

CHAPTER 1

Method

The following chapter is structured as follows. Initially, we shall introduce the basic concept of causality and *structural causal models* (SCMs) based on [?]. From that, we shall discuss the method proposed by [?] to infer such SCMs. In particular, we shall state the underlying assumptions of the method and discuss the implication of these. Furthermore, we shall see how based on observed similarities between pairs of random variables, the proposed method deconvolves a matrix of similarities and result in a set of proposed structures for the underlying SCM. Namely, how we from data can make predictions of which variables influence each other directly or through mediators - also known as in-between variable.

From the basic assumptions, we shall then discuss correlation and mutual information as similarity measures and how they differ. In particular, we shall see that mutual information and copulas, a statistical tool for isolating the joint behavior of random variables, are very related topics. In particular, mutual information can be computed from only the copula. Furthermore, we shall extend on the original methodology by considering mixed random variables as well. This is important, such that the delays T_u^D from ?? can be included later on in ??, where we shall use the methods obtained here to infer possible causal structures.

In Section 1.3, we will discuss the algorithms to be used later in detail. In

particular, we shall clarify a few results from the original paper [?] and extend on their results on robustness of the deconvolution algorithm. A major assumption regarding the *size* of the observed matrix of associations (G_{obs}) is removed, obtaining a more general and useful result. The Frobenius and maximum matrix norms are especially considered.

Finally, in Section 1.4, we shall discuss different methods for estimating mutual information as well as their drawbacks. Especially kernel based density estimators are considered discussed in terms of a simple dataset from [?] with observations related to suicide risk.

General introduction to method and causality

1.1 Causality and Causal Discovery

In this section, we shall discuss the method for network deconvolution, originally proposed by [?]. The underlying problem is inferring direct effects and dependencies. From this, using prior information on the production setup, we shall be able to infer causal dependencies by directing the resulting edges from the network deconvolution (ND) algorithm. Particularly, the framework and general algorithm proposed by Feizi et al. stems from a graph-theoretic approach to the problem of inferring direct dependencies. Namely, suppose that observations depend on quantities such as levels and sojourn times of in this case a chemical process. We shall represent these properties as vertices (nodes) V and dependencies between properties as edges. Initially, when observing the vertices, we observe both direct and indirect effects. Particularly, a vertex v_1 might influence some other vertex v_3 through another vertex v_2 if v_2 depends on v_1 and v_3 of v_2 . In this case, we will observe that v_1 influences v_3 , but actually it is v_2 that has a direct influence on v_3 . In graph-theoretical terms, we thus observe the transitive closure of the information that flows between vertices but want to infer the underlying network structure.

An important note on the algorithm to come is that we only use vertices that we have observed. Namely, the underlying structure might be as in Figure 1.1(a) with an unobserved node/variable (named U in this case). However, without any more assumptions or modelling choices we would (ideally) infer the network structure depicted in Figure 1.1(b). With these initial comments, we proceed with the general setup and assumptions for network deconvolution based on observations.

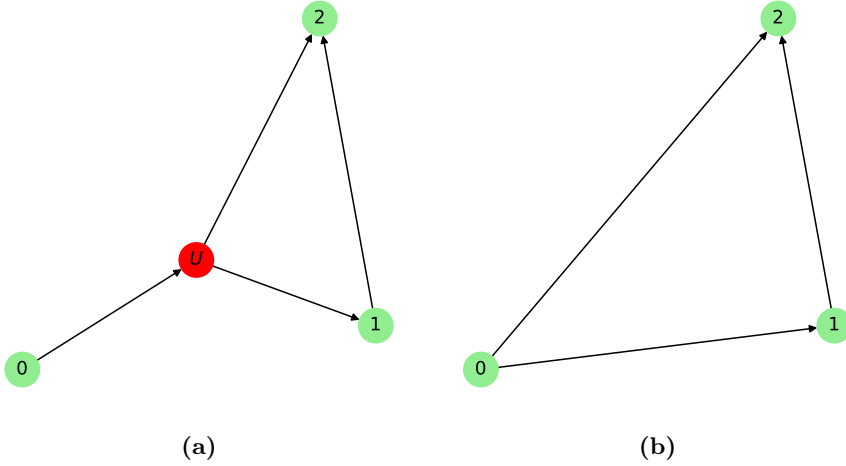


Figure 1.1: (a) An example of a causal structure depicted as a graph. When observing the network, only nodes 0, 1 and 2 are observed/recorded. (b) The resulting inferred graph from observational data. Although this is not a complete picture of the true underlying dynamics of the system, if only the observed variables are of interest, this will be an equally proper representation of the system. Furthermore, in practice this means no further assumptions are made which can and can not be of desire. Namely, if prior information is accessible one might introduce new nodes in the inferred network.

1.1.1 Setup and Assumptions

Suppose a set of d random variables (X_i) is given i.e. a d -dimensional random vector \mathbf{X} . The method presented in this section aims to discover direct relationships between pairs X_i and X_j for $i \neq j$. These relationships will be presented by a directed graph as in the previous section or an undirected graph in case the causal direction is either unknown or such an assumption on direction is not plausible. In particular, we shall let each random variable X_i be represented by a vertex in a graph. We will later discuss a way of directing edges such that a causal network may be discovered i.e. a directed acyclic graph that may be used for inference.

The method proposed by [?] then works as follows. Given an observed matrix $G_{obs} \in \mathbb{R}^{d \times d}$ of similarities between each pair of variables, we shall deduce a matrix $G_{dir} \in \mathbb{R}^{d \times d}$ of direct similarities between each pair of random variables X_i and X_j . In particular, we wish to filter out indirect effects which we will denote by G_{dir} defined as effects between pairs of variables that is the result

of effects propagating through other variables. The measure of similarity, can in practice be any desired measure such as correlation or mutual information which we will focus on in this thesis. See Section 1.2 for a further discussion on these two measures and Section 1.3 for how to obtain such a matrix. Note that the algorithm presented will in theory work for non-symmetric measures as well such as *Interaction information*, *Directed information* and *Normalized information*.



Figure 1.2: On overview of the methodology that we shall develop in this thesis. In particular, the estimation of mutual information - a possible measure of similarity - and deconvolution will be discussed in this section.

The (direct) network is then presented by the discovered G_{dir} containing only the direct effects i.e. interaction between pairs of variables which can be viewed as weights on the edges of the complete graph with nodes representing the random variables. As we shall see in Subsection 1.3.3, the algorithm is somewhat robust to noise in the sense that we can ensure accuracy depending on the level of noise observed present in G_{obs} and on the norm chosen (from a certain, although rather general, set of norms). Namely, if G_{obs} is subject to noise, we find a bound on how different the inferred directed effects can be to the true direct effects using different matrix norms to measure this difference. This hints to that a threshold on the inferred weights on the edges of the network might be a good idea to remove small inferred effects. This is further supported by the facts that often only the most influential variables are of importance when trying to control the process.

The first assumption is that the observed matrix of co-dependence G_{obs} may be expressed as

$$G_{obs} = G_{dir} + G_{indir} \quad (1.1)$$

Namely, that the direct and indirect effects can be added together to get the total and thus observed interdependence between each pair of variables. Often, this is not the case as we shall see later on. However, the error made from this assumption and the ones to be presented seem to be small enough that the discovered network accurately resemble the true underlying network.

The second and final assumption is that the indirect effects G_{indir} can be computed in terms of G_{dir} . Namely, that

$$G_{indir} = G_{dir}^2 + G_{dir}^3 + \dots = \sum_{k=2}^{\infty} G_{dir}^k \quad (1.2)$$

i.e. that the observed *information* exchanged on an edge e_{ij} between nodes X_i and X_j is the sum of the second, third etc. order effects, each given by the information on the n -path (where n is the order of the (diminishing) indirect effect) again assumed to be a sum of products. In other terms, the second order indirect effect between X_i and X_j (given as the (i, j) element of G_{dir}^2) is the sum of products on edges e_{ik} and e_{kj} for all k

$$[G_{dir}^2]_{ij} = \sum_{k=1}^d e_{ik} e_{kj}$$

where e_{ij} is the (i, j) element of G_{dir} . This is of course not true in general. However, through error analysis in Subsection 1.3.3 and controlled examples in ?? we shall see that this assumption is either true under some additional assumptions or only results in small numerical errors. Immediately, we observe that e_{ii} is of interest in terms of its physical meaning. The co-dependence between a random variable and itself might be somewhat ambiguous or even undefined depending on the measure. Thus, the notion of (non-existing) edges e_{ii} will be of interest later on when using the method on controlled cases. We note that in G_{obs} we shall in general set these elements to 0.

Thus, from the above assumptions, it follows that we can express G_{obs} as

$$G_{obs} = G_{dir} + G_{dir}^2 + G_{dir}^3 + \dots = G_{dir} + G_{dir} G_{obs} \quad (1.3)$$

Clearly, such a G_{dir} must have spectral radius at most 1 as otherwise, the above sum diverges and thus G_{obs} will not exist. I.e. $\rho(G_{dir}) < 1$, where $\rho(\cdot)$ denotes the spectral radius. Thus, assuming convergence we can rewrite the infinite series as

$$G_{obs} = G_{dir} (I - G_{dir})^{-1} \quad (1.4)$$

Multiplying the above by $(I - G_{dir})$ from the right and moving around terms, it immediately follows that

$$G_{obs} = G_{dir} (I + G_{obs}) \quad (1.5)$$

Thus, if we can show that $-1 \notin \sigma(G_{obs})$ (where $\sigma(\cdot)$ denotes the spectrum of an operator), we can isolate G_{dir} . Namely, we need -1 to not be an eigenvalue of G_{obs} . This is indeed true under the assumption that $\rho(G_{dir}) < 1$. In particular, assume that (λ, v) is an eigenpair of G_{dir} . Then, by assumption $|\lambda| < 1$ and by Equation 1.3:

$$G_{obs} v = \sum_{k=1}^{\infty} G_{dir}^k v = \sum_{k=1}^{\infty} \lambda^k v = \frac{\lambda}{1 - \lambda} v$$

where we used that v is an eigenvector of G_{dir} and the geometric series converges as $|\lambda| < 1$. In particular, $\left(\frac{\lambda}{1 - \lambda}, v\right)$ is an eigenpair of G_{obs} . In ??, we show

that λ is an eigenvalue of G_{dir} if and only if $\frac{\lambda}{1-\lambda}$ is an eigenvalue of G_{obs} i.e. there is a bijection between the eigenvalues of G_{dir} and G_{obs} . Thus, as the spectral radius of G_{dir} is less than one such that $\lambda \in (-1, 1)$, we conclude that $\sigma(G_{obs}) \subset (-\frac{1}{2}, \infty)$. Hence, $-1 \notin \sigma(G_{obs})$ and G_{dir} can easily be isolated in Equation 1.5 as

$$G_{dir} = G_{obs} (I + G_{obs})^{-1} \quad (1.6)$$

We note that from the above, we have that G_{obs} is a result of a G_{dir} only if the smallest eigenvalue of G_{obs} is larger than $-1/2$.

Furthermore, if the measure of dependence between pairs of variables is symmetric, then so is G_{obs} and hence diagonalizable by some orthogonal matrix U (such that $U^T = U^{-1}$) and diagonal matrix Λ_{obs} such that $G_{obs} = U\Lambda_{obs}U^T$ (with the columns of U being right eigenvectors of G_{obs}). This follows from the fact that any real symmetric matrix is diagonalizable. It follows that G_{dir} can be expressed in the following simple way (which is useful for computational efficiency)

$$\begin{aligned} G_{dir} &= U\Lambda_{obs}U^T (I + U\Lambda_{obs}U^T)^{-1} \\ &= U\Lambda_{obs}U^T (UU^T + U\Lambda_{obs}U^T)^{-1} \\ &= U\Lambda_{obs}U^T (U(I + \Lambda_{obs})U^T)^{-1} \\ &= U\Lambda_{obs}U^T U(I + \Lambda_{obs})^{-1}U^T \\ &= U\Lambda_{obs}(I + \Lambda_{obs})^{-1}U^T \\ &= U\Lambda_{dir}U^T \end{aligned}$$

where $\Lambda_{dir} = \Lambda_{obs}(I + \Lambda_{obs})^{-1}$ is also a diagonal matrix, with elements corresponding to the inverse of the mapping $\lambda \mapsto \frac{\lambda}{1-\lambda}$.

