Greedily add edges in decreasing order of weight such that
  - (a) the degree of any node is at most  $t + 1$ ;
  - (b) no cycles are formed.

The resulting forest is  $F' = \{T_1, T_2, \dots, T_m\}$ .

- 4: **Output**  $F_t = \cup_j \text{TreePartition}(T_j, t)$ .
- 

**Definition 15** A  $t$ -restricted forest of a graph  $G$  is a subgraph  $F_t$  such that

1.  $F_t$  is the disjoint union of connected components  $\{T_1, \dots, T_m\}$ , each of which is a tree;
2.  $|T_i| \leq t$  for each  $i \leq m$ , where  $|T_i|$  denotes the number of edges in the  $i$ th component.

Given a weight  $w_e$  assigned to each edge of  $G$ , an optimal  $t$ -restricted forest  $F_t^*$  satisfies

$$w(F_t^*) = \max_{F \in \mathcal{F}_t(G)} w(F)$$

where  $w(F) = \sum_{e \in F} w_e$  is the weight of a forest  $F$  and  $\mathcal{F}_t(G)$  denotes the collection of all  $t$ -restricted forests of  $G$ .

For  $t = 1$ , the problem is maximum weighted matching. However, for  $t \geq 7$ , we show that finding an optimal  $t$ -restricted forest is an NP-hard problem; however, this problem appears not to have been previously studied. Our reduction is from Exact 3-Cover (X3C), shown to be NP-complete by Garey and Johnson 1979). In X3C, we are given a set  $X$ , a family  $\mathcal{S}$  of 3-element subsets of  $X$ , and we must choose a subfamily of disjoint 3-element subsets to cover  $X$ . Our reduction constructs a graph with special tree-shaped subgraphs called *gadgets*, such that each gadget corresponds to a 3-element subset in  $\mathcal{S}$ . We show that finding a maximum weight  $t$ -restricted forest on this graph would allow us to then recover a solution to X3C by analyzing how the optimal forest must partition each of the gadgets.

Given the NP-hardness for finding optimal  $t$ -restricted forest, it is of interest to study approximation algorithms for the problem. Our first algorithm is Algorithm 2, which runs in two stages. In the first stage, a forest is greedily constructed in such a way that each node has degree no larger than  $t$  (a property that is satisfied by all  $t$ -restricted forests). However, the trees in the forest may have more than  $t$  edges; hence, in the second stage, each tree in the forest is partitioned in an optimal way by removing edges, resulting in a collection of trees, each of which has size at most  $t$ . The second stage employs a procedure we call `TreePartition` that takes a tree and returns the optimal  $t$ -restricted subforest. `TreePartition` is a divide-and-conquer procedure of Lukes (1974) that finds a carefully chosen set of forest partitions for each child subtree. It then merges these sets with the parent node one subtree at a time. The details of the `TreePartition` procedure are given in Appendix A.

**Theorem 16** *Let  $F_t$  be the output of Algorithm 2, and let  $F_t^*$  be the optimal  $t$ -restricted forest. Then  $w(F_t) \geq \frac{1}{4}w(F_t^*)$ .*

In Appendix A.7, we present a proof of the above result. In that section, we also present an improved approximation algorithm, one based on solving linear programs, that finds a  $t$ -restricted forest  $F'_t$  such that  $w(F'_t) \geq \frac{1}{2}w(F_t^*)$ . Although we cannot guarantee optimality in theory, algorithm 2 performs very well in practice. In Figure 1, we can see that the approximation picks out a  $t$ -restricted forest that is close to optimal among the set of all  $t$ -restricted forests.

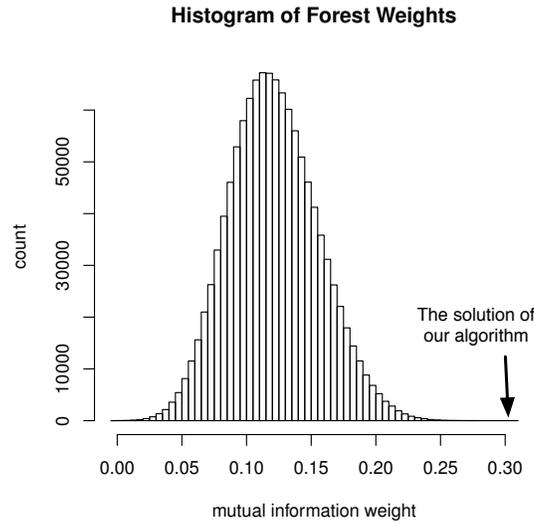


Figure 1: Histogram distribution of weights of all  $t$ -restricted forests on 11 nodes with  $t = 7$ . Edge weights are the mutual informations computed on the training data.

### 5.1 Pruning Based on $t$ -Restricted Forests

For a given  $t$ , after producing an approximate maximum weight  $t$ -restricted forest  $\hat{F}_t$  using  $\mathcal{D}_1$ , we prune away edges using  $\mathcal{D}_2$ . To do so, we first construct a new set of univariate and bivariate kernel density estimates using  $\mathcal{D}_2$ , as before,  $\hat{p}_{n_2}(x_i)$  and  $\hat{p}_{n_2}(x_i, x_j)$ . Recall that we define the “cross-entropies” of the kernel density estimates  $\hat{p}_{n_1}$  for each pair of variables as

$$\begin{aligned} \hat{I}_{n_2, n_1}(X_i, X_j) &= \int \hat{p}_{n_2}(x_i, x_j) \log \frac{\hat{p}_{n_1}(x_i, x_j)}{\hat{p}_{n_1}(x_i) \hat{p}_{n_1}(x_j)} dx_i dx_j \\ &\approx \frac{1}{m^2} \sum_{k=1}^m \sum_{\ell=1}^m \hat{p}_{n_2}(x_{ki}, x_{\ell j}) \log \frac{\hat{p}_{n_1}(x_{ki}, x_{\ell j})}{\hat{p}_{n_1}(x_{ki}) \hat{p}_{n_1}(x_{\ell j})}. \end{aligned}$$

We then eliminate all edges  $(i, j)$  in  $\hat{F}_t$  for which  $\hat{I}_{n_2, n_1}(X_i, X_j) \leq 0$ . For notational simplicity, we denote the resulting pruned forest again by  $\hat{F}_t$ .

---

**Algorithm 3**  $t$ -Restricted Forest Density Estimation

---

- 1: Divide data into two halves  $\mathcal{D}_1$  and  $\mathcal{D}_2$ .
  - 2: Compute kernel density estimators  $\hat{p}_{n_1}$  and  $\hat{p}_{n_2}$  for all pairs and single variable marginals.
  - 3: For all pairs  $(i, j)$  compute  $\hat{I}_{n_1}(X_i, X_j)$  according to (8) and  $\hat{I}_{n_2, n_1}(X_i, X_j)$  according to (10).
  - 4: For  $t = 0, \dots, t_{\text{final}}$  where  $t_{\text{final}}$  is chosen based on the application
    1. Compute or approximate (for  $t \geq 2$ ) the optimal  $t$ -restricted forest  $\hat{F}_t$  using  $\hat{I}_{n_1}$  as edge weights.
    2. Prune  $\hat{F}_t$  to eliminate all edges with negative weights  $\hat{I}_{n_2, n_1}$ .
  - 5: Among all pruned forests  $\hat{p}_{F^t}$ , select  $\hat{t} = \arg \min_{0 \leq t \leq t_{\text{final}}} \hat{R}_{n_2}(\hat{p}_{\hat{F}_t})$ .
- 

To estimate the risk, we simply use  $\hat{R}_{n_2}(\hat{p}_{\hat{F}_t})$  as defined in (9), and select the forest  $\hat{F}_{\hat{t}}$  according to

$$\hat{t} = \arg \min_{0 \leq t \leq d-1} \hat{R}_{n_2}(\hat{p}_{\hat{F}_t}).$$

The resulting procedure is summarized in Algorithm 3.

Using the approximation guarantee and our previous analysis, we have that the population weights of the approximate  $t$ -restricted forest and the optimal forest satisfy the following inequality. We state the result for a general  $c$ -approximation algorithm; for the algorithm given above,  $c = 4$ , but tighter approximations are possible.

**Theorem 17** *Assume the conditions of Theorem 9. For  $t \geq 2$ , let  $\hat{F}_t$  be the forest constructed using a  $c$ -approximation algorithm, and let  $F_t^*$  be the optimal forest; both constructed with respect to finite sample edge weights  $\hat{w}_{n_1} = \hat{I}_{n_1}$ . Then*

$$w(\hat{F}_t) \geq \frac{1}{c} w(F_t^*) + O_P \left( (k^* + \hat{k}) \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}} \right)$$

where  $\hat{k}$  and  $k^*$  are the number of edges in  $\hat{F}_t$  and  $F_t^*$ , respectively, and  $w$  denotes the population weights, given by the mutual information.

As seen below, although the approximation algorithm has weaker theoretical guarantees, it outperforms other approaches in experiments.

## 6. Experimental Results

In this section, we report numerical results on both synthetic data sets and microarray data. We mainly compare the forest density estimator with sparse Gaussian graphical models, fitting a multivariate Gaussian with a sparse inverse covariance matrix. The sparse Gaussian models are estimated using the graphical lasso algorithm (glasso) of Friedman et al. (2007), which is a refined version of an algorithm first derived by Banerjee et al. (2008). Since the glasso typically results in a large parameter bias as a consequence of the  $\ell_1$  regularization, we also compare with a method that we call

the *refit glasso*, which is a two-step procedure—in the first step, a sparse inverse covariance matrix is obtained by the glasso; in the second step, a Gaussian model is refit without  $\ell_1$  regularization, but enforcing the sparsity pattern obtained in the first step.

To quantitatively compare the performance of these estimators, we calculate the log-likelihood of all methods on a held-out data set  $\mathcal{D}_2$ . With  $\hat{\mu}_{n_1}$  and  $\hat{\Omega}_{n_1}$  denoting the estimates from the Gaussian model, the held-out log-likelihood can be explicitly evaluated as

$$\ell_{\text{gauss}} = -\frac{1}{n_2} \sum_{s \in \mathcal{D}_2} \left\{ \frac{1}{2} (X^{(s)} - \hat{\mu}_{n_1})^F \hat{\Omega}_{n_1} (X^{(s)} - \hat{\mu}_{n_1}) + \frac{1}{2} \log \left( \frac{|\hat{\Omega}_{n_1}|}{(2\pi)^d} \right) \right\}.$$

For a given tree structure  $\hat{F}$ , the held-out log-likelihood for the forest density estimator is

$$\ell_{\text{fde}} = \frac{1}{n_2} \sum_{s \in \mathcal{D}_2} \log \left( \prod_{(i,j) \in E(\hat{F})} \frac{\hat{p}_{n_1}(X_i^{(s)}, X_j^{(s)})}{\hat{p}_{n_1}(X_i^{(s)}) \hat{p}_{n_1}(X_j^{(s)})} \prod_{k \in V_{\hat{F}}} \hat{p}_{n_1}(X_k^{(s)}) \right),$$

where  $\hat{p}_{n_1}(\cdot)$  are the corresponding kernel density estimates, using a Gaussian kernel with plug-in bandwidths.

Since the held-out log-likelihood of the forest density estimator is indexed by the number of edges included in the tree, while the held-out log-likelihoods of the glasso and the refit glasso are indexed by a continuously varying regularization parameter, we need to find a way to calibrate them. To address this issue, we plot the held-out log-likelihood of the forest density estimator as a step function indexed by the tree size. We then run the full path of the glasso and discretize it according to the corresponding sparsity level, that is, how many edges are selected for each value of the regularization parameter. The size of the forest density estimator and the sparsity level of the glasso (and the refit glasso) can then be aligned for a fair comparison.

### 6.1 Synthetic Data

We use a procedure to generate high dimensional Gaussian and non-Gaussian data which are consistent with an undirected graph. We generate high dimensional graphs that contain cycles, and so are not forests. In dimension  $d = 100$ , we sample  $n_1 = n_2 = 400$  data points from a multivariate Gaussian distribution with mean vector  $\mu = (0.5, \dots, 0.5)$  and inverse covariance matrix  $\Omega$ . The diagonal elements of  $\Omega$  are all 62. We then randomly generate many connected subgraphs containing no more than eight nodes each, and set the corresponding non-diagonal elements in  $\Omega$  at random, drawing values uniformly from  $-30$  to  $-10$ . To obtain non-Gaussian data, we simply transform each dimension of the data by its empirical distribution function; such a transformation preserves the graph structure but the joint distribution is no longer Gaussian (see Liu et al., 2009).

