
A Critical View of the Structural Causal Model

Tomer Galanti¹ Ofir Nabati¹ Lior Wolf^{1,2}

Abstract

In the univariate case, we show that by comparing the individual complexities of univariate cause and effect, one can identify the cause and the effect, without considering their interaction at all. In our framework, complexities are captured by the reconstruction error of an autoencoder that operates on the quantiles of the distribution. Comparing the reconstruction errors of the two autoencoders, one for each variable, is shown to perform surprisingly well on the accepted causality directionality benchmarks. Hence, the decision as to which of the two is the cause and which is the effect may not be based on causality but on complexity.

In the multivariate case, where one can ensure that the complexities of the cause and effect are balanced, we propose a new adversarial training method that mimics the disentangled structure of the causal model. We prove that in the multidimensional case, such modeling is likely to fit the data only in the direction of causality. Furthermore, a uniqueness result shows that the learned model is able to identify the underlying causal and residual (noise) components. Our multidimensional method outperforms the literature methods on both synthetic and real world datasets.

1. Introduction

A long standing debate in the causality literature, is whether causality can be inferred without intervention (Pearl, 2009; Spirtes et al., 2000). The Structural Causal Model (SCM) (Spirtes et al., 2000) is a simple causative model for which many results demonstrate the possibility of such inference (Stegle et al., 2010; Bloebaum et al., 2018; Goudet et al., 2018; Lopez-Paz et al., 2017; 2015). In this model, the effect (Y) is a function of the cause (X) and some independent random noise (E).

In this work, we take a critical perspective of the univariate SCM. We demonstrate empirically that for the univariate case, which is the dominant case in the existing literature, the SCM leads to an effect that has a lower complexity than the cause. Therefore, one can identify the cause and the effect, by measuring their individual complexities, with no need to make the inference based on both variables simultaneously. Thus, the decision as to which of the two is the cause and which is the effect may not be based on causality but on complexity.

Since we are dealing with unordered univariate random variables, the complexity measure has to be based on the probability distribution function. As we show empirically, comparing the entropies of the distribution of two random variables is ineffective for inferring the causal direction. We, therefore, consider the quantiles, i.e. fixed sized vectors that are obtained as sub-sequences of the sorted sampled values of the variable.

We consider suitable complexity scores for these vectors. In our analysis, we show that the reconstruction error of an autoencoder of a multivariate random variable is a valid complexity measure. In addition, we link the reconstruction error based complexity, in the case of variational autoencoders, to the differential entropy of the input random variable. Hence, by computing the reconstruction errors of trained autoencoders on these vectors, we estimate the entropies of the quantile vectors of X and Y .

The challenges of measuring causality independently of complexity in the 1D case lead us to consider the multidimensional case, where the complexity can be controlled by, for example, manipulating the dimension of the noise signal in the SCM. Note that unlike (Goudet et al., 2018), we consider pairs of multivariate vectors and not many univariate variables in a graph structure. We demonstrate that for the multidimensional case, any method that is based on comparing the complexity of the individual random variables X and Y fails to infer causality of random variables. Furthermore, we extend a related univariate result by (Zhang & Hyvrinen, 2010) to the multidimensional case and prove that an SCM is unlikely to hold in both directions $X \rightarrow Y$ and $Y \rightarrow X$.

Based on our observations, we propose a new causality inference method for multidimensional cause and effect. The algorithm learns three networks in a way that mimics

¹School of Computer Science, Tel Aviv University, Israel

²Facebook AI Research (FAIR). Correspondence to: Tomer Galanti <tomerga2@post.tau.ac.il>, Lior Wolf <wolf@fb.com>.

the parts of the SCM. The noise part is unknown and is replaced by a function that is constrained to be independent of the cause, as captured by an adversarial loss. However, we show empirically that even without the explicit constraint, in several cases, such an independence emerges.

Our empirical results support our analysis and demonstrate that in the univariate case, assigning cause and effect based on complexity is competitive with the state of the art methods. In the multidimensional case, we show that the proposed method outperforms existing multivariate methods, as well as new extensions of univariate literature methods.

1.1. Problem Setup

We investigate the problem of causal inference from observational data. A non-linear structural causal model (SCM for short) is a generative process of the following form:

$$\begin{aligned} X &\sim \mathbb{P}_X \\ E &\sim \mathbb{P}_E \\ Y &\leftarrow g(f(X), E) \end{aligned} \quad (1)$$

The functions $g : \mathbb{R}^{d_f+d_e} \rightarrow \mathbb{R}^{d_y}$ and $f : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^{d_f}$ are fixed and unknown. In general, g and f are non-linear. Here, X is the input random variable and E is the environment random variable that is independent of X . We say that $X \in \mathbb{R}^{d_x}$ causes $Y \in \mathbb{R}^{d_y}$ if they satisfy a generative process, such as Eq. 1.

We present methods for inferring whether X causes Y (denoted by $X \rightarrow Y$) or Y causes X , or neither. The algorithm is provided with i.i.d samples $\{(x_i, y_i)\}_{i=1}^m \sim \mathbb{P}_{X,Y}^m$ (the distribution of m i.i.d samples from the joint distribution $\mathbb{P}_{X,Y}$) from the generative process of Eq. 1. In general, by (cf. Prop 4.8, (Peters et al., 2017)), for any joint distribution $\mathbb{P}_{X,Y}$ of two random variables X and Y , there is an SCM, $Y = g(f(X), E)$, where E is a noise variable, such that, $X \perp\!\!\!\perp E$ and f, g are some (measurable) functions. Therefore, in general, deciding whether X causes Y or vice versa is ill-posed when only provided with samples from the joint distribution. However, (Zhang & Hyvrinen, 2010) showed for the one dimensional case (i.e., $X, Y \in \mathbb{R}$) that under reasonable conditions, a representation $Y = g(f(X) + E)$ holds only in one direction. In Sec. 3.2, we extend this theorem and show that a representation $Y = g(f(X), E)$ holds only in one direction when g and f are assumed to be neural networks and X, Y are multidimensional (we call such SCMs neural SCMs).

Throughout the paper, we denote by $\mathbb{P}_U[u] := \mathbb{P}[U \leq u]$ the cumulative distribution function of a uni/multi-variate real valued random variable U and \mathbb{P} is a standard Lebesgue measure. Additionally, we denote by $p_U(u) = \frac{d}{du} \mathbb{P}_U[u]$ the probability density function of U (if exists, i.e., $\mathbb{P}_U[u]$ is absolutely continuous). We denote by $\mathbb{E}_{u \sim U}[f(u)]$ the

expected value of $f(u)$ for u that is distributed by $\mathbb{P}_U[u]$. The identity matrix of dimension $n \times n$ is denoted by I_n or I , when the dimension is obvious from the context.

1.2. Related Work

In causal inference, the algorithm is provided with a dataset of matched samples (x, y) of two random variables X and Y and decides whether X causes Y or vice versa. The early wisdom in this area asserted that this asymmetry of the data generating process (i.e., that Y is computed from X and not vice versa) is not apparent from looking at $\mathbb{P}_{X,Y}$ alone. That is, in general, provided with samples from the joint distribution $\mathbb{P}_{X,Y}$ of two variables X, Y does tell us whether it has been induced by an SCM from X to Y or from Y to X .

In publications, such as (Pearl, 2009; Spirtes et al., 2000), it is argued that in order to decide whether X causes Y or vice versa, one needs to observe the influence of interventions on the environment parameter. To avoid employing interventions, most publications assume prior knowledge on the generating process and/or independence between the cause and the mechanism.

Various methods for causal inference under the SCM have been suggested. Many of these methods are based on independence testing, where the algorithm models the data as $Y = g(f(X), E)$ (and vice versa) and decides upon the side that provides a better fitting in terms of mapping accuracy and independence between $f(X)$ and $E = r(X, Y)$. The LiNGAM (Shimizu et al., 2006) algorithm assumes that the SCM takes the form $Y = \beta X + E$, where $X \perp\!\!\!\perp E$, $\beta \in \mathbb{R}$ and E is non-Gaussian. The algorithm learns β , such that, X and $Y - \beta X$ are independent by applying independent component analysis (ICA). The Direct-LiNGAM (Shimizu et al., 2011) extends this method and replaces the mutual information minimization with a non-parametric kernel based loss (Bach & Jordan, 2003). However, the computation of this loss is of order $\Theta(m^2)$ in the the worst case (m is the number of samples).

The ANM approach (Hoyer et al., 2009) extends LiNGAM's modeling and assumes that $Y = f(X) + E$, where $X \perp\!\!\!\perp E$. A Gaussian Process is employed as the learned mechanism between the two random variables. The function f is trained to map between X and Y (and vice versa) and the method then tests whether, X and $f(X) - Y$ are independent. The independence test is based on kernels (Gretton et al., 2005).

A different extension of LiNGAM is the PNL algorithm by (Zhang & Hyvrinen, 2010). This algorithm learns a mapping between X and Y (and vice versa) of the form $Y = g(f(X) + E)$, where $f(X)$ and E are restricted to be independent. To do so, PNL trains two neural networks f and g to minimize the mutual information between $f(X)$

and $E = g^{-1}(Y) - f(X)$. The main disadvantage of this method is the reliance on the minimization of the mutual information. It is often hard to measure and optimize the mutual information directly, especially in higher dimensions. In many cases, it requires having an explicit modeling of the density functions, because of the computation of expected log-probability within the formulation of the entropy measure.

In our multivariate method, we take a similar approach to the above methods. However, our GAN-based independence constraint is non-parametric, is applied on the observations rather on an explicit modeling of the density functions, and the method is computationally efficient. In addition, we do not assume restrictive structural assumptions and treat the generic case, where the effect is of the form $Y = g(f(X), E)$.