As we shall later use some assumptions regarding causality leading G_{obs} to be a triangular matrix, we shall investigate the properties of the resulting G_{dir} . Namely, in the following, we show that given the existence of G_{dir} (with necessary and sufficient conditions on G_{obs} as given above), G_{obs} is triangular if and only if G_{dir} is triangular. Thus, by directing the observed similarity (by removing half the edge weights in G_{obs}), we also infer a directed graph G_{dir} .

Clearly, if G_{dir} is triangular, so are the powers G_{dir}^i for all $i \in \mathbb{N}$ and hence if the infinite sum $\sum_{i=1}^{\infty} G_{dir}^i$ converges, G_{obs} is triangular as well.

To show the other way, assume that G_{obs} is triangular and is the result of a G_{dir} with spectral radius smaller than 1. By Equation 1.6, G_{dir} is triangular if the inverse of $I + G_{obs}$ is triangular (upper triangular if G_{obs} is also upper triangular and similarly for lower triangular). This is indeed the case as in

general, the inverse of a triangular matrix is also triangular provided that the diagonal elements are non-zero. Note that $I + G_{obs}$ is never 0 in the diagonal, as $-1/2$ is the smallest possible eigenvalue of G_{obs} and hence smallest diagonal element. A simple proof is as follows. Assume without loss of generality, that a matrix T is upper triangular. Let D be the diagonal elements of T and T_u be the remaining strictly upper triangular part of T such that $T = D + T_u$. Then, assuming that D has non-zero diagonal elements, $T = D(I + D^{-1}T_u)$. Therefore,

$$\begin{aligned} T^{-1} &= (I + D^{-1}T_u)^{-1} D^{-1} \\ &= \sum_{i=0}^{\infty} (-D^{-1}T_u)^i D^{-1} \\ &= \sum_{i=0}^{d-1} (-D^{-1}T_u)^i D^{-1} \end{aligned}$$

which is clearly also upper triangular. The second and final equality follows from T_u being strictly upper triangular and thus nilpotent such that $\sigma(-D^{-1}T_u) = \{0\}$ and $T_u^d = 0$. We conclude that G_{obs} is triangular if and only if G_{dir} is (under the assumption G_{dir} exists and G_{obs}).

Finally, before discussing the implementation and analyzing the algorithm both analytically and through examples, we will take a closer look at the similarity measures that are to be used with this method and that in the end will make up the matrix G_{obs} . Namely, *mutual information* and *correlation*.

1.2 Information Measures and Computation

In this section we discuss two measures that can be used to construct the matrices of codependency from the previous section. Namely, we shall touch on correlation and discuss what one might choose to call Copula-based entropy. However, before discussing Copula entropy (CE) we first need to define what a copula is.

1.2.1 Copula

Given a set of d random variables X_1, \dots, X_d , a copula is loosely speaking a distribution function with domain $[0, 1]^d$ incorporating the dependence structure between the random variables. Given a joint distribution function F for

(X_1, \dots, X_d) and (invertible) marginals F_1, \dots, F_d we define a copula C as

$$\begin{aligned} F(x_1, \dots, x_d) &= \mathbb{P}(X_1 \leq x_1, \dots, X_d \leq x_d) \\ &= \mathbb{P}(F_1(X_1) \leq F_1(x_1), \dots, F_d(X_d) \leq F_d(x_d)) \\ &= C(F_1(x_1), \dots, F_d(x_d)) \end{aligned}$$

Letting $u_i = F_i(x_i) \in [0, 1]$ it is clear that C is a distribution function as described above [?]. Furthermore, it follows that the marginals of C are uniform as $F_i(X_i)$ is uniformly distributed. We thus define a copula in probabilistic terms as

Definition 1.1 (Copula). *A function $C : [0, 1]^d \rightarrow [0, 1]$ is called a copula if it has uniform marginals and is a distribution function for a d -dimensional random vector \mathbf{X} .*

An important and fundamental theorem of copulas for especially continuous random variables where the marginals are also continuous functions is stated by Sklar:

Theorem 1.2 (Sklar's theorem). *For a random vector \mathbf{X} with CDF F and univariate marginal CDFs F_1, \dots, F_d . There exists a copula C such that*

$$F(x_1, \dots, x_d) = C(F_1(x_1), \dots, F_d(x_d)) \quad (1.7)$$

If X is continuous, C is unique; otherwise C is uniquely determined on the Cartesian product of the ranges of distribution functions F_i , $\prod \text{Ran}(F_i)$.

Note that the last statement for non-continuous random variables can be made unique by instead using subcopulas, a generalization of copulas with domain I only a subdomain of the unit hypercube $\mathbb{I}^d = [0, 1]^d$ containing all faces of the unit hyper cube. However, there are infinitely many ways of extending such a subcopula to a copula C [?]. In our case, this means that for discrete and/or mixed variables, we will later have to work around this non-uniqueness when calculating mutual information. The example made by Geenens [?] is a bivariate random vector of independent variables $X \sim \text{Bern}(\pi_X)$ and $Y \sim \text{Bern}(\pi_Y)$. The support of F_X and F_Y is then $\{0, 1 - \pi_X\}$ and $\{0, 1 - \pi_Y\}$ respectively. Due to the restriction on the boundary of the unit square, the only unique point of a copula C is then $(1 - \pi_X, 1 - \pi_Y)$, and by independence we must have

$$C(1 - \pi_X, 1 - \pi_Y) = (1 - \pi_X)(1 - \pi_Y)$$

Geenens then proceed to define an uncountable set of copulas that fulfill the above criterion which further illustrates that the basic concepts of copulas are not well suited for discrete random vectors. Note that in the article it is however

argued how one can extend the concept to a more general concept that works for mixed variables.

From Equation 1.7 we see that a copula is thus simply just a function that *couples* the marginals of a random vector to the joint distribution. The following corollary follows immediately

Corollary 1.2.1 (Coordinate transformation). *Under the assumptions of Theorem 1.2, given any set (T_1, \dots, T_d) of strictly increasing functions, if C is a copula of (X_1, \dots, X_d) then it is also a copula of $(T_1(X_1), \dots, T_d(X_d))$.*

Proof. Suppose (X_1, \dots, X_d) admits a copula C and let T_i be given as stated. Consider coordinate wise the result of the transformation $Y_i = T_i(X_i)$ and consider the CDF $F_{Y_i}(y_i)$

$$\begin{aligned} F_{Y_i}(y_i) &= \mathbb{P}(Y_i \leq y_i) \\ &= \mathbb{P}(T_i^{-1}(Y_i) \leq T_i^{-1}(y_i)) \\ &= \mathbb{P}(X_i \leq T_i^{-1}(y_i)) \\ &= F_{X_i}(T_i^{-1}(y_i)) \end{aligned}$$

The above is easily generalized for a joint distribution as well. Thus, by the existence of a copula C for \mathbf{X}

$$\begin{aligned} F_{\mathbf{Y}}(y_1, \dots, y_d) &= F_{\mathbf{X}}(T_1^{-1}(y_1), \dots, T_d^{-1}(y_d)) \\ &= C(F_{X_1}(T_1^{-1}(y_1)), \dots, F_{X_d}(T_d^{-1}(y_d))) \\ &= C(F_{Y_1}(y_1), \dots, F_{Y_d}(y_d)) \end{aligned}$$

where Sklar's theorem have been used for the second equality. The above shows that C is indeed also a copula for $\mathbf{Y} = (T_1(X_1), \dots, T_d(X_d))$. \square

The above corollary is actually equivalent with a seemingly stronger statement and follows easily

Proposition 1.3. *Since T_i is strictly increasing, the inverse T_i^{-1} exists and is also strictly increasing. Thus, the above implication is bidirectional and hence for strictly increasing functions T_i , C is a copula of (X_1, \dots, X_d) if and only if it is a copula of $(T_1(X_1), \dots, T_d(X_d))$.*

1.2.2 Mutual Information and Copula Entropy

In this section we introduce copula entropy as done in [?] and see how it actually is equal to the well known mutual information (multiplied by -1) and hence as a corollary that mutual information is independent of marginals. Namely, under a coordinate transformation as in Corollary 1.2.1, the mutual information is constant. The name comes from the general definition of (differential) entropy as we shall see shortly. However, first we define mutual information between a set of random variables

Definition 1.4 (Mutual information). *For a random vector $\mathbf{X} = \{X_i\}$, we define the mutual information as*

$$I(\mathbf{X}) = \mathbb{E} \left[\log_b \left(\frac{f(\mathbf{X})}{\prod_i f_i(X_i)} \right) \right]$$

where f is the joint density function with marginals f_i of the random vector \mathbf{X} . The base of the logarithm b is often chosen to be 2, e or 10 although the choice is unimportant as all logarithms are equivalent up to a scaling factor. We shall in general choose $b = e$ and drop the base b from this point on.

We note that later on, as the choice of b will result in a scaling of G_{obs} , but we will also introduce a scaling parameter α for G_{obs} to both ensure the convergence of the algorithm and to control higher order effects, we shall in general choose $b = e$.

An important property of mutual information is that the continuous version is the limit of the discrete mutual information for random (continuous) vector discretized as the mesh size goes to zero i.e. recovering the continuity of the random vector. This is discussed in Subsection 1.2.3. For now, we proceed with the definition of (joint) entropy for both discrete and continuous random vectors.

Definition 1.5 (Entropy). *The (joint) entropy of a random vector \mathbf{X} is defined as*

$$H(\mathbf{X}) = -\mathbb{E}[\log f(\mathbf{X})]$$

In case of a discrete random vector, this is called the Shannon entropy while for continuous random vectors, this is called differential entropy and is often denoted as $h(\mathbf{X})$ instead of $H(\mathbf{X})$.

We note the need for two separate notations of entropy as differential entropy is not the limit of Shannon entropy in the way mutual information is. Again, this is further discussed in Subsection 1.2.3.

Before discussing copula entropy (CE), we note a very useful relation between entropy and mutual information. Indeed, we shall later use this to show that mutual information in the continuous version is the limit of the discretization.

Lemma 1.6 (Mutual information and entropy relation). *For a continuous random vector \mathbf{X} , the (joint) mutual information $I(\mathbf{X})$ can be decomposed into a sum of differential entropies as*

$$I(\mathbf{X}) = \sum_{i=1}^d h_i(X_i) - h(\mathbf{X})$$

where d is the dimension of \mathbf{X} . The same is true for discrete variables but with entropy H instead of differential entropy h .

Proof. This follows immediately from the definition of mutual information and entropy:

$$\mathbb{E} \left[\log \frac{f(\mathbf{X})}{\prod_{i=1}^d f_i(X_i)} \right] = - \sum_{i=1}^d \mathbb{E} [\log f_i(X_i)] + \mathbb{E} [\log f(\mathbf{X})]$$

□

With the definitions of mutual information and entropy we are finally ready to introduce copula entropy.

