

# [Re] Explaining Groups of Points in Low-Dimensional Representations

Damiaan J. W. Reijnaers<sup>1, ID</sup>, Daniël B. van de Pavert<sup>1, ID</sup>, Giguru Scheuer<sup>1, ID</sup>, and Liang Huang<sup>1, ID</sup>

<sup>1</sup>Faculty of Science, University of Amsterdam, Science Park 904, 1098 XH Amsterdam, the Netherlands

## Edited by

Koustuv Sinha,  
Jesse Dodge

## Reviewed by

Anonymous Reviewers

## Received

29 January 2021

## Published

27 May 2021

## DOI

10.5281/zenodo.4835056

## Reproducibility Summary

### Scope of Reproducibility

In this paper we present an analysis and elaboration of [1], in which an algorithm is posed by Plumb *et al.* for the purpose of finding human-understandable explanations in terms of given explainable features of input data for differences between groups of points occurring in a lower-dimensional representation of that input data.

### Methodology

We have upgraded the original code provided by the authors such that it is compatible with recent versions of popular deep learning frameworks, namely the TensorFlow 2.x- and PyTorch 1.7.x-libraries. Furthermore, we have created our own implementation of the algorithm in which we have incorporated additional experiments in order to evaluate the algorithm's relevance in the scope of different dimensionality reduction techniques and differently structured data. We have performed the same experiments as described in the original paper using both the upgraded version of the code and our own implementation taking the authors' code and paper as references.

### Results

The results presented in [1] were reproducible, both by using the provided code and our own implementation. Our additional experiments have highlighted several limitations of the explanatory algorithm in question: the algorithm severely relies on the shape and variance of the clusters present in the data (and, if applicable, the method used to label these clusters), and highly non-linear dimensionality reduction algorithms perform worse in terms of explainability.

### What was easy

The authors have provided an implementation<sup>1</sup> that cleanly separates different experiments on different datasets and the core functional methodology. Given a working environment, it is easy to reproduce the experiments performed in [1].

---

Copyright © 2021 D.J.W. Reijnaers et al., released under a Creative Commons Attribution 4.0 International license.

Correspondence should be addressed to Damiaan J. W. Reijnaers (info@damiaanreijnaers.nl)

The authors have declared that no competing interests exist.

Code is available at <https://github.com/damiaanr/fact-ai> – DOI 10.5281/zenodo.4686025. – SWH

swh:1:dir:df71f3541f882fb2474dc177b8eaba31f905d9d.

Open peer review is available at <https://openreview.net/forum?id=hq3TxQK5cox>.

## What was difficult

Minor difficulties were experienced in setting up the required environment for running the code provided by Plumb *et al.* locally (i.e. trivial changes in the code such as the usage of absolute paths and obtaining external dependencies). Evidently, it was time-consuming to rewrite all corresponding code, including the architecture for the variational auto-encoder provided by an external package, `scvis` 0.1.0<sup>2</sup>.

## Communication with original authors

No communication with the original authors was required to reproduce their work.

## 1 Introduction

As AI models are getting more integrated into applications with economic or social implications, the need for *explaining* decisions made by (potentially complex) models is increasing. In many AI applications, big datasets form the basis of a decision-making algorithm [2]. As these datasets often involve data of high dimensionality, the dimensionality of data is, in many different applications, often reduced using *dimensionality reduction* (DR) techniques [3].

Data, and consequently the decisions made by an algorithm that are (in)directly based on that data, often involve some sort of ‘grouping,’ either by utilizing a *clustering method* or by (manually) pre-defined ‘clusters of data.’ The ‘grouping process’ may happen before the dimensionality reduction, for which, assuming well-organized data in which the dimensions correspond to explainable ‘real-life’ *features* of the phenomena measured in the data, the distinguishing characteristics of a certain group becomes relatively explainable as these groups are marked by boundaries in terms of explainable dimensions. However, this process can also happen after the dimensionality of the data has been reduced, which results in the groups being defined in terms of dimensions in a *latent space*. Especially when grouping occurs after dimensionality reduction, different clusters in data often play a key role in the decision-making process of an algorithm – one could, for example, when regarding *credit risk modelling*, point out a group in latent space which ‘poses a high risk’ when provided a loan [4]. Moreover, such groupings arise naturally under the context of (variational) auto-encoders, which are the architecture of preference for many deep learning applications [1, p. 8].

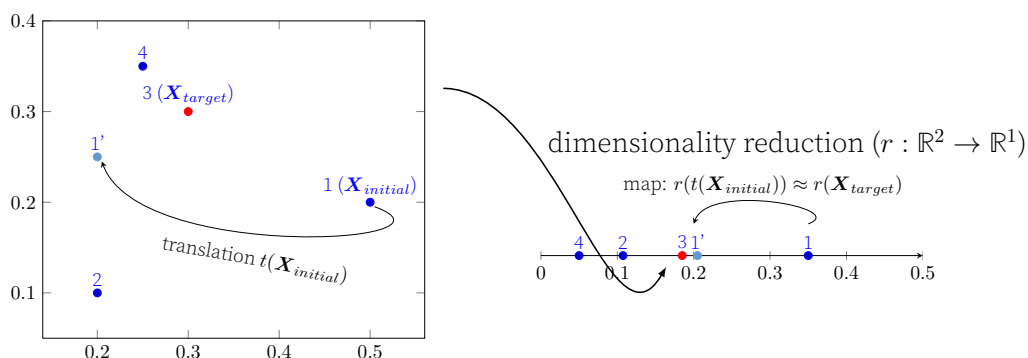
Thus, considering the above-mentioned need for explainable machine decisions, in combination with the increasing use of DR algorithms and the reliance on observed data clusters in abstract latent spaces which are not understandable to the human mind, it is of great relevance and importance to develop methods to explain differences between groups in the light of the application of a certain dimensionality reduction technique.

