
LOCAL EXPLANATION OF DIMENSIONALITY REDUCTION

A PREPRINT

Avraam Bardos
Aristotle University of
Thessaloniki, 54636, Greece
ampardos@csd.auth.gr

Ioannis Mollas
Aristotle University of
Thessaloniki, 54636, Greece
iamollas@csd.auth.gr

Nick Bassiliades
Aristotle University of
Thessaloniki, 54636, Greece
nbassili@csd.auth.gr

Grigorios Tsoumakas
Aristotle University of
Thessaloniki, 54636, Greece
greg@csd.auth.gr

May 2, 2022

ABSTRACT

Dimensionality reduction (DR) is a popular method for preparing and analyzing high-dimensional data. Reduced data representations are less computationally intensive and easier to manage and visualize, while retaining a significant percentage of their original information. Aside from these advantages, these reduced representations can be difficult or impossible to interpret in most circumstances, especially when the DR approach does not provide further information about which features of the original space led to their construction. This problem is addressed by Interpretable Machine Learning, a subfield of Explainable Artificial Intelligence that addresses the opacity of machine learning models. However, current research on Interpretable Machine Learning has been focused on supervised tasks, leaving unsupervised tasks like Dimensionality Reduction unexplored. In this paper, we introduce LXDR, a technique capable of providing local interpretations of the output of DR techniques. Experiment results and two LXDR use case examples are presented to evaluate its usefulness.

Keywords Dimensionality Reduction · Interpretable Machine Learning · Local Interpretations · Model-Agnostic · Black Box Models

1 Introduction

In recent years, the amount of data being produced on a regular basis, either through sensors or other technological means, has been rapidly increasing. Building on top of the availability of abundant data, Artificial Intelligence (AI) and Machine Learning (ML) have become an integral part of our daily lives. End-to-end transparency of AI and ML algorithms is crucial in high-risk applications such as disease diagnosis or criminal justice tools. Explainable AI (XAI) and Interpretable ML (IML) are fields that are researching how to make decision-making systems interpretable [1].

The data that are used for training AI systems are often high-dimensional, comprising thousands of features. This raises the computational complexity of building such systems on the one hand, while it also hurts their effectiveness on the other, as some features can be irrelevant, noisy or redundant. Dimensionality Reduction (DR) is a key solution for addressing these issues [2]. However, non-linear DR methods are opaque. The features extracted by such DR methods are difficult, if not impossible, to comprehend. By using them to solve supervised tasks, such as classification or regression, we are not able to identify which original feature contributed the most to a prediction.

Limited research has been conducted to address the interpretability of DR techniques. Related work includes the use of subspace projections [3], as well as adjusting existing methods designed to explain predictive models, such as LIME [4]. None of these methods however is general enough to be applicable to any DR technique. To fill this gap, we present

Local eXplanation of Dimensionality Reduction (LXDR), a model-agnostic technique for explaining the results of any non-linear DR technique, with the assumption that this technique will be able to reduce new instances. LXDR extracts interpretations via a local surrogate DR model, obtained by fitting a linear multi-target prediction model in the neighborhood of a reduced instance. Empirical results verify the interpretability capabilities of LXDR.

The rest of this paper is organized as follows. Section 2 presents background material and related work on interpretable dimensionality reduction. Section 3 presents LXDR and Section 4 our experiments. Finally, Section 5 presents the conclusions of this work, as well as future research directions.

2 Background and Related Work

This section starts by reviewing popular DR techniques. It then gives a brief introduction to the topic of interpretability, with a focus on surrogate models. Finally, it discusses related work on interpretable dimensionality reduction.

2.1 Dimensionality Reduction

A key use of DR is for enabling the visualization of high-dimensional datasets, as humans cannot comprehend a lot of dimensions. In addition, DR is beneficial when training ML models for supervised tasks, through the removal of noisy or redundant features that delay training and may even hurt generalization.

One of the most popular DR algorithms is Principal Component Analysis (PCA) [5]. Aiming to retain as much information from the data as possible through optimization, PCA constructs a reduced number of features, dubbed components, expressing linear relations among the original features. A variation of PCA is Sparse Principal Components Analysis (SPCA), which involves restricting as many coefficients as possible to being zero (sparsity constraint) [6]. In the same vein, SCoTLASS combines PCA and least absolute shrinkage and selection operator (LASSO) to produce sparser models [7].

Non-Negative Matrix Factorization (NMF) is a technique capable of both dimensionality reduction, and topic extraction, applicable to data with non-negative values [8]. Linear Discriminant Analysis (LDA) [9] attempts to project the data into components that maximize class separation. LDA is a supervised technique that requires ground truth information on the target (class or labels), unlike PCA and NMF.

Kernel PCA (KPCA) [10] is a variation of PCA that allows the representation of non-linear relations among the original features through the use of a kernel function. The t-distributed Stochastic Neighbor Embedding (t-SNE) [11] algorithm was designed for data visualization. Each high-dimensional data sample is represented by a two- or three-dimensional point in t-SNE, with a high likelihood of similar (dissimilar) samples being modelled by adjacent (distant) points. Other non-linear reduction techniques are Isometric mapping (IsoMap) [12] and Spectral Embedding [13]. Finally, neural Autoencoders (AE) are also used for dimensionality reduction [14].

2.2 Interpretability

Interpretability concerns the ability of ML models to provide reasons for their outputs. However, not every model is able to provide such information right out of the box. Such models are called *black box* models. As a result, a large body of research has focused on developing techniques that enable the interpretability of black box models.

Interpretability techniques can be either: a) global, offering explanations for the entire structure of the model, b) local, providing interpretations for a model’s output regarding a specific instance, or c) both. Furthermore, interpretability techniques can be classified as model-agnostic, applicable to any type of ML model, or model-specific, designed to interpret particular models and architectures [1].

