# LABEL-FREE EXPLAINABILITY FOR UNSUPERVISED MODELS

**Jonathan Crabbé**
DAMTP
University of Cambridge
jc2133@cam.ac.uk

**Mihaela van der Scaar**
University of Cambridge
The Alan Turing Institute
UCLA
mv472@cam.ac.uk

## ABSTRACT

Unsupervised black-box models are challenging to interpret. Indeed, most existing explainability methods require labels to select which component(s) of the black-box's output to interpret. In the absence of labels, black-box outputs often are representation vectors whose components do not correspond to any meaningful quantity. Hence, choosing which component(s) to interpret in a label-free unsupervised/self-supervised setting is an important, yet unsolved problem. To bridge this gap in the literature, we introduce two crucial extensions of post-hoc explanation techniques: (1) label-free feature importance and (2) label-free example importance that respectively highlight influential features and training examples for a black-box to construct representations at inference time. We demonstrate that our extensions can be successfully implemented as simple wrappers around many existing feature and example importance methods. We illustrate the utility of our label-free explainability paradigm through a qualitative and quantitative comparison of representation spaces learned by various autoencoders trained on distinct unsupervised tasks.

## 1 Introduction

Are machine learning models ready to be deployed in high-stakes applications? Recent years have witnessed a success of deep models on nontrivial tasks such as computer vision Voulodimos et al. [2018], natural language processing Young et al. [2018] and scientific discovery Jumper et al. [2021], Davies et al. [2021]. The success of these models comes at the cost of their complexity. Deep models typically involve millions to billions operations in order to turn their input data into a prediction. Since it is not possible for a human user to analyze each of these operations, the models appear as a *black-boxes*. When the deployment of these models impact critical areas such as healthcare, finance or justice, their opacity appears as a major obstruction Lipton [2016], Ching et al. [2018], Tjoa and Guan [2020].

**Post-Hoc Explainability** As a response to this transparency problem, the field of *explainable artificial intelligence* (XAI) received an increasing interest, see Adadi and Berrada [2018], Barredo Arrieta et al. [2020], Das and Rad [2020] for reviews. In trying to achieve transparency by retaining the approximation power of deep models, many *post-hoc explainability* methods were developed. These methods aim at complementing the predictions of black-box models with various explanations. In this way, models can be understood through the lens of explanations regardless their complexity. In this work, we focus on two types of such methods. (**1**) **Feature Importance** explanations highlight features that are crucial for the black-box to issue a prediction. Typical feature importance methods are Saliency Simonyan et al. [2013], Lime Ribeiro et al. [2016], Integrated Gradients Sundararajan et al. [2017], Shap Lundberg and Lee [2017] and DeepLift Shrikumar et al. [2017]. (**2**) **Example Importance** explanations highlight training examples that are crucial for the black-box to issue a prediction. Typical example importance methods are Influence Function Cook and Weisenberg [1982], Koh and Liang [2017], Deep K-Nearest Neighbours Papernot and McDaniel [2018], TraceIn Pruthi et al. [2020] and SimplEx Crabbé et al. [2021].

**Supervised Setting** The previous works almost exclusively focus on explaining models obtained in a *supervised setting*. In this setting, the black box $f : \mathcal{X} \to \mathcal{Y}$ connects two meaningful spaces: the feature space $\mathcal{X}$ and the label space $\mathcal{Y}$. By meaningful, we mean that each axis of these spaces corresponds to a quantity known by the model's user. This is illustrated at the top of Figure 1 with an idealized prostate cancer risk predictor. Each axis of $\mathcal{Y}$ corresponds to a clear

label: the probability of mortality with and without treatment. We can explain the predictions of this model with the importance of features/examples to make a prediction on each individual axis $y_1, y_2$ or for a combination of those axes (e.g in this case $y_1 - y_2$ is associated to the treatment effect). The key point is that, in the supervised setting, the user knows the meaning of the black-box output they try to interpret. This is not always the case in machine learning.

**Label-Free Setting** In the label-free setting, we are interested in black-boxes $\boldsymbol{f} : \mathcal{X} \to \mathcal{H}$ that connect a meaningful feature space $\mathcal{X}$ to a latent (or representation) space $\mathcal{H}$. Unlike the feature space $\mathcal{X}$ and the label space $\mathcal{Y}$, the axes of the latent space $\mathcal{H}$ do not correspond to labelled quantities known by the user. This is illustrated at the bottom of Figure 1, where the examples are mapped in a representation space for clustering purposes. Unlike in the supervised setting, there is no obvious way for the user to choose an axis among $h_1, h_2$ to interpret. This distinction goes well beyond the philosophical consideration. As we show in Sections 2 and 3, the aforementioned feature and example importance methods require labels to select an axis to interpret or evaluate a loss. The usage of these methods in the label-free setting hence requires a non-trivial extension.
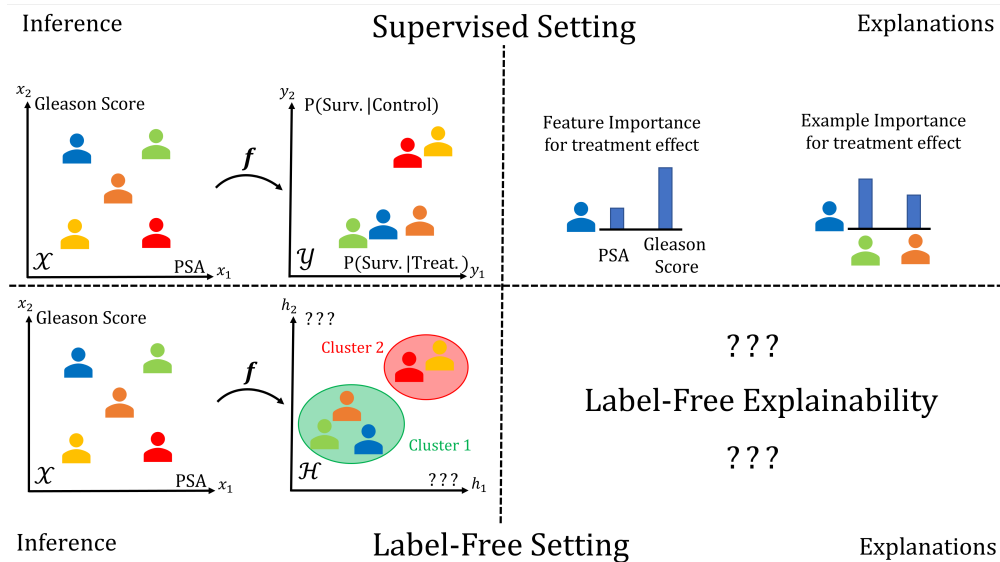


Figure 1: Supervised and Label-Free Settings.

**Motivation** To illustrate the significance of the label-free extension of explainability, we outline 2 widespread setups where it is required. Note that these 2 setups are not mutually exclusive. **(1) Unsupervised Learning:** Whenever we solve an unsupervised task such as clustering or density estimation, it is common to learn a function like $\boldsymbol{f}$ that projects the data onto a low dimensional representation space. Due to the unsupervised nature of the problem, any explanation of the latent space has to be done without label. **(2) Self-Supervised Learning:** Even in a supervised learning setting, we have often more unlabelled than labelled data available. When this is the case, self-supervised learning preconises to leverage unlabelled data to solve an unsupervised pretext tasks such as predicting the missing part of a masked image Jing and Tian [2021]. This yields a representation $\boldsymbol{f}$ of the data that we can use as part of a model to solve the downstream task. If we want to interpret the model's representations of unlabelled examples, label-free explainability is a must. Let us now review related work.

**Related Work** The majority of the relevant literature focuses on increasing the interpretability of representations spaces $\mathcal{H}$. Disentangled-VAEs constitute the best example Higgins et al. [2017], Burgess et al. [2018], Chen et al. [2018], Sarhan et al. [2019], we discuss them in more details in Section 4.3. When it comes to post-hoc approaches for explainability of latent spaces, *concept-based explanations* Kim et al. [2017], Brocki and Chung [2019] are central. These methods identify a dictionary between human concepts (like the presence of stripes on an image) and latent vectors. We note that these methods are only partially relevant here since concepts are typically learned with labelled examples, although some recent works challenge this assumption Ghorbani et al. [2019]. To the best of our knowledge, no existing work extends feature importance methods to the label-free setting. Regarding example importance, Kong and Chaudhuri [2021] adapt TraceIn for VAEs.

**Contribution (1) Label-Free Feature Importance:** We extend *linear* feature importance methods to the label-free setting (Section 2). Our extension is done by defining an auxiliary scalar function as a wrapper around the label-free black-box to interpret. This permits to compute feature importance in the label-free setting by retaining useful properties

of the original methods without increasing their complexity. **(2) Label-Free Example Importance:** We extend example importance methods to the label-free setting (Section 3). In this work, we treat two types of example importance methods that we call *loss-based* and *representation-based*. For the former, the extension requires to specify a label-free loss and a set of relevant model parameters to differentiate the loss with. For the latter, the extension is straightforward. Our feature and example importance extensions are validated experimentally (Section 4.1) and their practical utility is demonstrated with a use case motivated by self-supervised learning (Section 4.2). **(3) Challenging Interpretability of Disentangled Representations:** In testing the limits of our feature importance hypotheses with disentangled VAEs, we noticed that the interpretability of saliency maps associated to individual latent units seems unrelated to the strength of disentanglement between the units (Section 4.3). We analyze this phenomenon both qualitatively and quantitatively. This insight could be the seed of future developments in interpretable representation learning.

## 2  Feature Importance

In this section, we present our approach to extend linear feature importance methods to the label-free setting. We start by reviewing the typical setup with label to grasp some useful insights for our extension. With these insights, the extension to the label-free regime is immediate.

### 2.1  Feature Importance with Labels

We consider an input (or feature) space $\mathcal{X} \subset \mathbb{R}^{d_X}$ and an output (or label) space $\mathcal{Y} \subset \mathbb{R}^{d_Y}$, where $d_X \in \mathbb{N}^*$ and $d_Y \in \mathbb{N}^*$ are respectively the dimension of the input and output spaces. We are given a black-box model $\boldsymbol{f} : \mathcal{X} \to \mathcal{Y}$ that maps each input $\boldsymbol{x} \in \mathcal{X}$ to an output $\boldsymbol{y} = \boldsymbol{f}(\boldsymbol{x}) \in \mathcal{Y}$. Note that we use bold symbols to emphasize that those elements are typically vectors with more than one component ($d_X, d_Y > 1$). Feature importance methods explain the black-box prediction $\boldsymbol{f}(\boldsymbol{x})$ by attributing an importance score $a_i(\boldsymbol{f}, \boldsymbol{x})$ to each feature $x_i$ of $\boldsymbol{x}$ for[1] $i \in [d_X]$. Note that feature importance methods require to select one component $f_j(\boldsymbol{x}) \in \mathbb{R}$ for some $j \in [d_Y]$ of the output in order to compute these scores: $a_i(\boldsymbol{f}, \boldsymbol{x}) \equiv a_i(f_j, \boldsymbol{x})$. In a classification setting, $j$ typically corresponds to the ground-truth label (when it is known) $j = \arg\max_{k \in [d_Y]} [y_k]$ or to the label with maximal predicted probability $j = \arg\max_{k \in [d_Y]} [f_k(\boldsymbol{x})]$.

We now suggest an alternative approach: combining the importance scores for each component of $\boldsymbol{f}$ by weighting them with the associated class probability: $b_i(\boldsymbol{f}, \boldsymbol{x}) \equiv \sum_{j=1}^{d_y} f_j(\boldsymbol{x}) \cdot a_i(f_j, \boldsymbol{x})$. We note that, when a class probability $f_j(\boldsymbol{x}) \approx 1$ dominates, this reduces to the previous definition. However, when the class probabilities are balanced, this accounts for the contribution of each class. In the image classification context, this might be more appropriate than cherry-picking the saliency map of the appropriate class while disregarding the others Rudin [2019]. To the best of our knowledge, this approach has not been used in the literature. A likely reason for this is that this definition requires to compute $d_Y \cdot d_X$ importance scores per sample, which quickly becomes expensive as the number of classes grows. However, this limitation can easily be avoided when the importance scores are linear with respect to the black-box[2]. In this case, the weighted importance score can be rewritten as $b_i(\boldsymbol{f}, \boldsymbol{x}) = a_i(\sum_{j=1}^{d_Y} f_j(\boldsymbol{x}) \cdot f_j, \boldsymbol{x})$. With this trick, we can compute the weighted importance score by only calling the auxiliary function $g_{\boldsymbol{x}}(\tilde{\boldsymbol{x}}) = \sum_{j=1}^{d_Y} f_j(\boldsymbol{x}) \cdot f_j(\tilde{\boldsymbol{x}})$. We will now use the same reasoning to generalize feature importance scores to the label-free setting.

### 2.2  Label-Free Feature Importance

We now turn our setting of interest. We consider a latent (or representation) space $\mathcal{H} \subset \mathbb{R}^{d_H}$, where $d_H \in \mathbb{N}^*$ is the dimension of the latent space. We are given a black-box model $\boldsymbol{f} : \mathcal{X} \to \mathcal{H}$ that maps each input $\boldsymbol{x} \in \mathcal{X}$ to a representation $\boldsymbol{h} = \boldsymbol{f}(\boldsymbol{x}) \in \mathcal{H}$. As aforementioned, the latent space dimensions $h_j, j \in [d_H]$ are not related to labels with clear interpretations. Again, we would like to attribute an importance score $b_i(\boldsymbol{f}, \boldsymbol{x})$ to each feature $x_i$ of $\boldsymbol{x}$ for $i \in [d_X]$. Ideally, this score should reflect the importance of feature $x_i$ in constructing the latent representation $\boldsymbol{h} = \boldsymbol{f}(\boldsymbol{x})$. Unlike in the previous setting, we do not have a principled way of choosing a component $f_j$ for some $j \in [d_H]$. How can we compute these importance scores?

We can simply mimic the approach described in the previous section. For a feature importance method, we use the weighted importance score $b_i(\boldsymbol{f}, \boldsymbol{x}) = a_i(\sum_{j=1}^{d_H} f_j(\boldsymbol{x}) \cdot f_j, \boldsymbol{x})$. At this stage, we should stress that the individual components $f_j(\boldsymbol{x})$ do not correspond to probabilities in this case. Does it still make sense to compute a sum weighted by these components? In most cases, it does. Let us remind that in all the settings described above, the components will typically correspond to a neuron's activation function Glorot et al. [2011]. With typical activation functions such as

---

[1] We denote by $[N]$ the positive integers from 1 to $N \in \mathbb{N}^*$.
[2] Which is the case for most methods, including Shap.

ReLU or Sigmoid, inactive neurons will correspond to a vanishing component $f_j(\boldsymbol{x}) = 0$. From the above formula, this implies that these neurons will not contribute in the computation of $b_i(\boldsymbol{f}, \boldsymbol{x})$. Similarly, neurons that are more activated will have a bigger contribution in the weighted sum of $b_i(\boldsymbol{f}, \boldsymbol{x})$. By linearity of the feature importance method, this reasoning extends to linear combinations of neurons. Since the weighted sum can be interpreted as a latent space inner product, we introduce the following definition.

**Definition 2.1** (Label-Free Feature Importance). Let $\boldsymbol{f} : \mathcal{X} \to \mathcal{H}$ be a black-box latent map and for all $i \in [d_X]$ a $a_i(\cdot, \cdot) : \mathbb{R}^{\mathcal{X}} \times \mathcal{X} \to \mathbb{R}$ be a feature importance score linear w.r.t. its first argument[3]. We define the *label-free* feature importance as a score $b_i(\cdot, \cdot) : \mathcal{H}^{\mathcal{X}} \times \mathcal{X} \to \mathbb{R}$ such that

$$b_i(\boldsymbol{f}, \boldsymbol{x}) = a_i(g_{\boldsymbol{x}}, \boldsymbol{x}) \tag{1}$$

$$g_{\boldsymbol{x}} : \mathcal{X} \to \mathbb{R} \text{ s.t. } g_{\boldsymbol{x}}(\tilde{\boldsymbol{x}}) = \langle \boldsymbol{f}(\boldsymbol{x}), \boldsymbol{f}(\tilde{\boldsymbol{x}}) \rangle_{\mathcal{H}}, \tag{2}$$

where $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ denotes an inner product for the space $\mathcal{H}$.