Elaborating upon the earlier introduced example of credit risk, one could argue that a reasonable explanation for a machine decision is of the kind: “*Your loan was rejected, but if you would have earned €422,86 more per month, your loan would have been accepted.*” This type of explanation is a *counterfactual explanation* – a decision on the basis of a ‘fake,’ counterfactual, ‘world,’ in which features *would have* been shaped differently. The proposed explanatory technique in [1] does exactly this: it reverse-engineers obtained cluster labels in latent space to label corresponding data in original space; then ‘tweaks’

<sup>1</sup>GitHub, *Explaining Low Dimensional Representations*, <https://github.com/GDPlumb/ELDR>, accessed on January 29th, 2021

<sup>2</sup>GitHub, *shahcompbio/scvis: Python package for dimension reduction of high-dimensional biological data*, <https://github.com/shahcompbio/scvis>, accessed on January 22nd, 2021

an initial cluster by translating that cluster in original space, so that it is mapped to (approximately) the area of the target cluster in latent space. Since the dimensions in original space correspond to explainable features—the groups’ characteristics—which are comprehensible for a human mind; a translation that corresponds to merely adding or subtracting values to each of these features can be perfectly explained in ‘human language.’ Since it is desirable for the explanations to be concise, the translation should be as *sparse* as possible: the translation should ‘tweak’ as few dimensions as possible. The intuition behind this idea is visualized in figure 1.



**Figure 1.** Visualisation of the method described by Plumb *et al.* Each point represents an individual cluster. The difference between clusters 1 and 3 are explained. Note that although point 1 maps to 1', which is far from point 3 in original space, point 1' is almost identical to 3 in latent space – this is the goal, as we want to find a point 1' which maps 1 close to 3 in latent space while keeping the translation from 1 to 1' as sparse as possible. Note that the shown clusters (1–4) are determined in latent space and back-engineered to original space.  $r(\mathbf{X}_i)$  is also referred to in [1] by  $\mathbf{R}_i$ .

## 2 Scope of reproducibility

As follows from the introduction, the authors of [1] opted for a counterfactual, sparse explanation for key differences between (naturally arising) groups. Plumb *et al.* propose an algorithm that generalizes this idea and attempts to find *global* differences between all groups, constructed through a composition of simpler explanations and introduced as *Transitive Global Translations* (TGT). The main tested contributions are that “TGTs provide a *global* and *counterfactual* explanation between all groups which is mathematically consistent (i.e. symmetrical and transitive)” and that “TGTs overcome the shortcomings of statistical and manual interpretation which do not use the model that learned the low dimensional representation that was used to define the groups in the first place.”

In addition to replicating these claims using the provided code, we further evaluated these statements by testing compatibility with other DR methods than the originally used variational auto-encoder and tested the applicability of the algorithm on different datasets with different internal structures. All the latter mentioned experiments were performed by rewriting and implementing the algorithm in PyTorch (taking the authors’ code and paper as references), whereas the original code is written in TensorFlow, motivated by the development of PyTorch in becoming the preferred framework by researchers [5] and the assertion that it offers clearer coding structure<sup>3</sup>. For the sake of ensuring future reproducibility of the experiments originally presented in [1], as an addition, we have also provided an upgraded version of the original code, without further modifications, to make it compatible with recent versions of TensorFlow.

<sup>3</sup>Kirill Dubovikov, *PyTorch vs TensorFlow – spotting the difference*, <https://towardsdatascience.com/pytorch-vs-tensorflow-spotting-the-difference-25c75777377b>, accessed on January 25th, 2021

### 3 Methodology

In accordance with section 2, for the purpose of reproducing [1], several steps were taken in order to best cover the scope of the original paper with limited resources. Firstly, the provided artifacts were ‘modernized.’ Concretely – the code provided by Plumb *et al.* and the external library `scvis` were upgraded to be compatible with TensorFlow 2.4.1<sup>4</sup> and matplotlib 3.3.3<sup>5</sup>. Secondly, the provided code was rewritten from scratch – both the code complementing the original paper; and the `scvis` library, to make the computation of the explanations independent from the DR algorithms. Thirdly, we have run our implementation on different DR algorithms—both linear and non-linear (see section 3.1)—and on different datasets (see section 3.3). Both the reproduced version (using the original code) and the rewritten version are further explained in section 3.2. All relevant code complementing [1]—the chunks which we have chosen to rewrite—are explicitly stated in the paper, making the method by Plumb *et al.* reproducible, even if the code would not have been provided by the authors.

#### 3.1 Dimensionality reduction algorithms

Throughout [1], DR algorithms are only mentioned in the general sense. Only in section 4 the algorithm which Plumb *et al.* use for their experiments is introduced: a variational auto-encoder (VAE), which is a non-linear family of dimensionality reduction algorithms based on (deep) symmetrical neural networks with a *bottleneck* in the middle layers which form the latent representation of the input data. The implemented VAE is based on the architecture proposed by [6].

The presented explanatory method *should* work for any DR algorithm while effectively treating the algorithm as a ‘black box.’ Apart from testing this hypothesis, it is relevant to compare ‘explanatory performance’ on different DR algorithms as different algorithms suit different types of data. Furthermore, an explanation based solely on translations could possibly perform worse in situations wherein data is transformed in a non-linear manner, especially when methods inherently non-linearly transform (and warp) the input space. Therefore, as an addition to simply reproducing the implementation, experiments were done with several commonly known DR algorithms listed below.

**Linear methods** – We have opted for two different linear DR algorithms: truncated SVD (TSVD) and sparse PCA (SPCA) [7]. Both TSVD and SPCA reduce the dimensionality in a linear fashion. However, a key difference between both algorithms is that SPCA centers the data before computing the decomposition, where TSVD does not. As a result, TSVD handles sparse data more efficiently. The objective of SPCA is to find the sparse components that best reconstruct the data. While standard PCA, in most cases, extracts components using dense expressions, these are often hard to interpret. However, the sparse vectors extracted by sparse PCA naturally match the latent components, which increases explainability.

**Non-linear methods** – The implemented VAE is based on a Gaussian distribution, which results in a probabilistic generative model that preserves both local and global neighbor structures in the data [6]. As the VAE is solely used for encoding a constant latent space, the model is trained on the entire dataset. We further incorporate the use of three additional non-linear methods: kernel PCA (KPCA), locally linear embedding (LLE) [8] and isomaps. The latter two mentioned techniques are examples of *manifold learning* [9].

<sup>4</sup>TensorFlow, *Automatically upgrade code to TensorFlow 2*, <https://www.tensorflow.org/guide/upgrade>, accessed on January 22nd, 2021

<sup>5</sup>Matplotlib, *API Changes*, [https://matplotlib.org/3.3.3/api/api\\_changes.html](https://matplotlib.org/3.3.3/api/api_changes.html), accessed on January 22nd, 2021

To achieve a non-linear transformation, KPCA [10] extends standard PCA through the use of kernels. These kernels effectively mimic a complex function that projects the input data on a *higher* dimensional space in which, consequently, a lower-dimensional subspace is found in which the data is represented, resulting in a more efficient low-dimensional latent representation. The advantage is that KPCA is able to identify clusters which are not linear separable. For our experiments, we found that a sigmoid kernel provided the best performance (measured by the metrics proposed in section 3.2) when the latent space resulting from the KPCA was used to find TGTs.

