# Generative causal explanations of black-box classifiers

**Matthew O'Shaughnessy, Gregory Canal, Marissa Connor,**
**Mark Davenport, and Christopher Rozell**
School of Electrical & Computer Engineering
Georgia Institute of Technology

## Abstract

We develop a method for generating causal post-hoc explanations of black-box classifiers based on a learned low-dimensional representation of the data. The explanation is causal in the sense that changing learned latent factors produces a change in the classifier output statistics. To construct these explanations, we design a learning framework that leverages a generative model and information-theoretic measures of causal influence. Our objective function encourages both the generative model to faithfully represent the data distribution and the latent factors to have a large causal influence on the classifier output. Our method learns both global and local explanations, is compatible with any classifier that admits class probabilities and a gradient, and does not require labeled attributes or knowledge of causal structure. Using carefully controlled test cases, we provide intuition that illuminates the function of our objective. We then demonstrate the practical utility of our method on image recognition tasks.[1]

## 1 Introduction

There is a growing consensus among researchers, ethicists, and the public that machine learning models deployed in sensitive applications should be able to *explain* their decisions [1, 2]. A powerful way to make "explain" mathematically precise is to use the language of causality: explanations should identify *causal* relationships between certain data aspects — features which may or may not be semantically meaningful — and the classifier output [3–5]. In this conception, an aspect of the data helps explain the classifier if changing that aspect (while holding other data aspects fixed) produces a corresponding change in the classifier output.

Constructing causal explanations requires reasoning about how changing different aspects of the input data affects the classifier output, but these observed changes are only meaningful if the modified combination of aspects occurs naturally in the dataset. A challenge in constructing causal explanations is therefore the ability to change certain aspects of data samples without leaving the data distribution. In this paper we propose a novel learning-based framework that overcomes this challenge. Our framework has two fundamental components that we argue are necessary to operationalize a causal explanation: a method to *represent and move within the data distribution*, and a *rigorous metric for causal influence* of different data aspects on the classifier output.

To do this, we construct a generative model consisting of a disentangled representation of the data and a generative mapping from this representation to the data space (Figure 1(a)). We seek to learn this disentangled representation in such a way that each factor controls a different aspect of the data, and a subset of the factors have a large causal influence on the classifier output. To formalize this notion of causal influence, we define a structural causal model (SCM) [6] that relates independent
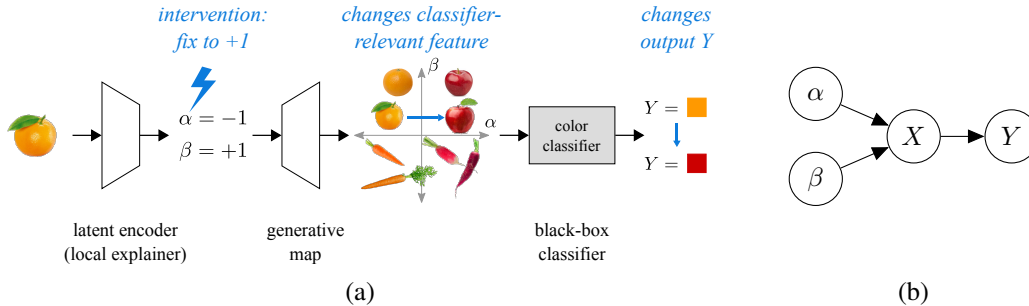
---

[1]Code is available at `https://github.com/siplab-gt/generative-causal-explanations`.

Figure 1: (a) Computational architecture used to learn explanations. Here, the low-dimensional representation $(\alpha, \beta)$ learns to describe the color and shape of inputs. Changing $\alpha$ (color) changes the output of the classifier, which detects the color of the data sample, while changing $\beta$ (shape) does not affect the classifier output. (b) DAG describing our causal model, satisfying principles in Section 3.1.

latent factors defining data aspects, the classifier inputs, and the classifier outputs. Leveraging recent work on information-theoretic measures of causal influence [7, 8], we use the independence of latent factors in the SCM to show that in our framework the causal influence of the latent factors on the classifier output can be quantified simply using mutual information. The crux of our approach is an optimization program for learning a mapping from the latent factors to the data space. The objective ensures that the learned disentangled representation represents the data distribution while simultaneously encouraging a subset of latent factors to have a large causal influence on the classifier output.

A natural benefit of our framework is that the learned disentangled representation provides a rich and flexible vocabulary for explanation. This vocabulary can be more expressive than feature selection or saliency map-based explanation methods: a latent factor, in its simplest form, could describe a single feature or mask of features in input space, but it can also describe much more complex patterns and relationships in the data. Crucially, unlike methods that crudely remove features directly in data space, the generative model enables us to construct explanations that respect the data distribution. This is important because an explanation is only meaningful if it describes combinations of data aspects that naturally occur in the dataset. For example, a loan applicant would not appreciate being told that his loan would have been approved if he had made a negative number of late payments, and a doctor would be displeased to learn that her automated diagnosis system depends on a biologically implausible attribute.

Once the disentangled representation is learned, explanations can be constructed using the generative mapping. Our framework can provide both global and local explanations: a practitioner can understand the aspects of the data that are important to the classifier at large by visualizing the effect in data space of changing each causal factor, and they can determine the aspects that dictated the classifier output for a specific input by observing its corresponding latent values. These visualizations can be much more descriptive than saliency maps, particularly in vision applications.

The major contributions of this work are a new conceptual framework for generating explanations using causal modeling and a generative model (Section 3), analysis of the framework in a simple setting where we can obtain analytical and intuitive understanding (Section 4), and a brief evaluation of our method applied to explaining image recognition models (Section 5).

## 2   Related work

We focus on methods that generate *post-hoc* explanations of black-box classifiers. While post-hoc explanations are typically categorized as either global (explaining the entire classifier mechanism) or local (explaining the classification of a particular datapoint) [9], our framework joins a smaller group of methods that globally learn a model that can be then used to generate local explanations [10–13].

**Forms of explanation.** Post-hoc explanations come in varying forms. Some methods learn an interpretable model such as a decision tree that *approximates the black-box* either globally [14–16] or locally [17–20]. A larger class of methods create local explanations directly in the data space,

2

performing *feature selection* or creating *saliency maps* using classifier gradients [21–25] or by training a new model [10]. A third category of methods generate *counterfactual data points* that describe how inputs would need to be altered to produce a different classifier output [26–32]. Other techniques identify the *points in the training set* most responsible for a particular classifier output [33, 34]. Our framework belongs to a separate class of methods whose explanations consist of a low-dimensional set of *latent factors* that describe different aspects (or "concepts") of the data. These latent factors form a rich and flexible vocabulary for both global and local explanations, and provide a means to represent the data distribution. Unlike some methods that learn concepts using labeled attributes [35, 36], we do not require side information defining data aspects; rather, we visualize the learned aspects using a generative mapping to the data space as in [37–39]. This type of latent factor explanation has also been used in the construction of self-explaining neural networks [37, 40].

**Causality in explanation.** Because explanation methods seek to answer "why" and "how" questions that use the language of cause and effect [3, 4], causal reasoning has played an increasingly important role in designing explanation frameworks [5]. (For similar reasons, causality has played a prominent part in designing metrics for fairness in machine learning [41–45].) Prior work has quantified the impact of features in data space by using Granger causality [13], a priori known causal structure [46, 36], an average or individual causal effect metric [47, 19], or by applying random valued-interventions [48]. Other work generates causal explanations by performing interventions in different network layers [49], using latent factors built into a modified network architecture [38], or using labeled examples of human-interpretable latent factors [50].

Generative models have been used to compute interventions that respect the data distribution [51, 36, 19, 52], a key idea in this paper. Our work, however, is most similar to methods using generative models whose explanations use notions of causality and are constructed directly from latent factors. Goyal et al. compute the average causal effect (ACE) of human-interpretable concepts on the classifier [50], but require labeled examples of the concepts and suffer from limitations of the ACE metric [8]. Harradon et al. construct explanations based on latent factors, but these explanations are specific to neural network classifiers and require knowledge of the classifier network architecture [38]. Our method is unique in constructing a framework from principles of causality that generates latent factor-based explanations of black-box classifiers without requiring side information.

**Disentanglement perspective.** Our method can also be interpreted as a *disentanglement* procedure [53, 54] supervised by classifier output probabilities. Unlike work that encourages a one-to-one correspondence between individual latent factors and semantically meaningful features (i.e., "data generating factors"), we aim to separate the latent factors that are relevant to the classifier's decision from those that are irrelevant. We outline connections to this literature in more detail in Section 3.5.

# 3   Methods

Our goal is to explain a black-box classifier $f\colon \mathcal{X} \to \mathcal{Y}$ that takes data samples $X \in \mathcal{X}$ and assigns a probability to each class $Y \in \{1, \ldots, M\}$ (i.e., $\mathcal{Y}$ is the $M$-dimensional probability simplex). We assume that the classifier also provides the gradient of each class probability with respect to the classifier input.

Our explanations take the form of a low-dimensional and independent set of "causal factors" $\alpha \in \mathbb{R}^K$ that, when changed, produce a corresponding change in the classifier output statistics. We also allow for additional independent latent factors $\beta \in \mathbb{R}^L$ that contribute to representing the data distribution but need not have a causal influence on the classifier output. Together, $(\alpha, \beta)$ constitute a low-dimensional representation of the data distribution $p(X)$ through the generative mapping $g\colon \mathbb{R}^{K+L} \to \mathcal{X}$. The generative mapping is learned so that the explanatory factors $\alpha$ have a large causal influence on $Y$, while $\alpha$ and $\beta$ together faithfully represent the data distribution (i.e., $p(g(\alpha, \beta)) \approx p(X)$). The $\alpha$ learned in this manner can be interpreted as aspects *causing $f$* to make classification decisions [6].

To learn a generative mapping with these characteristics, we need to define (i) a model of the causal relationship between $\alpha$, $\beta$, $X$, and $Y$, (ii) a metric to quantify the causal influence of $\alpha$ on $Y$, and (iii) a learning framework that maximizes this influence while ensuring that $p(g(\alpha, \beta)) \approx p(X)$.

### 3.1 Causal model

We first define a directed acyclic graph (DAG) describing the relationship between $(\alpha, \beta)$, $X$, and $Y$, which will allow us to derive a metric of causal influence of $\alpha$ on $Y$. We propose the following principles for selecting this DAG:

(1) **The DAG should describe the functional (causal) structure of the data, not simply the statistical (correlative) structure.** This principle allows us to interpret the DAG as a structural causal model (SCM) [6] and interpret our explanations causally.

(2) **The explanation should be derived from the classifier output $Y$, not the ground truth classes.** This principle affirms that we seek to understand the action of the classifier, not the ground truth classes.

(3) **The DAG should contain a (potentially indirect) causal link from $X$ to $Y$.** This principle ensures that our causal model adheres to the functional operation of $f \colon X \to Y$.

Based on these principles, we adopt the DAG shown in Figure 1(b). Note that the difference in the roles played by $\alpha$ and $\beta$ is subtle and not apparent from the DAG alone: the difference arises from the fact that the functional relationship defining the causal connection $X \to Y$ is $f$, which by construction uses only features of $X$ that are controlled by $\alpha$. In other words, interventions on both $\alpha$ and $\beta$ produce changes in $X$, but only interventions on $\alpha$ produce changes in $Y$. A key feature of this DAG is that the latent factors $(\alpha, \beta)$ are independent, which we enforce with an isotropic prior when learning the generative mapping. This independence improves the parsimony and interpretability of the learned disentangled representation (see Appendix A). It also results in our metric for causal influence simplifying to mutual information. Importantly, unlike methods that *assume* independence of features in data space (e.g., [48, 17, 23, 25]), our framework *intentionally learns* independent latent factors.

### 3.2 Metric for causal influence

We now derive a metric $\mathcal{C}(\alpha, Y)$ for the causal influence of $\alpha$ on $Y$ using the DAG in Figure 1(b). A satisfactory measure of causal influence in our application should satisfy the following principles:

(1) **The metric should completely capture functional dependencies.** This principle allows us to capture the complete causal influence of $\alpha$ on $Y$ through the generative mapping $g$ and classifier $f$, which may both be defined by complex and nonlinear functions such as neural networks.

(2) **The metric should quantify indirect causal relationships between variables.** This principle allows us to quantify the indirect causal relationship between $\alpha$ and $Y$.

Principle 1 eliminates common metrics such as the average causal effect (ACE) [55] and analysis of variance (ANOVA) [56], which capture only causal relationships between first- and second-order statistics, respectively [8]. Recent work has overcome these limitations by using information-theoretic measures [7, 8, 57]. Of these, we select the *information flow* measure of [7] to satisfy Principle 2 because it is node-based, naturally accommodating our goal of quantifying the causal influence of $\alpha$ on $Y$.

The information flow metric adapts the concept of mutual information typically used to quantify *statistical* influence to quantify *causal* influence by the observational distributions in the standard definition of conditional mutual information with interventional distributions:

**Definition 1** (Ay and Polani 2008 [7]). *Let $U$ and $V$ be disjoint subsets of nodes. The* information flow *from $U$ to $V$ is*

$$I(U \to V) := \int_U p(u) \int_V p(v \mid do(u)) \log \frac{p(v \mid do(u))}{\int_{u'} p(u')p(v \mid do(u'))du'} dV \, dU, \qquad (1)$$

*where $do(u)$ represents an intervention in a causal model that fixes $u$ to a specified value regardless of the values of its parents [6].*

The independence of $(\alpha, \beta)$ makes it simple to show that information flow and mutual information coincide in our DAG:

**Proposition 2** (Information flow in our DAG). *The information flow from $\alpha$ to $Y$ in the DAG of Figure 1(b) coincides with the mutual information between $\alpha$ and $Y$. That is, $I(\alpha \to Y) = I(\alpha; Y)$, where mutual information is defined as $I(\alpha; Y) = \mathbb{E}_{\alpha, Y} \left[ \log \frac{p(\alpha, Y)}{p(\alpha)p(Y)} \right]$.*

---

**Algorithm 1** Principled procedure for selecting $(K, L, \lambda)$.

---

1: Initialize $K, L, \lambda = 0$. Optimizing only $\mathcal{D}$, increase $L$ until objective plateaus.
2: **repeat** increment $K$ and decrement $L$. Increase $\lambda$ until $\mathcal{D}$ approaches value from Step 1.
3: **until** $\mathcal{C}$ reaches plateau. Use $(K, L, \lambda)$ from immediately before plateau was reached.

---

The proof, which follows easily from the rules of do-calculus [6, Thm. 3.4.1], is provided in Appendix C.1. Based on this result, we use

$$\mathcal{C}(\alpha, Y) = I(\alpha; Y) \tag{2}$$

to quantify the causal influence of $\alpha$ on $Y$. This metric, derived in our work from principles of causality using the DAG in Figure 1(b), has also been used to select informative features in other work on explanation [58, 11, 40, 59–61]. Our framework, then, generates explanations that benefit from both causal and information-theoretic perspectives. Note, however, that the validity of the causal interpretation is predicated on our modeling decisions; mutual information is in general a correlational, not causal, metric.

Other variants of (conditional) mutual information are also compatible with our development. These variants retain causal interpretations, but produce explanations of a slightly different character. For example, $\sum_{i=1}^{K} I(\alpha_i; Y)$ and $I(\alpha; Y \mid \beta)$ (the latter corresponding to the information flow of $\alpha$ on $Y$ when "imposing" $\beta$ in [7]) encourage interactions between the explanatory features to generate $X$. These variants are described and analyzed in more detail in Appendices A and B.

### 3.3 Optimization framework

We now turn to our goal of learning a generative mapping $g \colon (\alpha, \beta) \to X$ such that $p(g(\alpha, \beta)) \approx p(X)$, the $(\alpha, \beta)$ are independent, and $\alpha$ has a large causal influence on $Y$. We do so by solving

$$\underset{g \in G}{\arg\max} \quad \mathcal{C}(\alpha, Y) + \lambda \cdot \mathcal{D}\left(p(g(\alpha, \beta)), p(X)\right), \tag{3}$$

where $g$ is a function in some class $G$, $\mathcal{C}(\alpha, Y)$ is our metric for the causal influence of $\alpha$ on $Y$ from (2), and $\mathcal{D}(p(g(\alpha, \beta)), p(X))$ is a measure of the similarity between $p(g(\alpha, \beta))$ and $p(X)$.