A simple technique to interpret a black box is to use another, transparent model to explain it. The latter is called a *surrogate model*, and learns the data and the complex model’s predictions with the ultimate goal of extracting explanations. LIME, a state-of-the-art interpretability technique, follows this paradigm, but in local subspaces [15]. To interpret an instance’s prediction, LIME creates a synthetic local neighborhood around the instance. After using the black box to predict the generated neighbors, it trains a linear model on the neighborhood and the predictions, to extract weights for the features. These weights express the influence of a feature on the prediction of the instance, which is the final interpretation.

2.3 Interpretable Dimensionality Reduction

Because of their linear nature, PCA, NMF and LDA are intrinsically interpretable. As shown in Figure 1, internally, PCA learns a matrix of weights that linearly map the original dimensions N to the reduced dimensions $N_{reduced}$. The reduced data can therefore be obtained through a multiplication of the original data with this matrix.

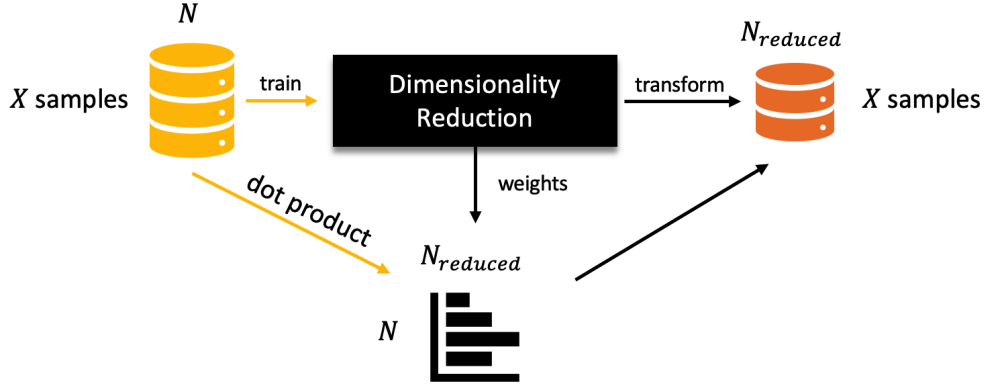


Figure 1: Reduced space data extraction through a combination of weights and original data

As an example, Table 1 shows the matrix of weights obtained when reducing the 4 dimensions of the Iris dataset [16] to 3 dimensions via PCA. These weights constitute an interpretation of the PCA transformation. We can see, for example, that among the original features, F_3 influences more the first component, and F_2 the second and third components.

Table 1: Components' weights extracted via PCA in the Iris dataset

Dimensions	F_1	F_2	F_3	F_4
1 st Component	0.361	-0.084	0.856	0.358
2 nd Component	0.656	0.73	-0.173	-0.075
3 rd Component	-0.582	0.597	0.076	0.545

These weights, however, are not provided by non-linear DR approaches. KPCA is uninterpretable when radial basis function or polynomial kernels are used. One solution, addressing the opaqueness of KPCA [3], employs a subspace on which features are projected before they are mapped into a high dimensional space using a kernel function. Hence, the mapping to the original features remains linear via a projection matrix, making kernel DR techniques interpretable. This technique, however, interferes with the algorithm of KPCA, which has a negative impact on the technique's performance. In addition, there is no publicly available implementation to experiment with it.

The non-parametric non-linear mapping of t-SNE is also difficult to interpret. Given that t-SNE uses neighborhoods to generate new representations for instances, known as embeddings, LIME was considered as a good fit for improving its interpretability in [4]. However, as LIME was designed for supervised models, a number of adjustments were explored. To begin, a LIME-inspired neighborhood creation process based on SMOTE was used. Afterwards, the neighbors were projected onto the reduced space. Lastly, a LASSO regression model was utilized to explain an instance's projection. Unfortunately, this work does not give details on how the second step is performed. This is an important issue, considering that t-SNE cannot project new samples and must create all the embeddings from scratch each time.

Manifold learning techniques, including IsoMap and Spectral Embedding, are also uninterpretable. Finally, an AE model can be interpreted intrinsically if it only contains an input layer, a latent layer, and an output layer. However, if it incorporates additional hidden layers, usually to increase the capacity, it becomes uninterpretable.

3 Our Approach

Given a black box DR model, the main idea of LXDR is to build a linear, hence interpretable, surrogate of that model. Instead of building a single global linear surrogate, LXDR follows a local approach, building a surrogate DR model around each particular instance for which an explanation of the black box model is needed. We consider two such uses cases: a) explaining decisions of supervised models operating on data that have been reduced by black box DR models (see Figure 7), and b) explaining outlier instances in the reduced space (see Figure 11) of a black box DR model.