Finally, manifold learning techniques are implemented to analyse whether TGTs can still perform on a higher-dimensional embedding of a low-dimensional manifold, especially for datasets of which its data might not lay on an underlying low-dimensional manifold. Isomaps can be seen as an extension of KPCA. Isomaps present a lower-dimensional space while maintaining distances between all points, while LLE aims to map to a lower-dimensional space while maintaining the distance in local neighborhoods. While Isomap could be seen as an extension of KPCA, LLE can be understood as a combination of PCAs run on local neighborhoods.

As shown in section 4.1, different latent spaces (resulting from applying different processes for the purpose of reducing dimensions) seem to perform better in terms of explainability, which is further discussed in section 5.

### 3.2 Model descriptions

The original paper aimed to find an explanation for the key differences between a pair of groups in latent space where the explanation is expressed in terms of the original dimensions. This explanation is represented by a counterfactual translation of one of the groups (in the pair) in original space: “*what if* all points in group A, thus  $\forall x \in A, x \in \mathbb{R}^d$ , would have been translated, so that  $\forall x \in A, \delta \in \mathbb{R}^d, x' = x + \delta$ ? Would group A have been roughly the same as group B *after* the groups have been mapped to latent space, so that  $\forall x \in A, \forall y \in B, x \in \mathbb{R}^d, y \in \mathbb{R}^m, r : \mathbb{R}^d \rightarrow \mathbb{R}^m, r(x') \approx r(y)$ ?”.

As both the original space ( $\mathbb{R}^d$ ) and latent space ( $\mathbb{R}^m$ ) are Euclidean spaces, translations between any pair of groups within a larger number of groups ( $> 2$ ) can be constructed by composing two translations (a vector addition in this context) or negating a translation (in this context equivalent to vector-scalar multiplication with  $\lambda = -1$ ). Using these operations and a *reference group*, explanations ( $\delta$ ) can be obtained for every possible pair of groups. The intuition behind this idea is illustrated in figure 2. Points in figure 2 do not directly correspond to the groups for which we want to find differences, since these locations will be approximated using sparse translations (see figure 1). Moreover, this linearity enables the possibility to measure the difference between the points in the two groups in latent space using the  $l_2$ -norm of the squared differences. At the same time, it is also possible to measure the *sparsity* of  $\delta$  using  $l_1$ -regularization, which directly relates to the comprehensibility of the explanation posed by  $\delta$ .

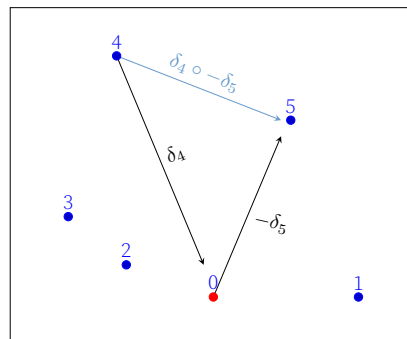
In order to obtain adequate components for the  $\delta$ -vector between a group and the reference group, all  $\delta$ -vectors can be initialized to zero and optimized by adjusting the components of  $\delta$  corresponding to two randomly chosen groups in the negative direction of the gradient of a loss function which accounts for the above-mentioned constraints. This loss function is formulated in equation 9 of the original paper. While Plumb *et al.* decide to incorporate the calculation of the gradient as an ‘extension’ of the `scvis` package, which we have reproduced using the provided code, we have decided to analytically compute the gradient using `SciPy 1.6.0` (using the `optimize.approx_fprime` method). This implementation choice is in line with our aim to rewrite all code from scratch to obtain a ‘stand-alone’ division between the explanatory algorithm (which is

the main idea of the original paper) and various ‘black-box’ dimensionality reduction algorithms, including the variational auto-encoder (the only algorithm Plumb *et al.* experiment with), as introduced in the beginning of this section.

The model’s performance is measured by two separate, but related (see section 4.1), metrics defined by [1, p. 3]: *correctness*, which measures the degree to which projected points are actually *in the vicinity* of points they should map to (which they, in a sense, should ‘imitate’) and *coverage*, which measures the degree to which projected points *cover* the target group. Both are defined in equation 3–4.

### 3.3 Datasets

In the upgraded-code model, we decided to reproduce the explanations on all provided datasets, including the synthetic and corrupted versions. The datasets used in the experiments in the original paper are the Heart Disease, Boston Housing, and Iris datasets; and the single-cell RNA dataset [11]. Our from-scratch model has incorporated all provided datasets, except for the latter, due to lack of computation power. However, experiments were run on three additional UCI datasets: Seeds, Wine and Glass.<sup>6</sup> The additional data allow us to understand how varying underlying structures in data influence explainability using different DR algorithms. Not all datasets can be reduced by all techniques: all three added datasets do not yield eigencomponents for a sigmoidal Kernel-PCA procedure.



**Figure 2.** Visualisation of the intuition behind the composition of two translations with regard to a reference group 0, of which one is negated. Note that in this example every data point represents a whole group.

### 3.4 Hyperparameters

A constraint on the resulting explanations is that they must be sparse. The sparsity level is formally defined by  $k$ , representing the number of dimensions in the original space (with ‘real-life’ features) used as the main components of the translation that forms the explanation. The hyperparameter  $k$  is enforced *after* the learning process by truncating the  $\delta$ -vector. To enforce sparsity *during* the learning process, the  $\delta$ -vector is being  $l_1$ -regularized by a term  $\lambda$ , as follows from [1, p. 5] and was inspired by the field of *compressed sensing*. The resulting equation, which is minimized, is stated in equation 1. As further explained in section 3.5, different values for  $\lambda$  are tested for optimization.

$$\text{loss}(\delta) = \|r(\bar{x}_{\text{initial}} + \delta) - \bar{r}_{\text{target}}\|_2^2 + \lambda \|\delta\|_1 \quad (1)$$

Following [1, p. 6], in equation 2 a *similarity measure* is defined to capture the degree of which the features of a  $k_1$ -sparse explanation correspond to  $k_2$ -sparse explanation where  $k_2 > k_1$ .