To calculate the pairwise mutual information  $\hat{I}(X_i; X_j)$ , we need to numerically evaluate two-dimensional integrals. We first rescale the data into  $[0, 1]^d$  and calculate the kernel density estimates on a grid of points; we choose  $m = 128$  evaluation points  $x_i^{(1)} < x_i^{(2)} < \dots < x_i^{(m)}$  for each dimension  $i$ , and then evaluate the bivariate and the univariate kernel density estimates on this grid.

There are three different kernel density estimates that we use—the bivariate kde, the univariate kde, and the marginalized bivariate kde. Specifically, the bivariate kernel density estimate on  $x_i, x_j$

based on the observations  $\{X_i^{(s)}, X_j^{(s)}\}_{s \in \mathcal{D}_1}$  is defined as

$$\widehat{p}(x_i, x_j) = \frac{1}{n_1} \sum_{s \in \mathcal{D}_1} \frac{1}{h_{2i} h_{2j}} K\left(\frac{X_i^{(s)} - x_i}{h_{2i}}\right) K\left(\frac{X_j^{(s)} - x_j}{h_{2j}}\right),$$

using a product kernel. The bandwidths  $h_{2i}, h_{2j}$  are chosen as

$$h_{2k} = 1.06 \cdot \min\left\{\widehat{\sigma}_k, \frac{\widehat{q}_{k,0.75} - \widehat{q}_{k,0.25}}{1.34}\right\} \cdot n^{-1/(2\beta+2)},$$

where  $\widehat{\sigma}_k$  is the sample standard deviation of  $\{X_k^{(s)}\}_{s \in \mathcal{D}_1}$  and  $\widehat{q}_{k,0.75}, \widehat{q}_{k,0.25}$  are the 75% and 25% sample quantiles of  $\{X_k^{(s)}\}_{s \in \mathcal{D}_1}$ .

In all the experiments, we set  $\beta = 2$ , such a choice of  $\beta$  and the ‘‘plug-in’’ bandwidth  $h_{2k}$  (and  $h_{1k}$  in the following) is a very common practice in nonparametric Statistics. For more details, see Fan and Gijbels (1996) and Tsybakov (2008).

Given an evaluation point  $x_k$ , the univariate kernel density estimate  $\widehat{p}(x_k)$  based on the observations  $\{X_k^{(s)}\}_{s \in \mathcal{D}_1}$  is defined as

$$\widehat{p}(x_k) = \frac{1}{n_1} \sum_{s \in \mathcal{D}_1} \frac{1}{h_{1k}} K\left(\frac{X_k^{(s)} - x_k}{h_{1k}}\right),$$

where  $h_{1k} > 0$  is defined as

$$h_{1k} = 1.06 \cdot \min\left\{\widehat{\sigma}_k, \frac{\widehat{q}_{k,0.75} - \widehat{q}_{k,0.25}}{1.34}\right\} \cdot n^{-1/(2\beta+1)}.$$

Finally, the marginal univariate kernel density estimate  $\widehat{p}_M(x_k)$  based on the observations  $\{X_k^{(s)}\}_{s \in \mathcal{D}_1}$  is defined by integrating the irrelevant dimension out of the bivariate kernel density estimates  $\widehat{p}(x_j, x_k)$  on the unit square  $[0, 1]^2$ . Thus,

$$\widehat{p}_M(x_k) = \frac{1}{m-1} \sum_{\ell=1}^m \widehat{p}(x_j^{(\ell)}, x_k).$$

With the above definitions of the bivariate and univariate kernel density estimates, we consider estimating the mutual information  $I(X_i; X_j)$  in three different ways, depending on which estimates for the univariate densities are employed.

$$\begin{aligned} \widehat{I}_{\text{fast}}(X_i, X_j) &= \frac{1}{(m-1)^2} \sum_{k'=1}^m \sum_{\ell=1}^m \widehat{p}(x_i^{(k')}, x_j^{(\ell)}) \log \widehat{p}(x_i^{(k')}, x_j^{(\ell)}) \\ &\quad - \frac{1}{m-1} \sum_{k'=1}^m \widehat{p}(x_i^{(k')}) \log \widehat{p}(x_i^{(k')}) - \frac{1}{m-1} \sum_{\ell=1}^m \widehat{p}(x_j^{(\ell)}) \log \widehat{p}(x_j^{(\ell)}) \\ \widehat{I}_{\text{medium}}(X_i, X_j) &= \frac{1}{(m-1)^2} \sum_{k'=1}^m \sum_{\ell=1}^m \widehat{p}(x_i^{(k')}, x_j^{(\ell)}) \log \frac{\widehat{p}(x_i^{(k')}, x_j^{(\ell)})}{\widehat{p}(x_i^{(k')}) \widehat{p}(x_j^{(\ell)})}. \\ \widehat{I}_{\text{slow}}(X_i, X_j) &= \frac{1}{(m-1)^2} \sum_{k'=1}^m \sum_{\ell=1}^m \widehat{p}(x_i^{(k')}, x_j^{(\ell)}) \log \widehat{p}(x_i^{(k')}, x_j^{(\ell)}) - \\ &\quad - \frac{1}{m-1} \sum_{k'=1}^m \widehat{p}_M(x_i^{(k')}) \log \widehat{p}_M(x_i^{(k')}) - \frac{1}{m-1} \sum_{\ell=1}^m \widehat{p}_M(x_j^{(\ell)}) \log \widehat{p}_M(x_j^{(\ell)}). \end{aligned}$$

The terms “fast,” “medium” and “slow” refer to the theoretical statistical rates of convergence of the estimators. The “fast” estimate uses one-dimensional univariate kernel density estimators wherever possible. The “medium” estimate uses the one-dimensional kernel density estimates in the denominator of  $p(x_i, x_j)/(p(x_i)p(x_j))$ , but averages with respect to the bivariate density. Finally, the “slow” estimate marginalizes the bivariate densities to estimate the univariate densities. While the rate of convergence is the two-dimensional rate, the “slow” estimate ensures the consistency of the bivariate and univariate densities.

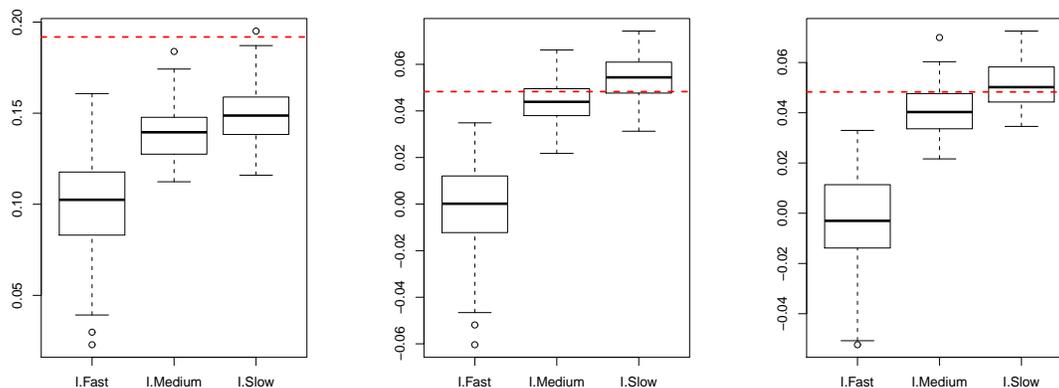


Figure 2: (Gaussian example) Boxplots of  $\hat{I}_{\text{fast}}$ ,  $\hat{I}_{\text{medium}}$ , and  $\hat{I}_{\text{slow}}$  on three different pairs of variables. The red-dashed horizontal lines represent the population values.

Figure 2 compares  $\hat{I}_{\text{fast}}$ ,  $\hat{I}_{\text{medium}}$ , and  $\hat{I}_{\text{slow}}$  on different pairs of variables. The boxplots are based on 100 trials. Compared to the ground truth, which can be computed exactly in the Gaussian case, we see that the performance of  $\hat{I}_{\text{medium}}$  and  $\hat{I}_{\text{slow}}$  is better than that of  $\hat{I}_{\text{fast}}$ . This is due to the fact that simply replacing the population density with a “plug-in” version can lead to biased estimates; in fact,  $\hat{I}_{\text{fast}}$  is not even guaranteed to be non-negative. In what follows, we employ  $\hat{I}_{\text{medium}}$  for all the calculations, due to its ease of computation and good finite sample performance. Figure 3 compares the bivariate fits of the kernel density estimates and the Gaussian models over four edges. For the Gaussian fits of each edge, we directly calculate the bivariate sample covariance and sample mean and plug them into the bivariate Gaussian density function. From the perspective and contour plots, we see that the bivariate kernel density estimates provide reasonable fits for these bivariate components.

A typical run showing the held-out log-likelihood and estimated graphs is provided in Figure 4. We see that for the Gaussian data, the refit glasso has a higher held-out log-likelihood than the forest density estimator and the glasso. This is expected, since the Gaussian model is correct. For very sparse models, however, the performance of the glasso is worse than that of the forest density estimator, due to the large parameter bias resulting from the  $\ell_1$  regularization. We also observe an efficiency loss in the nonparametric forest density estimator, compared to the refit glasso. The graphs are automatically selected using the held-out log-likelihood, and we see that the nonparametric forest-based kernel density estimator tends to select a sparser model, while the parametric

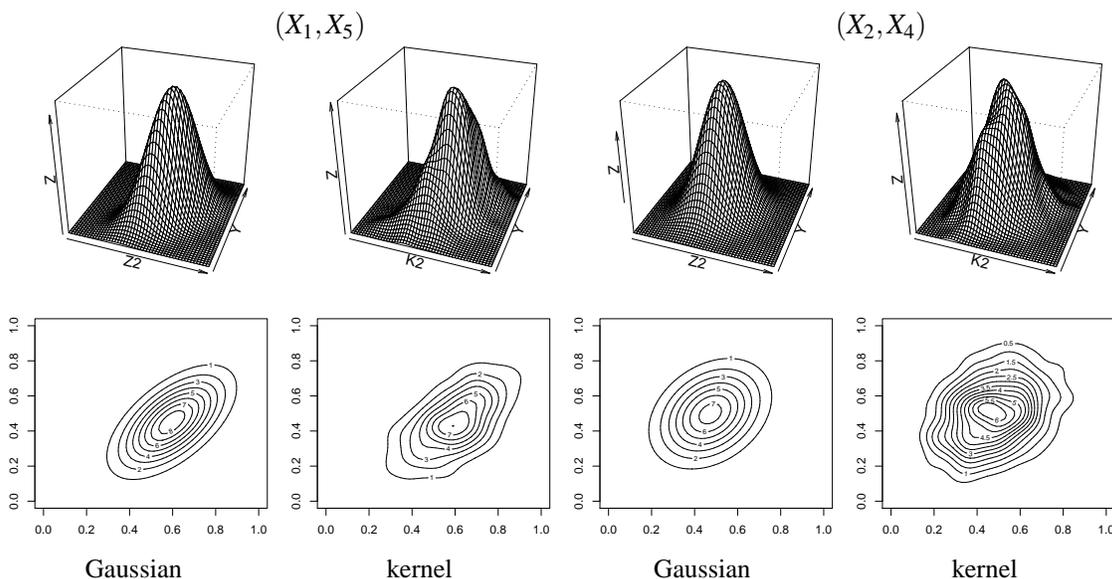


Figure 3: Perspective and contour plots of the bivariate Gaussian fits vs. the kernel density estimates for two edges of a Gaussian graphical model.

Gaussian models tend to overselect. This observation is new and is quite typical in our simulations. Another observation is that the held-out log-likelihood curve of the glasso becomes flat for less sparse models but never goes down. This suggests that the held-out log-likelihood is not a good model selection criterion for the glasso. For the non-Gaussian data, even though the refit glasso results in a reasonable graph, the forest density estimator performs much better in terms of held-out log-likelihood risk and graph estimation accuracy.

To compare with  $t$ -restricted forests, we generated additional Gaussian and non-Gaussian synthetic data as before except on a different graph structure. In Figure 5, we use 400 training examples while varying the size of heldout data to compare the log-likelihoods of four different methods; the log-likelihood is evaluated on a third large data set. In Figure 6, we consider only non-Gaussian data, use 400 training data and 400 heldout data, and generate graphs with best heldout log-likelihood across the four methods. We compute bandwidth, heldout log-likelihood, and mutual information same as before.