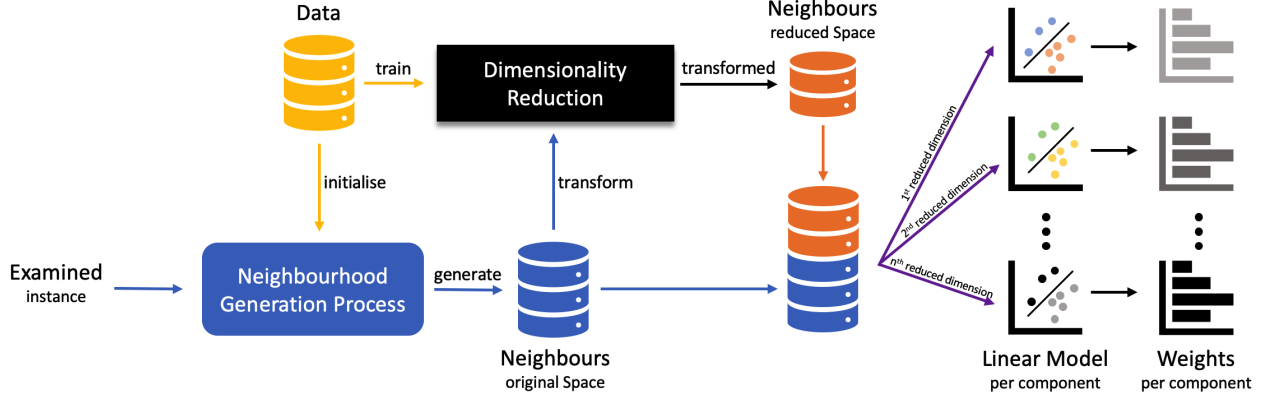


Figure 2: Workflow approach for the dimensionality reduction interpretation

To achieve its purpose, LXDR trains one local linear regression model per reduced dimension, using the original dimensions as predictor variables. The output of LXDR is therefore a matrix with the parameters of these linear models, representing the contribution of each original dimensions towards each of the reduced dimensions. This matrix can be used to compute a set of feature weights in the original space, through the inverse transform function as discussed later in Section 4.3.1, given a set of feature weights in the reduced space for the first use case. Moreover, utilizing LXDR to obtain insight on which original dimensions contribute more to the outlieriness of an instance in the reduced space is examined in the second use case. The local nature of LXDR is expected to lead to more faithful interpretations of non-linear black box DR models, compared to a single global linear surrogate. LXDR is model-agnostic, as it can be used in tandem with any black box DR model. Figure 2 presents the workflow of LXDR, while its pseudocode is given in Algorithm 1. The rest of this section describes LXDR in more detail.

Input : Instance x_i , Original Dataset D ,
Trained DR r , # of Neighbors K ,
Neighbourhood Generation Technique NG

Output : Coefficients $weights$

```

weights  $\leftarrow \emptyset$ 
# Generate a neighborhood for the given instance
 $B \leftarrow x_i \cup NG(D, K, x_i)$ 
# Apply DR to the neighbors
 $B' \leftarrow r(B)$ 
# For each dimension (component)
for dimension in reduced_neighbors do
    # Take the values of the selected dimension of the neighbors
     $rd \leftarrow B'[dimension]$ 
    # Train a linear model to predict those values
     $model \leftarrow train\_linear\_model(B, rd)$ 
    # Extract the weights as explanation for this component
     $weights \leftarrow weights \cup [model.get\_weights]$ 
end
return weights

```

Algorithm 1: The pseudocode of LXDR

We assume that the given black box model can be used to reduce an arbitrary instance from N dimensions to N_r , as we will use it to reduce the given instance x as well as a set of nearby instances. Note that some DR techniques, like t-SNE for example, do not offer this functionality. They act similarly to transductive models in supervised learning, reducing only the data they have been trained on. Retraining them on the original data expanded with the given instance and its neighbors would lead to a different DR model than the one we aim to approximate via a surrogate.

A crucial component of our technique is the Neighborhood Generation process NG . For a given instance x_i , we use an NG technique (Figure 3), in order to find or generate, depending on the technique, similar instances (neighbors). This utilization of a neighborhood, gives to our technique its local flavor.

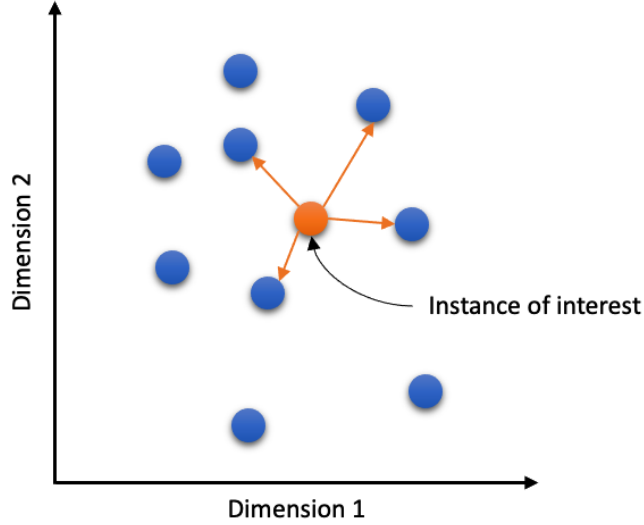


Figure 3: Instance neighborhood generation

For the neighborhood generation, we investigated two techniques. The first one is a very simple and commonly used technique (KNNs). The second one is LIME’s *NG* technique for tabular data, as presented in Section 2.

KNNs algorithm relies on the assumption that the most similar instances are close as well. Manhattan and Euclidean are the most often used metrics for calculating the distance between instances. The Euclidean distance was applied in our case. As inputs, this technique requires both the learned DR technique and the whole training dataset.

LIME was considered as a second *NG* technique. LIME was presented in Section 2.2 as a local model-agnostic interpretation technique that uses surrogate models to explain black box models. LIME produces these explanations by generating new synthetic data instances around a given instance. LIME does not directly provide us with these neighborhoods, and therefore, we made adjustments to its implementation to access the generated neighbors. LIME is a better option when only a portion of the training data is available, and it is also faster than KNNs.

Therefore, applying an *NG* technique, we are gathering the neighbors n^j of an instance x_i , in a set $B = \{x_i, n^1, n^2, \dots, n^K\}$ containing the given instance x and its neighbors $n^{(i)}$, $i = 1 \dots K$. However, in order to train our linear model, we must also have the target values, which in this case are the transformed representations of the neighbors. Therefore, we reduce the data set B using the black box DR model r : $B' \leftarrow r(B)$.

Data sets B and B' form a multi-target regression task, where B are the input values and B' the output values. To build the surrogate DR model, we simply train one linear regression model for each output variable, i.e. for each reduced dimension. We leave more sophisticated multi-target regression approaches [17, 18] for future work. However, we do employ a weighting mechanism, so that the linear models pay more attention to neighbors closer to the given instance x , as well as to x itself. Let $d(x, n^{(i)})$ denote the Euclidean distance of instance x from a neighbor $n^{(i)}$. The given instance takes a weight of 1, while each neighbor n^i is weighted with a weight $w^{(i)}$ given via the following equation:

$$w^{(i)} = e^{-2 \cdot d(x, n^{(i)})} \quad (1)$$

Finally, LXDR returns a matrix $weights \in \mathbf{R}^{N \times N_r}$ comprising the parameters learned by the linear models.