$$\text{similarity}(e_1, e_2) = \frac{\sum |e_1[i]| \mathbb{1}[e_2[i] \neq 0]}{\|e_1\|_1} \quad (2)$$

Lastly, hyperparameter  $\epsilon$  defines a threshold for the *correctness* and *coverage* metrics, as defined in equation 3–4. Although Plumb *et al.* do hint on the type of search they have performed to find values for  $\epsilon$ , it is not clearly expanded upon [1, p. 4]. For the purpose of analysing reproducibility, we have adopted the same values for  $\epsilon$  (which were not present

<sup>6</sup>UCI Machine Learning Repository, <https://archive.ics.uci.edu/ml/datasets/{heart+disease,iris,seeds,Wine,glass+identification}> and <https://archive.ics.uci.edu/ml/machine-learning-databases/housing/>, accessed on January 29th, 2021; note the set notation in the latter section of the URL.

in the paper) for the same datasets. For all new data, we have introduced similar static values for  $\epsilon$ . Plumb *et al.* do, however, propose an evaluation method for  $\epsilon$ , although it is again not mentioned in the paper. The authors take the mean, minimum and maximum of the diagonal values of the matrix of the correctness measures for all clusters (similar to the matrices shown in figure 4). These values can be misleading, as larger values for  $\epsilon$  would inherently score high on correctness. Therefore, we have constrained  $\epsilon \in [0, 2]$  and dynamically scale the latent spaces to force the data to be spread out (see section 3.5). We evaluated all values for  $\epsilon$  and outline the means in table 1.

$$\text{correctness}(t) = \frac{1}{|\mathbf{X}_{\text{initial}}|} \sum_{x \in \mathbf{X}_{\text{target}}} \mathbb{1}[\exists x' \in \mathbf{X}_{\text{target}} \|r(t(x)) - r(x')\|_2^2 \leq \epsilon] \quad (3)$$

$$\text{coverage}(t) = \frac{1}{|\mathbf{X}_{\text{target}}|} \sum_{x \in \mathbf{X}_{\text{initial}}} \mathbb{1}[\exists x' \in \mathbf{X}_{\text{initial}} \|r(x) - r(t(x'))\|_2^2 \leq \epsilon] \quad (4)$$

	VAE	PCA	TSVD	KPCA	SPCA	ISO	LLE	$ \text{clusters} $	$\epsilon$
Housing	0.99	1.00	1.00	1.00	1.00	1.00	0.83	6	1.50
Iris	0.98	0.96	0.98	0.94	0.96	0.85	0.99	3	0.75
Heart	0.96	0.99	1.00	0.99	0.99	0.98	1.00	8	1.00
Seeds	1.00	0.97	1.00	-	0.97	0.98	0.94	3	1.00
Wine	1.00	0.99	1.00	-	0.99	0.99	1.00	3	1.00
Glass	0.99	0.94	0.94	-	0.94	0.94	0.98	7	1.75

**Table 1.** Evaluation of fixed epsilon values present in equations 3–4: mean of the diagonal of the correctness matrix.

### 3.5 Experimental setup and computational requirements

We have reproduced the experiments in [1] by using the upgraded TensorFlow-code. Additionally, we have run our from-scratch code on three additional datasets and six other dimensionality reduction algorithms. All experiments were run for five trials for eleven different values of  $\lambda$  (evenly spaced between 0 and 5), for which the best set of  $\delta$ -vectors is chosen for all values of  $k$  ( $k$  is evenly spaced between 1 and the number of dimensions  $d$  in the original space, with a step size of 1 for  $d \leq 5$ , and 2 otherwise).

As pointed out in the previous section, an inadequately high value for  $\epsilon$  poses a problem if the data in latent space is of low variance (especially if variance  $< \epsilon$ ), as the performance measures, introduced in section 3.2 and shown in equations 3–4, would unjustly report very high scores. We have solved this problem by rescaling the data in latent space so that a variance of 10 is preserved (which amounts to a standard deviation of  $\approx 3.16$ ). We have deliberately chosen not to utilize a method for removing *outliers* prior to defining the factor with which to rescale the data in latent space, as to preserve ‘the spirit of the data,’ especially since the internal structure of some datasets contain outliers which potentially correspond to groups with a significantly different character, as opposed to outliers resulting from noisy measurements. We further discuss our choice in section 5. However, as we failed to preserve an adequate amount of variance when performing *TSVD* on the *Glass*-dataset using this method, we manually scaled the data in this latent space, for this particular dataset, with an additional factor of 20.

In the from-scratch implementation, K-means is used for clustering, while for the replicated experiments in the original notebooks, following the code provided by the authors, clusters were manually selected.

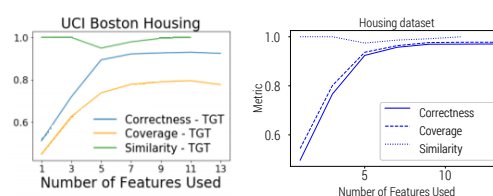
Experimental meta-results are shown in table 2. The original code, the ‘modernized’ version of the original code and the code for the from-scratch implementation can be found on our GitHub repository [12].

	VAE + model	PCA	TSVD	KPCA	SPCA	ISO	LLE	$ k $
Housing $\in \mathbb{R}^{13}$	3.7 + 0.1	0.3	0.3	4.4	0.4	4.5	4.7	7
Iris $\in \mathbb{R}^4$	0.8 + 0.1	0.1	0.1	0.7	0.1	1.0	1.1	4
Heart $\in \mathbb{R}^{14}$	2.6 + 0.1	0.6	0.5	4.9	0.5	4.3	3.6	7
Seeds $\in \mathbb{R}^7$	0.4 + 0.1	0.2	0.1	-	0.2	3.0	2.7	4
Wine $\in \mathbb{R}^{13}$	1.6 + 0.1	0.2	0.2	-	0.2	4.2	4.8	7
Glass $\in \mathbb{R}^{10}$	0.9 + 0.1	0.4	0.3	-	0.3	2.8	2.8	5

**Table 2.** Rounded TGT training time in hours per dataset and DR algorithm for all values of  $k$ , measured on an i7-4720HQ CPU @ $\approx$  2.6GHz. Measured for 5 trials per  $\lambda \in \{0, 0.5, \dots, 5\}$ . VAE models train on a minimum of 3000 iterations.