We observe that although creating a maximum spanning tree (MST) on the held-out data is asymptotically optimal; it can perform quite poorly. Unless there are copious amount of heldout data, held-out MST overfits on the heldout data and tend to give large graphs; in contrast,  $t$ -restricted forest has the weakest theoretical guarantee but it gives the best log-likelihood and produces sparser graphs. It is not surprising to note that MST on heldout data improves as heldout data size increases. Somewhat surprisingly though, Training-MST-with-pruning and  $t$ -restricted forest appear to be insensitive to the heldout data size.

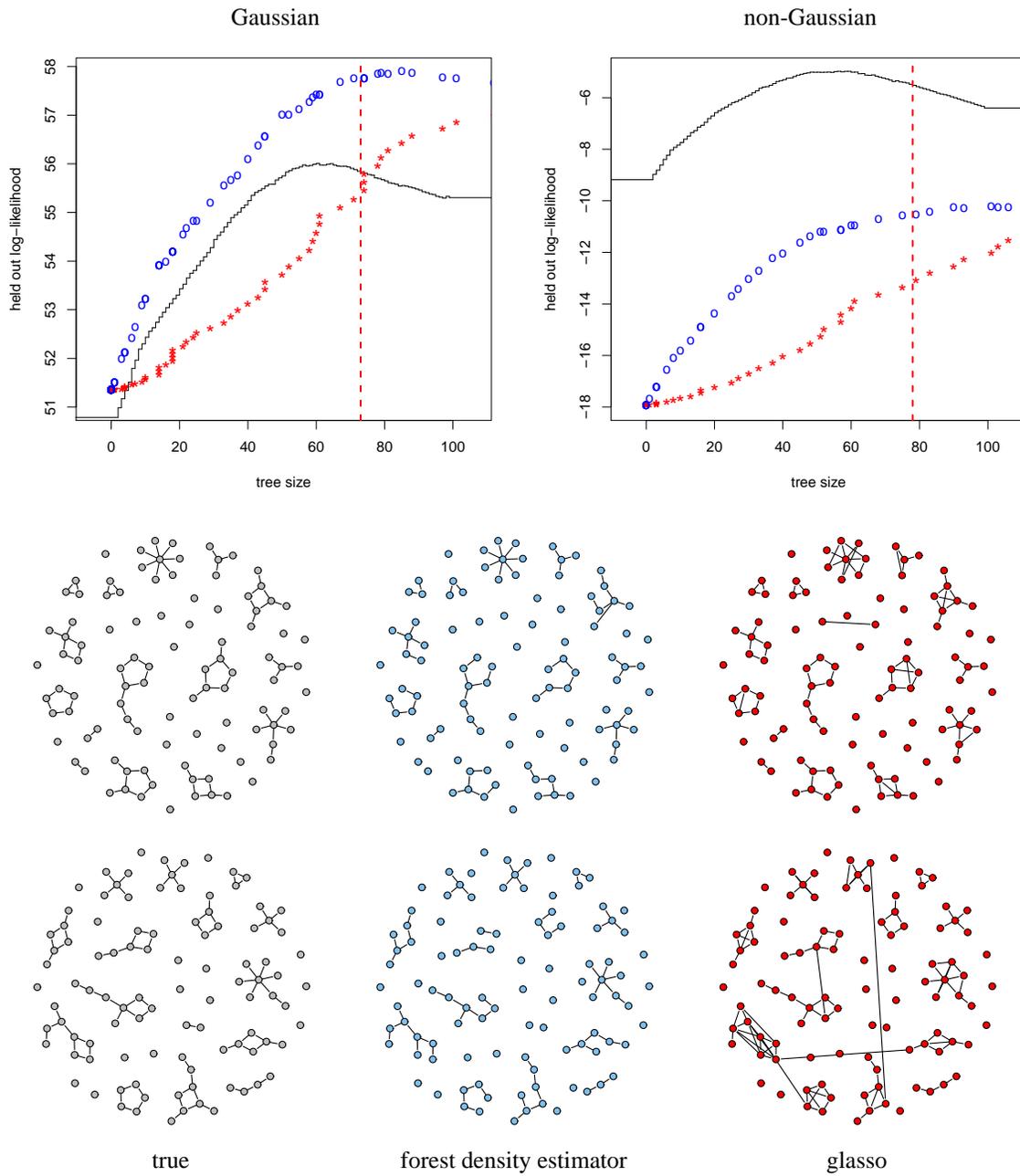


Figure 4: Synthetic data. Top-left Gaussian, and top-right non-Gaussian: Held-out log-likelihood plots of the forest density estimator (black step function), glasso (red stars), and refit glasso (blue circles), the vertical dashed red line indicates the size of the true graph. Bottom plots show the true and estimated graphs for the Gaussian (second row) and non-Gaussian data (third row).

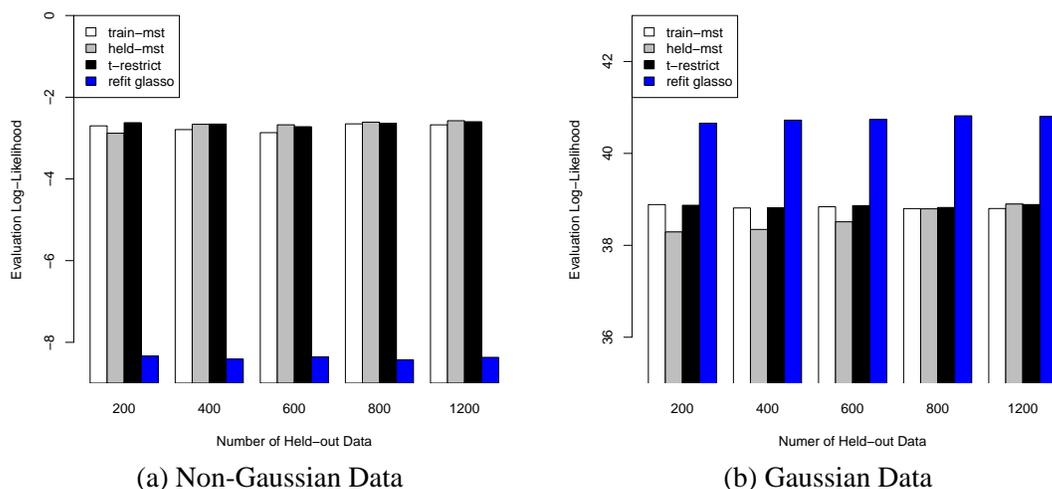


Figure 5: Log-likelihood comparison of various methods: (left white) MST on Training Data with Pruning (gray) MST on Heldout Data (black) t-Restricted Graph (blue) Refit Glasso

## 6.2 Microarray Data

In this example, we study the empirical performance of the algorithms on a microarray dataset.

### 6.2.1 ARABIDOPSIS THALIANA DATA

We consider a data set based on Affymetrix GeneChip microarrays for the plant *Arabidopsis thaliana*, (Wille et al., 2004). The sample size is  $n = 118$ . The expression levels for each chip are pre-processed by a log-transformation and standardization. A subset of 40 genes from the isoprenoid pathway are chosen, and we study the associations among them using the glasso, the refit glasso, and the forest-based kernel density estimator.

From the held-out log-likelihood curves in Figure 7, we see that the tree-based kernel density estimator has a better generalization performance than the glasso and the refit glasso. This is not surprising, given that the true distribution of the data is not Gaussian. Another observation is that for the tree-based kernel density estimator, the held-out log-likelihood curve achieves a maximum when there are only 35 edges in the model. In contrast, the held-out log-likelihood curves of the glasso and refit glasso achieve maxima when there are around 280 edges and 100 edges respectively, while their predictive estimates are still inferior to those of the tree-based kernel density estimator.

Figure 7 also shows the estimated graphs for the tree-based kernel density estimator and the glasso. The graphs are automatically selected based on held-out log-likelihood. The two graphs are clearly different; it appears that the nonparametric tree-based kernel density estimator has the potential to provide different biological insights than the parametric Gaussian graphical model.

### 6.2.2 HAPMAP DATA

This data set comes from Nayak et al. (2009). The data set contains Affymetrix chip measured expression levels of 4238 genes for 295 normal subjects in the *Centre d’Etude du Polymorphisme*

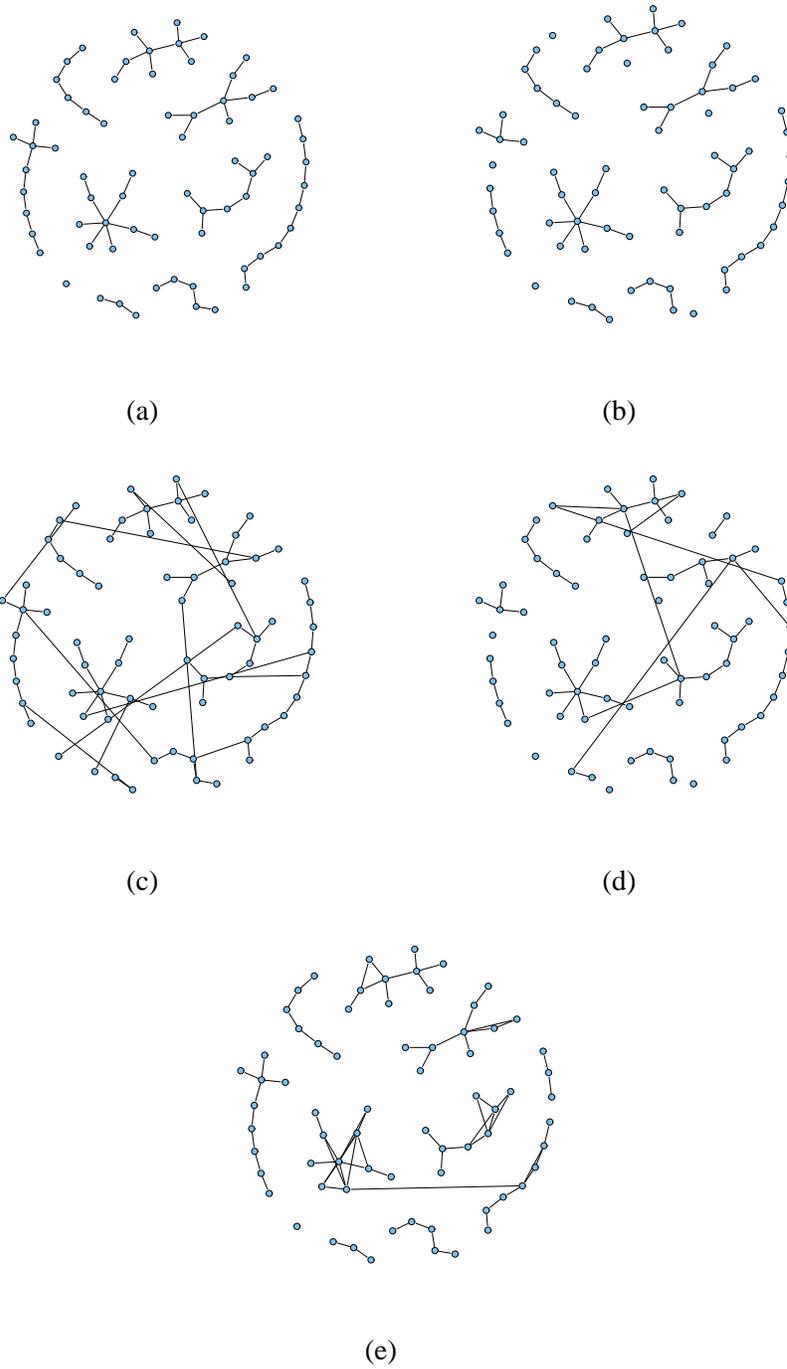


Figure 6: Graphs generated on non-Gaussian Data: (a) True Graph, (b) t-Restricted Forest (c) MST on Heldout Data (d) MST on Training Data with Pruning (e) Refit Glasso