The use of $\mathcal{D}$ is a crucial feature of our framework because it forces $g$ to produce samples that are in the data distribution $p(X)$. Without this property, the learned causal factors could specify combinations of aspects that do not occur in the dataset, providing little value for explanation. The specific form of $\mathcal{D}$ is dependent on the class of decoder models $G$. In this paper we focus on two specific instantiations of $G$. Section 4 takes $G$ to be the set of linear mappings with Gaussian additive noise, using negative KL divergence for $\mathcal{D}$. This setting allows us to provide more rigorous intuition for our model. Section 5 adopts the variational autoencoder (VAE) framework shown in Figure 1(a), parameterizing $G$ by a neural network and using a variational lower bound [62] as $\mathcal{D}$.

### 3.4 Training procedure

In practice, we maximize the objective (3) using Adam [63], computing a sample-based estimate of $\mathcal{C}$ at each iteration. The sampling procedure is detailed in Appendix D. Training our causal explanatory model requires selecting $K$ and $L$, which define the number of latent factors, and $\lambda$, which trades between causal influence and data fidelity in our objective. A proper selection of these parameters should set $\lambda$ sufficiently large so that the distributions $p(X \mid \alpha, \beta)$ used to visualize explanations lie in the data distribution $p(X)$, but not so high that the causal influence term is overwhelmed.

To properly navigate this trade-off it is instructive to view (3) as a constrained problem [64] in which $\mathcal{C}$ is maximized subject to an upper bound on $\mathcal{D}$. Algorithm 1 provides a principled method for parameter selection based on this idea. Step 1 selects the total number of latent factors needed to adequately represent $p(X)$ using only noncausal factors. Steps 2-3 then incrementally convert noncausal factors into causal factors until the total explanatory value of the causal factors (quantified by $\mathcal{C}$) plateaus. Because changing $K$ and $L$ affects the relative weights of the causal influence and data fidelity terms, $\lambda$ should be increased after each increment to ensure that the learned representation continues to satisfy the data fidelity constraint.
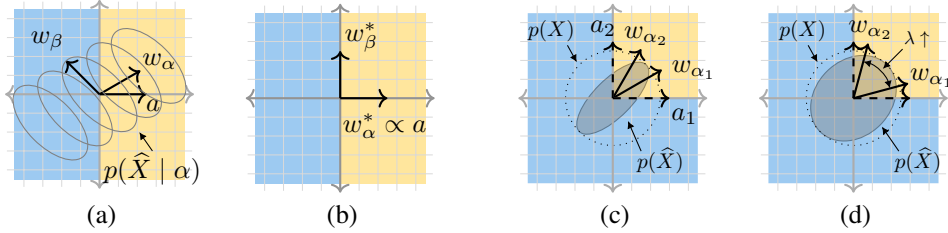
Figure 2: Explaining simple classifiers in $\mathbb{R}^2$. (a) Visualizing the conditional distribution $p(\widehat{X} \mid \alpha)$ provides intuition for the linear-Gaussian model. (b) Linear classifier with yellow encoding high probability of $y = 1$ (right side), and blue encoding high probability of $y = 0$ (left side). Proposition 3 shows that the optimal solution to (3) is $w_\alpha^* \propto a$ and $w_\beta^* \perp w_\alpha^*$ for $\lambda > 0$. (c-d) For the "and" classifier, varying $\lambda$ trades between causal alignment and data representation.

## 3.5 Disentanglement perspective

Disentanglement procedures seek to learn low-dimensional data representations in which latent factors correspond to data aspects that concisely and independently describe high dimensional data [53, 54]. Although some techniques perform unsupervised disentanglement [65–67], it is common to use side information as a supervisory signal.

Because our goal is explanation, our main objective is to separate classifier-relevant and classifier-irrelevant aspects. Our framework can be thought of as a disentanglement procedure with two distinguishing features:

First, we use classifier probabilities to aid disentanglement. This is similar in spirit to disentanglement methods that incorporate grouping or class labels as side information by modifying the VAE training procedure [68], probability model [69], or loss function [70]. Although these methods could be adapted for explanation using classifier-based groupings, our method intelligently uses classifier *probabilities* and gradients.

Second, we develop our framework from a causal perspective. Suter et al. also develop a disentanglement procedure from principles of causality [71], casting the disentanglement task as learning latent factors that correspond to parent-less causes in the generative structural causal model. Unlike this framework, we assume that the latent factors are independent based on properties of the VAE evidence lower bound. We then use this fact to show that the commonly-used MI metric measures *causal* influence of $\alpha$ on $Y$ using the information flow metric of [7].

This provides a causal interpretation for information-based disentanglement methods such as In-foGAN [66] (which adds a term similar to $I(\alpha; X)$ to the VAE objective). Encouragement of independence in latent factors plays an important role in much work on disentanglement (e.g., [65, 66, 72]); priors that better encourage independence could be applied in our framework to increase the validity of our proposed causal graph.

## 4 Analysis with linear-Gaussian generative map

We first consider the instructive setting in which a linear generative mapping is used to explain simple classifiers with decision boundaries defined by hyperplanes. This setting admits geometric intuition and basic analysis that illuminates the function of our objective.

In this section we define the data distribution as isotropic normal in $\mathbb{R}^N$, $X \sim \mathcal{N}(0, I)$ (but note that elsewhere in the paper we make no assumptions on the data distribution). Let $(\alpha, \beta) \sim \mathcal{N}(0, I)$, and consider the following generative model to be used for constructing explanations:

$$g(\alpha, \beta) = [W_\alpha \quad W_\beta] \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + \varepsilon,$$

where $W_\alpha \in \mathbb{R}^{N \times K}$, $W_\beta \in \mathbb{R}^{N \times L}$, and $\varepsilon \sim \mathcal{N}(0, \gamma I)$. We illustrate the behavior of our method applied with this generative model on two simple binary classifiers ($Y \in \{0, 1\}$).

**Linear classifier.** Consider first a linear separator $p(y = 1 \mid x) = \sigma(a^T x)$, where $a \in \mathbb{R}^N$ denotes the decision boundary normal and $\sigma$ is a sigmoid function (visualized in $\mathbb{R}^2$ in Figure 2(a)). With a single causal and single noncausal factor ($K = L = 1$), learning an explanation consists of finding the $w_\alpha, w_\beta \in \mathbb{R}^2$ that maximize (3). Intuitively, we expect $w_\alpha$ to align with $a$ because this direction allows $\alpha$ to produce the largest change in classifier output statistics. This can be seen by considering the distribution $p(\widehat{X} \mid \alpha)$ depicted in Figure 2(a), where we denote $\widehat{X} = g(\alpha, \beta)$ for convenience. Since the generative model is linear-Gaussian, varying $\alpha$ translates $p(\widehat{X} \mid \alpha)$ along the direction $w_\alpha$. When this direction is more aligned with the classifier normal $a$, interventions on $\alpha$ cause a larger change in classifier output by moving $p(\widehat{X} \mid \alpha)$ across the decision boundary. Because the data distribution is isotropic, we expect $\mathcal{D}$ to achieve its maximum when $w_\beta$ is orthogonal to $w_\alpha$, allowing $w_\alpha$ and $w_\beta$ to perfectly represent the data distribution. By combining these two insights, we see that the solution of (3) is given by $w_\alpha^* \propto a$ and $w_\beta^* \perp w_\alpha^*$ (Figure 2(b)).

This intuition is formalized in the following proposition, where for analytical convenience we use the (sigmoidal) normal cumulative distribution function as the classifier nonlinearity $\sigma$:

**Proposition 3.** *Let $\mathcal{X} = \mathbb{R}^N$, $K = 1$, $L = N - 1$, and $p(Y = 1 \mid x) = \sigma(a^T x)$, where $\sigma$ is the normal cumulative distribution function. Suppose that the columns of $W = [w_\alpha \; W_\beta]$ are normalized to magnitude $\sqrt{1 - \gamma}$ with $\gamma < 1$. Then for any $\lambda > 0$ and for $\mathcal{D}(p(\widehat{X}), p(X)) = -\mathrm{D}_{\mathrm{KL}}(p(X) \parallel p(\widehat{X}))$, the objective (3) is maximized when $w_\alpha \propto a$, $W_\beta^T a = 0$, and $W_\beta^T W_\beta = (1 - \gamma)I$.*

The proof, which is listed in Appendix C.2, follows geometric intuition for the behavior of $\mathcal{C}$. This result verifies our objective's ability to construct explanations with our desired properties: the causal factor learns the direction in which the classifier output changes, and the complete set of latent factors represent the data distribution.

**"And" classifier.** Now consider the slightly more complex "and" classifier parameterized by two orthogonal hyperplane normals $a_1, a_2 \in \mathbb{R}^2$ (Figure 2(c)) given by $p(Y = 1 \mid x) = \sigma(a_1^T x) \cdot \sigma(a_2^T x)$. This classifier assigns a high probability to $Y = 1$ when both $a_1^T x > 0$ and $a_2^T x > 0$. Here we use $K = 2$ causal factors and $L = 0$ noncausal factors to illustrate the role of $\lambda$ in trading between the terms in our objective. In this setting, learning an explanation entails finding the $w_{\alpha_1}, w_{\alpha_2} \in \mathbb{R}^2$ that maximize (3).

Figure 2(c-d) depicts the effect of $\lambda$ on the learned $w_{\alpha_1}, w_{\alpha_2}$ (see Appendix B for empirical visualizations). Unlike in the linear classifier case, when explaining the "and" classifier there is a tradeoff between the two terms in our objective: the causal influence term encourages both $w_{\alpha_1}$ and $w_{\alpha_2}$ to point towards the upper right-hand quadrant of the data space, the direction that produces the largest variation in class output probability. On the other hand, the isotropy of the data distribution results in the data fidelity term encouraging orthogonality between the factor directions. Therefore, when $\lambda$ is small the causal effect term dominates, aligning the causal factors to the upper right-hand quadrant of the data space (Figure 2(c)). As $\lambda$ increases (Figure 2(d)), the larger weight on the data fidelity term encourages orthogonality between the factor directions so that $p(\widehat{X})$ more closely approximates $p(X)$. This example illustrates how $\lambda$ must be selected carefully to represent the data distribution while learning meaningful explanatory directions (see Section 3.4).

## 5 Experiments with VAE architecture

In this section we generate explanations of CNN classifiers trained on image recognition tasks, letting $G$ be a set of neural networks and adopting the VAE architecture shown in Figure 1(a) to learn $g$.

**Qualitative results.** We train a CNN classifier with two convolutional layers followed by two fully connected layers on MNIST 3 and 8 digits, a common test setting for explanation methods [25, 13]. Using the parameter tuning procedure described in Algorithm 1, we select $K = 1$ causal factor, $L = 7$ noncausal factors, and $\lambda = 0.05$. Figure 3(a) shows the global explanation for this classifier and dataset, which visualizes how $g(\alpha, \beta)$ changes as $\alpha$ is modified. We observe that $\alpha$ controls the features that differentiate the digits 3 and 8, so changing $\alpha$ changes the classifier output while preserving stylistic features irrelevant to the classifier such as skew and thickness. By contrast, Figures 3(b-d) show that changing each $\beta_i$ affects stylistic aspects such as thickness and skew but not
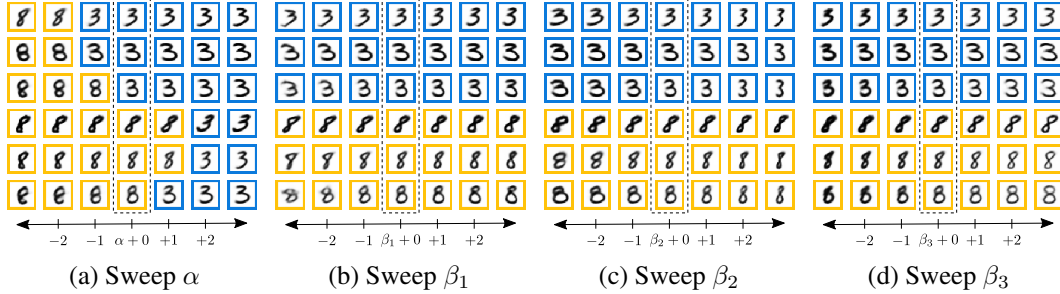
|  (a) Sweep $\alpha$ | (b) Sweep $\beta_1$ | (c) Sweep $\beta_2$ | (d) Sweep $\beta_3$ |

Figure 3: Visualizations of learned latent factors. (a) Changing the causal factor $\alpha$ provides the global explanation of the classifier. Images in the center column of each grid are reconstructed samples from the validation set; moving left or right in each row shows $g(\alpha, \beta)$ as a single latent factor is varied. Changing the learned causal factor $\alpha$ affects the classifier output (shown as colored outlines). (b-d) Changing the noncausal factors $\{\beta_i\}$ affects stylistic aspects such as thickness and skew but does not affect the classifier output.
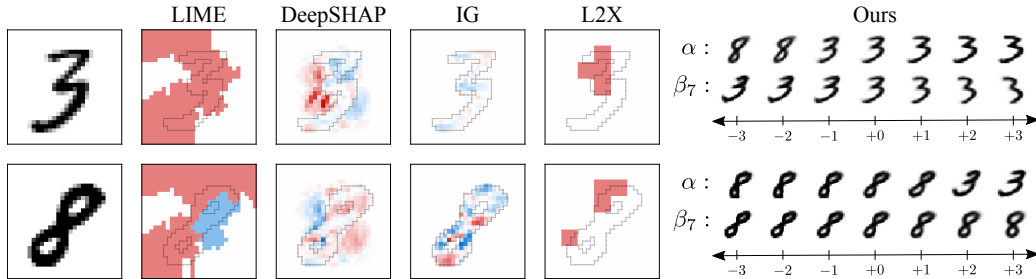


Figure 4: Compared to popular explanation techniques that generate saliency map-based explanations, our explanations consist of learned aspect(s) of the data, visualized by sweeping the associated latent factors (remaining latent factor sweeps are shown in Appendix E.2). Our explanations are able to differentiate causal aspects (pixels that define 3 from 8) from purely stylistic aspects (here, rotation).

the classifier output. Details of the experimental setup and training procedure are listed in Appendix E.1 along with additional results.

**Comparison to other methods.** Figure 4 shows the explanations generated by several popular competitors: LIME [17], DeepSHAP [25], Integrated Gradients (IG) [24], and L2X [11]. Each of these methods generates explanations that quantify a notion of relevance of (super)pixels to the classifier output, visualizing the result with a saliency map. While this form of explanation can be appealing for its simplicity, it fails to capture more complex relationships between pixels. For example, saliency map explanations cannot differentiate the "loops" that separate the digits 3 and 8 from other stylistic factors such as thickness and rotation present in the same (super)pixels. Our explanations overcome this limitation by instead visualizing latent factors that control different aspects of the data. This is demonstrated on the right of Figure 4, where latent factor sweeps show the difference between classifier-relevant and purely stylistic aspects of the data. Observe that $\alpha$ controls data aspects used by the classifier to differentiate between classes, while the noncausal factor controls rotation. Appendix E.2 visualizes the remaining noncausal factors and details the experimental setup.