## 4 Results

Using our upgraded TensorFlow-code and the provided pre-trained VAE models, the obtained results are identical to those provided by the paper’s code-base for all datasets and methods.<sup>7</sup> After re-training all VAE models and re-learning all explanations, using the same upgraded TensorFlow-code, we obtain very similar results.<sup>8</sup> The results are not identical, as the models learn slightly different representations, in which the manually selected clusters also differ. This indicates that the used cluster generation techniques influence the model’s ability to explain based on a translation. An example arises with the Iris dataset: Plumb *et al.* yield 0.833 coverage, while we get only 0.66. This difference caused by a different organization of clusters propagates further into the results for the corrupted data. Except for minor differences arising from this same issue, we successfully reproduced the results for all other datasets.



(a) Plumb *et al.* (b) Implemented from scratch

**Figure 3.** Comparison of VAE on Housing data

When running the Housing, Iris and Heart datasets on the from-scratch implementation of the explanation algorithm using our implementation of `scvis`, we obtain either similar (see figure 3) or even better results than Plumb *et al.* do. This again indicates that the method for dimensionality reduction does impact the results.

For different datasets and different dimensionality reduction algorithms that map the datasets to different latent spaces, we have encountered the same problem related to the inability of explaining pairs of groups of which the clusters have a different standard deviation, as shown in the original paper in figures 4-7 [1, p. 4]. This problem occurs, among others, in figure 5f in appendix A (and the corresponding measures in figure 6f).

For different datasets and different dimensionality reduction algorithms that map the datasets to different latent spaces, we have encountered the same problem related to the inability of explaining pairs of groups of which the clusters have a different standard deviation, as shown in the original paper in figures 4-7 [1, p. 4]. This problem occurs, among others, in figure 5f in appendix A (and the corresponding measures in figure 6f).

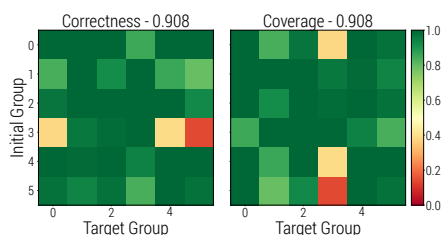
<sup>7</sup>To compare, open the `.ipynb` in all folders (except ‘code,’ ‘scvis,’ ‘MiscFigures,’ and ‘Integrated-Gradients-master’) on both repositories, and look at the resulting plots: <https://github.com/GDPlumb/ELDR> and [https://github.com/damiaa/r/fact-ai/tree/main/ELDR-TF2.x\\_\(pre\\_trained\\_models\)](https://github.com/damiaa/r/fact-ai/tree/main/ELDR-TF2.x_(pre_trained_models))

<sup>8</sup>Now, compare the plots with [https://github.com/damiaa/r/fact-ai/tree/main/ELDR-TF2.x\\_\(newly\\_trained\\_models\)](https://github.com/damiaa/r/fact-ai/tree/main/ELDR-TF2.x_(newly_trained_models))



#### 4.1 Additional results not present in the original paper

As [1, p. 6] pointed out, the *correctness* and *coverage* metrics are closely related. If the translations between arbitrary groups are symmetrical, the former translation—which is the technical representation of the explanation—is the negative of the latter (and vice-versa). Thus, when applying different **linear** dimensionality reduction algorithms, both metrics will exactly equal each other, as follows from the results shown below in table 3, and in figure 4 (note that the color map plot for correctness is the ‘transpose’ of the plot for coverage, and vice-versa) and figure 6 (along all different values for  $k$ , the graphs conveying correctness and coverage are on top of each other for all linear algorithms).



**Figure 4.** Metrics for Housing on PCA,  $k = 5$

An additional observation is that topology-based dimensionality reduction techniques are less suitable for data compression when the compressed data consequently needs to be explained using translation-based counterfactual explanations. Locally linear embedding and isomapping are methods based on manifold mapping: a one- or two-manifold is embedded in two-dimensional space, and points of a higher-dimensional space are mapped onto this manifold. This often yields problem-

atic data structures. For instance, when data is represented using a line (a one-manifold) and embedded in two-dimensional space in which the one-dimensional line is twirled. Since we are only considering translations (and, *e.g.*, no rotations), the method is unable to adequately explain differences between clusters present on manifolds in latent space.

Furthermore, we observed that certain datasets yield constant results when varying the sparsity-constraint. In this case, the compression from high-dimensional to two-dimensional space was (nearly) *lossless* – a single component suffices to explain the difference between any group. This phenomenon is shown in figure 6e in appendix A, indicating that the chemical composition of wine actually depends on only one latent variable: when using PCA, 99.98% of the variance is explained by only two components (of which 99.81% by the first component) out of 13 input dimensions. Note that the VAE seems to learn an approximation of LLE space (see figure 5e).

All results for all introduced performance measures, datasets and dimensionality reduction algorithms, are shown in table 3 below and illustrated in figures 5–7 in appendix A.

	Housing	Iris	Heart	Seeds	Wine	Glass
VAE	0.97; 0.98; 0.99	0.96; 0.94; 1.00	0.83; 0.82; 0.98	0.40; 0.40; 0.99	0.39; 0.39; 0.33	0.33; 0.33; 0.86
PCA	0.99; 0.99; 0.99	0.89; 0.89; 1.00	1.00; 1.00; 0.98	0.95; 0.95; 0.99	0.40; 0.40; 0.33	0.82; 0.82; 0.86
TSVD	1.00; 1.00; 0.99	0.95; 0.95; 1.00	1.00; 1.00; 0.98	0.84; 0.84; 0.99	0.52; 0.52; 0.33	0.94; 0.94; 0.86
KPCA	0.99; 0.98; 0.99	0.82; 0.86; 1.00	0.99; 0.99; 0.98	-	-	-
SPCA	0.99; 0.99; 0.99	0.88; 0.88; 1.00	1.00; 1.00; 0.98	0.95; 0.95; 1.00	0.39; 0.39; 0.33	0.82; 0.82; 0.86
ISO	0.22; 0.22; 0.99	0.67; 0.66; 1.00	0.16; 0.17; 0.98	0.93; 0.94; 1.00	0.39; 0.39; 0.33	0.58; 0.61; 1.00
LLE	0.20; 0.25; 1.00	0.54; 0.55; 1.00	0.40; 0.51; 1.00	0.88; 0.87; 1.00	0.37; 0.37; 0.33	0.33; 0.45; 1.00

**Table 3.** Correctness, coverage and similarity scores respectively for each model and dataset. Mean value is shown for the similarity scores while the score for the largest  $k$  is taken for the correctness and coverage measures.