Another independence constraint is applied by the Information Geometric Causal Inference (IGCI) (Danusis et al., 2012) approach, which determines the causal relationship in a deterministic setting $Y = f(X)$ under an independence assumption between the cause X and the mechanism f , $\text{Cov}(\log f'(x), p_X) = 0$.

The Conditional Distribution Similarity Statistic (CDS) (Fonollosa, 2016) measures the standard deviation of the values of Y (resp. X) after binning in the X (resp. Y) direction. The lower the standard deviation, the more likely the pair to be $X \rightarrow Y$. The CURE algorithm (Sgouritsa et al., 2015) compares between $X \rightarrow Y$ and $Y \rightarrow X$ directions in the following manner: if we can estimate $p_{X|Y}$ based on samples from p_Y more accurately than $p_{Y|X}$ based on samples from p_X , then $X \rightarrow Y$ is inferred.

The BivariateFit method learns a Gaussian Process regressor in both directions and decides upon the side that had the lowest error. The RECI method (Bloebaum et al., 2018) trains a regression model (a logistic function, polynomial functions, support vector regression, or a neural networks) in both directions, and returns the side that produced a lower MSE loss. The CGNN algorithm (Goudet et al., 2018) uses the Maximum Mean Discrepancy (MMD) distance between the distribution produced by modeling Y as an effect of X , $(X, g(X, E))$ (and vice versa), and the ground truth distribution. The algorithm compares the two distances and returns the direction that led to a smaller distance. The Gaussian Process Inference model (GPI) (Stegle et al., 2010) builds two generative models, one for $X \rightarrow Y$ and one for $Y \rightarrow X$. The distribution of the candidate cause variable is modelled as a Gaussian Mixture Model, and the mechanism f is a Gaussian Process. The causal direction is determined from the generative model that best fits the data.

Finally, it is worth mentioning that several other methods,

such as (Heinze-Deml et al., 2017; Zhang et al., 2011) assume a different type of SCM, where the algorithm is provided with separate datasets that correspond to different environments, i.e., sampled i.i.d from $\mathbb{P}_{X,Y|E}$, where the value of E is fixed for all samples in the dataset. In these publications, a different independence condition is assumed: Y is independent of E given X . This assumption fails in our setting, since we focus on the vanilla SCM, where the algorithm is provided only with observational i.i.d. samples of X and $Y = g(f(X), E)$ and the samples are not divided into subsets that are invariant w.r.t E .

2. The Univariate Case

In this section, we show that the univariate SCM does not necessarily capture causality. For this purpose, we describe a method for identifying the cause and the effect, which considers each of the two variables independently without considering the mapping between them. The success of this method, despite neglecting any interaction between the variables, indicates that univariate SCM challenges can be solved without considering causality.

The proposed method computes a complexity score for X and, independently, for Y . It then compares the scores and decides that the cause is the random variable with the larger score among them. Capturing the complexity of a univariate random variable without being able to anchor the observations in additional features is challenging. One can observe the probability distribution function and compute, for example, its entropy. As we show empirically, in Sec. 4, this is ineffective.

Our complexity scoring method, therefore, has a few stages. As a first step, it converts the random variable at hand (say, X) into a multivariate random variable. This is done by sorting the samples of the random variable, and then cutting the obtained list into fixed sized vectors of length k . We discard the largest measurements in the case, where the number of samples is not a multiple of k . We denote the random variable obtained this way by U . At the second stage, the method computes the complexity of the obtained random variable U using an autoencoder reconstruction error.

2.1. Reconstruction Errors as Complexity Measures

One can consider the complexity of a multivariate random variable in various ways. We consider non-negative complexity measures $C(X)$, which satisfy the weak assumption that when X and Y are independent then their complexities are lower than the complexity of their concatenation:

$$C(X, Y) \geq \max(C(X), C(Y)). \quad (2)$$

Examples of sample complexity measures that satisfy this condition are the Shannon Entropy and the Kolmogorov Complexity. The following lemma shows that a complexity that is based on autoencoder modeling is also in this family.

Let $\mathcal{F} = \{\mathcal{H}^d\}_{d=1}^\infty$ be a family of classes of autoencoders $A : \mathbb{R}^d \rightarrow \mathbb{R}^d$. Assume that the family \mathcal{F} is closed to fixations, i.e., for any autoencoder $A \in \mathcal{H}^{d_1+d_2}$ and any fixed vector $y^* \in \mathbb{R}^{d_2}$ ($x^* \in \mathbb{R}^{d_1}$), we have: $A(x, y^*)_{1:d_1} \in \mathcal{H}^{d_1}$ ($A(x^*, y)_{d_1+1:d_2} \in \mathcal{H}^{d_2}$). Here, $v_{i:j} = (v_i, \dots, v_j)$. Note that this is the typical situation when considering neural networks with biases.

Let X be a random variable. Let X be a multivariate random variable dimension d . We define the autoencoding complexity of X as follows:

$$C_{\mathcal{F}}(X) := \min_{A^* \in \mathbb{H}^d} \mathbb{E}_{x \sim X} [\ell(A^*(x), x)] \quad (3)$$

where $\ell(a, b)$ is some loss function.

Lemma 1. *Let $\{\mathcal{H}^d\}_{d=1}^\infty$ be a family of classes of autoencoders that is closed to fixations. The function $C_{\mathcal{F}}(X)$ is a proper complexity measure.*

2.2. The AEQ method

The AEQ method we propose estimates and compares the auto-encoder reconstruction error of the quantile vectors of X and Y . It is important to note that it does not imply that the AEQ method compares between the entropies of X and Y .

Once the random variable U is obtained as the quantiles of a random variable (either X or Y), our method trains an autoencoder $A : \mathbb{R}^k \rightarrow \mathbb{R}^k$ on U . A is trained to minimize the following objective:

$$\mathcal{L}_{\text{recon}}(A) := \mathbb{E}_{u \sim U} [\ell(A(u), u)] \quad (4)$$

where $\ell(a, b)$ is some loss function. In our implementation, we employ the L_2 -loss function, defined as $\ell(a, b) = \|a - b\|_2^2$. Finally, the method uses the value of $\mathcal{L}_{\text{recon}}(A)$, which we refer to as the AEQ score, as a proxy for the complexity of X (smaller loss means lower complexity). It decides that X or Y is the cause, based on which side provides a higher AEQ.

As we show in Sec. 4, the proposed causality-free method is as successful at solving SCM challenges as the leading literature methods. However, we do not propose it as a standalone method, and rather develop it to show the shortcoming of the univariate SCM setting and the associated literature datasets.

3. The Multivariate Case

For the univariate case, one can consider the complexity of the X and Y variables of the SCM and infer directionality.

We propose the AEQ complexity for this case, since more conventional complexities are ill-defined for unordered 1D data or, in the case of entropy, found to be ineffective.

The following technical lemma shows that for any complexity measure C , one cannot infer directionality in the multivariate SCM based on C .

Lemma 2. *Let C be a complexity measure of multivariate random variables (i.e, non-negative and satisfies Eq. 2). Then, there are triplets of random variables (X, E, Y) and (\hat{X}, E, Y) and functions g and g' , such that, $Y = g(X, E)$, $Y = g'(\hat{X}, E)$, $C(X) < C(Y)$ and $C(\hat{X}) > C(Y)$. Therefore, C cannot serve as a score for causal inference.*

We now turn our attention to a new multivariate causality inference method.

3.1. An Adversarial Method for Causal Inference

Our causality inference algorithm trains neural networks G, F, R and D . The success of fitting these networks serves as the score for the causality test. The function F models the function f , G models g and $R(Y)$ aims to model the environment parameter E . In general, our method aims at solving the following objective:

$$\begin{aligned} \min_{G, F, R} \mathcal{L}_{\text{err}}(G, F, R) &:= \frac{1}{m} \sum_{i=1}^m \|G(F(a_i), R(b_i)) - b_i\|_2^2 \\ \text{s.t: } &A \perp\!\!\!\perp R(B) \end{aligned} \quad (5)$$

where A is either X or Y and B is the other option, and $a_i = x_i, b_i = y_i$ or $a_i = y_i, b_i = x_i$ accordingly. To decide whether $X \rightarrow Y$ or vice versa, we train a different triplet G, F, R for each direction and see if we can minimize the mapping error \mathcal{L}_{err} subject to independence. We decide upon a specified direction, if the loss can be minimized subject to independence. In general, searching within the space of functions that satisfy $A \perp\!\!\!\perp R(B)$ is an intractable problem. However, we can replace it with a loss term that is minimized when $A \perp\!\!\!\perp R(B)$.

Independence loss We would like $R(B)$ to capture the information encoded in E . Therefore, restrict $R(B)$ and A to be independent in each other. We propose an adversarial loss for this purpose, which is a modified version of a loss proposed by (Brakel & Bengio, 2017) and later analyzed by (Press et al., 2019).

This loss measures the discrepancy between the joint distribution $\mathbb{P}_{A, R(B)}$ and the product of the marginal distributions $\mathbb{P}_A \times \mathbb{P}_{R(B)}$. Let d_F (d_R) be the dimension of F 's output (R). To measure the discrepancy, we make use of a discriminator $D : \mathbb{R}^{d_a+d_R} \rightarrow [0, 1]$ (d_a equals d_x or d_y depending

on $A = X$ or $A = Y$) that minimizes the following term:

$$\begin{aligned} \mathcal{L}_D(D; R) := & \frac{1}{m} \sum_{i=1}^m \ell(D(a_i, R(b_i)), 1) \\ & + \frac{1}{m} \sum_{i=1}^m \ell(D(\hat{a}_i, R(\hat{b}_i)), 0) \end{aligned} \quad (6)$$

where D is a discriminator network, and $l(p, q) = -(q \log(p) + (1 - q) \log(1 - p))$ is the binary cross entropy loss for $p \in [0, 1]$ and $q \in \{0, 1\}$. In addition, $\{(\hat{a}_i, \hat{b}_i)\}_{i=1}^m$ are i.i.d samples from $\mathbb{P}_A \times \mathbb{P}_B$. To create these samples, we sample independently \hat{a}_i and \hat{b}_i from the respective training sets $\{\hat{a}_i\}_{i=1}^m$ and $\{\hat{b}_i\}_{i=1}^m$ and then arbitrarily match them into couples (\hat{a}_i, \hat{b}_i) .