**Quantitative results.** We next learn explanations of a CNN trained to classify t-shirt, dress, and coat images from the Fashion MNIST dataset [73]. Following the parameter selection procedure of Algorithm 1, we select $K = 2$, $L = 4$, and $\lambda = 0.05$. We evaluate the efficacy of our explanations in this setting using two quantitative metrics. First, we compute the information flow (1) from each latent factor to the classifier output $Y$. Figure 5(a) shows that, as desired, the information flow from $\alpha$ to $Y$ is large while the information flow from $\beta$ to $Y$ is small. Second, we evaluate the reduction in classifier accuracy after individual aspects of the data are removed by fixing a single latent factor in each validation data sample to a different random value drawn from the prior $\mathcal{N}(0, 1)$. This test is frequently used as a metric for explanation quality; our method has the advantage of allowing us to remove certain data aspects while remaining in-distribution rather than crudely removing features
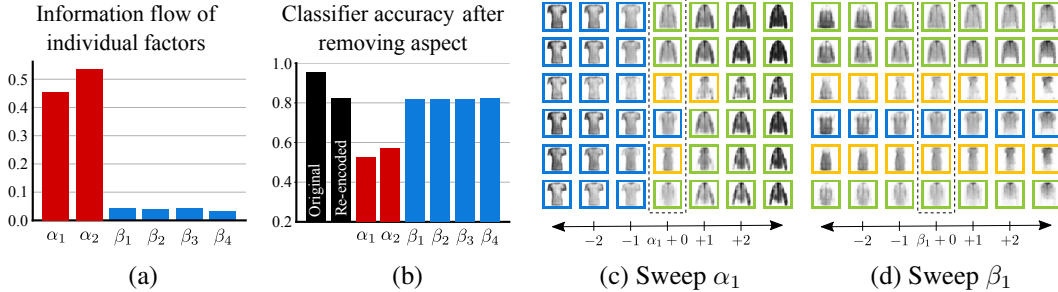
Figure 5: (a) Information flow (1) of each latent factor on the classifier output statistics. (b) Classifier accuracy when data aspects controlled by individual latent factors are removed (original: accuracy on validation set; re-encoded: classifier accuracy on validation set encoded and reconstructed by VAE), showing that learned causal factors (but not noncausal factors) control data aspects relevant to the classifier. (c-d) Modifying $\alpha_1$ changes the classifier output, while modifying $\beta_1$ does not.

by masking (super)pixels [74]. Figure 5(b) shows this reduction in classifier accuracy. Observe that changing aspects controlled by learned causal factors indeed significantly degrades the classifier accuracy, while removing aspects controlled by noncausal factors has only a negligible impact on the classifier accuracy. Figure 5(c-d) visualizes the aspects learned by $\alpha_1$ and $\beta_1$. As before, only the aspects of the data controlled by $\alpha$ are relevant to the classifier: changing $\alpha_1$ produces a change in the classifier output, while changing $\beta_1$ affects only aspects that do not modify the classifier output. Appendix E.3 contains details on the experimental setup and complete results.

## 6 Discussion

The central contribution of our paper is a generative framework for learning a rich and flexible vocabulary to explain a black-box classifier, and a method that uses this vocabulary and causal modeling to construct explanations. Our derivation from a causal model allows us to learn explanatory factors that have a causal, not correlational, relationship with the classifier, and the information-theoretic measure of causality that we adapt allows us to completely capture complex causal relationships. Our use of a generative framework to learn independent latent factors that describe different aspects of the data allows us to ensure that our explanations respect the data distribution.

Applying this framework to practical explanation tasks requires selecting a generative model architecture, and then training this generative model using data relevant to the classification task. The data used to train the explainer may be the original training set of the classifier, but more generally it can be any dataset; the resulting explanation will reveal the aspects in that specific dataset that are relevant to the classifier. The user must also select a generative model $g$ with appropriate capacity. Underestimating this capacity could reduce the effectiveness of the resulting explanations, while overestimating this capacity will needlessly increase the training cost. We explore this selection further in Appendix F both empirically and by using results from [75] to show how the value of $I(\alpha; Y)$ can be interpreted as a "certificate" of sufficient generative model capacity.

Our framework combining generative and causal modeling is quite general. Although we focused on the use of learned data aspects to generate explanations by visualizing the effect of modifying learned causal factors, the learned representation could also be used to generate counterfactual explanations — minimal perturbations of a data sample that change the classifier output [29, 3]. Our framework would address two common challenges in counterfactual explanation: because we can optimize over a low-dimensional set of latent factors, we avoid a computationally infeasible search in input space, and because each point in space maps to an in-distribution data sample, our model naturally ensures that perturbations result in a valid data point. Another promising avenue for future work is relaxing the independence structure of learned causal factors. Although this would result in a more complex expression for information flow, the sampling procedure we use to compute causal effect would generalize naturally; the more challenging obstacle would be learning latent factors with nontrivial causal structure. Finally, techniques that make the classifier-relevant latent factors more interpretable or better communicate the aspects controlled by each latent factor to humans would improve the quality of our generated explanations.

9

## Broader impacts

Explanation methods have the potential to play a major role in enabling the safe and fair deployment of machine learning systems [2, 76], and explainability is a oft-mentioned constraint in their legal and ethical analysis. Policy discussions about machine learning have increasingly turned to principles of transparency and fairness [77], with some legal scholars arguing that the 2016 European General Data Protection Regulation (GDPR) contains a "right to explanation" [78], and recent G20 and OECD recommendations both identifying "transparency and explainability" as important principles for the development of machine learning algorithms [79, 80].

The growing literature on explainability that our work contributes to has the potential to improve the transparency and fairness of machine learning systems and increase the level of trust users place in their decisions. Yet these explanation methods, often built from complex and nontransparent components and each proposing subtly different notions of explanation, also risk providing deceptively incomplete understanding of systems used in sensitive applications, or providing false assurances of fairness and lack of bias (see, e.g., [81]). This criticism may be especially true for our method, which constructs explanations using neural networks that are themselves difficult to understand. For the explanation literature to have a positive impact, it is necessary for explanations to be easily yet precisely understood by the nontechnical generalists deploying and regulating machine learning systems. We believe that causal perspective used in this work is valuable in this regard because causality has been identified as a vocabulary appropriate for translating technical concepts to psychological [3] and legal frameworks [2, 29]. We also believe our analysis with simple models is important because it endows our explanations with some theoretical grounding. However, a critical need remains for more interdisciplinary research examining how end users understand the outputs of explanation tools (e.g., [82]) and how technical tools can be brought to bear to address identified deficiencies.

## Acknowledgments and Disclosure of Funding

## References

[1] Finale Doshi-Velez, Ryan Budish, and Mason Kortz. The Role of Explanation in Algorithmic Trust. Technical report, Artificial Intelligence and Interpretability Working Group, Berkman Klein Center for Internet & Society, December 2017.

[2] Joshua Kroll, Joanna Huey, Solon Barocas, Edward Felten, Joel Reidenberg, David Robinson, and Harlan Yu. Accountable Algorithms. *Univ. Pa. Law Rev.*, 165(3):633, January 2017.

[3] Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial Intelligence*, 267:1–38, February 2019.

[4] Judea Pearl. The Seven Tools of Causal Inference with Reflections on Machine Learning. *Commun. ACM*, 62(3):54–60, March 2019.

[5] Raha Moraffah, Mansooreh Karami, Ruocheng Guo, Adrienne Raglin, and Huan Liu. Causal Interpretability for Machine Learning – Problems, Methods and Evaluation. *ArXiv200303934 Cs Stat*, March 2020.

[6] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, second edition, 2009.

[7] Nihat Ay and Daniel Polani. Information flows in causal networks. *Advs. Complex Syst.*, 11 (01):17–41, February 2008.

[8] Dominik Janzing, David Balduzzi, Moritz Grosse-Wentrup, and Bernhard Schölkopf. Quantifying causal influences. *Ann. Statist.*, 41(5):2324–2358, October 2013.

[9] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. A Survey of Methods for Explaining Black Box Models. *ACM Comput Surv*, 51(5):93:1–93:42, August 2018.

[10] Piotr Dabkowski and Yarin Gal. Real Time Image Saliency for Black Box Classifiers. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS)*, pages 6967–6976, Long Beach, CA, USA, 2017.

[11] Jianbo Chen, Le Song, Martin Wainwright, and Michael Jordan. Learning to Explain: An Information-Theoretic Perspective on Model Interpretation. In *Proc. Int. Conf. on Mach. Learn.*, pages 883–892, Stockholm, Sweden, July 2018.

[12] Seojin Bang, Pengtao Xie, Heewook Lee, Wei Wu, and Eric Xing. Explaining a black-box using deep variational information bottleneck approach. *arXiv:1902.06918*, 2019.

[13] Patrick Schwab and Walter Karlen. CXPlain: Causal Explanations for Model Interpretation under Uncertainty. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS)*, pages 10220–10230, Vancouver, BC, Canada, December 2019.

[14] Mark Craven and Jude W. Shavlik. Extracting Tree-Structured Representations of Trained Networks. In *Proc. Adv. in Neural Inf. Proc. Sys. 1996*, pages 24–30, Denver, CO, USA, 1996.

[15] Osbert Bastani, Carolyn Kim, and Hamsa Bastani. Interpretability via Model Extraction. In *Proc. KDD 2017 Work. on Fairness and Transparency in Machine Learning*, Halifax, NS, Canada, August 2017.

[16] Wenbo Guo, Sui Huang, Yunzhe Tao, Xinyu Xing, and Lin Lin. Explaining deep learning models–a bayesian non-parametric approach. In *Proc. Adv. in Neural Inf. Proc. Syst.*, pages 4514–4524, 2018.

[17] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In *Proc. of the SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD)*, pages 1135–1144, San Francisco, California, USA, 2016.

[18] Himabindu Lakkaraju, Ece Kamar, Rich Caruana, and Jure Leskovec. Interpretable & Explorable Approximations of Black Box Models. In *Proc. 2017 Workshop on Fairness, Accountability, and Transparency in Machine Learning (FAT)*, Halifax, NS, Canada, July 2017.

[19] Carolyn Kim and Osbert Bastani. Learning Interpretable Models with Causal Guarantees. *ArXiv190108576 Cs Stat*, January 2019.

[20] Jorg Wagner, Jan Mathias Kohler, Tobias Gindele, Leon Hetzel, Jakob Thaddaus Wiedemer, and Sven Behnke. Interpretable and Fine-Grained Visual Explanations for Convolutional Neural Networks. In *Proc. Computer Vision and Pattern Recognition*, pages 9089–9099, Long Beach, CA, USA, June 2019.

[21] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. In *Proc. 2014 Int. Conf. on Learning Representations Workshop Track*, Banff, AB, Canada, December 2013.

[22] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation. *PLOS ONE*, 10(7):e0130140, July 2015.

[23] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proc. Int. Conf. on Machine Learning*, pages 3145–3153, Sydney, NSW, Australia, August 2017.

[24] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proc. Int. Conf. on Machine Learning (ICML)*, ICML'17, pages 3319–3328, Sydney, NSW, Australia, August 2017.

[25] Scott M Lundberg and Su-In Lee. A Unified Approach to Interpreting Model Predictions. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS)*, pages 4765–4774, Long Beach, CA, USA, December 2017.

[26] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-Precision Model-Agnostic Explanations. In *Proc. AAAI Conf. on Artificial Intell.*, New Orleans, LA, USA, 2018.

[27] Adam White and Artur Garcez. Towards Providing Causal Explanations for the Predictions of any Classifier. In *Proc. Human-Like Computing Machine Intelligence Workshop (MI21-HLC)*, page 3, July 2019.

[28] Xin Zhang, Armando Solar-Lezama, and Rishabh Singh. Interpreting Neural Network Judgments via Minimal, Stable, and Symbolic Corrections. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS) 2018*, pages 4874–4885, Montréal, Quebec, Canada, December 2018. Curran Associates, Inc.

[29] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual Explanations Without Opening the Black Box: Automated Decisions and the GDPR. *Harv. J. Law Technol.*, 31(2), 2018.

[30] Brandon Carter, Jonas Mueller, Siddhartha Jain, and David Gifford. What made you do this? Understanding black-box decisions with sufficient input subsets. In *Proc. Int. Conf. on Artificial Intell. and Stat. (AISTATS)*, pages 567–576, Naha, Okinawa, Japan, April 2019.

[31] Ramaravind K. Mothilal, Amit Sharma, and Chenhao Tan. Explaining machine learning classifiers through diverse counterfactual explanations. In *Proc. Conf. on Fairness, Accountability, and Transparency (FAT*)*, FAT* '20, pages 607–617, Barcelona, Spain, January 2020.

[32] Arnaud Van Looveren and Janis Klaise. Interpretable Counterfactual Explanations Guided by Prototypes. *ArXiv190702584 Cs Stat*, February 2020.

[33] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 1885–1894, 2017.

[34] Rajiv Khanna, Been Kim, Joydeep Ghosh, and Sanmi Koyejo. Interpreting black box predictions using fisher kernels. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3382–3390, 2019.

[35] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, and Rory Sayres. Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV). In *Proc. Int. Conf. on Machine Learning (ICML)*, Stockholm, Sweden, July 2018.

[36] Álvaro Parafita and Jordi Vitrià. Explaining Visual Models by Causal Attribution. In *Proc. ICCV Work. on Interpretability and Explainability*, Seoul, Korea, November 2019.

[37] Oscar Li, Hao Liu, Chaofan Chen, and Cynthia Rudin. Deep Learning for Case-Based Reasoning through Prototypes: A Neural Network that Explains Its Predictions. In *Proc. AAAI Conf. on Artificial Intelligence*, New Orleans, LA, USA, February 2018.

[38] Michael Harradon, Jeff Druce, and Brian Ruttenberg. Causal Learning and Explanation of Deep Neural Networks via Autoencoded Activations. *ArXiv180200541 Cs Stat*, February 2018.

[39] David Alvarez Melis and Tommi Jaakkola. Towards Robust Interpretability with Self-Explaining Neural Networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 7775–7784. Curran Associates, Inc., 2018.

[40] Maruan Al-Shedivat, Avinava Dubey, and Eric P. Xing. Contextual Explanation Networks. *ArXiv170510301 Cs Stat*, December 2018.

[41] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual Fairness. In *Proc. Adv. in Neural Inf. Proc. Sys. (NeurIPS)*, pages 4066–4076, Long Beach, CA, USA, December 2017.

[42] Niki Kilbertus, Mateo Rojas-Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, and Bernhard Schölkopf. Avoiding Discrimination through Causal Reasoning. In *Proc. Adv. in Neural Inf. Proc. Sys. (NeurIPS)*, Long Beach, CA, USA, December 2017.

[43] Junzhe Zhang and Elias Bareinboim. Fairness in decision-making—the causal explanation formula. In *AAAI Conf. on Artificial Intelligence*, 2018.

[44] Junzhe Zhang and Elias Bareinboim. Equality of opportunity in classification: A causal approach. In *Proc. Adv. in Neural Inf. Proc. Sys. (NeurIPS)*, pages 3671–3681, 2018.

[45] Yongkai Wu, Lu Zhang, Xintao Wu, and Hanghang Tong. Pc-fairness: A unified framework for measuring causality-based fairness. In *Proc. Adv. Neural Inf. Proc. Syst.*, pages 3399–3409, 2019.

[46] Christopher Frye, Ilya Feige, and Colin Rowat. Asymmetric Shapley values: Incorporating causal knowledge into model-agnostic explainability. *ArXiv191006358 Cs Stat*, October 2019.

[47] Aditya Chattopadhyay, Piyushi Manupriya, Anirban Sarkar, and Vineeth N. Balasubramanian. Neural Network Attributions: A Causal Perspective. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 981–990, Long Beach, CA, USA, May 2019.

[48] Anupam Datta, Shayak Sen, and Yair Zick. Algorithmic Transparency via Quantitative Input Influence: Theory and Experiments with Learning Systems. In *2016 IEEE Symp. on Security and Privacy (SP)*, pages 598–617, May 2016.

[49] Tanmayee Narendra, Anush Sankaran, Deepak Vijaykeerthy, and Senthil Mani. Explaining Deep Learning Models using Causal Inference. *ArXiv181104376 Cs Stat*, November 2018.

[50] Yash Goyal, Amir Feder, Uri Shalit, and Been Kim. Explaining Classifiers with Causal Concept Effect (CaCE). *ArXiv190707165 Cs Stat*, February 2020.

[51] David Alvarez-Melis and Tommi Jaakkola. A causal framework for explaining the predictions of black-box sequence-to-sequence models. In *Proc. Conf. Empirical Methods in Natural Language Proc. (EMNLP)*, pages 412–421, Copenhagen, Denmark, 2017.

[52] Chun-Hao Chang, Elliot Creager, Anna Goldenberg, and David Duvenaud. Explaining Image Classifiers by Counterfactual Generation. In *Proc. Int. Conf. on Learning Representations (ICLR) 2019*, New Orleans, LA, USA, May 2019.

[53] Y. Bengio, A. Courville, and P. Vincent. Representation Learning: A Review and New Perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, August 2013.

