
Explaining Low Dimensional Representation, a reproduction

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 **Scope of Reproducibility**

3 This report covers our reproduction of the paper 'Explaining Low dimensional Representation' [8] by Plumb et al.
4 In this paper, a method (Transitive Global Translations, TGT) is proposed for explaining different clusters in low
5 dimensional representations of high dimensional data. They show their method outperforms the Difference Between the
6 Means (DBM) method, is consistent in explaining differences with few features and matches real patterns in data. We
7 verify these claims by reproducing their experiments and testing their method on new data. We also investigate the use
8 of more complex transformations to explain differences between clusters.

9 **Methodology**

10 We reproduce the original experiments using their source code. We also replicate their findings by re-implementing the
11 authors' method in PyTorch [7] and evaluating on two of the dataset used in the paper and two new ones. Furthermore,
12 we compare TGT with our own extension of TGT, which uses a larger class of transformations.

13 **Results**

14 We were able to reproduce their results using their code, yielding mostly similar results. TGT generally outperforms
15 DBM, especially when explanations use few features. TGT is consistent in terms of the features to which it attributes
16 cluster differences, across different sparsity levels. TGT matches real patterns in data. When extending the types of
17 functions used for explanations, performance did not improve significantly, suggesting translations make for adequate
18 explanations. However, the scaling extension shows promising performance on the modified synthetic data to recover
19 the original signal.

20 **What was easy**

21 The easiest part was running the existing code with the pre-trained model files. The original authors had set up their
22 code base in an organized manner with clear instructions.

23 **What was difficult**

24 The first difficulty that we encounter was finding the right environment. The source code depends on deprecated
25 functionality. The clustering method they used, had to be re-implemented for us to use it in our replication. Another
26 difficulty was the selection of clusters. The authors did not prove a consistent method for selecting clusters in a latent
27 space representation. When retraining the provided models, we get a latent space representation different to the original
28 experiments. The clusters have to be manually selected. The metrics that they used to evaluate their explanations are
29 also depend on the clustering. This means that there is some variability in the exact verification of reproducibility.

30 **Communication with original authors**

31 We asked the original authors for clarification on how to choose the ϵ hyper-parameter. However, it became apparent
32 that we had misread, and the procedure is indeed adequately reported in the paper.

33 1 Introduction

34 The *curse of dimensionality* [3] is a long-standing problem in Machine Learning. Data in many domains and applications
 35 (e.g. Bioinformatics) has high-dimensional representations. Finding patterns in such high-dimensional data is a
 36 challenging task. To this end, *dimensionality reduction* [12] techniques have greatly helped in data-analysis, information
 37 extraction, building computational models, and in doing inference. Given an input $x \in \mathbb{R}^d$, dimensionality reduction
 38 learns a function $r : x \mapsto r(x)$, $r(x) \in \mathbb{R}^m$, where $m \ll d$. Such a dimensionality reduction function r naturally
 39 arises in deep learning due to the expressivity and representational power of neural networks. The goal of r is to encode
 40 useful knowledge about the input space, thus providing distinctive information in the transformed output $r(x)$. This
 41 results in “clusters” or “groups of points” in the transformation space. The downside of this exercise, however, is that
 42 the output space is usually non-interpretable. There is usually no easy way to know what information is present in the
 43 transformed points $r(x)$ and what sort of distinctive knowledge they contain.

44 In this work, we reproduce the paper ‘Explaining Groups of Points in Low-Dimensional Representations’ by Plumb
 45 et al. [8]. This paper proposes a method for explaining different clusters in latent space representation. They look at
 46 the problem of explaining the points in the latent space representation through the lens of Interpretability in Machine
 47 Learning. We reproduce their findings and expand upon their work with our extension. We extend their research by
 48 applying their method to a larger class of explanation functions and testing their method on new dataset. We further
 49 investigate the efficacy of the explanations using a probing classifier [2].

50 2 Methodology

51 Counterfactual Explanations [11] have emerged as an active research area in the field of Interpretable Machine Learning.
 52 A counterfactual explanation is defined as the smallest perturbation to the input that would change the output of a
 53 machine learning model. As such, these explanations are promising as they can provide suggestive recourse to the
 54 beneficiary in a machine learning based decision system. As an interpretable machine learning problem, Plumb et al.
 55 [8] aim to find such counterfactual explanations in order to explain the differences between the groups in latent space.
 56 To this end, they employ the function r itself to find what perturbation δ needs to be made to the input $x \in \mathbb{R}^d$ so that
 57 $r(x + \delta)$ belongs to the different target group. The goal is to find the *global* explanations that apply to the whole group
 58 as opposed to the *local* explanations which explain only individual examples [4]. Furthermore, the explanations need
 59 to be *sparse* for them to be interpretable by practitioners. Finally, these explanations should be both *symmetric*
 60 and *transitive*. To obtain these *Global Counterfactual Explanations*(GCE), the authors propose the algorithm called,
 61 Transitive Global Translations (TGT), explained hereafter.

62 Following the previous notation, let $r : \mathbb{R}^d \rightarrow \mathbb{R}^m$ denote our dimensionality reduction function, where d is the
 63 dimensionality of the input space and m is the latent space’s dimensionality. Suppose $X_i, X_j \subset \mathbb{R}^d$ get mapped
 64 to the clusters $R_i, R_j \subset \mathbb{R}^m$ respectively. The goal is to define the transformation $t_{i \rightarrow j} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ on $x \in X_i$ as
 65 $x' = t_{i \rightarrow j}(x)$, so that $r(x') \in R_j$, or equivalently $x' \in X_j$.