## 5 Discussion

We have been able to successfully reproduce and upgrade the code provided by Plumb *et al.*, which is the implementation of the technique presented in [1]. We were able to reproduce results by running the pre-trained models; after re-training these models; and by rewriting the algorithm from scratch.

The performance depends on the mapping to latent space. A limitation of the algorithm is the lack of variable freedom: only translations can be used to explain differences between groups. By utilizing different DR methods, we have shown that not all algorithms produce latent spaces in which translations suffice for an explanation between clusters. Especially manifold-based algorithms require a more sophisticated type of explanation using, for instance, rotation and scaling.

A possible solution would be to generate an explanation in terms of a matrix (so that  $\mathbf{x}' = \mathbf{M}\mathbf{x} + \boldsymbol{\delta}$  instead of  $\mathbf{x}' = \mathbf{x} + \boldsymbol{\delta}$ ). However, explanations using translations are directly interpretable for humans as the dimensions of the tested datasets correspond to explainable features. Considering, for example, rotations, would come at the cost of explainability. However, forcing  $\mathbf{M}$  to be diagonal (which leads to multiplying all data features by the corresponding diagonal factors) might yield proper explanations (of the type “*If your monthly income would have been twice as big...*”).

Furthermore, by using different datasets, we have shown that some DR techniques do a better job in terms of explainability between clusters. Data might be structured in different ways and produce different types of clusters. It directly follows from our results that the explanation method heavily relies on the generated clusters, and more importantly, its shapes and variance in latent space, which is another limitation.

In section 3.5, we opted for an approach in which the latent spaces are dynamically scaled for the purpose of achieving a certain amount of variance in the data, as to obey the fixed hyperparameter  $\epsilon$  which was introduced in section 3.4. A more robust method would be to dynamically define  $\epsilon$  based on the variance in the latent space depending on the dataset and, potentially, proper handling of its outliers. We view the lack of expansion on this hyperparameter in the original paper as a (minor) limitation.

As VAEs are sampling points in latent space from a distribution, an identical point  $\mathbf{x}$  in the input space which is repeatedly mapped to latent space will yield different mappings. It would be interesting to further investigate to what extent explanatory algorithms can ‘handle’ this variance – to what extent such algorithms could find a stable and unchanging explanation in an ever-changing counterfactual world.

### 5.1 What was easy

After having asserted the dependencies needed to run the code provided by the authors on their GitHub repository, replicating the experiments performed in the paper was easy; the code was cleanly written, and it was easy to understand the architectural choices which the authors’ made in their model.

### 5.2 What was difficult

Although the code required for computing the explanations was separated from the implementations in other folders in which it is applied on the different datasets, Plumb *et al.* decided to perform experiments using a VAE and intertwined all code for generating the explanations with the `scvis`-library. This caused the explanatory model to rely

on TensorFlow, while the explanation method itself does not require any deep learning. It is preferred to dissect the model into independent components, considering that the explanatory model should work with other DR algorithms, including those which do not require neural pipelines. Because of the authors' choice to use a VAE, the entire repository relied on TensorFlow. Mainly for this reason, we rewrote everything from scratch, as to provide the code for the dimensionality reduction methods on a 'stand-alone' basis. As we were aiming to reproduce the results provided by the original paper, we also had to rewrite the external `scvis`-library, which provided the VAE architecture, in PyTorch (for the sake of keeping clean, separated and object-oriented code). This was very time-consuming as both frameworks are fundamentally different (dynamic vs static graph definition). The process of rewriting took approximately two weeks.

Lastly, the results for the Bipolar dataset could only be replicated with the provided model configuration file. Although we have re-trained the model using the upgraded originally provided code, we have chosen to exclude this dataset for the additional experiments due to lack of computation power and time.

### 5.3 Communication with original authors

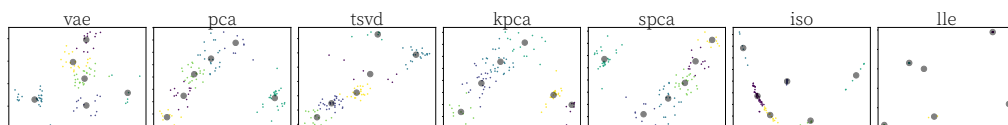
A short interaction occurred to gain more insight into the required external libraries as these were not listed in any documentation. Although the authors responded quickly, the issue had already been solved in the meantime.

