

UNIVERSIDADE FEDERAL DE SÃO CARLOS

COMPUTER SCIENCE

FINAL REPORT

A Study of the Isomap Algorithm and Its Applications in Machine Learning

Author:

Lucas O. DAVID

Supervisor:

Dr. Alexandre M. LEVADA

November 16, 2015

Contents

1 Abstract	3
2 Introduction	4
3 Relevant Background	5
3.1 Data set	5
3.1.1 Example of a canonical data set	5
3.1.2 Data set as a collection of vectors in the \mathbb{R}^n	6
3.1.3 Modern Problems and Applications	7
3.2 Probability Theory	7
3.2.1 Feature Normalization	7
3.2.2 Centering Matrix	8
3.2.3 Variance	8
3.2.4 Covariance	9
3.2.5 Correlation	10
3.3 Numerical Analysis	11
3.3.1 Eigenvalues and Eigenvectors of a Matrix	11
3.3.2 Spectral Decomposition of a Matrix	11
3.3.3 Singular Value Decomposition	12
3.4 Topology	12
3.4.1 Manifolds	12
3.5 Graph Theory	13
3.5.1 Graphs	13
3.5.2 Related Problems	15
3.6 Machine Learning	17
3.6.1 Machine Learning Algorithms	18
3.6.2 Multi-class Classification	22
3.6.3 Evaluating learners	23
3.6.4 Examples of learning	25
4 Linear Dimensionality Reduction	26
4.1 Principal Component Analysis	27

4.1.1	Study of the PCA Algorithm	28
4.1.2	Formalization of the PCA Algorithm	29
4.2	Multidimensional Scaling	29
4.2.1	Study of the MDS	30
4.2.2	Formalization of the Multidimensional Scaling Method	32
4.3	Classification and Regression Over Linearly Reduced Data Sets	32
4.3.1	K data set	33
4.3.2	The Iris flower data set	34
4.3.3	The Digits data set	34
5	Non-linear Dimensionality Reduction	35
5.1	The Isomap Algorithm	37
5.1.1	Study of the Isomap Algorithm	37
5.1.2	Formalization of the Isomap Algorithm	39
5.1.3	Computational Complexity	39
5.2	Classification and Regression Over Data Sets Reduced with Isomap	41
5.2.1	The Swiss Roll Data Set	41
5.2.2	The Digits Data Set	42
5.2.3	The Leukemia Data Set	43
5.3	Applicability and Limitations of Isomap	44
5.3.1	Infeasibility on Highly Dense Data Sets	44
5.3.2	Necessary Settings for Convergence	44
6	Final Considerations	46

1 Abstract

This work aims to study the foundations of nonlinear dimensionality reduction with the algorithm known as Isometric Mapping (or simply Isomap).

All the experiments presented here were developed in computational environment, which can be found here: <https://github.com/lucasdavid/manifold-learning>.

2 Introduction

Throughout the years, machine learning and pattern recognition techniques have grown popular between both academical and the corporative sector. Their vast application fields and promising results indubitably contributed to our current scenario where not only computer scientists or mathematicians, but engineers, psychologists and many other groups have taken interest on how to adapt and apply these studies to their own problems.

Machine learning can help us to understand, analyze and process large amount of data and even take decisions based on it. The size of data, however, is not necessarily the predominant factor which dictates how machine learning models perform. We are, of course, interested in the information that lies within the data and how to efficiently extract it.

In the last 80 years, there were many different machine learning algorithms developed. Between those, many could successfully generalize low dimensional data [1]. In the other hand, problems of our world are often too complex and may be represented by high dimensional data. For example, images, sounds or text documents can be expressed as vectors of the \mathbb{R}^n , where each element corresponds to a pixel, wave signal or character, respectively. When analyzing these problems, we observed that many of the algorithms would often become unstable. **Dimensionality reduction** (or DR) then quickly became a key concept for minimizing the data size while maintaining its meaning.

Finally, dimensionality reduction has evolved into an extensive area. Nowadays, DR is not only applied to data reduction, specifically, but often employed in order to improve visualization or as preprocessing step for noise reduction.

3 Relevant Background

3.1 Data set

In the context of Computer Science, very often our goal is to develop machines that can assist or automate the process of solving real-world problems. Firstly, however, we must find ways to express these problems numerically.

Although the recurrent usage of the term in scientific work, there is not a clear definition established. It is possible, however, to observe the regular presence of four related features: grouping, content, relatedness e purpose [2]. For the scope of this work, the term data set is invariably associated with the idea of a collection of samples. Each sample is a sequence of features, where the i -th feature of all instances belong to a same set of symbols f_i .

Succinctly, consider S a sequence of samples and $F := \{f_i \mid f_i \text{ is a set of symbols}\}$. Then, the dataset D is defined as:

$$D := [d]_{ij} \mid d_{ij} \in f_j, \forall i \in [1, |S|], \forall j \in [1, |F|]$$

3.1.1 Example of a canonical data set

The table bellow illustrates an examples of data set, where each row represents a **sample**, and each column a **feature**.

	Sepal length	Sepal width	Petal length	Petal width	Species
1	5.1	3.5	1.4	0.2	I. setosa
2	4.9	3.0	1.4	0.2	I. setosa
3	4.7	3.2	1.3	0.2	I. setosa
...

Table 1: The first three samples of the Iris flower data set.

Iris flower is an example of data set broadly used in the machine learning demonstrations, being usually interpreted as a classification problem where the feature *Species* will be learned from its adjacent features. In that scenario, *Species* is denominated **target feature**.

3.1.2 Data set as a collection of vectors in the \mathbb{R}^n

A data set can have each one of its nominal features enumerated. I.e., mapped to an element of \mathbb{N} . Such set could then be expressed as a collections of vectors in \mathbb{R}^n . Consider the data set bellow:

	Age	Gen.	TB	DB	Alk.	Sgpt	Sgot	TP	ALB	A/G	S
1	65	Female	0.7	0.1	187	16	18	6.8	3.3	0.9	1
2	62	Male	10.9	5.5	699	64	100	7.5	3.2	0.74	1
3	62	Male	7.3	4.1	490	60	68	7	3.3	0.89	1
...

Table 2: The first three samples of the Indian Liver Patient Dataset (ILPD)

Composed by 583 samples and 11 features, the data set ILPD has a nominal feature $Gender := \{Male, Female\}$. $Gender$ can, of course, be mapped on $\{0, 1\}$. ILPD can finally be expressed by the figure bellow:

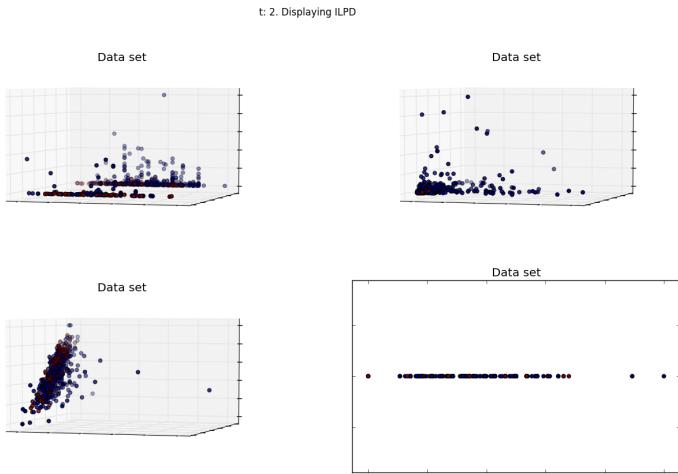


Figure 1: The data set ILPD mapped onto \mathbb{R}^n , where each of its features is an axis, except for $S := \{1, 2\}$, which was represented by the vertex's color.

Considering the many graphs required to display the data set, it is quite difficult to identify a plausible distribution for ILPD. We define here our first encouragement towards the study of dimensionality reduction: the identification of the most

significant features and plotting of those might yield simpler and more intuitive representations.

3.1.3 Modern Problems and Applications

Differently from Iris flower or even the Glass data set, modern problems are, in many cases, associated with large data sets. I.e., data sets that contain many samples and features. Although the high number of samples is essentially benefic, a high number of features might be irrelevant or even unconstructive to the learning process [3].

	A	B	...	AAAV
1	0.0111486888670454	-0.01541263850539861	...	0.007440367302352156
2	0.03016080450207878	0.1772161342899135	...	0.01309011094914101
...
8200	0.02680808496910305	-0.0320375843317954	...	0.1772161342899135

Table 3: A data set with 8200 samples and 100 features.

In many cases, there are indicatives that the data set lie near a lower-dimensional manifold embedded in the \mathbb{R}^n [4]. For the data set above, in special, the \mathbb{R}^{100} . A second encouragement can then be set: it is possible that the data set might be shrunk by combining similar (linearly dependent) features or eliminating the ones that poorly contribute towards the learning process. In order to do this, one must be able to qualify the “contribution” of each feature or even identify dependencies between features.

3.2 Probability Theory

3.2.1 Feature Normalization

Many of the methods ahead will require the data set to be centered in the origin. To center a data set X , we build a data set X' s.t. each column has zero mean and it is contracted by its standard deviation:

$$X'_{.j} = \frac{X_{.j} - \mu}{\sigma}, \text{ where}$$

1. $X_{\cdot j}$ is the j -th column of the matrix X .
2. μ is the mean of $X_{\cdot j}$.
3. σ is the standard deviation of $X_{\cdot j}$.

3.2.2 Centering Matrix

The symmetric matrix H is named the **centering matrix** when the multiplication of it by a vector X produces the same effect of subtracting the mean of the components from each component of X . H is defined as:

$$H = I_n - \frac{1}{n} \mathbf{1} \mathbf{1}^T, \text{ where:}$$

1. I_n is the identity matrix of order n .
2. $\mathbf{1}$ is the column vector of 1's.

3.2.3 Variance

Variance is the measure which describes how far the samples in a given set X vary. For this work scope, will consider only discrete probabilities. That is, if X represents a variate with known distribution $P(x)$, where $\sum P(x) = 1, \forall x \in X$ and population mean μ , then

$$\text{var}(X) = \frac{1}{n} (X - \mu) \cdot (X - \mu) = \frac{1}{n} \sum (x - \mu)^2, \forall x \in X$$

Remark 3.1 If $\text{var}(X) = 0$, then it is easy to assert that all variables in X assume the exact same value by using the elementary properties of the inner product.

Example 3.1 If $X = \{1, 2, -2, 4\}$ and $\mu = \frac{1}{4} \sum X_i = \frac{1+2-2+4}{4} = 1.25$, then

$$\begin{aligned} \text{var}(X) &= \frac{1}{4} \sum (X_i - \mu)^2 \\ &= \frac{(1 - 1.25)^2 + (2 - 1.25)^2 + (-2 - 1.25)^2 + (4 - 1.25)^2}{4} \\ &= 4.6875 \end{aligned}$$

Example 3.2 The variance of the Sepal length feature in the Iris flower data set can be calculated as:

$$\begin{aligned} \text{var}(X) &= \frac{1}{150} \sum (X_i - \mu)^2 \\ &= \frac{1}{150} [(5.1 - 5.84)^2 + (4.9 - 5.84)^2 + \dots + (5.9 - 5.84)^2] \\ &= \frac{102.17}{150} = .681122 \end{aligned}$$

3.2.4 Covariance

The covariance measures the variance of two random variates in respect to each other. Formally, if X and Y are two given random variates with known mean population distribution μ_X and μ_Y , respectively, then

$$\sigma(X, Y) = \frac{1}{n} (X - \mu_X) \cdot (Y - \mu_Y)$$

Simply putting, the covariance of two random variables X and Y can be interpreted as one of the following behaviors:

$$\sigma(X, Y) = \begin{cases} \sigma_{xy} > 0 \implies X \text{ tends to increase as } Y \text{ increases.} \\ \sigma_{xy} < 0 \implies X \text{ tends to increase as } Y \text{ decreases.} \\ \sigma_{xy} = 0 \implies X \text{ and } Y \text{ are completely unrelated.} \end{cases}$$

Remark 3.2 For a random variate X , $\text{var}(X) = \sigma(X, X)$.