*Remark* 2.2. This definition gives a simple recipe to extend any linear feature importance method $a_i$ to the label-free setting. In practice, this is implemented by defining the auxiliary function $g_{\boldsymbol{x}}$ as a simple wrapper around the black-box function $\boldsymbol{f}$.

Arguably one of the most important property shared by many feature importance methods is *completeness*. Feature importance methods endowed with this property produce importance scores whose sum equals the black-box prediction up to a constant baseline $a_0 \in \mathbb{R}$: $\sum_{i=1}^{d_x} a_i(g, \boldsymbol{x}) = g(\boldsymbol{x}) - a_0$. This provides a meaningful connection between importance scores and black-box predictions. Typical examples of baselines are $a_0 = 0$ for Lime; the expected prediction $a_0 = \mathbb{E}_{\boldsymbol{X}}[g(\boldsymbol{X})]$ for Shap and a baseline prediction $a_0 = g(\bar{\boldsymbol{x}})$ for Integrated Gradients. We show that our label-free feature importance scores are endowed with an analogous property in higher dimension.

**Proposition 2.3** (Label-Free Completeness). *If the feature importance scores $a_i(\cdot, \cdot) : \mathbb{R}^{\mathcal{X}} \times \mathcal{X} \to \mathbb{R}$ are linear and satisfy completeness, then the label-free importance scores $b_i(\boldsymbol{f}, \boldsymbol{x}), i \in [d_X]$ sum to the black-box representation's norm $\|f(\boldsymbol{x})\|_{\mathcal{H}}^2 = \langle \boldsymbol{f}(\boldsymbol{x}), \boldsymbol{f}(\boldsymbol{x}) \rangle_{\mathcal{H}}$ up to a constant baseline $b_0 \in \mathbb{R}$ for all $\boldsymbol{x} \in \mathcal{X}$:*

$$\sum_{i=1}^{d_X} b_i(\boldsymbol{f}, \boldsymbol{x}) = \|\boldsymbol{f}(\boldsymbol{x})\|_{\mathcal{H}}^2 - b_0. \tag{3}$$

*Proof.* The proof is provided in Appendix A. $\qquad\square$

In Appendix A.2, we demonstrate that our label-free extension of feature importance verifies crucial invariance properties with respect to latent space symmetries.

# 3 Example Importance

In this section, we present our approach to extend example importance methods to the label-free setting. Since example importance methods are harder to unify, the structure of this section differs from the previous one. We split the example importance methods in two families: the loss-based and representation-based methods. The extension to the label-free setting works differently for these two families. Hence, we treat them in two distinct subsections. In both cases, we work with an input space $\mathcal{X}$ and a latent space $\mathcal{H}$. We are given a training set of $N \in \mathbb{N}^*$ examples $\mathcal{D}_{\text{train}} = \{\boldsymbol{x}^n \mid n \in [N]\}$. This training set is used to fit a black-box latent map $\boldsymbol{f} : \mathcal{X} \to \mathcal{H}$. We want to assign an example importance score $c^n(\boldsymbol{f}, \boldsymbol{x})$ to each training example $\boldsymbol{x}^n \in \mathcal{D}_{\text{train}}$ for the black-box $\boldsymbol{f}$ to build a representation $\boldsymbol{f}(\boldsymbol{x}) \in \mathcal{H}$ of a test example $\boldsymbol{x} \in \mathcal{X}$. Note that we use upper indices for examples in contrast with the lower indices for the features. Hence, $x_i^n$ denotes feature $i$ of training example $n$. Similarly, $c^n$ denotes an example importance in contrast with $b_i$ that is used for feature importance.

## 3.1 Loss-Based Example Importance

Loss-based example importance methods assign a score to each training example $\boldsymbol{x}^n$ by simulating the effect of their removal on the loss for a testing example. To make this more explicit, we denote by $\boldsymbol{z}$ the data of an example that is required to evaluate the loss. In a supervised setting, this typically correspond to a couple $\boldsymbol{z} = (\boldsymbol{x}, \boldsymbol{y})$ with an input $\boldsymbol{x} \in \mathcal{X}$ and a label $\boldsymbol{y} \in \mathcal{Y}$. Similarly, the training set is of the form $\mathcal{D}_{\text{train}} = \{\boldsymbol{z}^n \mid n \in [N]\}$. This training set is used to fit a black-box model $\boldsymbol{f}_\theta$ parametrized by $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^P$, where $P \in \mathbb{N}^*$ is the number of model parameters. Training

---

[3]Gradient-based methods like Integrated Gradients require differentiable functions, a restriction of $\mathbb{R}^{\mathcal{X}}$.

is done by solving $\boldsymbol{\theta}_* = \arg\min_{\boldsymbol{\theta}\in\Theta} \sum_{n=1}^{N} L(\boldsymbol{z}^n, \boldsymbol{\theta})$ with the loss $L : \mathcal{Z} \times \Theta \to \mathbb{R}$. This yields a model $\boldsymbol{f}_{\boldsymbol{\theta}_*}$. If we remove an example $\boldsymbol{z}^n$ from $\mathcal{D}_{\text{train}}$, the optimization problem turns into $\boldsymbol{\theta}_*^{-n} = \arg\min_{\boldsymbol{\theta}\in\Theta} \sum_{m=1, m\neq n}^{N} L(\boldsymbol{z}^m, \boldsymbol{\theta})$. This creates a parameter shift $\boldsymbol{\delta}^n\boldsymbol{\theta} \equiv \boldsymbol{\theta}_*^{-n} - \boldsymbol{\theta}_*$. This parameter shift, in turns, impacts the loss $L(\boldsymbol{z}, \boldsymbol{\theta})$ on a test example $\boldsymbol{z}$. This shift is reflected by the quantity $\boldsymbol{\delta}^n L(\boldsymbol{z}, \boldsymbol{\theta}_*) \equiv L(\boldsymbol{z}, \boldsymbol{\theta}_*^{-n}) - L(\boldsymbol{z}, \boldsymbol{\theta}_*)$. This loss shift permits to draw a distinction between proponents ($\boldsymbol{z}^n$ that decrease the loss: $\boldsymbol{\delta}^n L(\boldsymbol{z}, \boldsymbol{\theta}_*) > 0$) and opponents ($\boldsymbol{z}^n$ that increase the loss: $\boldsymbol{\delta}^n L(\boldsymbol{z}, \boldsymbol{\theta}_*) < 0$). Hence, it provides a meaningful measure of example importance. In order to estimate the loss shift without retraining the model, Koh and Liang [2017] propose to evaluate the influence function:

$$\boldsymbol{\delta}^n L(\boldsymbol{z}, \boldsymbol{\theta}_*) \approx \frac{1}{N} \left\langle \boldsymbol{\nabla}_{\boldsymbol{\theta}} L(\boldsymbol{z}, \boldsymbol{\theta}_*), \boldsymbol{H}_{\boldsymbol{\theta}_*}^{-1} \boldsymbol{\nabla}_{\boldsymbol{\theta}} L(\boldsymbol{z}^n, \boldsymbol{\theta}_*) \right\rangle_{\Theta},$$

where $\langle\cdot,\cdot\rangle_{\Theta}$ is an inner product on the parameter space $\Theta$ and $\boldsymbol{H}_{\boldsymbol{\theta}_*} \equiv \sum_{n=1}^{N} N^{-1} \boldsymbol{\nabla}_{\boldsymbol{\theta}}^2 L(\boldsymbol{z}^n, \boldsymbol{\theta}_*) \in \mathbb{R}^{P \times P}$ is the training loss Hessian matrix. Similarly, Pruthi et al. [2020] propose to use checkpoints during the model training to evaluate the loss shift:

$$\boldsymbol{\delta}^n L(\boldsymbol{z}, \boldsymbol{\theta}_*) \approx \sum_{t=1}^{T} \eta_t \left\langle \boldsymbol{\nabla}_{\boldsymbol{\theta}} L(\boldsymbol{z}, \boldsymbol{\theta}_t), \boldsymbol{\nabla}_{\boldsymbol{\theta}} L(\boldsymbol{z}^n, \boldsymbol{\theta}_t) \right\rangle_{\Theta},$$

where $\eta_t$ and $\boldsymbol{\theta}_t$ are respectively the learning rate and the model's parameters at checkpoint $t \in [T]$, $T \in \mathbb{N}^*$ is the total number of checkpoints. Similar approaches building on the theory of Shapley values exist Ghorbani and Zou [2019], Ghorbani et al. [2020].

We now turn to the label-free setting. In this case, we train our model with a label-free loss $L : \mathcal{X} \times \Theta \to \mathbb{R}$. Is it enough to drop the label and fix $\boldsymbol{z} = \boldsymbol{x}$ in all the above expressions? Most of the time, no. It is important to notice that the latent map $\boldsymbol{f} : \mathcal{X} \to \mathcal{H}$ that we wish to interpret is not necessarily equal to the model that we use to evaluate the loss $L(\boldsymbol{z}, \boldsymbol{\theta})$. To understand, it is instructive to consider a concrete example. Let us assume that we are in a self-supervised setting and that we train an autoencoder $\boldsymbol{f}_d \circ \boldsymbol{f}_e : \mathcal{X} \to \mathcal{X}$ that consists in an encoder $\boldsymbol{f}_e : \mathcal{X} \to \mathcal{H}$ and a decoder $\boldsymbol{f}_d : \mathcal{H} \to \mathcal{X}$ on a pretext task such as denoising. At the end of the pretraining, we typically keep the encoder $\boldsymbol{f}_e$ and drop the decoder $\boldsymbol{f}_d$. If we want to compute the loss-based example importance, we are facing a difficulty. On the one hand, we would like to draw interpretations that rely solely on the encoder $\boldsymbol{f}_e$ that is going to be used in a downstream task. On the other hand, loss-based example importance scores are computed with the autoencoder $\boldsymbol{f}_d \circ \boldsymbol{f}_e$ that also involves the irrelevant decoder $\boldsymbol{f}_d$. Whenever the loss evaluation involves more than the black-box latent map $\boldsymbol{f}$ we wish to interpret, replacing $\boldsymbol{z}$ by $\boldsymbol{x}$ in the above expressions is therefore insufficient.

To provide a more satisfactory solution, we split the parameter space $\Theta = \Theta_{\text{r}} \times \Theta_{\text{irr}}$ into a relevant component $\Theta_{\text{r}} \subseteq \mathbb{R}^{P_{\text{r}}}$ and an irrelevant component $\Theta_{\text{irr}} \subseteq \mathbb{R}^{P_{\text{irr}}}$. The black-box $\boldsymbol{f}$ to interpret is parametrized only by the relevant parameters $\boldsymbol{\theta}_{\text{r}}$ and can be denoted $\boldsymbol{f}_{\boldsymbol{\theta}_{\text{r}}}$. Concretely, we are interested in isolating the part of loss shift $\boldsymbol{\delta}^n L(\boldsymbol{z}, \boldsymbol{\theta}_*)$ caused by the variation of these relevant parameters. We note that the above estimators for this quantity involve gradients with respect to the parameters $\boldsymbol{\theta}$. We decompose the overall parameters gradient in terms of relevant and irrelevant parameters gradients: $\boldsymbol{\nabla}_{\boldsymbol{\theta}} = \boldsymbol{\nabla}_{\boldsymbol{\theta}_{\text{r}}} + \boldsymbol{\nabla}_{\boldsymbol{\theta}_{\text{irr}}}$. Ignoring the variation of irrelevant parameters is equivalent to setting $\boldsymbol{\nabla}_{\boldsymbol{\theta}_{\text{irr}}} = 0$ in the above expressions. This is trivially equivalent to the replacement of $\boldsymbol{\theta}$ by $\boldsymbol{\theta}_{\text{r}}$. This motivates the following definition.

**Definition 3.1** (Label-Free Loss-Based Importance). Let $f_{\boldsymbol{\theta}_{\text{r}}} : \mathcal{X} \to \mathcal{H}$ be a black-box latent map trained to minimize a loss $L : \mathcal{X} \times \Theta \to \mathbb{R}$ on a training set $\mathcal{D}_{\text{train}} = \{\boldsymbol{x}^n \mid n \in [N]\}$. To measure the impact of removing example $\boldsymbol{x}^n$ from $\mathcal{D}_{\text{train}}$ with $n \in [N]$, we define the *Label-Free Loss-Based Example Importance* as a score $c^n(\cdot, \cdot) : \mathcal{H}^{\mathcal{X}} \times \mathcal{X} \to \mathbb{R}$ such that

$$c^n(\boldsymbol{f}_{\boldsymbol{\theta}_{\text{r}}}, \boldsymbol{x}) = \boldsymbol{\delta}^n L(\boldsymbol{x}, \boldsymbol{\theta}_{\text{r},*}) \tag{4}$$

*Remark* 3.2. This definition gives a simple recipe to extend any loss-based example importance method to the label-free setting. In practice, this is implemented by using the unsupervised loss $L$ trained to fit the model and differentiating $L$ with respect to parameters of the model we wish to interpret.

## 3.2   Representation-Based Example Importance

Although representation-based example importance methods are introduced in a supervised context, their extension to the label-free setting is straightforward. These methods assign a score to each training example $\boldsymbol{x}^n$ by analyzing the latent representations of these examples. To make this more concrete, we start with a typical supervised setting. Consider a model $\boldsymbol{f}_l \circ \boldsymbol{f}_e : \mathcal{X} \to \mathcal{Y}$, where $\boldsymbol{f}_e : \mathcal{X} \to \mathcal{H}$ maps inputs to latent representations and $\boldsymbol{f}_l : \mathcal{H} \to \mathcal{Y}$ maps representations to labels. In the case of deep neural networks, the representation space typically corresponds to the output of an intermediate layer. We would like to see how the representation map relates a test example $\boldsymbol{x} \in \mathcal{X}$ to the training set examples. This can be done by mapping the training set inputs into the representation space $\boldsymbol{f}_e(\mathcal{D}_{\text{train}}) =$

$\{\boldsymbol{f}_e(\boldsymbol{x}^n) \mid n \in [N]\}$. To quantify the affinity between $\boldsymbol{x}$ and the training set examples, we attempt a reconstruction of $\boldsymbol{f}_e(\boldsymbol{x})$ with training representations from $\boldsymbol{f}_e(\mathcal{D}_{\text{train}})$: $\boldsymbol{f}_e(\boldsymbol{x}) \approx \sum_{n=1}^{N} w^n(\boldsymbol{x}) \cdot \boldsymbol{f}_e(\boldsymbol{x}^n)$. If this reconstruction is good, it also provides an approximation of the predicted output via $\boldsymbol{f}_l \circ \boldsymbol{f}_e(\boldsymbol{x}) \approx \boldsymbol{f}_l(\sum_{n=1}^{N} w^n(\boldsymbol{x}) \cdot \boldsymbol{f}_e(\boldsymbol{x}^n))$. Following Papernot and McDaniel [2018], the most obvious approach to define weights $w^n(\boldsymbol{x})$ is to identify the indices $\text{KNN}(\boldsymbol{x}) \subset [N]$ of the $K$ nearest neighbours (DKNN) of $\boldsymbol{f}_e(\boldsymbol{x})$ in $\boldsymbol{f}_e(\mathcal{D}_{\text{train}})$ and weigh them according to a Kernel function $\kappa : \mathcal{H}^2 \to \mathbb{R}^+$:

$$w^n(\boldsymbol{x}) = \mathbf{1}\left[n \in \text{KNN}(\boldsymbol{x})\right] \cdot \kappa\left[\boldsymbol{f}_e(\boldsymbol{x}^n), \boldsymbol{f}_e(\boldsymbol{x})\right], \tag{5}$$

where $\mathbf{1}$ denotes the indicator function. Similarly, Crabbé et al. [2021] propose to learn the weights by solving

$$\boldsymbol{w}(\boldsymbol{x}) = \underset{\boldsymbol{\lambda} \in [0,1]^N}{\arg\min} \left\| \boldsymbol{f}_e(\boldsymbol{x}) - \sum_{n=1}^{N} \lambda^n \boldsymbol{f}_e(\boldsymbol{x}^n) \right\|_{\mathcal{H}} \tag{6}$$

such that $\sum_{n=1}^{N} \lambda^n = 1$. Similar approaches building on the representer theorem exist Yeh et al. [2018].