66 The proposed algorithm TGT considers the transformations of the form $t_{i \rightarrow j}(x) = x + \delta_{i \rightarrow j}$. To find the optimal
 67 parameters of the transformation function, authors imply a compressed-sensing based objective function as below:

$$l(\delta_{i \rightarrow j}) = \|r(t_{i \rightarrow j}(\bar{X}_i)) - \bar{R}_j\|_2^2 + \lambda \|\delta_{i \rightarrow j}\|_1 \quad (1)$$

68 where $\lambda \|\delta_{i \rightarrow j}\|_1$ is a regularization term to incentivize sparser explanations, and $\bar{X}_i \in \mathbb{R}^d$ and $\bar{R}_j \in \mathbb{R}^m$ denote the
 69 means of the clusters in the input space and latent space respectively. Given clusters $0, 1, \dots, n$, we get a total of
 70 $\frac{1}{2}n(n+1)$ transformations. To further increase sparsity, we can truncate $\delta_{i \rightarrow j}$ to only the k features with the largest
 71 absolute value, for some k . An issue with this is that the translation using the truncated $\delta_{i \rightarrow j}$ might no longer correctly
 72 transform inputs that get mapped to R_i into inputs that get mapped to R_j .

73 Furthermore, the transformations $t_{i \rightarrow j}$ have to adhere to several mathematical properties. Namely, for any clusters
 74 i, j, k these transformations should be : a) Symmetric, i.e. $t_{i \rightarrow j} = t_{j \rightarrow i}^{-1}$ and b) Transitive, i.e. $t_{j \rightarrow k} \circ t_{i \rightarrow j} = t_{i \rightarrow k}$.
 75 From these properties it follows that $t_{i \rightarrow i}$ is the identity function \mathcal{I} as

$$t_{i \rightarrow i} = t_{i \rightarrow 0} \circ t_{0 \rightarrow i} = t_{i \rightarrow 0} \circ t_{i \rightarrow 0}^{-1} = \mathcal{I} \quad (2)$$

76 We define this condition as *self-similarity*. Furthermore, the group of translations is uniquely defined by $t_{0 \rightarrow 1}, \dots, t_{0 \rightarrow n}$,
 77 because for any i, j :

$$t_{i \rightarrow j} = t_{0 \rightarrow j} \circ t_{i \rightarrow 0} = t_{0 \rightarrow j} \circ t_{0 \rightarrow i}^{-1} \quad (3)$$

78 Plumb et al. [8] compare their method against the naive baseline of Difference Between the Means (DBM). With DBM,
 79 each transformation is still a translation: $t_{i \rightarrow j}(x) = x + \delta_{i \rightarrow j}$. However, now $\delta_{i \rightarrow j} = (\bar{X}_j - \bar{X}_i)$. We also use this as a
 80 baseline for comparison in this report.

81 Since translations are a very narrow class of functions, we expanded upon the research by investigating other trans-
 82 formations that still satisfy the GCE requirements. We investigate the transformations of the form $t_{0 \rightarrow i}(x) =$
 83 $\exp(\gamma_{0 \rightarrow i}) \odot x + \delta_{0 \rightarrow i}$. These always have a well defined inverse, given by $t_{0 \rightarrow i}^{-1}(x) = \exp(-\gamma_{0 \rightarrow i}) \odot (x - \delta_{0 \rightarrow i})$
 84 and only have $\mathcal{O}(d)$ parameters. The inclusion of scaling could enhance performance, while the necessary components
 85 of GCE are maintained.

86 2.1 Metrics to evaluate Global Counterfactual Explanations

87 To measure the efficacy of the transformation function $t_{i \rightarrow j}$, the authors propose two metrics, *Coverage* and *Correctness*.

88 1. The **Coverage** ($cv(t_{i \rightarrow j})$) is the fraction of points $a \in R_j$ for which there is a point $b \in X_i$ such that
 89 $\|r(t_{i \rightarrow j}(b)) - a\|_2 < \epsilon$, i.e.

$$cv(t_{i \rightarrow j}) = \frac{1}{|R_j|} \sum_{a \in R_j} \mathbb{I}[\exists b \in X_i \|r(t_{i \rightarrow j}(b)) - a\|_2 < \epsilon] \quad (4)$$

90 2. The **Correctness** ($cr(t_{i \rightarrow j})$) is the fraction of points $b \in X_i$ for which there is some $a \in R_j$ such that
 91 $\|r(t_{i \rightarrow j}(b)) - a\|_2 < \epsilon$, i.e.

$$cr(t_{i \rightarrow j}) = \frac{1}{|X_i|} \sum_{a \in R_j} \mathbb{I}[\exists b \in X_i \|r(t_{i \rightarrow j}(b)) - a\|_2 < \epsilon] \quad (5)$$

92 Note that both these metrics have the hyperparameter ϵ which is to be chosen carefully. When $i = j$ we do not count
 93 the point itself, there must be some other point within distance ϵ .¹

94 Furthermore, the **Similarity** metric measures the consistency of the explanations at different sparsity levels. Given two
 95 explanations e_1, e_2 where e_1 is more sparse than e_2 , the similarity of e_1 and e_2 is defined as

$$sim(e_1, e_2) = \frac{\sum_i |e_1[i]| \mathbb{I}(e_2[i] \neq 0)}{\|e_1\|_1} \quad (6)$$

96 This is equal to 1 if e_1 uses a subset of the features that e_2 uses. By definition, DBM always has similarity 1.

97 3 Scope of Reproducibility

98 We investigate the following claims from the original paper:

- 99 1. In terms of the average correctness and coverage, TGT performs equally well or better than the DBM method.
 100 This remains true, especially for sparser explanations.
- 101 2. TGT explanations have similarity close to 1. It is consistent in which features it uses for explanations across
 102 different sparsities.
- 103 3. TGT correctly identifies known causal structure in data.
- 104 4. Furthermore, TGT explanations are consistent. When altering the dataset by adding a copy of a cluster with a
 105 specific feature altered, TGT recovers the modification with little change to the other explanations.

106 4 Methodology of Reproducibility

107 We make use of the code made available by the original authors² for our pilot investigative study. We first verify that
 108 the provided models and explanations stay true to the claims made in the paper. We further retrain their models on the
 109 provided dataset. We also made our own PyTorch [7] implementation to to further verify the claims, and to perform
 110 experiments with the proposed extension.

¹We use this definition to set the value for epsilon, as explained in the *Methodology* section of the original Paper.