Covariance Matrix of Features in a Data Set

For a data set D , where its i -th feature column is represented by $D_{\cdot i}$, the covariance between each one of its features can be represented by the matrix:

$$\begin{aligned}\Sigma &= [\sigma_{xy}]_{n \times n} \\ &= \begin{bmatrix} \sigma(D_{\cdot 0}, D_{\cdot 0}) & \sigma(D_{\cdot 0}, D_{\cdot 1}) & \cdots & \sigma(D_{\cdot 0}, D_{\cdot n-1}) \\ \sigma(D_{\cdot 1}, D_{\cdot 0}) & \sigma(D_{\cdot 1}, D_{\cdot 1}) & & \sigma(D_{\cdot 1}, D_{\cdot n-1}) \\ \vdots & & & \vdots \\ \sigma(D_{\cdot n-1}, D_{\cdot 0}) & \sigma(D_{\cdot n-1}, D_{\cdot 1}) & \cdots & \sigma(D_{\cdot n-1}, D_{\cdot n-1}) \end{bmatrix} \\ &= \frac{1}{n}(HD)^T HD \\ &= \frac{1}{n}D^T H^T HD \\ &= \frac{1}{n}D^T HD\end{aligned}$$

Where H is the **centering matrix**.

3.2.5 Correlation

“The correlation is a measure of the direction and strength of a linear relationship among variables.”

Let X and Y be two random variables, $r(X, Y)$, i.e., the correlation between X and Y is defined as:

$$r(X, Y) = \frac{\sigma(X, Y)}{\sigma_X \sigma_Y},$$

where σ_X and σ_Y are the standard deviations of the variables X and Y , respectively.

Correlation Matrix of Features in a Data Set

For a data set D , where its i -th feature column is represented by $D_{\cdot i}$, the correlation between each one of its features can be represented by the matrix:

$$corr(D) = \begin{bmatrix} \frac{\sigma(D_{\cdot 0}, D_{\cdot 0})}{\sigma_{D_{\cdot 0}}^2} & \frac{\sigma(D_{\cdot 0}, D_{\cdot 1})}{\sigma_{D_{\cdot 0}} \sigma_{D_{\cdot 1}}} & \dots & \frac{\sigma(D_{\cdot 0}, D_{\cdot n-1})}{\sigma_{D_{\cdot 0}} \sigma_{D_{\cdot n-1}}} \\ \frac{\sigma(D_{\cdot 1}, D_{\cdot 0})}{\sigma_{D_{\cdot 1}} \sigma_{D_{\cdot 0}}} & \frac{\sigma(D_{\cdot 1}, D_{\cdot 1})}{\sigma_{D_{\cdot 1}}^2} & \dots & \frac{\sigma(D_{\cdot 1}, D_{\cdot n-1})}{\sigma_{D_{\cdot 1}} \sigma_{D_{\cdot n-1}}} \\ \vdots & & & \vdots \\ \frac{\sigma(D_{\cdot n-1}, D_{\cdot 0})}{\sigma_{D_{\cdot n-1}} \sigma_{D_{\cdot 0}}} & \frac{\sigma(D_{\cdot n-1}, D_{\cdot 1})}{\sigma_{D_{\cdot n-1}} \sigma_{D_{\cdot 1}}} & \dots & \frac{\sigma(D_{\cdot n-1}, D_{\cdot n-1})}{\sigma_{D_{\cdot n-1}}^2} \end{bmatrix}$$

3.3 Numerical Analysis

3.3.1 Eigenvalues and Eigenvectors of a Matrix

Given a matrix $A \neq 0 \in \mathbb{R}^{2n}$, a vector $v \in \mathbb{R}^n$ is said to be an **eigenvector** of A if the multiplication Av does not change the direction of v ; that is:

$$\exists \lambda \in \mathbb{R} \mid Av = \lambda v, \text{ where}$$

λ is the **eigenvalue** associated to the eigenvector v .

3.3.2 Spectral Decomposition of a Matrix

If $[M]_{n \times n}$ is a symmetric matrix of rank n and admits n pairs of eigenvalues λ and eigenvectors $[V]_{n \times n} = [v_0, v_1, v_2, \dots, v_{n-1}]$, such that $V^T V = 1$, then

$$MV = V\lambda$$

$\lambda = diag(\sigma_0, \sigma_1, \dots, \sigma_{n-1})$ is the diagonal matrix, where σ_i is the eigenvalue associated to the eigenvector v_i . [5] Furthermore, V 's columns are linear independent, hence V is invertible.

$$\begin{aligned} MV &= V\lambda \\ MVV^T &= V\lambda V^T \\ M &= V\lambda V^T \end{aligned}$$

3.3.3 Singular Value Decomposition

If $M \in \mathbb{R}^{m \times n}$, then $\exists U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ and $\Sigma = \text{diag}(\sigma_0, \dots, \sigma_{n-1})$ conditioned to $\sigma_i \geq \sigma_{i+1} \geq 0, \forall \sigma \in [0, n)$ s.t. [6]

$$M = U\Sigma V^T$$

Theorem 3.1 If $A = U\Sigma V^T$, $AA^T = U\Sigma^2 U$ and $A^T A = V\Sigma^2 V$.

Proof

$$\begin{aligned} A^T A &= (U\Sigma V^T)^T (U\Sigma V^T) \\ &= V\Sigma^T U^T U\Sigma V^T \\ &= V\Sigma\Sigma V^T \\ &= V\Sigma^2 V^T \end{aligned}$$

Proving $AA^T = U\Sigma^2 U$ is analogous to the above.

3.4 Topology

3.4.1 Manifolds

Intuitively, n-dimensional topological manifolds are sets that are “locally Euclidean” [7]. In other words, they can be decomposed into sub sets that can be mapped to the \mathbb{R}^n .

Formally, a set M is said to be a **n-dimensional topological manifold** $\iff M$ is a paracompact Hausdorff topological space $| \forall p \in M, p \in U_p$, where U_p is an open set that is homeomorphic to an open set V_p of the Euclidean space \mathbb{R}^n [7].

From here, we will use the word manifold to refer to the n-dimensional topological manifold.

Charts

The pair (U_i, ϕ_i) is called a **coordinate chart** or **chart** on M if $U \in M$ and ϕ_i is a **homeomorphism** such that $\phi_i(U_i) = V_i \subseteq \mathbb{R}^n$ [?].

Atlas

A set $A = \{(U_i, \phi_i)\}_{i \in A}$ is said to be an **atlas** on a manifold M if $\cup_{i \in A} U_i = M$ [?].

Example 3.3 *The \mathbb{R}^n is, directly, a manifold.*

Example 3.4 *A n -dimensional sphere is a manifold. Earth, in special, is a 3-dimensional sphere and its stereographic projection is its mapping to the \mathbb{R}^2 [8].*

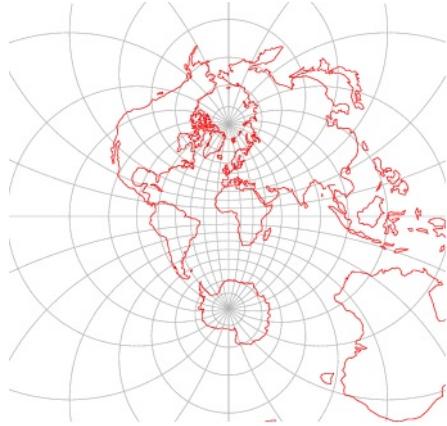


Figure 2: Stereographic projection applied to Earth.

3.5 Graph Theory

3.5.1 Graphs

Let G be the pair (X, U) . G is defined as a **graph** [9] where

1. X is a set of objects called **vertices**.
2. U is a family of elements $u_i \in X \times X$ called arcs.

Basic Concepts [9]

Multiplicity If $G = (X, U)$ and $(x, y) \in X \times X$, the multiplicity $m_g^+(x, y)$ of x, y is defined to be the number of arcs with initial endpoint x and terminal endpoint y . Furthermore:

1. $m_g^-(x, y) = m_g^+(y, x)$

$$2. \ m_G(x, y) = m_G^+(x, y) + m_G^-(x, y)$$

Degree If $G = (X, U)$ is a graph and $x \in X$, the degree $d(x)$ of x is defined [10] as $d(x) = 2n_s + n_n$, where n_s is the number of arcs self-incident at x (i.e., $\{(x, x)\}$) and n_n is the number of arcs incident at x .

Adjacency Matrix If $G = (X, U)$, $X = \{x_1, x_2, \dots, x_n\}$, define the adjacency matrix $A = [a_{ij}]_{n \times n}$ associated with graph G , where $a_{ij} = m_G^+(x_i, x_j)$.

Further Specifications

Undirected graph “All graphs are directed, but sometimes the direction need not be specified [9]”.

Let $G = (X, U)$ be a graph and $u_k \in U \mid u_k = (a, b)$. The element $e_i = [a, b]$ can be defined as the **edge** that links a to b without specifying direction. Finally, define the **undirected graph** H as (X, E) , where $E = \{e_i\}$ is the set of edges created from U .

Figure 3 shows the graph **Les Miserables**, where each node is a character and each arc links two characters that have shared stage at some point during the play.

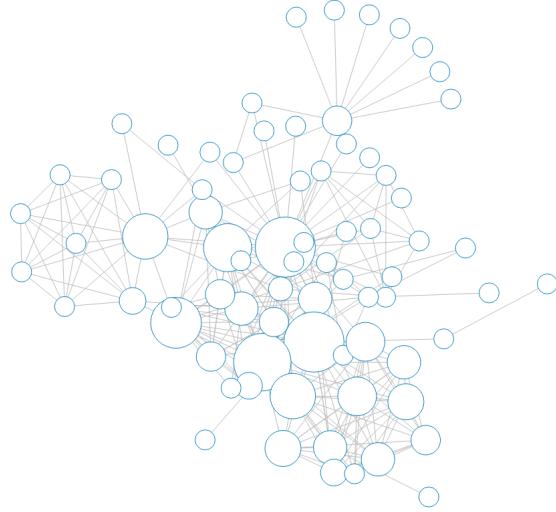


Figure 3: The **Les Miserables** graph.

Complete graph A graph $G = (X, U)$ is said to be complete if

$$\forall(x, y) \in X, x \neq y, m_G(x, y) \geq 1$$

Remark 3.3 Let $G = (X, U), G_1 = (X, E)$ and $n = |X|$. G is called the complete graph K_n if

$$\forall(x, y) \in X, \exists e \in E \mid e = [x, y]$$

Weighted graph Let $G = (X, U)$ be a graph and $w: U \rightarrow \mathbb{R} \mid w(u) = w_u$ be the weighted associated with arc u . G is said to be a weighted graph.