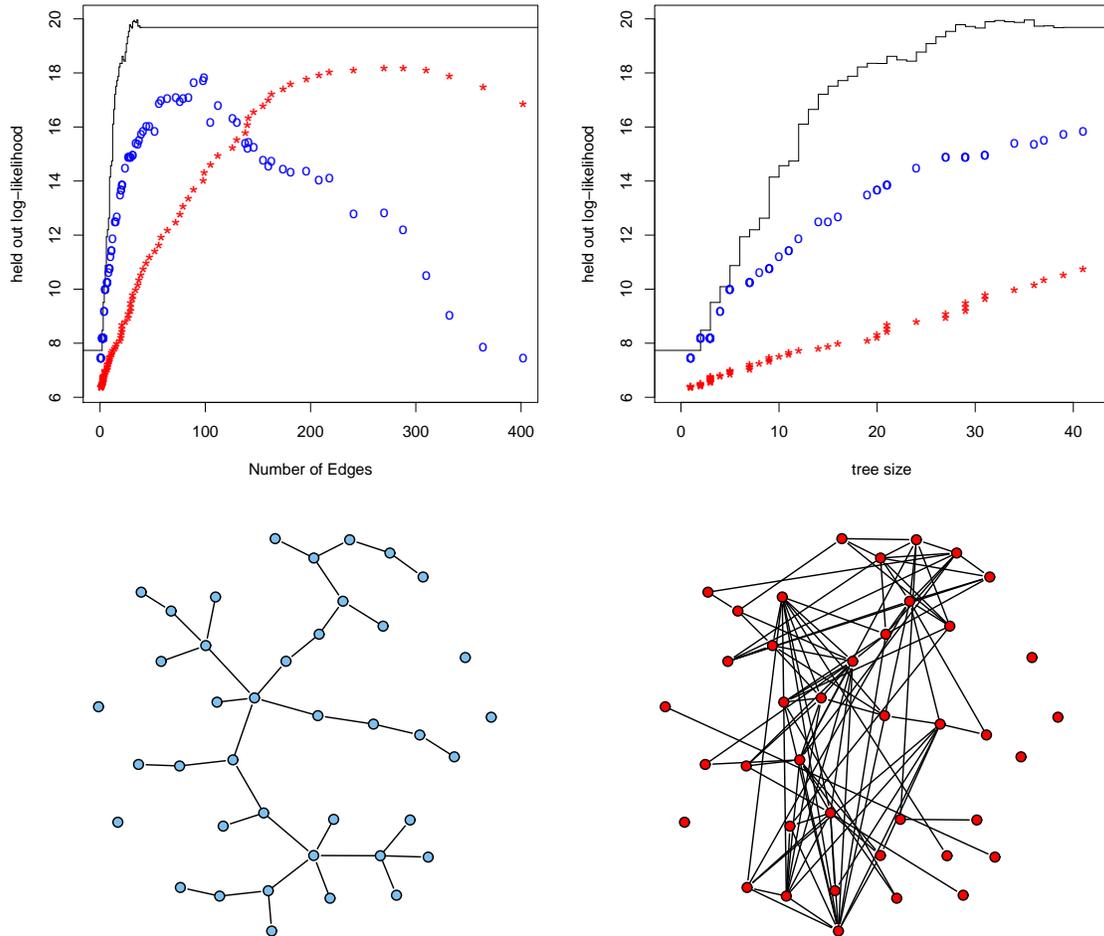


Figure 7: Results on microarray data. Top: held-out log-likelihood (left) and its zoom-in (right) of the tree-based kernel density estimator (black step function), glasso (red stars), and refit glasso (blue circles). Bottom: estimated graphs using the tree-based estimator (left) and glasso (right).

*Humain* (CEPH) and the International HapMap collections. The 295 subjects come from four different groups: 148 unrelated grandparents in the CEPH-Utah pedigrees, 43 Han Chinese in Beijing, 44 Japanese in Tokyo, and 60 Yoruba in Ibadan, Nigeria. Since we want to find common network patterns across different groups of subjects, we pooled the data together into a  $n = 295$  by  $d = 4238$  numerical matrix.

We estimate the full 4238 node graph using both the forest density estimator (described in Section 3.1 and 3.2) and the Meinshausen-Bühlmann neighborhood search method as proposed in Meinshausen and Bühlmann (2006) with regularization parameter chosen to give it about same number as edges as the forest graph.

To construct the kernel density estimates  $\hat{p}(x_i, x_j)$  we use an array of Nvidia graphical processing units (GPU) to parallelize the computation over the pairs of variables  $X_i$  and  $X_j$ . We discretize the domain of  $(X_i, X_j)$  into a  $128 \times 128$  grid, and correspondingly employ  $128 \times 128$  parallel cells in the GPU array, taking advantage of shared memory in CUDA. Parallelizing in this way increases the total performance by approximately a factor of 40, allowing the experiment to complete in a day.

The forest density estimated graph reveals one strongly connected component of more than 3000 genes and various isolated genes; this is consistent with the analysis in Nayak et al. (2009) and is realistic for the regulatory system of humans. The Gaussian graph contains similar component structure, but the set of edges differs significantly. We also ran the  $t$ -restricted forest algorithm for  $t = 2000$  and it successfully separates the giant component into three smaller components. For visualization purposes, in Figure 8, we show only a 934 gene subgraph of the strongly connected component among the full 4238 node graphs we estimated. More detailed analysis of the biological implications of this work will left as a future study.

## 7. Conclusion

We have studied forest density estimation for high dimensional data. Forest density estimation skirts the curse of dimensionality by restricting to undirected graphs without cycles, while allowing fully nonparametric marginal densities. The method is computationally simple, and the optimal size of the forest can be robustly selected by a data-splitting scheme. We have established oracle properties and rates of convergence for function estimation in this setting. Our experimental results compared the forest density estimator to the sparse Gaussian graphical model in terms of both predictive risk and the qualitative properties of the estimated graphs for human gene expression array data. Together, these results indicate that forest density estimation can be a useful tool for relaxing the normality assumption in graphical modeling.

## Acknowledgments

The research reported here was carried out at Carnegie Mellon University and was supported in part by NSF grant CCF-0625879, AFOSR contract FA9550-09-1-0373, and a grant from Google.

## Appendix A. Proofs

In the following, we present the detailed proofs of all the technical results.

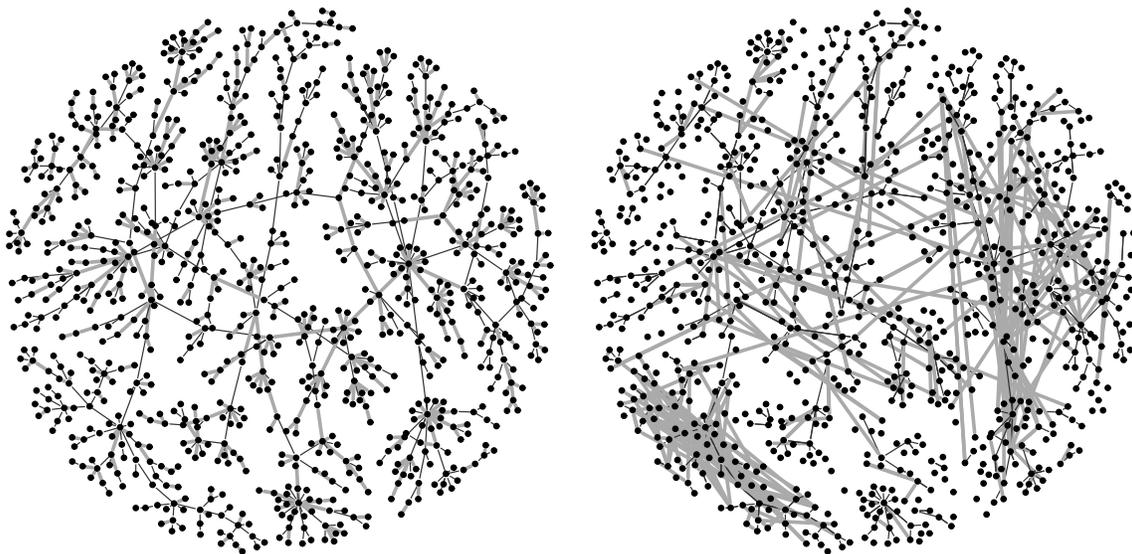


Figure 8: A 934 gene subgraph of the full estimated 4238 gene network. Left: estimated forest graph. Right: estimated Gaussian graph. The bold gray edges in the forest graph are missing from the Gaussian graph and vice versa; the thin black edges are shared by both graphs. Note that the layout of the genes is the same for both graphs.

### A.1 Proof of Lemma 8

We only need to consider the more complicated bivariate case (14); the result in (15) follows from the same line of proof. First, given the assumptions, the following lemma can be obtained by an application of Corollary 2.2 of Giné and Guillou (2002). For a detailed proof, see Rinaldo and Wasserman (2010).

**Lemma 18** (Giné and Guillou, 2002) *Let  $\hat{p}$  be a bivariate kernel density estimate using a kernel  $K(\cdot)$  for which Assumption 2 holds and suppose that*

$$\sup_{t \in \mathcal{X}^2} \sup_{h_2 > 0} \int_{\mathcal{X}^2} K_2^2(u) p^*(t - uh_2) du \leq D < \infty. \quad (19)$$

1. *Let the bandwidth  $h_2$  be fixed. Then there exist constants  $L > 0$  and  $C > 0$ , which depend only on the VC characteristics of  $\mathcal{F}_2$  in (11), such that for any  $c_1 \geq C$  and  $0 < \varepsilon \leq c_1 D / \|K_2\|_\infty$ , there exists  $n_0 > 0$  which depends on  $\varepsilon$ ,  $D$ ,  $\|K_2\|_\infty$  and the VC characteristics of  $K_2$ , such that for all  $n \geq n_0$ ,*

$$\mathbb{P} \left( \sup_{u \in \mathcal{X}^2} |\hat{p}(u) - \mathbb{E}\hat{p}(u)| > 2\varepsilon \right) \leq L \exp \left\{ - \frac{1}{L} \frac{\log(1 + c_1/(4L))}{c_1} \frac{nh_2^2 \varepsilon^2}{D} \right\}. \quad (20)$$

2. Let  $h_2 \rightarrow 0$  in such a way that  $nh_2^2/\log h_2 \rightarrow \infty$ , and let  $\varepsilon \rightarrow 0$  so that

$$\varepsilon = \Omega \left( \sqrt{\frac{\log r_n}{nh_2^2}} \right), \quad (21)$$

where  $r_n = \Omega(h_2^{-1})$ . Then (20) holds for sufficiently large  $n$ .

From (D2) in Assumption 1 and (K1) in Assumption 2, it is easy to see that (19) is satisfied. Also, since

$$h_2 \asymp \left( \frac{\log n}{n} \right)^{\frac{1}{2+2\beta}},$$

it is clear that  $nh_2^2/\log h_2 \rightarrow \infty$ . Part 2 of Lemma 18 shows that there exist  $c_2$  and  $c_3$  such that

$$\mathbb{P} \left( \sup_{(x_i, x_j) \in \mathcal{X}_i \times \mathcal{X}_j} |\widehat{p}(x_i, x_j) - \mathbb{E}\widehat{p}(x_i, x_j)| \geq \frac{\varepsilon}{2} \right) \leq c_2 \exp \left( -c_3 n^{\frac{\beta}{1+\beta}} (\log n)^{\frac{1}{1+\beta}} \varepsilon^2 \right) \quad (22)$$

for all  $\varepsilon$  satisfying (21).

This shows that for any  $i, j \in \{1, \dots, d\}$  with  $i \neq j$ , the bivariate kernel density estimate  $\widehat{p}(x_i, x_j)$  is uniformly close to  $\mathbb{E}\widehat{p}(x_i, x_j)$ . Note that  $\mathbb{E}\widehat{p}(x_i, x_j)$  can be written as

$$\mathbb{E}\widehat{p}(x_i, x_j) = \int \frac{1}{h_2^2} K \left( \frac{u_i - x_i}{h_2} \right) K \left( \frac{v_j - x_j}{h_2} \right) p^*(u_i, v_j) du_i dv_j.$$

The next lemma, from Rigollet and Vert (2009), provides a uniform deviation bound on the bias term  $\mathbb{E}\widehat{p}(x_i, x_j) - p^*(x_i, x_j)$ .

**Lemma 19** (Rigollet and Vert, 2009) *Under (D1) in Assumption 1 and (K3) in Assumption 2, we have*

$$\sup_{(x_i, x_j) \in \mathcal{X}_i \times \mathcal{X}_j} |\mathbb{E}\widehat{p}(x_i, x_j) - p^*(x_i, x_j)| \leq L_1 h_2^\beta \int_{\mathcal{X}^2} (u^2 + v^2)^{\beta/2} K(u)K(v) dudv.$$

where  $L$  is defined in (D1) of Assumption 1.