²<https://github.com/GDPlumb/ELDR>

111 4.1 Model description

112 We identify that the scope of the original paper is to explain clusters in the low-dimensional representations. However,
113 obtaining meaningful and discernible low-dimensional representations is an active area of research. The original authors
114 employ a t-SNE [10] objective based Variational Autoencoder (called, henceforth, as scVIS) [5] as the r function. They
115 make use of library³ by the original scVIS authors in their implementation. We also implement this model in Pytorch
116 for our experiments. However, we deliberately decide not to match the model implementation exactly. This is done
117 to study the model-agnosticism of the TGT algorithm. By design, TGT should be able to explain the clusters for any
118 differentiable r function. However, we maintain that r should give discernible latent representations with preserved
119 global structure in the data. In our implementation of the scVIS library, we therefore do not employ the hyperparameters
120 and the training settings from the original library.

121 4.2 Dataset Description

122 We reproduce the findings of the authors on four datasets that they used. We use two of these datasets as well as two
123 new ones to test our PyTorch implementation.

- 124 1. **Single cell RNA** [9]: This dataset has 13166 features. We use the same number of clusters at the original
125 authors, 18 in this case.
- 126 2. **UCI Boston housing** This dataset has 506 entries with 13 features. We use 6 clusters for both reproduction
127 and replication.
- 128 3. **UCI Heart disease** This dataset has 303 entries with 13 features and 1 binary label. We used 8 clusters in the
129 reproduction and 4 in the replication. The data was normalized to be in the range [0, 1].
- 130 4. **UCI Iris** This dataset has 150 entries with 4 features and 1 ternary label. Ran in the reproduction with 3
131 clusters. $N = 150$
- 132 5. **Breast Cancer Wisconsin (Diagnostic)** ⁴ This dataset has 569 entries with 30 features and 1 binary label. We
133 use 3 clusters in the replication.
- 134 6. **Pima Indians Diabetes Database** ⁵ This dataset has 768 entries with 8 features and 1 binary label. We used
135 3 clusters in the replication. The data was normalized to be in the range [0, 1].
136 Note that the number of clusters depend on the latent-space representation, and thus, are user dependent.

137 4.3 Hyperparameters

138 **Tensorflow [1] Experiments** For the reproduction of the original experiments, we use the same hyperparameters as
139 the original authors.

140 **Pytorch** For our implementation of the scVIS model, we use l2 regularization of 0.001, learning rate 0.01, and
141 perplexity of 10. Furthermore, the degree-of-freedom for the studentT distribution is set to 2.0. Perplexity and the
142 degree-of-freedom is used same as the original scVIS implementation. We use validation set to monitor the training
143 process of the scVIS model, and stop training when the ELBO(Evidence Lower BOUND)[6] stops improving. For
144 training the TGT explanations, we closely follow the settings from Plumb et al. [8]. We initialize the δ_{s_s} as
145 zero vectors. We tune the regularization parameter λ by grid search over a fixed range [0.0, 5.0] incremented by 0.5.
146 Defining the metrics for TGT requires careful setting of the ϵ hyper-parameter. We follow the *self-similarity* condition
147 (transformations of clusters to themselves should, theoretically, have correctness and coverage to be 1.0), and increase
148 the ϵ in the range [0.0, 2.0] with increments of 0.02 until the correctness and coverage metrics are greater than 0.95.
149 Furthermore, we use the truncation values(TV)(refer Table 1) to evaluate on the sparsity of the explanations. For the
150 Pima Indians Diabetes Database and Breast Cancer Wisconsin(Diagnostic) dataset, we use the same truncation values
151 as for UCI Boston Housing dataset.

152 4.4 Experimental setup and code

153 We closely follow the experimental setup in the original paper for our experiments. We make our Pytorch code available
154 ⁶ to further support the reproducible research. We reran the code of the original authors with new clustering models

³<https://github.com/shahcompbio/scvis>

⁴<https://www.kaggle.com/uciml/breast-cancer-wisconsin-data>

⁵<https://www.kaggle.com/uciml/pima-indians-diabetes-database>

⁶<https://github.com/elfrink1/FACT>

Dataset	Truncation Values(TV)	ϵ
Single Cell RNA	50, 100, 250, 500, 1000, 15000	0.75
Heart Disease	1, 3, 5, 7, 9, 11, 13	1.0
Housing	1, 3, 5, 7, 9, 11, 13	1.5
Iris	1, 2, 3, 4	0.75

Table 1: Truncation Values (TV) and ϵ value used for each of the dataset.

Explanation	x_1	x_2	x_3	x_4
$0 \rightarrow 1$	-1.01	-0.02	0.00	-0.88
$0 \rightarrow 1$	-1.05	0.99	0.00	-0.88
$0 \rightarrow 3$	0.00	0.89	0.00	0.00

Table 2: Explanations for the synthetic dataset as given by our implementation. Note that both DBM and TGT are able to infer that the x_3 is not causing any cluster. However, the authors’ claim that TGT also discovers that x_4 doesn’t cause any cluster cannot be verified.

155 and new explanations. We optimize the compressed-sensing based objective function for the TGT algorithm using the
 156 gradient descent algorithm. Our scaling extension is easily integrated in the source code, and can be optimized in a
 157 similar way. We train the scVIS models on the Lisa computing cluster ⁷. We use approximately 30 hours of GPU time.
 158 We train the TGT explanations on CPU (Intel i5).

159 5 Results

160 For the reproduction of the authors’ experiments, we achieve approximately similar results to the original paper. The
 161 TGT method does seem to outperform DBM method. The TGT explanations also have high similarity across sparsity
 162 levels. However, the TGT algorithm is unable to identify known causal structure in synthetic data with as good precision
 163 as reported in the original paper. We are also unable to match the results on the modified and corrupted data to a good
 164 precision. We describe the results in the following sections:

165 5.1 Results reproducing original paper

166 5.1.1 Coverage, Correctness and similarity

167 In figure 1, we can see a comparison between the correctness, coverage and similarity of the TGT and DBM methods.
 168 Note that the DBM always has similarity 1. The similarity of TGT stays between 1 and 0.9, which supports claim 2.

169 We see that the coverage and correctness are similar for the UCI Heart disease dataset. On the UCI Iris dataset, the
 170 coverage is comparable but the correctness is better for TGT. In both housing and RNA, the coverage and correctness
 171 are better at less features and similar for more features. Overall, these results support claim 1, especially for a small
 172 amount of features.