To restrict that $R(B)$ and A are independent, R is trained to confuse the discriminator D such that the two sets of samples are indistinguishable by D ,

$$\begin{aligned} \mathcal{L}_{\text{indep}}(R; D) := & \frac{1}{m} \sum_{i=1}^m \ell(D(a_i, R(b_i)), 1) \\ & + \frac{1}{m} \sum_{i=1}^m \ell(D(\hat{a}_i, R(\hat{b}_i)), 1) \end{aligned} \quad (7)$$

Full objective The full objective of our method is then translated into the following program:

$$\begin{aligned} \min_{G, F, R} \mathcal{L}_{\text{err}}(G, F, R) + \lambda \cdot \mathcal{L}_{\text{indep}}(R; D) \\ \min_D \mathcal{L}_D(D; R) \end{aligned} \quad (8)$$

Where λ is some positive constant. The discriminator D minimizes the loss $\mathcal{L}_D(D; R)$ concurrently with the other networks. Our method decides if X causes Y or vice versa, by comparing the score $\mathcal{L}_{\text{err}}(G, F, R)$. A lower error means a better fit. The full description of the architecture employed for the encoders, generator and discriminator is given in Appendix A. A sensitivity experiment for the parameter λ is provided in Appendix B.

In addition to the success in fitting, we also measure the degree of independence between A and $R(B)$. We denote by c_{real} the percentage of samples (a_i, b_i) that the discriminator classifies as 1 and by c_{fake} the percentage of samples (\hat{a}_i, \hat{b}_i) that are classified as 0. We note that when $c_{\text{real}} \approx 1 - c_{\text{fake}}$, the discriminator is unable to discriminate between the two distributions, i.e., it is wrong in classifying half of the samples. We, therefore, use $|c_{\text{real}} + c_{\text{fake}} - 1|$ as a measure of independence.

3.2. Analysis

In this section, we analyze the proposed method. In Thm. 1, we show that if X and Y admit a SCM in one direction,

then it admits a SCM in the opposite direction, only if the involved functions satisfy a specific partial differential equation.

Theorem 1 (Identifiability of neural SCMs). *Let $\mathbb{P}_{X, Y}$ admit a neural SCM from X to Y as in Eq. 1, such that p_X , and the activation functions of f and g are three-times differentiable. Then it admits a neural SCM from Y to X , only if p_X , f , g satisfy Eq. 27 in the appendix.*

This result generalizes the one-dimensional case presented in (Zhang & Hyvrinen, 2010), where a one-dimensional version of this differential equation is shown to hold in the analog case.

In the following theorem, we show that minimizing the proposed losses is sufficient to recover the different components, i.e., $F(X) \propto f(X)$ and $R(Y) \propto E$, where $A \propto B$ means that $A = f(B)$ for some invertible function f .

Theorem 2 (Uniqueness of Representation). *Let $\mathbb{P}_{X, Y}$ admit a nonlinear model from X to Y as in Eq. 1, i.e., $Y = g(f(X), E)$ for some random variable $E \perp\!\!\!\perp X$. Assume that f and g are invertible. Let G, F and R be functions, such that, $\mathcal{L}_{\text{err}} := \mathbb{E}_{(x, y) \sim (X, Y)} [\|G(F(x), R(y)) - y\|_2^2] = 0$ and G and F are invertible functions and $X \perp\!\!\!\perp R(Y)$. Then, $F(X) \propto f(X)$ and $R(Y) \propto E$.*

where, \mathcal{L}_{err} is the mapping error proposed in Eq. 5. In addition, the assumption $X \perp\!\!\!\perp R(Y)$ is sufficed by the independence loss.

A more general results, but which requires additional terminology, is stated as Thm. 3 in Appendix C. It extends Thm. 2 to the case, where the mapping loss is not necessarily zero and the independence $X \perp\!\!\!\perp R(Y)$ is replaced by a discriminator-based independence measure. Thm. 3 also gets rid of the assumption that the various mappings f, g and F, G are invertible. In this case, instead of showing that $R(Y) \propto E$, we provide an upper bound on the reconstruction of E out of $R(Y)$ (and vice versa) that improves as the training loss of G, F and R decreases.

To conclude our analysis, by Thm. 1, under reasonable assumptions, if X and Y admit a multivariate SCM in direction $X \rightarrow Y$, then, there is no such representation in the other direction. By Thm. 2, by training our method in both directions, one is able to capture the causal model in the correct direction. This is something that is impossible to do in the other direction by Thm. 1.

4. Experiments

This section is divided into two parts. The first is devoted to showing that causal inference in the one-dimensional case highly depends on the complexities of the distributions of X and Y . In the second part of this section, we show that our multivariate causal inference method outperforms existing

baselines. Most of the baseline implementations were taken from the Causality Discovery Toolbox of (Kalainathan & Goudet, 2019). The experiments with PNL (Zhang & Hyvrinen, 2010), LiNGAM (Shimizu et al., 2006) and GPI (Stegle et al., 2010) are based on their original matlab code.

4.1. One-Dimensional Data

We compared the autoencoder method on several well-known one dimensional cause-effect pairs datasets. Each dataset consists of a list of pairs of real valued random variables (X, Y) with their direction 1 or 0, depending on $X \rightarrow Y$ or $Y \rightarrow X$ (resp.). For each pair, we have a dataset of samples $\{(x_i, y_i)\}_{i=1}^m$.

Five cause-effect inference datasets, covering a wide range of associations, are used. CE-Net (Goudet et al., 2018) contains 300 artificial cause-effect pairs generated using random distributions as causes, and neural networks as causal mechanisms. CE-Gauss contains 300 artificial cause-effect pairs as generated by (Mooij et al., 2016), using random mixtures of Gaussians as causes, and Gaussian Process priors as causal mechanisms. CE-Multi (Goudet et al., 2018) contains 300 artificial cause-effect pairs built with random linear and polynomial causal mechanisms. In this dataset, simulated additive or multiplicative noise is applied before or after the causal mechanism.

The real-world datasets include the diabetes dataset by (Frank & Asuncion, 2010), where causality is from Insulin \rightarrow Glucose. Glucose curves and Insulin doses were analysed for 69 patients, each serves as a separate dataset. To match the literature protocols, the pairs are taken in an orderless manner, ignoring the time series aspect of the problem. Finally, the Tübingen cause-effect pairs dataset by (Mooij et al., 2016) is employed. This dataset is a collection of 100 heterogeneous, hand-collected, real-world cause-effect samples.

The autoencoder A employed in our method, Eq. 4, is a fully-connected five-layered neural network with three layers for the encoder and two layers for the decoder. The hyperparameters of this algorithm are the sizes of each layer, the activation function and the input dimension, i.e., length of sorted cuts (denoted by k in Sec. 2). Throughout the experiments, we noticed that the hyperparameter with the highest influence is the input dimension. For all datasets, results are stable in the range of $200 \leq k \leq 300$, and we, therefore, use $k = 250$ throughout the experiments. For all datasets, we employed the ReLU activation function, except the Tübingen dataset, where the sigmoid activation function produced better results (results are also reasonable with ReLU, but not state of the art).

In addition to our method, we also present results obtained with the entropy of each individual variable as a complexity

Table 1. Mean AUC rates of various baselines on different one dimensional cause-effect pairs datasets. Our interaction-less AEQ algorithm achieves competitive results on most datasets.

Method	CE-Net	CE-Gauss	CE-Multi	Tübingen	Diabetes
BivariateFit	77.6	36.3	55.4	58.4	0.0
LiNGAM(Shimizu et al., 2006)	43.7	66.5	59.3	39.7	100.0
CDS (Fonollosa, 2016)	89.5	84.3	37.2	59.8	12.0
IGCI (Danusis et al., 2012)	57.2	33.2	80.7	62.2	100.0
ANM (Hoyer et al., 2009)	85.1	88.9	35.5	53.7	22.2
PNL(Zhang & Hyvrinen, 2010)	75.5	83.0	49.0	68.1	28.1
GPI (Stegle et al., 2010)	88.4	89.1	65.8	66.4	92.9
RECI (Bloebaum et al., 2018)	60.0	64.2	85.3	62.6	95.4
CGNN (Goudet et al., 2018)	89.6	82.9	96.6	79.8	34.1
Entropy as complexity	49.6	49.7	50.8	54.5	53.4
Our AEQ comparison	62.5	71.0	96.0	82.8	95.0

measure. This is done by binning the values of the variables into 50 bins. Other numbers of bins produce similar results.

Tab. 1 presents the mean AUC for each literature benchmark. As can be seen, the AEQ complexity measure produces reasonable results in comparison to the state of the art methods, indicating that the 1D SCM can be overcome by comparing per-variable scores. On the popular Tübingen dataset, the AEQ computation outperforms all literature methods.

Tab. 2 presents accuracy rates for various methods on the Tübingen dataset, where such results are often reported in the literature. As can be seen, our interaction-less method outperforms almost all other methods, including methods that employ supervised learning of the cause-effect relation.

4.2. Multivariate Data

We first compare our method on several synthetic datasets. Each dataset consists of a list of pairs of real multivariate random variables (X, Y) with their direction 1 or 0, depending on $X \rightarrow Y$ or $Y \rightarrow X$ (resp.). For each pair, we have a dataset of samples $\{(x_i, y_i)\}_{i=1}^m$.