[54] Irina Higgins, David Amos, David Pfau, Sebastien Racaniere, Loic Matthey, Danilo Rezende, and Alexander Lerchner. Towards a Definition of Disentangled Representations. *ArXiv181202230 Cs Stat*, December 2018.

[55] Paul W. Holland. Causal Inference, Path Analysis and Recursive Structural Equations Models. *ETS Res. Rep. Ser.*, 1988(1):i–50, 1988.

[56] R C Lewontin. The analysis of variance and the analysis of causes. *Am. J. Hum. Genet.*, 26(3): 400–411, May 1974.

[57] Gabriel Schamberg and Todd P Coleman. Quantifying Context-Dependent Causal Influences. In *Proc. NeurIPS 2018 Work. on Causal Learning*, page 10, Montréal, Quebec, Canada, December 2018.

[58] Shuyang Gao, Greg Ver Steeg, and Aram Galstyan. Variational Information Maximization for Feature Selection. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS)*, pages 487–495, Barcelona, Spain, December 2016.

[59] Atsushi Kanehira and Tatsuya Harada. Learning to Explain With Complemental Examples. In *Proc. Conf. on Comp. Vision and Pattern Recognition (CVPR)*, pages 8595–8603, Long Beach, CA, USA, June 2019.

[60] Shiyu Chang, Yang Zhang, Mo Yu, and Tommi Jaakkola. A Game Theoretic Approach to Class-wise Selective Rationalization. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS)*, pages 10055–10065, Vancouver, BC, Canada, 2019.

[61] Tameem Adel, Zoubin Ghahramani, and Adrian Weller. Discovering Interpretable Representations for Both Deep Generative and Discriminative Models. In *Proc. Int. Conf. on Machine Learning (ICML)*, pages 50–59, Stockholm, Sweden, July 2018.

[62] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *Proc. Int. Conf. on Learning Representations (ICLR)*, Banff, AB, Canada, April 2014.

[63] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *ArXiv14126980 Cs*, January 2017.

[64] Stephen P. Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, UK ; New York, 2004. ISBN 978-0-521-83378-3.

[65] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. Beta-VAE: Learning basic visual concepts with a constrained variational framework. In *Proc. Int. Conf. on Learning Representations (ICLR) 2017*, Toulon, France, April 2017.

[66] Xi Chen, Yan Duan, Rein Houthooft, John Schulman, Ilya Sutskever, and Peter Abbeel. InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS) 2016*, Barcelona, Spain, December 2016.

[67] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In *Proc. Int. Conf. on Mach. Learn. (ICML) 2018*, Stockholm, Sweden, June 2018.

[68] Tejas Kulkarni, William Whitney, Pushmeet Kohli, and Josh Tenenbaum. Deep convolutional inverse graphics network. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS)*, Montréal, Quebec, Canada, December 2015.

[69] Diane Bouchacourt, Ryota Tomioka, and Sebastian Sebastian. Multi-level variational autoencoder: Learning disentangled representations from grouped observations. In *Proc. AAAI Conf. on Artificial Intell.*, New Orleans, LA, USA, February 2018.

[70] Karl Ridgeway and Michael Mozer. Learning deep disentangled embeddings with the F-statistic loss. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS)*, Montréal, Quebec, Canada, December 2018.

[71] Raphael Suter, Dorde Miladinovic, Bernhard Schölkopf, and Stefan Bauer. Robustly Disentangled Causal Mechanisms: Validating Deep Representations for Interventional Robustness. In *Proc. Int. Conf. on Mach. Learn. (ICML) 2019*, Long Beach, CA, USA, May 2019.

[72] Ricky T. Q. Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating sources of disentanglement in variational autoencoders. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS)*, Montréal, Quebec, Canada, December 2018.

[73] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *arXiv:1708.07747*, 2017.

[74] Sara Hooker, Dumitru Erhan, Pieter-Jan Kindermans, and Been Kim. A Benchmark for Interpretability Methods in Deep Neural Networks. In *Proc. Adv. Neural Inf. Proc. Syst. (NeurIPS)*, Vancouver, BC, Canada, 2019.

[75] M. Feder and N. Merhav. Relations between entropy and error probability. *IEEE Trans. Inform. Theory*, 40(1):259–266, Jan./1994.

[76] David Danks and Alex John London. Regulating Autonomous Systems: Beyond Standards. *IEEE Intell. Syst.*, 32(1):88–91, January 2017.

[77] Jack Karsten. New White House AI principles reach beyond economic and security considerations, Brookings Institution, January 2020.

[78] Gianclaudio Malgieri and Giovanni Comandé. Why a Right to Legibility of Automated Decision-Making Exists in the General Data Protection Regulation. *International Data Privacy Law*, 7 (4):243–265, November 2017.

[79] G20. G20 Ministerial Statement on Trade and Digital Economy. Technical report, Tsukuba, Japan, June 2019.

[80] OECD. Recommendation of the Council on Artificial Intelligence. Technical Report OECD/LEGAL/0449, OECD, 2020.

[81] Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat Mach Intell*, 1(5):206–215, May 2019.

[82] Sana Tonekaboni, Shalmali Joshi, Melissa D. McCradden, and Anna Goldenberg. What Clinicians Want: Contextualizing Explainable Machine Learning for Clinical End Use. In *Machine Learning for Healthcare Conf.*, pages 359–380, Ann Arbor, MI, USA, August 2019.

[83] Thomas M Cover and Joy A Thomas. Elements of Information Theory. 2006.

# A Intuition for and variants of causal influence metric

**Intuition for causal influence objective.** To better understand the causal portion of our objective (4), we use standard identities to decompose it as

$$\mathcal{C} = I(Y;\alpha) = H(Y) - \mathbb{E}_\alpha[H(Y \mid \alpha)], \tag{4}$$

where

$$p(y \mid \alpha) = \int_\beta \int_x p(y \mid x)p(x \mid \alpha, \beta)p(\beta)dxd\beta. \tag{5}$$

The conditional distribution (5) can be interpreted as the probability of $Y = y$ for a fixed value of $\alpha$, averaged over the values of $\beta$. The decomposition in (4) therefore shows that $\mathcal{C}$ is the *reduction in uncertainty about $Y$ provided by knowledge of* $\alpha$, where this reduction is measured in a global sense in that the effect of $\beta$ is averaged together to produce a single probability estimate for $Y$ and fixed $\alpha$.

As an example, consider the color classifier and generative mapping shown in Figure 1(a), in which $f$ classifies based on color. The first term in (4) represents how similar the classifier output is for all objects in the training set. The second term represents how similar the classifier output for groups of objects is, on average, after being grouped by $\alpha$. A large $\mathcal{C} = I(\alpha; Y)$ means that grouping by $\alpha$ significantly increases the confidence the classifier has that objects in each group are of the same class. In this case, grouping by $\alpha =$ 'color' has a much larger effect on the classifier output — and therefore results in a larger $\mathcal{C}$ — than grouping by $\alpha =$ 'shape' would, since grouping the objects by color results in the classifier gaining much more confidence that each group shares the same class.

**Variants of causal objective.** Consider the following variants of the *joint, unconditional* objective $\mathcal{C} = I(\alpha; Y)$, our measure of causal influence from Section 3.2:

1. *Independent, unconditional:* $\mathcal{C}_{iu} = \frac{1}{K} \sum_i I(\alpha_i; Y)$
2. *Independent, conditional:* $\mathcal{C}_{ic} = \frac{1}{K} \sum_i I(\alpha_i; Y \mid \alpha_{\neg i}, \beta)$, where $\alpha_{\neg i} = \{\alpha_j\}_{j \neq i}$
3. *Joint, conditional:* $\mathcal{C}_{jc} = I(\alpha; Y \mid \beta)$

Each objective variant gives rise to a classifier explanation that has a causal interpretation, but as we will show, the *character* of each is subtly different. The following proposition begins to explore these differences by relating them using information-theoretic quantities.

**Proposition 4** (Relationship between candidate causal objectives)**.** *The following hold in the DAG of Figure 1(b):*

(a) $\mathcal{C} = \mathcal{C}_{iu} + \frac{1}{K} \sum_{i=1}^{K} I(\alpha_{\neg i}; Y \mid \alpha_i)$.

(b) $\mathcal{C}_{jc} = \mathcal{C}_{ic} + \frac{1}{K} \sum_{i=1}^{K} I(\alpha_{\neg i}; Y \mid \beta)$.

(c) $\mathcal{C}_{jc} = \mathcal{C} + I(\alpha; \beta \mid Y)$.

(d) $\mathcal{C}_{ic} = \mathcal{C}_{iu} + \frac{1}{K} \sum_i I(\alpha_i; \alpha_{\neg i}, \beta \mid Y)$.

These relationships are depicted visually in Figure 6 and proved in Appendix C.3. Note that only (c) and (d) use the independence of the latent variables in our DAG. The "adjustment factors" that relate the objective variants can be interpreted as follows:

1. By conditioning on other latent factors (i.e., using $\mathcal{C}_{ic}$, $\mathcal{C}_{iu}$, or $\mathcal{C}_{jc}$ rather than $\mathcal{C}$) we include the "adjustment factor" $\frac{1}{K} \sum_i I(\alpha_i; \alpha_{\neg i}, \beta \mid Y)$ (in the "independent" case) or $I(\alpha; \beta \mid Y)$ (in the "joint" case) in the objective. These terms encourage complex interactions between latent factors within each group of similarly-classified points. On the one hand, the stastistical pattern that these terms encourage arises naturally from the DAG in Figure 1(b): although the latent factors are independent, conditioning on $Y$ renders them dependent. This conditional dependence pattern is often referred to as Berkson's paradox or the "explaining away" phenomenon. To illustrate this concept, consider a classifier that classifies paintings at an auction as $Y \in \{$'sold', 'not sold'$\}$ based on the learned latent factors $z_1 =$ 'beautiful' and $z_2 =$ 'historical value', which we assume to be independent. Once $Y$ is known, however, $z_1$ and $z_2$ are rendered dependent: learning that a sold painting does not have historical

16

| independent, unconditional $\mathcal{C}_{iu} = \frac{1}{K}\sum_i I(\alpha_i; Y)$ | $+\frac{1}{K}\sum_i I(\alpha_{\neg i}; Y \mid \alpha_i)$ | joint, unconditional $\mathcal{C} = I(\alpha; Y)$ |

$+\frac{1}{K}\sum_i I(\alpha_i; \alpha_{\neg i}, \beta \mid Y)$  $+I(\alpha; \beta \mid Y)$

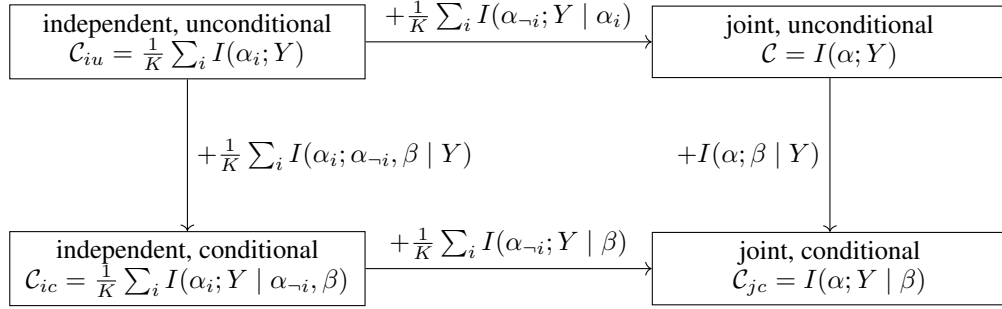| independent, conditional $\mathcal{C}_{ic} = \frac{1}{K}\sum_i I(\alpha_i; Y \mid \alpha_{\neg i}, \beta)$ | $+\frac{1}{K}\sum_i I(\alpha_{\neg i}; Y \mid \beta)$ | joint, conditional $\mathcal{C}_{jc} = I(\alpha; Y \mid \beta)$ |

Figure 6: Graphical representation of relationships between causal objective variants derived from Proposition 4.

value would allow us to infer that it is likely to be beautiful. On the other hand, we do not in general expect that our learned latent factors, which we encourage to be independent, will correspond to semantically meaningful features, so we may not expect them to fit this "explaining away" conditional dependence pattern.

2. By jointly considering the causal factors $\alpha$ rather than summing the causal influence of each $\alpha_i$ (i.e., by using $\mathcal{C}$ rather than $\mathcal{C}_{iu}$, or $\mathcal{C}_{jc}$ rather than $\mathcal{C}_{ic}$) we include the "adjustment factor" $\frac{1}{K}\sum_{i=1}^{K} I(\alpha_{\neg i}; Y \mid \alpha_i)$ in the objective. This term encourages each learned causal factor to make the remaining causal factors more predictable given the classifier output $Y$, encouraging *interactions* between latent factors to have an effect on the classifier output probability. We consider this to be positive, but using an independent objective might aid in visualizing the relationship between the latent space and data space.

The next section provides more intuition for these objectives in the context of the linear-Gaussian generative map and simple classifiers introduced in Section 4.
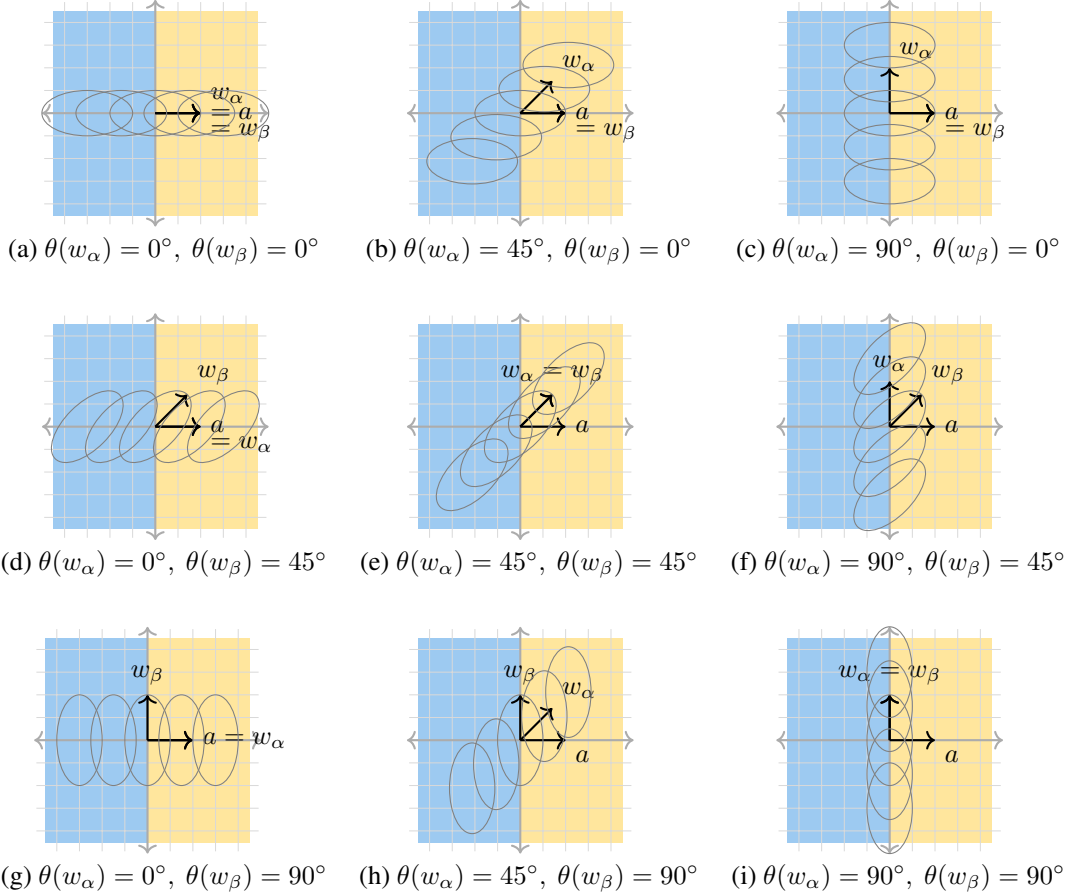
(a) $\theta(w_\alpha) = 0°$, $\theta(w_\beta) = 0°$    (b) $\theta(w_\alpha) = 45°$, $\theta(w_\beta) = 0°$    (c) $\theta(w_\alpha) = 90°$, $\theta(w_\beta) = 0°$

(d) $\theta(w_\alpha) = 0°$, $\theta(w_\beta) = 45°$    (e) $\theta(w_\alpha) = 45°$, $\theta(w_\beta) = 45°$    (f) $\theta(w_\alpha) = 90°$, $\theta(w_\beta) = 45°$

(g) $\theta(w_\alpha) = 0°$, $\theta(w_\beta) = 90°$    (h) $\theta(w_\alpha) = 45°$, $\theta(w_\beta) = 90°$    (i) $\theta(w_\alpha) = 90°$, $\theta(w_\beta) = 90°$

Figure 7: Distributions $p(x \mid \alpha)$ for the linear-Gaussian generative map and single hyperplane classifier when $a = [1,\ 0]^T$. The orientation of $w_\alpha$ controls the direction in which the probability mass of $p(x \mid \alpha)$ shifts as $\alpha$ is varied, while the orientation of $w_\beta$ controls the rotation of each distribution $p(x \mid \alpha)$.