## References

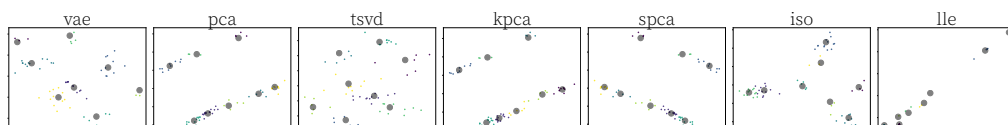
1. G. Plumb, J. Terhorst, S. Sankararaman, and A. Talwalkar. "Explaining Groups of Points in Low-Dimensional Representations." In: **Proceedings of the 37th International Conference on Machine Learning**. Ed. by H. D. III and A. Singh. Vol. 119. Proceedings of Machine Learning Research. PMLR, 13–18 Jul 2020, pp. 7762–7771. URL: <http://proceedings.mlr.press/v119/plumb20a.html>.
2. D. E. O'Leary. "Artificial Intelligence and Big Data." In: **IEEE Intelligent Systems** 28.2 (2013), pp. 96–99. doi: 10.1109/MIS.2013.39.
3. L. Van Der Maaten, E. Postma, and J. Van den Herik. "Dimensionality reduction: a comparative review." In: **J Mach Learn Res** 10 (2009), pp. 66–71.
4. R. A. Mancisidor, M. Kampffmeyer, K. Aas, and R. Jenssen. "Learning latent representations of bank customers with the Variational Autoencoder." In: **Expert Systems with Applications** 164 (2021), p. 114020. doi: <https://doi.org/10.1016/j.eswa.2020.114020>. URL: <https://www.sciencedirect.com/science/article/pii/S0957417420307910>.
5. H. He. "The State of Machine Learning Frameworks in 2019." In: **The Gradient** (2019).
6. J. Ding, A. Condon, and S. P. Shah. "Interpretable dimensionality reduction of single cell transcriptome data with deep generative models." en. In: **Nature Communications** 9.1 (May 2018), p. 2002. doi: 10.1038/s41467-018-04368-5. URL: <https://www.nature.com/articles/s41467-018-04368-5> (visited on 04/13/2021).
7. H. Zou, T. Hastie, and R. Tibshirani. "Sparse Principal Component Analysis." In: **Journal of Computational and Graphical Statistics** 15.2 (2006), pp. 265–286. doi: 10.1198/106186006X113430. eprint: <https://doi.org/10.1198/106186006X113430>. URL: <https://doi.org/10.1198/106186006X113430>.
8. S. T. Roweis and L. K. Saul. "Nonlinear Dimensionality Reduction by Locally Linear Embedding." In: **Science** 290.5500 (2000), pp. 2323–2326. doi: 10.1126/science.290.5500.2323. eprint: <https://science.sciencemag.org/content/290/5500/2323.full.pdf>. URL: <https://science.sciencemag.org/content/290/5500/2323>.
9. L. Cayton. "Algorithms for manifold learning." In: **Univ. of California at San Diego Tech. Rep** 12.1-17 (2005), p. 1.
10. B. Schölkopf, A. Smola, and K.-R. Müller. "Kernel principal component analysis." In: **Artificial Neural Networks – ICANN'97**. Ed. by W. Gerstner, A. Germond, M. Hasler, and J.-D. Nicoud. Berlin, Heidelberg: Springer Berlin Heidelberg, 1997, pp. 583–588.
11. K. Shekhar et al. "Comprehensive Classification of Retinal Bipolar Neurons by Single-Cell Transcriptomics." eng. In: **Cell** 166.5 (Aug. 2016), 1308–1323.e30. doi: 10.1016/j.cell.2016.07.054.
12. D. J. W. Reijnaers, D. B. van de Pavert, and G. Scheuer. **damiaanr/fact-ai: Reproduction of (Plumb et al., 2020)**. Version 1.0.0. Apr. 2021. doi: 10.5281/zenodo.4686025. URL: <https://doi.org/10.5281/zenodo.4686025>.

# Appendices

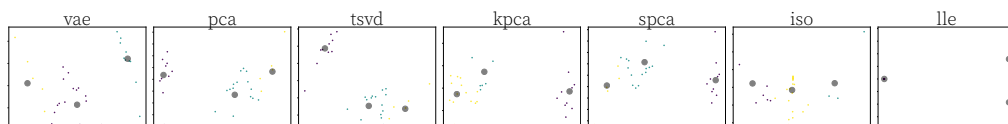
## A Additional latent spaces and corresponding measures



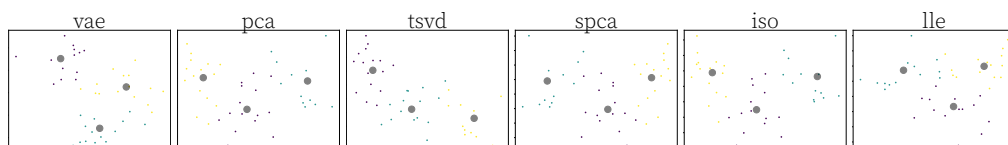
(a) Latent spaces for the Housing dataset



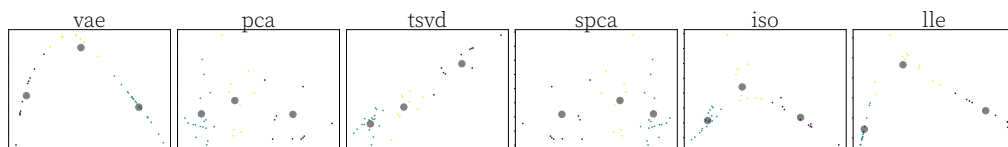
(b) Latent spaces for the Heart dataset



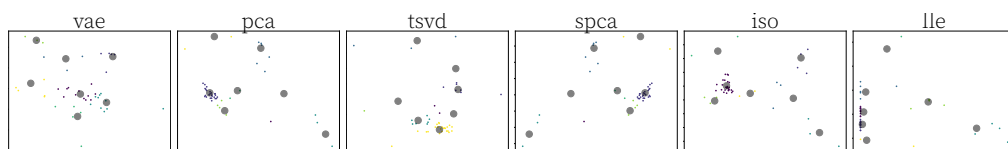
(c) Latent spaces for the Iris dataset



(d) Latent spaces for the Seeds dataset

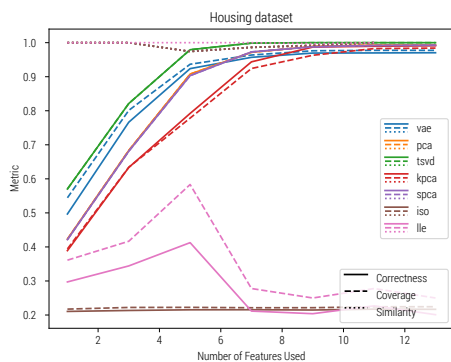


(e) Latent spaces for the Wine dataset

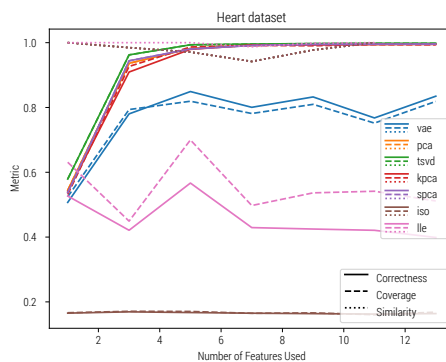


(f) Latent spaces for the Glass dataset

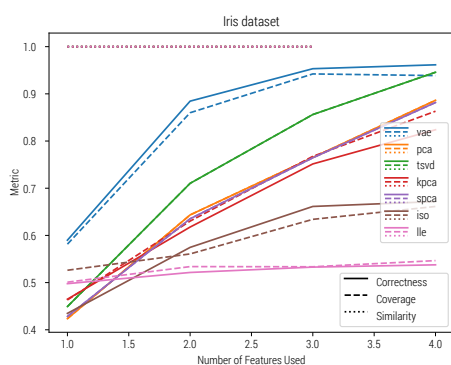
**Figure 5.** Visualization of two-dimensional latent spaces in which higher dimensional data is compressed after applying various dimensionality reduction techniques for different datasets. This figure illustrates how the underlying structure of the data, in combination with the involved method for reducing dimensionality, influences the resulting representational spaces. Every circle represents a cluster to which the (cluster-colored) data points belong.