4 Experiments

This section focuses on the basic parameters of our technique, the setup, the quantitative experiments, and the qualitative experiments. The experiments’ code will be publicly available on the LXDR’s GitHub repo: (<https://github.com/avrambaldas/Interpretable-Unsupervised-Learning.git>).

In the setup section, we introduce the datasets utilized in the experiments, as well as the features we take into account for each one. The parameters of our technique and how they are employed, as well as the metrics used to evaluate LXDR,

are also presented. In the quantitative experiments section, we provide all the results supporting the effectiveness of our technique for the various datasets. We also perform a scalability study regarding the time performance of LXDR. Finally, the qualitative experiments’ section demonstrates the usefulness of our technique on the component explanation in two different use cases.

4.1 Setup

We experimented with three datasets, which demonstrate the universality of our technique to both regression and classification tasks. Therefore, two classification related datasets (Iris¹, Digits²) and one regression related (Diabetes³) were considered. More information can be found in Table 2.

Table 2: Information about datasets in terms of size, features, task, and classes

Dataset	Instances	Features	Task	Classes
Iris	150	4	Classification	3
Digits	1797	64	Classification	10
Diabetes	442	10	Regression	Regression

We applied DR to each dataset, trying to preserve at least 95% of the original information. Iris, for example, was reduced from 4 to 3 dimensions, Diabetes from 10 to 8 dimensions, and Digits from 64 to 25 dimensions. Furthermore, only 25% (449 out of 1,797 instances) of the entire data was used for Digits, in order to achieve faster execution time.

Additional parameters were introduced and employed for the experiments to improve the results, boosting our technique to more accurate results with fewer mistakes. These parameters, as well as their descriptions, are presented below:

K: K is the # of neighbors, and is set by the user, default = 10% of the total instances.

NG: NG can be either LIME or KNNs, default = KNNs.

AA: Auto Alpha (AA) is the linear model’s regularization strength (alpha). If AA is true, our approach would select the optimum alpha for the specific instance by experimenting with several alpha values. The selection is based on a self-evaluation procedure, which measures the performance of the linear model to the original data. Otherwise, if AA is false, the linear model’s default value of alpha is utilized.

4.1.1 Evaluation Metrics

We employ two different evaluation methods to assess the performance of our technique. The first one, *instance difference*, uses the extracted by LXDR weights $weights_{LXDR}$ to reduce the dimensions of an instance x and compare the reduced representation $x'_{LXDR} = x \cdot weights_{LXDR}$, with the reduced representation of the instance as provided by the DR technique $x' = r(x)$. The second metric, *weights difference*, assesses how much different the weights $weights_{LXDR}$ produced by LXDR are from the original weights of a DR technique $weights$. However, this metric is applicable only to intrinsically interpretable DR techniques.

$$ED(y, \hat{y}) = \sqrt{\sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

For these two metrics, we measure the Euclidean Distance (ED). In more depth, ED calculates the distance between two points. We consider that y and \hat{y} correspond to x' and x'_{LXDR} in the first metric, while $weights$ and $weights_{LXDR}$ in the second.

4.2 Quantitative Experiments

In this section, we evaluate the performance of LXDR. We compare LXDR to a global variation of LXDR, GXDR, which does not employ neighbors but instead is trained on the whole dataset. Three DR techniques, PCA, RBF KPCA, and AE, are utilized in the experiments to demonstrate our approach’s applicability and generalizability to a variety of

¹<https://bit.ly/3rMKxpA>

²<https://bit.ly/3jr6h6P>

³<https://bit.ly/3xfuTnS>

dimensionality reduction techniques. In the first set of experiments, we measure the *weights difference* between the weights provided by LXDR and GXDR applied in PCA, which is interpretable, and we can extract the original weights. The second set of experiments compare the performance of LXDR and GXDR based on the *instance difference* using RBF KPCA and AE as targeted DR techniques. We measure ED for both metrics across the three datasets.

The LXDR parameters we use are $AA = True$, $NG = KNNs$, and $K = \{50, 150, 750\}$, for the Iris, Diabetes, and Digits datasets, respectively. In our evaluations, the KNNs NG technique outperformed the LIME technique, in terms of the aforementioned metrics. As a result, KNNs NG is used in the experiments presented below. However, LIME is still a good alternative, particularly when the training dataset is unavailable or just partially available. Finally, we generate 25 synthetically generated datasets with 1K samples and $\{|F| = 10l | l \in \mathbb{N}, n \leq 25\}$ number of features, to measure the weights difference and instance difference of LXDR with different number of neighbors $K = \{|F|/1000, 2|F|/1000, 3|F|/1000\}$, as well as to measure its time performance.

4.2.1 Weights Difference

We measure the average *weights difference* between the weights produced by LXDR for each instance and the interpretable DR approach (PCA) in this set of experiments. We also compare LXDR’s performance with that of its global variant, GXDR. In Table 3, we compare LXDR and GXDR, and we see that LXDR outperforms its variation in all cases based on ED. We can also see that the performance difference is significant in two of the three examples.

Table 3: LXDR and GXDR performance based on the *weights difference* metric across datasets. For visualization purposes, all values have been multiplied by 100

Dataset	ED	
	LXDR	GXDR
Iris	0.9743	80.7851
Diabetes	0.1669	9.7474
Digits	5.3172	5.3961

We also evaluate LXDR performance using the synthetic datasets we produced based on this metric. Figure 4 depicts LXDR with different values for K to perform in these datasets. Using greater numbers in K , such as $1000 \cdot 3/4 = 750$, clearly resulted in more accurate weights. However, it is still local because LXDR weights the neighbors as shown in Eq. 1, and this is why it outperforms its global counterpart. GXDR performance in this experiment ranged from $ED = 0.029 \cdot 100 = 2.9$ to $ED = 0.204 \cdot 100 = 20.4$, for 10 to 250 features, respectively, and was not included in the figure for visualization purposes.