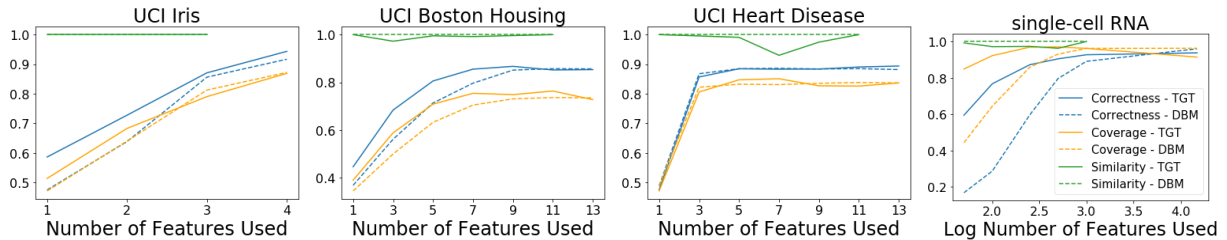


Figure 1: Comparison of the metrics(Correctness, Coverage, and Similarity) across different datasets for reproduction experiments.

⁷<https://userinfo.surfsara.nl/systems/lisa>

173 **5.2 Explaining Causal Structure in the Synthetic Data**

174 We verify the claim that TGT identifies the causal structure in the data (claim 3). The synthetic dataset is generated
 175 same as the original paper, i.e. $x_1, x_2 \sim \mathcal{N}(0, 0.2) + \text{Bern}(0.5)$, $x_2 \sim \mathcal{N}(0, 0.05)$, $x_4 \sim x_1 + \mathcal{N}(0, 0.05)$. Note that
 176 this dataset has four different clusters, caused by the first two dimensions x_1 and x_2 , x_3 is noise, and x_4 is correlated
 177 with x_1 and x_2 . The authors claim that for this synthetic data, TGT is able to find that x_4 is not the cause for any group.
 178 However, the said claim cannot be re-verified. Interestingly, the re-run of their code doesn't provide the justification
 179 either to the degree as mentioned in the paper. We observe that both TGT and DBM are able to identify x_3 is not
 180 causing any groups. Thus, in this scenario, both TGT and DBM are comparable. Refer table 2 for the explanations
 181 obtained. We, hereby note, that the explanations vary across multiple runs and we use the experimental setup same as
 182 the original authors. However, the values across the third dimension are consistently approximately 0.

183 **5.2.1 Feature modifications**

184 For each of the UCI datasets, the original authors add a 'corrupted' version where an extra cluster is added with artificial
 185 feature modification. With the exception of the modified features, the corrupted class is a copy of a chosen target class.
 186 They train TGT explanations using both the original scVIS model for the respective dataset and a model retrained on the
 187 corrupted dataset. We reproduce these experiments to see if TGT correctly attributes the difference between the target
 188 and corrupted class to the right features. Refer to Appendix A.1 for the illustrated figures and description. Overall, we
 189 observe that TGT is unable to identify the modifications to as good a precision as reported in the original paper. TGT is
 190 able to identify the modification for the UCI Iris dataset. For UCI Heart Disease Dataset (figure 7), it does not identify
 191 the features modified and on the UCI Boston Housing Dataset (figure 6), it identifies noisy modifications. However,
 192 with the retrained scVIS model and new representations, TGT is consistent in identifying the modifications across all
 193 the datasets.

194 **5.3 Results beyond original paper**

195 **5.3.1 PyTorch replication**

196 We also replicate the TGT algorithm in PyTorch. Our Pytorch implementation includes the entire method along with
 197 the scVIS clustering method. In our implementation, we use the Scikit learn ⁸ kmeans module for our cluster selection
 198 as opposed to the manual clustering in the Tensorflow implementation. However, our number of clusters argument to
 199 the kmeans algorithm was informed by the learned low-dimensional representations for each dataset. Due to differences
 200 in the clustering model and cluster selection, we cannot directly compare the coverage and correctness metrics between
 201 our Pytorch replication and the TensorFlow reproduction. We additionally experiment with our scaling extension to the
 202 TGT algorithm. In the scaling extension of the TGT algorithm, along with the δ (δ) parameters, each cluster now has a
 203 γ (γ) parameters. The transformation from cluster 0 to i is now given by: $t_{0 \rightarrow i} = e^{\gamma_i} \odot x + \delta_i$. The gammas (γ_s) are
 204 truncated just like the deltas and their $L1$ norm is added to the regularization term. Note that these transformations are
 205 strictly more expressive. If γ is the zero vector, these transformations reduce to regular TGT.

206 **5.3.2 UCI Heart Disease and UCI Boston Housing Dataset**

207 In figure 2 we see the results of our replication on the UCI Boston Housing and UCI Heart Disease dataset. For the UCI
 208 Boston Housing data, the TGT method seems to slightly outperform DBM both with and without scaling. This supports
 209 claim 1. The deltas (δ_s) and gammas (γ_s) show high similarity, supporting claim 2. For the UCI Heart Disease dataset,
 210 we do not see a difference in performance without scaling while TGT with scaling performs slightly worse.

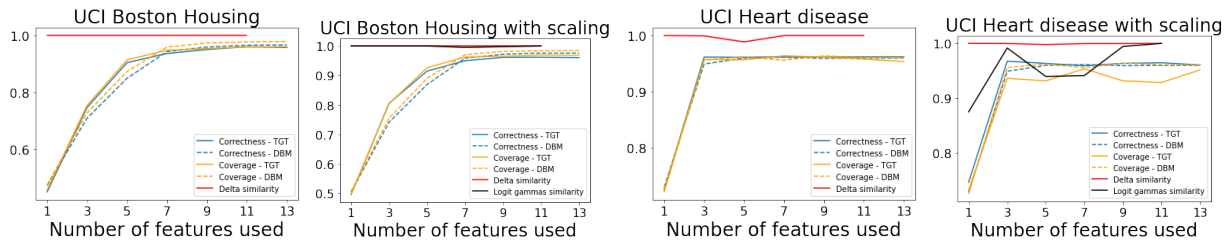


Figure 2: Results for the PyTorch replication for UCI Boston Housing and UCI Heart Disease dataset.