Euclidean graph If $G = (X, U)$ and $W = \{w_u, \forall u \in U\} \subset \mathbb{R}$, G is said to be an euclidean graph if w_u corresponds to the euclidean distance between the vertices connected by u in a specified embedding.

Tree Let G be the graph (V, E) such that G is connected and admits no cycles. G is said to be a **tree** [9].

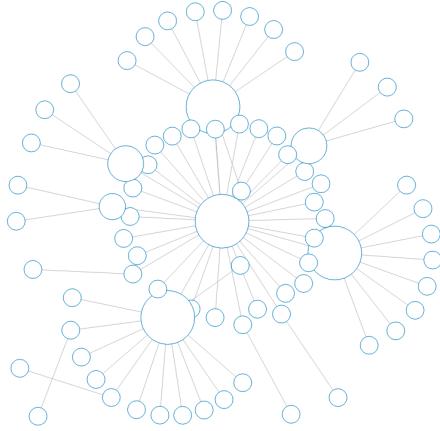


Figure 4: A tree extracted (a subgraph) from the **Les Misérables** graph.

3.5.2 Related Problems

Nearest-Neighbor Search

Let $G = (X, U)$ be a weighted graph, where $\forall u \in U, \exists w_u \in \mathbb{R}$, i.e., the weight or **length** of the arc u , and $n: U \rightarrow \{0, 1\}$ a definition of **nearness** in G . The nearest-neighbor search is a optimization problem that consists of finding a subgraph $H = (X, F \subseteq U) \mid f \in F \iff n(f) = 1$.

A Generic Algorithm for NN Search

1. $F_x := \emptyset, \forall x \in X$
2. $F_x := F_x \bigcup_{\forall p \in U_x} \begin{cases} \{p\}, & \text{if } n(p) = 1 \\ \emptyset, & \text{if } n(p) = 0 \end{cases}$
3. $H := (X, F), F = \bigcup_{x \in X} F_x$

Let's consider two (between many) specifications of this algorithm:

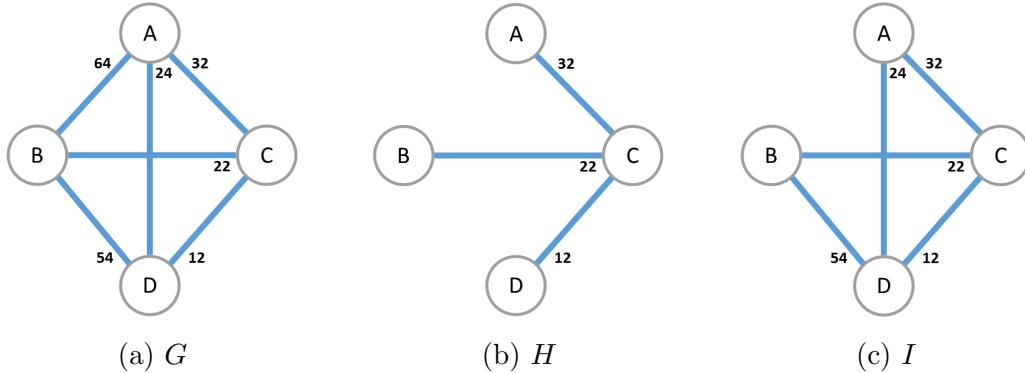
K-Nearest Neighbor Search (K-NN) Fixed $k \in \mathbb{N}$ and $U_x \subseteq U$, where U_x is the set of all arcs with x as initial endpoint, define:

$$n(u_x) := \begin{cases} 1, & \text{if } w_{u_x} \leq w_p, \forall p \in U_x - F_x \text{ and } |F_x| \leq k \\ 0, & \text{otherwise.} \end{cases}$$

ϵ -Nearest Neighbor (ϵ -NN) Fixed $\epsilon \in \mathbb{R}$, define:

$$n(f) := \begin{cases} 1, & \text{if } w_f \leq \epsilon \\ 0, & \text{otherwise} \end{cases}$$

Example 3.5 If $k = 1$ and $\epsilon = 60$, the graph G , the sub-graph H found from K-Nearest neighbor algorithm and the sub-graph I found from the ϵ -Nearest neighbor are defined as follows:



Shortest-path Problem [11]

Let $G = (X, U)$ be a weighted graph, a path $p(x, y) = (u_0, u_1, u_2, \dots, u_{p-2}, u_{p-1}, u_p)$ | $u_0 = (x, -), u_p = (-, y), u_i \in U, \forall i \in [0, p]$, and the weight of p be given by $w(p) = \sum_i^p w(u_i)$. The shortest-path weight between two vertices x and y is defined as:

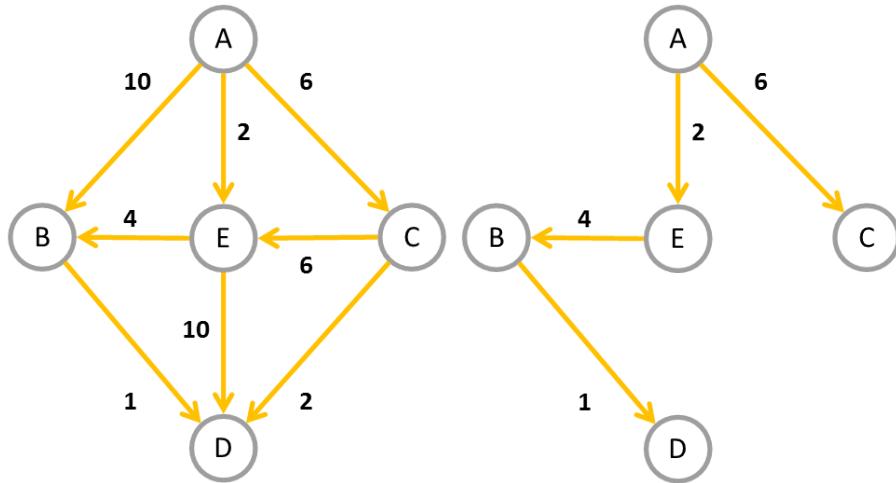
$$\sigma(x, y) = \begin{cases} \min\{w(p(x, y))\}, & \text{if } p(x, y) \text{ exists,} \\ \infty, & \text{otherwise.} \end{cases}$$

while $p(x, y)$ is the **shortest path** from x to y .

The shortest-path between a vertex $x_0 \in X$ and all other vertices can be expressed by the tree $S = (X, F), F \subseteq U$, denominated **shortest-path tree**.

Dijkstra's Algorithm [11]

Example 3.6 Let G be the weighted graph as defined in figure 6a. The shortest-path between A and all other vertices is described by the tree S illustrated in figure 6b.



(a) The weighted graph G .

(b) The shortest-path tree S .

3.6 Machine Learning

“Learning is the improvement of performance in some environment through the acquisition of knowledge resulting from experience in that environment.” [12]

“Machine learning is the area in Computer Science that has as goal the projection and implementation of machines that have the ability to learn autonomously” [13]. In other words, the creation of machines that are able to recognize patterns in an environment and interpret those using concepts related to artificial intelligence. Such interpretation can create a model which might eventually be used to predict new patterns, take actions and/or solve domain problems.

3.6.1 Machine Learning Algorithms

Based on how the learning phase of a problem is, most of the ML algorithms may be divided into one of the following categories:

Supervised A direct feedback is presented during the learning phase [12]. For the instances where the ME problem relies on a data set, the learning task uses the labeled training samples (i.e., samples that present the **target feature**) to synthesize the model that attempts to generalize the relationship between the feature vectors and the target variable [14].

Unsupervised Infer hidden structures from the data set without direct feedback, such as known labels for the samples [14].

Semi-supervised Most commonly, it is given by the extension of either supervised or unsupervised learning to include the other paradigm, resulting in a combination of both [15].

Reinforcement Through iterative exploration, the learner is positively or negatively reinforced for its actions. The learner’s goal is, ultimately, maximize the cumulative reward gained [14].

Support Vector Machine: An Example of Supervised Learning

The Support Vector Machine algorithm (i.e., SVM) is a powerful tool, often used in classification problems.

Intuitively, the SVM algorithm attempts to find a hyperplane that separates the samples in a data set into two different groups: positives and negatives. Furthermore, the hyperplane is placed such that the distance between the **support vectors** (the closest samples) and it are maximized.

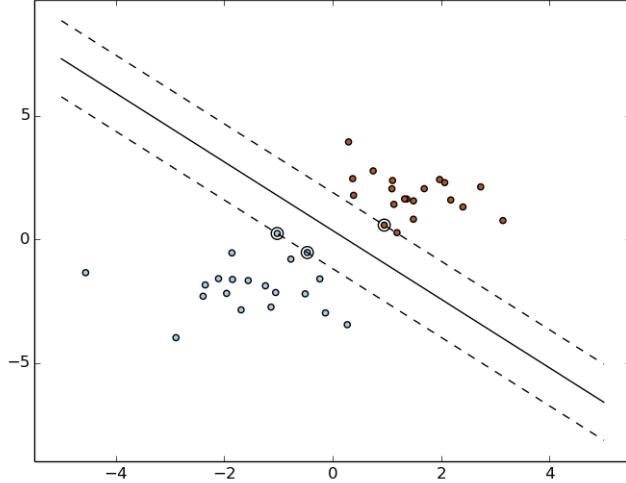


Figure 7: A SVM classifier projecting a hyperplane that perfectly separates two classes of samples [16].

More elaborately, given any linearly separable data set X containing samples from two distinguished classes $\{-1, +1\}$, consider the vector w a hyperplane $d = \{x \mid w \cdot x + b = 0\}$ (represented in fig. 7 by the contiguous line) s.t.

$$\text{decision rule} \begin{cases} w \cdot u + b \geq 0 \implies y_u = +1 \\ w \cdot u + b < 0 \implies y_u = -1 \end{cases}$$

To prevent samples from falling into the margin or being misclassified [17], reinforce that for any positive sample x_+ , $w \cdot x_+ + b \geq 1$. Similarly for x_- samples, $w \cdot x_- + b \leq -1$. These both constraints can be expressed as

$$y_i(w \cdot x_i + b) - 1 \geq 0 \quad (1)$$

Notice that $y_i(w \cdot x_i + b) - 1 = 0 \iff x_i$ is a **support vector**.

The width (the distance between the two margins) of the street is the vector $(x_+^0 - x_-^0)$ projected onto the vector w , where x_-^0 is a positive support vector and

x_+^0 is a negative one.

$$\begin{aligned}
width &= (x_+^0 - x_-^0) \cdot \frac{w}{\|w\|} \\
&= \frac{x_+^0 \cdot w - x_-^0 \cdot w}{\|w\|} \\
&= \frac{1 - b - (-1 - b)}{\|w\|} = \frac{2}{\|w\|}
\end{aligned} \tag{2}$$

As the goal is to maximize the width, while still respecting the constraint (1).

$$\max width = \max \frac{2}{\|w\|} \equiv \min \frac{1}{2} \|w\|^2 \tag{3}$$

Which can be solved using standard quadratic programming.

SVM for non-separable data sets (soft margins)

To deal with non-separable data sets, it is possible to introduce the variables ξ_i [17], which represent a trade-off between maximum margin/distance of the misclassified samples from the decision boundary.

$$\begin{aligned}
&\min \frac{1}{2} \|w\|^2 + C \sum_1^m \xi_i, \text{ constrained to:} \\
&y_i(w \cdot x_i + b) \geq 1 - \xi_i, \xi_i \geq 0
\end{aligned}$$

Such trade-off can be adjusted through the parameter C . Notice that small values for C might result in misclassification, whereas high values can produce overfitting.