## B    Detailed analysis with linear-Gaussian generative map

In this section we provide empirical simulations supporting the analysis with a linear-Gaussian generative map in Section 4. Recall that we use the isotropic data distribution $X \sim \mathcal{N}(0, I)$, latent space prior $(\alpha, \beta) \sim \mathcal{N}(0, I)$, and

$$g(\alpha, \beta) = \begin{bmatrix} W_\alpha & W_\beta \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} + \varepsilon,$$

where $W_\alpha \in \mathbb{R}^{N \times K}$, $W_\beta \in \mathbb{R}^{N \times L}$, and $\varepsilon \sim \mathcal{N}(0, \gamma I)$.

**Linear classifier.** Consider first the linear separator in $\mathbb{R}^2$ from Section 4, $p(Y = 1 \mid x) = \sigma(a^T x)$. With $K = L = 1$, learning an explanation entails learning the $w_\alpha, w_\beta \in \mathbb{R}^2$ that maximize the objective (3). As shown in Proposition 3, the data representation term $\mathcal{D}$ encourages $w_\alpha \perp w_\beta$; here we focus on the causal influence term $\mathcal{C}$. The decomposition in (4) shows that $\mathcal{C}$ depends on both $p(Y)$ and $p(Y \mid \alpha)$; Figure 7 visualizes how the distributions $p(Y \mid \alpha)$ change with $\alpha$ (gray ellipses) and $w_\alpha, w_\beta$ (subplots). Note first that the isotropy of $p(\alpha)$ means that $p(Y)$ has equal probability mass on either side of the classifier decision boundary, regardless of $w_\alpha$ and $w_\beta$. This implies that $H(Y)$ is invariant to $w_\alpha$ and $w_\beta$ for this classifier, a fact formalized in the proof of Proposition 3.

We next explore the role of $w_\alpha$ and $w_\beta$ in $p(x \mid \alpha)$ (and therefore $p(y \mid \alpha)$). Our causal objective $\mathcal{C}$ is large when the $p(y \mid \alpha)$ have low entropy in expectation over $\alpha$. Note from Figure 7 that $w_\alpha$ controls the direction in which the probability mass of $p(x \mid \alpha)$ shifts as $\alpha$ is varied, while $w_\beta$ controls the
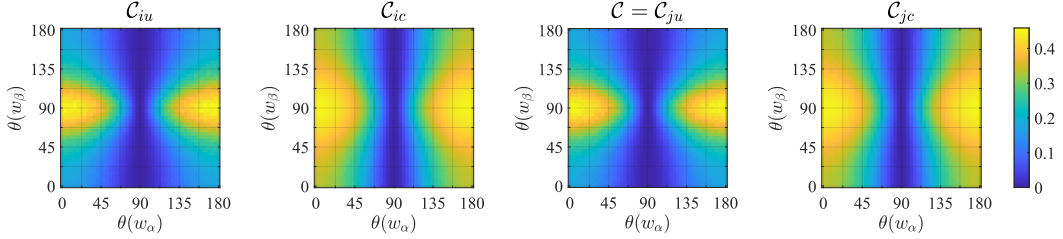
Figure 8: Value of each causal objective variant in the linear-Gaussian generative map, linear classifier setting described in Section 4, as the orientations of $w_\alpha$ and $w_\beta$ are varied. The classifier decision boundary normal is $\theta(a) = 0°$. Each variant is maximized when $w_\alpha \propto a$ (i.e., $\theta(w_\alpha) = 0°$) and $w_\beta \perp a$ (i.e., $\theta(w_\beta) = 90°$). $\mathcal{C} = \mathcal{C}_{ju}$ refers to the causal objective (2) used in the main text.
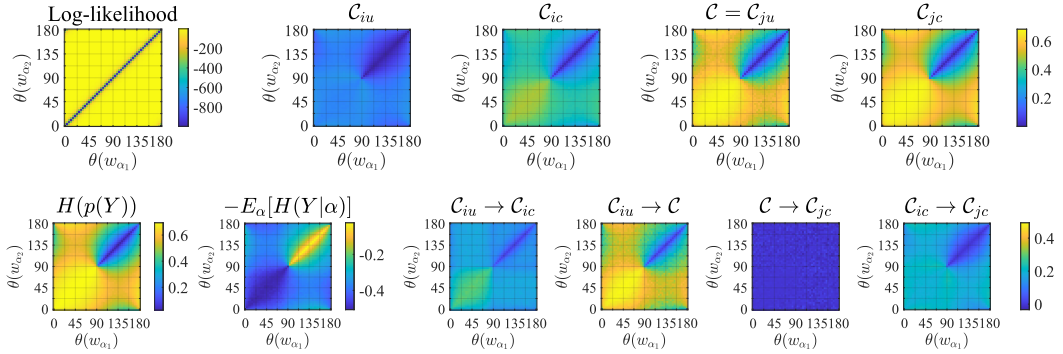


Figure 9: Empirically-computed values of terms relevant to the causal objective variants in the linear-Gaussian generative map, "and" classifier setting described in Section 4. The angles of the classifier decision boundary normals are $\theta(a_1) = 0°$ and $\theta(a_2) = 90°$. Top row: log-likelihood used as $\mathcal{D}$; causal objective variants from Appendix A. $\mathcal{C} = \mathcal{C}_{ju}$ refers to the causal objective (2). Bottom row: terms in decomposition (4); "adjustment factors" from Proposition 4.

rotation of each distribution $p(x \mid \alpha)$. The causal objective $\mathcal{C}$ is maximized when the entropy of $p(y \mid \alpha)$ (in expectation over $\alpha$) is smallest — in other words, when the distributions $p(x \mid \alpha)$ have as little overlap possible with the classifier decision boundary. From Figure 7, we observe that this occurs when $w_\alpha$ is aligned with the decision boundary normal ($w_\alpha \propto a$) and when $w_\beta$ is orthogonal to the decision boundary normal ($w_\beta \perp a$). This selection of $w_\alpha$ and $w_\beta$ minimizes the range of $\alpha$ for which $p(x \mid \alpha)$ contains mass on both sides of the decision boundary.

Figure 8 shows the value of each of the causal objective variants described in Appendix A as the orientation of $w_\alpha$ and $w_\beta$ with respect to the classifier decision boundary normal $a$ are varied. For each combination of angles, we compute the causal objective using the sample-based estimate described in Appendix D with $N_\alpha = 2500$, $N_\beta = 500$, and the logistic sigmoid function $\sigma$ with steepness 5. (Note that in the training procedure we achieve satisfactory results with much lower $N_\alpha, N_\beta$.) These results verify the intuition presented above and formalized in Proposition 3: the causal effect is greatest when $w_\alpha \propto a$ and $w_\beta \perp a$. As noted in Section 4, in this setting both $\mathcal{C}$ and $\mathcal{D}$ encourage $w_\alpha$ and $w_\beta$ to be orthogonal.

**"And" classifier.** We now consider the "and" classifier in $\mathbb{R}^2$ from Section 4, $p(Y = 1 \mid x) = \sigma(a_1^T x) \cdot \sigma(a_2^T x)$, where we learn $K = 2$ causal explanatory factors and $L = 0$ noncausal factors. In this setting learning an explanation consists of learning $w_{\alpha_1}, w_{\alpha_2} \in \mathbb{R}^2$ maximizing (3).

Figure 9 shows how the value of the causal objective changes with the learned generative mapping in the linear-Gaussian setting of Section 4. The top row shows the terms in the objective (3): the likelihood and the causal objective variants described in Appendix A. The bottom row shows the components of these causal objective variants, which provide further intuition for their differences:
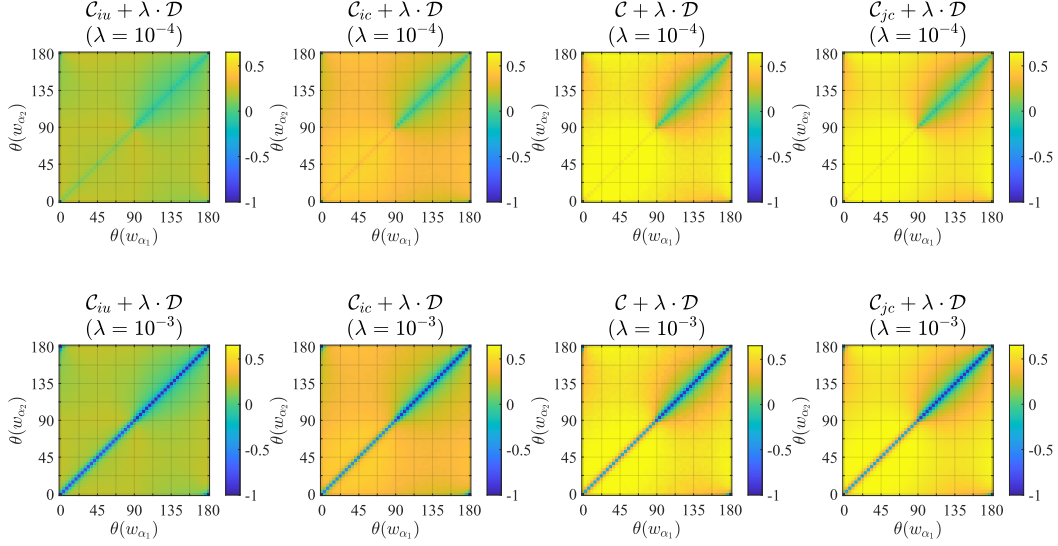
Figure 10: Empirically-computed value of combined objective (3) for the causal objective variants in the linear-Gaussian generative map, "and" classifier setting described in Section 4. The angles of the classifier decision boundary normals are $\theta(a_1) = 0°$ and $\theta(a_2) = 90°$. As $\lambda$ increases, the increased weight of the data representation term in the objective encourages the learned $w_{\alpha_1}$ and $w_{\alpha_2}$ to be more orthogonal to better represent the isotropic distribution of the data.

the first two plots show the decomposition of $\mathcal{C} = \mathcal{C}_{ju}$ from (4), and the remaining plots show the "adjustment factors" from Proposition 4 and Figure 6 that describe the differences between the causal influence objective variants. The logistic sigmoid with steepness 100 is used to implement the classifier, and the causal influence objective variants are computed with $N_\alpha = 2500$ and $N_\beta = 500$.

With the exception of the variant $\mathcal{C}_{iu}$, each of these causal objectives is maximized when $w_{\alpha_1}$ and $w_{\alpha_2}$ are aligned in the direction of maximum classifier change: $\theta(w_{\alpha_1}) = \theta(w_{\alpha_2})$ when $a_1 = [1, \ 0]^T$ and $a_2 = [0, \ 1]^T$ as in our example (see Figure 2(c-d)). Because with this classifier $\mathcal{C}$ does not encourage $w_{\alpha_1} \perp w_{\alpha_2}$, here the data representation term $\mathcal{D}$ serves to regularize $\mathcal{C}$. Figure 10 shows the value of the combined objective (3) for each causal influence variant and two different values of $\lambda$. We observe that as $\lambda$ increases and the weight of the data representation term increases, the optimal angles of $w_{\alpha_1}$ and $w_{\alpha_2}$ move in opposing directions from $45°$ (the angle of normal bisecting $a_1$ and $a_2$). This supports the intuition described in Section 4 and stylized in Figure 2(c-d).

## C  Proofs

### C.1  Proof of Proposition 2

Proposition 2 states that information flow coincides with mutual information in our DAG. Here we prove a generalization of the proposition that is also helpful when considering the conditional causal influence objective variants in Appendix A. Specifically, we consider the information flow from $U$ to $V$ imposing $W$:

**Definition 5** (Ay and Polani 2008 [7])**.** *Let $U$, $V$, and $W$ be disjoint subsets of nodes. The* information flow from $U$ to $V$ imposing $W$, *denoted $I(U \to V \mid W)$, is*

$$\mathbb{E}_{w \sim W} \left[ \int_U p(u \mid do(w)) \int_V p(v \mid do(u), do(w)) \log \frac{p(v \mid do(u), do(w))}{\int_{u'} p(u' \mid do(w)) p(v \mid do(u'), do(w))} dV \, dU \right],$$

*where $do(w)$ represents an intervention in a model that fixes $w$ to a specified value regardless of the values of its parents [6].*

**Proposition 6** (Information flow in our DAG)**.** *The information flow from $\alpha$ to $Y$ imposing $\beta$ in the DAG of Figure 1(b) coincides with the mutual information of $\alpha$ and $Y$ conditioned on $\beta$,*

$$I(\alpha \to Y \mid do(\beta)) = I(\alpha; Y \mid \beta),$$

*where conditional mutual information is defined as $I(X; Y \mid Z) = \mathbb{E}_{X,Y,Z} \left[ \log \frac{p(x,y|z)}{p(x|z)p(y|z)} \right]$.*

*Proof.* The proof follows from the "action/observation exchange" rule of the *do*-calculus [6, Thm. 3.4.1]. This rule asserts that $p(y \mid do(x), do(z), w) = p(y \mid do(x), z, w)$ if $Y \perp Z \mid X, W$ in $\mathcal{G}_{\overline{X}\underline{Z}}$, the causal model modified to remove connections entering $X$ and leaving $Z$. When applied to our model, it yields

1. $p(Y \mid do(\alpha)) = p(Y \mid \alpha)$ (because $Y \perp \alpha$ in $\mathcal{G}_{\underline{\alpha}}$);

2. $p(\alpha \mid do(\beta)) = p(\alpha \mid \beta)$ (because $\alpha \perp \beta$ in $\mathcal{G}_{\underline{\beta}}$); and

3. $p(Y \mid do(\alpha), do(\beta)) = p(Y \mid \alpha, \beta)$ (because $Y \perp (\alpha, \beta)$ in $\mathcal{G}_{\underline{\alpha},\underline{\beta}}$).

Starting with the definition of the information flow from $\alpha$ to $Y$ imposing $\beta$, we have that

$$I(\alpha \to Y \mid do(\beta)) = \mathbb{E}_\beta \left[ \int_\alpha p(\alpha \mid do(\beta)) \int_Y p(Y \mid do(\alpha), do(\beta)) \right.$$
$$\left. \times \log \frac{p(Y \mid do(\alpha), do(\beta))}{\int_{a'} p(\alpha = a' \mid do(\beta)) p(Y \mid do(\alpha = a'), do(\beta))} \right] dY \, d\alpha$$

$$= \mathbb{E}_\beta \left[ \int_\alpha p(\alpha \mid \beta) \int_Y p(Y \mid \alpha, \beta) \right.$$
$$\left. \times \log \frac{p(Y \mid \alpha, \beta)}{\int_{a'} p(\alpha = a' \mid \beta) p(Y \mid \alpha = a', \beta)} \right] dY \, d\alpha$$

$$= \mathbb{E}_\beta \left[ \int_{\alpha,Y} p(Y, \alpha \mid \beta) \log \frac{p(Y \mid \alpha, \beta)}{p(Y \mid \beta)} \right] dY \, d\alpha$$

$$= \int_\beta p(\beta) \int_{\alpha,Y} p(Y, \alpha \mid \beta) \log \frac{p(Y \mid \alpha, \beta)}{p(Y \mid \beta)} dY \, d\alpha \, d\beta$$

$$= \int_\beta p(\beta) \int_{\alpha,Y} p(Y, \alpha \mid \beta) \log \frac{p(Y \mid \alpha, \beta) p(\alpha \mid \beta)}{p(Y \mid \beta) p(\alpha \mid \beta)} dY \, d\alpha \, d\beta$$

$$= \int_\beta p(\beta) \int_{\alpha,Y} p(Y, \alpha \mid \beta) \log \frac{p(Y, \alpha \mid \beta)}{p(Y \mid \beta) p(\alpha \mid \beta)} dY \, d\alpha \, d\beta$$

$$= I(\alpha; Y \mid \beta).$$

$\square$

Proposition 2 follows from Proposition 6 by imposing the null set.