We employ five datasets, covering multiple associations. Each dataset contains 300 artificial cause-effect pairs. The cause random variable is of the form $X = h(z)$, where h is some function and $z \sim \mathcal{N}(0, \sigma_1^2 \cdot I_n)$. The effect is of the form $Y = g(u(X, E))$, where $E \sim \mathcal{N}(0, \sigma_2^2 \cdot I_n)$ is independent of X , u is a fixed function that combined the cause X and the noise term E and g is the causal mechanism. For each dataset, the functions h and g are taken from the same family of causal mechanisms \mathcal{H} . Each pair of random variables is specified by randomly selected functions h and

Table 2. Accuracy rates of various baselines on the CE-Tüb dataset. Our interaction-less algorithm AEQ achieves almost SOTA accuracy.

Method	Supervised	Acc
LiNGAM (Shimizu et al., 2006)	-	44.3%
BivariateFit	-	44.9%
Entropy as a complexity measure	-	52.5%
IGCI (Daniusis et al., 2012)	-	62.6%
CDS (Fonollosa, 2016)	-	65.5%
ANM (Hoyer et al., 2009)	-	59.5%
CURE (Sgouritsa et al., 2015)	-	60.0% ^a
GPI (Stegle et al., 2010)	-	62.6%
PNL (Zhang & Hyvrinen, 2010)	-	66.2%
CGNN (Goudet et al., 2018)	-	74.4%
RECI (Bloebaum et al., 2018)	-	77.5%
SLOPE (Marx & Vreeken, 2017)	-	81.0%
Our AEQ comparison	-	80.0%
Jarfo (Fonollosa, 2016)	+	59.5%
RCC (Lopez-Paz et al., 2015)	+	75.0% ^b
NCC (Lopez-Paz et al., 2017)	+	79.0%

^aThe accuracy of CURE is reported on version 0.8 of the dataset in (Sgouritsa et al., 2015) as 75%. In (Bloebaum et al., 2018) they re-ran this algorithm and achieved an accuracy rate of around 60%.

^bThe accuracy scores reported in (Lopez-Paz et al., 2015) are for version 0.8 of the dataset, in (Lopez-Paz et al., 2017) they re-ran RCC (Lopez-Paz et al., 2015) on version 1.0 of the dataset.

g.

The synthetic datasets extend the standard synthetic data generators of (Kalainathan & Goudet, 2019) to the multivariate causal pairs. MCE-Poly is generated element-wise polynomials composed on linear transformations as mechanisms and $u(X, E) = X + E$. MCE-Net pairs are generated using neural networks as causal mechanisms and u is the concatenation operator. The mechanism in MCE-SigMix consists of linear transformation followed by element wise application of $q_{a,b,c}(x) := ab(\tilde{x} + c)/(1 + |b \cdot (\tilde{x} + c)|)$, where a, b, c are random real valued numbers, which are sampled for each pair and $\tilde{x} = x + e$, where e is the environment random variable. In this case, $u(X, E) = X + E$. We noticed that a-priori, the produced datasets are imbalanced in a way that the reconstruction error of a standard autoencoder on each random variable can be employed as a score that predicts the cause variable with a high accuracy. Therefore, in order to create balanced datasets, we varied the amount of noise dimensions and their intensity, until the autoencoder reconstruction error of both X and Y became similar. Note that for these multivariate variables, we do not use quantiles and use the variables themselves. As the AutoEncoder reconstruction results in Tab. 3 show, in the MCE-SigMix dataset, balancing was only partly successful.

Table 3. Mean AUC rates of various baselines on different multivariate cause-effect pairs datasets. The datasets are designed and balanced, such that an autoencoder method would fail. Our method achieves SOTA results.

Method	MCE- Poly	MCE- Net	MCE- SigMix	MOUS- MEG
AE reconstruction	57.2	42.4	22.3	41.2
BivariateFit	54.7	48.4	48.2	44.2
IGCI (Daniusis et al., 2012)	41.9	49.3	59.8	56.0
CDS (Fonollosa, 2016)	63.8	57.0	62.1	89.9
ANM (Hoyer et al., 2009)	52.2	51.1	46.4	52.4
PNL (Zhang & Hyvrinen, 2010)	76.4	54.7	16.8	56.3
CGNN (Goudet et al., 2018)	47.8	67.8	58.8	40.9
Our method	95.3	84.2	98.5	97.7

Table 4. Results of various methods on different variations of the MOUS-MEG dataset. R stands for the MEG scan at rest, W stands for the word presented to the subject and A stands for the MEG scan, when the subject is active.

Method	R + W → A	R → A	W → A
Expected to be causal	Yes	No	No
AE reconstruction	41.2	51.7	98.6
BivariateFit	44.2	58.1	0.0
IGCI (Daniusis et al., 2012)	56.0	50.6	42.2
CDS (Fonollosa, 2016)	89.9	52.1	90.2
ANM (Hoyer et al., 2009)	52.4	49.3	0.0
PNL (Zhang & Hyvrinen, 2010)	56.3	43.7	0.0
CGNN (Goudet et al., 2018)	40.9	52.2	100.0
Our method	97.7	44.4	0.0

We compare our results to two types of baseline methods: (i) BivariateFit and ANM (Hoyer et al., 2009) are methods that were designed (also) for the multivariate case, (ii) CGNN (Goudet et al., 2018) and PNL (Zhang & Hyvrinen, 2010) are naturally extended to this case. To extend the CDS (Fonollosa, 2016) and IGCI (Daniusis et al., 2012) methods to higher dimension, we applied quantizations over the data samples, i.e., cluster the samples $\{x_i\}_{i=1}^m$ and $\{y_i\}_{i=1}^m$ using two distinct k-means with $k = 10$, and then, each sample is replaced with its corresponding cluster to obtain a univariate representation of the data. After pre-processing the data, we apply the corresponding method. To select the hyperparameter k , we varied its value between 5 to 500 for different scales and found 10 to provide the best results. RECI (Bloebaum et al., 2018) could be extended. However, RECI’s runtime is of order $\mathcal{O}(n^3)$, where n is

Table 5. Emergence of independence. Ind C (Ind E) is the mean of $|c_{\text{real}} + c_{\text{fake}} - 1|$ over all pairs of random variables, epochs and samples, when training the method from X to Y (vice versa). w/o backprop means without backpropagating gradients from D to R .

Dataset	Full method			w/o backprop		
	AUC	Ind C	Ind E	AUC	Ind C	Ind E
MCE-Poly	95.3	0.06	0.05	95.1	0.10	0.10
MCE-Net	84.2	0.28	0.31	65.1	0.55	0.55
MCE-SigMix	98.5	0.05	0.06	98.8	0.16	0.20
MOUS-MEG	97.7	0.14	0.14	80.7	0.74	0.75

the input dimension. Other methods cannot be extended, or require significant modifications. For example, the SLOPE method (Marx & Vreeken, 2017) heavily relies on the ability to order the samples of the random variables X and Y . However, it is impossible to do so in the multivariate case. We could not find any open source implementation of the CURE algorithm (Sgouritsa et al., 2015).

The results, given in Tab. 3 show a clear advantage over the literature methods across the four datasets. The exact same architecture is used throughout all experiments, with the same λ parameter. See Sec. 1 of the supplementary material. A sensitivity analysis (see supplementary Sec. 2) shows that our results are better than all baseline methods, regardless of the parameter λ .

In addition to the synthetic datasets, we also employ the MOUS-MEG real world dataset, provided to us by the authors of (King et al., 2020). This dataset is part of Mother Of Unification Studies (MOUS) dataset (Schoffelen et al., 2019). This dataset contains magneto-encephalography (MEG) recordings of 102 healthy Dutch-speaking subjects performing a reading task (9 of them were excluded due to corrupted data). Each subject was asked to read 120 sentences in Dutch, both in the right order and randomly mixed order, which adds up to a total of over 1000 words. Each word was presented on the computer screen for 351ms on average and was separated from the next word by 3-4 seconds. Each time step consists of 301 MEG readings of the magnetometers, attached to different parts of the head. For more information see (Schoffelen et al., 2019). For each pair (X, Y) , X is the interval $[-1.5s, -0.5s]$ relative to the word onset concatenated with the word embedding (using the spaCy python module with the Dutch language model), this presents the subject in his “rest” state (i.e. the cause). Y is the interval $[0, 1.0s]$ relative to the word onset, which presents the subject in his “active” state (i.e. the effect).

To validate the soundness of the dataset, we ran a few experiments on variations of the dataset and report the results as additional columns in Tab. 4. As can be seen, a dataset where the cause consists of the word embedding and the

effect consists of the subject’s “active” state is highly imbalanced. This is reasonable, since the word embedding and the MEG readings are encoded differently and are of different dimensions. In addition, when the cause is selected to be the “rest” state and the effect is the “active” state, the various algorithms are unable to infer which side is the cause and which one is the effect, since the word is missing. Finally, when considering the Rest+Word \rightarrow Active variation, the relationship is expected to be causal, the AE reconstruction indicates that the dataset is balanced, and our method is the only one to achieve a high AUC rate.

Emergence of independence To check the importance of our adversarial loss in identifying the direction of causality and capturing the implicit independent representation $f(X)$ and E , we applied our method without training R against the discriminator. Therefore, in this case, the discriminator only serves as a test whether X and $R(Y)$ are independent or not and does not contribute to the training loss of R ($\lambda = 0$).

As mentioned in Sec. 3.1, the distance between $c_{\text{real}} + c_{\text{fake}}$ to 1 indicates the amount of dependence between X and $R(Y)$. We denote by Ind C the mean values of $|c_{\text{real}} + c_{\text{fake}} - 1|$ over all pairs of random variables and samples when training our method in the causal direction. The same mean score when training in the anti-causal direction is denoted Ind E. As is evident from Tab. 5, the independence is similar between the two directions, emphasizing the importance of the reconstruction error in the score.