Dependency over the dot product

By the Lagrange multipliers method [18]:

$$L = \frac{1}{2} \|w\|^2 - \sum \alpha_i [y_i(w \cdot x_i + b) - 1] \tag{4}$$

$$\frac{\partial L}{\partial w} = w - \sum \alpha_i y_i x_i = 0 \implies w = \sum \alpha_i y_i x_i \tag{5}$$

$$\frac{\partial L}{\partial b} = - \sum \alpha_i y_i = 0 \implies \sum \alpha_i y_i = 0 \tag{6}$$

Applying (5) and (6) on (4), our problem of minimizing (3) constrained by (1) becomes maximizing L , subject to (6):

$$\begin{aligned} L &= \frac{1}{2} \sum \alpha_i y_i x_i \cdot \sum \alpha_j y_j x_j - \sum \alpha_i y_i x_i \cdot \sum \alpha_j y_j x_j - \sum \alpha_i y_i b + \sum \alpha_i \\ &= \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) \end{aligned}$$

w and b can easily be found from (5) and $\alpha_i[y_i(w \cdot x_i + b) - 1] = 0$ (for any $i \mid \alpha_i \neq 0$), respectively.

Now, as we finally plug (5) back into our decision rule, it becomes clear that both training and prediction phases depend only on the dot product between the sample vectors [18]:

$$\text{decision rule } \begin{cases} \sum \alpha_i y_i x_i \cdot u + b \geq 0 \implies y_u = +1 \\ \sum \alpha_i y_i x_i \cdot u + b < 0 \implies y_u = -1 \end{cases}$$

The implications of such fact will be discussed in the next section.

Kernel functions

”The kernel function represent the dot product of both vectors projected onto the new space” [18].

Some data sets are not linearly separable, as mentioned in a few sections above. They can, however, be projected to a different vector space, where they hopefully will be. To achieve this, a transformation ϕ is applied on both vectors. The dot product between those is then calculated in the new space and the value is used during training (the maximizing L) and prediction (the decision rule). Practically, let $k(u, v) := f : (\mathbb{R}^n, \mathbb{R}^n) \rightarrow \mathbb{R} \mid k(u, v) = \phi(u) \cdot \phi(v)$,

$$\begin{aligned} L &= \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j k(x_i, x_j) \\ \text{decision rule } &\begin{cases} \sum \alpha_i y_i k(x_i, u) + b \geq 0 \implies y_u = +1 \\ \sum \alpha_i y_i k(x_i, u) + b < 0 \implies y_u = -1 \end{cases} \end{aligned}$$

The figure bellow illustrates a non-linearly separable data set defined in the \mathbb{R} being projected to the \mathbb{R}^2 .

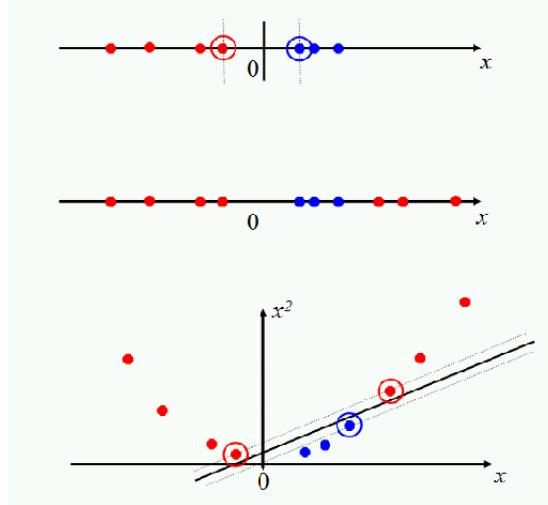


Figure 8: Projection of samples from the \mathbb{R} to the \mathbb{R}^2 , allowing SVM to find a hyperplane that perfectly separates both classes [19].

Between many different kernels, two often used are the RBF: $\exp(-\frac{\|x-x'\|^2}{2\sigma^2})$ and the Polynomial: $(u^T v + c)^d$. [19]

3.6.2 Multi-class Classification

Not all problems are binary. In classification, this reflects data sets that have their samples separated into more than two classes. The Iris flower data set is an example of this, when predicting the samples' *species*.

There are many different approaches for multi-class classification. [20] For the scope of this work, consider only the following:

One-vs-All (OVA) n binary classifiers are built, where n is also the number of classes in the data set. For the classifier $c^j, j \in [1, n]$, all samples of the j th class are taken as positive examples, at the same time that all the other samples are considered negative.

If

$$c^j(x) = \sum_{i=1}^m y_i \alpha_i^j x \cdot x_i + b$$

Then positive values for $c^j(x)$ indicate that the sample x belongs to the j th class. Additionally, greater $c^j(x)$ values imply on further distance from the

hyperplane (i.e., $c^j(x)$ can also be interpreted as a **confidence value**) and the sample x should be assigned to the class which holds greatest confidence. [21] Shortly, classification is given by

$$f(x) = \arg \max_j c^j(x)$$

All-vs-All (AVA) Also known as all-pairs or one-vs-one, a classifier c_{ij} is built for each pair of classes (i, j) , resulting in a total of $n(n - 1)$ classifiers. c_{ij} responds with positive values for samples of the i th class and negative values for samples of the j th class. Classification can be done by simply counting the class most frequently associated with x :

$$f(x) = \arg \max_i \sum_{j=1}^n c_{ij}(x)$$

Or equivalently, but only using $\frac{n(n-1)}{2}$ classifiers,

$$f(x) = \arg \max_i \sum_{j=1}^n \frac{j - i}{|j - i|} c_{\min(i,j) \max(i,j)}(x)$$

3.6.3 Evaluating learners

Machine Learning algorithms might be susceptible to data noise, incorrect configuration or even random factors, which would eventually decrease the generated model's accuracy. In order to evaluate this same accuracy, models are quite often tested after trained.

Considering that testing with the same data used for training will most likely produce unreliable results, a simple way to test a learner is to separate the labeled data set into two chunks, where the first is used for training. The second chunk is then given to the learner, which attempts to predict the samples. Finally, the predictions made by the learner would be compared with the actual labels.

Confusion Matrix

When testing classification models, one way to visualize the wrong predictions made by the learner is a confusion matrix, where the item a_{ij} is the number of times that a sample of the class i was classified as being of the class j .

	a	b	c	d
a	12	3	2	0
b	7	12	2	4
c	0	4	54	8
d	6	0	1	23

Table 4: Example of confusion matrix for a data-set with four different classes.

Naturally, a diagonal matrix represents that all samples of the class i were classified as i , which is the best possible outcome (no errors).

Cross Validation

Sometimes (for example, when the data set does not have too many samples), partitioning of the data set into subsets might be benign to the learning process, as the model will be constructed only considering a small, random portion of the samples. This event is known as **underfitting**.

When k -fold cross-validating, [22] the data set can be partitioned into k folds. For each fold k_i , a model is trained with all folds, except for k_i . The model is then tested over k_i . Finally, the score reported by the cross-validation method is the average accuracy when testing over all folds.

Grid Search

Machine Learning algorithms may require specific parameters to run. For instance, SVM requires C and which *kernel* it should use. Additionally, kernels might require their own parameters. As these parameters strongly affect how the generated model will be, one cannot choose them arbitrarily.

Grid Search is a traditional method used to find the parameters that optimize the generalization of learning model over a specific data set. [23] Given a set of possible parameters, it will exhaustively search for the combination of those that generate the best outcome, which might be evaluated by a validation method such as cross-validation.

3.6.4 Examples of learning

Coffee Selling Rate

The figure below illustrates a data set which express the relationship between the variables **time of day** and the **selling rate** of coffee in a particular coffee-shop. Clearly being a regression problem, a machine learning algorithm must create a model that appropriately generalizes the distribution observed. Such model will eventually be used to predict the selling rate in the following days.

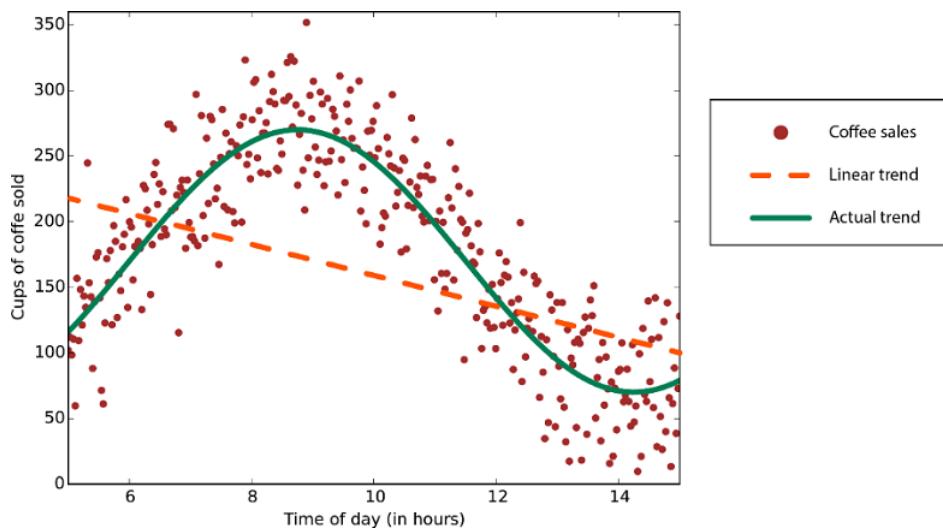


Figure 9: graphic representation of a data set generalization by a linear (orange) and a non-linear model (green)

The orange line and the green arc represent two different models. The orange line, which represents the linear model, clearly does not generalize the data set appropriately, once it induces an error much larger than necessary. [24] The non-linear model, i.e., the green arc, was capable of generalizing the data inducing a smaller error.

Iris Flower

Consider the Iris flower data set as defined in 3.1.1. In order to predict the feature *Species* of a given sample, one could train a classification model using the **Support Vector Machine** algorithm.

Through GridSearch, it was found that the SVM algorithm with $C = 100$, $gamma = .01$ and rbf kernel is capable of finding a model yielding .99% accuracy. The samples misclassified during the test phase were represented by the confusion matrix bellow.

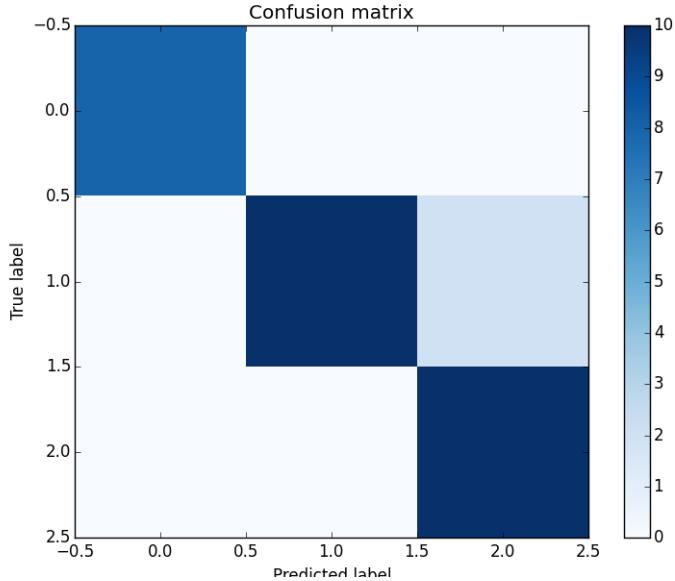


Figure 10: Confusion matrix of a SVM with $C = 100$, $gamma = .01$ and rbf kernel when predicting samples from the Iris flower data set.