⁸<https://scikit-learn.org/stable/index.html>

	0 to 1		1 to 0		δ_1	δ_2	δ_3	δ_4	γ_1	γ_2	γ_3	γ_4
	C_r	C_v	C_r	C_v								
TGT	1	0.529	0.333	1.000	2.581	0.014	0.001	0.842	-	-	-	-
Scaling	1	0.529	0.520	1.000	2.608	0.023	-0.001	0.927	0.879	0.009	0.002	0.007

Table 3: The deltas(δ_s) and gammas(γ_s) for the mapping from group 0 to group 1 on the modified synthetic dataset for regular TGT and TGT with scaling. C_r and C_v indicate correctness and coverage, respectively.

211 5.3.3 Breast Cancer Wisconsin (Diagnostic) and Pima Indians Diabetes Database

212 In figure 3 we see the results for the Pima Indians Diabetes and Breast Cancer Wisconsin (Diagnostic) dataset. For the
 213 diabetes dataset, TGT with and without scaling outperforms DBM when more than one feature is used. This supports
 214 claim 1. Since the delta (δ) similarity is close to 1, claim 2 is also supported. For the Breast Cancer dataset, we see
 215 similar performance for DBM and TGT and slightly worse performance with scaling. The deltas (δ_s) still have high
 216 similarity, supporting claim 2.

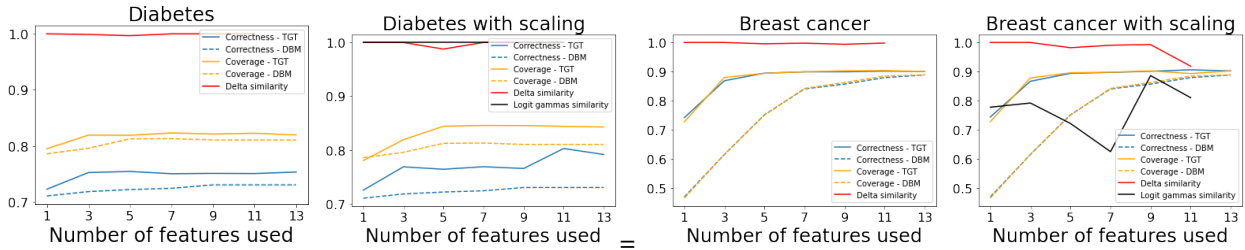


Figure 3: Results of the PyTorch replication on PIMA Indians Diabetes and Breast Cancer Wisconsin (Diagnostic) Dataset.

217 5.3.4 Scaling extension

218 In figure 2, we can see the difference in performance on two dataset included in the original experiments. Scaling does
 219 not seem to improve performance on the UCI Boston Housing dataset and slightly decreases performance on the UCI
 220 Heart Disease dataset. The similarity of the gammas (γ_s) is mostly above 0.9.

221 In figure 3, we see the same metrics for the Breast Cancer and Pima Indians Diabetes Dataset. For the Diabetes
 222 dataset, the performance improves slightly and the gammas(γ_s) show high similarity. For the Breast Cancer dataset, the
 223 performance is about the same but the gammas(γ_s) show relatively low similarity.

224 Altogether, these results suggest that the addition of scaling does not significantly improve the accuracy and correctness
 225 while making the transformations more complex. Based on our experiments, we do not recommend the addition of
 226 scaling in the explanation functions, and conclude that the original TGT is expressive enough.

227 5.3.5 Experiment with Modified Synthetic Data

228 In order to study the efficacy of the proposed scaling function, we perform experiments on the synthetic dataset. We
 229 modify one of the groups of points G by performing the operation $ax_i^k + b$, where i corresponds to the group number
 230 and k denotes the feature dimension which we modify. We define $a \sim \mathcal{U}(1.0, 2.0)$ and $b \sim \mathcal{U}(-0.5, 1.0)$. We add
 231 modified group G' into the original data D to get the new data D' . We follow the experiment setup from the original
 232 paper as: a) $r(G')$ should form a different group of its own. b) G' should be within the distribution of the original D .
 233 In this study, we want to investigate whether the TGT with scaling is able to recover the modifications, and if in doing
 234 so it affects the explanations between other groups. The sampling procedure gave $a=2.0$, $b=0.60$ and we keep $k=0$. We
 235 observe that the explanations with scaling are able to recover the modification to an approximate degree (scaling factor
 236 $e^\gamma \approx 2.38$, actual $a=2.0$), and give better correctness as compared to the regular TGT (refer figure 3). Interestingly, the
 237 translations explanations of the scaling extension are approximately equal to the deltas of the regular TGT. The exact
 238 results can be found in Table 3. Figures 8 and 9 in the appendix show the data spread and resulting translations.

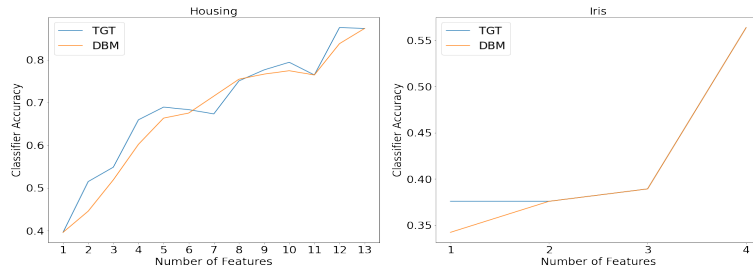


Figure 4: Classification accuracy of probing classifier at different sparsity levels for Housing (left) and Iris (right) dataset.

239 5.3.6 Experiments with Probing Classifier

240 To further investigate the efficacy of TGT explanations, we use a probing-classifier [2] as a proxy to study the qualitative
 241 differences of the features selected by TGT and DBM. For each cluster, we train a binary classifier with features ranked
 242 highest by TGT and DBM at different sparsity levels K . We compute the overall accuracy at each sparsity level using
 243 the ensemble of these binary classifiers. As can be seen in Figure 4, the results demonstrate that for sparser explanations,
 244 TGT selects features that lead to higher accuracy of the ensemble classifier than those selected by DBM. This further
 245 validates the paper’s claim that TGT leads to better sparse explanations as compared to DBM. Furthermore, we also
 246 use the probing classifier to understand the differences between the groups. For each pair of group, we train a Binary
 247 Linear Classifier to predict the group of a test point. We, then, investigate the feature importances of the classifier
 248 towards decision making. We ascertain that the features classifier give more importance to while decision making are
 249 the defining property of the class. Interestingly, we find that the more important features according to the classifier
 250 correspond to the explanations provided by the TGT algorithm. Refer to figure 11. This provides further evidence that
 251 TGT is able to find real distinctive signals as explanations.