As can be seen in Tab. 5, the adversarial loss improves the results when there is no implicit emergence of independence. However, in cases where there is emergence of independence, the results are similar. We noticed that the values of Ind C and Ind E are smaller for the full method. However, in MCE-Poly and MCE-SigMix they are still very small and, therefore, there is implicit emergence of independence between X and $R(Y)$, even without explicitly training $R(Y)$ to be independent of X .

5. Summary

We discover an imbalance in the complexities of cause and effect in the univariate SCM and suggest a method to exploit it. Since the method does not consider the interactions between the variables, its success in predicting cause and effect indicates an inherent bias in the univariate datasets. Turning our attention to the multivariate case, where the complexity can be actively balanced, we propose a new method in which the learned networks model the underlying SCM itself. Since the noise term E is unknown, we replace it by a function of Y that is enforced to be independent of X . We also show that under reasonable conditions, the independence emerges, even without explicitly enforcing it.

6. Acknowledgements

This project has received funding from the European Research Council (ERC) under the European Unions Horizon 2020 research and innovation programme (grant ERC CoG 725974). The authors would like to thank Dimitry Shaiderman for insightful discussions. The contribution of Tomer Galanti is part of Ph.D. thesis research conducted at Tel Aviv University.

References

- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017*, pp. 214–223, 2017.
- Bach, F. R. and Jordan, M. I. Kernel independent component analysis. *J. Mach. Learn. Res.*, 3:1–48, March 2003. ISSN 1532-4435. doi: 10.1162/153244303768966085. URL <https://doi.org/10.1162/153244303768966085>.
- Bloebaum, P., Janzing, D., Washio, T., Shimizu, S., and Schoelkopf, B. Cause-effect inference by comparing regression errors. In Storkey, A. and Perez-Cruz, F. (eds.), *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pp. 900–909, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR.
- Brakel, P. and Bengio, Y. Learning independent features with adversarial nets for non-linear ica, 2017.
- Chazelle, B. *The Discrepancy Method: Randomness and Complexity*. Cambridge University Press, USA, 2000. ISBN 0521003571.
- Daniusis, P., Janzing, D., Mooij, J. M., Zscheischler, J., Steudel, B., Zhang, K., and Schölkopf, B. Inferring deterministic causal relations. *CoRR*, 2012.
- Fonollosa, J. A. R. Conditional distribution variability measures for causality detection. *ArXiv*, abs/1601.06680, 2016.
- Frank, A. and Asuncion, A. UCI machine learning repository, 2010. <http://archive.ics.uci.edu/ml>.
- Goudet, O., Kalainathan, D., Caillou, P., Lopez-Paz, D., Guyon, I., and Sebag, M. Learning functional causal models with generative neural networks. In *Explainable and Interpretable Models in Computer Vision and Machine Learning*, Springer Series on Challenges in Machine Learning. Springer International Publishing, 2018.
- Gretton, A., Herbrich, R., Smola, A., Bousquet, O., and Schölkopf, B. Kernel methods for measuring independence. *J. Mach. Learn. Res.*, 6:2075–2129, December 2005. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1046920.1194914>.
- Heinze-Deml, C., Peters, J., and Meinshausen, N. Invariant causal prediction for nonlinear models. *Journal of Causal Inference*, 6, 2017.

- Hoyer, P. O., Janzing, D., Mooij, J. M., Peters, J., and Schölkopf, B. Nonlinear causal discovery with additive noise models. In Koller, D., Schuurmans, D., Bengio, Y., and Bottou, L. (eds.), *Advances in Neural Information Processing Systems 21*, pp. 689–696. Curran Associates, Inc., 2009.
- Kalainathan, D. and Goudet, O. Causal discovery toolbox: Uncover causal relationships in python, 2019.
- King, J.-R., Charton, F., Oquab, M., and Lopez-Paz, D. Measuring causal influence with back-to-back regression: the linear case, 2020. URL <https://openreview.net/forum?id=B1lKd1HtwS>.
- Lopez-Paz, D., Muandet, K., Schölkopf, B., and Tolstikhin, I. Towards a learning theory of cause-effect inference. In Bach, F. and Blei, D. (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 1452–1461, Lille, France, 07–09 Jul 2015. PMLR.
- Lopez-Paz, D., Nishihara, R., Chintala, S., Schölkopf, B., and Bottou, L. Discovering causal signals in images. In *Proceedings IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 2017*, pp. 58–66, Piscataway, NJ, USA, July 2017. IEEE.
- Mansour, Y., Mohri, M., and Rostamizadeh, A. Domain adaptation: Learning bounds and algorithms. In *COLT*, 2009.
- Marx, A. and Vreeken, J. Telling cause from effect using mdl-based local and global regression. In *2017 IEEE International Conference on Data Mining (ICDM)*, pp. 307–316, Nov 2017.
- Mooij, J. M., Peters, J., Janzing, D., Zscheischler, J., and Schölkopf, B. Distinguishing cause from effect using observational data: Methods and benchmarks. *Journal of Machine Learning Research*, 17(32):1–102, 2016.
- Müller, A. Integral probability metrics and their generating classes of functions advances in applied probability. In *Advances in Applied Probability*, pp. 429–443, 1997.
- Pearl, J. *Causality: Models, Reasoning and Inference*. Cambridge University Press, New York, NY, USA, 2nd edition, 2009. ISBN 052189560X, 9780521895606.
- Peters, J., Janzing, D., and Schölkopf, B. *Elements of Causal Inference - Foundations and Learning Algorithms*. Adaptive Computation and Machine Learning Series. The MIT Press, Cambridge, MA, USA, 2017.
- Press, O., Galanti, T., Benaim, S., and Wolf, L. Emerging disentanglement in auto-encoder based unsupervised image content transfer. In *International Conference on Learning Representations*, 2019.
- Schoffelen, J.-M., Oostenveld, R., Lam, N. H., Uddén, J., Hultén, A., and Hagoort, P. A 204-subject multimodal neuroimaging dataset to study language processing. *Scientific data*, 6(1):17, 2019.
- Sgouritsa, E., Janzing, D., Hennig, P., and Schölkopf, B. Inference of Cause and Effect with Unsupervised Inverse Regression. In Lebanon, G. and Vishwanathan, S. V. N. (eds.), *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, volume 38 of *Proceedings of Machine Learning Research*, pp. 847–855, San Diego, California, USA, 09–12 May 2015. PMLR.
- Shimizu, S., Hoyer, P. O., Hyvärinen, A., and Kerminen, A. A linear non-gaussian acyclic model for causal discovery. *J. Mach. Learn. Res.*, 7:2003–2030, December 2006. ISSN 1532-4435.
- Shimizu, S., Inazumi, T., Sogawa, Y., Hyvärinen, A., Kawahara, Y., Washio, T., Hoyer, P. O., and Bollen, K. Directlingam: A direct method for learning a linear non-gaussian structural equation model. *J. Mach. Learn. Res.*, 12:1225–1248, July 2011. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1953048.2021040>.
- Silvester, J. R. Determinants of block matrices. *The Mathematical Gazette*, pp. 2000, 1999.
- Spirtes, P., Glymour, C., and Scheines, R. *Causation, Prediction, and Search*. MIT press, 2nd edition, 2000.
- Stegle, O., Janzing, D., Zhang, K., Mooij, J. M., and Schölkopf, B. Probabilistic latent variable models for distinguishing between cause and effect. In Lafferty, J. D., Williams, C. K. I., Shawe-Taylor, J., Zemel, R. S., and Culotta, A. (eds.), *Advances in Neural Information Processing Systems 23*, pp. 1687–1695. Curran Associates, Inc., 2010.
- Zhang, K. and Hyvrinen, A. Distinguishing causes from effects using nonlinear acyclic causal models. In Guyon, I., Janzing, D., and Schölkopf, B. (eds.), *Proceedings of Workshop on Causality: Objectives and Assessment at NIPS 2008*, volume 6 of *Proceedings of Machine Learning Research*, pp. 157–164, Whistler, Canada, 12 Dec 2010. PMLR.
- Zhang, K., Peters, J., Janzing, D., and Schölkopf, B. Kernel-based conditional independence test and application in causal discovery. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, UAI’11, pp. 804–813, Arlington, Virginia, United States, 2011. AUAI Press. ISBN 978-0-9749039-7-2. URL <http://dl.acm.org/citation.cfm?id=3020548.3020641>.

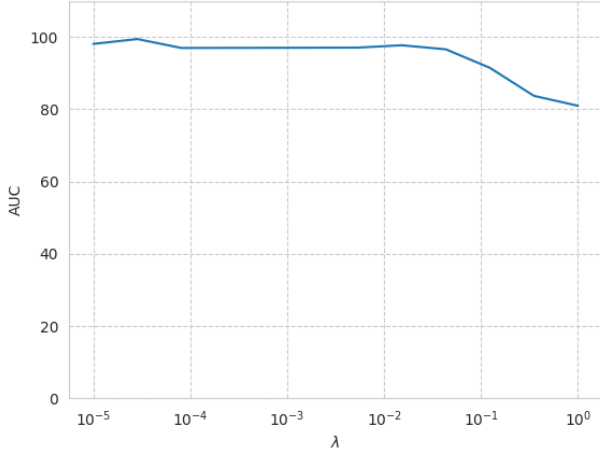


Figure 1. Sensitivity experiment. The graph presents the AUC of our algorithm on MOUS-MEG dataset with λ , which varies between 10^{-5} to 1 in a logarithmic scale.