4 Linear Dimensionality Reduction

As presented in the previous sections, data sets with many features may present a series of issues: difficult visualization, high performance requirements, noise etc. In this section, it will be discussed methods related with linear dimensionality reduction, i.e., the shrinking of data sets by transformation and/or removal of features, while minimizing information loss.

Consider the data set K . K has its samples expressed by two similarly scaled dimensions. It is clear, however, that the samples follow a very particular distribution:

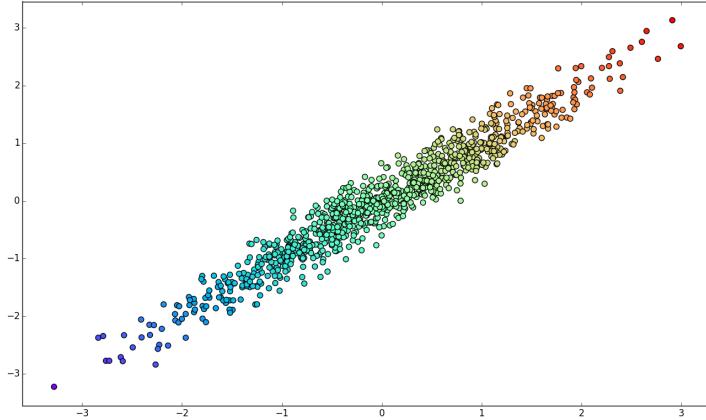


Figure 11: The data set $K \in \mathbb{R}^2$.

Additionally, something interesting can be observed when analyzing the covariance matrix of K : as it is not a diagonal matrix, the variance of x from its mean somehow correlates with the variance of y . [25]

	x	y
x	1.26682132	1.29158697
y	1.29158697	1.40358478

Table 5: Covariance between the components of K .

4.1 Principal Component Analysis

As in K , some data sets follow certain distributions that are majorly contained in a few orthogonal components, where a component is the result of a linear combination of the original features.

Principal Component Analysis (PCA) is a statistical technique that attempts to transform a n -dimensional data set X into a m -dimensional data set Y , where, hopefully, $k \ll n$. Furthermore, the dimensions of Y will necessarily be orthogonal components aligned with the direction in which the variance of samples in X is maximum, commonly referred to as **principal components**. [26] In figure 12, the orange and purple vectors are the principal components of the data set K .

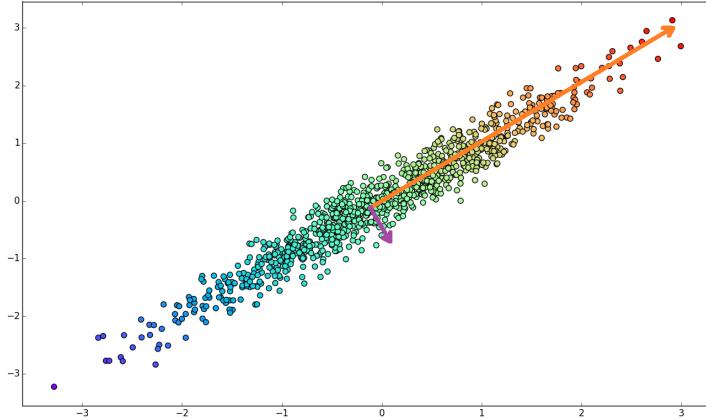


Figure 12: The principal components of K .

4.1.1 Study of the PCA Algorithm

Let D be a dataset with n samples and f features and $X = HD$, where H is the centering matrix. Our goal is to find which are the principal components of the covariance matrix Σ_X .

$$\Sigma_X = \frac{1}{n} X^T X \quad (7)$$

Using the **Singular Value Decomposition** method described in section 3.3.3, we know that

$$X = U\Sigma V^T \quad (8)$$

Needless to say, Σ is the diagonal matrix of singular values, not to be mistaken by the covariance matrix Σ_X .

From 7 and 8:

$$\begin{aligned} \Sigma_X &= \frac{1}{n} X^T X \\ &= \frac{1}{n} (U\Sigma V^T)^T U\Sigma V^T \\ &= \frac{1}{n} V\Sigma^2 V^T \end{aligned}$$

Which entails that V is the orthonormal matrix with Σ_X 's eigenvectors as columns, whereas Σ contains the correspondent eigenvalues σ_{ii}^2 associated with

$v_i \in V$. $v_i \in V$ is, in fact, a principal component of X and its associated eigenvalue σ_i module gives X spectral radius. As we are interested in the dimensions that give most variance, keep only the $m \in \mathbb{R}$ most significant eigenvalues and their correspondent eigenvectors.

Finally, it also worth remarking once again that the principal components are linear combinations of the original features (the canonical base). V^{-1} is, therefore, a change-of-basis matrix from the original feature basis (the canonical base) to a new one generated by the principal components. As V is orthogonal, V^{-1} exists and it is equal to V^T . Formally, if x is a sample from the X data set,

$$y = V^T x$$

4.1.2 Formalization of the PCA Algorithm

Let D be a data set with n samples and f features and $m \in \mathbb{R}$ the number of dimensions desired for the reduced data set. [27] [28]

1. Find $X = HD$, where H is the centering matrix.
2. Calculate the covariance matrix Σ_X .
3. Use singular value decomposition to find the eigenvalues $\lambda = \{\lambda_i\}$ and eigenvectors $V = \{v_i\}$ of Σ_X .
4. Sort the eigenvalues by their absolute value in descending order and select the first m ones and their respective eigenvectors.

4.2 Multidimensional Scaling

Alternatively to PCA, Multidimensional Scaling (or simply MDS) can be used to reduce the dimensionality of a data set. The method has, however, an extensive application domain and often appears in the literature in different contexts. An example of this is the problem of, given a set of objects O and a dissimilarity measurement $\delta_{rs}, \forall(r, s) \in O \times O$, finding a suitable representation in the \mathbb{R}^n for the objects in O . [5]

For this project, we study the **classic MDS**. That is, when the dissimilarities considered are the euclidean distances between coordinates in the \mathbb{R}^n .

4.2.1 Study of the MDS

If $\delta = [\delta_{rs}]_{n \times n}$ is the dissimilarity matrix, where δ_{rs} represents the euclidean distances between two samples $x_r, x_s \in \mathbb{R}^m$ from the data set $[X]_{n \times m}$ induced by the L2-norm. In other words,

$$\begin{aligned}\delta_{rs} &= \sqrt{\sum_i (x_{ri} - x_{si})^2} \\ &\iff \\ \delta_{rs}^2 &= \sum_i (x_{ri} - x_{si})^2 \\ &= (x_r - x_s) \cdot (x_r - x_s) \\ &= x_r \cdot x_r + x_s \cdot x_s - 2x_r \cdot x_s\end{aligned}\tag{9}$$

Now consider the inner product matrix $B = XX^T$, where $b_{rs} = x_r \cdot x_s$. Given that B can be decomposed as $U\Sigma U^T = U\Sigma^{\frac{1}{2}}\Sigma^{\frac{1}{2}}U^T = U\Sigma^{\frac{1}{2}}(U\Sigma^{\frac{1}{2}})^T = XX^T \iff X = U\Sigma^{\frac{1}{2}}$, if we can derive B from 9, it will be possible to apply the same decomposition considered in PCA to find a new data set Y which is a reduction of X . [5]

Firstly, we will assume that Y is centered in the origin (i.e., Y has its features' means equal to zero):

$$\sigma_f = \sum_i y_{if} = 0, \forall f \in [0, m)\tag{10}$$

Now, 9 \implies

$$\begin{aligned}\frac{1}{n} \sum_r \delta_{rs}^2 &= \frac{1}{n} \sum_r (x_r \cdot x_r + x_s \cdot x_s - 2x_r \cdot x_s) \\ &= \frac{1}{n} \sum_r x_r \cdot x_r + \sum_r x_s \cdot x_s - 2 \sum_r x_r \cdot x_s \\ &= \frac{1}{n} \sum_r x_r \cdot x_r + nx_s \cdot x_s - 2 \sum_r 0 \cdot x_s \\ &= \frac{1}{n} \sum_r x_r \cdot x_r + x_s \cdot x_s \\ &\iff \\ x_s \cdot x_s &= \frac{1}{n} \left(\sum_r \delta_{rs}^2 - \sum_r x_r \cdot x_r \right)\end{aligned}\tag{11}$$

Similarly to 11,

$$x_r \cdot x_r = \frac{1}{n} \left(\sum_s \delta_{rs}^2 - \sum_s x_s \cdot x_s \right) \quad (12)$$

Putting 11 and 12 back in 9:

$$\begin{aligned} \delta_{rs}^2 &= \frac{1}{n} \left(\sum_s \delta_{rs}^2 - \sum_s x_s \cdot x_s + \sum_r \delta_{rs}^2 - \sum_r x_r \cdot x_r \right) - 2x_r \cdot x_s \\ &\implies \\ x_r \cdot x_s &= -\frac{1}{2} (\delta_{rs}^2 - \frac{1}{n} [\sum_s \delta_{rs}^2 - \sum_s x_s \cdot x_s + \sum_r \delta_{rs}^2 - \sum_r x_r \cdot x_r]) \\ &= -\frac{1}{2} (\delta_{rs}^2 - \frac{1}{n} [\sum_s \delta_{rs}^2 + \sum_r \delta_{rs}^2 - 2 \sum_r x_r \cdot x_r]) \end{aligned} \quad (13)$$

To eliminate the $x_r \cdot x_r$ term from 13:

$$\begin{aligned} \frac{1}{n^2} \sum_s \sum_r \delta_{rs}^2 &= \frac{1}{n^2} \sum_s \sum_r (x_r \cdot x_r + x_s \cdot x_s - 2x_r \cdot x_s) \\ &= \frac{1}{n^2} \sum_s (\sum_r x_r \cdot x_r + \sum_r x_s \cdot x_s - 2 \sum_r x_r \cdot x_s) \\ &= \frac{1}{n^2} \sum_s (\sum_r x_r \cdot x_r + n x_s \cdot x_s) \\ &= \frac{1}{n^2} (n \sum_r x_r \cdot x_r + n \sum_s x_s \cdot x_s) \\ &= \frac{1}{n^2} 2n \sum_r x_r \cdot x_r \\ &= \frac{2}{n} \sum_r x_r \cdot x_r \end{aligned} \quad (14)$$

Finally, applying 14 on 11:

$$B_{rs} = x_r \cdot x_s = -\frac{1}{2} (\delta_{rs}^2 - \frac{1}{n} [\sum_s \delta_{rs}^2 + \sum_r \delta_{rs}^2 - \frac{1}{n} \sum_s \sum_r \delta_{rs}^2]) \quad (15)$$

From 15, it becomes clear that B is, in fact, the double centering of the matrix $A = -\frac{1}{2}\delta^2$. I.e., $B = HAH$. Singular value decomposition can now be performed onto B , resulting in the matrices U and Σ .

Finally, we can sort the eigenvalues (and their respective eigenvectors, the columns of U) in decrease order and keep only the ones that offer greater variance.

Remark 4.1 As euclidean distances were used to build the dissimilarity matrix δ , B is indubitably positive semidefinite, hence $\Sigma_i \geq 0, \forall i \in [0, n]$. However, negative eigenvalues might appear if other dissimilarity measurement were to be used. In these cases, one might consider to simply ignore such components.