Let  $c_4 = L_1 \int_{\mathcal{X}^2} (u^2 + v^2)^{\beta/2} K(u)K(v) dudv$ . From the discussion of Example 6.1 in Rigollet and Vert (2009) and (K1) in Assumption 2, we know that  $c_4 < \infty$  and only depends on  $K$  and  $\beta$ . Therefore

$$\mathbb{P} \left( \sup_{(x_i, x_j) \in \mathcal{X}_i \times \mathcal{X}_j} |p^*(x_i, x_j) - \mathbb{E}\widehat{p}(x_i, x_j)| \geq \frac{\varepsilon}{2} \right) = 0 \quad (23)$$

for  $\varepsilon \geq 4c_4 h_2^\beta$ .

The desired result in Lemma 8 is an exponential probability inequality showing that  $\widehat{p}(x_i, x_j)$  is close to  $p^*(x_i, x_j)$ . To obtain this, we use a union bound:

$$\begin{aligned} & \mathbb{P} \left( \max_{(i,j) \in \{1, \dots, d\} \times \{1, \dots, d\}} \sup_{(x_i, x_j) \in \mathcal{X}_i \times \mathcal{X}_j} |\widehat{p}(x_i, x_j) - p^*(x_i, x_j)| \geq \varepsilon \right) \\ & \leq d^2 \mathbb{P} \left( \sup_{(x_i, x_j) \in \mathcal{X}_i \times \mathcal{X}_j} |\widehat{p}(x_i, x_j) - \mathbb{E}\widehat{p}(x_i, x_j)| \geq \frac{\varepsilon}{2} \right) \\ & \quad + d^2 \mathbb{P} \left( \sup_{(x_i, x_j) \in \mathcal{X}_i \times \mathcal{X}_j} |p^*(x_i, x_j) - \mathbb{E}\widehat{p}(x_i, x_j)| \geq \frac{\varepsilon}{2} \right). \end{aligned} \quad (24)$$

The first result follows from (22) and (24).

Choosing

$$\varepsilon = \Omega \left( 4c_4 \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}} \right),$$

the result directly follows by combining (22) and (23)

### A.2 Proof of Theorem 9

First, from (D2) in Assumption 1 and Lemma 8, we have for any  $i \neq j$ ,

$$\max_{(i,j) \in \{1, \dots, d\} \times \{1, \dots, d\}} \sup_{(x_i, x_j) \in \mathcal{X}_i \times \mathcal{X}_j} \left( \frac{\widehat{p}(x_i, x_j)}{p^*(x_i, x_j)} - 1 \right) = O_P \left( \sqrt{\frac{\log n + \log d}{n^{\beta/(\beta+1)}}} \right).$$

The next lemma bounds the deviation of  $\widehat{R}(\widehat{p}_F)$  from  $R(p_F^*)$  over different choices of  $F \in \mathcal{F}_d$  with  $|E(F)| \leq k$ . In the following, we let

$$\mathcal{F}_d^{(k)} = \{F \in \mathcal{F}_d : |E(F)| \leq k\}$$

denote the family of  $d$ -node forests with no more than  $k$  edges.

**Lemma 20** *Under the assumptions of Theorem 9, we have*

$$\sup_{F \in \mathcal{F}_d^{(k)}} |\widehat{R}(\widehat{p}_F) - R(p_F^*)| = O_P \left( k \sqrt{\frac{\log n + \log d}{n^{\beta/(\beta+1)}}} + d \sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}} \right).$$

**Proof** For any  $F \in \mathcal{F}_d^{(k)}$ , we have

$$\begin{aligned} & |\widehat{R}(\widehat{p}_F) - R(p_F^*)| \\ & \leq \underbrace{\left| \sum_{(i,j) \in E(F)} \left( \int_{\mathcal{X}_i \times \mathcal{X}_j} p^*(x_i, x_j) \log p^*(x_i, x_j) dx_i dx_j - \int_{\mathcal{X}_i \times \mathcal{X}_j} \widehat{p}(x_i, x_j) \log \widehat{p}(x_i, x_j) dx_i dx_j \right) \right|}_{A_1(F)} \\ & \quad + \underbrace{\left| \sum_{k \in V} (\deg_F(k) - 1) \left( \int_{\mathcal{X}_k} p^*(x_k) \log p^*(x_k) dx_k - \int_{\mathcal{X}_k} \widehat{p}(x_k) \log \widehat{p}(x_k) dx_k \right) \right|}_{A_2(F)} \end{aligned}$$

where  $\deg_F(k)$  is the degree of node  $k$  in  $F$ . Let  $\varepsilon \geq 4c_4 h_2^\beta$  and let  $\Omega_n$  be the event that

$$\max_{(i,j) \in \{1, \dots, d\} \times \{1, \dots, d\}} \sup_{(x_i, x_j) \in \mathcal{X}_i \times \mathcal{X}_j} |\widehat{p}(x_i, x_j) - p^*(x_i, x_j)| \leq \varepsilon.$$

By Lemma 8,  $\Omega_n$  holds except on a set of probability at most

$$c_2 d^2 \exp \left( -c_3 n^{\frac{\beta}{1+\beta}} (\log n)^{\frac{1}{1+\beta}} \varepsilon^2 \right).$$

From (D2) in Assumption 1, and from the fact that  $|\log(1 + u)| \leq 2|u|$  for all small  $u$ , we have that, on the event  $\Omega_n$ ,

$$\sup_{F \in \mathcal{F}_d^{(k)}} A_1(F) \leq ck\varepsilon.$$

By choosing  $\varepsilon = \Omega\left(4c_4 \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}}\right)$  we conclude that

$$\sup_{F \in \mathcal{F}_d^{(k)}} A_1(F) = O_P\left(k \sqrt{\frac{\log n + \log d}{n^{\beta/(\beta+1)}}}\right).$$

By a similar argument and using the fact that  $\sum_k |\deg_F(k) - 1| = O(d)$ , we have

$$\sup_{F \in \mathcal{F}_d^{(k)}} A_2(F) = O_P\left(d \sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}}\right).$$

■

The next auxiliary lemma is also needed to obtain the main result. It shows that  $\widehat{R}(\widehat{p}_F)$  does not deviate much from  $R(\widehat{p}_F)$  uniformly over different choices of  $F \in \mathcal{F}_d^{(k)}$ .

**Lemma 21** *Under the assumptions of Theorem 9, we have*

$$\sup_{F \in \mathcal{F}_d^{(k)}} |R(\widehat{p}_F) - \widehat{R}(\widehat{p}_F)| = O_P\left(k \sqrt{\frac{\log n + \log d}{n^{\beta/(\beta+1)}}} + d \sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}}\right).$$

**Proof** The argument is similar to the proof of Lemma 20. ■

The proof of the main theorem follows by repeatedly applying the previous two lemmas. As in Proposition 2, with

$$p_{F_d}^* = \arg \min_{q_F \in \mathcal{P}_d^{(k)}} R(q_F),$$

we have

$$\begin{aligned}
 & R(\widehat{p}_{\widehat{F}_d^{(k)}}) - R(p_{F_d^{(k)}}^*) \\
 &= R(\widehat{p}_{\widehat{F}_d^{(k)}}) - \widehat{R}(\widehat{p}_{\widehat{F}_d^{(k)}}) + \widehat{R}(\widehat{p}_{\widehat{F}_d^{(k)}}) - R(p_{F_d^{(k)}}^*) \\
 &= \widehat{R}(\widehat{p}_{\widehat{F}_d^{(k)}}) - R(p_{F_d^{(k)}}^*) + O_P\left(k\sqrt{\frac{\log n + \log d}{n^{\beta/(\beta+1)}}} + d\sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}}\right) \tag{25}
 \end{aligned}$$

$$\leq \widehat{R}(\widehat{p}_{F_d^{(k)}}) - R(p_{F_d^{(k)}}^*) + O_P\left(k\sqrt{\frac{\log n + \log d}{n^{\beta/(\beta+1)}}} + d\sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}}\right) \tag{26}$$

$$\begin{aligned}
 &= R(p_{F_d^{(k)}}^*) - R(p_{F_d^{(k)}}^*) + O_P\left(k\sqrt{\frac{\log n + \log d}{n^{\beta/(\beta+1)}}} + d\sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}}\right) \tag{27} \\
 &= O_P\left(k\sqrt{\frac{\log n + \log d}{n^{\beta/(\beta+1)}}} + d\sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}}\right).
 \end{aligned}$$

where (25) follows from Lemma 21, (26) follows from the fact that  $\widehat{p}_{\widehat{F}_d^{(k)}}$  is the minimizer of  $\widehat{R}(\cdot)$ , and (27) follows from Lemma 20.

### A.3 Proof of Theorem 10

To simplify notation, we denote

$$\begin{aligned}
 \phi_n(k) &= k\sqrt{\frac{\log n + \log d}{n^{\beta/(\beta+1)}}} \\
 \psi_n(d) &= d\sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}}.
 \end{aligned}$$

Following the same proof as Lemma 21, we obtain the following.

**Lemma 22** *Under the assumptions of Theorem 9, we have*

$$\sup_{F \in \mathcal{F}_d^{(k)}} |R(\widehat{p}_F) - \widehat{R}_{n_2}(\widehat{p}_F)| = O_P\left(\phi_n(k) + \psi_n(d)\right).$$

where  $\widehat{R}_{n_2}$  is the held out risk.

To prove Theorem 10, we now have

$$\begin{aligned}
 R(\widehat{p}_{\widehat{F}_d^{(\widehat{k})}}) - R(\widehat{p}_{\widehat{F}_d^{(k^*)}}) &= R(\widehat{p}_{\widehat{F}_d^{(\widehat{k})}}) - \widehat{R}_{n_2}(\widehat{p}_{\widehat{F}_d^{(\widehat{k})}}) + \widehat{R}_{n_2}(\widehat{p}_{\widehat{F}_d^{(\widehat{k})}}) - R(\widehat{p}_{\widehat{F}_d^{(k^*)}}) \\
 &= O_P(\phi_n(\widehat{k}) + \psi_n(d)) + \widehat{R}_{n_2}(\widehat{p}_{\widehat{F}_d^{(\widehat{k})}}) - R(\widehat{p}_{\widehat{F}_d^{(k^*)}}) \\
 &\leq O_P(\phi_n(\widehat{k}) + \psi_n(d)) + \widehat{R}_{n_2}(\widehat{p}_{\widehat{F}_d^{(k^*)}}) - R(\widehat{p}_{\widehat{F}_d^{(k^*)}}) \tag{28} \\
 &= O_P\left(\phi_n(\widehat{k}) + \phi_n(k^*) + \psi_n(d)\right).
 \end{aligned}$$

where (28) follows from the fact that  $\widehat{k}$  is the minimizer of  $\widehat{R}_{n_2}(\cdot)$ .

#### A.4 Proof of Theorem 11

Using the shorthand

$$\begin{aligned}\phi_n(k) &= k \sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}} \\ \psi_n(d) &= d \sqrt{\frac{\log n + \log d}{n^{2\beta/(1+2\beta)}}}\end{aligned}$$

We have that

$$\begin{aligned}R(\widehat{p}_{\widehat{F}_{n_2}}) - R(\widehat{p}_{F^*}) &= R(\widehat{p}_{\widehat{F}_{n_2}}) - \widehat{R}_{n_2}(\widehat{p}_{\widehat{F}_{n_2}}) + \widehat{R}_{n_2}(\widehat{p}_{\widehat{F}_{n_2}}) - R(\widehat{p}_{F^*}) \\ &= O_P(\phi_n(\widehat{k}) + \psi_n(d)) + \widehat{R}_{n_2}(\widehat{p}_{\widehat{F}_{n_2}}) - R(\widehat{p}_{F^*}) \\ &\leq O_P(\phi_n(\widehat{k}) + \psi_n(d)) + \widehat{R}_{n_2}(\widehat{p}_{F^*}) - R(\widehat{p}_{F^*}) \\ &= O_P(\phi_n(\widehat{k}) + \phi_n(k^*) + \psi_n(d))\end{aligned}\tag{29}$$

where (29) follows because  $\widehat{F}_{n_2}$  is the minimizer of  $\widehat{R}_{n_2}(\cdot)$ .

#### A.5 Proof of Theorem 12

We begin by showing an exponential probability inequality on the difference between the empirical and population mutual informations.