We now turn to the label-free setting. The above discussion remains valid if we replace the supervised representation map $\boldsymbol{f}_e$ by an unsupervised representation map $\boldsymbol{f}$. The only part that we are loosing is the reconstruction of the label prediction $\boldsymbol{f}_l \circ \boldsymbol{f}_e(\boldsymbol{x})$ in terms of the training set predictions. In short, we can take $c^n = w^n$ without any additional work for representation-based methods.

## 4 Experiments

In this section, we conduct quantitative evaluations of the label-free extensions of various explanation methods. We start with simple consistency checks to ensure that these methods provide sensible explanations for unsupervised models. Then, we demonstrate how our label-free explanation paradigm makes it possible to compare representations learned from different pretext tasks. Finally, we challenge Definition 2.1 by studying saliency maps of VAEs. For each experiment, a more detailed description can be found in Appendix C. Supplementary experiments with medical time series can be found in Appendix D.

### 4.1 Consistency Checks

We would like to asses whether the approaches described in Sections 2 and 3 provide a sensible way to extend feature and example importance scores to the unsupervised setting.

**Setup** We explain the predictions of a denoising autoencoder $\boldsymbol{f}_d \circ \boldsymbol{f}_e : \mathcal{X} \to \mathcal{X}$ that consists in an encoder $\boldsymbol{f}_e : \mathcal{X} \to \mathcal{H}$ and a decoder $\boldsymbol{f}_d : \mathcal{H} \to \mathcal{X}$ with $d_X = 28^2$, $d_H = 4$. This model is trained on the MNIST dataset LeCun et al. [1998] with 60k training images and 10k testing images. We corrupt each training image $\boldsymbol{x} \in \mathcal{D}_{\text{train}}$ with random noise $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I}/3)$ where $\boldsymbol{I}$ is the identity matrix on $\mathcal{X}$. The model is trained to minimize the denoising loss $L_{\text{den}}(\boldsymbol{x}) = \mathbb{E}_{\boldsymbol{\epsilon}}[\boldsymbol{x} - \boldsymbol{f}_d \circ \boldsymbol{f}_e(\boldsymbol{x} + \boldsymbol{\epsilon})]^2$. The autoencoder is trained for 100 epochs with patience 10. ▶ **Feature Importance:** We compute the associated label-free importance score $b_i(\boldsymbol{f}_e, \boldsymbol{x})$ of each pixel $x_i$ for building the latent representation of the test images $\boldsymbol{x} \in \mathcal{D}_{\text{test}}$. We mask the $M$ most important pixels of the image with a mask $\boldsymbol{m} \in \{0,1\}^{d_X}$. For each image, we measure the latent shift $\|\boldsymbol{f}_e(\boldsymbol{x}) - \boldsymbol{f}_e(\boldsymbol{m} \odot \boldsymbol{x})\|_{\mathcal{H}}$ induced by masking the most important pixels, where $\odot$ denotes the Hadamard product. We expect this shift to increase with the importance of masked pixels. We report the average shift over the testing set (10k images) for several values of $M$ and feature importance methods in Figure 3. ▶ **Example Importance:** The consistency check for our label-free example importance methods is inspired by Kong and Chaudhuri [2021]. We sample $N = 1000$ training images $\boldsymbol{x}^n \in \mathcal{D}_{\text{train}}, n \in [N]$ without replacement. We compute the importance score $c^n(\boldsymbol{f}_e, \boldsymbol{x})$ of each training example $\boldsymbol{x}^n$ for predicting the latent representation of the test images $\boldsymbol{x} \in \mathcal{D}_{\text{test}}$. We select the $M$ most important training examples $\boldsymbol{x}^{n_1}, \dots, \boldsymbol{x}^{n_M}$. We compare their ground truth label $y^{n_1}, \dots, y^{n_M}$ to the label $y$ of $\boldsymbol{x}$. We compute the similarity rates $\sum_{m=1}^{M} \delta_{y,y^{n_m}}/M$, where $\delta$ denotes the Kronecker delta. We reproduce the above steps for the $M$ least important examples. If the encoder meaningfully represents the data, we expect the similarity rate of the most important examples to be higher than for the least important examples. We report the distribution of similarity rates across 1,000 test images for various values of $M$ and example importance methods in Figure 4.

**Results** ▶ **Feature Importance:** The label-free version of all the feature importance methods exhibit a similar behaviour: the latent shift increases sharply as we perturb the few most important pixels. This increase decelerates when we start perturbing pixels that are less relevant. Furthermore, selecting the perturbed pixels according to the various importance scores $b_i(\boldsymbol{f}_e, \boldsymbol{x})$ yields latent shifts that are significantly bigger than the shift induced by perturbing random pixels. These observations confirm that the label-free importance scores $b_i(\boldsymbol{f}_e, \boldsymbol{x})$ allow us to identify the
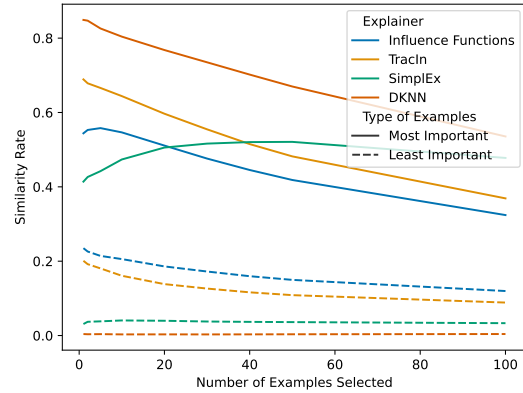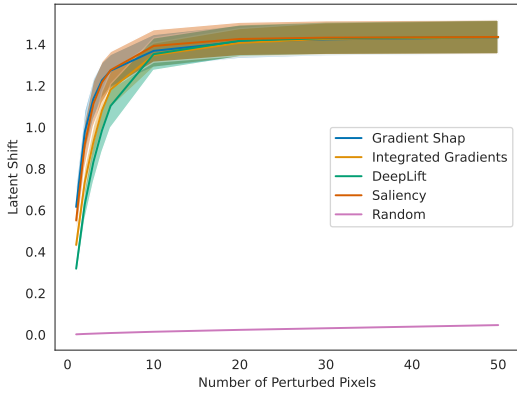
Figure 3: Shift induced by masking important pixels.



Figure 4: Similarity between labels of selected examples.

features that are the most relevant for the encoder $\boldsymbol{f}_e$ to build a latent representation for the image $\boldsymbol{x}$. ▶ **Example Importance:** For all example importance method, we observe that the similarity rate is significantly higher among the most similar examples than among the least similar examples. This observation confirms that the label-free importance scores $c^n(\boldsymbol{f}_e, \boldsymbol{x})$ allow us to identify training examples that are related to the test example we wish to explain. In this case, the verification also validates our autoencoder $\boldsymbol{f}$ since no label was used during training.

## 4.2 Use case: comparing the representations learned with different pretext tasks
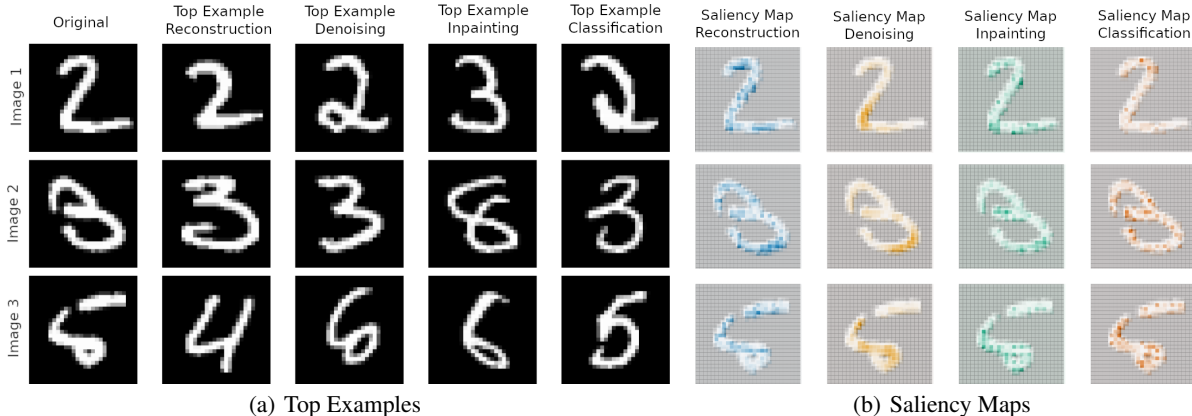


(a) Top Examples

(b) Saliency Maps

Figure 4: Explanations for different pretext tasks.

In a self-supervised learning setting, many unsupervised pretext task can be chosen to learn a meaningful representation of the data. How do the representations from different pretext tasks compare to each other? We show that label-free explainability permits to answer this question through qualitative and quantitative comparisons of representations.

**Setup** The setup is similar to Section 4.1. We consider two extra pretext tasks on top of denoising: reconstruction and inpainting Pandit et al. [2019]. The reconstruction autoencoder is trained to minimize the reconstruction loss $L_{\text{rec}}(\boldsymbol{x}) = [\boldsymbol{x} - \boldsymbol{f}_d \circ \boldsymbol{f}_e(\boldsymbol{x})]^2$. The inpainting autoencoder is trained to minimize the inptaiting loss $L_{\text{in}}(\boldsymbol{x}) = \mathbb{E}_{\boldsymbol{M}}[\boldsymbol{x} - \boldsymbol{f}_d \circ \boldsymbol{f}_e(\boldsymbol{M} \odot \boldsymbol{x})]^2$, where $\boldsymbol{M}$ is a random mask with $M_i \sim \text{Bernoulli}(0.8)$ for all $i \in [d_X]$. On top of that, we use the labelled training set to fit a MNIST classifier $\boldsymbol{f}_l \circ \boldsymbol{f}_e$, where $\boldsymbol{f}_l : \mathcal{X} \to \mathcal{Y}$ is an extra linear layer followed by a softmax activation. We are interested in comparing the representation spaces learned by the encoder $\boldsymbol{f}_e$ for the various tasks. Each autoencoder is trained for 100 epochs with patience 10. ▶ **Feature Importance:** For each encoder $\boldsymbol{f}_e$, we use our label-free Gradient Shap to produce saliency maps $b_i(\boldsymbol{f}_e, \boldsymbol{x})$ for the test images $\boldsymbol{x} \in \mathcal{D}_{\text{test}}$. To compare the saliency maps obtained by different models, a common approach is to compute their Pearson correlation coefficient Le

Meur and Baccino [2013]. We report the average Pearson coefficients between the encoders across 5 runs in Table 1.
▶ **Example Importance:** For each encoder $\boldsymbol{f}_e$, we use our label-free Deep-KNN to produce example importance $c^n(\boldsymbol{f}_e, \boldsymbol{x})$ of 1,000 training examples $\boldsymbol{x}^n \in \mathcal{D}_{\text{train}}$ for 1,000 test images $\boldsymbol{x} \in \mathcal{D}_{\text{test}}$. Again, we use the Pearson correlation coefficient to compare different encoders. We report the average Pearson coefficients between the encoders across 5 runs in Table 2.

Table 1: Pearson correlation for saliency maps (avg +/- std).

| PEARSON | RECON. | DENOIS. | INPAINT. | CLASSIF. |
|---|---|---|---|---|
| RECON. | $1.0 \pm 0.0$ | $.39 \pm .01$ | $.31 \pm .02$ | $.44 \pm .02$ |
| DENOIS. | $.39 \pm .01$ | $1.0 \pm .00$ | $.32 \pm .01$ | $.40 \pm .00$ |
| INPAINT. | $.31 \pm .02$ | $.32 \pm .01$ | $1.0 \pm 0.0$ | $.32 \pm .02$ |
| CLASSIF. | $.44 \pm .02$ | $.40 \pm .00$ | $.32 \pm .02$ | $1.0 \pm 0.0$ |

Table 2: Pearson correlation for example importance (avg +/- std).

| PEARSON | RECON. | DENOIS. | INPAINT. | CLASSIF. |
|---|---|---|---|---|
| RECON. | $1.0 \pm 0.0$ | $.10 \pm .04$ | $.11 \pm .03$ | $.07 \pm .03$ |
| DENOIS. | $.10 \pm .04$ | $1.0 \pm 0.0$ | $.12 \pm .03$ | $.06 \pm .02$ |
| INPAINT. | $.11 \pm .03$ | $.12 \pm .03$ | $1.0 \pm 0.0$ | $.07 \pm .01$ |
| CLASSIF. | $.07 \pm .03$ | $.06 \pm .02$ | $.07 \pm .01$ | $1.0 \pm 0.0$ |

**Results** ▶ **Not all representations are created equal.** For saliency maps, the Pearson correlation coefficients range from .31 to .44. This corresponds to moderate positive correlations. A good baseline to interpret these results is provided by Ouerhani et al. [2003]: the correlation between the fixation of two human subjects (human saliency maps) are typically in the same range. Hence, two encoders trained on distinct pretext tasks pay attention to different parts of the image like two separate human subjects typically do. For example importance scores, the Pearson correlation coefficients range from .06 to .12, which correspond to weak correlations. For both explanation types, these quantitative results strongly suggest that distinct pretext tasks do not yield interchangeable representations. ▶ **What makes classification special?** For saliency maps, the autoencoder-classifier correlations are comparable to those of the autoencoder-autoencoder. This shows that using labels creates a shift in the saliency maps comparable to changing the unsupervised pretext task. Hence, classification does not appear as a privileged task in terms of feature importance. Things are different for example importance: the autoencoder-classifier correlations are substantially lower than those of the autoencoder-autoencoder. One likely reason is that the classifier groups examples together according to an external label that is unknown to the autoencoders.

**Qualitative Analysis** Beyond quantitative analysis, label-free explainability allows us to appreciate qualitative differences between different encoders. To illustrate, we plot the most important DKNN and the saliency maps for the various encoders in Figure 4. ▶ **Feature Importance:** In accordance with our quantitative analysis, the saliency maps between different tasks look different. For instance, the denoising encoder seems to focus on small contiguous parts of the images. In contrast, the classifier seems to focus on a few isolated pixels. ▶ **Example Importance:** The top examples are rarely similar across various pretext tasks, as suggested by the quantitative analysis. The classifier is the only one that associates an example of the same class given an ambiguous image like Image 3. ▶ **Synergies:** Sometimes, saliency maps permit to better understand example importance. For instance, let us consider Image 3. In comparison to the other encoders, the reconstruction encoder pays less attention to the crucial loop at the image bottom. Hence, it is not surprising that the corresponding top example is less relevant than those selected by the other encoders.