4.2.2 Formalization of the Multidimensional Scaling Method

Let X be a data set with n samples and f features and $m \in \mathbb{R}$ the number of dimensions desired for the reduced data set. [5]

1. Calculate the dissimilarity matrix $[\delta]_{rs}$, where $\delta_{rs} = \sqrt{\sum_i (x_{ri} - x_{si})^2}$
2. Calculate the matrix $B = H - \frac{1}{2}\delta_{rs}^2 H$.
3. Use singular value decomposition to find the matrices Σ and U .
4. Select the m greatest eigenvalues in Σ . From these, create the matrices $\Sigma' = [\sigma'_{m \times m}]$ and $U' = [u'_{n \times m}]$ (where each column i contains the eigenvector associated with σ'_i).
5. Construct $Y = U'\Sigma'$

4.3 Classification and Regression Over Linearly Reduced Data Sets

Consider that all tests bellow were done in following order:

1. The data set was retrieved.
2. Grid Search was executed, where the parameters grid was given by

kernel	linear	rbf	sigmoid		
C	1	10	100	1000	
gamma	.001	.01	.1	1	10

Table 6: Parameters Grid. Note: the *gamma* parameter is not used in the linear kernel, hence the search space does not contain multiple combinations of these two parameters.

3. Reduction was performed.
4. Grid Search was re-executed and the results compared.

4.3.1 K data set

The figure bellow illustrates the results of PCA algorithm application over the artificial K data set. Notice that, for the second application, it correctly chose to discard the vertical dimension, as the samples offer less variability in this component.

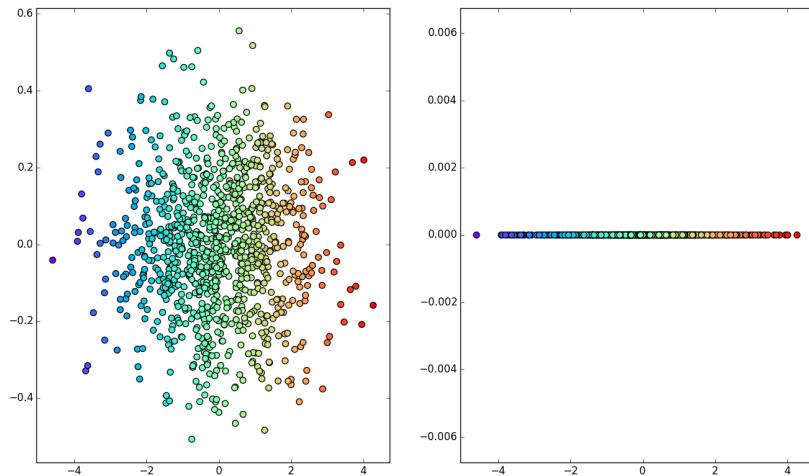


Figure 13: The PCA algorithm reducing K to 2 and 1 dimension, respectively.

	Original data	Reduced data (\mathbb{R}^2)	Reduced data (\mathbb{R})
Pred. accuracy	.97	.98	.99
GridSearch time	1.98s	2s	2.26s
Reduction time	—	0.995ms	1.118ms
Data size	1600 bytes	1600 bytes	8000 bytes

Table 7: Description of predictions and reduction performance for k .

4.3.2 The Iris flower data set

The Iris flower being reduced from 4 to 2 features. Although the data set became non-linearly separable, classes are still somewhat organized in different clusters.

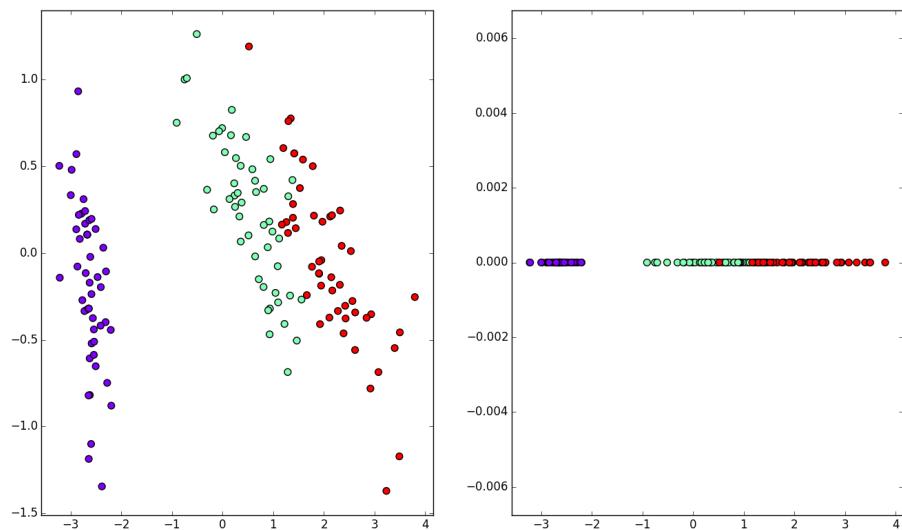


Figure 14: The PCA algorithm reducing the Iris flower to 2 and 1 dimension, respectively.

	Original data	Reduced data (\mathbb{R}^2)	Reduced data (\mathbb{R})
Pred. accuracy	.99	.97	.94
GridSearch time	1.71s	1.64s	1.78s
Reduction time	—	0.995ms	1.118ms
Data size	4800 bytes	2400 bytes	1200 bytes

Table 8: Description of predictions and reduction performance for Iris flower.

4.3.3 The Digits data set

Digits data set is composed by 1.797 samples, 64 features and 10 classes. Each sample is a 8x8 image of a hand-written digit from 0 to 9.

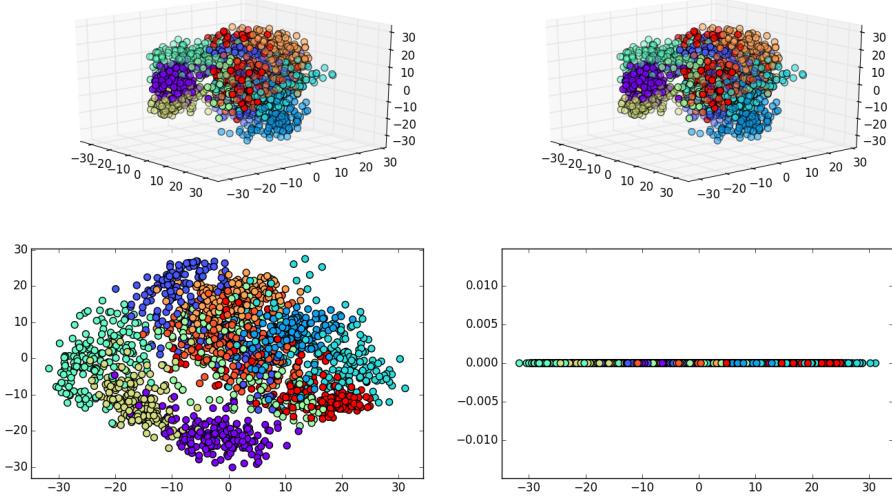


Figure 15: Digits data set reduced to 10, 3, 2 and 1 dimension, respectively.

	\mathbb{R}^{64}	\mathbb{R}^{10}	\mathbb{R}^3	\mathbb{R}^2	\mathbb{R}
Pred. accuracy	.98	.95	.74	.64	.39
GridSearch time	12.82s	30.25s	186.73s	154.56s	122.47s
Reduction time	—	0.02s	0.01s	0.01s	0.01s
Data size	920064 B	143760 B	43128 B	28752 B	14376 B

Table 9: Description of predictions and reduction performance for Digits.

Notice that it was possible to eliminate 54 dimensions, consistently reducing the data set size, and only suffering 3% of prediction accuracy loss. The score drastically decreased, however, when more dimensions were removed.

5 Non-linear Dimensionality Reduction

Although PCA has presented promising results in the previous section, hence its great popularity in dimensionality reduction problems, there are many examples in which PCA will fail in its task. Consider the example below:

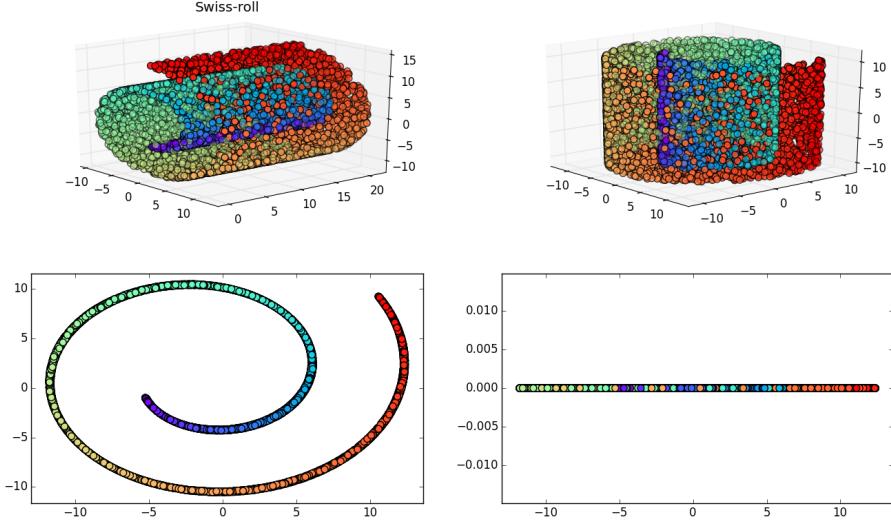


Figure 16: The **Swiss-roll manifold** and its reductions to 3, 2 and 1 dimensions, respectively, using the PCA algorithm.

Specially in the last reduction, to a single dimension, a big drawback of the algorithm PCA is clear: it assumes that the data lies roughly on a liner subspace, causing it to incorrect extract the underlying structure when this is not the case. [3]

The table bellow describes a regression attempt over the data sets illustrated in figure 16, were the feature being predicted is the contiguous value represented by the vertexes' colors:

	Original data	Reduced (\mathbb{R}^2)	Reduced (\mathbb{R}^1)
Accuracy	1.	.68	.54
GS time	3.84s	2.26s	3.26s
Reduction time	—	0.00ms	0.00ms
Data size	23.44 KB	15.62 KB	7.81 KB

Table 10: Regression accuracy and reduction performance for the Swiss-roll data set.

5.1 The Isomap Algorithm

Firstly suggested by Tenenbaum, de Silva and Langford, **Isometric Mapping** (or Isomap) assumes that the data lies near a smooth manifold. If the assumption is reasonable, it is possible to explore concepts such as neighborhood and local linearity to map the manifold to a linear structure before reducing it with a linear algorithm.

In this section we will discuss the Isomap algorithm and its inner workings. We will then proceed to formalize it. Finally, empirical tests results achieved during the project will be shown.

5.1.1 Study of the Isomap Algorithm

As the original data set might be folded, twisted or curved, [29] we must first find a suitable linear representation for it.

Let S be our original data set, as illustrated in figure 17.