**Lemma 23** *Under Assumptions 1, 2, there exist generic constants  $c_5$  and  $c_6$  satisfying*

$$\mathbb{P}\left(|I(X_i; X_j) - \widehat{I}(X_i; X_j)| > \varepsilon\right) \leq c_5 \exp\left(-c_6 n^{\frac{\beta}{1+\beta}} (\log n)^{\frac{1}{1+\beta}} \varepsilon^2\right).$$

for arbitrary  $i, j \in \{1, \dots, d\}$  with  $i \neq j$ , and  $\varepsilon \rightarrow 0$  so that

$$\varepsilon = \Omega\left(\sqrt{\frac{\log r_n}{nh_2^2}}\right),$$

where  $r_n = \Omega(h_2^{-1})$ .

**Proof** For any  $\varepsilon = \Omega\left(\sqrt{\frac{\log r_n}{nh_2^2}}\right)$ , we have

$$\begin{aligned}&\mathbb{P}\left(|I(X_i; X_j) - \widehat{I}(X_i; X_j)| > \varepsilon\right) \\ &= \mathbb{P}\left(\left|\int_{\mathcal{X}_i \times \mathcal{X}_j} p^*(x_i, x_j) \log \frac{p^*(x_i, x_j)}{p^*(x_i)p^*(x_j)} dx_i dx_j - \int_{\mathcal{X}_i \times \mathcal{X}_j} \widehat{p}(x_i, x_j) \log \frac{\widehat{p}(x_i, x_j)}{\widehat{p}(x_i)\widehat{p}(x_j)} dx_i dx_j\right| > \varepsilon\right) \\ &\leq \mathbb{P}\left(\left|\int_{\mathcal{X}_i \times \mathcal{X}_j} (p^*(x_i, x_j) \log p^*(x_i, x_j) - \widehat{p}(x_i, x_j) \log \widehat{p}(x_i, x_j)) dx_i dx_j\right| > \frac{\varepsilon}{2}\right) \\ &\quad + \mathbb{P}\left(\left|\int_{\mathcal{X}_i \times \mathcal{X}_j} (p^*(x_i, x_j) \log p^*(x_i)p^*(x_j) - \widehat{p}(x_i, x_j) \log \widehat{p}(x_i)\widehat{p}(x_j)) dx_i dx_j\right| > \frac{\varepsilon}{2}\right)\end{aligned}\tag{30}$$

Since the second term of (30) only involves univariate kernel density estimates, this term is dominated by the first term, and we only need to analyze

$$\mathbb{P}\left(\left|\int_{\mathcal{X}_i \times \mathcal{X}_j} (p^*(x_i, x_j) \log p^*(x_i, x_j) - \widehat{p}(x_i, x_j) \log \widehat{p}(x_i, x_j)) dx_i dx_j\right| > \frac{\varepsilon}{2}\right).$$

The desired result then follows from the same analysis as in Lemma 20. ■

Let

$$L_n = \Omega\left(\sqrt{\frac{\log n + \log d}{n^{\beta/(1+\beta)}}}\right)$$

be defined as in Assumption 3. To prove the main theorem, we see the event  $\widehat{F}_d^{(k)} \neq F_d^{(k)}$  implies that there must be at least exist two pairs of edges  $(i, j)$  and  $(k, \ell)$ , such that

$$\text{sign}\left(I(X_i, X_j) - I(X_k, X_\ell)\right) \neq \text{sign}\left(\widehat{I}(X_i, X_j) - \widehat{I}(X_k, X_\ell)\right). \quad (31)$$

Therefore, we have

$$\begin{aligned} & \mathbb{P}\left(\widehat{F}_d^{(k)} \neq F_d^{(k)}\right) \\ & \leq \mathbb{P}\left(\left(I(X_i, X_j) - I(X_k, X_\ell)\right) \cdot \left(\widehat{I}(X_i, X_j) - \widehat{I}(X_k, X_\ell)\right) \leq 0, \text{ for some } (i, j), (k, \ell)\right). \end{aligned}$$

With  $d$  nodes, there can be no more than  $d^4/2$  pairs of edges; thus, applying a union bound yields

$$\begin{aligned} & \mathbb{P}\left(\left(I(X_i, X_j) - I(X_k, X_\ell)\right) \cdot \left(\widehat{I}(X_i, X_j) - \widehat{I}(X_k, X_\ell)\right) \leq 0, \text{ for some } (i, j), (k, \ell)\right) \\ & \leq \frac{d^4}{2} \max_{((i,j),(k,\ell)) \in \mathcal{J}} \mathbb{P}\left(\left(I(X_i, X_j) - I(X_k, X_\ell)\right) \cdot \left(\widehat{I}(X_i, X_j) - \widehat{I}(X_k, X_\ell)\right) \leq 0\right). \end{aligned}$$

Assumption 3 specifies that

$$\min_{((i,j),(k,\ell)) \in \mathcal{J}} |I(X_i, X_j) - I(X_k, X_\ell)| > 2L_n.$$

Therefore, in order for (31) hold, there must exist an edge  $(i, j) \in \mathcal{J}$  such that

$$|I(X_i, X_j) - \widehat{I}(X_i, X_j)| > L_n.$$

Thus, we have

$$\begin{aligned} & \max_{((i,j),(k,\ell)) \in \mathcal{J}} \mathbb{P}\left(\left(I(X_i, X_j) - I(X_k, X_\ell)\right) \cdot \left(\widehat{I}(X_i, X_j) - \widehat{I}(X_k, X_\ell)\right) \leq 0\right) \\ & \leq \max_{i,j \in \{1, \dots, d\}, i \neq j} \mathbb{P}\left(|I(X_i, X_j) - \widehat{I}(X_i, X_j)| > L_n\right) \\ & \leq c_5 \exp\left(-c_6 n^{\frac{\beta}{1+\beta}} (\log n)^{\frac{1}{1+\beta}} L_n^2\right). \end{aligned} \quad (32)$$

where (32) follows from Lemma 23.

Chaining together the above arguments, we obtain

$$\begin{aligned}
 & \mathbb{P}\left(\widehat{F}_d^{(k)} \neq F_d^{(k)}\right) \\
 & \leq \mathbb{P}\left(\left(I(X_i, X_j) - I(X_k, X_\ell)\right) \cdot \left(\widehat{I}(X_i, X_j) - \widehat{I}(X_k, X_\ell)\right) \leq 0, \text{ for some } (i, j), (k, \ell)\right) \\
 & \leq \frac{d^4}{2} \max_{((i,j),(k,\ell)) \in \mathcal{J}} \mathbb{P}\left(\left(I(X_i, X_j) - I(X_k, X_\ell)\right) \cdot \left(\widehat{I}(X_i, X_j) - \widehat{I}(X_k, X_\ell)\right) \leq 0\right) \\
 & \leq d^4 \max_{i,j \in \{1, \dots, d\}, i \neq j} \mathbb{P}\left(|I(X_i, X_j) - \widehat{I}(X_i, X_j)| > L_n\right) \\
 & \leq d^4 c_5 \exp\left(-c_6 n^{\frac{\beta}{1+\beta}} (\log n)^{\frac{1}{1+\beta}} L_n^2\right) \\
 & = o\left(c_5 \exp\left(4 \log d - c_6 (\log n)^{\frac{1}{1+\beta}} \log d\right)\right) \\
 & = o(1).
 \end{aligned}$$

The conclusion of the theorem now directly follows.

### A.6 Proof of NP-hardness of $t$ -Restricted Forest

We will reduce an instance of exact 3-cover (X3C) to an instance of finding a maximum weight  $t$ -restricted forest ( $t$ -RF).

Recall that in X3C, we are given a finite set  $X$  with  $|X| = 3q$  and a family of 3-element subsets of  $X$ ,  $\mathcal{S} = \{S \subset X : |S| = 3\}$ . The objective is to find a subfamily  $\mathcal{S}' \subset \mathcal{S}$  such that every element of  $X$  occurs in exactly one member of  $\mathcal{S}'$ , or to determine that no such subfamily exists.

Suppose then we are given  $X = \{x_1, \dots, x_n\}$  and  $\mathcal{S} = \{S \subset X : |S| = 3\}$ , with  $m = |\mathcal{S}|$ . We construct the graph  $G$  in an instance of  $t$ -RF as follows, and as illustrated in Figure 9.

For each  $x \in X$ , add an *element node* to  $G$ . For each  $S \in \mathcal{S}$ , construct a *gadget*, which is a subgraph comprised of a *nexus node*, three *junction nodes*, and three *lure nodes*; see Figure 9. We assign weights to the edges in a gadget in the following manner:

$$\begin{aligned}
 w(\text{element, junction}) &= 2 \\
 w(\text{nexus, lure}_1) &= 5 \\
 w(\text{lure}_1, \text{lure}_2) &= 10 \\
 w(\text{lure}_2, \text{lure}_3) &= 10 \\
 w(\text{nexus, junction}) &= N > 31m.
 \end{aligned}$$

Note that the weight  $N$  is chosen to be strictly greater than the weight all of the non-nexus-junction edges in the graph combined. To complete the instance of  $t$ -RF, let  $t = 7$ .

**Lemma 24** *Suppose  $G$  is a graph constructed in the transformation from X3C described above. Then  $F_t^*$  must contain all the nexus-junction edges.*

**Proof** The set of all nexus-junction edges together form a well-defined  $t$ -restricted forest, since each subtree has a nexus node and 3 junction nodes. Call this forest  $F$ . If some forest  $F'$  is missing a nexus-junction edge, then  $F'$  must have weight strictly less than  $F$ , since  $N$  is larger than the sum of all of the non-nexus-junction edges.  $\blacksquare$

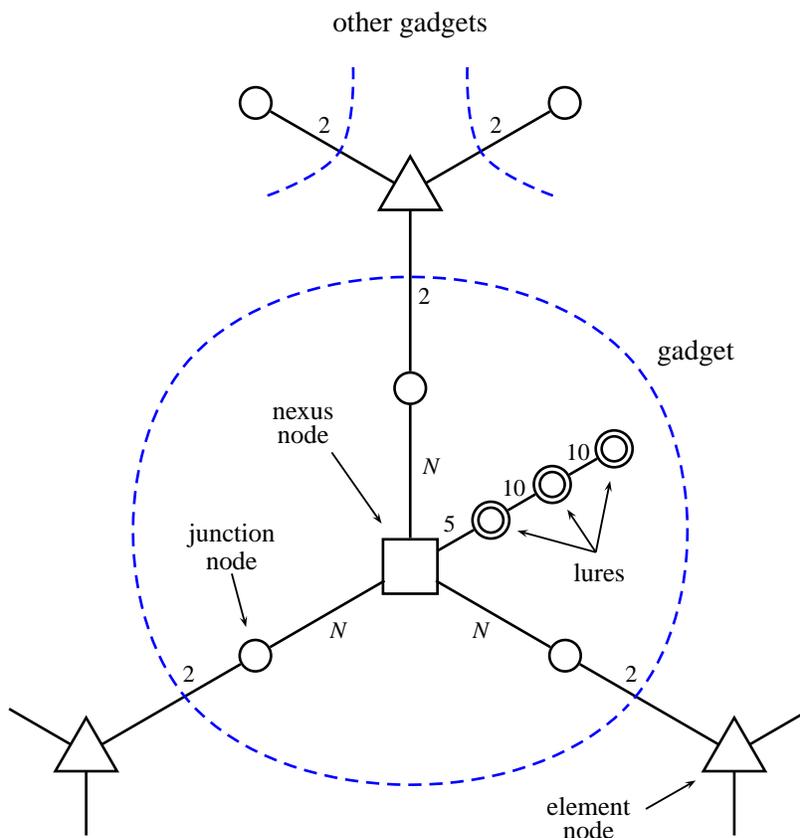


Figure 9: Gadget constructed in the reduction from X3C

**Lemma 25** *Each subtree in  $F_t^*$  can contain at most one nexus node.*

**Proof** Suppose a subtree  $T$  in  $F_t^*$  contains two nexus nodes. Then it must contain 6 junction nodes by Lemma 24. Thus,  $T$  contains at least 8 nodes, and therefore violates the  $t$ -restriction constraint. ■

**Lemma 26** *For each nexus node contained in  $F_t^*$ , the corresponding three junction nodes are either connected to all or none of the three neighboring element nodes.*

**Proof** By the previous two Lemmas 24 and 25, each subtree is associated with at most one gadget, and hence at most one  $S \in \mathcal{S}$ , and moreover each gadget has at least one associated subtree.

Without loss of generality, we consider a region of the graph local to some arbitrary subtree. By the size constraint, a subtree cannot contain all the adjacent element nodes and all the lure nodes.