Definition 1.7 (Copula entropy). *For a continuous random vector \mathbf{X} with a uniquely defined copula C , and copula density c , we define the copula entropy CE of \mathbf{X} as*

$$CE(\mathbf{X}) = h(\mathbf{U})$$

where \mathbf{U} has density c . In particular,

$$CE(\mathbf{X}) = -\mathbb{E} [\log c(\mathbf{U})]$$

As stated above, copula entropy is actually equal to the negative mutual information which we state as a theorem

Theorem 1.8 (Equality of Copula entropy). *For a continuous random vector \mathbf{X} , the copula entropy CE is equal to the negative joint mutual information of \mathbf{X}*

$$CE(\mathbf{X}) = -I(\mathbf{X})$$

Proof. By Theorem 1.2, letting $x_i = F_i^{-1}(u_i)$, we can relate the copula density to the joint density of \mathbf{X} and its marginals

$$\begin{aligned}
 c(u_1, \dots, u_d) &= \frac{\partial}{\partial \mathbf{u}} C(u_1, \dots, u_d) \\
 &= \frac{\partial}{\partial \mathbf{u}} F(F_1^{-1}(u_1), \dots, F_d^{-1}(u_d)) \\
 &= \frac{\partial^d}{\prod_{i=1}^d \partial u_i} F(F_1^{-1}(u_1), \dots, F_d^{-1}(u_d)) \\
 &= \frac{\partial^{d-1}}{\prod_{i=2}^d \partial u_i} \left(\frac{\partial}{\partial x_1} F \right) (F_1^{-1}(u_1), \dots, F_d^{-1}(u_d)) \cdot \frac{1}{f_1(F_1^{-1}(u_1))} \\
 &\vdots \\
 &= \frac{\partial}{\partial \mathbf{x}} F(F_1^{-1}(u_1), \dots, F_d^{-1}(u_d)) \cdot \frac{1}{\prod_{i=1}^d f_i(F_i^{-1}(u_i))} \\
 &= f(F_1^{-1}(u_1), \dots, F_d^{-1}(u_d)) \frac{1}{\prod_{i=1}^d f_i(F_i^{-1}(u_i))}
 \end{aligned}$$

Where for the fourth equality, we have applied the chain rule and that $\frac{d}{dx} f^{-1}(x) = \frac{1}{f'(f^{-1}(x))}$ for any differentiable and invertible function f . The final equality follows by definition of the joint probability density function f in terms of F . We note that there is no need for the Jacobian as $x_i = F_i^{-1}(u_i)$ and hence no need for non-diagonal partial derivatives $\frac{\partial x_i}{\partial u_j}$. It follows directly that

$$\begin{aligned}
 -CE(\mathbf{X}) &= \int_{[0,1]^d} c(\mathbf{u}) \log c(\mathbf{u}) \, d\mathbf{u} \\
 &= \int_{\mathcal{X}} \frac{f(\mathbf{x})}{\prod_{i=1}^d f_i(x_i)} \log \left(\frac{f(\mathbf{x})}{\prod_{i=1}^d f_i(x_i)} \right) \prod_{i=1}^d f_i(x_i) \, d\mathbf{x} \\
 &= \int_{\mathcal{X}} f(\mathbf{x}) \log \left(\frac{f(\mathbf{x})}{\prod_{i=1}^d f_i(x_i)} \right) \, d\mathbf{x} \\
 &= I(\mathbf{X})
 \end{aligned}$$

where $\mathcal{X} = \prod_{i=1}^d \text{dom } F_i$ is the domain of the random vector \mathbf{X} and the third equality follows from a change of variables with the trivial substitution $u_i = F_i(x_i)$ such that $du_i = f_i(x_i) dx_i$ and $x_i = F_i^{-1}(u_i)$. This concludes the proof. \square

Finally, before moving on to correlation as a measure of similarity, we discuss what happens in the limit of mutual information and entropy as we shall later need this as arguments for numerical stability.

1.2.3 Entropy and Mutual Information in the Limit

In this section, we shall discuss the differences between entropy and differential entropy and observe how this difference cancels when computing mutual information. In fact, we shall see that mutual information defined for continuous random vectors is the limit of the discrete version which will be useful later when implementing the algorithm.

First, although one may think differential entropy is the limit of (discrete) entropy, this is not the case. Namely, consider the support of $f(x)$ (here assumed to be the entire real line) binned into intervals i.e. a discretization of the continuous random variable X , which we shall denote X^Δ . To make notation simpler, we shall bin into equal-sized intervals of width Δ . Then, for each interval $[i\Delta, (i+1)\Delta)$ for $i \in \mathbb{Z}$, there exists an x_i such that the probability mass on this interval is represented by this x_i :

$$\mathbb{P}(X^\Delta = x_i) = f(x_i)\Delta = \int_{i\Delta}^{(i+1)\Delta} f(x) dx \quad (1.8)$$

Clearly, this discretization generates a valid distribution as

$$\sum_{i \in \mathbb{Z}} f(x_i)\Delta = \int_{\mathbb{R}} f(x) dx = 1$$

and in the limit, as $\Delta \rightarrow 0$ we recover the original distribution $f(x)$. However, if we try to calculate the entropy of this discretization, denoted by H^Δ , we get a diverging limit

$$\begin{aligned} H^\Delta &= - \sum_{i \in \mathbb{Z}} f(x_i)\Delta \log(f(x_i)\Delta) \\ &= - \sum_{i \in \mathbb{Z}} f(x_i)\Delta \log f(x_i) - \sum_{i \in \mathbb{Z}} f(x_i)\Delta \log \Delta \\ &= - \sum_{i \in \mathbb{Z}} f(x_i)\Delta \log f(x_i) - \log \Delta \end{aligned}$$

Clearly, the first term in the above expression converges to the differential entropy $h(X)$ as $\Delta \rightarrow 0$ whereas $\log \Delta \rightarrow -\infty$ i.e. the expression diverges altogether when differential entropy is well-defined.

A similar argument for the joint entropy between the discretization of X_1 and X_2 (and in principle to any number of dimensions), denoted by H_{12}^Δ , results in

$$H_{12}^\Delta = - \sum_{i,j \in \mathbb{Z}} f(x_1^{(i)}, x_2^{(j)}) \Delta_1 \Delta_2 \log f(x_1^{(i)}, x_2^{(j)}) - \log \Delta_1 - \log \Delta_2$$

where $x_1^{(i)} \in [i\Delta_1, (i+1)\Delta_1)$ and $x_2^{(j)} \in [j\Delta_2, (j+1)\Delta_2)$ are defined such that

$$f(x_1^{(i)}, x_2^{(j)}) \Delta_1 \Delta_2 = \int_{j\Delta_2}^{(j+1)\Delta_2} \int_{i\Delta_1}^{(i+1)\Delta_1} f(x_1, x_2) dx_1 dx_2, \quad \forall i, j \in \mathbb{Z}$$

Note that clearly $(x_1^{(i)}, x_2^{(j)})$ exists for all $i, j \in \mathbb{Z}$. Again, the joint entropy diverges however, when computing the mutual information, we see that the diverging terms cancel. Namely, from Lemma 1.6

$$\begin{aligned} I_{12}^\Delta &= H_1^\Delta + H_2^\Delta - H_{12}^\Delta \\ &= - \sum_{i \in \mathbb{Z}} f_1(\tilde{x}_1^{(i)}) \Delta_1 \log f_1(\tilde{x}_1^{(i)}) - \log \Delta_1 \\ &\quad - \sum_{j \in \mathbb{Z}} f_2(\tilde{x}_2^{(j)}) \Delta_2 \log f_2(\tilde{x}_2^{(j)}) - \log \Delta_2 \\ &\quad + \sum_{i, j \in \mathbb{Z}} f(x_1^{(i)}, x_2^{(j)}) \Delta_1 \Delta_2 \log f(x_1^{(i)}, x_2^{(j)}) + \log \Delta_1 \Delta_2 \\ &= - \sum_{i \in \mathbb{Z}} f_1(\tilde{x}_1^{(i)}) \log f_1(\tilde{x}_1^{(i)}) \Delta_1 - \sum_{j \in \mathbb{Z}} f_2(\tilde{x}_2^{(j)}) \log f_2(\tilde{x}_2^{(j)}) \Delta_2 \\ &\quad + \sum_{i, j \in \mathbb{Z}} f(x_1^{(i)}, x_2^{(j)}) \log f(x_1^{(i)}, x_2^{(j)}) \Delta_1 \Delta_2 \\ &\rightarrow h(X_1) + h(X_2) - h(X_1, X_2) \text{ as } \Delta_1, \Delta_2 \rightarrow 0 \end{aligned}$$

Thus, the limit of the mutual information for discrete random variables is indeed the mutual information defined for continuous random variables and can be computed either as the limit of discretizing the probability density function and then computing entropies or just using the initial definition for (discrete) mutual information in Definition 1.4. In particular, mutual information for discrete and random variables are comparable such that it makes sense define mutual information between mixed random variables. For a more rigorous treatment of this, we refer to [?] where they define mutual information between discrete and continuous random variables from a measure theoretical point of view.

Before continuing, we discuss the case where X_1 is equal to X_2 . In this case, discretizing with a common Δ we have that

$$f(x_1^{(i)}, x_2^{(j)}) \Delta^2 = \int_{j\Delta}^{(j+1)\Delta} \int_{i\Delta}^{(i+1)\Delta} f(x_1, x_2) dx_1 dx_2, \quad \forall i, j \in \mathbb{Z}$$

Clearly, the above integral is 0 for $i \neq j$. Although $f(x_1, x_2)$ is not well-defined in the usual functional sense, extending to distribution, we might write $f(x_1, x_2) = f(x_2|x_1)f(x_1)$. In terms of distributions, it works to put $f(x_2|x_1) = \delta(x_2 - x_1)$ where δ is the *Dirac delta* distribution, as then $\int_{\mathbb{R}} f(x_1, x_2) dx_2 = f(x_1)$ and

$f(x_1, x_2)$ is "0" when $x_1 \neq x_2$. I.e. the marginals and probability mass are correct. Then, when calculating the above integral, we get that

$$\begin{aligned} f\left(x_1^{(i)}, x_1^{(i)}\right) \Delta^2 &= \int_{i\Delta}^{(i+1)\Delta} \int_{i\Delta}^{(i+1)\Delta} f(x_1, x_2) dx_1 dx_2 \\ &= \int_{i\Delta}^{(i+1)\Delta} f(x_1) dx_1 \\ &= f\left(\tilde{x}_1^{(i)}\right) \Delta \end{aligned}$$

Thus, when calculating I^Δ for two identical variables, we obtain

$$\begin{aligned} I^\Delta &= - \sum_{i \in \mathbb{Z}} f\left(\tilde{x}_1^{(i)}\right) \log f\left(\tilde{x}_1^{(i)}\right) \Delta - \sum_{j \in \mathbb{Z}} f\left(\tilde{x}_2^{(j)}\right) \log f\left(\tilde{x}_2^{(j)}\right) \Delta \\ &\quad + \sum_{i \in \mathbb{Z}} f\left(\tilde{x}_1^{(i)}\right) \log f\left(\tilde{x}_1^{(i)}\right) \Delta - \log \Delta \\ &\rightarrow \infty \text{ as } \Delta \rightarrow 0 \end{aligned}$$

Thus in practice, it would not make much sense to compare equal variables or even a random vector only defined on a lower dimensional manifold as we would get an infinite copula entropy.

1.2.4 Correlation

At this point, we have a good understanding of copula entropy/mutual information for calculations later on. However, another typical measure of similarity is correlation which is easily estimated from sample data. However, in this section we show that in general, we can not compute the correlation coefficient from a copula which we saw above is the case for mutual information. Namely, given a copula C for some set of random variables $\{X_i\}_{i \in I}$ indexed by finite I , one can not calculate ρ between any pair (X_i, X_j) , $i \neq j$ from the copula. This is easily shown by the following argument.

First, note that from Corollary 1.2.1, C is also a copula for $Z_i := (X_i - \mu_i) / \sigma_i$ for $i \in I$ where $\mu_i = \mathbb{E}[X_i]$ and $\sigma_i = \sqrt{\text{Var } X_i}$ (assuming that these exist). Clearly, the correlation coefficient for Z_i and Z_j is the same as between X_i and X_j . We thus proceed trying to calculate the correlation between any pair Z_i and Z_j .

$$\begin{aligned} \rho_{ij} &= \int \int_{\mathbb{R}^2} z_i z_j f_{ij}(z_i, z_j) dz_i dz_j \\ &= \int \int_{[0,1]^2} F_i^{-1}(u_i) F_j^{-1}(u_j) c_{ij}(u_i, u_j) du_i du_j \end{aligned} \tag{1.9}$$

where c_{ij} density version of the copula defined for X_i and X_j and F_i and F_j are the marginals of Z_i and Z_j with mean 0 and variance 1. From the above, it is then clear for a fixed, non-constant copula C , the correlation depends on the marginals of X_i and X_j . Also, we see that a constant copula density (only admissible if $c \equiv 1$ on $[0, 1]^2$ and 0 elsewhere) always results in $\rho_{ij} = 0$ as

$$\int_0^1 F^{-1}(u) du = \int_{\mathbb{R}} z f(z) dz = 0$$

where the final equality follows from the construction of Z_i .