A. Architecture for All Multivariate Experiments

The functions G , F , R and D in the adversarial multivariate method are fully connected neural networks and their architectures are as follows: F is a 2-layered network with dimensions $100 \rightarrow 60 \rightarrow 50$, R is a 3-layered network with dimensions $100 \rightarrow 50 \rightarrow 50 \rightarrow 20$, G is a 2-layers neural network with dimensions $50 + 20 \rightarrow 80 \rightarrow 100$ (the input has 50 dimensions for $F(X)$ and 20 for $R(Y)$). The discriminator is a 3-layers network with dimensions $100 + 20 \rightarrow 60 \rightarrow 50 \rightarrow 2$ (the input is the concatenation of X and $R(Y)$). The activation function in all networks is the sigmoid function except the discriminator that applies the leaky ReLU activation. For all networks, the activation is not applied at the output layer.

Throughout the experiments the learning rate for training G , F and R is 0.01 and the learning rate of D is 0.001.

B. Sensitivity Experiment

To check that our results are robust with respect to λ , we conducted a sensitivity analysis. In this experiment we ran our algorithm on the MOUS-MEG dataset (i.e., Rest + Word \rightarrow Active variation) with λ that varies between 10^{-5} to 1 in a logarithmic scale. As can be seen in Fig. 1, our algorithm is highly stable to the selection of $\lambda \in [10^{-5}, 10^{-1}]$. The performance decays (gradually) only for $\lambda \geq 0.1$.

C. Analysis

C.1. Terminology and Notations

We recall some relevant notations and terminology. For a vector $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ we denote $\|x\|_2 := \sqrt{\sum_{i=1}^n x_i^2}$ the Euclidean norm of x . For a differentiable function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ and $x \in \mathbb{R}^m$, we denote by

$$J(f(x)) := \left(\frac{\partial f_i}{\partial \zeta_j}(x) \right)_{i \in [n], j \in [m]} \quad (9)$$

the Jacobian matrix of f in x . For a twice differentiable function $f : \mathbb{R}^m \rightarrow \mathbb{R}$, we denote by

$$H(f(x)) := \left(\frac{\partial^2 f}{\partial \zeta_i \partial \zeta_j}(x) \right)_{i, j \in [m]} \quad (10)$$

the Hessian matrix of f in x . Additionally, for a twice differentiable function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$, $f(x) = (f_1(x), \dots, f_n(x))$, we denote the Hessian of f by $H(f(x)) := (H(f_1(x)), \dots, H(f_n(x)))$. For a scalar function $f : \mathbb{R}^m \rightarrow \mathbb{R}$ instead of using the Jacobian notation, the gradient notation will be employed, $\nabla(f(x)) := J(f(x))$. For two positive functions $f(x)$ and $g(x)$, we denote, $f(x) \lesssim g(x)$ if there is a constant $C > 0$, such that, $f(x) \leq C \cdot g(x)$.

C.2. Proofs for the Results

In this section we provide the proofs of the main results in the paper.

Lemma 1. *Let $\{\mathcal{H}^d\}_{d=1}^\infty$ be a family of classes of autoencoders that is closed to fixations. The function $C_{\mathcal{F}}(X)$ is a proper complexity measure.*

Proof. First, since $\ell(a, b) \geq 0$ for all $a, b \in \mathbb{R}^k$, this function is non-negative. Next, we would like to show that $C_{\mathcal{F}}(X, Y) \geq \max(C_{\mathcal{F}}(X), C_{\mathcal{F}}(Y))$. Let A^* be the minimizer of $\mathbb{E}_{x \sim X} [\ell(A^*(x), x)]$ within $\mathcal{H}^{d_1+d_2}$. We consider that there is a vector y^* , such that,

$$\begin{aligned} \mathbb{E}_{(x,y) \sim (X,Y)} [\ell(A(x, y), (x, y))] & \\ & \geq \mathbb{E}_{y \sim Y} \mathbb{E}_{x \sim X} [\ell(A(x, y), (x, y))] \\ & \geq \mathbb{E}_{x \sim X} [\ell(A(x, y^*), (x, y^*))] \\ & \geq \mathbb{E}_{x \sim X} [\ell(A(x, y^*)_{1:d_1}, x)] \end{aligned} \quad (11)$$

We note that $A(x, y^*)_{1:d_1} \in \mathcal{H}^{d_1}$. Therefore,

$$\begin{aligned} \mathbb{E}_{(x,y) \sim (X,Y)} [\ell(A(x, y), (x, y))] & \\ & \geq \min_{A^* \in \mathcal{H}^{d_1}} \mathbb{E}_{x \sim X} [\ell(A^*(x), x)] = C_{\mathcal{F}}(X) \end{aligned} \quad (12)$$

By similar considerations, $C_{\mathcal{F}}(X, Y)$. \square

Lemma 2. *Let C be a complexity measure of multivariate random variables (i.e, non-negative and satisfies Eq. 2).*

Then, there are triplets of random variables (X, E, Y) and (\hat{X}, E, Y) and functions g and g' , such that $Y = g(X, E)$, $Y = g'(\hat{X}, E)$, $C(X) < C(Y)$ and $C(\hat{X}) > C(Y)$. Therefore, C cannot serve as a score for causal inference.

Proof. Let X be a random variable and $E \perp\!\!\!\perp X$, such that $Y = g(X, E)$. Assume that $C(X) < C(Y)$. Then, let X' be a random variable independent of X , such that, $C(X') > C(Y)$. Then, according to the definition of a complexity measure, we have: $C(X, X') > C(Y)$ and we have: $Y = g'(X, X', E)$, for $g'(a, b, c) = g(a, c)$. \square

The following lemma is an extension of Thm. 1 in (Zhang & Hyvrinen, 2010) to real valued random variables of dimension > 1 .

Lemma 3. Assume that (X, Y) can be described by both:

$$Y = g_1(f_1(X) + E_1), \text{ s.t. } X \perp\!\!\!\perp E_1 \text{ and } g_1 \text{ is invertible} \quad (13)$$

and

$$X = g_2(f_2(Y) + E_2), \text{ s.t. } Y \perp\!\!\!\perp E_2 \text{ and } g_2 \text{ is invertible} \quad (14)$$

Assume that g_1 and g_2 are invertible and let:

$$\begin{aligned} T_1 &:= g_1^{-1}(Y) \text{ and } h_1 := f_2 \circ g_1 \\ T_2 &:= g_2^{-1}(X) \text{ and } h_2 := f_1 \circ g_2 \end{aligned} \quad (15)$$

Assume that the involved densities p_{T_2} , p_{E_1} and nonlinear functions f_1, g_1 and f_2, g_2 are third order differentiable. We then have the following equations for all (X, Y) satisfying:

$$\begin{aligned} & \mathbf{H}(\eta_1(t_2)) \cdot \mathbf{J}(h_1(t_1)) - \mathbf{H}(\eta_2(e_1)) \cdot \mathbf{J}(h_2(t_2)) \\ & + \mathbf{H}(\eta_2(e_1)) \cdot \mathbf{J}(h_2(t_2)) \cdot \mathbf{J}(h_1(t_1)) \cdot \mathbf{J}(h_2(t_2)) \\ & - \nabla(\eta_2(e_1)) \cdot \mathbf{H}(h_2(t_2)) \cdot \mathbf{J}(h_1(t_1)) = 0 \end{aligned} \quad (16)$$

where $\eta_1(t_2) := \log p_{T_2}(t_2)$ and $\eta_2(e_1) := \log p_{E_1}(e_1)$.

Proof. The proof is an extension of the proof of Thm. 1 in (Zhang & Hyvrinen, 2010). We define:

$$\begin{aligned} T_1 &:= g_1^{-1}(Y) \text{ and } h_1 := f_2 \circ g_1 \\ T_2 &:= g_2^{-1}(X) \text{ and } h_2 := f_1 \circ g_2 \end{aligned} \quad (17)$$

Since g_2 is invertible, the independence between X and E_1 is equivalent to the independence between T_2 and E_1 . Similarly, the independence between Y and E_2 is equivalent to the independence between T_1 and E_2 . Consider the transformation $F : (E_2, T_1) \mapsto (E_1, T_2)$:

$$\begin{aligned} E_1 &= T_1 - f_1(X) = T_1 - f_1(g_2(T_2)) \\ T_2 &= f_2(Y) + E_2 = f_2(g_1(T_1)) + E_2 \end{aligned} \quad (18)$$

The Jacobian matrix of this transformation is given by:

$$\begin{aligned} \mathbf{J} &:= \mathbf{J}(F(e_2, t_1)) \\ &= \left[\begin{array}{c|c} -\mathbf{J}(h_2(t_2)) & I - \mathbf{J}(h_2(t_2)) \cdot \mathbf{J}(h_1(t_1)) \\ \hline I & \mathbf{J}(h_1(t_1)) \end{array} \right] \quad (19) \end{aligned}$$

Since I commutes with any matrix, by Thm. 3 in (Silvester, 1999), we have:

$$\begin{aligned} & \left| \det(\mathbf{J}(F(E_2, T_1))) \right| \\ &= \left| \det \left(-\mathbf{J}(h_2(T_2)) \cdot \mathbf{J}(h_1(T_1)) \right. \right. \\ & \quad \left. \left. - I \cdot (I - \mathbf{J}(h_2(T_2)) \cdot \mathbf{J}(h_1(T_1))) \right) \right| = 1 \end{aligned} \quad (20)$$

Therefore, we have: $p_{T_2}(t_2) \cdot p_{E_1}(e_1) = p_{T_1, E_2}(t_1, e_2) / |\det \mathbf{J}| = p_{T_1, E_2}(t_1, e_2)$. Hence, $\log(p_{T_1, E_2}(t_1, e_2)) = \eta_1(t_2) + \eta_2(e_1)$ and we have:

$$\frac{\partial \log(p_{T_1, E_2}(t_1, e_2))}{\partial e_2} = \nabla \eta_1(t_2) - \nabla \eta_2(e_1) \cdot \mathbf{J}(h_2(t_2)) \quad (21)$$