We now perform a case analysis:

1. If a subtree contains no element nodes and all the lure nodes, then it has weight  $3N + 25$ . Call this an OFF configuration.

2. If a subtree contains two element nodes, and a second subtree of three nodes contains all the lure nodes, then the total weight of both subtrees is  $3N + 24$ . This is suboptimal because we can convert to an OFF configuration and gain additional weight without affecting any other subtrees. Hence, such a configuration cannot exist in  $F_t^*$ .
3. If a subtree contains two element nodes and  $\text{lure}_1$ , and a second subtree contains just  $\text{lure}_2$  and  $\text{lure}_3$ , then the total weight of the two subtrees is  $3N + 19$ . This is again suboptimal.
4. If a subtree contains an element node and both  $\text{lure}_1$  and  $\text{lure}_2$ , then there cannot be a second subtree in region local to the gadget. The weight of this one subtree is  $(3N + 2 + 5 + 10) = 3N + 17$ , which is suboptimal.
5. If a subtree contains all three element nodes and no lure nodes, and a second subtree contains all the lure nodes, then the total weight is  $(3N + 6) + 20 = 3N + 26$ . Call this an ON configuration.

Thus, we see that each gadget in  $F_t^*$  must be either an ON or an OFF configuration. ■

Recall that each gadget corresponds to a 3-element subset  $S$  in the family  $\mathcal{S}$ . Since a gadget in an ON configuration has greater weight than a gadget in an OFF configuration, an optimal  $t$ -RF will have as many gadgets in the ON configuration as possible. Thus, to solve X3C we can find the optimal  $t$ -RF and, to obtain a subcover  $S'$ , we place all  $S$  into  $S'$  that correspond to ON gadgets in the forest. By Lemma 25 each subtree can contain at most one nexus node, which implies that each ON gadget is connected to element nodes that are not connected to any other ON gadgets. Thus, this results in a subcover for which each element of  $X$  appears in at most one  $S \in S'$ .

### A.7 Proof of Theorem 16

Recall that we want to show that Algorithm 2 returns a forest with weight that is at least a quarter of the weight of the optimal  $t$ -restricted forest. Let us distinguish two types of constraints:

- (a) the degree of any node is at most  $t$ ;
- (b) the graph is acyclic.

Note that the optimal  $t$ -restricted forest  $F_t^*$  satisfies both the constraints above, and hence the maximum weight set of edges that satisfy both the constraints above has weight at least  $w(F_t^*)$ . Recall that the first stage of Algorithm 2 greedily adds edges subject to these two constraints—the next two lemmas show that the resulting forest has weight at least  $\frac{1}{2}w(F_t^*)$ .

**Lemma 27** *The family of subgraphs satisfying the constraints (a) and (b) form a 2-independence family. That is, for any subgraph  $T$  satisfying (a) and (b), and for any edge  $e \in G$ , there exist at most two edges  $\{e_1, e_2\}$  in  $T$  such that  $T \cup \{e\} - \{e_1, e_2\}$  also satisfies constraints (a) and (b).*

**Proof** Let  $T$  be a subgraph satisfying (a) and (b) and suppose we add  $e = (u, v)$  in  $T$ . Then the degrees of both  $u$  and  $v$  are at most  $t + 1$ . If no cycles were created, then we can simply remove an edge in  $T$  containing  $u$  (if any) and an edge in  $T$  containing  $v$  (if any) to satisfy the degree constraint (a) as well. If adding  $e$  created a cycle of the form  $\{\dots, (u', u), (u, v), (v, v')\}$ , then the edges  $(u', u)$  and  $(v, v')$  can be removed to satisfy both constraints (a) and (b). ■

**Lemma 28** *Let  $F_1$  be the forest output after Step 1 of algorithm 2. Then  $w(F_1) \geq \frac{1}{2}w(F_t^*)$ .*

**Proof** Let  $F^{**}$  be a maximum weight forest that obeys both constraints (a) and (b). Since the optimal  $t$ -restricted forest  $F_t^*$  obeys both these constraints, we have  $w(F_t^*) \leq w(F^{**})$ . By a theorem of Hausmann et al. (1980), in a  $p$ -independence family the greedy algorithm is a  $\frac{1}{p}$ -approximation to the maximum weight  $p$ -independent set. By Lemma 27, we know that the set of all subgraphs satisfying constraints (a) and (b) is a 2-independent family. Hence,  $w(F_1) \geq \frac{1}{2}w(F^{**}) \geq \frac{1}{2}w(F_t^*)$ . ■

We can now turn to the proof of Theorem 16.

**Proof** Given a graph  $G$ , let  $F_1$  be the forest output by first step of Algorithm 2, and let  $F_A$  be the forest outputted by the second step. We claim that  $w(F_A) \geq \frac{1}{2}w(F_1)$ ; combined with Lemma 28, this will complete the proof of the theorem.

To prove the claim, we first show that given any tree  $T$  with edge weights and maximum degree  $t \geq 2$ , we can obtain a sub-forest  $F$  with total weight  $w(F) \geq \frac{1}{2}w(T)$ , and where the number of edges in each tree in the forest  $F$  is at most  $t - 1$ . Indeed, root the tree  $T$  at an arbitrary node of degree-1, and call an edge  $e$  *odd* or *even* depending on the parity of the number of edges in the unique path between  $e$  and the root. Note that the set of odd edges and the set of even edges partition  $T$  into sub-forests composed entirely of stars of maximum degree  $t - 1$ , and one of these sub-forests contains half the weight of  $T$ , which is what we wanted to show.

Applying this procedure to each tree  $T$  in the forest  $F_1$ , we get the existence of a  $t - 1$ -restricted subforest  $F'_1 \subseteq F_1$  that has weight at least  $\frac{1}{2}w(F_1)$ . Observe that a  $t - 1$ -restricted subforest is *a fortiori* a  $k$ -restricted subforest, and since  $w(F_A)$  is the best  $t$ -restricted subforest of  $F_1$ , we have

$$w(F_A) \geq w(F'_1) \geq \frac{1}{2}w(F_1) \geq \frac{1}{4}w(F_t^*),$$

completing the proof. ■

### A.7.1 AN IMPROVED APPROXIMATION ALGORITHM

We can get an improved approximation algorithm based on a linear programming approach. Recall that  $F^{**}$  is a maximum weight forest satisfying both (a) and (b). A result of Singh and Lau (2007) implies that given any graph  $G$  with non-negative edge weights, one can find in polynomial time a forest  $F_{SL}$  such that

$$w(F_{SL}) \geq w(F^{**}) \geq w(F_t^*), \tag{33}$$

but where the maximum degree in  $F_{SL}$  is  $t + 1$ . Now applying the procedure from the proof of Theorem 16, we get a  $t$ -restricted forest  $F'_{SL}$  whose weight is at least half of  $w(F_{SL})$ . Combining this with (33) implies that  $w(F'_{SL}) \geq w(F_t^*)$ , and completes the proof of the claimed improved approximation algorithm. We remark that the procedure of Singh and Lau (2007) to find the forest  $F_{SL}$  is somewhat computationally intensive, since it requires solving vertex solutions to large linear programs.

### A.8 Proof of Theorem 17

Proceeding as in the proof of Theorem 10, we have that

$$\begin{aligned} \left| R(\widehat{p}_{\widehat{F}_t}) - R(\widehat{p}_{F_t^*}) \right| &\leq R(\widehat{p}_{\widehat{F}_t}) - \widehat{R}_{n_1}(\widehat{p}_{\widehat{F}_t}) + \left| \widehat{R}_{n_1}(\widehat{p}_{\widehat{F}_t}) - R(\widehat{p}_{F_t^*}) \right| \\ &= O_P(k\phi_n(d) + d\psi_n(d)) + \left| \widehat{R}_{n_1}(\widehat{p}_{\widehat{F}_t}) - R(\widehat{p}_{F_t^*}) \right|. \end{aligned}$$

Now, let  $\widehat{H}_{n_1}$  denote the estimated entropy  $H(X) = \sum_k H(X_k)$ , constructed using the kernel density estimates  $\widehat{p}_{n_1}(x_k)$ . Since the risk is the negative expected log-likelihood, we have using the approximation guarantee that

$$\begin{aligned} \widehat{R}_{n_1}(\widehat{p}_{\widehat{F}_t}) - R(\widehat{p}_{F_t^*}) &= -\widehat{w}_{n_1}(\widehat{F}_t) + \widehat{H}_{n_1} - R(\widehat{p}_{F_t^*}) \\ &\leq -\frac{1}{c}\widehat{w}_{n_1}(F_t^*) + \widehat{H}_{n_1} - R(\widehat{p}_{F_t^*}) \\ &= \widehat{R}_{n_1}(\widehat{p}_{F_t^*}) + \frac{c-1}{c}\widehat{w}_{n_1}(F_t^*) - R(\widehat{p}_{F_t^*}) \\ &= O_P\left(k^*\phi_n(d) + d\psi_n(d) + \frac{c-1}{c}w(F_t^*)\right) \end{aligned}$$

and the result follows.

### A.9 The TreePartition Subroutine

To produce the best  $t$ -restricted subforest of the forest  $F_1$ , we use a divide-and-conquer forest partitioning algorithm described by Lukes (1974), which we now describe in more detail.

To begin, note that finding an optimal subforest is equivalent to finding a partition of the nodes in the forest, where each disjoint tree in the subforest is a cluster in the partition. Since a forest contains a disjoint set of trees, it suffices to find the optimal  $t$ -restricted partition of each of the trees.

For every subtree  $T$ , with root  $v$ , we will find a *list of partitions*  $v.P = \{v.P_0, v.P_1, \dots, v.P_k\}$  such that

1. for  $i \neq 0$ ,  $v.P_i$  is a partition whose cluster containing root  $v$  has size  $i$ ;
2.  $v.P_i$  has the maximum weight among all partitions satisfying the above condition.

We define  $v.P_0$  to be  $\arg \max\{w(v.P_1), \dots, w(v.P_k)\}$ . The Merge subroutine used in TreePartition takes two lists of partitions  $\{v.P, u_i.P\}$ , where  $v$  is the parent of  $u_i$ ,  $v.P$  is a partition of node  $v$  unioned with subtrees of children  $\{u_1, \dots, u_{i-1}\}$ , and  $u_i.P$  is a partition of the subtree of child  $u_i$ ; refer to Figure 10.

Since a partition is a list of clusters of nodes, we denote by  $\text{Concat}(v.P_2, u.P_{k-2})$  the concatenation of clusters of partitions  $v.P_2, u.P_{k-2}$ . Note that the concatenation forms a partition if  $v.P_2$  and  $u.P_{k-2}$  are respectively partitions of two disjoint sets of vertices. The weight of a partition is denoted  $w(v.P_2)$ , that is, the weight of all edges between nodes of the same cluster in the partition  $v.P_2$ .

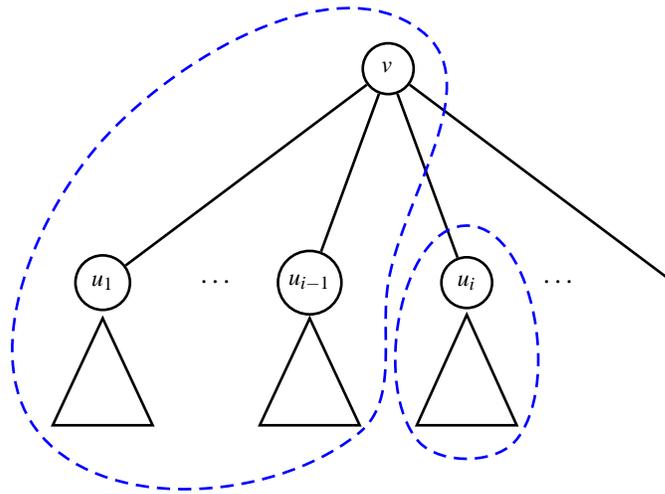


Figure 10: The TreePartition procedure to merge two subproblems.