Thus, we conclude that indeed mutual information and correlation are very different measures of codependency. Namely, mutual information does not depend on the marginal distributions whereas from Equation 1.9 we see that correlation does. Thus, it does not make much sense to introduce copulas in the setting of correlation albeit at this point we do not favor one measure above the other. Only if marginals should be insignificant to the network, copula entropy is at this point preferred.

1.3 Copula Based Network Discovery

In this section, we will present the general algorithm and discuss some of its properties regarding uncertainty and convergence. We will focus on using mutual information i.e. copula entropy as the measure of similarity but other measures such as correlation can be interchanged at will in the general algorithm.

By Theorem 1.8 we can compute the mutual information from observed data from the copula. Namely, let CE_{ij} denote the (pairwise) copula entropy of variables X_i and X_j . We shall then set

$$G_{obs} = \begin{bmatrix} 0 & -CE_{12} & \dots & -CE_{1n} \\ -CE_{21} & 0 & \dots & -CE_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ -CE_{n1} & -CE_{n2} & \dots & 0 \end{bmatrix} \quad (1.10)$$

where n is the number of nodes in the graph i.e. random variables that we have observed. Notice that we have chosen the diagonal elements as 0 since information between a random variable X and itself is not really well-defined and when trying to compute this numerically, we observe diverging results as also discussed in the previous section. Furthermore, only the information that propagates through the network is of interest and so setting 0 in the diagonal avoids a bias when deconvolving the information or any similarity in general.

Especially for mutual information where the information between a variable and itself diverges to ∞ thus in the limit, from Equation 1.6, we would get the identity matrix which does not tell us much about the direct dependencies.

Algorithm 1 then follows immediately from Equation 1.10

Algorithm 1 G_{obs} computation

Require: $n > 0$ \triangleright Number of variables
 $G_{obs} \leftarrow \mathbf{0}$
for $1 \leq i, j \leq n \mid i \neq j$ **do**
 Estimate F_i and F_j from $x_i^{\mathcal{D}}$ and $x_j^{\mathcal{D}}$
 $u_i^{\mathcal{D}} \leftarrow F_i(x_i^{\mathcal{D}})$
 $u_j^{\mathcal{D}} \leftarrow F_j(x_j^{\mathcal{D}})$
 Estimate c_{ij} from $u_i^{\mathcal{D}}$ and $u_j^{\mathcal{D}}$
 Compute NCE_{ij}
 $[G_{obs}]_{ij} \leftarrow -NCE_{ij}$
end for
return G_{obs}

Namely, for each entry in G_{obs} , except for the diagonal elements, first estimate the cumulative distributions of X_i and X_j based on samples $x_i^{\mathcal{D}}$. Then, transform the samples by the estimated distribution function to obtain corresponding uniform samples. This may be done outside the loop to increase computational efficiency. From the paired samples $(x_i^{\mathcal{D}}, x_j^{\mathcal{D}})$, estimate the copula density c_{ij} and finally use this to compute the mutual information/copula entropy. Methods for estimating the densities and in continuation hereof the distribution functions are presented in Section 1.4. The negative copula entropy is then recorded in (i, j) entry of G_{obs} . We note that the algorithm can be optimized for symmetric measures such as copula entropy itself, to only loop through $i < j$ and saving the computed entropy in the (j, i) entry as well. Also, as copula entropy diverges as X_i and X_j are jointly distributed closer to a one-dimensional manifold, ideally there should be a check for such or the user should check the paired observations to exclude such variable combinations.

From Subsection 1.2.3, to calculate the (joint) copula entropy of a continuous random vector, we simply discretize the domain of each random variable and use the estimated copula density evaluated at these points to estimate the total copula entropy. Furthermore, if one or more elements of the random vector are mixed random variables, we choose the discrete events to be their own bins and discretize the rest or in the context of Algorithm 1 only estimate the distribution functions for the continuous component of the random variable. This works due to the copula entropy for continuous random variables being the limit of the

discretization and as such, the copula entropy is well-defined for mixed random variables as well.

We continue with an example of how this discretization of a mixed random variable would work. Notice that we only have a discrete event (an atom) at 0 as this resembles the observed behavior of the delays, although the example could be extended to more complex discrete distributions.

Example 1.1 (Discretization of mixed random variable). *Let X be a mixture of an atom in e.g. 0 and an exponential with parameter λ with proportions p and $1-p$. Then, a discretization of X is 0 with probability mass p and the remaining support $(0, \infty)$ discretized in some way with total probability mass $1-p$ and each bin having probability according to Equation 1.8 scaled with $1-p$. If the bin size is a constant Δ , then for the discretized variable X^Δ , we have $\mathbb{P}(X^\Delta = 0) = p$ and $\mathbb{P}(X^\Delta = x_i) = (1-p) \exp(-\lambda i \Delta) (1 - \exp(-\lambda \Delta))$, where x_i is given by*

$$x_i = i\Delta + \frac{1}{\lambda} (\log(\lambda\Delta) - \log(1 - e^{-\lambda\Delta})), \quad i \in \mathbb{N}_0$$

1.3.1 Network Deconvolution

At this point, we have obtained a convolved matrix of information G_{obs} and are ready to use Equation 1.6. We present the original algorithm from [?] in the case G_{obs} is symmetric and hence diagonalizable by an orthogonal matrix U . The original `Matlab` implementation was translated to `Python` and is summarized in the following pseudocode.

Algorithm 2 (ND) Network Deconvolution

Require: G_{obs}, α, β
 $[G_{obs}]_{ii} \leftarrow 0, \forall i \in \{1, \dots, d\}$ \triangleright ensure zero-diagonal
 $[G_{obs}]_{ij} \leftarrow 0$, when $[G_{obs}]_{ij} < Q_\alpha(G_{obs})$
 Compute eigendecomposition U, Λ of G_{obs}
 $\lambda^+ \leftarrow \max(\lambda^{\max}, 0)$
 $\lambda^- \leftarrow \min(\lambda^{\min}, 0)$
 $k^+ \leftarrow \frac{1-\beta}{\beta} \lambda^+$
 $k^- \leftarrow \frac{1+\beta}{\beta} \lambda^-$
 $c_s^{-1} \leftarrow \max(k^+, -k^-)$
 $\hat{\Lambda} \leftarrow \Lambda (c_s^{-1} I + \Lambda)^{-1}$
return $U \hat{\Lambda} U^T$

where $Q_\alpha(G_{obs})$ denotes the α quantile of the strictly upper (or lower due to

symmetry) triangular part of G_{obs} . We note the two extra parameters *alpha* and β which we will discuss shortly. In particular, the paper contains conflicting information on how to find β from how it is defined. Furthermore, they include some analysis on the robustness of the above deconvolution algorithm but only in a somewhat particular case and with some confusion on matrix norms and spectral radius. This analysis on robustness, we will extend and clarify in the following Subsection 1.3.3.

From the definition of $Q_\alpha(G_{obs})$ it is clear that the α parameter is a filter on the observed edges and is useful if one wants to filter out insignificant observations. However, in practice, as we will see, it is often not very influential except for large α (corresponding to many edges set to 0) as small perturbations from e.g. imperfect calculations should not influence the results for fairly conditioned matrices as we shall observe in ???. Thus, setting $\alpha = 0$ retains all values in G_{obs} after setting the diagonal equal to 0. As a technical detail, we note that the `quantile` function from NumPy (v. 1.26.4) has been used to find this quantile as quantiles can be defined in many ways from a data set.

Finally, we note that the $\beta \in (0, 1)$ parameter corresponds to a scaling of G_{obs} such that the resulting spectral norm of G_{dir} is β . From Algorithm 2 it is seen that it serves as a regularization on the eigenvalues of G_{obs} and although this is discussed in [?], their results do not conform with their implementation, and we thus comment on this and what else could be done to ensure convergence of the algorithm in the following section. Also, in practice we choose a threshold t on the elements of G_{dir} returned from Algorithm 2 to further filter out insignificant direct dependencies.

1.3.2 Ensuring Convergence and the Effect of β

In this section we will further discuss the effect of β and how the steps for rescaling the observed similarity matrix G_{obs} are derived. In particular, we will reformulate the original derivation from [?] as there is a discrepancy between their code¹ and their proof of choosing a scaling parameter c_s of G_{obs} . Namely, denote \tilde{G}_{obs} as the rescaled G_{obs} such that $\tilde{G}_{obs} = c_s G_{obs}$. Choosing c_s as in Algorithm 2 i.e. $c_s^{-1} = \max\left(\frac{1-\beta}{\beta}\lambda^+, -\frac{1+\beta}{\beta}\lambda^-\right)$ where λ^+ is the largest positive eigenvalue of G_{obs} (and 0 if no eigenvalue is positive) and λ^- is the most negative eigenvalue of G_{obs} (and 0 if no eigenvalue is negative) then implies \tilde{G}_{dir} obtained from the new \tilde{G}_{obs} has spectral radius $\beta < 1$ i.e. a proper G_{dir} with the largest numerical eigenvalue equal to β . This holds in general and not only for symmetric G_{obs} as we will see in the following. However, when G_{obs} is symmetric

¹<https://compbio.mit.edu/nd/>

the resulting \tilde{G}_{dir} can easily be expressed through the eigendecomposition of G_{obs} , U , Λ as

$$\begin{aligned}\tilde{G}_{dir} &= \tilde{G}_{obs} \left(I + \tilde{G}_{obs} \right)^{-1} \\ &= c_s G_{obs} (I + c_s G_{obs})^{-1} \\ &= U c_s \Lambda U^T (U U^T + U c_s \Lambda U^T)^{-1} \\ &= U c_s \Lambda U^T U (I + c_s \Lambda)^{-1} U^T \\ &= U \Lambda (c_s^{-1} I + \Lambda)^{-1} U^T\end{aligned}$$

which can also be seen in Algorithm 2. Thus, with everything else explained about the algorithm, we show that the resulting \tilde{G}_{dir} in general have spectral radius β .

Let (λ, v) be an eigenpair of G_{obs} with $\lambda \neq 0$, it then follows that $\left(\frac{\lambda}{c_s^{-1} + \lambda}, v \right)$ is an eigenpair of \tilde{G}_{dir} . Then, following the arguments in [?] (which we have redone to know why the original implementation and derivation differs), we obtain that for a λ in $[0, \infty)$, we must have that

$$c_s^{-1} \geq \frac{1 - \beta}{\beta} \lambda^+$$

and similarly for $\lambda \in (-c_s^{-1}, 0)$

$$c_s^{-1} \geq -\frac{1 + \beta}{\beta} \lambda^-$$

for $\lambda < -c_s^{-1}$ we obtain that the resulting eigenvalue is larger than 1 hence we must also have that $c_s^{-1} \geq -\lambda^-$ which is covered by the above constraint on c_s^{-1} . Thus, the smallest c_s^{-1} we can choose to ensure that $\rho(\tilde{G}_{dir}) \leq \beta$ is by $c_s^{-1} = \max\left(\frac{1 - \beta}{\beta} \lambda^+, -\frac{1 + \beta}{\beta} \lambda^-\right)$ which also implies that $\rho(\tilde{G}_{dir}) = \beta$ as either the most negative or most positive eigenvalue is mapped to β or $-\beta$ respectively. This coincides with the original implementation, noting that some error has been made in the original discussion of the parameter β in [?]. Furthermore, we note that if we just want the algorithm to converge, as we discussed before, this is equivalent to $\sigma(\tilde{G}_{obs}) \subseteq (-1/2, \infty)$, so really, we can just choose $c_s^{-1} = -(2 + \delta)\lambda^-$ for some small δ if $\lambda^- < -1/2$ and otherwise not scale G_{obs} to preserve the structure. Finally, we note that as β tends to 0, higher order interactions become less significant as can clearly be seen from Equation 1.4. Thus, β also allows us to tune how much influence higher order interactions should have and one should try different β to see how influenced results are to higher order effects.