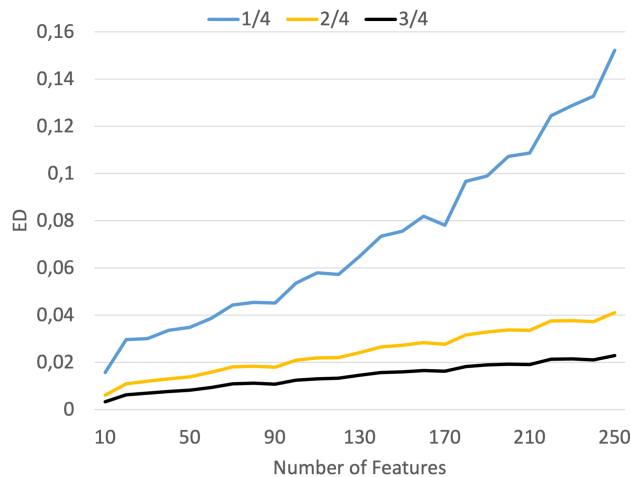


Figure 4: LXDR performance based on *weights difference* on synthetic datasets for three distinct *number of neighbors* choices. For visualization purposes, all values have been multiplied by 100

4.2.2 Instance Difference

In this set of experiments, we measure the performance of LXDR and GXDR, in terms of instance difference, on reducing instances based on the weights produced by these techniques. In Table 4, we notice that LXDR outperforms

Table 4: LXDR and GXDR performance based on the *instance difference* metric across datasets. For visualization purposes, all values have been multiplied by 100

Dataset	DR	ED	
		LXDR	GXDR
Iris	KPCA	3.9099	7.1388
	AE	4.9070	8.3904
Diabetes	KPCA	2.0346	8.0876
	AE	7.8590	10.6608
Digits	KPCA	3.9550	4.4079
	AE	29.4543	39.4562

GXDR in every dataset. An interesting observation here is that both techniques produce better results when interpreting the output of KPCA in contrast to AE. This can be due to the higher complexity of the AE model.

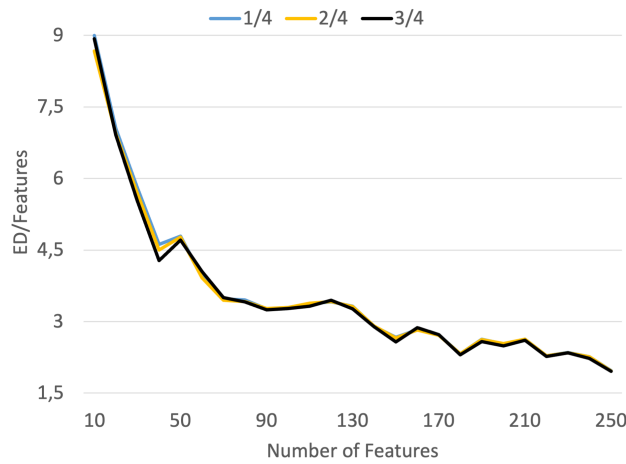


Figure 5: LXDR performance based on *instance difference* on synthetic datasets for three distinct *number of neighbors* choices. For visualization purposes, all values have been multiplied by 100

Furthermore, we assess LXDR’s performance on this metric using the synthetic datasets. Figure 5 illustrates how LXDR performs on these datasets with varying K values. In this metric, it is quite ambiguous which number of neighbors is preferable, but $1000 \cdot 3/4 = 750$ appears to be the best option. In particular, the average performance across the datasets is $1000 \cdot 1/4 = 3.56$, $1000 \cdot 2/4 = 3.52$, and $1000 \cdot 3/4 = 3.51$, while GXDR has a mean value of 3.55.

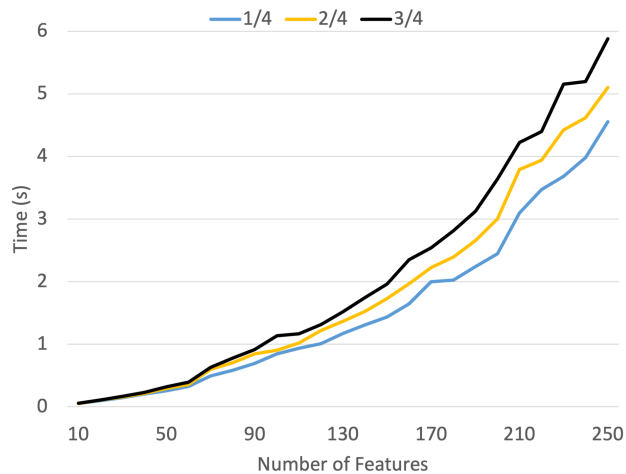


Figure 6: LXDR *time performance* on synthetic datasets for three distinct *number of neighbors* choices

4.2.3 Scalability

The final experiment in this section examines the LXDR’s time performance. Figure 6 shows the average time required by LXDR to generate an explanation for an instance for different values of K across the datasets. There is no noticeable difference in LXDR’s time performance for datasets with fewer features. However, with higher-dimensional datasets, we can see the technique’s response time increases as the number of neighbors increases. In the dataset with 250 features, for example, there is a 1-second difference between the three K alternatives. However, we did not optimize or parallelize this process, thus further optimization would result in improved time performance.

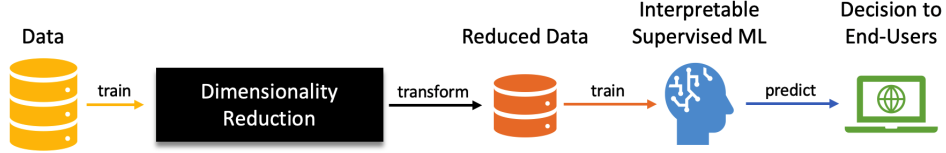


Figure 7: This workflow illustrates a pipeline containing a supervised decision-making system that makes use of data processed by a DR approach

4.3 Qualitative Experiments

We also demonstrate the use and effectiveness of LXDR in two use cases. The first one is through a step-by-step process, as presented in Figure 7. We use the Diabetes dataset to solve a regression problem with PCA as the DR technique. We choose an interpretable DR technique, in order to be able to compare LXDR to the ground truth weights of PCA. The second one concerns using DR techniques for visualization.