21

## C.2 Proof of Proposition 3

With $K = 1$ we can decompose $\mathcal{C}$ as

$$\mathcal{C} = I(Y; \alpha) = H(Y) - H(Y \mid \alpha). \tag{6}$$

where $H$ denotes entropy of a discrete random variable [83]. First consider the entropy term $H(Y)$. From the illustrations of $p(\widehat{X} \mid \alpha)$ in Figure 7, we can see in $\mathbb{R}^2$ that this entropy is constant for all values of $w_\alpha$ and $w_\beta$: regardless of their angle and offsets, the aggregate set of distributions $p(\widehat{X} \mid \alpha)$ is symmetric about the origin and so the probability mass of $p(\widehat{X})$ is spread symmetrically across both sides of the decision boundary. This idea is generalized in the following lemma, which shows that $H(Y)$ is equal to $\log(2) \approx 0.69$ nats for all values of $W$:

**Lemma 7.** *Under the conditions of Propsition 3, $H(Y) = \log(2)$ nats for all $W \in \mathbb{R}^{N \times N}$.*

*Proof.* Since $(\alpha, \beta) \sim \mathcal{N}(0, I)$, we have $\widehat{X} \sim \mathcal{N}(0, WW^T + \gamma I)$. Letting $U = a^T X$, we have $U \sim \mathcal{N}(0, a^T (WW^T + \gamma I) a)$ which we note has an even probability density function. Considering the classifier output probability marginalized over the generated inputs $\widehat{X}$, we have

$$
\begin{aligned}
p(Y = 1) &= \mathbb{E}_{\widehat{X}}[p(Y = 1 \mid \widehat{X})] \\
&= \mathbb{E}_{\widehat{X}}[\sigma(a^T \widehat{X})] \\
&= \mathbb{E}_U[\sigma(U)] \\
&= \mathbb{E}_U[\sigma(U) - 0.5] + 0.5 \\
&\overset{(\star)}{=} 0.5
\end{aligned}
$$

where in $(\star)$ we use the fact that since $U$ has an even probability density and $\sigma(U) - 0.5$ is an odd function, we have that $\mathbb{E}_U[\sigma(U) - 0.5] = 0$. Letting $h_b(p) = -(p \log p + (1 - p) \log(1 - p))$ denote the binary entropy function, we have that $H(\widehat{Y}) = h_b(p(\widehat{Y} = 1)) = h_b(0.5) = \log(2)$ nats. $\qquad\square$

We now consider the second term in (6), the conditional entropy $H(Y \mid \alpha)$. In $\mathbb{R}^2$ (Figure 7), this term corresponds to the average over $\alpha$ of the classification entropies for each distribution $p(\widehat{X} \mid \alpha)$ (depicted as individual ellipses). Intuitively, this entropy is small when many of the conditional distributions $p(\widehat{X} \mid \alpha)$ lie almost entirely on a single side of the decision boundary (corresponding to high classifier output agreement within each distribution, and therefore low entropy). The orientation of $w_\beta$ can reduce this term by rotating the data distributions so that their *minor*, not *major* axes cross the classifier, reducing the variance of classifier outputs in $\widehat{X} \mid \alpha$ for each unique $\alpha$. The orientation of $w_\alpha$ can reduce this term by moving the distributions $p(\widehat{X} \mid \alpha)$ away from the decision boundary (where disagreement in corresponding $Y$ values is lower) as quickly as possible as $|\alpha|$ increases.

**Lemma 8.** *Let $W = [w_\alpha \quad W_\beta]$, for $w_\alpha \in \mathbb{R}^N$ and $W_\beta \in \mathbb{R}^{N \times (N-1)}$. Suppose that each column $w_i$ of $W$ is bounded by $c > 0$, i.e., $\|w_i\|_2 \leq c$. Then under the conditions of Proposition 3, $H(Y \mid \alpha)$ is minimized when $w_\alpha = \pm c \frac{a}{\|a\|_2}$ and $W_\beta^T a = 0$.*

*Proof.* We have $\widehat{X} = w_\alpha \alpha + W_\beta \beta + \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \gamma I)$. For fixed $\alpha$, $p(\widehat{X} \mid \alpha) = \mathcal{N}(w_\alpha \alpha, W_\beta W_\beta^T + \gamma I)$. Defining $U = a^T X$, we have $U \mid \alpha \sim \mathcal{N}(\alpha a^T w_\alpha, a^T W_\beta W_\beta^T a + \gamma \|a\|_2^2)$. Then,

$$
\begin{aligned}
p(\widehat{Y} = 1 \mid \alpha) &= \mathbb{E}_{\widehat{X} \mid \alpha}[p(\widehat{Y} = 1 \mid \widehat{X}, \alpha)] \\
&= \mathbb{E}_{\widehat{X} \mid \alpha}[p(\widehat{Y} = 1 \mid \widehat{X})] \\
&= \mathbb{E}_{\widehat{X} \mid \alpha}[\sigma(a^T X)] \\
&= \mathbb{E}_{U \mid \alpha}[\sigma(U)] \\
&\overset{(\star)}{=} \sigma\left( \frac{\alpha \langle a, w_\alpha \rangle}{\sqrt{1 + a^T W_\beta W_\beta^T a + \gamma \|a\|_2^2}} \right),
\end{aligned}
$$

where ($\star$) follows from the fact that for $Z \sim \mathcal{N}(\mu, \sigma^2)$, $\mathbb{E}_Z[\sigma(Z)] = \sigma\left(\frac{\mu}{\sqrt{1+\sigma^2}}\right)$.

We can now evaluate the entropy $H(Y \mid \alpha) = \mathbb{E}_{t \sim \alpha}[H(Y \mid \alpha = t)]$. Again denoting the binary entropy function by $h_b$, we have

$$H(Y \mid \alpha = t) = h_b(p(Y = 1 \mid \alpha = t))$$

$$= h_b(\sigma(s)) \quad \text{where } s := \frac{t\langle a, w_\alpha \rangle}{\sqrt{1 + a^T W_\beta W_\beta^T a + \gamma \|a\|_2^2}}$$

$$= h_b((\sigma(s) - 0.5) + 0.5).$$

Let $q := p - 0.5$ and define $\widetilde{h}_b(q) = h_b(q + 0.5)$ for $q \in [-0.5, 0.5]$ so that $\widetilde{h}_b$ is an even function. Therefore, $\widetilde{h}_b(q) = \widetilde{h}_b(|q|)$, and we have $h_b(p) = \widetilde{h}_b(p - 0.5) = \widetilde{h}_b(|p - 0.5|)$. Applying this fact yields

$$= \widetilde{h}_b(\sigma(s) - 0.5)$$

$$= \widetilde{h}_b(|\sigma(s) - 0.5|)$$

$$\overset{(\dagger)}{=} \widetilde{h}_b(|\sigma(|s|) - 0.5|)$$

where ($\dagger$) follows since $|\sigma(s) - 0.5|$ is an even function of $s$. On $\mathbb{R}_{\geq 0}$ we have that $\widetilde{h}_b(\cdot)$ is a monotonically decreasing function and $|\sigma(\cdot) - 0.5|$ is a monotonically increasing function, and therefore $H(Y \mid \alpha = t) = \widetilde{h}_b(|\sigma(|s|) - 0.5|)$ is a monotonically decreasing function of $|s|$ where

$$|s| = \frac{|t|\,|\langle a, w_\alpha \rangle|}{\sqrt{1 + a^T W_\beta W_\beta^T a + \gamma \|a\|_2^2}}. \tag{7}$$

For any value of $t$, it is clear that the expression in (7) is maximized (and therefore $H(Y \mid \alpha = t)$ is minimized) with respect to $w_\alpha$ and $W_\beta$ when both $|\langle a, w_\alpha \rangle|$ is maximized and $a^T W_\beta W_\beta^T a$ is minimized. By the Cauchy-Schwarz inequality and from boundedness of the column magnitudes of $W$ by $c$, we have that $|\langle a, w_\alpha \rangle|$ is maximized at $w_\alpha = \pm c\frac{a}{\|a\|_2}$. Since $a^T W_\beta W_\beta^T a \geq 0$, this quadratic term is minimized at $W_\beta^T a = 0$ in which case $a^T W_\beta W_\beta^T a = 0$.

Since choosing $w_\alpha$ and $W_\beta$ in this way minimizes $H(Y \mid \alpha = t)$ for any $t$, we have that $H(Y \mid \alpha) = \mathbb{E}_{t \sim \alpha}[H(Y \mid \alpha = t)]$ is also minimized with this choice of $w_\alpha$ and $W_\beta$. $\qquad \square$

Since $H(Y)$ is constant for any $W$ (Lemma 7), we have that the conditions on $W$ described in Lemma 8 maximize $\mathcal{C} = I(\alpha; Y)$. We combine this result with the following lemma to characterize the minimum of the entire objective (3):

**Lemma 9.** *Suppose that $\varepsilon < 1$, $W \in \mathbb{R}^{N \times N}$, and that $X, \varepsilon, z \sim \mathcal{N}(0, I)$ in $\mathbb{R}^N$. With $U = Wz + \gamma\varepsilon$, $\mathrm{D}_{\mathrm{KL}}(p(X) \| p(U))$ is minimized by any orthogonal $W$ with columns normalized to magnitude $\sqrt{(1-\gamma)}$.*

*Proof.* Noting that $U \sim \mathcal{N}(0, WW^T + \gamma I)$, we have from a standard result on KL divergence between multivariate normal distributions that

$$\underset{W}{\arg\min}\, \mathrm{D}_{\mathrm{KL}}(p(X) \| p(U)) = \underset{W}{\arg\min}\, \log\left|WW^T + \gamma I\right| + \mathrm{tr}((WW^T + \gamma I)^{-1}). \tag{8}$$

Since $WW^T + \gamma I$ is positive definite, there exists orthogonal $V$ and diagonal $\Lambda$ with positive entries $\{\lambda_i\}_{i=1}^N$ such that $WW^T + \gamma I = V\Lambda V^T$. We then have

$$\log\left|WW^T + \gamma I\right| + \mathrm{tr}((WW^T + \gamma I)^{-1}) = \log\left|V\Lambda V^T\right| + \mathrm{tr}((V\Lambda V^T)^{-1})$$

$$= \log|\Lambda| + \mathrm{tr}(\Lambda^{-1})$$

$$= \sum_i \log \lambda_i + \frac{1}{\lambda_i}. \tag{9}$$

(9) is minimized at $\lambda_i = 1$ for all $i$. Therefore, the minimizer of (8) is characterized by $WW^T = VV^T - \gamma I = (1 - \gamma)I$. Any orthogonal $W$ with column magnitudes equal to $\sqrt{1 - \gamma}$ satisfies this condition. $\qquad \square$

Combining these lemmas, consider the solution $w_\alpha = \sqrt{1 - \gamma} \frac{a}{\|a\|_2}$, and $W_\beta$ with orthogonal, $\sqrt{1 - \gamma}$-norm columns satisfying $W_\beta^T a = 0$. From Lemma 8 we have that this solution minimizes $H(Y \mid \alpha)$ within the class of $N \times N$ matrices whose column magnitudes are bounded by $\sqrt{1 - \gamma}$. Combined with the invariance of $H(Y)$ to $W$ (Lemma 7), we have that $I(\alpha; Y)$ is maximized by this choice of $w_\alpha$ and $W_\beta$. From Lemma 9 we have that this solution also minimizes $\mathcal{D} = \mathrm{D}_{\mathrm{KL}}(p(X) \parallel p(U))$, and thus this solution minimizes the objective (3) for any $\lambda > 0$.

### C.3  Proof of Proposition 4

Proposition 4 states the relationships between information flow-based objectives depicted graphically in Figure 6.

*Proof of (a).* We have that

$$
\begin{aligned}
I(Y; \alpha) &= \frac{1}{K} \sum_{i=1}^{K} I(Y; \alpha_1, \ldots, \alpha_K) \\
&= \frac{1}{K} \sum_{i=1}^{K} [I(Y; \alpha_i) + I(Y; \alpha_{\neg i} \mid \alpha_i)] \\
&= \frac{1}{K} \sum_{i=1}^{K} I(Y; \alpha_i) + \frac{1}{K} \sum_{i=1}^{K} I(Y; \alpha_{\neg i} \mid \alpha_i).
\end{aligned}
$$

$\square$

*Proof of (b).* First, note that

$$
\begin{aligned}
I(X; Y \mid Z, W) &= \int_x \int_y \int_z \int_w p(x, y, z, w) \log \frac{p(x, y \mid z, w)}{p(x \mid z, w)p(y \mid z, w)} dx\,dy\,dz\,dw \\
&= \int_x \int_y \int_z \int_w p(x, y, z, w) \log \frac{p(x, y, z, w)/p(z, w)}{p(x, z, w)/p(z, w)p(y, z, w)/p(z, w)} dx\,dy\,dz\,dw \\
&= \int_x \int_y \int_z \int_w p(x, y, z, w) \log \frac{p(x, y, z, w)p(z, w)p(x, w)}{p(x, z, w)p(y, z, w)p(x, w)} dx\,dy\,dz\,dw \\
&= \int_x \int_y \int_z \int_w p(x, y, z, w) \log \frac{p(x, y, z \mid w)p(z \mid w)p(x \mid w)}{p(x, z \mid w)p(y, z \mid w)p(x \mid w)} dx\,dy\,dz\,dw \\
&= \int_x \int_y \int_z \int_w p(x, y, z, w) \left( \log \frac{p(x, y, z \mid w)}{p(x \mid w)p(y, z \mid w)} \right. \\
&\qquad \left. - \log \frac{p(x, z \mid w)}{p(x \mid w)p(z \mid w)} \right) dx\,dy\,dz\,dw \\
&= I(X; Y, Z \mid W) - \mathbb{E}_Y \left[ I(X; Z \mid W) \right] \\
&= I(X; Y, Z \mid W) - I(X; Z \mid W).
\end{aligned}
$$

Applying this identity,

$$
\begin{aligned}
I(Y; \alpha \mid \beta) &= \frac{1}{K} \sum_{i=1}^{K} I(Y; \alpha_i, \alpha_{\neg i} \mid \beta) \\
&= \frac{1}{K} \sum_{i=1}^{K} [I(Y; \alpha_i \mid \alpha_{\neg i}, \beta) + I(Y; \alpha_{\neg i} \mid \beta)] \\
&= \frac{1}{K} \sum_{i=1}^{K} I(Y; \alpha_i \mid \alpha_{\neg i}, \beta) + \frac{1}{K} \sum_{i=1}^{K} I(Y; \alpha_{\neg i} \mid \beta).
\end{aligned}
$$