Therefore,

$$\begin{aligned} & \frac{\partial^2 \log(p_{T_1, E_2}(t_1, e_2))}{\partial e_2 \partial t_1} \\ &= \mathbf{H}(\eta_1(t_2)) \cdot \mathbf{J}(h_1(t_1)) \\ & \quad - \mathbf{H}(\eta_2(e_1)) \cdot (I - \mathbf{J}(h_2(t_2)) \cdot \mathbf{J}(h_1(t_1))) \cdot \mathbf{J}(h_2(t_2)) \\ & \quad - \nabla(\eta_2(e_1)) \cdot \mathbf{H}(h_2(t_2)) \cdot \mathbf{J}(h_1(t_1)) \\ &= \mathbf{H}(\eta_1(t_2)) \cdot \mathbf{J}(h_1(t_1)) - \mathbf{H}(\eta_2(e_1)) \cdot \mathbf{J}(h_2(t_2)) \\ & \quad + \mathbf{H}(\eta_2(e_1)) \cdot \mathbf{J}(h_2(t_2)) \cdot \mathbf{J}(h_1(t_1)) \cdot \mathbf{J}(h_2(t_2)) \\ & \quad - \nabla(\eta_2(e_1)) \cdot \mathbf{H}(h_2(t_2)) \cdot \mathbf{J}(h_1(t_1)) \end{aligned} \quad (22)$$

The independence between T_1 and E_2 implies that for every possible (t_1, e_2) , we have: $\frac{\partial^2 \log p_{T_1, E_2}(t_1, e_2)}{\partial e_2 \partial t_1} = 0$. \square

Lemma 4 (Reduction to post-linear models). Let $f(x) = \sigma_1(W_d \dots \sigma_1(W_1 x))$ and $g(u, v) = \sigma_2(U_k \dots \sigma_2(U_1(u, v)))$ be two neural networks. Then, if $Y = g(f(X), E)$ for some $E \perp\!\!\!\perp X$, we can represent $Y = \hat{g}(\hat{f}(X) + N)$ for some $N \perp\!\!\!\perp X$.

Proof. Let $f(x) = \sigma_1(W_d \dots \sigma_1(W_1 x))$ and $g(u, v) = \sigma_2(U_k \dots \sigma_2(U_1(u, v)))$ be two neural networks. Here, (u, v) is the concatenation of the vectors u and v . We consider that $U_1(f(X), E) = U_1^1 f(X) + U_1^2 E$. We define a noise variable $N := U_1^2 E$ and have: $X \perp\!\!\!\perp N$. In addition, let $\hat{f}(x) := U_1^1 f(x)$ and $\hat{g}(z) := \sigma_2(U_k \dots \sigma_2(U_2 \sigma_2(z)))$. We consider that: $Y = \hat{g}(\hat{f}(X) + N)$ as desired. \square

Theorem 1 (Identifiability of neural SCMs). Let $\mathbb{P}_{X, Y}$ admit a neural SCM from X to Y as in Eq. 1, such that p_X ,

and the activation functions of f and g are three-times differentiable. Then it admits a neural SCM from Y to X , only if p_X , f , g satisfy Eq. 27 in the appendix.

Proof. Let $f_i(z) = \sigma_1(W_{i,d} \dots \sigma_1(W_{i,1}z))$ and $g_i(u, v) = \sigma_2(U_{i,k} \dots \sigma_2(U_{i,1}(u, v)))$ (where $i = 1, 2$) be pairs of neural networks, such that, σ_1 and σ_2 are three-times differentiable. Assume that:

$$Y = g(f(X), E_1) \text{ and } X = g(f(Y), E_2) \quad (23)$$

for some $E_1 \perp\!\!\!\perp X$ and $E_2 \perp\!\!\!\perp Y$. By Lem. 4, we can represent

$$\begin{aligned} Y &= \hat{g}_1(\hat{f}_1(X) + N_1), \\ \text{where } N_1 &= U_{1,1}^2 E_1, \hat{f}_1 = U_{1,1}^1 f_1(X) \\ \text{and } \hat{g}_1(z) &= \sigma_2(U_{1,k} \dots \sigma_2(U_{1,2} \sigma_2(z))) \end{aligned} \quad (24)$$

and also,

$$\begin{aligned} X &= \hat{g}_2(\hat{f}_2(Y) + N_2), \\ \text{where } N_2 &= U_{2,1}^2 E_2, \hat{f}_2 = U_{2,1}^1 f_2(X) \\ \text{and } \hat{g}_2(z) &= \sigma_2(U_{2,k} \dots \sigma_2(U_{2,2} \sigma_2(z))) \end{aligned} \quad (25)$$

Here, for each $i = 1, 2$ and $j = 1, 2$, $U_{i,1}^j$ are the submatrices of $U_{i,1}$ that satisfy:

$$U_{i,1}(f_i(X), E_i) = U_{i,1}^1 f_i(X) + U_{i,1}^2 E_i \quad (26)$$

From the proof of Lem. 4, it is evident that the constructed \hat{g}_1 , \hat{f}_1 and \hat{g}_2 , \hat{f}_2 are three-times differentiable whenever σ_1 and σ_2 are. Therefore, by Lem. 3, the following differential equation holds:

$$\begin{aligned} & \mathbf{H}(\eta_1(t_2)) \cdot \mathbf{J}(h_1(t_1)) - \mathbf{H}(\eta_2(n_1)) \cdot \mathbf{J}(h_2(t_2)) \\ & + \mathbf{H}(\eta_2(n_1)) \cdot \mathbf{J}(h_2(t_2)) \cdot \mathbf{J}(h_1(t_1)) \cdot \mathbf{J}(h_2(t_2)) \\ & - \nabla(\eta_2(n_1)) \cdot \mathbf{H}(h_2(t_2)) \cdot \mathbf{J}(h_1(t_1)) = 0 \end{aligned} \quad (27)$$

where

$$\begin{aligned} T_1 &:= \hat{g}_1^{-1}(Y) \text{ and } h_1 := \hat{f}_2 \circ \hat{g}_1 \\ T_2 &:= \hat{g}_2^{-1}(X) \text{ and } h_2 := \hat{f}_1 \circ \hat{g}_2 \end{aligned} \quad (28)$$

and $\eta_1(t_2) := \log p_{T_2}(t_2)$ and $\eta_2(n_1) := \log p_{N_1}(n_1)$. \square

Theorem 2 (Uniqueness of Representation). *Let $\mathbb{P}_{X,Y}$ admit a nonlinear model from X to Y as in Eq. 1, i.e., $Y = g(f(X), E)$ for some random variable $E \perp\!\!\!\perp X$. Assume that f and g are invertible. Let G , F and R be functions, such that, $\mathcal{L}_{\text{err}} := \mathbb{E}_{(x,y) \sim (X,Y)} [\|G(F(x), R(y)) - y\|_2^2] = 0$ and G and F are invertible functions and $X \perp\!\!\!\perp R(Y)$. Then, $F(X) \propto f(X)$ and $R(Y) \propto E$.*

Proof. Since F and f are invertible, one can represent: $F(X) = F(f^{-1}(f(X)))$ and $f(X) = f(F^{-1}(F(X)))$. Similarly, since G and g are invertible, we also have: $(F(X), R(Y)) \propto (f(X), E)$. Since $(F(X), R(Y)) \propto (f(X), E)$ and $F(X) \propto f(X)$, we have: $R(Y) = Q(F(X), E)$. However, $R(Y) \perp\!\!\!\perp F(X)$ and therefore, we can represent $R(Y) = P(E)$ and vice versa. \square

C.3. An Extension of Thm. 2

In this section we extend Thm. 2. As a reminder, in our method, we employ two losses: a mapping loss $\mathcal{L}_{\text{err}}(G, F, R)$ and a GAN-like independence loss $\mathcal{L}_{\text{indep}}(R; D)$.

Informally, in similar fashion to Thm. 2, we would like to claim that when the algorithm successfully minimizes the losses, the information present in $r(Y) := E$ can be recovered from $R(Y)$. In Thm. 2, it is shown that whenever the losses are optimal, we have: $R(Y) \propto r(Y)$. In Thm. 3, we relax the optimality assumption and we would like to express the recoverability of $r(Y)$ given $R(Y)$ in terms of the success of the algorithm in minimizing the losses. By similar arguments we can also show that $f(X)$ can be recovered from $F(X)$.

To define a measure of recoverability of one random variable given another random variable we consider a class \mathcal{T} of transformations $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$. The reconstruction of a given random variable V out of U is given by:

$$\text{Rec}_{\mathcal{T}}(V; U) := \inf_{T \in \mathcal{T}} \mathbb{E}_{(u,v) \sim (U,V)} [\|T(u) - v\|_2^2] \quad (29)$$

The class \mathcal{T} of transformations serves as the set of possible candidate mappings from U to V .

In our case, we are interested in measuring the ability to recover the information present in $r(Y)$ given $R(Y)$. Therefore, we would like to show that our algorithm implicitly minimizes:

$$\text{Rec}_{\mathcal{T}}(r(y); R(y)) = \inf_{T \in \mathcal{T}} \mathbb{E}_{y \sim Y} [\|T(R(y)) - r(y)\|_2^2] \quad (30)$$

To do so, we upper bound the recoverability using the mapping error and a discriminator based divergence. In our bound, instead of employing $\mathcal{L}_{\text{indep}}(R; D)$ directly, we make use of a different discriminator based measure of independence. For simplicity, we will assume that \mathcal{T} consists of functions $T : \cup_{n \in \mathbb{N}} \mathbb{R}^n \rightarrow \mathbb{R}^{d_e}$ and for every fixed $u \in \mathbb{R}^k$, we have: $T_u(x) := T(x, u) \in \mathcal{T}$. This is the case of $\mathcal{T} = \cup_{n \in \mathbb{N}} \mathcal{T}_n$, where \mathcal{T}_n is the class of fully-connected neural networks (with biases) with input dimension n and fixed hidden dimensions.