### 4.3   Challenging our assumptions with disentangled VAEs

In Definition 2.1, the inner product appearing in the label-free importance expression corresponds to a sum over the latent space dimensions. In practice, this has the effect of mixing the feature importance for each latent unit to compute an overall feature importance. While this is reasonable when no particular meaning is attached to each latent unit, it might be inappropriate when the units are designed to be interpretable. Disantangled VAEs, for instance, involve latent units that are sensitive to change in single data generative factors, while invariants to other factors. This selective sensitivity permits to assign a meaning to each unit. An important question ensues: can we identify the generative factor associated to each latent unit by using their saliency maps? To answer, we propose a qualitative and quantitative analysis of the saliency maps from disentangled latent units.
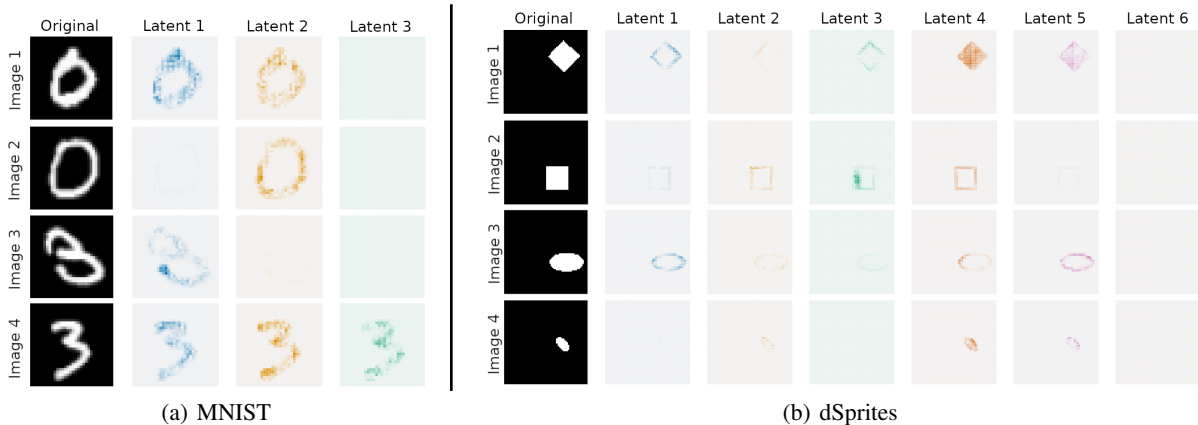
(a) MNIST

(b) dSprites

Figure 5: Saliency maps for each unit of the disentangled VAEs. The scale is constant for each image.
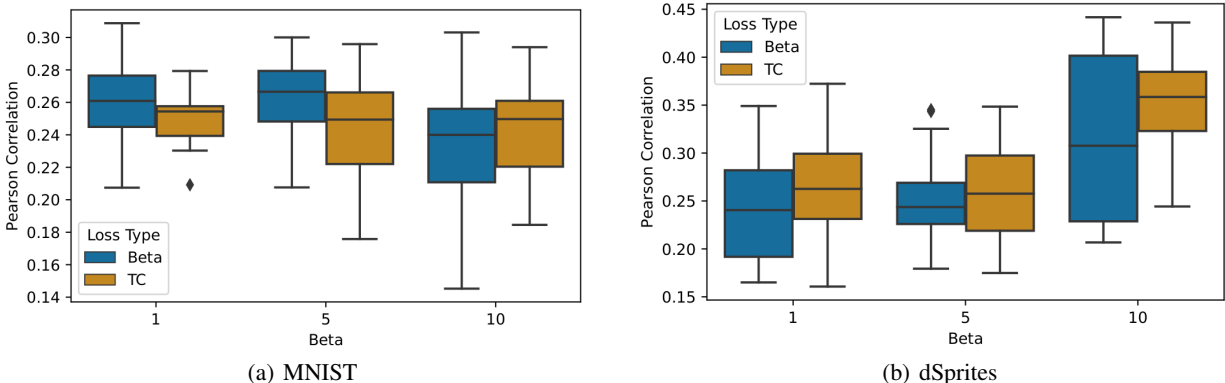


(a) MNIST

(b) dSprites

Figure 6: Pearson correlation between saliency maps for different values of $\beta$.

**Setup** We study the interpretability of two popular disentangled VAEs : the $\beta$-VAE Higgins et al. [2017] and the TC-VAE Chen et al. [2018]. Those two VAEs involve a variational encoder computing the expected representation $\boldsymbol{\mu} : \mathcal{X} \to \mathcal{H}$ as well as its standard deviation $\boldsymbol{\sigma} : \mathcal{X} \to \mathcal{H}$ and a decoder $\boldsymbol{f}_d : \mathcal{H} \to \mathcal{X}$. Latent samples are obtained with the reparametrization trick Kingma and Welling [2013]: $\boldsymbol{h} = \boldsymbol{\mu}(\boldsymbol{x}) + \boldsymbol{\sigma}(\boldsymbol{x}) \odot \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$. These VAEs are trained on the MNIST and dSprites datasets Matthey et al. [2017] ($90\% - 10\%$ train-test split) to minimize their objective. We use $d_H = 3$ latent units for MNIST and $d_H = 6$ for dSprites. We train 20 disentangled VAEs of each type for $\beta \in \{1, 5, 10\}$ and for 100 epochs with patience 10.

**Qualitative Analysis** We use Gradient Shap to evaluate the importance $a_i(\mu_j, \boldsymbol{x})$ of each pixel $x_i$ from the image $\boldsymbol{x}$ to predict each latent unit $j \in [d_H]$. Again, we use the Pearson correlation to compare different saliency maps associated to each latent unit [4]. In this case, a low Pearson correlation corresponds to latent units paying attention to distinct parts of the images. Clearly, this is desirable if we want to easily identify the specialization of each latent unit. Therefore, use this criterion to select a VAE to analyse among the 120 VAEs we trained on each dataset. This corresponds to a $\beta$-VAE with $\beta = 10$ for MNIST and a TC-VAE with $\beta = 1$ for dSprites. We show the various saliency maps for 4 test images on Figure 5. The saliency maps appear to follow patterns that make the interpretation difficult. Here are a couple of examples that we can observe: **(1)** A latent unit is sensitive to a given image while insensitive to a similar image (e.g. Latent Unit 1 of the MNIST VAE is sensitive to Image 1 but not to Image 2). **(2)** The focus of a latent unit changes completely between two similar images (e.g. Latent Unit 4 of the dSprites VAE focuses on the interior of the square from Image 1 but only on the edges of the square from Image 2). **(3)** Several latent units focus on the same part of the image (e.g. Image 4 of MNIST and Image 3 of dSprites). Additional examples can be found in Appendix C.

---

[4]We average the correlation over pairs of latent units.

**Quantitative Analysis** If we increase the disentanglement of a representation space, does it imply that distinct latent units are going to focus on more distinct features? For the disantangled VAEs we consider, the strength of disentanglement increases with $\beta$. Hence, the above question can be formulated more precisely : does the correlation between the latent units saliency map decrease with $\beta$? To answer this question, we make box-plots of the Pearson correlation coefficients for both VAE types (Beta and TC) and various values of $\beta$. Results can be observed in Figure 6. For MNIST, the Pearson correlation slightly decreases with $\beta$ (Spearman $\rho = -.15$). For dSprites, the Pearson correlation moderately increases with $\beta$ (Spearman $\rho = .43$). This analysis shows that increasing $\beta$ *does not* imply that latent units are going to pay attention to distinct part of the images. If this results is surprising at first, it can be understood by thinking about disentanglement. As aforementioned, increasing disentanglement encourages distinct latent units to pay attention to distinct generative factors. There is no guarantee that distinct generative factors are associated to distinct features. To illustrate, let us consider two generative factors of the dSprites dataset: the position of a shape and its scale. Clearly, these two generative factors can be identified by paying attention to the edges of the shape appearing on the image, as various latent units do in Figure 5(b). Whenever generative factors are not unambiguously associated to a subset of features, latent units can identify distinct generative factors by looking at similar features. In this case, increasing disentanglement of latent units does not necessarily make their saliency maps more decorrelated. We conclude two things: **(1)** If we want to identify the role of each latent unit with their saliency maps, disentanglement might not be the right approach. Perhaps it is possible to introduce priors on the saliency map to control the features the model pays attention to, like it was done by Erion et al. [2021] in a supervised setting. We leave this for future works. **(2)** In the absence of interpretability for individual units, our label-free feature importance extension combines the saliency of individual units by extending crucial properties such as completeness and invariance with respect to latent symmetries.

# References

Amina Adadi and Mohammed Berrada. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access*, 6:52138–52160, 2018. ISSN 21693536. doi:10.1109/ACCESS.2018.2870052.

Naman Agarwal, Brian Bullins, and Elad Hazan. Second-Order Stochastic Optimization for Machine Learning in Linear Time. *Journal of Machine Learning Research*, 18:1–40, 2016. ISSN 15337928.

Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador Garcia, Sergio Gil-Lopez, Daniel Molina, Richard Benjamins, Raja Chatila, and Francisco Herrera. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58:82–115, 2020. ISSN 15662535. doi:10.1016/j.inffus.2019.12.012.

Lennart Brocki and Neo Christopher Chung. Concept saliency maps to visualize relevant features in deep generative models. *Proceedings - 18th IEEE International Conference on Machine Learning and Applications, ICMLA 2019*, pages 1771–1778, 2019. doi:10.1109/ICMLA.2019.00287.

Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, Alexander Lerchner, and Deepmind London. Understanding disentangling in $\beta$-VAE. *arXiv*, 2018.

Tian Qi Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating Sources of Disentanglement in Variational Autoencoders. *6th International Conference on Learning Representations, ICLR 2018 - Workshop Track Proceedings*, 2018.

Travers Ching, Daniel S. Himmelstein, Brett K. Beaulieu-Jones, Alexandr A. Kalinin, Brian T. Do, Gregory P. Way, Enrico Ferrero, Paul-Michael Agapow, Michael Zietz, Michael M. Hoffman, Wei Xie, Gail L. Rosen, Benjamin J. Lengerich, Johnny Israeli, Jack Lanchantin, Stephen Woloszynek, Anne E. Carpenter, Avanti Shrikumar, Jinbo Xu, Evan M. Cofer, Christopher A. Lavender, Srinivas C. Turaga, Amr M. Alexandari, Zhiyong Lu, David J. Harris, Dave DeCaprio, Yanjun Qi, Anshul Kundaje, Yifan Peng, Laura K. Wiley, Marwin H. S. Segler, Simina M. Boca, S. Joshua Swamidass, Austin Huang, Anthony Gitter, and Casey S. Greene. Opportunities and obstacles for deep learning in biology and medicine. *Journal of The Royal Society Interface*, 15(141):20170387, 2018. ISSN 1742-5689. doi:10.1098/rsif.2017.0387.

R. Dennis Cook and Sanford Weisenberg. *Residuals and influence in regression*. Chapman and Hall, New York, 1982.

Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley, 2005. ISBN 9780471241959. doi:10.1002/047174882X.

Jonathan Crabbé, Zhaozhi Qian, Fergus Imrie, and Mihaela van der Schaar. Explaining Latent Representations with a Corpus of Examples. In *Advances in Neural Information Processing Systems*, 2021.

Arun Das and Paul Rad. Opportunities and Challenges in Explainable Artificial Intelligence (XAI): A Survey. *arXiv*, 2020.

Alex Davies, Petar Veličković, Lars Buesing, Sam Blackwell, Daniel Zheng, Nenad Tomašev, Richard Tanburn, Peter Battaglia, Charles Blundell, András Juhász, Marc Lackenby, Geordie Williamson, Demis Hassabis, and Pushmeet Kohli. Advancing mathematics by guiding human intuition with AI. *Nature 2021 600:7887*, 600(7887):70–74, 2021. ISSN 1476-4687. doi:10.1038/s41586-021-04086-x.

Gabriel Erion, Joseph D. Janizek, Pascal Sturmfels, Scott M. Lundberg, and Su-In Lee. Improving performance of deep learning models with axiomatic attribution priors and expected gradients. *Nature Machine Intelligence 2021 3:7*, 3(7):620–631, 2021. ISSN 2522-5839. doi:10.1038/s42256-021-00343-w.

Amirata Ghorbani and James Zou. Data Shapley: Equitable Valuation of Data for Machine Learning. *36th International Conference on Machine Learning, ICML 2019*, 2019-June:4053–4065, 2019.

Amirata Ghorbani, James Wexler, James Zou, and Been Kim. Towards Automatic Concept-based Explanations. *Advances in Neural Information Processing Systems*, 32, 2019. ISSN 10495258.

Amirata Ghorbani, Michael P Kim, and James Zou. A Distributional Framework for Data Valuation. *arXiv*, 2020.

Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep Sparse Rectifier Neural Networks, 2011. ISSN 1938-7228.

A. L. Goldberger, L. A. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C. K. Peng, and H. E. Stanley. PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals. *Circulation*, 101(23), 2000. doi:10.1161/01.CIR.101.23.E215.

Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, Alexander Lerchner, and Google Deepmind. beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. In *International Conference on Learning Representations*, 2017.

Longlong Jing and Yingli Tian. Self-Supervised Visual Feature Learning with Deep Neural Networks: A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(11):4037–4058, 2021. doi:10.1109/TPAMI.2020.2992393.

John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with AlphaFold. *Nature 2021 596:7873*, 596(7873):583–589, 2021. ISSN 1476-4687. doi:10.1038/s41586-021-03819-2.

Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised Contrastive Learning. *Advances in Neural Information Processing Systems*, 2020-December, 2020. ISSN 10495258.

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). *35th International Conference on Machine Learning*, 6:4186–4195, 2017.

Diederik P. Kingma and Jimmy Lei Ba. Adam: A Method for Stochastic Optimization. *3rd International Conference on Learning Representations, ICLR 2015 - Conference Track Proceedings*, 2014.

Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. *2nd International Conference on Learning Representations, ICLR 2014 - Conference Track Proceedings*, 2013.

Pang Wei Koh and Percy Liang. Understanding Black-box Predictions via Influence Functions. In *International Conference on Machine Learning*, volume 34, pages 2976–2987. PMLR, 2017.

Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. Captum: A unified and generic model interpretability library for pytorch, 2020.

Zhifeng Kong and Kamalika Chaudhuri. Understanding Instance-based Interpretability of Variational Auto-Encoders. In *Advances in Neural Information Processing Systems*, 2021.

Olivier Le Meur and Thierry Baccino. Methods for comparing scanpaths and saliency maps: Strengths and weaknesses. *Behavior Research Methods*, 45(1):251–266, 2013. ISSN 1554351X. doi:10.3758/S13428-012-0226-9/TABLES/2.

Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2323, 1998. ISSN 00189219. doi:10.1109/5.726791.

Zachary C. Lipton. The Mythos of Model Interpretability. *Communications of the ACM*, 61(10):35–43, 2016.

Scott Lundberg and Su-In Lee. A Unified Approach to Interpreting Model Predictions. *Advances in Neural Information Processing Systems*, pages 4766–4775, 2017.

Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. dSprites: Disentanglement testing Sprites dataset. https://github.com/deepmind/dsprites-dataset/, 2017.

Nabil Ouerhani, Heinz Hügli, René Müri, and Roman Von Wartburg. Empirical Validation of the Saliency-based Model of Visual Attention. *Electronic Letters on Computer Vision and Image Analysis*, 3(1):13–23, 2003.

Parthe Pandit, Mojtaba Sahraee-Ardakan, Sundeep Rangan, Philip Schniter, Alyson K Fletcher, P Pandit, M Sahraee-Ardakan, and A K Fletcher. Inference with Deep Generative Priors in High Dimensions. *IEEE Journal on Selected Areas in Information Theory*, 1(1):336–347, 2019. doi:10.1109/jsait.2020.2986321.

Nicolas Papernot and Patrick McDaniel. Deep k-Nearest Neighbors: Towards Confident, Interpretable and Robust Deep Learning. *arXiv*, 2018.

Barak A. Pearlmutter. Fast Exact Multiplication by the Hessian. *Neural Computation*, 6:160, 1994.

Garima Pruthi, Frederick Liu, Mukund Sundararajan, and Satyen Kale. Estimating Training Data Influence by Tracing Gradient Descent. In *Advances in Neural Information Processing Systems*, pages 19920–19930. Curran Associates, Inc., 2020.

Arun Rai. Explainable AI: from black box to glass box. *Journal of the Academy of Marketing Science 2019 48:1*, 48(1):137–141, 2019. ISSN 1552-7824. doi:10.1007/S11747-019-00710-5.

Alvin Rajkomar, Jeffrey Dean, and Isaac Kohane. Machine Learning in Medicine. *New England Journal of Medicine*, 380(14):1347–1358, 2019. ISSN 0028-4793. doi:10.1056/NEJMRA1814259/SUPPL_FILE/NEJMRA1814259_DISCLOSURES.PDF.

Carl Edward Rasmussen. Gaussian Processes in Machine Learning. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 3176:63–71, 2003. ISSN 16113349. doi:10.1007/978-3-540-28650-9_4.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, volume 13-17, pages 1135–1144. Association for Computing Machinery, 2016. ISBN 9781450342322. doi:10.1145/2939672.2939778.

Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence 2019 1:5*, 1(5):206–215, 2019. ISSN 2522-5839. doi:10.1038/s42256-019-0048-x.

Mhd Hasan Sarhan, Abouzar Eslami, Nassir Navab, and Shadi Albarqouni. Learning Interpretable Disentangled Representations Using Adversarial VAEs. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 11795 LNCS:37–44, 2019. doi:10.1007/978-3-030-33391-1_5.

C. E. Shannon. A Mathematical Theory of Communication. *Bell System Technical Journal*, 27(4):623–656, 1948. ISSN 00058580. doi:10.1002/j.1538-7305.1948.tb00917.x.

Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning Important Features Through Propagating Activation Differences. *34th International Conference on Machine Learning, ICML 2017*, 7:4844–4866, 2017.

Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. *2nd International Conference on Learning Representations, ICLR 2014 - Workshop Track Proceedings*, 2013.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic Attribution for Deep Networks. *34th International Conference on Machine Learning, ICML 2017*, 7:5109–5118, 2017.

Erico Tjoa and Cuntai Guan. A Survey on Explainable Artificial Intelligence (XAI): Toward Medical XAI. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–21, 2020. ISSN 2162-237X. doi:10.1109/tnnls.2020.3027314.