$\square$

*Proof of (c).* We have that

$$I(Y;\alpha \mid \beta) = \int_Y \int_\alpha \int_\beta p(Y,\alpha,\beta) \log \frac{p(Y,\alpha \mid \beta)}{p(Y \mid \beta)p(\alpha \mid \beta)} dY d\alpha d\beta$$

$$\overset{(\star)}{=} \int_Y \int_\alpha \int_\beta p(Y,\alpha,\beta) \log \frac{p(\beta \mid Y,\alpha)p(Y,\alpha)}{p(\beta)} \frac{p(\beta)}{p(\beta \mid Y)p(Y)} \frac{p(\beta)}{p(\beta \mid \alpha)p(\alpha)} dY d\alpha d\beta$$

$$\overset{(\star\star)}{=} \int_Y \int_\alpha \int_\beta p(Y,\alpha,\beta) \log \frac{p(\beta \mid Y,\alpha)p(Y,\alpha)}{p(\beta)} \frac{p(\beta)}{p(\beta \mid Y)p(Y)} \frac{p(\beta)}{p(\beta)p(\alpha)} dY d\alpha d\beta$$

$$= \int_Y \int_\alpha \int_\beta p(Y,\alpha,\beta) \log \frac{p(Y,\alpha)}{p(Y)p(\alpha)} \frac{p(\beta \mid Y,\alpha)}{p(\beta \mid Y)} dY d\alpha d\beta$$

$$= \int_Y \int_\alpha \int_\beta p(Y,\alpha,\beta) \log \frac{p(Y,\alpha)}{p(Y)p(\alpha)} \frac{p(\beta \mid Y,\alpha)p(\alpha \mid Y)}{p(\beta \mid Y)p(\alpha \mid Y)} dY d\alpha d\beta$$

$$= \int_Y \int_\alpha \int_\beta p(Y,\alpha,\beta) \log \frac{p(Y,\alpha)}{p(Y)p(\alpha)} \frac{p(\alpha,\beta \mid Y)}{p(\alpha \mid Y)p(\beta \mid Y)} dY d\alpha d\beta$$

$$= \int_Y \int_\alpha \int_\beta p(Y,\alpha,\beta) \left( \log \frac{p(Y,\alpha)}{p(Y)p(\alpha)} + \log \frac{p(\alpha,\beta \mid Y)}{p(\alpha \mid Y)p(\beta \mid Y)} \right) dY d\alpha d\beta$$

$$= I(Y;\alpha) + I(\alpha;\beta \mid Y),$$

where $(\star)$ follows from Bayes' rule and $(\star\star)$ follows from the independence of $\alpha$ and $\beta$ in our model. $\square$

*Proof of (d).* Similar to (c). $\square$

---
**Algorithm 2** Sample-based estimate of $\mathcal{C}(\alpha; Y)$
---
**Input:** number of samples $N_\alpha$ and $N_\beta$, number of latent factors $K$ and $L$, number of classes $M$

   $I \leftarrow 0$
   $\boldsymbol{q}_y \leftarrow \text{zeros}(M)$
   **for** $i = 1$ **to** $N_\alpha$ **do**
     $\alpha \leftarrow K$-dimensional vector sampled from $\mathcal{N}(0, I)$
     $\boldsymbol{p}_{y|\alpha} \leftarrow \text{zeros}(M)$
     **for** $j = 1$ **to** $N_\beta$ **do**
       $\beta \leftarrow L$-dimensional vector sampled from $\mathcal{N}(0, I)$
       $x \leftarrow$ sample from $p(x \mid \alpha, \beta)$
       $\boldsymbol{p}_{y|\alpha} \leftarrow \boldsymbol{p}_{y|\alpha} + \frac{1}{N_\beta} p(y \mid x)$ (where $p(y \mid x) \in \mathbb{R}^M$ is the classifier probability for each class)
     **end for**
     $I \leftarrow I + \frac{1}{N_\alpha} \sum_{m=1}^M \boldsymbol{p}_{y|\alpha}[m] \log \boldsymbol{p}_{y|\alpha}[m]$
     $\boldsymbol{q}_y \leftarrow \boldsymbol{q}_y + \frac{1}{N_\alpha} \boldsymbol{p}_{y|\alpha}$
   **end for**
   $I \leftarrow I - \sum_{m=1}^M \boldsymbol{q}_y[m] \log \boldsymbol{q}_y[m]$
**Output:** $I$ (sample-based estimate of $I(\alpha; Y)$)

---

## D Sample-based estimate of causal influence

Here we detail the sampling procedure for approximating the causal objective in (2). (The variants described in Appendix A can be approximated in similar fashion.) We have

$$\mathcal{C}(\alpha; Y) = I(\alpha; Y) = \int_\alpha p(\alpha) \left( \sum_y p(y \mid \alpha) \log p(y \mid \alpha) \right) d\alpha - \sum_y p(y) \log p(y)$$

where

$$p(y \mid \alpha) = \int_\beta \int_x p(y \mid x) p(x \mid \alpha, \beta) p(\beta) dx d\beta \qquad (10)$$

and

$$p(y) = \int_{\alpha,\beta} \int_x p(y \mid x) p(x \mid \alpha, \beta) p(\alpha) p(\beta) dx d\alpha d\beta. \qquad (11)$$

For fixed $\alpha$, we approximate (10) with $N_x$ and $N_\beta$ samples of $x$ and $\beta$, respectively, as

$$p(y \mid \alpha) \approx \frac{1}{N_\beta N_x} \sum_{j=1}^{N_\beta} \sum_{n=1}^{N_x} p(y \mid x^{(n)}),$$

where each $x^{(n)} \sim p(x \mid \alpha, \beta^{(j)})$ and $\beta^{(j)} \sim p(\beta)$. Similarly, we approximate (11) with $N_x$, $N_\alpha$, and $N_\beta$ samples of $x$, $\alpha$, and $\beta$, respectively, as

$$p(y) \approx \frac{1}{N_\alpha N_\beta N_x} \sum_{j=1}^{N_\beta} \sum_{i=1}^{N_\alpha} \sum_{n=1}^{N_x} p(y \mid x^{(n)}),$$

where each $x^{(n)} \sim p(x \mid \alpha^{(i)}, \beta^{(j)})$, $\alpha^{(i)} \sim p(\alpha)$, and $\beta^{(j)} \sim p(\beta)$. Therefore,

$$I(\alpha_i; y) \approx \frac{1}{N_\alpha N_\beta N_x} \left[ \sum_{i=1}^{N_\alpha} \sum_y \left( \sum_{j=1}^{N_\beta} \sum_{n=1}^{N_x} p(y \mid x^{(n)}) \right) \log \left( \frac{1}{N_\beta N_x} \sum_{j=1}^{N_\beta} \sum_{n=1}^{N_x} p(y \mid x^{(n)}) \right) \right.$$
$$\left. - \sum_y \left( \sum_{j=1}^{N_\beta} \sum_{i=1}^{N_\alpha} \sum_{n=1}^{N_x} p(y \mid x^{(n)}) \log \left( \frac{1}{N_\alpha N_\beta N_x} \sum_{j=1}^{N_\beta} \sum_{i=1}^{N_\alpha} \sum_{n=1}^{N_x} p(y \mid x^{(n)}) \right) \right) \right]$$

where each $x^{(n)} \sim p(x \mid \alpha^{(i)}, \beta^{(j)})$, $\alpha^{(i)} \sim p(\alpha)$, and $\beta^{(j)} \sim p(\beta)$.

The complete procedure is described algorithmically in Algorithm 2 with $N_x = 1$.

| Classifier Architecture |
| --- |
| Input (28×28) |
| Conv2 (32 channels, 3×3 kernels, stride 1, pad 0) |
| ReLU |
| Conv2 (64 channels, 3×3 kernels, stride 1, pad 0) |
| ReLU |
| MaxPool (2×2 kernel) |
| Dropout ($p = 0.5$) |
| Linear (128 units) |
| ReLU |
| Dropout ($p = 0.5$) |
| Linear ($M$ units) |
| Softmax |

Table 1: Network architecture for MNIST Classifier



Figure 11: Partial details of parameter tuning procedure used to select $K$, $L$, and $\lambda$ for explaining MNIST 3/8 classifier using Algorithm 1. *Left:* In Step 1 we select the total number of latent factors $K + L$ needed to adequately represent the data distribution. *Center:* In Steps 2-3 we iteratively convert noncausal latent factors to causal latent factors until $\mathcal{C}$ plateaus. *Right:* After each increment of $K$, we adjust $\lambda$ to approximately achieve the value of $\mathcal{D}$ from Step 1.

# E    VAE experimental details and additional results

## E.1    Details and additional results for MNIST experiments

All experiments were run using a single Nvidia GeForce GTX 1080 GPU. The traditional MNIST training set was split into training and validation sets composed of the first 50,000 and remaining 10,000 images, respectively. The testing set was the same as the traditional MNIST testing set, composed of 10,000 images. These sets were down-selected to include only samples with the labels of interest. Input images were scaled so that the network inputs are in $[0, 1]^{28 \times 28}$.

The network architecture for the classifier used in the MNIST experiments is shown in Table 1 where $M$, the number of class outputs, varies depending on the classification task. The classifier was trained with a batch size of 64 and a stochastic gradient descent optimizer with momentum 0.5 and learning rate 0.1. The 3/8 classifier was trained for 20 epochs and the 1/4/9 classifier was trained for 30 epochs. The test accuracy of the classifier trained on both the 3/8 and 1/4/9 datasets was 99.6%.

The VAE architecture used to learn the generative map $g$ is shown in Table 2. The objective (3) was maximized with 8000 training steps, batch size 64, and learning rate $5 \times 10^{-4}$. At each training step, the causal influence term 2 was estimated using the sampling procedure in Appendix D with $N_\alpha = 100$ and $N_\beta = 25$. For experiments with digits 3 and 8, we selected $K = 1$, $L = 7$, and $\lambda = 0.05$ using the parameter selection procedure in Algorithm 1; Figure 11 shows intermediate results from this procedure.

| VAE Encoder Architecture | VAE Decoder Architecture |
|---|---|
| Input (28×28) | Input ($K + L$) |
| Conv2 (64 chan., 4×4 kernels, stride 2, pad 1) | Linear (3136 units) |
| ReLU | ReLU |
| Conv2 (64 chan., 4×4 kernels, stride 2, pad 1) | Conv2Transp (64 chan., 4×4 kernels, stride 1, pad 1) |
| ReLU | ReLU |
| Conv2 (64 chan., 4×4 kernels, stride 1, pad 0) | Conv2Transp (64 chan., 4×4 kernels, stride 2, pad 2) |
| ReLU | ReLU |
| Linear ($K + L$ units for both $\mu$ and $\sigma$) | Conv2Transp (1 chan., 4×4 kernel, stride 2, pad 1) |
| | Sigmoid |

Table 2: VAE network architecture used for MNIST and Fashion MNIST experiments.



Figure 12: Visualizations for learned latent factors for MNIST 3/8 classifier. Images in the center column of each grid are reconstructed samples from the validation set; moving left or right in each row shows $g(\alpha, \beta)$ as a single latent factor is varied. This plot shows the complete results from Figure 3; it includes sweeps for two additional samples and visualizations of all $L = 7$ noncausal factors.

Figure 12 shows additional results for the experiment of Figure 3, which visualizes the learned latent factors that explain the MNIST 3/8 classifier. Here we show latent factor sweeps from this experiment with additional data samples and all $K + L = 8$ latent factors.

Figure 13 shows an explanation of the same classifier architecture detailed in Table 1 trained on the MNIST digits 1, 4, and 9. We use the VAE architecture of Table 2 with $K = 2$ causal factors, $L = 2$ noncausal factors, and $\lambda = 0.1$, and estimated the causal influence portion of the objective using the sampling procedure in Appendix D with $N_\alpha = 75$ and $N_\beta = 25$. While the factor sweeps in Figure 13 provide a high-level indication of the data features each factor corresponds to, a practitioner may also wish to visualize the fine-grained transitions between each class. This can be achieved by sweeping each factor on a finer scale, as visualized by the zoomed in regions of Figure 13 as well as the more comprehensive sweeps in Figure 14.
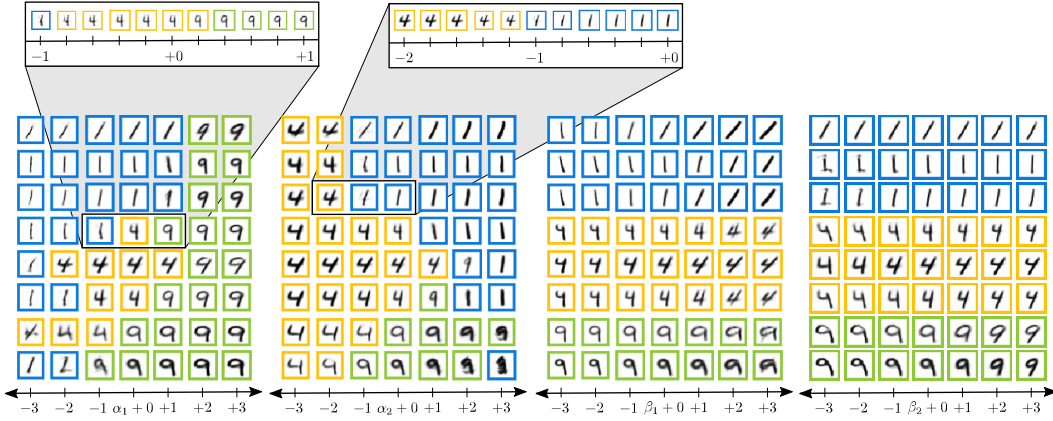
Figure 13: Visualizations of learned latent factors for MNIST 1/4/9 classifier. Images in the center column of each grid are reconstructed samples from the validation set; moving left or right in each row shows $g(\alpha, \beta)$ as a single latent factor is varied. Varying the causal factors $\alpha_1$ and $\alpha_2$ control aspects that affect the classifier output (colored borders); varying the noncausal factors $\beta_1$ and $\beta_2$ affect only stylistic aspects such as rotation and thickness.



Figure 14: High-resolution transition regions of the first causal factor in explaining the MNIST 1/4/9 classifier. Visualizing high-resolution latent factor sweeps can allow a practitioner to more easily identify which data features correspond to each underlying factor. For example, one can observe in the second row from the bottom how increasing $\alpha_1$ causes the left branch of the digit '4' to smoothly transition into completing the loop of the digit '9' while the digit stem remains fixed.

## E.2   Details and additional results for comparison experiments

Figure 4 compares our latent factor-based local explanations to the local explanations of four popular explanation methods. We generate explanations of the same CNN classifier trained on MNIST 3 and 8 digits described in Appendix E.1. The data samples explained in Figure 4 are the first example of each class in the MNIST validation set.

**Implementation details of other methods.** The following procedures were used to generate the results for LIME, DeepSHAP, IG, and L2X shown in Figure 4 (left):

- **LIME [17].** The LIME framework trains a sparse linear model using superpixel features. Following the recommendation in the authors' code, we generate superpixels using the Quickshift segmentation algorithm from scikit-image with kernel size 1, maximum distance 200, and color/image-space proximity ratio 0 (as the MNIST digits are grayscale). The LIME local approximation is fit using the default kernel width of 0.25, 10,000 samples, and $K = 10$ features. Figure 4 show superpixels identified as contributing positively (red) and or negatively (blue) to the classification decision.

- **DeepSHAP [25].** The DeepSHAP method uses the structure of the classifier network to efficiently approximate Shapley values, a game-theoretic formulation for how to optimally distribute rewards to players of a cooperative game. The Shapley values displayed in Figure
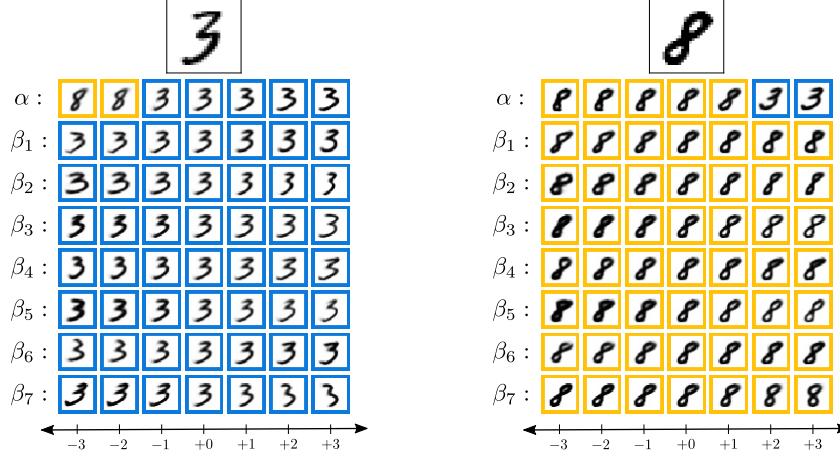
29

Figure 15: Complete results for local explanations of '3' and '8' from Figure 4. Our explanations are able to differentiate causal aspects (pixels that define 3 from 8) from purely stylistic aspects. Only the causal factor $\alpha$ controls changes in data space that result in a change in classifier output.