4.3.1 Use Case: Regression

We divide the dataset into a train set (80%) and a test set (20%). After training PCA, we reduce the dimensions of both the train and test sets from 10 to 8. In this manner, we preserve 95% of the original data. Then, we train a Ridge regression model, which serves as the decision-making mechanism. Trained and evaluated on the reduced datasets, the model achieves a MAE score of 42.08.

We then choose a random reduced instance $x'_i = [-0.73, -0.86, -0.07, -0.88, 0.17, -0.07, -0.1, -0.06]$ to explain the prediction of the model. The prediction of ridge for this instance is 140.03. Ridge regression, because it is intrinsically interpretable, can provide a set of coefficients. Because our reduced set has 8 dimensions, we also have 8 coefficients $coef = [8.73, -72.23, 28.48, 74.1, -27.07, -55.25, -15.1, 25.95]$. $coef$ presents the global weights of ridge, while $coef_{x'_i} = coef \times x'_i = [-6.38, 62.08, -1.9, -65.04, -4.51, 3.85, 1.47, -1.51]$ showcases the local interpretation for Ridge’s prediction regarding the examined instance x'_i .

$$x'_i = (x_i - mean) \cdot weights^T \quad (2)$$

$$x_i = x'_i \cdot weights + mean \quad (3)$$

Table 5: Contribution of features for each one of the 8 components from PCA

	C_1	C_2	C_3	C_4	C_5	C_6	C_7	C_8
F_1	0.1	-0.35	-0.82	-0.18	-0.37	-0.08	-0.09	-0.03
F_2	0.96	0.24	-0.01	-0.03	0.06	0.01	0.03	0.09
F_3	0.05	-0.3	0.18	0.29	-0.13	-0.2	-0.49	0.7
F_4	0.1	-0.34	-0.19	0.57	0.59	-0.26	-0.02	-0.32
F_5	0.04	-0.36	0.1	-0.53	0.39	0.02	-0.01	0.1
F_6	0.05	-0.23	0.13	-0.41	0.23	0.03	-0.35	-0.14
F_7	-0.11	0.11	-0.29	-0.15	0.47	0.16	0.23	0.49
F_8	0.12	-0.26	0.29	-0.16	-0.18	-0.12	-0.15	-0.32
F_9	0.09	-0.46	0.22	-0.0	-0.18	-0.23	0.74	0.17
F_{10}	0.1	-0.37	0.06	0.25	-0.06	0.89	0.02	-0.02

In PCA, an instance is reduced based on Eq. 2, while it is inverse transformed through Eq. 3. Influenced by this inverse transformation process of PCA, we will propagate this local interpretation $coef_{x'_i}$ to the original space to produce $coef_{x_i}$, through $coef_{x_i} = coef_{x'_i} \cdot weights + mean$ (Eq.4). PCA’s coefficients are presented in Table 5.

By doing this multiplication, we have $coef_{x_i}$, which is the local interpretation of the Ridge model’s prediction of the random instance in the original feature space. Therefore, $coef_{x_i} = [-6.38, 62.08, -1.9, -65.04, -4.51, 3.85, 1.47, -1.51]$. We also present this information in Figure 8.

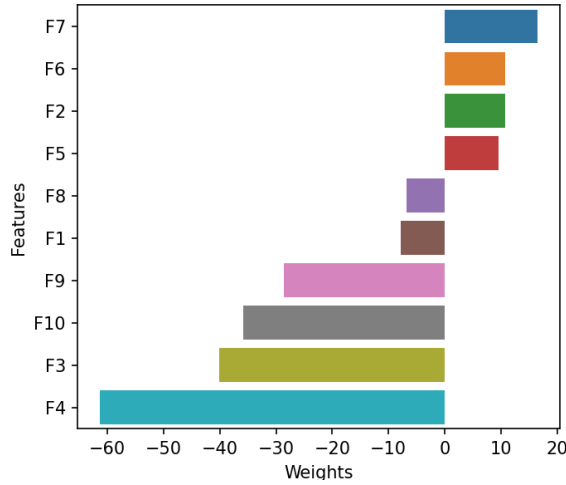


Figure 8: Local interpretation of x_i 's prediction in the original feature space

F_4 is the most important feature in the ridge’s prediction. As a result, we perform the following test. We decrease the value of the 4th feature from -0.121 to 0.0021 (the feature’s mean value). Because the importance is negative, we anticipate that the final prediction will be increased. We transform the modified instance with PCA before querying Ridge for another prediction. Ridge predicts a score of 146.86 , which is higher than the initial prediction of 140.03 . Hence, our hypothesis is correct.

Table 6: Contribution of features for each one of the 8 components from LXDR

	C_1	C_2	C_3	C_4	C_5	C_6	C_7	C_8
F_1	0.1	-0.35	-0.82	-0.18	-0.37	-0.08	-0.09	-0.03
F_2	0.96	0.24	-0.01	-0.03	0.06	0.01	0.03	0.09
F_3	0.05	-0.3	0.18	0.29	-0.13	-0.2	-0.49	0.7
F_4	0.1	-0.34	-0.19	0.57	0.59	-0.26	-0.02	-0.32
F_5	0.04	-0.36	0.1	-0.53	0.39	0.02	-0.01	0.1
F_6	0.05	-0.23	0.13	-0.41	0.23	0.03	-0.35	-0.14
F_7	-0.11	0.11	-0.29	-0.15	0.47	0.16	0.23	0.49
F_8	0.12	-0.26	0.29	-0.16	-0.18	-0.12	-0.15	-0.32
F_9	0.09	-0.46	0.22	-0.0	-0.18	-0.23	0.74	0.17
F_{10}	0.1	-0.37	0.06	0.25	-0.06	0.89	0.02	-0.02

Nevertheless, a lot of DR techniques, mostly non-linear, cannot provide information about the *weights* out of the box. Hence, a technique like LXDR can provide this kind of interpretation. Consider the PCA as a black box, and therefore we do not have the weights of Table 5. We are using LXDR to provide a local explanation of how the reduced representation of the examined instance was influenced by the original features. The parameters we select in this example are: AA is True, NG is KNNs, K is 150. By initializing LXDR, we request an interpretation. The output of LXDR is visible in Table 6.