Athanasios Voulodimos, Nikolaos Doulamis, Anastasios Doulamis, and Eftychios Protopapadakis. Deep Learning for Computer Vision: A Brief Review. *Computational Intelligence and Neuroscience*, 2018, 2018. ISSN 16875273. doi:10.1155/2018/7068349.

Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual Explanations without Opening the Black Box: Automated Decisions and the GDPR. *SSRN Electronic Journal*, 2017. ISSN 1556-5068. doi:10.2139/ssrn.3063289.

Chih-Kuan Yeh, Joon Sik Kim, Ian E. H. Yen, and Pradeep Ravikumar. Representer Point Selection for Explaining Deep Neural Networks. *Advances in Neural Information Processing Systems*, 31:9291–9301, 2018.

Tom Young, Devamanyu Hazarika, Soujanya Poria, and Erik Cambria. Recent trends in deep learning based natural language processing. *IEEE Computational Intelligence Magazine*, 13(3):55–75, 2018. ISSN 15566048. doi:10.1109/MCI.2018.2840738.

# A   Properties of the Label-Free Extensions

In this appendix, we prove the completeness property. Next, we motivate and prove the orthogonal invariance of our label-free extensions.

## A.1   Completeness

Let us write a proof for Proposition 2.3.

*Proof.* The proof is an immediate consequence of Definition 2.1 and the completeness property of the feature importance score $a_i$:

$$\sum_{i=1}^{d_X} b_i(\boldsymbol{f}, \boldsymbol{x}) \stackrel{(1)}{=} \sum_{i=1}^{d_X} a_i(g_{\boldsymbol{x}}, \boldsymbol{x})$$
$$= g_{\boldsymbol{x}}(\boldsymbol{x}) - a_0,$$

where we used the completeness property to obtain the second equality and $a_0 \in \mathbb{R}$ is the baseline for the importance score $a_i$. By noting that $g_{\boldsymbol{x}}(\boldsymbol{x}) = \langle \boldsymbol{f}(\boldsymbol{x}), \boldsymbol{f}(\boldsymbol{x}) \rangle_{\mathcal{H}} = \|\boldsymbol{f}(\boldsymbol{x})\|_{\mathcal{H}}^2$, we obtain the desired equality (3) with the identification $b_0 = a_0$. $\qquad\square$

## A.2   Invariance with respect to latent symmetries

In Section 1, we described the ambiguity associated to the axes of the latent space $\mathcal{H}$. This line of reasoning can be made more formal with symmetries. Due to the fact that each axis of the latent space is not associated with a fixed and predetermined label, there exists many latent spaces that are equivalent to each other. For instance, if we swap two axes of the latent space by relabelling $h_1 \mapsto h_2$ and $h_2 \mapsto h_1$, we do not change the latent space structure. More generally, given an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$, the set of transformations that leave the geometry of the latent space $\mathcal{H}$ invariant is the set of orthogonal transformations.

**Definition A.1** (Orthogonal Transformations). Let $\mathcal{H}$ be a real vector space equipped with an inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H}^2 \to \mathbb{R}$. An orthogonal transformation is a linear map $\boldsymbol{T} : \mathcal{H} \to \mathcal{H}$ such that for all $\boldsymbol{h}_1, \boldsymbol{h}_2 \in \mathcal{H}$, we have:

$$\langle \boldsymbol{T}(\boldsymbol{h}_1), \boldsymbol{T}(\boldsymbol{h}_2) \rangle_{\mathcal{H}} = \langle \boldsymbol{h}_1, \boldsymbol{h}_2 \rangle_{\mathcal{H}}.$$

*Remark* A.2. In the case where $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is the standard euclidean inner product $\langle \boldsymbol{h}_1, \boldsymbol{h}_2 \rangle_{\mathcal{H}} = \boldsymbol{h}_1^{\mathsf{T}} \boldsymbol{h}_2$, the orthogonal transformations are represented by orthogonal matrices in $O(d_H) = \{\boldsymbol{M} \in \mathbb{R}^{d_H \times d_H} \mid \boldsymbol{M}^{\mathsf{T}} \boldsymbol{M}\}$. These transformations include rotations, axes permutations and mirror symmetries.

Of course, since these transformations leave the geometry of the latent space invariant, we would expect the same for the explanations. We verify that this is indeed the case for our label-free extension of feature importance.

**Proposition A.3** (Label-Free Feature Importance Invariance). *The label-free importance scores $b_i(\cdot, \cdot), i \in [d_X]$ are invariant with respect to orthogonal transformations in the latent space $\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}}$. More formally, for all $\boldsymbol{f} \in \mathcal{H}^{\mathcal{X}}$, $\boldsymbol{x} \in \mathcal{X}$ and $i \in [d_X]$:*

$$b_i(\boldsymbol{T} \circ \boldsymbol{f}, \boldsymbol{x}) = b_i(\boldsymbol{f}, \boldsymbol{x})$$

*Remark* A.4. This property is a further motivation for the usage of an inner product in Definition 2.1.

*Proof.* This proposition is a trivial consequence of the inner product appearing in Definition 2.1. Let $g_{\boldsymbol{x}}^{\boldsymbol{T}}$ be the auxiliary function associated to $\boldsymbol{T} \circ \boldsymbol{f}$ for some $\boldsymbol{x} \in \mathcal{X}$. We note that for all $\tilde{\boldsymbol{x}} \in \mathcal{X}$:

$$g_{\boldsymbol{x}}^{\boldsymbol{T}}(\tilde{\boldsymbol{x}}) \stackrel{(2)}{=} \langle \boldsymbol{T} \circ \boldsymbol{f}(\boldsymbol{x}), \boldsymbol{T} \circ \boldsymbol{f}(\tilde{\boldsymbol{x}}) \rangle_{\mathcal{H}}$$
$$= \langle \boldsymbol{f}(\boldsymbol{x}), \boldsymbol{f}(\tilde{\boldsymbol{x}}) \rangle_{\mathcal{H}}$$
$$\stackrel{(2)}{=} g_{\boldsymbol{x}}(\tilde{\boldsymbol{x}}),$$

where we used the fact that $\boldsymbol{T}$ is orthogonal in the second equality and $g_{\boldsymbol{x}}(\tilde{\boldsymbol{x}})$ is the auxiliary function associated to $\boldsymbol{f}$. Since this holds for any $\tilde{\boldsymbol{x}} \in \mathcal{X}$, we have that $g_{\boldsymbol{x}}^{\boldsymbol{T}} = g_{\boldsymbol{x}}$ for all $\boldsymbol{x} \in \mathcal{X}$. This allows us to write:

$$b_i(\boldsymbol{T} \circ \boldsymbol{f}, \boldsymbol{x}) \stackrel{(1)}{=} a_i(g_{\boldsymbol{x}}^{\boldsymbol{T}}, \boldsymbol{x})$$
$$= a_i(g_{\boldsymbol{x}}, \boldsymbol{x})$$
$$\stackrel{(1)}{=} b_i(\boldsymbol{f}, \boldsymbol{x})$$

for all $\boldsymbol{x} \in \mathcal{X}$ and $i \in [d_X]$. This is the desired identity $\qquad\square$

When it comes to example importance methods, the same guarantee holds for representation-based methods:

**Proposition A.5** (Representation-Based Example Importance Invariance). *Let $\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}}$ be a latent space. The label-free importance scores $c^n(\cdot, \cdot), n \in [N]$ outputted by DKNN (5) are invariant with respect to orthogonal transformations of $\mathcal{H}$ if they are defined with a kernel $\kappa : \mathcal{H}^2 \to \mathbb{R}^+$ that is invariant with respect to orthogonal transformations:*

$$\kappa(\boldsymbol{T}(\boldsymbol{h}_1), \boldsymbol{T}(\boldsymbol{h}_2)) = \kappa(\boldsymbol{h}_1, \boldsymbol{h}_2)$$

*for all orthogonal transformation $\boldsymbol{T}$ and $\boldsymbol{h}_1, \boldsymbol{h}_2 \in \mathcal{H}$. Similarly, the importance scores outputted by SimplEx (6) are invariant with respect to orthogonal transformations of $\mathcal{H}$. In both cases, the invariance property can be written more formally: for all $\boldsymbol{f} \in \mathcal{H}^{\mathcal{X}}$, $\boldsymbol{x} \in \mathcal{X}$ and $n \in [N]$:*

$$c^n(\boldsymbol{T} \circ \boldsymbol{f}, \boldsymbol{x}) = c^n(\boldsymbol{f}, \boldsymbol{x})$$

*Remark* A.6. The invariance property for the kernel function is verified for kernels that involve the inner-product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ in their definition. This includes RBF, Matern and Polynomial Kernels Rasmussen [2003]. Note that constant kernels trivially verify this property. Finally, replacing the kernel function by the inverse-distance in latent space (as it is done in our implementation) $\kappa(\boldsymbol{h}_1, \boldsymbol{h}_2) = \|\boldsymbol{h}_1 - \boldsymbol{h}_2\|_{\mathcal{H}}^{-1}$ also preserves the invariance property.

*Proof.* We start by noting that the latent space distance is invariant under orthogonal transformations. for all $\boldsymbol{h}_1, \boldsymbol{h}_2 \in \mathcal{H}$:

$$\begin{aligned}
\|\boldsymbol{T}(\boldsymbol{h}_1) - \boldsymbol{T}(\boldsymbol{h}_2)\|_{\mathcal{H}}^2 &= \|\boldsymbol{T}(\boldsymbol{h}_1 - \boldsymbol{h}_2)\|_{\mathcal{H}}^2 \\
&= \langle \boldsymbol{T}(\boldsymbol{h}_1 - \boldsymbol{h}_2), \boldsymbol{T}(\boldsymbol{h}_1 - \boldsymbol{h}_2) \rangle_{\mathcal{H}} \\
&= \langle \boldsymbol{h}_1 - \boldsymbol{h}_2, \boldsymbol{h}_1 - \boldsymbol{h}_2 \rangle_{\mathcal{H}} \\
&= \|\boldsymbol{h}_1 - \boldsymbol{h}_2\|_{\mathcal{H}}^2,
\end{aligned}$$

where we successively used the linearity and orthogonality of $\boldsymbol{T}$. Note that this equation is equivalent to $\|\boldsymbol{T}(\boldsymbol{h}_1) - \boldsymbol{T}(\boldsymbol{h}_2)\|_{\mathcal{H}} = \|\boldsymbol{h}_1 - \boldsymbol{h}_2\|_{\mathcal{H}}$ since both norms are positive. Since the latent KNNs in (5) are computed with this latent space distance, we deduce their invariance under orthogonal transformations. By combining this to the invariance of the kernel, we obtain the desired invariance for the DKNN importance scores (5). For all $\boldsymbol{x} \in \mathcal{X}, \boldsymbol{f} \in \mathcal{H}^{\mathcal{X}}$ and orthogonal transformation $\boldsymbol{T}$:

$$\begin{aligned}
c_{\text{DKNN}}^n(\boldsymbol{T} \circ \boldsymbol{f}, \boldsymbol{x}) &\overset{(5)}{=} \mathbf{1}\left[n \in \text{KNN}(\boldsymbol{x})\right] \cdot \kappa\left[\boldsymbol{T} \circ \boldsymbol{f}(\boldsymbol{x}^n), \boldsymbol{T} \circ \boldsymbol{f}(\boldsymbol{x})\right] \\
&= \mathbf{1}\left[n \in \text{KNN}(\boldsymbol{x})\right] \cdot \kappa\left[\boldsymbol{f}(\boldsymbol{x}^n), \boldsymbol{f}(\boldsymbol{x})\right] \\
&\overset{(5)}{=} c_{\text{DKNN}}^n(\boldsymbol{f}, \boldsymbol{x}),
\end{aligned}$$

where we have used the invariance property to obtain the second equality. We can proceed similarly for SimplEx (6):

$$\begin{aligned}
c_{\text{SimplEx}}^n(\boldsymbol{T} \circ \boldsymbol{f}, \boldsymbol{x}) &\overset{(6)}{=} \underset{\boldsymbol{\lambda} \in [0,1]^N}{\arg\min} \left\| \boldsymbol{T} \circ \boldsymbol{f}(\boldsymbol{x}) - \sum_{n=1}^{N} \lambda^n \boldsymbol{T} \circ \boldsymbol{f}(\boldsymbol{x}^n) \right\|_{\mathcal{H}} \\
&= \underset{\boldsymbol{\lambda} \in [0,1]^N}{\arg\min} \left\| \boldsymbol{T} \left[ \boldsymbol{f}(\boldsymbol{x}) - \sum_{n=1}^{N} \lambda^n \boldsymbol{f}(\boldsymbol{x}^n) \right] \right\|_{\mathcal{H}} \\
&= \underset{\boldsymbol{\lambda} \in [0,1]^N}{\arg\min} \left\| \boldsymbol{f}(\boldsymbol{x}) - \sum_{n=1}^{N} \lambda^n \boldsymbol{f}(\boldsymbol{x}^n) \right\|_{\mathcal{H}} \\
&\overset{(6)}{=} c_{\text{SimplEx}}^n(\boldsymbol{f}, \boldsymbol{x}),
\end{aligned}$$

where we have successively used the linearity and orthogonality of $\boldsymbol{T}$. Those are the desired identities.  $\square$

The only label-free extension that we have not yet discussed are the loss-based example importance methods from Section 3.1. Unfortunately, due to the fact that the black-box $\boldsymbol{f}$ is only a component required in the evaluation of the loss $L$, it is not possible to provide a general guarantee like in the previous examples. If we take the example of the autoencoder $\boldsymbol{f}_d \circ \boldsymbol{f}_e$, we note that applying an orthogonal transformation $\boldsymbol{f}_e \mapsto \boldsymbol{T} \circ \boldsymbol{f}_e$ to the encoder leaves the autoencoder invariant only if this transformation is undone by the decoder $\boldsymbol{f}_d \mapsto \boldsymbol{f}_d \circ \boldsymbol{T}^{-1}$. Unlike the other methods, the invariance of loss-based example importance scores therefore requires restrictive assumptions. If invariance of the explanations under orthogonal transformations is required, this might be an argument in favour of representation-based methods.

# B    Implementation Details

In this appendix, we detail the implementation of our label-free extensions.

## B.1    Label-Free Feature Importance

The label-free feature importance methods used in our experiments are described in Table 3:

Table 3: Feature Importance Methods.

| Method | Ref. | Linearity | Completeness | Label-Free Expression |
|---|---|---|---|---|
| Saliency | Simonyan et al. [2013] | ✓ | ✗ | $b_i(\boldsymbol{f},\boldsymbol{x}) = \frac{\partial g_{\boldsymbol{x}}}{\partial x_i}(\boldsymbol{x})$ |
| Integrated Gradients | Sundararajan et al. [2017] | ✓ | ✓ | $b_i(\boldsymbol{f},\boldsymbol{x}) = (x_i - \bar{x}_i)\int_0^1 \frac{\partial g_{\boldsymbol{x}}}{\partial x_i}(\boldsymbol{x} + \tau(\boldsymbol{x} - \bar{\boldsymbol{x}}))d\tau$ |
| Gradient Shap | Lundberg and Lee [2017] | ✓ | ✓ | $b_i(\boldsymbol{f},\boldsymbol{x}) = (x_i - \bar{x}_i)\mathbb{E}_{\boldsymbol{\epsilon},U}\left[\frac{\partial g_{\boldsymbol{x}}}{\partial x_i}(\boldsymbol{x} + \boldsymbol{\epsilon} + U(\boldsymbol{x} - \bar{\boldsymbol{x}}))\right]$ |
| DeepLift | Shrikumar et al. [2017] | ✓ | ✓ | $b_i(\boldsymbol{f},\boldsymbol{x}) = C_{\Delta x_i \Delta g_{\boldsymbol{x}}}$ |

where $\bar{\boldsymbol{x}} \in \mathcal{X}$ is a baseline input, $\boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$, $U \sim \text{Uniform}(0,1)$ and $C_{\Delta x_i \Delta g_{\boldsymbol{x}}}$ is used by propagating the Deeplift rules along the computational graph of $g_{\boldsymbol{x}}$. Note that, in each case, partial derivatives are computed with respect to the argument of $g_{\boldsymbol{x}}$ only (hence we do not consider derivatives of the form $g_{\frac{\partial \boldsymbol{x}}{\partial x_i}}$). We use the Captum Kokhlikyan et al. [2020] implementation of each method.