4 can be interpreted as the (averaged) importance of each pixel for explaining the difference between $f(x)$ and $\mathbb{E}_{x \sim X}[f(x)]$. We train the explanation model using 1000 randomly chosen samples from the training set. The DeepSHAP method produces explanations for each possible class; we display the Shapley values corresponding to the classifier class (i.e., the top image shows the explanation for ground truth class 3 and the bottom image shows the explanation for ground truth class 8).

- **IG [24].** The integrated gradients (IG) method integrates the gradient of the classifier probabilities with respect to the input as the input changes from a "baseline." We use an all-zero image as the baseline and the trapezoid rule with 50 steps to approximate the integral. The output in Figure 4 shows the integrated gradient explanation for each input image.

- **L2X [11].** The learning to explain (L2X) algorithm learns a mask of features $S$ that (approximately) maximizes $I(Y; X \odot S)$. Following [11, Sec. 4.3], we find a mask with $k = 4$ active superpixels, each of size $4 \times 4$. The neural network parameterizing the "explainer" model $p(S \mid X)$ consists of two convolutional layers (32 filters of size $2 \times 2$ each with relu activation, each followed by a max pooling layer with a $2 \times 2$ pool size), followed by a single $2 \times 2$ convolutional filter. This explainer network learns a $7 \times 7$ mask, with each element corresponding to a $2 \times 2$ superpixel in data space. The neural network parameterizing the variational bound $q(Y \mid X \odot S)$ consists of two convolutional layers, each containing 32 filters of size $2 \times 2$, using relu activation, and followed by a max pooling layer with $2 \times 2$ pool size; followed by a dense layer. The networks parameterizing $p(S \mid X)$ and $q(Y \mid X \odot S)$ were trained together with 10 epochs of the 9943 MNIST training samples of 3's and 8's and the outputs $Y$ of the convolutional neural network classifier described in Appendix E.1.

**Complete results for our method.** In Figure 4 (right) we show only latent factor sweeps for the the causal factor $\alpha$ and a single noncausal factor $\beta_7$. Figure 15 shows complete local explanations with each noncausal factor. Our explanations use the VAE framework described in Appendix E.1.

### E.3 Details and additional results for fashion MNIST experiments

Our training set was the same as the traditional Fashion MNIST training set, composed of 60,000 images. The Fashion MNIST testing set was split into validation and testing sets composed of the first 6,000 and last 4,000 images, respectively. These sets were down-selected to include only samples with the labels of interest — in our experiment, classes 0 ('t-shirt/top'), 3 ('dress'), and 4 ('coat'). Input images were scaled so that the input images were in $[0, 1]^{28 \times 28}$.
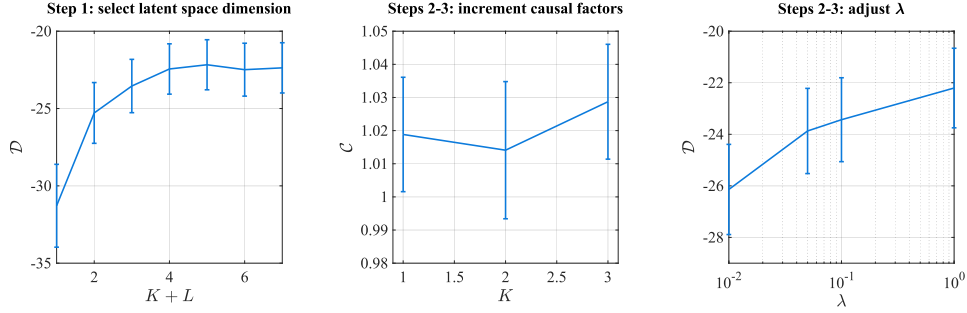
30

Figure 16: Partial details of parameter tuning procedure used to select $K$, $L$, and $\lambda$ for explaining a classifier trained on classes 0, 3, and 4 of the fashion MNIST dataset using Algorithm 1. *Left:* in Step 1 we select the total number of latent factors $K + L$ needed to adequately represent the data distribution. *Center:* In Steps 2-3 we iteratively convert noncausal latent factors to causal latent factors until $\mathcal{C}$ (shown in nats) plateaus. *Right:* After each increment of $K$, we adjust $\lambda$ to approximately achieve the value of $\mathcal{D}$ from Step 1.

The same classifier architecture described in Table 1 was used in this experiment. The classifier was trained with 50 epochs, a batch size of 64, a stochastic gradient descent optimizer with momentum 0.5 and learning rate 0.1. Because the classes used ('t-shirt/top,' 'dress,' and 'coat') are similar, this classifier task is more challenging than the MNIST digit classification task; the test accuracy of the classifier was 95.2%.

The same VAE architecture described in Table 2 was used to learn the generative map $g$. The objective (3) was maximized with 8000 training steps, batch size 32, and learning rate $10^{-4}$. At each training step, the causal influence term (2) was estimated using the sampling procedure in Appendix D with $N_\alpha = 100$ and $N_\beta = 25$. Using the parameter selection procedure in Algorithm 1, we selected $K = 2$, $L = 4$, and $\lambda = 0.05$; Figure 16 shows intermediate results from this procedure.

Figure 17 contains the complete results from the experiment in Figure 5 (right), showing a complete visualization of the global explanation learned for this classifier.
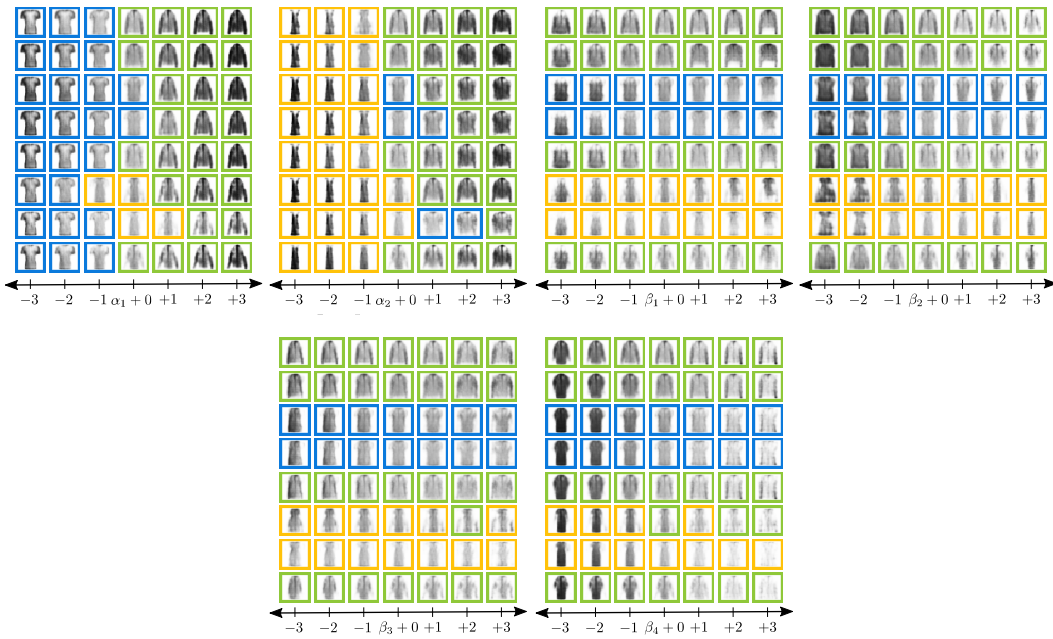
Figure 17: Visualizations of learned latent factors for Fashion MNIST classifier trained on classes 't-shirt-top,' 'dress,' and 'coat.' Images in the center column column of each grid are reconstructed samples from the validation set; moving left or right in each row shows $g(\alpha, \beta)$ as a single latent factor is varied. This plot shows the complete results from Figure 5 (right); it includes sweeps for two additional samples and visualizations of all $K + L = 6$ latent factors.

# F   Selecting generative model capacity

One practical decision to make when constructing explanations using our method is selecting the *capacity* of the generative model $g$. Set too low, the generative model will have insufficient capacity to represent the data distribution and classifier, reducing the quality of the explanation. Set too high, the generative model will require a more time- and energy-intensive training procedure.

We can use results from [75] to bound the capacity mismatch of our explainer (i.e., explainer error in predicting classifier outputs) with the $I(\alpha; Y)$ part of our objective. In practice, this result means that a sufficiently large value of $I(\alpha; Y)$ serves as a certificate that the explainer complexity is sufficient to explain the classifier. Below, we show details of this analysis and empirically demonstrate how $I(\alpha; Y)$ can be used to select an architecture with sufficient capacity.

## F.1   $I(\alpha; Y)$ serves as a certificate of sufficient explainer capacity

One reasonable measure for the quality of an explanation method is how accurately the black-box's classifications can be predicted from the explanation alone. If this prediction is accurate, then in a predictive sense the explanation has captured the relevant information about the classifier's behavior. In our model, the estimator that minimizes prediction error is the MAP estimate of the classifier's output from $p(Y \mid \alpha)$, where $p(Y \mid \alpha)$ is determined by marginalizing $p(Y \mid X)p(X \mid \alpha, \beta)p(\beta)$ over $\beta$ and $X$. As we show below, we can upper bound the error of this MAP estimator *directly* by the causal effect $I(\alpha; Y)$ of $\alpha$ on $Y$, the quantity our method explicitly optimizes.

Specifically, let $\pi(Y \mid \alpha) := \int_\alpha [1 - \max_y p(y \mid \alpha)] p(\alpha)d\alpha$ denote the expected error of this MAP estimator, averaged over the prior distribution on causal factor $\alpha$. From [75], we have

$$\phi^*(\pi(Y \mid \alpha)) \leq H(Y \mid \alpha),$$

where $H(Y \mid \alpha)$ is the conditional entropy of $Y$ given $\alpha$, and $\phi^*$ is a monotonically increasing, invertible function. Define $\widetilde{\phi} = (\phi^*)^{-1}$. Since $H(Y \mid \alpha) = H(Y) - I(Y; \alpha) \leq \log M - I(Y; \alpha)$ [83], we have

$$\pi(Y \mid \alpha) \leq \widetilde{\phi}(\log_2 M - I(Y; \alpha)) \tag{12}$$

where $I(Y; \alpha)$ is measured in bits.

If we take the prediction error of $Y$ from $\alpha$ as a measure of "mismatch" between our trained model and the blackbox classifier, (12) bounds this mismatch by the causal effect term in our objective and can serve as a certificate for having sufficient network capacity. For example, in 3-class Fashion MNIST ($M = 3$), a value of $I(\alpha; Y) = 1.03$ nats as in Figure 16 results in a bound of $\pi(Y \mid \alpha) \leq 0.05$. This translates to a MAP estimator of $Y$ from $\alpha$ having a black-box output prediction error of less than $5\%$, or that the causal factors can explain at least $95\%$ of the black-box's behavior. If this prediction accuracy is satisfactory, then the capacity of the generator $g$ is sufficient to learn appropriate latent factors and their mapping to the data space. If this prediction accuracy is not satisfactory, a class $G$ of generative models $g$ with higher capacity can be used. This will provide the model with more flexibility to optimize $I(\alpha; Y)$ and reduce prediction error.

## F.2   Empirical results

The drawback of a VAE with insufficient capacity can be seen in Figure 18, which shows the causal effect and data fidelity terms of the objective (3) as the VAE capacity and tuning parameter $\lambda$ are modified. The VAE in each trial, which is applied to explain the 3-class Fashion-MNIST classifier considered in the quantitative experiments of Section 5, uses the architecture described in Table 2 with $K = 2$ and $L = 4$ but with a variable number of convolutional filters in each layer of the encoder and decoder (see Table 3). The values of $\mathcal{C}$ and $\mathcal{D}$ reported in Figure 18 are the average values in the last 50 training steps for each model. The dotted line in Figure 18 represents the maximum achievable value of $I(\alpha; Y)$ in this three class setting, $\log(3) \approx 1.1$ nats.

As discussed in Section 3.4, the tuning parameter $\lambda$ dictates the trade-off between the objective's causal effect term $\mathcal{C}$ and data fidelity term $\mathcal{D}$. When the number of filters per layer is small, however, the model has insufficient capacity to simultaneously achieve a satisfactory value of both $\mathcal{C}$ and $\mathcal{D}$.

Figure 19 shows partial resulting explanations generated by an explainer with insufficient capacity (8 filters per convolutional layer; Figure 19(a–b)). Although the causal and noncausal factors do
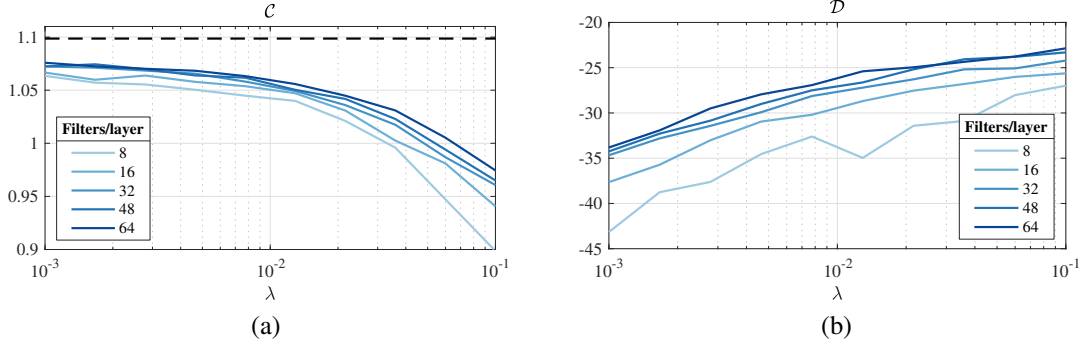
Figure 18: Post-training value of the (a) causal effect and (b) data fidelity terms in the objective (3) for various capacities of VAE. The capacity is modified by changing the number of convolutional filters in each layer.

| Filters per convolutional layer | Encoder parameters | Decoder parameters |
|:---:|:---:|:---:|
| 8 | 6,916 | 4,937 |
| 16 | 17,916 | 13,969 |
| 32 | 52,204 | 44,321 |
| 48 | 102,876 | 91,057 |
| 64 | 169,932 | 154,177 |

Table 3: Number of VAE parameters when $K + L = 6$.

indeed roughly correspond to classifier-relevant and classifier-irrelevant data aspects in the sense that changing $\alpha_1$, but not $\beta_1$, produces changes in the classifier output, the effect of the model's limited ability to represent the data distribution is evident in the weak correspondence of the generated samples to training samples. Meanwhile, the same explanation generated by an explainer with sufficient capacity (64 filters per convolutional layer; Figure 19(c–d)) shows both effectively disentangled classifier-relevant/irrelevant data aspects and generated samples that appear to lie in the training data distribution.
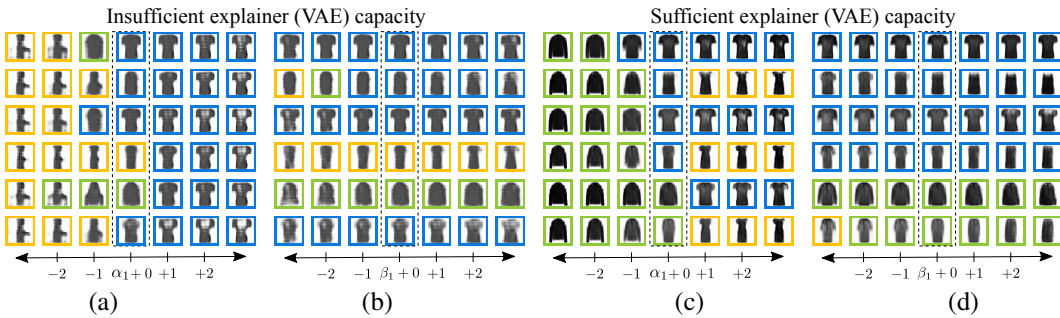


Figure 19: Global explanations with $\lambda \approx 0.013$ and varying VAE model capacity. (a–b) 8 filters per convolutional layer, defining a VAE with insufficient capacity to represent the data distribution. (c–d) 64 filters per convolutional layer, defining a VAE with sufficient capacity to represent the data distribution.