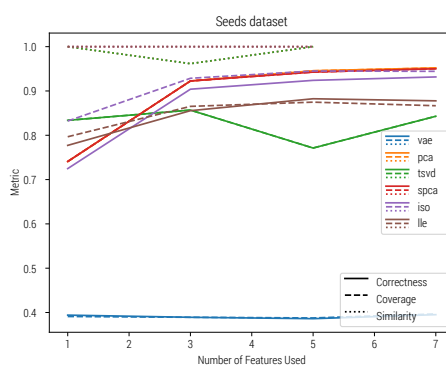
(a) Measures for Housing dataset (all algorithms)



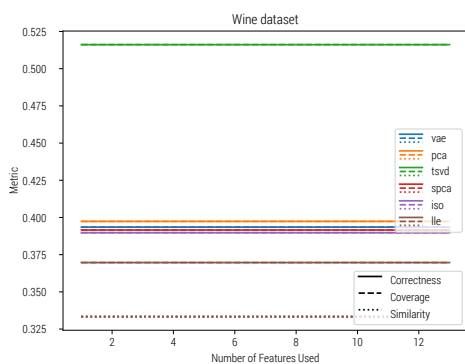
(b) Measures for Heart dataset (all algorithms)



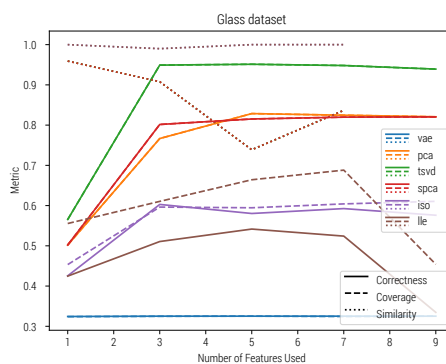
(c) Measures for Iris dataset (all algorithms)



(d) Measures for Seeds dataset (note: no KPCA)

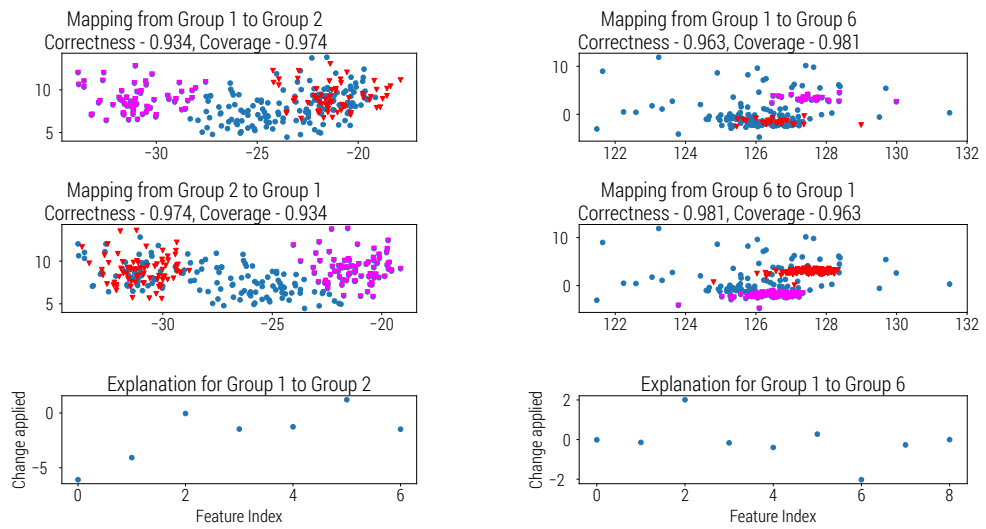


(e) Measures for Wine dataset (note: no KPCA)



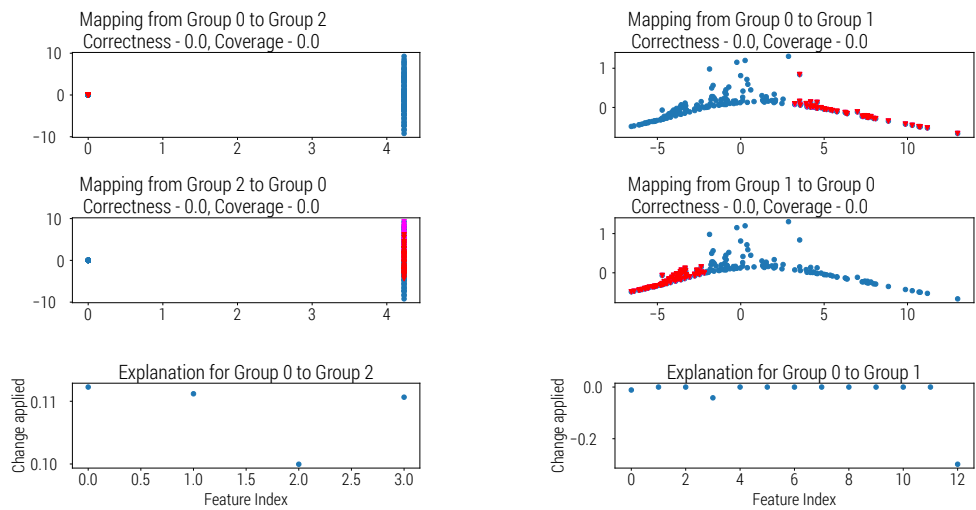
(f) Measures for Glass dataset (note: no KPCA)

**Figure 6.** Measures for coverage, correctness and similarity for all involved dimensionality reduction algorithms on all datasets.



(a) A proper explanation for clusters in the Seeds dataset mapped to a latent space generated by SPCA.

(b) A well-formed explanation for groups within the Glass dataset transformed by TSVD.



(c) Iris dataset mapped with LLE. Note how the model is unable to learn, caused by the shapes of the clusters.

(d) Wine dataset mapped with an isomap. The model is unable to learn, indicating that a translation cannot suffice as an explanation here. Note that  $k$ -sparsity does not influence the explanation, as the Wine dataset has only one significant latent variable.

Figure 7. Four sample explanations between two arbitrary groups for four different datasets projected onto four different latent spaces.