To extend this implementation to the label-free setting, it is necessary to define an auxiliary function $g_{\boldsymbol{x}}$ associated to the vectorial black-box function $\boldsymbol{f}$ for each testing example $\boldsymbol{x} \in \mathcal{D}_{\text{test}}$. With libraries such as Pytorch, it is possible to define an auxiliary function as a wrapper around the module that represents $\boldsymbol{f}$. This allows us to compute the importance scores with a single batch call of the original feature importance method, as described in Algorithm 1.

---

**Algorithm 1** Label-Free Feature Importance

---

**Input:** Batch of inputs $\boldsymbol{X} \in \mathcal{X}^B$ of size $B$, Black-box $\boldsymbol{f} : \mathcal{X} \to \mathcal{H}$, Feature importance method $a_i(\cdot,\cdot) : \mathcal{H}^{\mathcal{X}} \times \mathcal{X}^B \to \mathbb{R}^B$
**Output:** Batch label-free feature importance $b_i(\boldsymbol{f}, \boldsymbol{X})$.
Define batch auxiliary function $g_{\boldsymbol{X}}$ as a wrapper around $\boldsymbol{f}$ according to (2).
Compute the label-free importance scores for the batch $b_i(\boldsymbol{f}, \boldsymbol{X}) \leftarrow a_i(g_{\boldsymbol{X}}, \boldsymbol{X})$.

---

## B.2    Label-Free Example Importance

We detail the label-free example importance methods used in our experiments in Table 4:

Table 4: Example Importance Methods.

| Method Type | Method | Ref. | Label-Free Expression |
|---|---|---|---|
| Loss-Based | Influence Functions | Koh and Liang [2017] | $c^n(\boldsymbol{f},\boldsymbol{x}) = \frac{1}{N}\left\langle \boldsymbol{\nabla}_{\boldsymbol{\theta}_r}L(\boldsymbol{x},\boldsymbol{\theta}_*), \boldsymbol{H}_{\boldsymbol{\theta}_r}^{-1}\boldsymbol{\nabla}_{\boldsymbol{\theta}_r}L(\boldsymbol{x}^n,\boldsymbol{\theta}_*)\right\rangle_{\Theta_r}$ |
|  | TracIn | Pruthi et al. [2020] | $c^n(\boldsymbol{f},\boldsymbol{x}) = \sum_{t=1}^T \eta_t \left\langle \boldsymbol{\nabla}_{\boldsymbol{\theta}_t}L(\boldsymbol{x},\boldsymbol{\theta}_t), \boldsymbol{\nabla}_{\boldsymbol{\theta}_t}L(\boldsymbol{x}^n,\boldsymbol{\theta}_t)\right\rangle_{\Theta_r}$ |
| Representation-Based | Deep K-Nearest Neighbours | Papernot and McDaniel [2018] | $c^n(\boldsymbol{f},\boldsymbol{x}) = \boldsymbol{1}\left[n \in \text{KNN}(\boldsymbol{x})\right] \cdot \kappa\left[\boldsymbol{f}(\boldsymbol{x}^n), \boldsymbol{f}(\boldsymbol{x})\right]$ |
|  | SimplEx | Crabbé et al. [2021] | $c^n(\boldsymbol{f},\boldsymbol{x}) = \arg\min_{\boldsymbol{\lambda}\in[0,1]^N}\left\|\boldsymbol{f}(\boldsymbol{x}) - \sum_{n=1}^N \lambda^n \boldsymbol{f}(\boldsymbol{x}^n)\right\|_{\mathcal{H}}$ |
|  |  |  | s.t. $\sum_{n=1}^N \lambda^n = 1$ |

where $\boldsymbol{\theta}_r$ are the parameters of the black-box $\boldsymbol{f}$. Our implementation closely follows the above references with some subtle differences. For completeness, we detail the algorithm used for each method. We start with Influence Functions in Algorithm 2.

---

**Algorithm 2** Label-Free Influence Functions

---

**Input:** Test input $\boldsymbol{x} \in \mathcal{X}$, Black-box $\boldsymbol{f} : \mathcal{X} \to \mathcal{H}$, Optimized parameters $\boldsymbol{\theta}_* \in \Theta$, Relevant black-box parameters $\boldsymbol{\theta}_{\mathrm{r}} \in \Theta_{\mathrm{r}}$, Loss function $L : \mathcal{X} \times \Theta \to \mathbb{R}$ used to train the black-box, Training set $\mathcal{D}_{\mathrm{train}} = \{\boldsymbol{x}^n \mid n \in [N]\}$, Number of samples $S \in \mathbb{N}^*$, Number of recursions $R \in \mathbb{N}^*$.
**Output:** Label-free influence functions $c^n(\boldsymbol{f}, \boldsymbol{x})$.
Initialize $\boldsymbol{v}_0 \leftarrow \boldsymbol{\nabla}_{\boldsymbol{\theta}_{\mathrm{r}}} L(\boldsymbol{x}, \boldsymbol{\theta}_*)$.
Initialize $\boldsymbol{v} \leftarrow \boldsymbol{v}_0$.
**for** recursion in $[R]$ **do**
    Sample $S$ training points $\boldsymbol{x}^{n_1}, \ldots, \boldsymbol{x}^{n_S}$ from the training set $\mathcal{D}_{\mathrm{train}}$.
    Make a Monte-Carlo estimation of the training loss Hessian $\boldsymbol{H} \leftarrow \frac{1}{S} \boldsymbol{\nabla}^2_{\boldsymbol{\theta}_{\mathrm{r}}} \sum_{s=1}^{S} L(\boldsymbol{x}^{n_s}, \boldsymbol{\theta}_*)$.
    Update the estimate for the inverse Hessian-vector product $\boldsymbol{v} \leftarrow \boldsymbol{v}_0 + (\boldsymbol{I} - \boldsymbol{H})\boldsymbol{v}$.
**end for**
Compute the influence function $c^n(\boldsymbol{f}, \boldsymbol{x}) \leftarrow \boldsymbol{v}^{\intercal} \boldsymbol{\nabla}_{\boldsymbol{\theta}_{\mathrm{r}}} L(\boldsymbol{x}^n, \boldsymbol{\theta}_*)$

---

This implementation follows the original implementation by Koh and Liang [2017] that leverages the literature on second-order approximation techniques Pearlmutter [1994], Agarwal et al. [2016]. Note that Monte-Carlo estimations of the Hessian quickly become expensive when the number of model parameters grows. Due to our limited computing infrastructure, we limit the number of recursions to $R = 100$. Furthermore, we only compute influence functions for smaller subsets of the training and testing set. The label-free version of TracIn is described in Algorithm 3.

---

**Algorithm 3** Label-Free TracIn

---

**Input:** Test input $\boldsymbol{x} \in \mathcal{X}$, Black-box $\boldsymbol{f} : \mathcal{X} \to \mathcal{H}$, Checkpoint parameters $\{\boldsymbol{\theta}_t \in \Theta \mid t \in [T]\}$, Relevant black-box parameters $\boldsymbol{\theta}_{\mathrm{r}} \in \Theta_{\mathrm{r}}$, Loss function $L : \mathcal{X} \times \Theta \to \mathbb{R}$ used to train the black-box, Training set $\mathcal{D}_{\mathrm{train}} = \{\boldsymbol{x}^n \mid n \in [N]\}$, Checkpoint learning rates $\{\eta_t \in \mathbb{R} \mid t \in [T]\}$.
**Output:** Label-free TraceIn scores $c^n(\boldsymbol{f}, \boldsymbol{x})$.
Initialize $c \leftarrow 0$.
**for** $t$ in $[T]$ **do**
    Update the estimate $c \leftarrow c + \eta_t \boldsymbol{\nabla}^{\intercal}_{\boldsymbol{\theta}_{\mathrm{r}}} L(\boldsymbol{x}^n, \boldsymbol{\theta}_t) \boldsymbol{\nabla}_{\boldsymbol{\theta}_{\mathrm{r}}} L(\boldsymbol{x}, \boldsymbol{\theta}_t)$.
**end for**
Return the score $c^n(\boldsymbol{f}, \boldsymbol{x}) \leftarrow c$

---

In our implementation, we create a checkpoint after each interval of 10 epochs during training. When it comes to DKNN, the formula (5) can be computed explicitly without following a particular procedure. In our implementation, we replaced the kernel function by an inverse distance $\kappa(\boldsymbol{h}_1, \boldsymbol{h}_2) = \|\boldsymbol{h}_1 - \boldsymbol{h}_2\|_{\mathcal{H}}^{-1}$. Further, to make it more fair with other baselines that assign a score to each examples (and not only to $K \in \mathbb{N}^*$ examples), we removed the indicator in (5): $\mathbf{1}(n \in \mathrm{KNN}(x)) \mapsto 1$. In this way, the $K$ most important examples always correspond to the $K$ nearest neighbours. Finally, the label-free version of SimplEx is described in Algorithm 4.

---

**Algorithm 4** SimplEx

---

**Input:** Test input $\boldsymbol{x} \in \mathcal{X}$, Black-box $\boldsymbol{f} : \mathcal{X} \to \mathcal{H}$, Training set $\mathcal{D}_{\mathrm{train}} = \{\boldsymbol{x}^n \mid n \in [N]\}$, Number of epochs $E \in \mathbb{N}^*$.

**Output:** SimplEx scores $c^n(\boldsymbol{f}, \boldsymbol{x})$.
Initialize weights $(w^n)_{n=1}^N \leftarrow \mathbf{0}$.
**for** epoch in $[E]$ **do**
    Normalize weights $(\lambda^n)_{n=1}^N \leftarrow \mathrm{Softmax}\left[(w^n)_{n=1}^N\right]$
    Estimate Error $\mathcal{L} = \left\| \boldsymbol{f}(\boldsymbol{x}) - \sum_{n=1}^N \lambda^n \boldsymbol{f}(\boldsymbol{x}^n) \right\|_{\mathcal{H}}$.
    Update weights with Adam $(w^n)_{n=1}^N \leftarrow \mathrm{Adam\ Step}(\mathcal{L})$.
**end for**
Return the score $c^n(\boldsymbol{f}, \boldsymbol{x}) \leftarrow \lambda^n$

---

Note that our implementation of SimplEx is identical to the original one. It relies on an Adam Kingma and Ba [2014] with the default Pytorch parameters (learning rate $= .001, \beta_1 = .9, \beta_2 = .999, \epsilon = 10^{-8}$, weight decay $= 0$).

In this section, we have presented many label-free implementations of feature and example importance methods. For some types of explanations, like counterfactual explanations Wachter et al. [2017], the label plays an essential role. Hence, it does not always make sense to extend an explanation to the label-free setting.

# C   Experiments Details

In this appendix, we provide further details to support the experiments described in Section 4. All our experiments have been performed on a machine with Intel(R) Core(TM) i5-8600K CPU @ 3.60GHz [6 cores] and Nvidia GeForce RTX 2080 Ti GPU. Our implementation is done with Python 3.8 and Pytorch 1.10.0.

## C.1   Consistency Checks

We provide some details for the experiments in Section 4.1.

**Model**   The architecture for the autoencoder is described in Table 5. The autoencoder is trained to minimize the denoising loss described in Section 4.1 for 100 epochs with patience 10 by using Pytorch's Adam with hyperparameters: learning rate $= .001, \beta_1 = .9, \beta_2 = .999, \epsilon = 10^{-8}$, weight decay $= 10^{-5}$.

Table 5: MNIST Autoencoder Architecture.

| Component | Layer Type | Hyperparameters | Activation Function |
|---|---|---|---|
| | Conv2d | Input Channels:1 ; Output Channels:8 ; Kernel Size:3 ; Stride:2 ; Padding:1 | ReLU |
| | Conv2d | Input Channels:8 ; Output Channels:16 ; Kernel Size:3 ; Stride:2 ; Padding:1 | ReLU |
| | BatchNorm | Input Channels:16 | ReLU |
| Encoder | Conv2d | Input Channels:16 ; Output Channels:32 ; Kernel Size:3 ; Stride:2 ; Padding:0 | ReLU |
| | Flatten | Start Dimension:1 | |
| | Linear | Input Dimension: 288 ; Output Dimension: 128 | ReLU |
| | Linear | Input Dimension: 128 ; Output Dimension: 4 | |
| | Linear | Input Dimension: 4 ; Output Dimension: 128 | ReLU |
| | Linear | Input Dimension: 128 ; Output Dimension: 288 | ReLU |
| | Unflatten | Dimension:1 ; Unflatten Size:(32, 3, 3) | |
| Decoder | ConvTranspose2d | Input Channels:32 ; Output Channels:16 ; Kernel Size:3 ; Stride:2 ; Output Padding:0 | |
| | BatchNorm | Input Channels:16 | ReLU |
| | ConvTranspose2d | Input Channels:16 ; Output Channels:8 ; Kernel Size:3 ; Stride:2 ; Output Padding:1 | |
| | BatchNorm | Input Channels:8 | ReLU |
| | ConvTranspose2d | Input Channels:8 ; Output Channels:1 ; Kernel Size:3 ; Stride:2 ; Output Padding:1 | Sigmoid |

**Feature Importance**   As a baseline for the feature importance methods, we use a black image $\bar{x} = 0$.

## C.2   Pretext Task Sensitivity

We provide some details for the experiments in Section 4.2.

**Models**   All the autoencoders have the architecture described in Table 5. The classifier has all the layers from the encoder in Table 5 with an extra linear layer (Input Dimension:4 ; Output Dimension:10 ; Activation:Softmax) that converts the latent representations to class probabilities. The classifier is trained to minimize the cross-entropy loss $L_{CE}(\boldsymbol{x}, \boldsymbol{y}) = -\boldsymbol{y} \odot \log\left[\boldsymbol{f}_l \circ \boldsymbol{f}_e(\boldsymbol{x})\right]$, where $\boldsymbol{y}$ is the one-hot encoded label associated to the training example $\boldsymbol{x}$. All the models are trained to minimize their objective for 100 epochs with patience 10 by using Pytorch's Adam with hyperparameters: learning rate $= .001, \beta_1 = .9, \beta_2 = .999, \epsilon = 10^{-8}$, weight decay $= 10^{-5}$.

**Feature Importance**   As a baseline for the feature importance methods, we use a black image $\bar{x} = 0$.

**Metrics**   We use the Pearson coefficient to measure the correlation between two importance scores given a random test example and a random feature/training example. In our experiment, we compute the Pearson correlation between the label-free feature importance scores $b_i$ outputted by two different encoder $\boldsymbol{f}_{e1}, \boldsymbol{f}_{e2} : \mathcal{X} \to \mathcal{H}$:

$$r_{\text{feat.}}(\boldsymbol{f}_{e1}, \boldsymbol{f}_{e2}) = \frac{\text{cov}_{\boldsymbol{X}, I}\left[b_I(\boldsymbol{f}_{e1}, \boldsymbol{X}), b_I(\boldsymbol{f}_{e2}, \boldsymbol{X})\right]}{\sigma_{\boldsymbol{X}, I}\left[b_I(\boldsymbol{f}_{e1}, \boldsymbol{X})\right] \sigma_{\boldsymbol{X}, I}\left[b_I(\boldsymbol{f}_{e2}, \boldsymbol{X})\right]}$$
$$\boldsymbol{X} \sim \text{Empirical Distribution}(\mathcal{D}_{\text{test}}), I \sim \text{Uniform}([d_X]),$$

where cov denotes the covariance between two random variables and and $\sigma$ denotes the standard deviation of a random variable. Similarly, for label-free example importance scores $c^n$:

$$r_{\text{ex.}}(\boldsymbol{f}_{e1}, \boldsymbol{f}_{e2}) = \frac{\text{cov}_{\boldsymbol{X}, I}\left[c^I(\boldsymbol{f}_{e1}, \boldsymbol{X}), c^I(\boldsymbol{f}_{e2}, \boldsymbol{X})\right]}{\sigma_{\boldsymbol{X}, I}\left[c^I(\boldsymbol{f}_{e1}, \boldsymbol{X})\right] \sigma_{\boldsymbol{X}, I}\left[c^I(\boldsymbol{f}_{e2}, \boldsymbol{X})\right]}$$
$$\boldsymbol{X} \sim \text{Empirical Distribution}(\mathcal{D}_{\text{test}}), I \sim \text{Uniform}(\mathcal{J}),$$

where $\mathcal{J} \subset [N]$ is the indices of the sampled training examples for which the example importance is computed. Those two Pearson correlation coefficients are the one that we report in Tables 1 and 2.