1.3.3 Robustness to Noise

Finally, before discussing how to compute and estimate the mutual information between two random variables based on observations, we turn our heads to error analysis of the deconvolution algorithm. It is important to understand how well the algorithm performs subject to noise and errors. Namely, in the case of mutual information, the assumption that higher order effects can be calculated as a sum of matrix powers of the direct effects does not hold. Thus, if we can quantify the error in G_{obs} , we can from the following analysis quantify the resulting error in G_{dir} . We shall first discuss the original result from [?], correcting some errors in terms of definitions and see how their result can also be expressed as an absolute upper bound on the error instead of only how this error behaves for small perturbations. Furthermore, we shall extend their result to not only hold when $\rho(G_{obs}) < 1$ and $\rho(G_{obs} + N) < 1$ where N is some noise e.g. from computation or assumptions that does not completely hold.

The original result states that $\|G_{dir} - \tilde{G}_{dir}\|_2 \leq \gamma + \mathcal{O}(\delta^2 + \gamma^2 + \delta\gamma)$ where $\|\cdot\|_2$ is the Euclidean norm also known as the spectral norm as this is equal to the largest singular value of the input matrix. However, they note that the Euclidean norm of a matrix M is equal to $\sqrt{\sum_{i,j} m_{ij}^2}$ which is incorrect. This is the Frobenius norm, and instead it should have been defined as

$$\|M\|_2 = \sup_{\|x\|_2=1} \|Mx\|_2 = \sigma_{\max}(M)$$

They then proceed to let γ be the largest absolute eigenvalue of N and δ the largest absolute eigenvalue of $\tilde{G}_{obs} = G_{obs} + N$ however as the noise may be both positive and negative, it is easier to define δ as the largest absolute eigenvalue of G_{obs} instead which we will do in the following. We note that γ and δ are not the spectral/Euclidean norm of N and G_{obs} respectively as in general, we only have $\rho(M) \leq \|M\|_2$. However, if G_{obs} and N are both (real) symmetric matrices, then the spectral norms are equal to the largest absolute eigenvalues of G_{obs} and N respectively. Thus, if instead one wanted to measure the difference in the direct dependency matrices in terms of e.g. the Frobenius norm, it is important to differentiate between the spectral radius and the norm that is actually being used. Finally, before constructing the actual upper bound on the error instead of quantizing the asymptotic behavior for small γ , we note that $\|\cdot\|_2$ is a sub-multiplicative matrix norm defined as bellow ([?]), and that we shall assume that $\rho(G_{obs}), \rho(\tilde{G}_{obs}) < 1$.

Definition 1.9 (Sub-multiplicative Matrix norm). *A matrix norm $\|\cdot\|$ is said to be sub-multiplicative, if for every $A, B \in \mathbb{F}^{n \times n}$ where \mathbb{F} is either the real or complex field:*

$$\|AB\| \leq \|A\| \cdot \|B\|$$

As we do not use any property of the spectral norm except that it is sub-multiplicative, we shall consider any norm $\|\cdot\|$ in general that is also sub-multiplicative. Thus, consider the norm of the difference $G_{dir} - \tilde{G}_{dir}$:

$$\begin{aligned}
\|G_{dir} - \hat{G}_{dir}\| &= \left\| G_{obs} (I + G_{obs})^{-1} - \hat{G}_{obs} (I + \hat{G}_{obs})^{-1} \right\| \\
&= \left\| -\sum_{k \geq 1} (-G_{obs})^k + \sum_{k \geq 1} (-\hat{G}_{obs})^k \right\| \\
&\leq \sum_{k \geq 1} \left\| G_{obs}^k - (\hat{G}_{obs})^k \right\| \\
&\leq \sum_{k \geq 1} \sum_{i=1}^k \binom{k}{i} \|N\|^i \|G_{obs}\|^{k-i} \\
&= \sum_{k \geq 1} \sum_{i=1}^k \binom{k}{i} \gamma^i \delta^{k-i} \\
&= \sum_{k \geq 1} \left((\gamma + \delta)^k - \delta^k \right) \\
&= \frac{\gamma + \delta}{1 - \gamma - \delta} - \frac{\delta}{1 - \delta} \\
&= \frac{\gamma}{(1 - \gamma - \delta)(1 - \delta)}
\end{aligned} \tag{1.11}$$

where in the second to last inequality, we assume that $\gamma + \delta < 1$ as then both $\sum (\gamma + \delta)^k$ and $\sum \delta^k$ converges as $\gamma + \delta \geq \delta \geq 0$ and hence also the difference of the sums converges. Also, the second equality uses that the spectral norm of G_{obs} and \tilde{G}_{obs} is less than 1 in order to express the inverses as infinite series. Thus, the above bound on the difference $G_{dir} - \tilde{G}_{dir}$ does not hold in every case. Namely, for fixed γ , the bound tends to ∞ as $\delta \rightarrow 1$. Furthermore, we note that the final infinite sum diverges whenever $\gamma + \delta > 1$ through the following argument using the ratio test for infinite sums which is needed because we can not conclude on the convergence of a difference of diverging sums solely from

the fact that the individual sums diverge:

$$\begin{aligned}
\lim_{n \rightarrow \infty} \left| \frac{(\gamma + \delta)^{n+1} - \delta^{n+1}}{(\gamma + \delta)^n - \delta^n} \right| &= \lim_{n \rightarrow \infty} \left| \frac{(\gamma + \delta) \left(1 + \frac{\gamma}{\delta}\right)^n - \delta}{\left(1 + \frac{\gamma}{\delta}\right)^n - 1} \right| \\
&= \lim_{n \rightarrow \infty} \left| \delta + \gamma \frac{\left(1 + \frac{\gamma}{\delta}\right)^n}{\left(1 + \frac{\gamma}{\delta}\right)^n - 1} \right| \\
&= \lim_{n \rightarrow \infty} \left| \delta + \gamma \frac{1}{1 - \left(1 + \frac{\gamma}{\delta}\right)^{-n}} \right| \\
&= |\gamma + \delta| = \gamma + \delta
\end{aligned}$$

assuming that $\gamma, \delta > 0$ corresponding to neither N nor G_{obs} is the zero matrix in which case the above analysis is nonsensical.

Before continuing with a more general bound on the error, we first note that examples of sub-multiplicative matrix norms are every induced norm such as the spectral norm and the Frobenius norm which is often useful when interpreting error. Also, the max norm is *not* sub-multiplicative, but a scaled version is (which is true for any matrix norm from the fact that all matrix norms are equivalent).

Now, consider the general case, where we do not restrict the spectral radius of either G_{obs} or N except such that G_{obs} and \tilde{G}_{obs} admits direct similarity matrices G_{dir} and \tilde{G}_{dir} (with spectral radius less than 1). To obtain a more general result, we shall use the following result from [?], which is very useful when doing matrix perturbation analysis.

Theorem 1.10 (Inverse of sum of matrices). *Let $A, B \in \mathbb{R}^{n \times n}$ such that A and $A + B$ are invertible. Then the inverse of $A + B$ can be expressed as*

$$(A + B)^{-1} = A^{-1} - A^{-1}B(A + B)^{-1}$$

The proof of the above is simple through direct computation. Hence, we continue to once again consider the difference $G_{dir} - \tilde{G}_{dir}$

$$\begin{aligned}
G_{dir} - \tilde{G}_{dir} &= G_{obs} (I + G_{obs})^{-1} - (G_{obs} + N) (I + G_{obs} + N)^{-1} \\
&= G_{obs} \left((I + G_{obs})^{-1} - (I + G_{obs} + N)^{-1} \right) - N (I + G_{obs} + N)^{-1} \\
&= G_{obs} (I + G_{obs})^{-1} N (I + G_{obs} + N)^{-1} - N (I + G_{obs} + N)^{-1} \\
&= - (I + G_{obs})^{-1} N (I + G_{obs} + N)^{-1}
\end{aligned}$$

where the third equality follows from Theorem 1.10. This way, we have a simple exact expression for the difference without any further assumptions on G_{obs} and

N . Now, under a sub-multiplicative norm $\|\cdot\|$ we can bound the norm of the difference in the following way.

$$\left\|G_{dir} - \tilde{G}_{dir}\right\| \leq \|N\| \left\|(I + G_{obs})^{-1}\right\| \left\|(I + G_{obs} + N)^{-1}\right\| \quad (1.12)$$

We note that if once again, we assume that the spectral radius of G_{obs} and $G_{obs} + N$ are smaller than 1, we rediscover Equation 1.11. Equation 1.12 also shows that in general, if N is small or G_{obs} is large we should observe small errors which is also what we would expect intuitively. The above result is also very useful when later on in ?? we discuss the error from using mutual information in the case of a multi-variate Gaussian.

From Equation 1.12, by another application of Theorem 1.10, we find the relative error in general to be bounded as follows

$$\frac{\left\|G_{dir} - \tilde{G}_{dir}\right\|}{\left\|G_{dir}\right\|} \leq \|N\| \left|1 - \frac{\|I\|}{\left\|G_{obs} (I + G_{obs})^{-1}\right\|}\right| \left\|(I + G_{obs} + N)^{-1}\right\|$$

Finally, before discussing the methods for estimating the copula density, we comment on some frequently used matrix norms and show some explicit bounds on the error only using the difference of G_{obs} and \tilde{G}_{obs} , N . Namely, we shall consider the max norm and Frobenius norm of the difference $G_{dir} - \tilde{G}_{dir}$ and note that from [?], we can relate the Euclidean norm to the Frobenius and max norm in the following way. Namely, for any matrix $A \in \mathbb{R}^{n \times n}$ it holds that

$$\|A\|_2 \leq \|A\|_F \leq \sqrt{n} \|A\|_2$$

$$\|A\|_{\max} \leq \|A\|_2 \leq n \|A\|_{\max}$$

Finally, if G_{obs} and N are symmetric, the singular values are equal to the absolute eigenvalues for G_{obs} and \tilde{G}_{obs} and because $\sigma(I + G_{obs}), \sigma(I + \tilde{G}_{obs}) \subseteq (1/2, \infty)$ implies $\sigma\left((I + G_{obs})^{-1}\right), \sigma\left((I + \tilde{G}_{obs})^{-1}\right) \subseteq (0, 2)$ we infer that $\left\|(I + G_{obs})^{-1}\right\|_2, \left\|(I + \tilde{G}_{obs})^{-1}\right\|_2 \leq 2$. Using this with the above equivalence on the Euclidean norm with Equation 1.12, we conclude that

$$\begin{aligned} \left\|G_{dir} - \tilde{G}_{dir}\right\|_F &\leq 4\sqrt{d} \|N\|_F \\ \left\|G_{dir} - \tilde{G}_{dir}\right\|_{\max} &\leq 4d \|N\|_{\max} \end{aligned} \quad (1.13)$$

This clearly shows us that for small networks (thus small d) we risk smaller errors in terms of the Frobenius and max norm (which is not surprising) which

are clearly interpreted through the difference of individual element of G_{dir} and \tilde{G}_{dir} and that the max norm scales linearly with the number of nodes while the Frobenius difference only scales with the square root of the number of nodes.

1.4 Estimating Mutual Information

For Algorithm 1 to work, we need a good and preferably fast estimator of mutual information. [?] proposes to use B-splines for this based on [?] which we shall describe in the following section. It is however quickly apparent that this estimator has some problems when computing mutual information based on the Copula representation. We extend the method to other splines but end up using kernel density estimators (KDE) as they can be regularized in a continuous manner and as a result of this in general show great performance.