Furthermore, we provide the *weights difference* between PCA and LXDR weights. The error in terms of ED is 0.0001 . The weight differences can also be seen in Figure 9, where the differences are extremely small. It is clear that the weights are nearly identical. As a result, even if PCA could not provide this interpretation intrinsically, we could have computed such weights accurately with LXDR. The last comparison is to reduce the original instance using the PCA and LXDR weights. As seen in Figure 10, the reduced representations of the instance are almost identical. Indeed, the ED *instance difference* is $5.54e - 05$.

Researchers and end-users who want to evaluate and implement LXDR in their solutions should use the provided metrics for parameter tuning. When the DR technique is uninterpretable, the parameters should be tuned and the *instance difference* should be measured. On the other hand, if the DR is interpretable, both the *instance* and the *weights differences* should be used to tune the parameters.

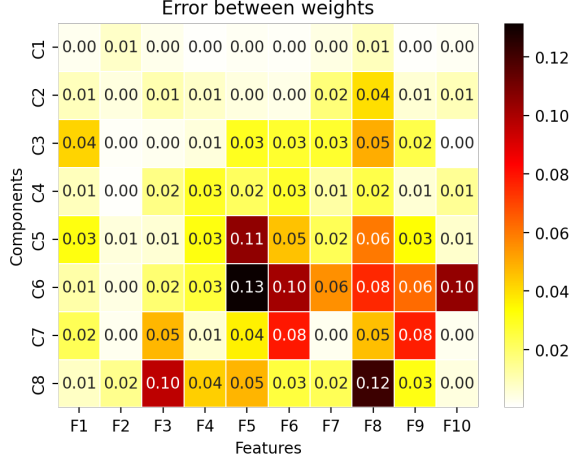


Figure 9: Weights difference between PCA and LXDR. For visualization purposes, all values have been multiplied by 1000

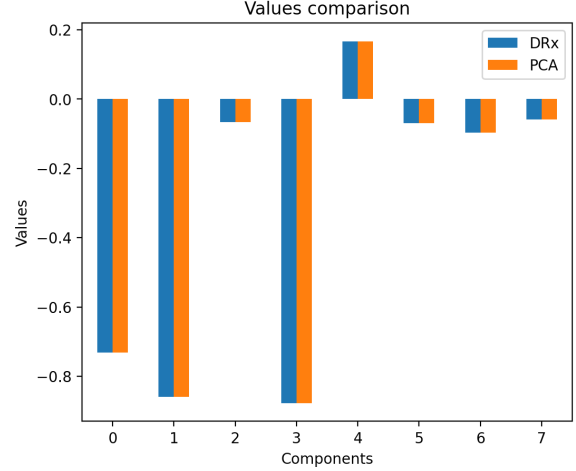


Figure 10: Instance difference between PCA and LXDR

4.3.2 Use Case: Visualization

DR techniques are frequently utilized to facilitate data visualization by reducing data to two or three dimensions. In this use case, we are using the same dataset, but we reduce the dimensions to two to visualize them using RBF KPCA.

In Figure 11, we visualize the data. The dataset concerns a regression problem, and the color map represents the target value. It is visible that on the right part of the figure the target values tend to be higher. However, there is a particular point (circled in red) that seems to be an outlier; it has a high target value while lying far from samples with similar targets. Employing LXDR we can investigate why this point has a high value on the first component (C_1).

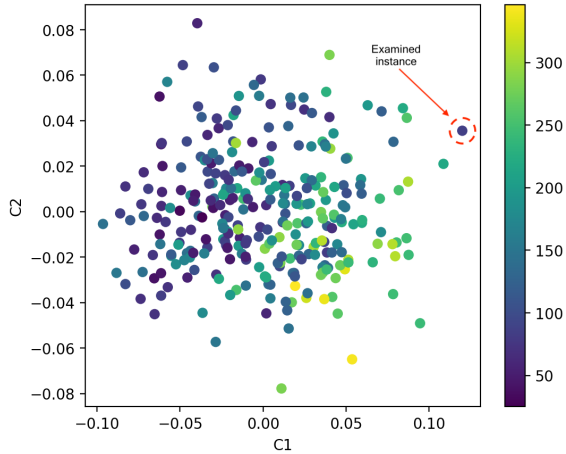


Figure 11: Visualized data after KPCA reduction to 2 dimensions

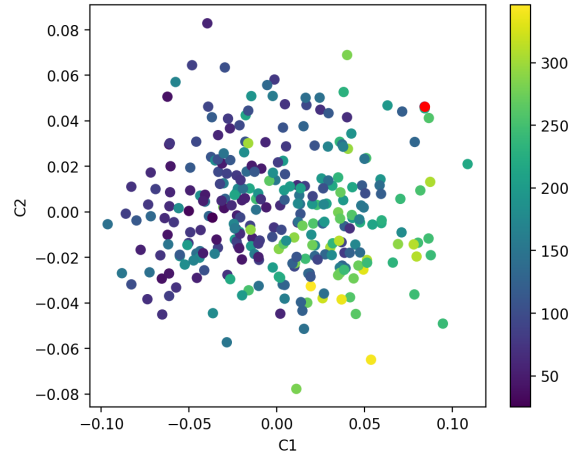


Figure 12: Visualized data after KPCA reduction to 2 dimensions with tweaked x_e instance