**Supplementary Examples**   To check that the qualitative analysis from Section 4.2 extends beyond the examples showed in the main paper, the reader can refer to Figures 9 and 10.

### C.3   Challenging our assumptions with disentangled VAEs

We provide some details for the experiments in Section 4.3.

Table 6: MNIST Variational Autoencoder Architecture.

| Component | Layer Type | Hyperparameters | Activation Function |
|---|---|---|---|
| | Conv2d | Input Channels:1 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Padding:1 | ReLU |
| | Conv2d | Input Channels:32 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Padding:1 | ReLU |
| | Conv2d | Input Channels:32 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Padding:1 | ReLU |
| Encoder | Flatten | Start Dimension:1 | |
| | Linear | Input Dimension: 512 ; Output Dimension: 256 | ReLU |
| | Linear | Input Dimension: 256 ; Output Dimension: 256 | ReLU |
| | Linear | Input Dimension: 256 ; Output Dimension: 6 | ReLU |
| Reparametrization Trick | | The output of the encoder contains $\boldsymbol{\mu}$ and $\log \boldsymbol{\sigma}$. | |
| | | The latent representation is then generated via $\boldsymbol{h} = \boldsymbol{\mu}(\boldsymbol{x}) + \boldsymbol{\sigma}(\boldsymbol{x}) \odot \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim \mathcal{N}(0, I)$ | |
| | Linear | Input Dimension: 3 ; Output Dimension: 256 | ReLU |
| | Linear | Input Dimension: 256 ; Output Dimension: 256 | ReLU |
| | Linear | Input Dimension: 256 ; Output Dimension: 512 | ReLU |
| Decoder | Unflatten | Dimension:1 ; Unflatten Size:(32, 4, 4) | |
| | ConvTranspose2d | Input Channels:32 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Output Padding:1 | ReLU |
| | ConvTranspose2d | Input Channels:32 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Output Padding:1 | ReLU |
| | ConvTranspose2d | Input Channels:32 ; Output Channels:1 ; Kernel Size:4 ; Stride:2 ; Output Padding:1 | Sigmoid |

Table 7: dSprites Variational Autoencoder Architecture.

| Component | Layer Type | Hyperparameters | Activation Function |
|---|---|---|---|
| | Conv2d | Input Channels:1 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Padding:1 | ReLU |
| | Conv2d | Input Channels:32 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Padding:1 | ReLU |
| | Conv2d | Input Channels:32 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Padding:1 | ReLU |
| Encoder | Conv2d | Input Channels:32 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Padding:1 | ReLU |
| | Flatten | Start Dimension:1 | |
| | Linear | Input Dimension: 512 ; Output Dimension: 256 | ReLU |
| | Linear | Input Dimension: 256 ; Output Dimension: 256 | ReLU |
| | Linear | Input Dimension: 256 ; Output Dimension: 12 | ReLU |
| Reparametrization Trick | | The output of the encoder contains $\boldsymbol{\mu}$ and $\log \boldsymbol{\sigma}$. | |
| | | The latent representation is then generated via $\boldsymbol{h} = \boldsymbol{\mu}(\boldsymbol{x}) + \boldsymbol{\sigma}(\boldsymbol{x}) \odot \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim \mathcal{N}(0, I)$ | |
| | Linear | Input Dimension: 6 ; Output Dimension: 256 | ReLU |
| | Linear | Input Dimension: 256 ; Output Dimension: 256 | ReLU |
| | Linear | Input Dimension: 256 ; Output Dimension: 512 | ReLU |
| Decoder | Unflatten | Dimension:1 ; Unflatten Size:(32, 4, 4) | |
| | Conv2d | Input Channels:32 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Padding:1 | ReLU |
| | ConvTranspose2d | Input Channels:32 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Output Padding:1 | ReLU |
| | ConvTranspose2d | Input Channels:32 ; Output Channels:32 ; Kernel Size:4 ; Stride:2 ; Output Padding:1 | ReLU |
| | ConvTranspose2d | Input Channels:32 ; Output Channels:1 ; Kernel Size:4 ; Stride:2 ; Output Padding:1 | Sigmoid |

**Model**   The architecture of the MNIST VAE is described in Table 6 and those of the dSprites VAE is described in Table 7. Both of these architectures are reproductions of the VAEs from Burgess et al. [2018]. The $\beta$-VAE is trained to minimize the objective $L_\beta(\boldsymbol{x}, \theta, \phi) = \mathbb{E}_{q_\phi(\boldsymbol{h}|\boldsymbol{x})} [\log p_\theta(\boldsymbol{x} \mid \boldsymbol{h})] - \beta D_{\mathrm{KL}} [q_\phi(\boldsymbol{h} \mid \boldsymbol{x}) \mid\mid p(\boldsymbol{h})]$, where $q_\phi(\boldsymbol{h} \mid \boldsymbol{x})$ is the distribution underlying the reparametrized encoder output, $p_\theta(\boldsymbol{x} \mid \boldsymbol{h})$ is the distribution underlying the decoder output, $p(\boldsymbol{h})$ is the density associated to isotropic unit Gaussian $\mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$ underlying $\boldsymbol{\epsilon}$ and $D_{\mathrm{KL}}$ is the KL-divergence. The objective of the TC-VAE is th same as in Chen et al. [2018]. We refer the reader to the original paper for the details. All the VAEs are trained to minimize their objective for 100 epochs with patience 10 by using Pytorch's Adam with hyperparameters: learning rate $= .001, \beta_1 = .9, \beta_2 = .999, \epsilon = 10^{-8}$, weight decay $= 10^{-5}$.

**Feature Importance**   As a baseline for the feature importance methods, we use a black image $\bar{\boldsymbol{x}} = \boldsymbol{0}$.

**Metrics** We use the Pearson coefficient to measure the correlation between two importance scores given a random test example and a random feature. In this case, one correlation coefficient can be computed for each couple $(i, j) \in [d_H]^2$ of latent units:

$$r_{ij} = \frac{\text{cov}_{\boldsymbol{X},I}\left[a_I(\mu_i, \boldsymbol{X}), a_I(\mu_j, \boldsymbol{X})\right]}{\sigma_{\boldsymbol{X},I}\left[a_I(\mu_i, \boldsymbol{X})\right]\sigma_{\boldsymbol{X},I}\left[a_I(\mu_j, \boldsymbol{X})\right]}$$

$$\boldsymbol{X} \sim \text{Empirical Distribution}(\mathcal{D}_{\text{test}}), I \sim \text{Uniform}([d_X]),$$

where $\mu_i$ is the $i$-th component of the expected representation computed by the encoder for all $i \in [d_H]$. To have an overall measure of correlation between the VAE units, we sum over all pairs of distinct latent units:

$$r = \frac{1}{d_H(d_H - 1)} \sum_{i=1, i \neq j}^{d_H} \sum_{j=1}^{d_H} r_{ij}$$

This averaged correlation coefficient is the one that we report in Figure 6. In our quantitative analysis, we also report the Spearman rank correlation between $\beta$ and $r$. Concretely, this is done by performing the experiment $M \in \mathbb{N}^*$ times for different values $\beta_1, \ldots, \beta_M$ of $\beta$. We then measure the correlation coefficients $r_1, \ldots, r_M$ associated to each experiment. The Spearman rank correlation coefficient can be computed form this data:

$$\rho = \frac{\text{cov}_{M_1, M_2}\left[\text{rank}(r_{M_1}), \text{rank}(\beta_{M_2})\right]}{\sigma_{M_1}[\text{rank}(r_{M_1})]\sigma_{M_2}[\text{rank}(\beta_{M_2})]}$$

$$M_1, M_2 \sim \text{Uniform}([M]).$$

This coefficient $\rho$ ranges from $-1$ to $1$, where $\rho = -1$ corresponds to a perfect monotonically decreasing relation, $\rho = 0$ corresponds to the absence of monotonic relation and $\rho = 1$ corresponds to a perfect monotonically increasing relation.

**Entropy** To further our quantitative analysis of the VAEs, we introduce a new metric called *entropy*. The purpose of this metric is to measure how the saliency for each feature is distributed across the different latent units. In particular, we would like to be able to distinguish the case where all latent units are sensitive to a feature and the case where only one latent unit is sensitive to a feature. As we have done previously, we can compute the importance score $a_i(\mu_j, \boldsymbol{x})$ of each feature $x_i$ from $\boldsymbol{x} \in \mathcal{X}$ for a latent unit $j \in [d_H]$. For each latent unit $j \in [d_H]$, we define the proportion of attribution as

$$p_j(i, \boldsymbol{x}) = \frac{|a_i(\mu_j, \boldsymbol{x})|}{\sum_{k=1}^{d_H} |a_i(\mu_k, \boldsymbol{x})|}.$$

This corresponds to the fraction of the importance score attributed to $j$ for feature $i$ and example $\boldsymbol{x}$. Note that this quantity is well defined if at least one of the $a_i(\mu_k, \boldsymbol{x}), k \in [d_H]$ is non-vanishing. Hence, we only consider the features $i \in [d_X]$ that are salient for at least one latent unit. We can easily check that $\sum_{j=1}^{d_H} p_j(i, \boldsymbol{x}) = 1$ by construction. This means that the proportions of attributions can be interpreted as probabilities of saliency. This allows us to define an entropy that summarizes the distribution over the latent units:

$$S(i, \boldsymbol{x}) = -\sum_{j=1}^{d_H} p_j(i, \boldsymbol{x}) \ln p_j(i, \boldsymbol{x}).$$

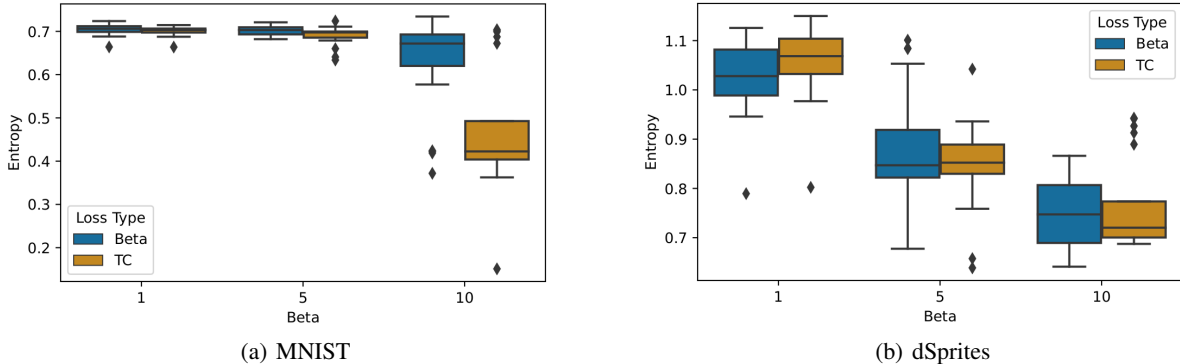This entropy is analogous to Shannon's entropy Shannon [1948]. It can be checked easily that this entropy is minimal ($S^{\min} = 0$) whenever only one latent unit $j \in [d_H]$ is sensitive to feature $i$: $p_j(i, \boldsymbol{x}) = 1$. Conversely, it is well known Cover and Thomas [2005] that the entropy is maximal ($S^{\max} = \ln d_H$) whenever the distribution is uniform over the latent units: $p_j(i, \boldsymbol{x}) = 1/d_H$ for all $j \in [d_H]$. In short: the entropy is low when mostly one latent unit is sensitive to the feature of interest and high when several latent units are sensitive to the feature of interest. Clearly, the former situation is more desirable if we want to distinguish the different latent units. For each VAE, we evaluate the average entropy

$$S = \mathbb{E}_{\boldsymbol{X},I}\left[S(I, \boldsymbol{X})\right]$$

$$\boldsymbol{X} \sim \text{Empirical Distribution}(\mathcal{D}_{\text{test}}), I \sim \text{Uniform}[\mathcal{I}(\boldsymbol{X})],$$

where $\mathcal{I}(\boldsymbol{x}) = \{i \in [d_X] \mid a_i(\mu_k, \boldsymbol{x}) \neq 0 \text{ for at least one } k \in [d_H]\}$ is the set of features that are salient for at least one latent unit. We measure the average entropy for each VAE and report the results as a function of $\beta$ in Figure 7.

We clearly see that the entropy decreases as the disentanglement strength $\beta$ increases for both MNIST (Spearman $\rho = -.56$) an dSprites (Spearman $\rho = -.76$). This means that disentangling has the effect of distributing the saliency

|     |     |
| --- | --- |
| (a) MNIST | (b) dSprites |

Figure 7: Entropy of saliency maps for different values of $\beta$.

over fewer units. This brings a nice complement to the quantitative analysis that we have conducted in Section 4.3: although increasing disentanglement does not make the latent units focus on different parts of the image (since the correlation does not decrease significantly), it does decrease the number of latent units that are simultaneously sensitive to a given part of the image (since the entropy decreases substantially). These two phenomena are not incompatible with each other. For instance, we see that the 6-th latent unit seems inactive in comparison with the other latent units in Figure 12. In fact, this latent unit might perfectly pay attention to the same parts of the image as the other units and, hence, be correlated. What distinguishes this unit from the others is that the feature importance scores on its saliency map are significantly smaller (we cannot appreciate it by plotting the saliency maps on the same scale) and, hence, reduces the entropy. Finally, we note that the entropies from Figure 7 remain fairly close to their maximal value ($S^{max} = \ln 3 \approx 1.1$ for MNIST and $S^{max} = \ln 6 \approx 1.8$ for dSprites). This means that the VAEs have several active units for each pixel.

**Supplementary Examples** To check that the qualitative analysis from Section 4.3 extends beyond the examples showed in the main paper, the reader can refer to Figures 11 and 12. These saliency maps are produced with the same VAEs we selected for our qualitative analysis.

## D    Supplementary Experiments with Medical Time Series

In this appendix, we perform further consistency checks with medical time series data. Those consistency checks mirror the one described in Section 4.1.

**Dataset** In this experiment, we work with the ECG5000 dataset Goldberger et al. [2000]. This dataset $\mathcal{D}$ contains 5000 univariate time series[5] $(x_t)_{t=1}^T \in \mathcal{X}^T$ describing the heartbeat of a patient. Each time series describes a single heartbeat with a resolution of $T = 140$ time steps. For the sake of notation, we will represent univariate time series by vectors: $\boldsymbol{x} = (x_t)_{t=1}^T$. Each time series comes with a label $y \in \{0, 1\}$ indicating if the heartbeat is normal ($y = 0$) or abnormal ($y = 1$). Since it is laborious to manually annotate 5000 time series, those labels were generated automatically. Of course, those labels are not going to be used in training our model. We only use the labels to perform consistency checks once the model has been trained.

### D.1    Feature Importance

**Model** We train a reconstruction autoencoder $\boldsymbol{f}_d \circ \boldsymbol{f}_e : \mathcal{X}^T \to \mathcal{X}^T$ that consists in an encoder $\boldsymbol{f}_e : \mathcal{X}^T \to \mathcal{H}$ and a decoder $\boldsymbol{f}_d : \mathcal{H} \to \mathcal{X}^T \ d_H = 64$. This model is trained with a training set $\mathcal{D}_{\text{train}}$ of 2919 time series from $\mathcal{D}$ that correspond to normal heartbeats: $y = 0$. The model is trained to minimize the reconstruction loss $L_{\text{rec}}(\boldsymbol{x}) = \sum_{t=1}^T |x_t - [\boldsymbol{f}_d \circ \boldsymbol{f}_e(\boldsymbol{x})]_t|$. The autoencoder is trained for 150 epochs with patience 10 by using Pytorch's Adam with hyperparameters: learning rate $= .001, \beta_1 = .9, \beta_2 = .999, \epsilon = 10^{-8}$, weight decay $= 0$. Its detailed architecture is presented in Table 8.

---

[5]Note that $t$ is used to index the time series steps, as opposed to model checkpoints in Section 3.1.