1.4.1 B-splines

A spline is in general a piecewise polynomial [?]. We say that a spline is of order $p + 1$ if the piecewise polynomials are of order p . A particular and widely used type of splines are B-splines which are a basis for all splines such that any spline can be expressed as a linear combination of B-splines. Five B-splines of degree 3 (order 4) can be seen in Figure 1.3(a) where we denote $B_{i,p}$ as a B-spline of degree p . The index i comes from the following. Namely, let $i \in \{1, \dots, m + p + 1\}$, we define knots t_i as where pieces of polynomials meet such that

$$B_{i,p}(t) = \begin{cases} \text{non-zero}, & t \in [t_i, t_{i+p+1}) \\ 0, & \text{otherwise} \end{cases}$$

uniquely defines m splines on $[t_{p+1}, t_m]$. If we furthermore constrain the splines such that

$$\sum_{i=1}^m B_{i,p}(t) = 1$$

The B-splines $B_{i,p}$ can then be evaluated at some t through recursion by the Cox-deBoor recursion formula:

$$B_{i,0}(t) = 1, \quad t \in [t_i, t_{i+1}) \quad (1.14)$$

$$B_{i,k}(t) = \frac{t - t_i}{t_{i+k} - t_i} B_{i,k-1}(t) + \frac{t_{i+k+1} - t}{t_{i+k+1} - t_{i+1}} B_{i+1,k-1}(t) \quad (1.15)$$

The fact that the B-splines sum to 1 and that we have m splines of degree p

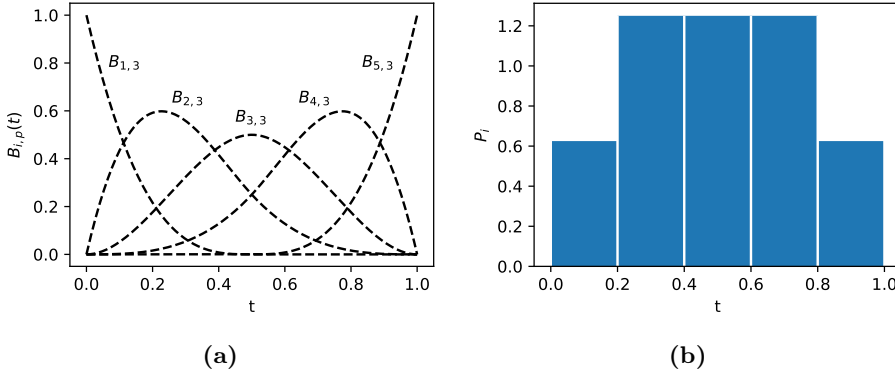


Figure 1.3

is what we will use to estimate the density and in return the mutual information between two random variables based on observations. Namely, suppose we want to compute the mutual information, this can then be done by discretizing the random variables by binning. We assign the probability mass $P_{i,j}$ to bin $(i,j) \in \{1, \dots, m\} \times \{1, \dots, m\}$ as the fraction of observations in the domain corresponding to that bin. In particular, using Copula entropy, we would divide the unit interval into m^2 equal bins. As we saw earlier, in theory, increasing the number of bins will result in a more and more exact estimate of the mutual information. However, with limited observations, this is not the case as in the limit, the bins would not represent the true underlying distribution due to the finite number of observations. Namely, we would observe a few bins with 1 observation and many with none. As an example, suppose we have n observations and m bins in both dimensions. Then, for m large enough, $P_{i,j} = \frac{1}{n}$ for a hundred distinct bins as well as the marginal probability masses $P_i = \frac{1}{n}$ and $P_j = \frac{1}{n}$. But then,

$$\begin{aligned}
 I(X_1^\Delta, X_2^\Delta) &= -\sum_{i=1}^m P_i \log P_i - \sum_{j=1}^m P_j \log P_j + \sum_{i,j=1}^m P_{i,j} \log P_{i,j} \\
 &= -n \left(\frac{1}{n} \log \frac{1}{n} \right) \\
 &= \log n
 \end{aligned}$$

which is clearly independent of the true underlying distribution. However, if for each observation x_j we assign it to bin i with probability mass $B_{i,p}(x_j)$ this problem is mitigated to some extent as long as m is not too large. It follows that in total each x_j is assigned to all m bins with a combined mass 1 as $\sum_i B_{i,p}(x_j) = 1$. Thus, let $(\mathbf{X}^{(1)}, \mathbf{X}^{(2)})$ be a pair of random vectors each of d i.i.d variables representing observations drawn from a joint distribution, we

then define the B-spline density estimator (i.e. a random variable) for $\mathbf{X}^{(1)}$ as

$$P_i^{(1)} = \frac{1}{d} \sum_{j=1}^d B_{i,p} \left(X_j^{(1)} \right), \quad i \in \{1, \dots, m\}$$

and similarly for $\mathbf{X}^{(2)}$. Furthermore, the B-spline joint density estimator is given by

$$P_{i,j} = \frac{1}{d} \sum_{k=1}^m B_{i,p} \left(X_k^{(1)} \right) B_{j,p} \left(X_k^{(2)} \right)$$

i.e. a product of the probability masses such that $\sum_{i,j=1}^d P_{i,j} = 1$ and the marginal probability masses are given by $P_i^{(1)}$ and $P_j^{(2)}$ respectively.

However, there is still one problem. Namely, the marginals are not uniform when $p > 0$ when $\mathbf{X}^{(k)}$ is. This can also be seen from Figure 1.3(b). This is especially bad when computing Copula entropy as for e.g. a Gaussian random vector we would dramatically underestimate the mutual information as we shall also see in the next chapter. To see that this is the case, consider the expectation of $P_i^{(k)}$ for $k \in \{1, 2\}$:

$$\mathbb{E} \left[P_i^{(k)} \right] = \frac{1}{d} \sum_{j=1}^d \mathbb{E} \left[B_{i,p} \left(X_j^{(k)} \right) \right] = \int_0^1 B_{i,p}(x) dx$$

Indeed, we would want this to be $\frac{1}{m}$ but from the Cox-deBoor recursion formula, we see that this is not the case. Namely, by choosing the knots as in [?] we see that the bins close to the boundary have too little probability mass. Thus, we turn our head to another family of splines called M-splines.

1.4.2 M-splines

Another known family of splines called M-splines have exactly the desired property of equal integrals. Namely, the M-splines $M_{i,p}$ all have unit integrals. Thus, rescaling with $\frac{1}{m}$ results in a family of splines $\tilde{M}_{i,p}$ such that on average we have that $P_i = \frac{1}{m}$. The M-splines equivalent to those of Figure 1.3 are shown in Figure 1.4. Indeed, we see that the marginals are uniform. M-splines can similarly to B-splines be computed recursively by the following [?]

$$M_{i,0}(t) = \frac{1}{t_{i+1} - t_i}, \quad t \in [t_i, t_{i+1}) \quad (1.16)$$

$$M_{i,k}(t) = \frac{k((t - t_i) M_{i,k-1}(t) + (t_{i+k} - t) M_{i+1,k-1}(t))}{(k-1)(t_{i+k} - t_i)} \quad (1.17)$$

However, now the problem is that the splines no longer sum to 1 (after rescaling with $\frac{1}{n}$). Namely, we can no longer guarantee that the probability masses $P_{i,j}^M$ based on the M-spline sum to 1. Thus, a renormalization is needed to ensure a proper probability mass function. From Figure 1.5 we see that the effect of $\sum_{i=1}^m \tilde{M}_{i,p}(x) \neq 1$ in this case is that observations on the interior are smoothed more than those near the boundary.

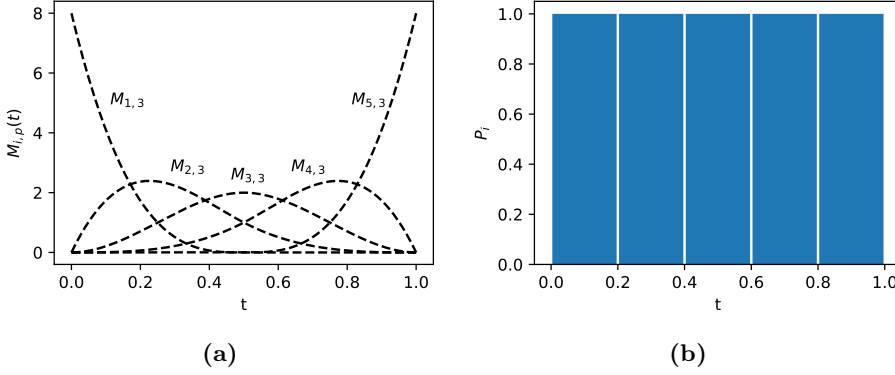


Figure 1.4: P_i is area of each rectangle i.e. 0.2.

This can however be a useful property especially for a bivariate Gaussian for which the Copula density, as we shall see, have peaks at $(0, 0)$ and $(1, 1)$ while being relatively smooth elsewhere.

We note that through construction, one can create a family of splines that both sum to 1 have integrals $\frac{1}{m}$. However, as these spline-based smoothing methods does not perform well in general perform very differently for different m with the lack of continuously varying this parameter that acts as a regularizing parameter, as we shall see in ??, we only present the method for constructing such splines in ?? and instead consider a more general way of estimating the Copula density, namely kernel density estimators.

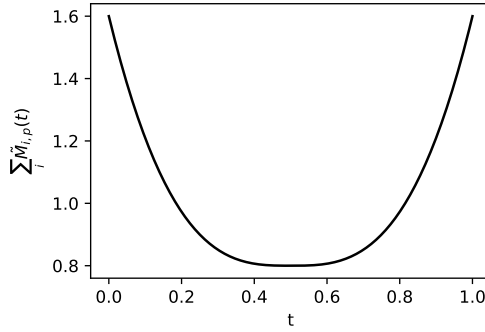


Figure 1.5

1.4.3 Naïve KDE

As we shall see in ??, if one can accurately determine the Copula density of X_1 and X_2 , then using an approximation of the integral, one can calculate the mutual information to any precision wanted. This clearly follows from the above analysis regarding the behavior of the discretization of X_1 and X_2 in the limit as the mesh gets more fine. Thus, if we can estimate the joint Copula density well, we obtain a good estimate of the mutual information of X_1 and X_2 . A widely used non-parametric method for density estimation is kernel density estimation. Namely, if $\{\mathbf{x}_i\}$ is a set of n , d -dimensional observations from a population, i.e. \mathbf{x}_i can be both scalars and vectors in the case of a multidimensional distribution, the kernel density estimator (KDE) of the probability density function is in general given as

$$\hat{f}(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\prod_{j=1}^d \mathbf{h}_{i,j}} K\left(\frac{\mathbf{x} - \mathbf{x}_i}{\mathbf{h}_i}\right) \quad (1.18)$$

where \mathbf{h}_i is the bandwidth (vector) associated with observation \mathbf{x}_i and K is the kernel (function), defined on the domain of \mathbf{X} which is often \mathbb{R}^d . Often, the bandwidths \mathbf{h}_i are taken to be equal and initially, we shall do so as well. Furthermore, the kernel K is a non-negative function, and is in itself a density function i.e. integrates to 1 as shown below. This ensures that \hat{f} in Equation 1.18 integrates to 1 and is non-negative i.e. a proper distribution.

$$\int_{\mathbb{R}^d} K(\mathbf{x}) d\mathbf{x} = 1$$

In one dimension, a particular useful kernel is the Gaussian kernel given by $K(x) = \phi(x)$ where ϕ is the density function for the standard Normal distribution. ϕ is chosen due to its simple behavior and mathematical properties. In particular, we shall see in the following section, that the properties of the

Gaussian kernel allows for simple expressions when correcting for a boundary such that computation is quick and efficient. For multiple dimensions, we often consider product kernels, which are kernels K of the form

$$K(\mathbf{x}) = \prod_{i=1}^d K_i(x_i) \quad (1.19)$$

I.e. just a product of kernels. In particular, we choose $K_i = \phi$ again due to the numerical properties. Thus, initially, we have a KDE \hat{f} of the following form where we once again note that $\mathbf{h}_i = \mathbf{h}$ for all $i \in \{1, \dots, n\}$ such that h_j denotes the bandwidth associated with the j th dimension.

$$\hat{f}(\mathbf{x}) = \frac{1}{n \prod_{j=1}^d h_j} \sum_{i=1}^n \prod_{j=1}^d \phi\left(\frac{x_j - x_{i,j}}{h_j}\right)$$

The choice of bandwidth \mathbf{h} is important regarding a trade-off between the variance and bias of the KDE. In general, we want to choose h as small as possible resulting in the least bias but a too small \mathbf{h} will result in large variance of the estimator. In particular, \mathbf{h} acts as a smoothing parameter like the number of bins from the previous section, but here, we can choose any $h > 0$ making the KDE a much more versatile tool. Often the *Mean Integrated Square Error* (MISE) is used which is the expected L^2 -norm of $\hat{f} - f$ i.e.