252 6 Discussion

253 Based on the reproduction of the original experiments, claims 1 and 2 seem to hold, the experiments for claim 4 do not
 254 all support it, but the claim does seem to hold. Claim 1 and 2 seem to hold in particular for sparse explanations. The
 255 evidence for claim 3 is inconclusive. The coverage and correctness in our reproduction were not always the same as in
 256 the original paper. It is difficult to compare these metrics for different clustering outcomes, as they depend on the ϵ
 257 parameter which depends on the clustering.

258 A major difficulty in reproduction is the cluster selection. When retraining the scVIS model, the latent space repre-
 259 sentation structure changes. The authors provide no method as to determine the different clusters other than visual
 260 inspection. Cluster selection could be an explanation for the differences in results between the original experiments and
 261 our reproduction. To verify the results with more confidence, a robust method for cluster selection might be required.

262 References

- 263 [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado,
 264 Andy Davis, Jeffrey Dean, Matthieu Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous
 265 distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- 266 [2] Yonatan Belinkov, Sebastian Gehrmann, and Ellie Pavlick. Interpretability and analysis in neural NLP. In
 267 *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*,
 268 pages 1–5, Online, July 2020. Association for Computational Linguistics.
- 269 [3] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- 270 [4] Diogo V Carvalho, Eduardo M Pereira, and Jaime S Cardoso. Machine learning interpretability: A survey on
 271 methods and metrics. *Electronics*, 8(8):832, 2019.
- 272 [5] Jiarui Ding, Anne Condon, and Sohrab P. Shah. Interpretable dimensionality reduction of single cell transcriptome
 273 data with deep generative models. *bioRxiv*, 2017.
- 274 [6] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2014.

- 275 [7] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,
276 Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep
277 learning library. *arXiv preprint arXiv:1912.01703*, 2019.
- 278 [8] Gregory Plumb, Jonathan Terhorst, Sriram Sankararaman, and Ameet Talwalkar. Explaining groups of points in
279 low-dimensional representations, 2020.
- 280 [9] Karthik Shekhar, Sylvain W Lapan, Irene E Whitney, Nicholas M Tran, Evan Macosko, Monika Kowalczyk, Xian
281 Adiconis, Joshua Z Levin, James Nemesh, Melissa Goldman, Steven Mccarroll, Constance L Cepko, Aviv Regev,
282 and Joshua R Sanes. Comprehensive classification of retinal bipolar neurons by single-cell transcriptomics. *Cell*,
283 166:1308–1323.e30, 08 2016.
- 284 [10] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*,
285 9(11), 2008.
- 286 [11] Sahil Verma, John Dickerson, and Keegan Hines. Counterfactual explanations for machine learning: A review,
287 2020.
- 288 [12] Haozhe Xie, Jie Li, and Hanqing Xue. A survey of dimensionality reduction techniques based on random
289 projection, 2018.

290 **A Extra figures**

291 **A.1 Experiments on Corrupted datasets**

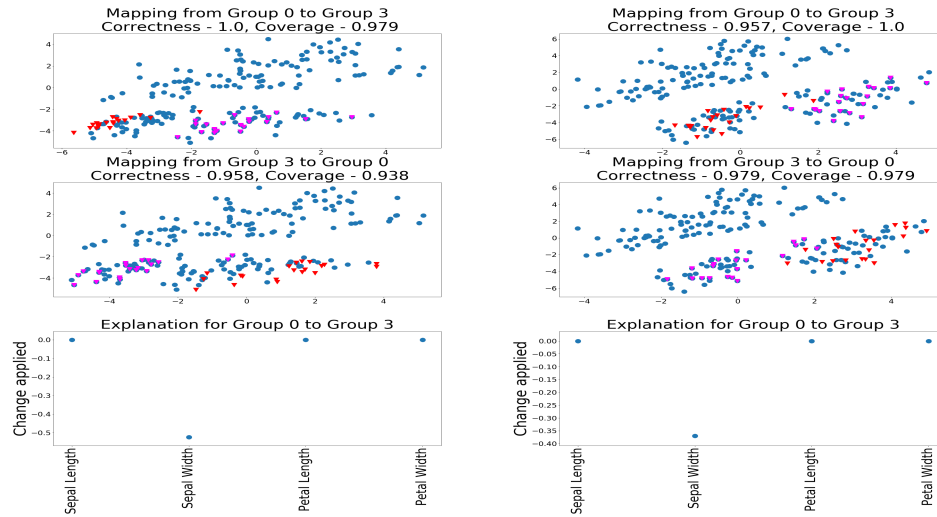


Figure 5: Explanation for corrupted features on UCI Iris dataset. Feature modified is 1(Sepal Width). Left: Visualization of the TGT explanations on the modified dataset. Right: Visualization of the TGT explanations with scVIS retrained on the modified dataset. We observe that the TGT explanations are robust to the modifications.