This instance’s original representation is $x_e = [0.01, 0.05, 0.03, -0.0, 0.15, 0.2, -0.06, 0.19, 0.02, 0.07]$. The reduced representation through KPCA is $x'_e = [0.12, 0.04]$. Figure 11 reveals that the abnormal value of x_e concerns C_1 ’s 0.12 value. Using LXDR we acquire the weights representing the influence of the input features to each component. The weights regarding the component we examine, C_1 , are $[0.08, 0.08, 0.14, 0.11, 0.16, 0.14, -0.11, 0.21, 0.16, 0.14]$. The 8th feature appears to have the highest impact on the value of C_1 . We repeat the experiment as earlier. We change the value of F_8 from 0.19 to 0 (the feature’s mean value), and then we reduce the instance again through KPCA.

The modified instance’s representation is $x'_{e,tweaked} = [0.08, 0.05]$. In Figure 12, we notice how the instance got closer to the center, and it is no longer considered an outlier. If we wanted to use this data for training, we should either exclude this point from the training set or investigate whether its target value is incorrect.

5 Conclusions

We introduced LXDR in this paper, a technique for providing local interpretations for DR techniques. We were able to investigate and determine which original features contributed the most to a specific reduced component by combining neighboring data with the predictions of a linear model. Experiments were conducted to support the effectiveness of LXDR. In three different datasets, we compared the results of LXDR to those of its global variant, GXDR, using two different metrics. LXDR outperforms GXDR in every experiment. Furthermore, we investigated how datasets with larger dimension spaces affect LXDR performance.

We also performed a scalability analysis to assess LXDR’s response time. While for smaller datasets applying LXDR is fast, in larger datasets the response time exceeds the one second. Therefore, in order to utilize LXDR in an online, production application further optimization of the process can be applied. For example, parallelization of the training process of the different linear models, or precomputed neighborhoods.

Regarding the qualitative experiments, we presented two use cases. The first use case concerns the usefulness of LXDR in interpreting the results of a regression model trained on data reduced by a DR technique. The second involves inspecting individual points in a plot that appear to be anomalous after being projected by a DR technique.

In the future, we can investigate new parameters for LXDR to improve performance. We intend to adapt LXDR to be applicable to DR techniques that cannot reduce new instances, such as t-SNE. The application of LXDR to image and textual datasets should be investigated as well. Ways to improve the evaluation of LXDR performance, inspired by the evaluation metrics used in the IML research on supervised tasks, will also be examined. Finally, a user study can be conducted to determine how effective a technique like LXDR may be for (non)expert users.

Acknowledgments

The research work was supported by the Hellenic Foundation for Research and Innovation (H.F.R.I.) under the “First Call for H.F.R.I. Research Projects to support Faculty members and Researchers and the procurement of high-cost research equipment grant” (Project Number: 514)

References

- [1] Amina Adadi and Mohammed Berrada. Peeking inside the black-box: A survey on explainable artificial intelligence (XAI). *IEEE Access*, 6:52138–52160, 2018.
- [2] Carlos Oscar Sánchez Sorzano, Javier Vargas, and A Pascual Montano. A survey of dimensionality reduction techniques. *arXiv preprint arXiv:1403.2877*, 2014.
- [3] Chieh Wu, Jared Miller, Yale Chang, Mario Sznajder, and Jennifer G. Dy. Solving interpretable kernel dimensionality reduction. In *NeurIPS*, pages 7913–7923, 2019.
- [4] Adrien Bibal, Viet Minh Vu, Géraldine Nanfack, and Benoît Fréney. Explaining t-sne embeddings locally by adapting LIME. In *28th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning, ESANN 2020, Bruges, Belgium, October 2-4, 2020*, pages 393–398, 2020.
- [5] Michael E. Tipping and Christopher M. Bishop. Mixtures of probabilistic principal component analysers. *Neural Comput.*, 11(2):443–482, 1999.
- [6] Hui Zou, Trevor Hastie, and Robert Tibshirani. Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, 15(2):265–286, 2006.
- [7] Doyo G Enki, Nickolay T Trendafilov, and Ian T Jolliffe. A clustering approach to interpretable principal components. *Journal of Applied Statistics*, 40(3):583–599, 2013.
- [8] A. Cichocki, R. Zdunek, A.H. Phan, and S. Amari. *Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation*. Wiley, 2009.
- [9] Petros Xanthopoulos, Panos M. Pardalos, and Theodore B. Trafalis. *Linear Discriminant Analysis*, pages 27–33. Springer New York, New York, NY, 2013.
- [10] Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Kernel principal component analysis. In Wulfram Gerstner, Alain Germond, Martin Hasler, and Jean-Daniel Nicoud, editors, *Artificial Neural Networks — ICANN’97*, pages 583–588, Berlin, Heidelberg, 1997. Springer Berlin Heidelberg.
- [11] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008.

- [12] Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319–2323, 2000.
- [13] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Comput.*, 15(6):1373–1396, 2003.
- [14] Yasi Wang, Hongxun Yao, and Sicheng Zhao. Auto-encoder based dimensionality reduction. *Neurocomputing*, 184:232–242, 2016.
- [15] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144. ACM, 2016.
- [16] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [17] G. Tsoumakas, E. Spyromitros-Xioufis, A. Vrekou, and I. Vlahavas. Multi-target regression via random linear target combinations. In Toon Calders, Floriana Esposito, Eyke H{ü}llermeier, and Rosa Meo, editors, *Proceedings of ECML PKDD 2014*, volume 8726 LNAI, pages 225–240. Springer Berlin Heidelberg, 2014.
- [18] E. Spyromitros-Xioufis, G. Tsoumakas, W. Groves, and I. Vlahavas. Multi-target regression via input space expansion: treating targets as inputs. *Machine Learning*, 104, 2016.