---

**Algorithm 4** TreePartition( $T, t$ )

---

- 1: **Input** a tree  $T$ , a positive integer  $t$
  - 2: **Returns** an optimal partition into trees of size  $\leq t$ .
  - 3: Initialize  $v.P_1 = [\{v\}]$  where  $v$  is root of  $T$ , if  $v$  has no children, return  $v.P_1$
  - 4: For all children  $\{u_1, \dots, u_s\}$  of  $v$ , recursively call TreePartition( $u_i, t$ ) to get a collection of lists of partitions  $\{u_1.P, u_2.P, \dots, u_s.P\}$
  - 5: For each child  $u_i \in \{u_1, \dots, u_s\}$  of  $v$   
Update  $v.P \leftarrow \text{Merge}(u_i.P, v.P)$
  - 6: **Output**  $v.P_0$
- 

---

**Algorithm 5** Merge( $v.P, u.P$ )

---

- 1: **Input** a list of partitions  $v.P$  and  $u.P$ , where  $v$  is a parent of  $u$ .
  - 2: **Returns** a single list of partitions  $v.P'$ .
  - 3: For  $i = 1, \dots, t$ :
    1. Let  $(s^*, t^*) = \arg \max_{(s,t):s+t=i} w(\text{Concat}(v.P_s, u.P_t))$
    2. Let  $v.P'_i = \text{Concat}(v.P_{s^*}, u.P_{t^*})$
  - 4: Select  $v.P'_0 = \arg \max_{v.P'_i} w(v.P'_i)$
  - 5: **Output**  $\{v.P'_0, v.P'_1, \dots, v.P'_n\}$
-

## Appendix B. Computation of the Mutual Information Matrix

In this appendix we explain different methods for computing the mutual information matrix, and making the tree estimation more efficient. One way to evaluate the empirical mutual information is to use

$$\widehat{I}(X_i; X_j) = \frac{1}{n_1} \sum_{s \in \mathcal{D}_1} \log \frac{\widehat{p}_{n_1}(X_i^{(s)}, X_j^{(s)})}{\widehat{p}_{n_1}(X_i^{(s)}) \widehat{p}_{n_1}(X_j^{(s)})}. \quad (34)$$

Compared with our proposed method

$$\widehat{I}_{n_1}(X_i, X_j) = \frac{1}{m^2} \sum_{k=1}^m \sum_{\ell=1}^m \widehat{p}_{n_1}(x_{ki}, x_{\ell j}) \log \frac{\widehat{p}_{n_1}(x_{ki}, x_{\ell j})}{\widehat{p}_{n_1}(x_{ki}) \widehat{p}_{n_1}(x_{\ell j})}, \quad (35)$$

(34) is somewhat easier to calculate. However, if the sample size in  $\mathcal{D}_1$  is small, the approximation error can be large. A different analysis is needed to provide justification of the method based on (34), which would be more difficult since  $\widehat{p}_{n_1}(\cdot)$  is dependent on  $\mathcal{D}_1$ . For these reasons we use the method in (35).

Also, note that instead of using the grid based method to evaluate the numerical integral, one could use sampling. If we can obtain  $m_1$  i.i.d. samples from the bivariate density  $\widehat{p}(X_i, X_j)$ ,

$$\left\{ (X_i^{(s)}, X_j^{(s)}) \right\}_{s=1}^{m_1} \stackrel{\text{i.i.d.}}{\sim} \widehat{p}_{n_1}(x_i, x_j),$$

then the empirical mutual information can be evaluated as

$$\widehat{I}(X_i; X_j) = \frac{1}{m_1} \sum_{s=1}^{m_1} \log \frac{\widehat{p}(X_i^{(s)}, X_j^{(s)})}{\widehat{p}(X_i^{(s)}) \widehat{p}(X_j^{(s)})}.$$

Compared with (34), the main advantage of this approach is that the estimate can be arbitrarily close to (8) for large enough  $m_1$  and  $m$ . Also, the computation can be easier compared to Algorithm 1. Let  $\widehat{p}_{n_1}(X_i, X_j)$  be the bivariate kernel density estimator on  $\mathcal{D}_1$ . To sample a point from  $\widehat{p}_{n_1}(X_i, X_j)$ , we first random draw a sample  $(X_i^{(k')}, X_j^{(\ell')})$  from  $\mathcal{D}_1$ , and then sample a point  $(X, Y)$  from the bivariate distribution

$$(X, Y) \sim \frac{1}{h_2^2} K \left( \frac{X_i^{(k')} - \cdot}{h_2} \right) K \left( \frac{X_j^{(\ell')} - \cdot}{h_2} \right).$$

Though this sampling strategy is superior to Algorithm 1, it requires evaluation of the bivariate kernel density estimates on many random points, which is time consuming; the grid-based method is preferred.

In our two-stage procedure, the stage requires calculation of the empirical mutual information  $\widehat{I}(X_i; X_j)$  for  $\binom{d}{2}$  entries. Each requires  $O(m^2 n_1)$  work to evaluate the bivariate and univariate kernel density estimates on the  $m \times m$  grid, in a naive implementation. Therefore, the total time to calculate the empirical mutual information matrix  $M$  is  $O(m^2 n_1 d^2)$ . In the second stage, the time complexity of the Chow-Liu algorithm is dominated by the first step. Therefore the total time complexity is  $O(m^2 n_1 d^2)$ . The first stage requires  $O(d^2)$  space to store the matrix  $M$  and  $O(m^2 n_1)$  space to

---

**Algorithm 6** More efficient calculation of the mutual information matrix  $M$ .
 

---

```

1: Initialize  $M = \mathbf{0}_{d \times d}$  and  $H^{(i)} = \mathbf{0}_{n_1 \times m}$  for  $i = 1, \dots, d$ .
2: % calculate and pre-store the univariate KDE
3: for  $k = 1, \dots, d$  do
4:   for  $k' = 1, \dots, m$  do
5:      $\hat{p}(x_k^{(k')}) \leftarrow \frac{1}{n_1} \sum_{s \in \mathcal{D}_1} \frac{1}{h_1} K \left( \frac{X_k^{(s)} - x_k^{(k')}}{h_1} \right)$ 
6:   for  $k' = 1, \dots, m$  do
7:     % calculate the components used for the bivariate KDE
8:     for  $i' = 1, \dots, n_1$  do
9:       for  $i = 1, \dots, d$  do
10:         $H^{(i)}(i', k') \leftarrow \frac{1}{h_2} K \left( \frac{X_{i'}^{(i)} - x_{i'}^{(k')}}{h_2} \right)$ 
11:      % calculate the mutual information matrix
12:      for  $\ell' = 1, \dots, m$  do
13:        for  $i = 1, \dots, d - 1$  do
14:          for  $j = i + 1, \dots, d$  do
15:             $\hat{p}(x_i^{(k')}, x_j^{(\ell')}) \leftarrow 0$ 
16:          for  $i' = 1, \dots, n_1$  do
17:             $\hat{p}(x_i^{(k')}, x_j^{(\ell')}) \leftarrow \hat{p}(x_i^{(k')}, x_j^{(\ell')}) + H^{(i)}(i', k') \cdot H^{(j)}(i', \ell')$ 
18:             $\hat{p}(x_i^{(k')}, x_j^{(\ell')}) \leftarrow \hat{p}(x_i^{(k')}, x_j^{(\ell')}) / n_1$ 
19:             $M(i, j) \leftarrow M(i, j) + \frac{1}{m^2} \hat{p}(x_i^{(k')}, x_j^{(\ell')}) \cdot \log \left( \hat{p}(x_i^{(k')}, x_j^{(\ell')}) / (\hat{p}(x_i^{(k')}) \cdot \hat{p}(x_j^{(\ell')})) \right)$ 

```

---

evaluate the kernel density estimates on  $\mathcal{D}_1$ . The space complexity for the Chow-Liu algorithm is  $O(d^2)$ , and thus the total space complexity is  $O(d^2 + m^2 n_1)$ .

The quadratic time and space complexity in the number of variables  $d$  is acceptable for many practical applications but can be prohibitive when the dimension  $d$  is large. The main bottleneck is to calculate the empirical mutual information matrix  $M$ . Due to the use of the kernel density estimate, the time complexity is  $O(d^2 m^2 n_1)$ . The straightforward implementation in Algorithm 1 is conceptually easy but computationally inefficient, due to many redundant operations. For example, in the nested for loop, many components of the bivariate and univariate kernel density estimates are repeatedly evaluated. In Algorithm 6, we suggest an alternative method which can significantly reduce such redundancy at the price of increased but still affordable space complexity.

The main technique used in Algorithm 6 is to change the order of the multiple nested for loops, combined with some pre-calculation. This algorithm can significantly boost the empirical performance, although the worst case time complexity remains the same. An alternative suggested by Bach and Jordan (2003) is to approximate the mutual information, although this would require further analysis and justification.

## References

- Martin Aigner and Günter Ziegler. *Proofs from THE BOOK*. Springer-Verlag, 1998.
- Francis R. Bach and Michael I. Jordan. Beyond independent components: Trees and clusters. *Journal of Machine Learning Research*, 4:1205–1233, 2003.
- Onureena Banerjee, Laurent El Ghaoui, and Alexandre d’Aspremont. Model selection through sparse maximum likelihood estimation. *Journal of Machine Learning Research*, 9:485–516, March 2008.
- Arthur Cayley. A theorem on trees. *Quart. J. Math.*, 23:376–378, 1889.
- Anton Chechetka and Carlos Guestrin. Efficient principled learning of thin junction trees. In *In Advances in Neural Information Processing Systems (NIPS)*, Vancouver, Canada, December 2007.
- C. K. Chow. and C. N. Liu. Approximating discrete probability distributions with dependence trees. *IEEE Transactions on Information Theory*, 14(3):462–467, 1968.
- Jianqing Fan and Irène Gijbels. *Local Polynomial Modelling and Its Applications*. Chapman and Hall, 1996.
- Jerome H. Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2007.
- Michael Garey and David Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman, 1979. ISBN 0-7167-1044-7.
- Evarist Giné and Armell Guillou. Rates of strong uniform consistency for multivariate kernel density estimators. *Annales de l’institut Henri Poincaré (B), Probabilités et Statistiques*, 38:907–921, 2002.
- Dirk Hausmann, Bernhard Korte, and Tom Jenkyns. Worst case analysis of greedy type algorithms for independence systems. *Math. Programming Studies*, 12:120–131, 1980.
- Joseph B. Kruskal. On the shortest spanning subtree of a graph and the traveling salesman problem. *Proceedings of the American Mathematical Society*, 7(1):48–50, 1956.
- Steffen L. Lauritzen. *Graphical Models*. Clarendon Press, 1996.
- Han Liu, John Lafferty, and Larry Wasserman. The nonparanormal: Semiparametric estimation of high dimensional undirected graphs. *Journal of Machine Learning Research*, 10:2295–2328, October 2009.
- Joseph A. Lukes. Efficient algorithm for the partitioning of trees. *IBM Jour. of Res. and Dev.*, 18(3):274, 1974.
- Nicolai Meinshausen and Peter Bühlmann. High dimensional graphs and variable selection with the Lasso. *The Annals of Statistics*, 34:1436–1462, 2006.

- Renuka Nayak, Michael Kearns, Richard Spielman, and Vivian Cheung. Coexpression network based on natural variation in human gene expression reveals gene interactions and functions. *Genome Research*, 19:1953–1962, 2009.
- Deborah Nolan and David Pollard. U-processes: Rates of convergence. *The Annals of Statistics*, 15(2):780 – 799, 1987.
- Philippe Rigollet and Régis Vert. Optimal rates for plug-in estimators of density level sets. *Bernoulli*, 15(4):1154–1178, 2009.
- Alessandro Rinaldo and Larry Wasserman. Generalized density clustering. *Ann. Statist.*, 38(5): 2678–2722, 2010.
- Mohit Singh and Lap Chi Lau. Approximating minimum bounded degree spanning trees to within one of optimal. In *STOC'07—Proceedings of the 39th Annual ACM Symposium on Theory of Computing*, pages 661–670. ACM, New York, 2007.
- Vincent Tan, Animashree Anandkumar, and A. Willsky. Learning Gaussian tree models: Analysis of error exponents and extremal structures. *IEEE Trans. on Signal Processing*, 58(5):2701–2714, 2010.
- Vincent Tan, Animashree Anandkumar, Lang Tong, and Alan Willsky. A large-deviation analysis for the maximum likelihood learning of tree structures. *IEEE Trans. on Info. Theory*, 57(3): 1714–1735, 2011.
- Alexandre B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Publishing Company, Incorporated, 2008.
- Anja Wille, Philip Zimmermann, Eva Vranová, Andreas Fürholz, Oliver Laule, Stefan Bleuler, Lars Hennig, Amela Prelić, Peter von Rohr, Lothar Thiele, Eckart Zitzler, Wilhelm Grussem, and Peter Bühlmann. Sparse Gaussian graphical modelling of the isoprenoid gene network in *Arabidopsis thaliana*. *Genome Biology*, 5:R92, 2004.