Table 8: ECG5000 Autoencoder Architecture.

| Component | Layer Type | Hyperparameters | Activation Function |
|---|---|---|---|
| Encoder | LSTM | Input Size:1 ; Hidden Size: $2 \cdot d_H$ | |
| | LSTM | Input Size: $2 \cdot d_H$ ; Hidden Size:$d_H$ | |
| Representation | | The latent representation $h$ is given by the output of the second LSTM at the last time step. | |
| | | This repreentation is copied at each time step to be a valid imput for the first decoder LSTM. | |
| Decoder | LSTM | Input Size: $d_H$ ; Hidden Size: $d_H$ | |
| | LSTM | Input Size: $d_H$ ; Hidden Size: $2 \cdot d_H$ | |
| | Linear | Input Dimension: $2 \cdot d_H$ ; Output Dimension: 1 | |

**Setup**    We test our model with the abnormal heartbeats $\mathcal{D}_{\text{test}} = \mathcal{D} \setminus \mathcal{D}_{\text{train}}$. We compute the label-free importance score $b_t(\boldsymbol{f}_e, \boldsymbol{x})$ of each time step $t \in [T]$ for building the latent representation of heartbeat $\boldsymbol{x} \in \mathcal{D}_{\text{test}}$. We use the average normal heartbeat as a baseline for the feature importance methods: $\bar{\boldsymbol{x}} = \sum_{\boldsymbol{x} \in \mathcal{D}_{\text{train}}} \boldsymbol{x} / |\mathcal{D}_{\text{train}}|$. We replace the $M$ most important time steps by their baseline value with a mask $\boldsymbol{m} \in \{0, 1\}^T$. For each time series, we measure the latent shift $\|\boldsymbol{f}_e(\boldsymbol{x}) - \boldsymbol{f}_e(\boldsymbol{m} \odot \boldsymbol{x} + (\mathbf{1} - \boldsymbol{m}) \odot \bar{\boldsymbol{x}})\|_{\mathcal{H}}$ induced by masking the most important time steps. We expect this shift to increase with the importance of masked time steps. We report the average shift over the testing set (2081 time series) for several values of $M$ and feature importance methods in Figure 8(a).
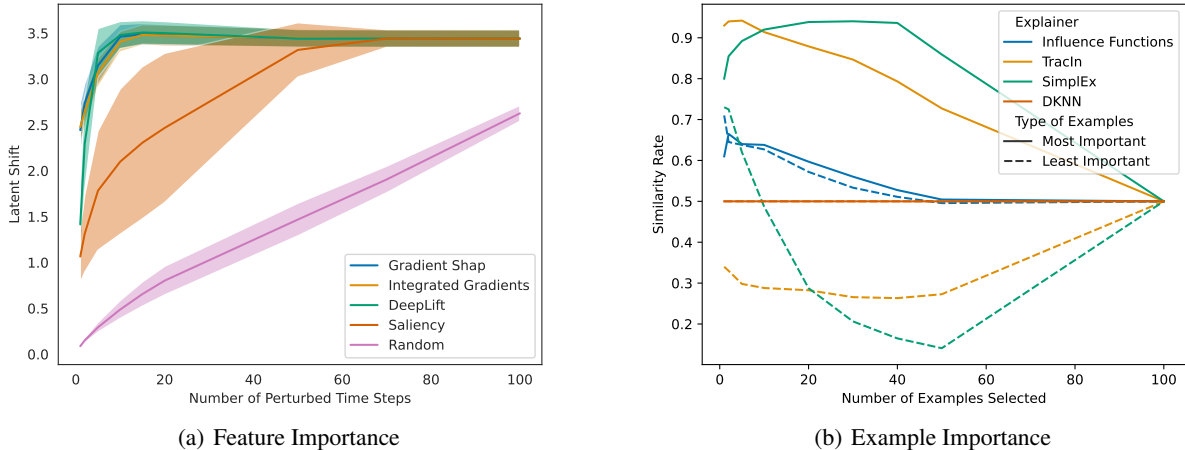


(a) Feature Importance

(b) Example Importance

Figure 8: Consistency checks for the ECG5000 dataset.

**Results**    Again, all the label-free extensions pass the consistency check. This is because the relevant time steps have a significantly larger impact than random time steps when perturbed. In this case, more sophisticate feature importance methods offer better performances than Saliency.

### D.2    Example Importance

**Model**    In this experiment, we use the autoencoder described in Table 8 with $d_H = 32$. The whole training process is identical to the one from the previous section with one difference: $\mathcal{D}_{\text{train}}$ and $\mathcal{D}_{\text{test}}$ are now obtained with a random split of $\mathcal{D}$ ($80\% - 20\%$). This means that both subsets contain normal and abnormal heartbeats.

**Setup**    We sample $N = 100$ training time series $\boldsymbol{x}^n \in \mathcal{D}_{\text{train}}, n \in [N]$ without replacement and by respecting a $50\% - 50\%$ proportion for each label $y = 0, 1$. We compute the importance score $c^n(\boldsymbol{f}_e, \boldsymbol{x})$ of each training example $\boldsymbol{x}^n$ for predicting the latent representation of the test time series $\boldsymbol{x} \in \mathcal{D}_{\text{test}}$. We select the $M$ most important training examples $\boldsymbol{x}^{n_1}, \ldots, \boldsymbol{x}^{n_M}$. We compare their ground truth label $y^{n_1}, \ldots, y^{n_M}$ to the label $y$ of $\boldsymbol{x}$. We compute the similarity rates $\sum_{m=1}^{M} \delta_{y, y^{n_m}} / M$, where $\delta$ denotes the Kronecker delta. We reproduce the above steps for the $M$ least

important examples. If the encoder meaningfully represents the data, we expect the similarity rate of the most important examples to be higher than for the least important examples. We report the distribution of similarity rates across 100 test time-series for various values of $M$ and example importance methods in Figure 8(b).

**Results**    In this case, only SimplEx and TraceIn pass the consistency check. Other methods don't show a significant difference between the most and least important examples. We note that each successful method is a refinement over a method that failed the consistency check (TraceIn was introduced as a refinement over Influence functions and SimplEx as a refinement over DKNN). A study involving more than 100 examples would have been necessary to obtain more significant results. We had to restrict our analysis to few examples due to the high computational cost of approximating Influence Functions for time series models.
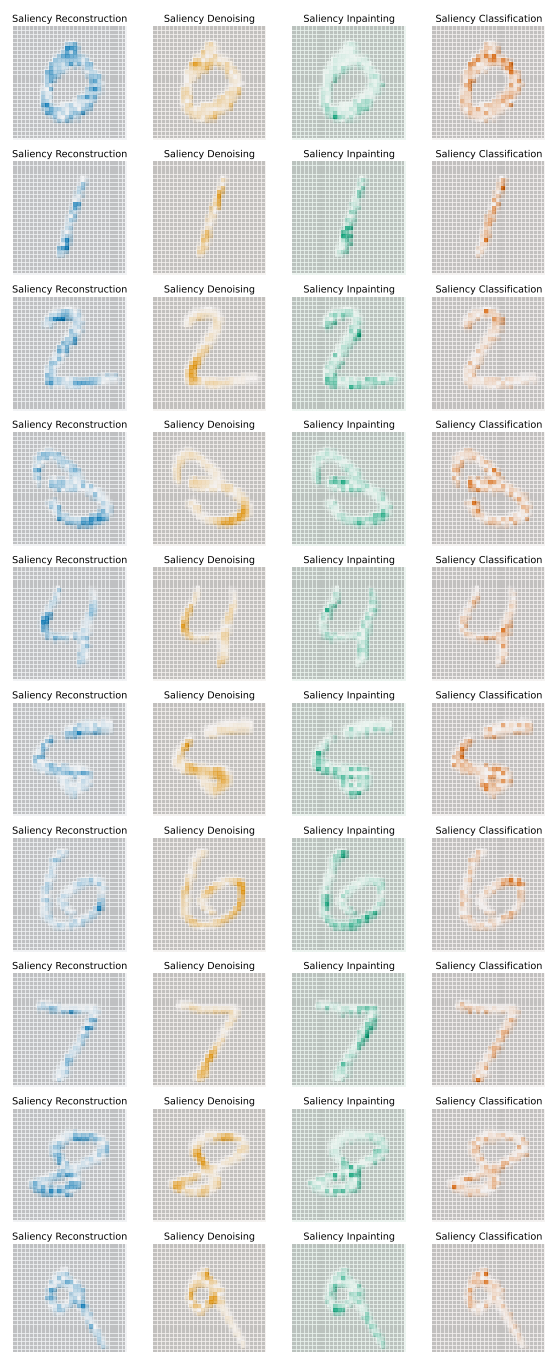
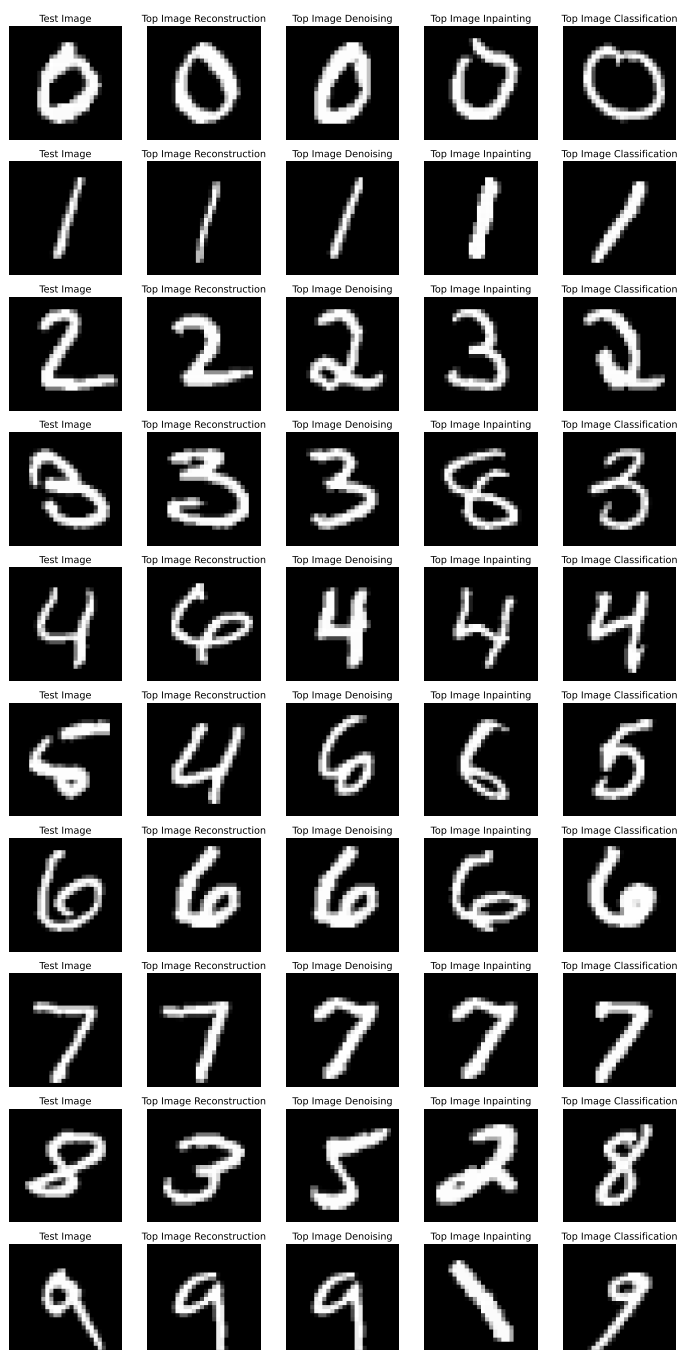Figure 9: Label-free saliency for various pretext tasks.

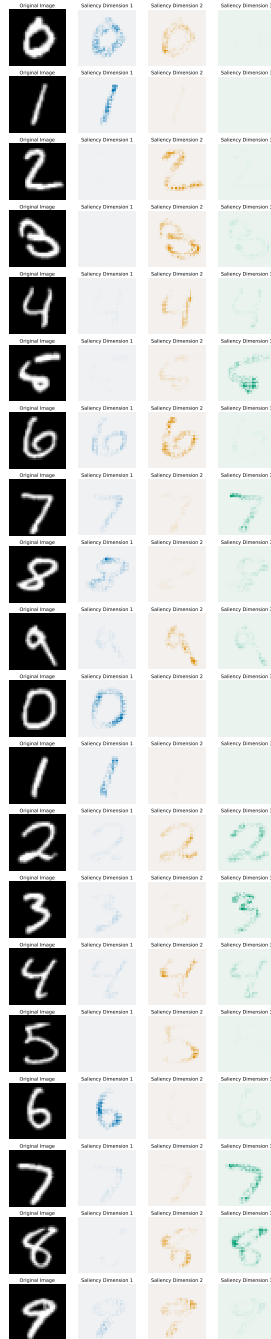Figure 10: Label-free top example for various pretext tasks.
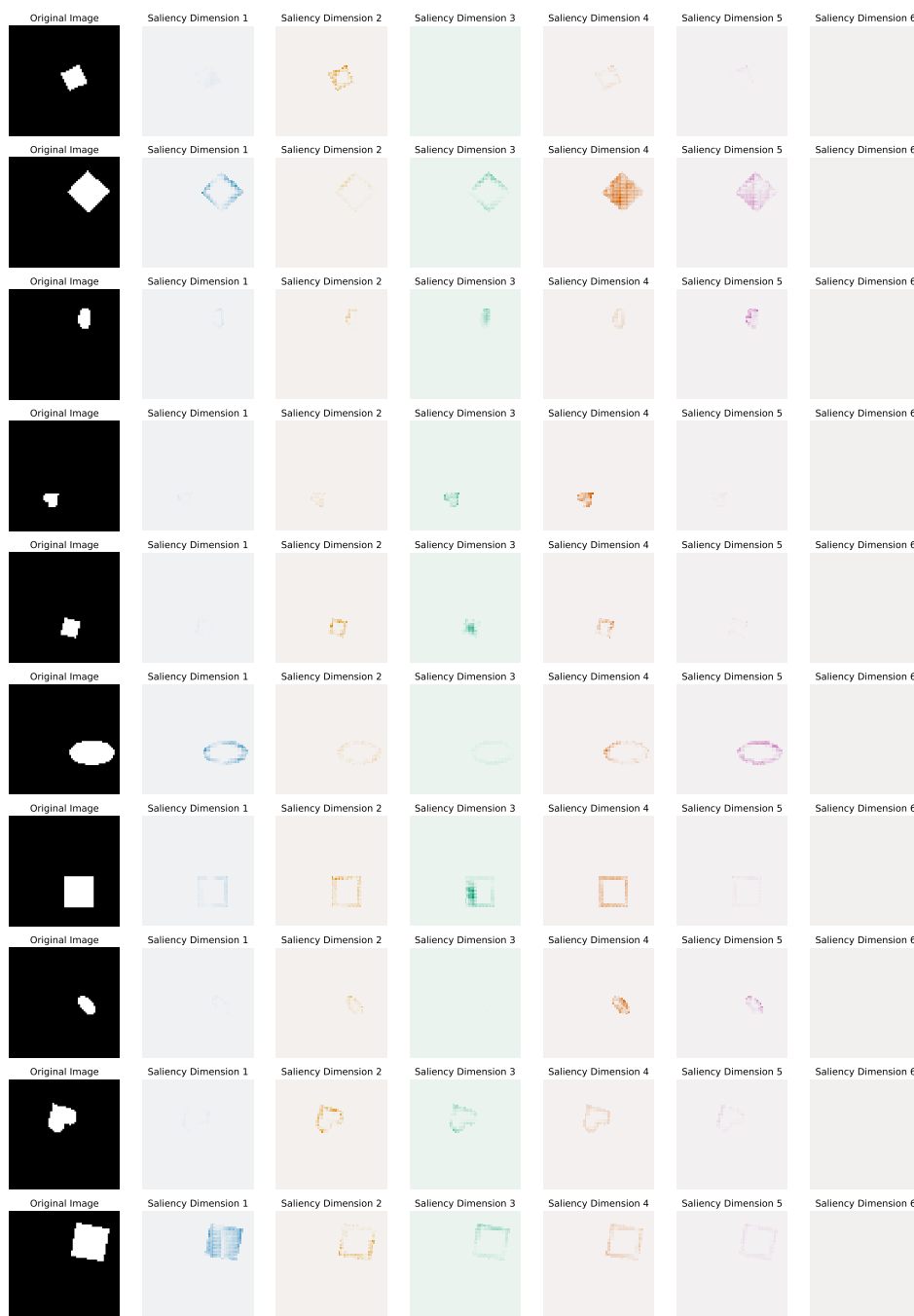
Figure 11: Saliency maps for the latent units of a MNIST VAE.

Figure 12: Saliency maps for the latent units of a dSprites VAE.