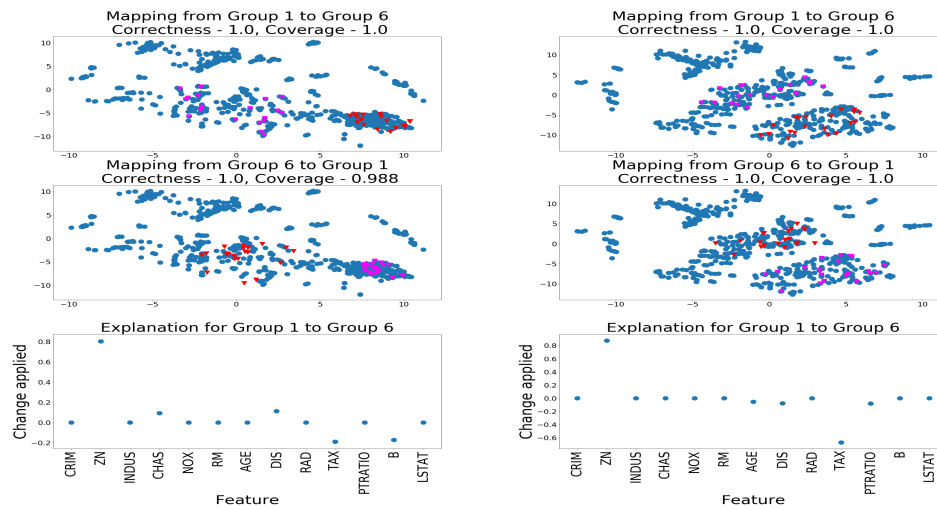


Figure 6: Explanation for corrupted features on UCI Boston Housing dataset. We modify the features 1(ZN) and 9(TAX). Left: Visualization of the TGT explanations on the modified dataset. We observe that TGT returns noisy explanations in this case. Right: Visualization of the TGT explanations with scVIS retrained on the modified dataset. With retrained scVIS model, TGT is able to recover the modifications.

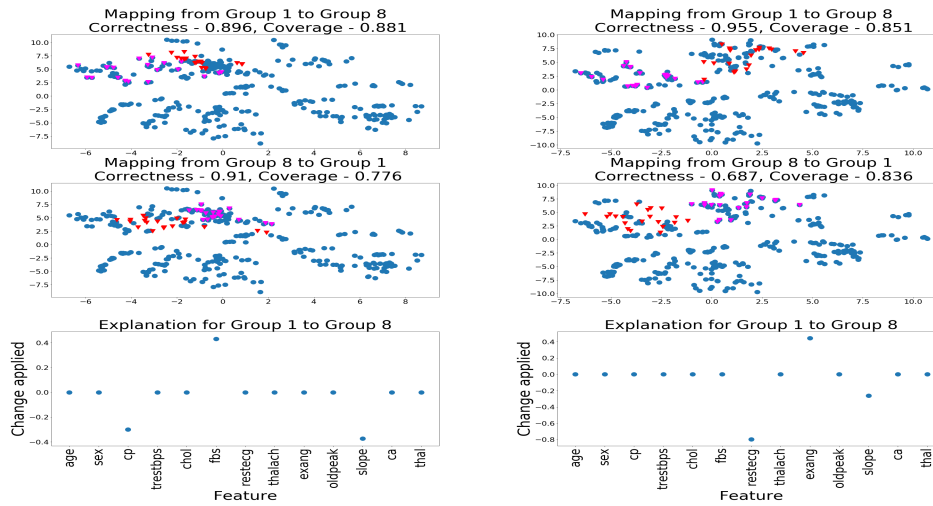


Figure 7: Explanation for corrupted features on UCI Heart Disease dataset. Left: Visualization of the TGT explanations on the modified dataset. We modified the 6(restecg) and 8(exang). However, the TGT recovers modifications in features 2(cp), 5(fbs), and 10(slope) instead. Right: Visualization of the TGT explanations with scVIS retrained on the modified dataset. With retrained scVIS model, TGT recovers the modified features along with 10(slope) feature. This observation does not entirely support the claim 4.

292 A.2 Synthetic data

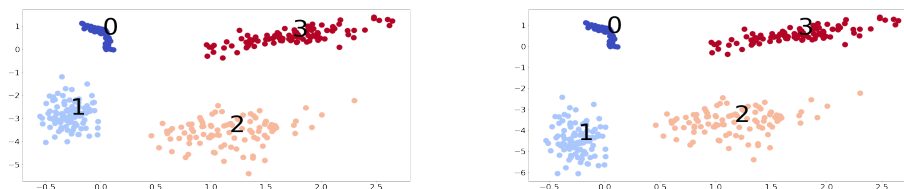


Figure 8: a) Synthetic data b) Synthetic data with the modification applied. We modify the data from group 1 across the 0^{th} dimension by $ax_1^0 + b$. Here a and b are 2.0, 0.60 respectively.

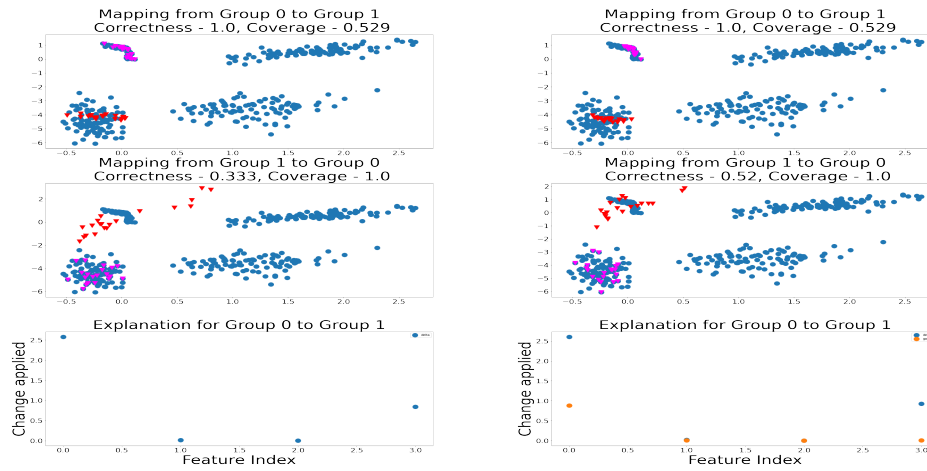


Figure 9: We compare the explanations from the TGT algorithm (left) and the TGT with scaling extension algorithm(right) on the modified synthetic data. We can observe that the TGT with scaling extension has better correctness, and is able to identify the scaling we have applied across the first dimension (i.e. $k=0$). The γ for this dimension is 0.87, which means the scaling factor is $e^\gamma \approx 2.38$. Moreover, the translation parameters are approximately same in both the variants of the TGT.