The proposed measure of independence will be based on the discrepancy measure (Chazelle, 2000; Mansour et al., 2009). For a given class \mathcal{D} of discriminator functions $D : \mathcal{X} \rightarrow \mathbb{R}$, we define the \mathcal{D} -discrepancy, also known as Integral Probability Metric (Müller, 1997), between two random variables X_1 and X_2 over \mathcal{X} by:

$$\text{disc}_{\mathcal{D}}[X_1 \| X_2] := \sup_{D \in \mathcal{D}} \{ \mathbb{E}_{x_1 \sim X_1} [D(x_1)] - \mathbb{E}_{x_2 \sim X_2} [D(x_2)] \} \quad (31)$$

A well known example of this measure is the WGAN divergence (Arjovsky et al., 2017) that is specified by a class \mathcal{D} of neural networks of Lipschitzness ≤ 1 .

In our bound, to measure the independence between $F(X)$ and $R(Y)$, we make use of the term:

$$\text{disc}_{\mathcal{D}}[(F(X), R(Y), Y) \parallel (F(X'), R(Y), Y)] \quad (32)$$

for some class of discriminators \mathcal{D} . Even though we do not use the original measure of independence, the idea is very similar. Instead of using a GAN-like divergence between $(X, R(Y))$ and $(X', R(Y))$, we employ a WGAN-like divergence between $(F(X), R(Y))$ and $(F(X'), R(Y))$. From a theoretical standpoint, it is easier to work with the discrepancy measure since it resembles a distance measure.

The selection of \mathcal{D} is a technical by-product of the proof of the theorem and one can treat it as an ‘‘expressive enough’’ class of functions. Specifically, each discriminator $D \in \mathcal{D}$ takes the following form:

$$D(u_1, u_2, u_3) = \|T(u_1, u_2) - Q(u_3)\|_2^2 \quad (33)$$

where $T \in \mathcal{T}$ and $Q \in \mathcal{Q}$. Here, $u_1 \in \mathbb{R}^{d_f}$, $u_2 \in \mathbb{R}^{d_e}$ and $u_3 \in \mathbb{R}^{d_y}$. In particular, the discrepancy measure is:

$$\begin{aligned} & \text{disc}_{\mathcal{D}}[(F(X), R(Y), Y) \parallel (F(X'), R(Y), Y)] \\ &= \sup_{T \in \mathcal{T}, Q \in \mathcal{Q}} \left\{ \mathbb{E}_{(x,y)} [\|T(F(x), R(y)) - Q(y)\|_2^2] \right. \\ & \quad \left. - \mathbb{E}_{(x',y)} [\|T(F(x'), R(y)) - Q(y)\|_2^2] \right\} \quad (34) \end{aligned}$$

where $(x, y) \sim (X, Y)$ and $x' \sim X$ is an independent copy of x . A small discrepancy indicates that there is no discriminator $D \in \mathcal{D}$ that is able to separate between $(F(X), R(Y), Y)$ and $(F(X'), R(Y), Y)$. In particular, if $F(X) \perp\!\!\!\perp R(Y)$, then, $\text{disc}_{\mathcal{D}}[(F(X), R(Y), Y) \parallel (F(X'), R(Y), Y)] = 0$.

Theorem 3. *Let $\mathbb{P}_{X,Y}$ admits a nonlinear model from X to Y , i.e., $Y = g(f(X), E)$ for some random variable $E \perp\!\!\!\perp X$. We denote by \mathcal{G} , \mathcal{F} and \mathcal{R} the classes from which the algorithm selects the mappings G , F , R (resp.). Let \mathcal{Q} be a class of L -Lipschitz continuous functions $Q : \mathbb{R}^{d_y} \rightarrow \mathbb{R}^{d_e}$. Let \mathcal{T} be a class of functions that satisfies $\mathcal{Q} \circ \mathcal{G} \subset \mathcal{T}$. Let $\mathcal{D} = \{D(u_1, u_2, u_3) := \|T(u_1, u_2) - Q(u_3)\|_2^2\}_{Q \in \mathcal{Q}, T \in \mathcal{T}}$ be the class of discriminators. Then, for any $G \in \mathcal{G}$, $F \in \mathcal{F}$ and $R \in \mathcal{R}$, we have:*

$$\begin{aligned} & \text{Rec}_{\mathcal{T}}(r(Y); R(Y)) \\ & \lesssim \mathcal{L}_{\text{err}}(G, F, R) + \lambda \\ & \quad + \text{disc}_{\mathcal{D}}[(F(X), R(Y), Y) \parallel (F(X'), R(Y), Y)] \quad (35) \end{aligned}$$

where $\lambda := \inf_{Q \in \mathcal{Q}} \mathbb{E}_{y \sim Y} [\|Q(y) - r(y)\|_2^2]$.

As can be seen from Thm. 3, when \mathcal{Q} is expressive enough, such that, λ is small and \mathcal{T} is expressive enough to satisfy $\mathcal{Q} \circ \mathcal{G} \subset \mathcal{T}$, for any functions G, F, R , the recoverability of $r(Y)$ given $R(Y)$ is upper bounded by the sum of the mapping error and the discriminator based independence

measure. Hence, when selecting G, F, R that minimize both losses, one implicitly learns a modeling $G(F(X), R(Y))$, such that, $r(Y)$ can be recovered from $R(Y)$. By a similar argument, the same relation holds for $f(X)$ and $F(X)$.

Proof. Let $Q^* \in \arg \min_{Q \in \mathcal{Q}} \mathbb{E}_{y \sim Y} [\|Q(y) - r(y)\|_2^2]$. We consider that:

$$\begin{aligned} & \inf_{T \in \mathcal{T}} \mathbb{E}_{(x,y) \sim (X,Y)} \|T(R(y)) - r(y)\|_2^2 \\ & \leq 3 \inf_{T \in \mathcal{T}} \mathbb{E}_{(x,y) \sim (X,Y)} \|T(R(y)) - Q^*(y)\|_2^2 \\ & \quad + 3 \inf_{Q \in \mathcal{Q}} \|Q(y) - r(y)\|_2^2 \\ & = 3 \inf_{T \in \mathcal{T}} \mathbb{E}_{(x,y) \sim (X,Y)} \|T(R(y)) - Q^*(y)\|_2^2 + 3\lambda \\ & = 3 \inf_{T \in \mathcal{T}} \mathbb{E}_{\substack{x' \sim X \\ (x,y) \sim (X,Y)}} \|T(F(x'), R(y)) - Q^*(y)\|_2^2 + 3\lambda \quad (36) \end{aligned}$$

where x' and x are two independent copies of X . The last equation follows from the fact that x' and y are independent and from the definition of \mathcal{T} ,

$$\begin{aligned} & \inf_{T \in \mathcal{T}} \mathbb{E}_{\substack{x' \sim X \\ (x,y) \sim (X,Y)}} \|T(F(x'), R(y)) - Q^*(y)\|_2^2 \\ & \geq \inf_{T \in \mathcal{T}} \mathbb{E}_{x'} \mathbb{E}_{(x,y) \sim (X,Y)} \|T_{F(x')}(R(y)) - Q^*(y)\|_2^2 \\ & \geq \mathbb{E}_{x'} \inf_{T \in \mathcal{T}} \mathbb{E}_{(x,y) \sim (X,Y)} \|T_{F(x')}(R(y)) - Q^*(y)\|_2^2 \quad (37) \\ & \geq \mathbb{E}_{x'} \inf_{T \in \mathcal{T}} \mathbb{E}_{(x,y) \sim (X,Y)} \|T(R(y)) - Q^*(y)\|_2^2 \\ & = \inf_{T \in \mathcal{T}} \mathbb{E}_{(x,y) \sim (X,Y)} \|T(R(y)) - Q^*(y)\|_2^2 \end{aligned}$$

Next we consider that for any $T \in \mathcal{T}$, we can rewrite:

$$\begin{aligned} & \mathbb{E}_{\substack{x' \sim X \\ (x,y) \sim (X,Y)}} \|T(F(x'), R(y)) - Q^*(y)\|_2^2 \\ & = \mathbb{E}_{(x,y) \sim (X,Y)} \|T(F(x), R(y)) - Q^*(y)\|_2^2 \\ & \quad + \left\{ \mathbb{E}_{\substack{x' \sim X \\ (x,y) \sim (X,Y)}} \|T(F(x'), R(y)) - Q^*(y)\|_2^2 \right. \\ & \quad \left. - \mathbb{E}_{(x,y) \sim (X,Y)} \|T(F(x), R(y)) - Q^*(y)\|_2^2 \right\} \\ & \leq \mathbb{E}_{(x,y) \sim (X,Y)} \|T(F(x), R(y)) - Q^*(y)\|_2^2 \\ & \quad + \text{disc}_{\mathcal{D}}[(F(X), R(Y), Y) \parallel (F(X'), R(Y), Y)] \quad (38) \end{aligned}$$

Since the class \mathcal{T} includes $Q^* \circ G$, we have:

$$\begin{aligned} & \inf_T \mathbb{E}_{(x,y) \sim (X,Y)} \|T(R(y)) - r(y)\|_2^2 \\ & \leq 3 \mathbb{E}_{(x,y) \sim (X,Y)} \|Q^*(G(F(x), R(y))) - Q^*(y)\|_2^2 \\ & \quad + \text{disc}_{\mathcal{D}}[(F(X), R(Y), Y) \parallel (F(X'), R(Y), Y)] + 3\lambda \quad (39) \end{aligned}$$

Since Q^* is a L -Lipschitz function for some constant $L > 0$, we have the desired inequality. \square