$$\text{MISE}(\hat{f}) = \mathbb{E}_f \left[\int_{\mathbb{R}^d} \left| \hat{f}(\mathbf{x}) - f(\mathbf{x}) \right|^2 d\mathbf{x} \right]$$

which of course depends on \mathbf{h} . The expectation \mathbb{E}_f denotes the expectation with respect to the samples $\{\mathbf{x}_i\}$ of \mathbf{X} with (true) density distribution function f . Expanding the above expression, we obtain a simple expression relating MISE to the integrated squared bias and integrated variance as shown below

$$\text{MISE}(\hat{f}) = \int_{\mathbb{R}^d} \left| \mathbb{E}_f [\hat{f}(\mathbf{x})] - f(\mathbf{x}) \right|^2 d\mathbf{x} + \int_{\mathbb{R}^d} \text{Var} [\hat{f}(\mathbf{x})] d\mathbf{x}$$

It is however quite complicated to optimize the above, and we shall thus often use a simple rule of thumb known as Scott's rule [?] for choosing \mathbf{h} . Namely, for product kernels, we let the bandwidths of each dimension j equal the following where $\hat{\sigma}_j$ is the standard deviation estimated from the observations of X_j

$$h_j^{\text{Scott}} = \hat{\sigma}_j n^{-1/(d+4)}, \quad j \in \{1, \dots, d\}$$

In Figure 1.6, we have shown a basic example in one dimension with two different manual choices of the bandwidth h and h chosen by Scott's rule. We have used

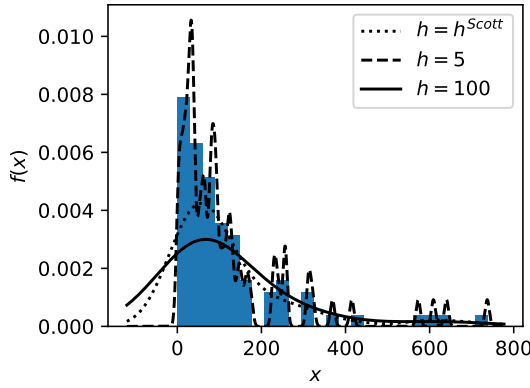


Figure 1.6: The suicide data from MHE. The 86 observations are shown as a histogram of densities along with Gaussian KDEs with bandwidth $h = 5$ and $h = 100$ and h chosen from Scott's rule $h = \hat{\sigma}n^{-1/5} \approx 59.86$.

data from [?], tabulated in [?] which has been used in [?] which propose a method for correcting the KDE near a boundary which we shall discuss in the following section. The data consists of 86 observations regarding suicide and is known to be non-negative. In consideration of the reader, we have included the observations in ?? . It is clear that using h^{Scott} results in what we qualitatively would deem a good estimate for the probability density function as $h = 100$ seen to be overly smoothed whereas $h = 100$ too under-smoothed. In particular, from repeated samples we would expect the estimator using $h = 100$ would have large bias but small variance whereas for $h = 5$ would have much larger variance but smaller bias. However, a problem the estimators, $h = 100$ and $h = h^{Scott}$ especially, have is that they have probability mass below 0 which in this case is unwanted. I.e. when restricting \hat{f} to $[0, \infty)$ they are no longer proper probability distribution functions as they do not integrate to 1. A Simple fix could be to simply rescale \hat{f} such that this is the case, but as seen from the example in Figure 1.7 where this method is applied to the same example as above, this tends to underestimate the peaks especially near the boundary.

We note that using a non-constant h would improve on this behavior, but simpler methods exists, and we thus proceed in the next section with a method that shows great promise regarding this seemingly fundamental issue with KDE. In particular, we refer to a systematic way of letting the shape of each of the kernels depend on the associated observation x_i .

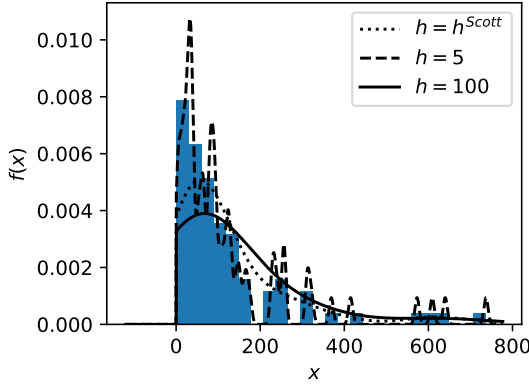


Figure 1.7: Using a rescaled version of \hat{f} on the interval $[0, \infty)$ and disregarding any probability mass below $x = 0$ we obtain proper probability distributions once again. However, neither of the methods capture the peak near the boundary $x = 0$. In particular, although h^{Scott} still seem to be a good choice for h , the KDE does not capture the tendency observed in the data.

1.4.4 Boundary Corrected KDE

Before introducing the boundary corrected kernels presented by [?], we mention another simple method of boundary correction called reflection. Namely, suppose without loss of generality $x = 0$ is the lower boundary of the domain of \mathbf{X} and let \hat{f} be KDE as from the previous section. Then, the reflection boundary corrected KDE denoted \hat{f}_R is defined as

$$\hat{f}_R(x) = \hat{f}(x) + \hat{f}(-x)$$

Clearly, \hat{f}_R is non-negative, and it follows from the below that it is also a proper density function as the probability mass is 1

$$\int_0^\infty \hat{f}_R(x) dx = \int_0^\infty \hat{f}(x) + \hat{f}(-x) dx = \int_0^\infty \hat{f}(x) dx + \int_{-\infty}^0 \hat{f}(x) dx = 1$$

Also, the above is easily extended to two boundaries. Namely, if the domain is $[a, b]$, the reflection boundary corrected KDE is given by

$$\hat{f}_R(x) = \hat{f}(x) + \hat{f}(2a - x) + \hat{f}(2b - x), \quad x \in [a, b]$$

If we once again apply this to the suicide data, comparing to Figure 1.7 we see a big improvement near the boundary as shown in Figure 1.8(a). However, we still proceed with the method originally presented in [?]. The reason for

this is that when testing for distribution type using the Kolmogorov Smirnov test (on a 5% significance level) we reject that the observations originate from \hat{f}_R with $h = 100$. For $h = 5$ we do not but due to the above considerations

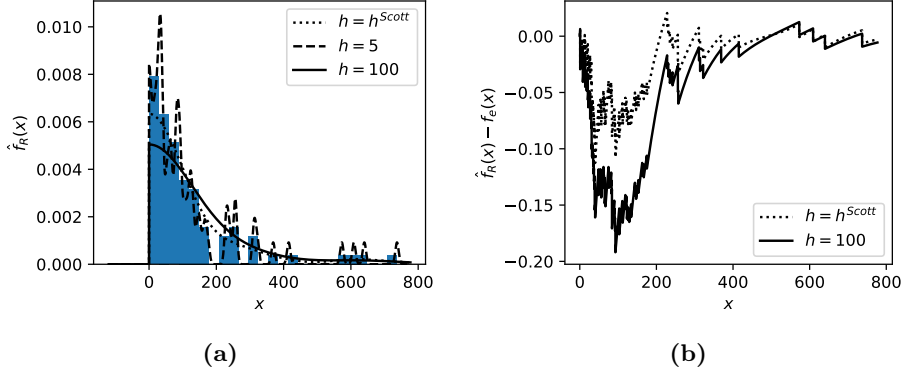


Figure 1.8

regarding the integrated point wise variance of the estimator, this choice of h is undesired in either case. For $h = h^{Scott}$ we do not reject the distribution but as the shape of the error resembles the error for $h = 100$ we suspect that there is some systematic error which we also see from Figure 1.8(b). Furthermore, the largest deviation from the empirical distribution to \hat{f}_R is close to the boundary (as expected). The test statistics (the largest absolute difference D between the distribution functions and the adjusted statistic) is shown in Table 1.1 where the adjusted test statistic should be compared to the critical value 1.358 on a 5% significance level. Furthermore, to really test the kernel density estimators, one should compute the MISE based on bootstrap and/or cross-validation. In particular, this way larger h i.e. more smoothing would be more favorable compared to the Kolmogorov Smirnov test results which shows small h best represent the empirical distribution.

h	D	Adjusted D
5	0.029042	0.27315
h^{Scott}	0.11262	1.0593
100	0.19194	1.8053

Table 1.1: $(\sqrt{n} + 0.12 + 0.11/\sqrt{n}) D$

We now turn our attention to the boundary corrected KDE from [?]. They shortly described this in [?] where it was used to estimate the mutual information in terms of Copula entropy. However, issues arise when using this KDE in terms

of non-negativity of the KDE and the facts that it does not integrate to unity. All of these issues can however be handled in a general way without effecting the results regarding bias of the estimator as we shall also see from the above example on suicide data. In particular, [?] show that the bias of their estimator \hat{f}_L is of order $\mathcal{O}(h^2)$ whereas the reflection method discussed above has $\mathcal{O}(h)$ bias. The basic idea of \hat{f}_L is that it is a linear combination of a symmetric kernel K and xK i.e. a first order kernel. They do however only give explicit expressions when a lower bound at $x = 0$ is enforced but the math generalize nicely to 2 boundaries. In particular, we shall let x_u denote the upper bound of the domain and keep $x = 0$ as the lower bound. Before continuing with the derivation and results regarding implementation for the Gaussian kernel we note that in \mathbb{R}^2 , there is a library `evmix` which implements the boundary corrected KDE from [?] but only for 1 dimension and only with a lower bound at $x = 0$. We thus expand on this library (although in `Python`) to include both lower and upper bounds and furthermore, generalize to multiple dimensions using a product kernel as noted in Equation 1.19.

To expand on the boundary corrected KDE to include an upper bound on the domain, we define the functions $a_m(x)$. Note that in [?], they define a_m as a function of $p = x/h$ where h is the bandwidth, but to keep the expression later on easier to understand in terms of kernel centers etc. we instead define them as a function of $x \in [0, x_u]$. The reason for initially defining a_m in terms of h is to keep the bandwidth out of the expression, which we then define as follows, when there is no upper bound

$$a_m(x) = \int_{-\infty}^{\frac{x}{h}} u^m K(u) du, \quad x \in [0, \infty)$$

The above is actually equal to the part of the m th moment of the kernel centered at x (and width bandwidth h), that is inside the interval $[0, \infty)$ up to a difference in sign. In particular, using the change of variables $z = x - hu$ we have that

$$a_m(x) = (-1)^m \int_0^\infty \left(\frac{z-x}{h} \right)^m \frac{1}{h} K\left(\frac{z-x}{h} \right) dz$$

where we have used that K assumed to be symmetric. Indeed, the above is as described the part of the m th moment that comes from $[0, \infty)$ (up to a difference in sign) of the kernel that is centered at x with bandwidth h . From this, it is a natural extension to replace the upper bound of the integral with x_u which, when expanding on the initial definition in [?] as done in e.g. [?] turns out to be the correct adjustment. Thus, a_m can be understood as part of the moments which we shall then use as normalizing functions. Thus, for an upper bound x_u

²<https://search.r-project.org/CRAN/refmans/evmix/html/bckden.html>

we find that instead a_m is defined as

$$a_m(x) = \int_{\frac{x-x_u}{h}}^{\frac{x}{h}} u^m K(u) du, \quad x \in [0, x_u]$$

The boundary adjusted KDE which we shall denote K^L and index depending on the i th kernel center is then given by

$$K_i^L(x) = \frac{1}{h} \frac{a_2(x) - a_1(x) \frac{x-x_i}{h}}{a_0(x) a_2(x) - a_1^2(x)} K\left(\frac{x-x_i}{h}\right) \quad (1.20)$$

For the Gaussian kernel, $a_m(x)$ for $m \in \{0, 1, 2\}$ are easily computed and implemented in code through standard routines as they have simple closed forms in terms of ϕ and Φ , i.e. the standard Gaussian density and distribution functions. Namely, for a_0 we simply have that by definition of Φ

$$\begin{aligned} a_0(x) &= \int_{\frac{x-x_u}{h}}^{\frac{x}{h}} \phi(u) du \\ &= \Phi\left(\frac{x}{h}\right) - \Phi\left(\frac{x-x_u}{h}\right) \end{aligned}$$