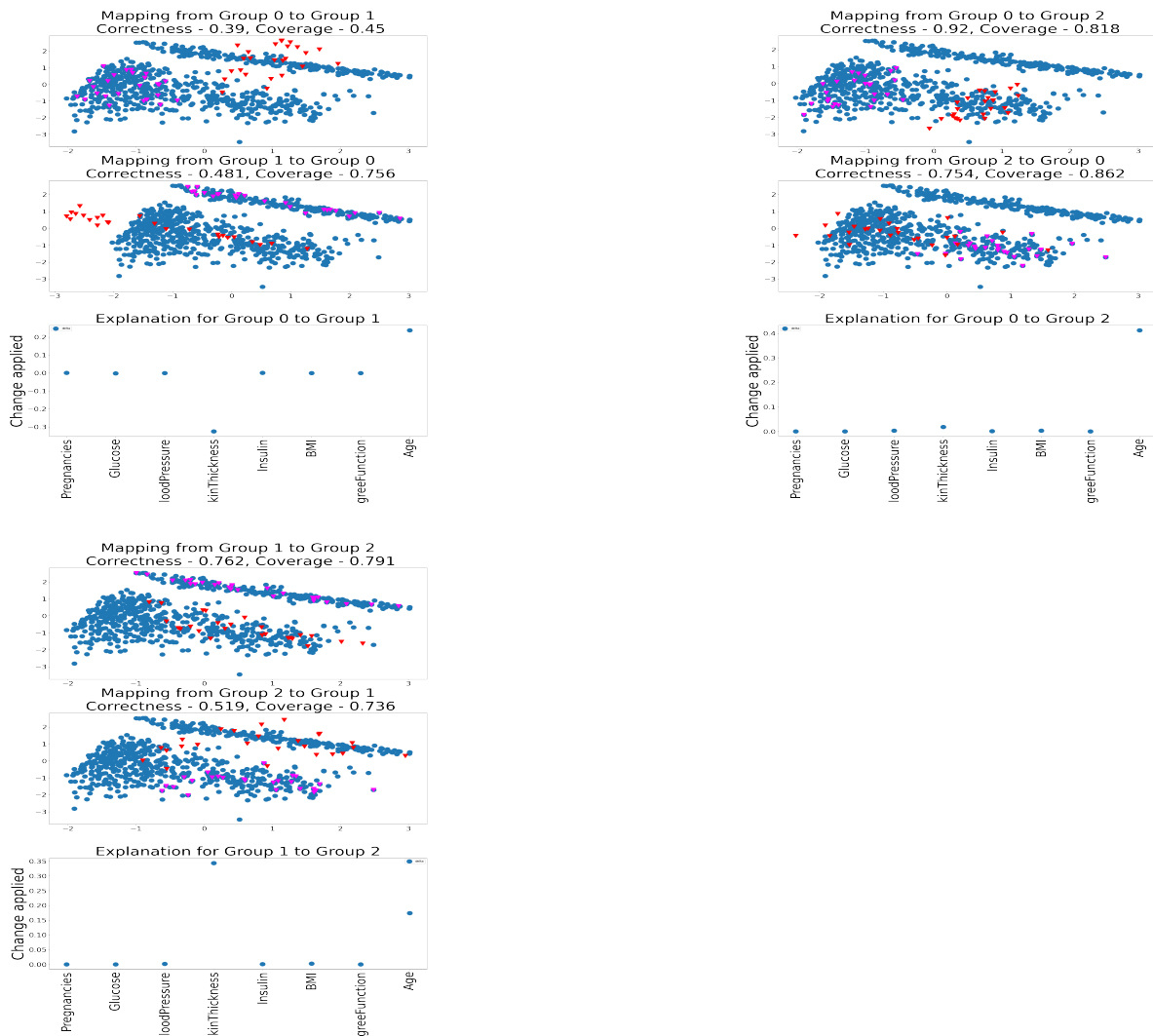


Figure 10: Explanations between different groups for the Pima Indians Diabetes Database.

293 **A.3 Probing Classifier and Feature Importance**

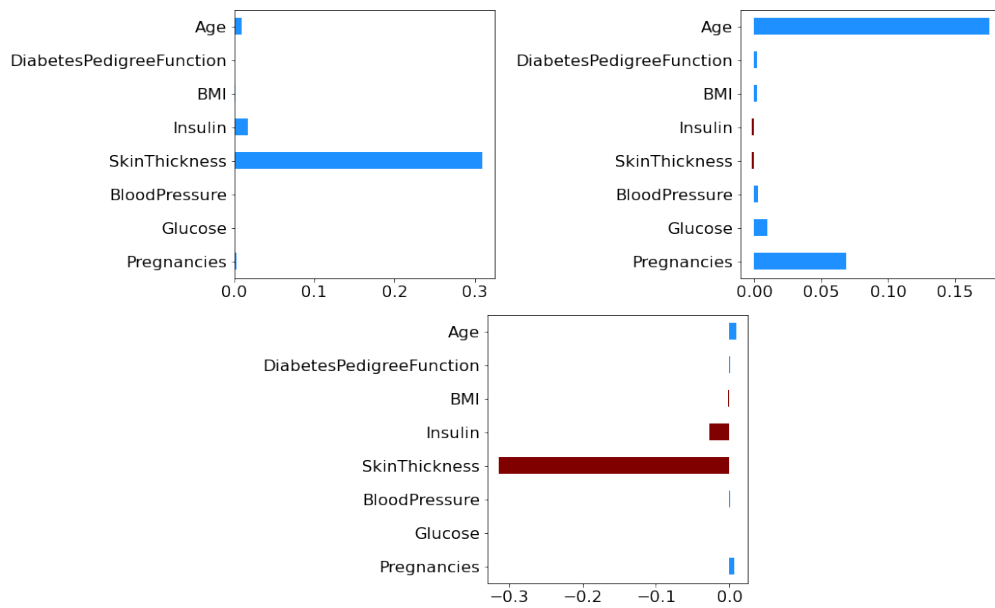


Figure 11: Feature Importance by the binary classifier for the Pima Indians Diabetes Database. a) (top-left): Feature importance for the classifier between groups 0 and 1. b) (top-right): Feature importance for the classifier between groups 0 and 2. c) (bottom-left): Feature importance for the classifier between groups 1 and 2. We note that the classifiers give significant feature importances to the features which correspond to the deltas (refer fig. 10).