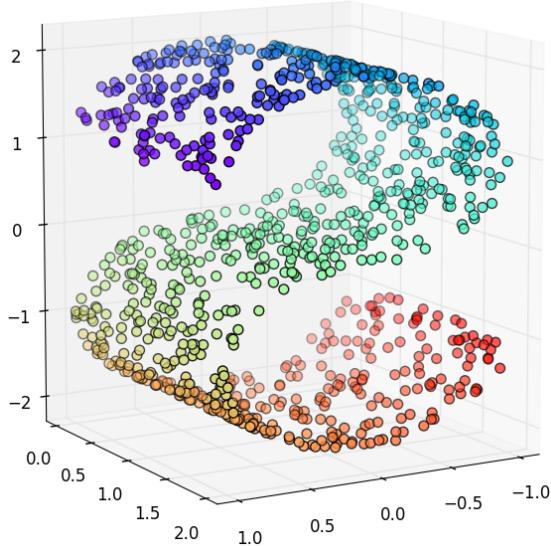


Figure 17: The data set S , consisting of 1000 samples and 3 features.

Additionally, consider the undirected weighted graph $G = (V, E)$. If $e_{xy} \in E$, $w(e_{xy}) = \delta_{xy}$, where δ_{xy} is the euclidean distance between the samples x and y .

in S . That is,

$$\delta_{xy} = \sqrt{\sum_i (x_i - y_i)^2}, \forall (x, y) \in S \times S \mid x \neq y$$

Now that only distances were kept, an infinite number of n-dimensional embeddings can be found with **MDS**, as every solution can be transposed, rotated or reflected. This does not fix the non-linearity of the data, though, as the original distances strictly constraint the samples to their original pattern. In order to achieve this, **Nearest neighbor search** can be performed over G , resulting in the graph H . Nearest-neighbor search will preserve edges connecting closer nodes, hence preserving local (linear) distances, but erase edges connecting nodes which are far from each other (non necessarily linear). Obviously, the search parameters (K or ϵ) must be carefully chosen to limit the connectivity of the nodes to a small (linear) neighborhood while maintaining the graph completely connected. Ideally, H will be a **mesh graph**.

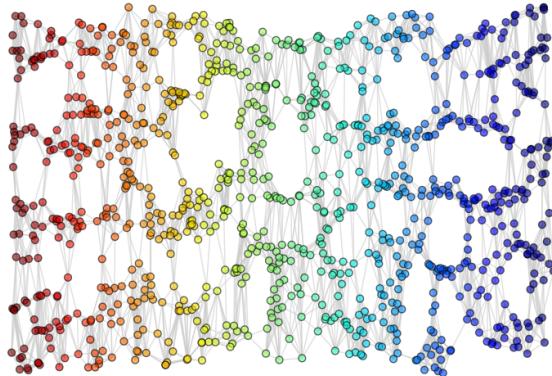
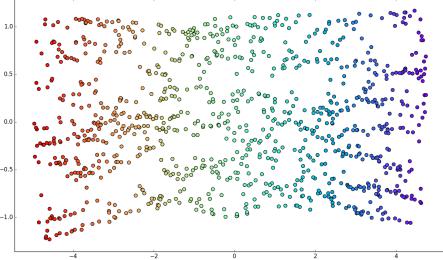


Figure 18: The graph H .

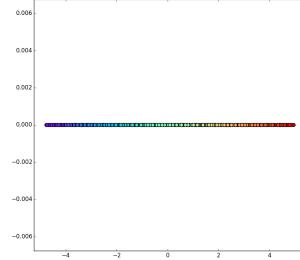
At this moment, not all distances in H are defined. This can be easily handled, though, by performing the **Floyd-Warshall** algorithm over H . Alternatively, M can be achieved by performing **Dijkstra's algorithm** for all nodes and joining all shortest-path trees found.

Finally, M is a euclidean graph which roughly lies on a hyperplane. Furthermore, the adjacency matrix associated to M contains not the distance induced by the L_2 norm, but the geodesic pairwise distances. [4] The **MDS** method can now

be used to construct a representation in sub-spaces of the \mathbb{R}^n . S , specifically, can be reduced to the \mathbb{R}^2 or \mathbb{R}^1 :



(a) S reduced to two dimensions.



(b) S reduced to one dimension.

5.1.2 Formalization of the Isomap Algorithm

Let X be the original data set and $p \in \mathbb{R}$ the number of dimensions desired for the reduced data set, [29]

1. Construct the weighted graph G from the distances pairwise $\delta_{xy}, \forall(x, y) \in S \times S, x \neq y$ and find the graph H by applying the **nearest-neighbor algorithm** on the graph G .
2. Compute the shortest path graph M between all pairs of nodes from graph H . This might be done by the **all-pairs Dijkstra's** or by the **Floyd-Warshall algorithm**.
3. Use M to construct the p -dimensional embedding using the **MDS** algorithm.

5.1.3 Computational Complexity

If n is the number of training samples, $f \in \mathbb{R}$ the number of features and $k \in \mathbb{R}$ the number of nearest neighbors:

1. The time complexity associated with building the neighborhood graph is $O(n^2)$.
2. For calculating the shortest path graph:

- (a) Dijkstra's algorithm implementation using Fibonacci Heaps have time complexity $O(nk + n \log n)$. As the algorithm must be calculated for each node, this step has time complexity $O[n^2(k + \log n)]$.
 - (b) Alternatively, using the Floyd-Warshall algorithm, this step has time complexity equals to $O(n^3)$.
3. MDS: taking $O(n^3)$ time steps to calculate the SVD decomposition, it is the bottleneck of the entire algorithm. [3]

The time complexity of the Isomap algorithm (when using Dijkstra's) is, therefore, $O[n^2 + n^2(k + \log n) + n^3]$. Better implementations exist, of course. For example, **Ball Tree** can be used for efficient neighbor search, requiring only $O(fn \log k \log n)$ time steps. [30]

5.2 Classification and Regression Over Data Sets Reduced with Isomap

5.2.1 The Swiss Roll Data Set

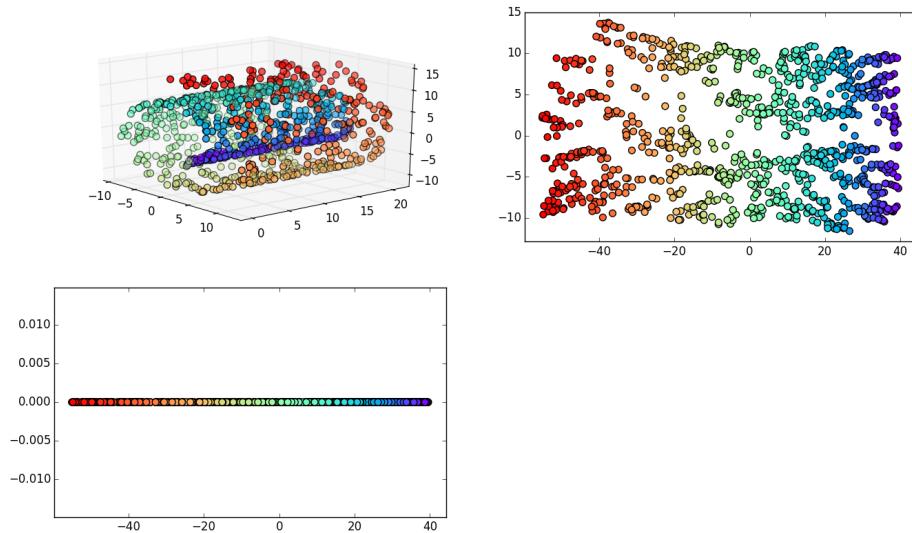


Figure 20: The Swiss Roll data set and its reductions to two and one dimensions, respectively.

	\mathbb{R}^3	\mathbb{R}^2	\mathbb{R}
Prediction accuracy	1	1	1
GridSearch time	15.53 s	415.72 s	387.98 s
Reduction time	—	0.51 s	0.48 s
Data size	23.44 KB	15.62 KB	7.81 KB

Table 11: Description of predictions and reduction performance for Swiss Roll.

5.2.2 The Digits Data Set

We will analyze Digits once again, only this time it will be reduced with the Isomap algorithm.

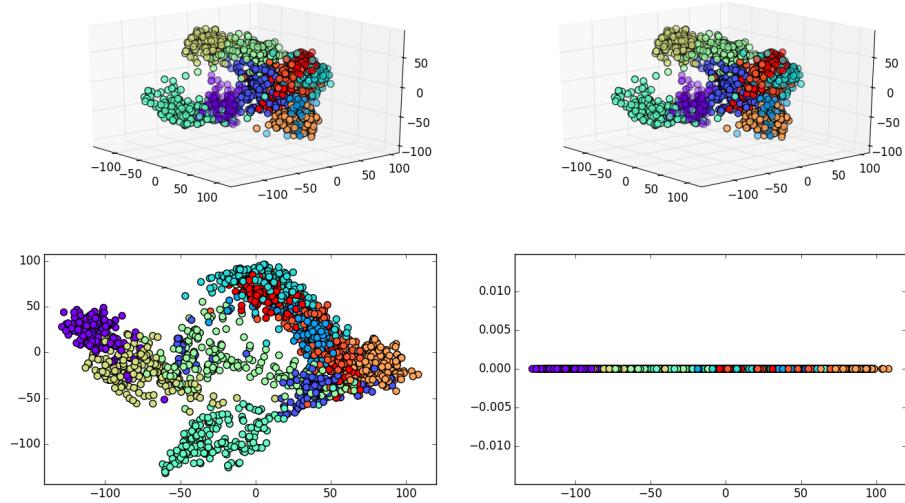


Figure 21: Digits data set reduced to 10, 3, 2 and 1 dimension, respectively.

	\mathbb{R}^{64}	\mathbb{R}^{10}	\mathbb{R}^3	\mathbb{R}^2	\mathbb{R}
Pred. accuracy	.98	.96	.91	.69	.45
GridSearch time	12.22 s	98.08 s	287.46 s	630.47 s	949.88 s
Reduction time	—	5.43 s	2.72 s	2.71 s	2.04 s
Data size	898.50 KB	140.39 KB	42.12 KB	28.08 KB	14.04 KB

Table 12: Description of predictions and reduction performance for Digits.

Notice that we managed to reduce the data set to only 3 dimensions while maintaining 91% of accuracy (remember that dimensionality reduction with PCA would reduce accuracy to 74%). Accuracy loss is still observable when reducing it to two or one dimensions, although it is less intense than losses caused by linear reduction.

5.2.3 The Leukemia Data Set

Extracted from mldata, the Leukemia data set contains 72 samples and 7130 features. The first 7129 features are expression levels of the genes in a given patient. The 7130th feature $t \in \{-1, 1\}$ indicate which of two variants of leukemia is present in the sample (AML, 25 samples, or ALL, 47 samples). [31]

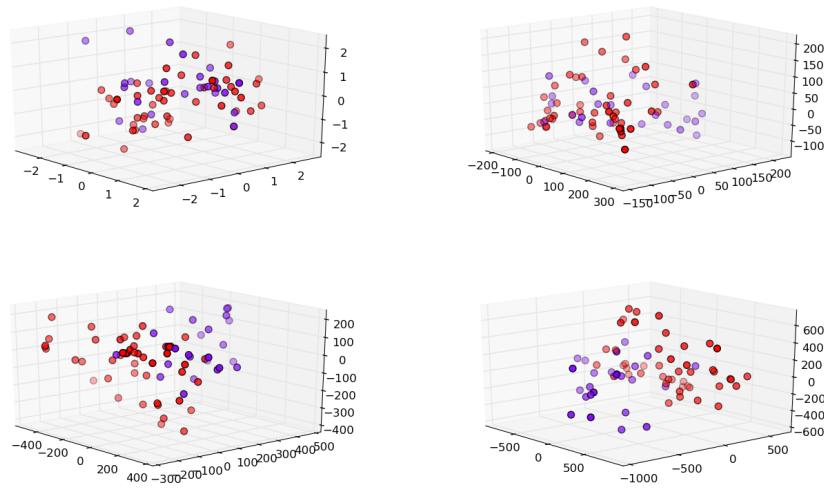


Figure 22: The Leukemia data set and its reduction to 30, 20 and 10 dimensions, respectively.