And similarly, for a_1 , using that $\int u \phi(u) du = -\phi(u) + C$

$$\begin{aligned} a_1(x) &= \int_{\frac{x-x_u}{h}}^{\frac{x}{h}} u \phi(u) du \\ &= \phi\left(\frac{x-x_u}{h}\right) - \phi\left(\frac{x}{h}\right) \end{aligned}$$

Finally, for a_2 we have, using integration by parts for the first step

$$\begin{aligned} a_2(x) &= \int_{\frac{x-x_u}{h}}^{\frac{x}{h}} u^2 \phi(u) du \\ &= [-u\phi(u)]_{\frac{x-x_u}{h}}^{\frac{x}{h}} + \int_{\frac{x-x_u}{h}}^{\frac{x}{h}} \phi(u) du \\ &= \left(\frac{x-x_u}{h} \phi\left(\frac{x-x_u}{h}\right) - \frac{x}{h} \phi\left(\frac{x}{h}\right)\right) + \left(\Phi\left(\frac{x}{h}\right) - \Phi\left(\frac{x-x_u}{h}\right)\right) \\ &= \frac{x}{h} a_1(x) - \frac{x_u}{h} \phi\left(\frac{x-x_u}{h}\right) + a_0(x) \end{aligned}$$

From the improved bias of $\mathcal{O}(h^2)$ for K_i^L , we have then obtained an estimator that should perform better in terms of a smaller MISE as we have less bias and by choosing h optimally, we should expect small variance as well. Figure 1.9(a) shows the KDE \hat{f}_L using K^L as the kernel (with $x_u = \infty$) and is compared to

the reflected kernel from above using the same bandwidth. We note that \hat{f}_L has been rescaled such that it integrates to unity on $[0, \infty)$. Also, from Figure 1.9(b) we see that the density is not statistically significant, and although we can not conclude from this alone, we observe that the deviation from the empirical distribution functions is less than that of the reflected KDE.

As noted above, we need to rescale \hat{f}_L to unity as the kernels in Equation 1.20 does not have unit integrals. Furthermore, \hat{f}_L is not even ensured to be non-negative. This can be handled in multiple ways. One way is to simply take the maximum of \hat{f}_L and 0 effectively cutting off any negative density however in [?], a method is given for KDE based on K_i^L specifically, ensuring the $\mathcal{O}(h^2)$ bias of the estimator. Namely, let K_i^N be given as

$$K_i^N(x) = \frac{1}{h a_0(x)} K\left(\frac{x - x_i}{h}\right)$$

which is then the kernel, locally renormalized using a_0 . We then define the related KDE as

$$\hat{f}_N(x) = \frac{1}{n} \sum_{i=1}^n K_i^N(x)$$

which is then non-negative everywhere as $a_0(x), K(x) > 0$. The non-negative boundary corrected KDE, denoted $\hat{f}_P(x)$ is then defined as

$$\hat{f}_P(x) = \hat{f}_N(x) e^{\frac{\hat{f}_L(x)}{\hat{f}_N(x)} - 1}$$

This \hat{f}_P is also shown in Figure 1.9 resulting in identical density functions. We note that \hat{f}_P works by multiplying the non-negative \hat{f}_N by a (large positive) constant when \hat{f}_L is larger than \hat{f}_N and thus drives \hat{f}_N towards \hat{f}_L . However, from implementation and trying on different distributions, sometimes we observe some odd behavior that can easily be derived from the definition of \hat{f}_P . Thus, we propose a modification to overcome these odd properties of \hat{f}_P . Namely, suppose both \hat{f}_L and \hat{f}_N are close to 0 at some point x . If \hat{f}_L is a magnitude of 10 larger than \hat{f}_N , then \hat{f}_P is approximately $8000 \cdot \hat{f}_N$ which even if \hat{f}_N is small may be large. Thus, it is possible even if both are close to 0, that $\hat{f}_P \gg 0$ which is in contrast to what we would want from the estimator. Thus, we propose a regularized KDE version of \hat{f}_P denoted \hat{f}_{regP} which is obtained from first rewriting \hat{f}_P as follows

$$\hat{f}_P(x) = \bar{f}(x) e^{\frac{\hat{f}_L(x) - \hat{f}_N(x)}{\hat{f}_N(x)}}$$

then, we introduce the regularizing parameters $\lambda \geq 0$ such that

$$\hat{f}_{regP}(x) = \bar{f}(x) e^{\frac{\hat{f}_L(x) - \hat{f}_N(x)}{\hat{f}_N(x) + \lambda}}$$

In practice, we have found that $\lambda = 0.001$ is a sufficient regularization whilst preserving the shape. A small λ is preferred as then we are ensured $\mathcal{O}(h^2)$ behavior of the bias. In Figure 1.9 we have also shown this regularized version with $\lambda = 0.001$ and observe that it is basically identical to \hat{f}_P on the domain while also behaving well numerically for large x , where the issue discussed above arise for \hat{f}_P . Before continuing with a few more interesting methods of density

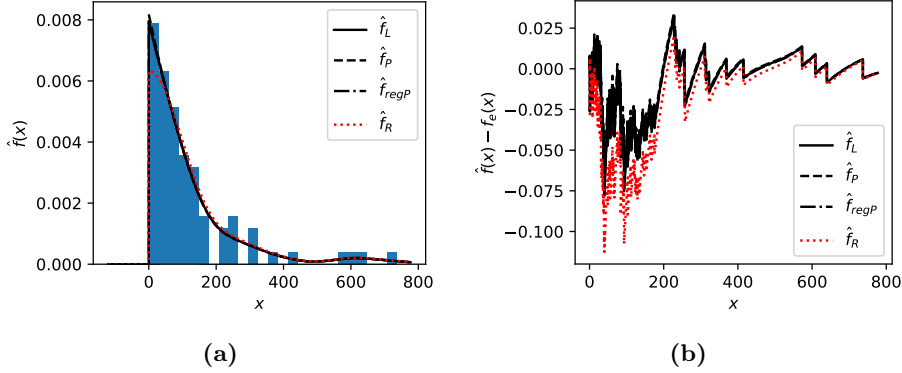


Figure 1.9

estimation and mutual information estimation in particular, we note that in the same R package there is a KDE based on the Gaussian Copula density which at first glance might seem a good choice especially if we are trying to estimate the density of a Copula that we know to be Gaussian. The KDE is based on [?], and trivially we get that the domain of the basis functions is already $[0, 1]$. In particular, the kernel for a given bandwidth h and kernel center x_i is given by

$$K_i^{GC}(x) = \frac{1}{h\sqrt{2-h^2}} e^{-\frac{(1-h^2)}{2h^2(2-h^2)}((\Phi^{-1}(x))^2 + (\Phi^{-1}(x_i))^2) - 2\Phi^{-1}(x)\phi^{-1}(x_i)}$$

which is obtained by letting $h^2 = 1 - \rho$ in the original expression for the Gaussian Copula density. However, in practice, we shall see that this choice of kernel has some numerical instabilities especially for large correlations i.e. small h and does not perform as well as \hat{f}_{regP} even when using their optimal bandwidth.

Furthermore, we note that as observations are initially transformed such that they are approximately uniform on $[0, 1]$, the variance of these uniform observations $u_i = \hat{F}(x_i)$ will also be approximately $1/12$ and hence the Scott rule of thumb will give a near constant bandwidth. However, as we shall see later, the performance is drastically improved by choosing smaller bandwidth for subdomains with high density. We thus note that a local bandwidth should be considered in the future. This could e.g. by using k -nearest neighbor and the

distances of these which is the basic idea of [?]. Although we will not be using this method as it is out of scope for this paper, we note that their estimator seem to result in good estimates when $n \geq 300$. Using the core ideas of this paper, one could perhaps deduce a better algorithm for choosing h globally or locally.

Finally, we note that a recent method based on diffusion [?] [?] which shows great potential. Without going into too much detail, it works by considering the following PDE

$$\frac{\partial}{\partial t} \hat{f}(x; t) = \frac{1}{2} \frac{\partial^2}{\partial x^2} \hat{f}(x; t)$$

where $t = h^2$. The Gaussian kernel is then the unique solution when the domain is \mathbb{R} with condition $\hat{f}[x; 0] = \frac{1}{n} \sum_{i=1}^n \delta(x - x_i)$ i.e. the observations is the boundary condition at $h = 0$. If instead we add that the support should be e.g. $[0, 1]$ the Neumann boundary conditions on the boundary results in an analytical solution. Namely, from the boundary conditions

$$\left. \frac{\partial}{\partial x} \hat{f}(x; t) \right|_{x=0} = \left. \frac{\partial}{\partial x} \hat{f}(x; t) \right|_{x=1} = 0$$

the analytical solution to \hat{f} is

$$\hat{f}_D(x; t) = \frac{1}{n} \sum_{i=1}^n K_i^D(x)$$

where

$$K_i^D(x) = \frac{1}{h} \sum_{k=-\infty}^{\infty} \phi\left(\frac{x - (2k + x_i)}{h}\right) + \phi\left(\frac{x - (2k - x_i)}{h}\right)$$

which is similar to the reflection method from above. They even give an algorithm for calculating t which from the example below appear to work very well. Before presenting this example we note that although it is out of scope to use this method in the following chapter, it would be interesting to see how well it would perform as it is the only KDE here that does not use any transformation near the boundary while still being consistent at the boundary as they show in [?].

Example 1.2 (Diffusion based KDE). *Let $X \sim \text{Beta}(1, 4)$ such that $f(x) = 4(1 - x)^3$ for $x \in [0, 1]$ and thus $F(x) = 1 - (1 - x)^4$. Generating 1000 samples from the distribution and using the algorithm in [?] to estimate the bandwidth, we get a bandwidth $h^D = 0.05461$ whereas with Scott's rule we get $h^{\text{Scott}} = 0.04145$. From Figure 1.10 we see that diffusion and the regularized boundary corrected KDE \hat{f}_{regP} from above agree almost everywhere. Only near the lower*

boundary at $x = 0$ there is a significant difference and here \hat{f}_{regP} seem to fit the true distribution better. This is due to the Neumann boundary condition which enforce \hat{f}_D to have horizontal derivative at both boundaries. Note that we have used both the h^D and h^{Scott} bandwidth for the regularized boundary corrected KDE resulting in basically no difference.

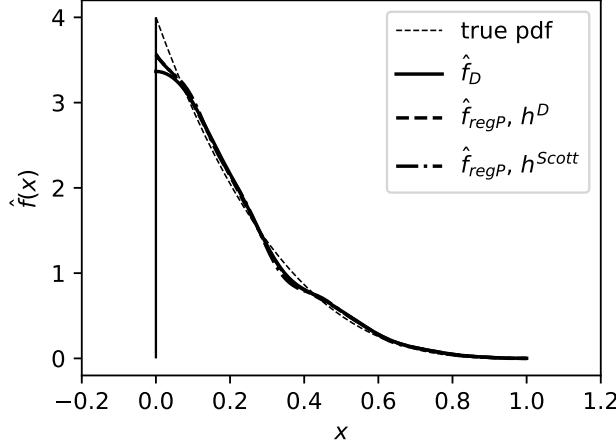


Figure 1.10: Using the diffusion based KDE on 1000 samples from a Beta (1, 4) distribution, we see that in general it fits very well. Only at the boundary $x = 0$ there seems to be a problem which on the other hand, \hat{f}_{regP} does not seem to suffer from to the same extent. In particular, repeating the experiment, \hat{f}_{regP} is on average 4 at $x = 0$ while \hat{f}_D constantly seem to undershoot.

At this point, we thus have a complete set of methods for both estimating the mutual information from observations and from these, algorithms to estimate the (causal) structure depending on assumptions regarding direction i.e. a topological ordering of the variables. We thus proceed with numerical results in the following chapter and apply the framework on both fully controlled systems and the observations from ??.