	\mathbb{R}^{7129}	\mathbb{R}^{30}	\mathbb{R}^{20}	\mathbb{R}^{10}
Pred. accuracy	.99	.88	.85	.88
GridSearch time	2.42 s	.26 s	.26 s	141.13 s
Reduction time	—	.06 s	.04 s	.04 s
Data size	4010.06 KB	16.88 KB	11.25 KB	5.62 KB

Table 13: Description of predictions and reduction performance for Leukemia data set.

5.3 Applicability and Limitations of Isomap

5.3.1 Infeasibility on Highly Dense Data Sets

In practice, Isomap's time complexity of $O[n^2 + n^2(k + \log n) + n^3]$ (see section 5.1.3) makes it unsuitable for data sets with great number of samples. This is illustrated in Shi and Gu's experiments: Isomap could not reduce data sets with more than 6000 samples in reasonable time. [32]

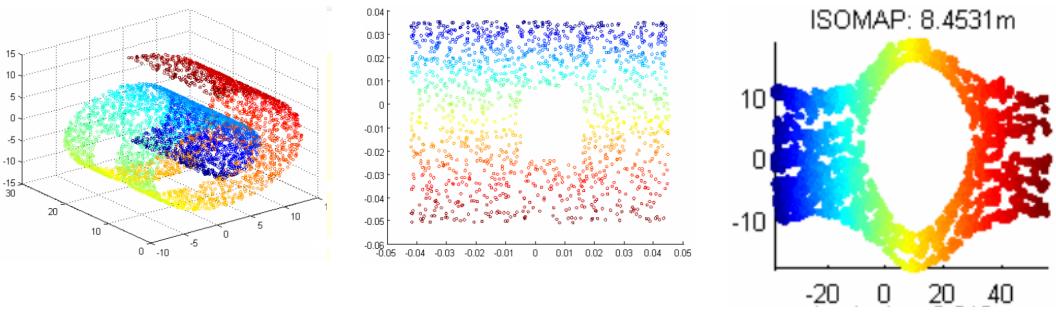
In order to visualize this issue, an experiment was made. Consider Spam the data set with 4601 samples and 57 features and the implementations MDS, Isomap, sk-PCA and sk-Isomap (the last two being native algorithms from the scikit-learn library).

5.3.2 Necessary Settings for Convergence

The experiments in the previous section made it quite clear that Isomap outperforms PCA on the selected data sets. Unfortunately, this trend cannot be projected for a generic case considering our experiments were strictly controlled. In real-world data sets, the application of Isomap may lead to poor low-dimensional embeddings. [33] Below are listed some the issues that have great influence over Isomap's results:

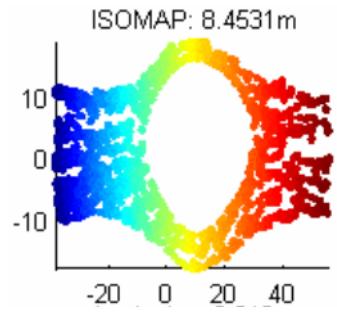
Manifold Assumption It refers to the initial assumption that the data lies on a low-dimensional manifold. Although wildly exploited by many authors, it is difficult to assert whether such assumption holds or not in real-world data sets. [34] Furthermore, even if the data roughly lies on a manifold, discontinuities in the data pattern can characterize the manifold as non-smooth. In these situations, graphs with edges that disrespect locality would be extracted and, hence, poor low-dimensional representations would be produced.

Convexity "Isomap relies on the data being geodesically convex." [35] This issue becomes clear when dealing with data sets with "holes" in it that are too big, resulting in great disconnected areas that require paths with great curvature to get around. Lerman demonstrated how reducing non-convex data sets can easily yield distorted results, when these have big enough "holes" on them: [36]



(a) The non-convex
Swiss-roll.

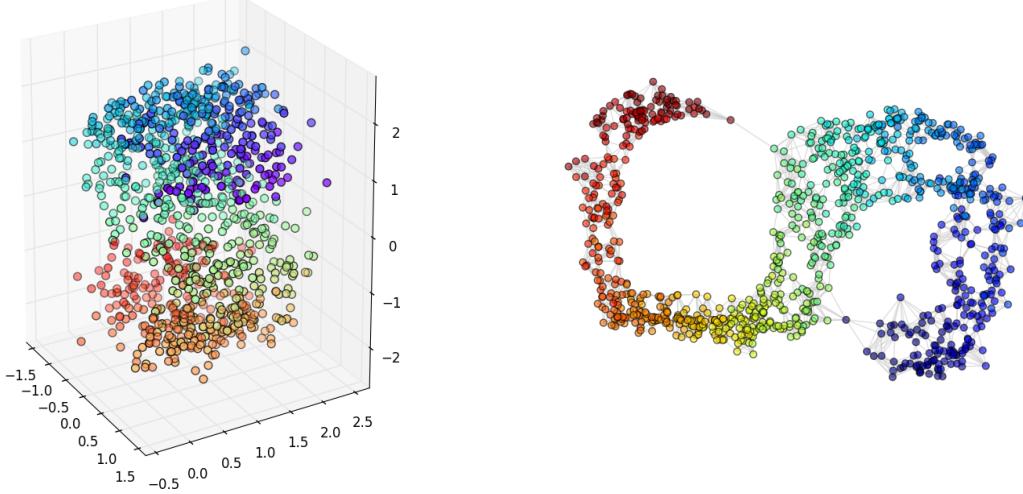
(b) The expected
reduction.



(c) The actual reduction.

Noise Noise in data is expressed by samples that somehow diverge from their neighborhood (outliers). As these samples become farthest from its original neighborhood, the chances of being linked to samples from other neighborhood increase, possibly decreasing the quality of the solution. A possible solution is to remove these samples during the pre-processing stage [33].

The image bellow illustrates an attempt to use Isomap to reduce the $S_{n=.4}$, which is the S manifold subjected to a noise factor of .4:



(a) The data set $S_{n=.4}$.

(b) The neighborhood graph extracted from $S_{n=.4}$.

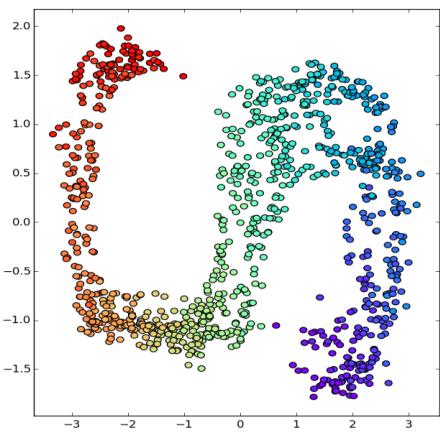


Figure 25: $S_{n=.4}$ reduced.

6 Final Considerations

References

- [1] J. Wang, *Geometric Structure of High-Dimensional Data and Dimensionality Reduction*. Springer Berlin Heidelberg, 2012.
- [2] S. S. A. H. Renear and K. M. Wickett, “Definitions of dataset in the scientific and technical literature,” *ASIS&T 2010*, 2010.
- [3] L. Cayton, “Algorithms for manifold learning,” *Univ. of California at San Diego Tech. Rep*, pp. 1–17, 2005.
- [4] A. Ghodsi, “Dimensionality reduction a short tutorial,” *Department of Statistics and Actuarial Science, Univ. of Waterloo, Ontario, Canada*, 2006.
- [5] T. Cox and M. A. A. Cox, *Multidimensional scaling*. Boca Raton, FL, USA: CRC Press, 2000.
- [6] W. Gander, *The Singular Value Decomposition*, December 2008.
- [7] J. M. Lee, *Manifolds and differential geometry*, vol. 107. American Mathematical Society Providence, 2009.
- [8] E. W. Weisstein, “Stereographic projection.”
- [9] C. Berge and E. Minieka, *Graphs and hypergraphs*, vol. 7. North-Holland publishing company Amsterdam, 1973.
- [10] W. Mayeda, *Graph Theory*. John Wiley & Sons, 1972.
- [11] T. Cormen, *Introduction to Algorithms*. Introduction to Algorithms, MIT Press, 2001.
- [12] P. Langley, *Elements of Machine Learning*. Machine Learning Series, Morgan Kaufmann, 1996.
- [13] L. W. Hosch, “Machine learning.”
- [14] M. Awad and R. Khanna, *Efficient Learning Machines: Theories, Concepts, and Applications for Engineers and System Designers*. New York, NY, USA: Apress, 2015.

- [15] X. Zhu and A. B. Goldberg, “Introduction to semi-supervised learning,” *Synthesis lectures on artificial intelligence and machine learning*, vol. 3, no. 1, pp. 1–130, 2009.
- [16] scikit learn, “Svm: Maximum margin separating hyperplane.”
- [17] J. Weston, “Support vector machine (and statistical learning theory) tutorial.”.
- [18] P. Winston, “Learning: Support vector machine.” 2010.
- [19] P. R. C. D. Manning and H. Schutze, *Introduction to information retrieval*, vol. 1. Cambridge university press Cambridge, 2008.
- [20] R. Rifkin, “Multiclass classification.” 2008.
- [21] “Svm multi-class classification.” 2008.
- [22] scikit learn, “Cross-validation: evaluating estimator performance.”
- [23] scikit learn, “Grid search: Searching for estimator parameterse.”
- [24] B. Rohrer, “How to choose algorithms for microsoft azure machine learning,” 2015.
- [25] J. Shlens, “A tutorial on principal component analysis,” *arXiv preprint arXiv:1404.1100*, 2014.
- [26] G. Dunteman, *Principal components analysis*. No. 69, Newbury Road, CA, USA: Sage, 1989.
- [27] L. I. Smith, “A tutorial on principal components analysis,” *Cornell University, USA*, vol. 51, p. 52, 2002.
- [28] S. Raschka, “Implementing a principal component analysis (pca) in python step by step,” April 2014.
- [29] V. S. T. B. Joshua and J. C. Langford, “A global geometric framework for nonlinear dimensionality reduction,” *Science*, vol. 290, no. 5500, pp. 2319–2323, 2000.

- [30] scikit learn, “Manifold learning: Isomap.”
- [31] F. Ducatelle, “Datasets for data mining.”
- [32] L. Shi and J. Gu, “A fast manifold learning algorithm,” *Information Technology Journal*, vol. 11, no. 3, pp. 380–383, 2012.
- [33] E. P. L. Maaten and J. Herik, “Dimensionality reduction: A comparative review,” *Journal of Machine Learning Research*, vol. 10, no. 1-41, pp. 66–71, 2009.
- [34] T. Lin and H. Zha, “Riemannian manifold learning,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 30, no. 5, pp. 796–809, 2008.
- [35] D. Donoho and C. Grimes, *When Does ISOMAP Recover the Natural Parameterization of Families of Articulated Images?* Technical report (Stanford University. Dept. of Statistics), Department of Statistics, Stanford University, 2002.
- [36] G. Lerman, “Manifold learning techniques: so which is the best?.” 2005.