

Lecture 7: Dimensionality Reduction

Machine learning models are well suited to approach problems with thousands to millions of features. This huge amount of features makes the training very slow and gives rise to selection almost impossible or impractical. This problem is known as the curse of dimensionality.

Simulon fallacies occurs when we attempt to deal with high geometrical dimensions or a huge amount of information. Our comprehension is limited to a 3D world and its projections. Thus, we cannot clearly imagine the features of the 4th time dimension of special relativity or the 11+ dimensions of string theory. The same with the 6N-dimensional plan above of a classical many particle system. In this case we have the problem with the 6-dimensions of momentum and space, needed to specify the state of a physical system and the N-dimensional position of a N-particle. Resulting in a good approximation for a dataset. This problem is also called "Curse's Demon".

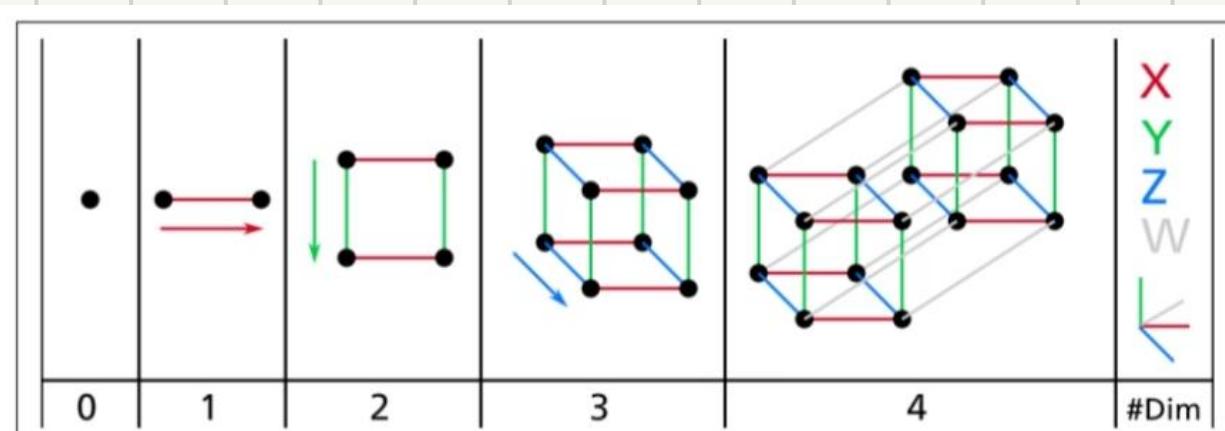


Figure 8-1. Point, segment, square, cube, and tesseract (0D to 4D hypercubes)²

Large dimensions makes the problem to be intractable. Not only due to computational complexity. But because due to our geometrical limitations, we can not take in account how basic facts scales with geometry:

- If you pick a random point in a unit square, only 0.4% of the points are at the very edge. However, in a 10^4 -hyperspace, 99.9% of the points lies at the border.
- If you pick a point randomly in a unit square, the distance between these two points will be on average of 0.52. In a 3D unit cube, about 0.66. In a 10^6 D unit cube, 408.25.

The counterintuitiveness of these basic facts arise from our limitations in interpret dimensions superior to 3. In theory, the solution to these problems, like only about ML, could be more datapoints and more computational power. However, the number of training instances to reach a given density grows exponentially with the number of dimensions. With 100 features, you would need more training instances than the number of atoms in the universe in order to each feature to about 0.1 of each other in a uniform universe.

6 Approaches to dimensionality reduction

There are two main families of approaches for dimensionality reduction: Projection and Manifold learning.

A) Projection

In a real dataset, most of the points are not in a uniform spread. There are extremes that are almost constant and others that are very highly correlated. The result is that most of the samples of a training instance lie on a subspace of a high-dimensional space.

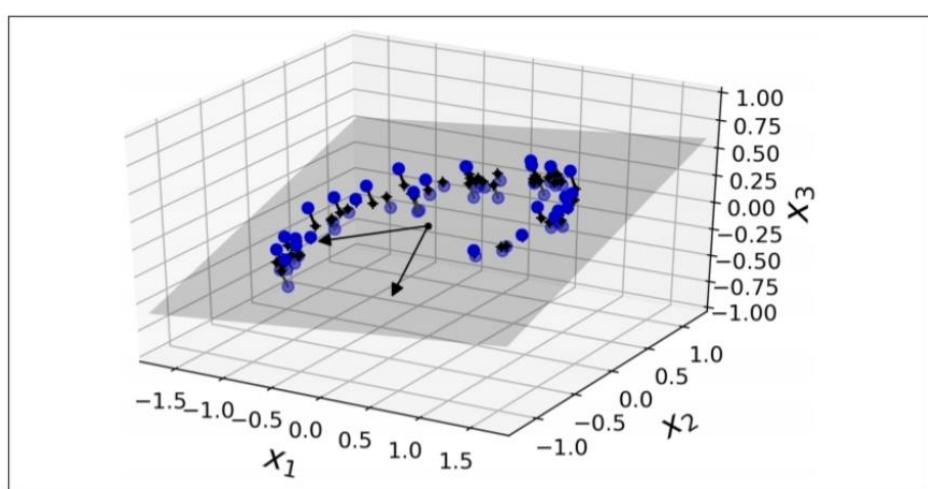


Figure 8-2. A 3D dataset lying close to a 2D subspace

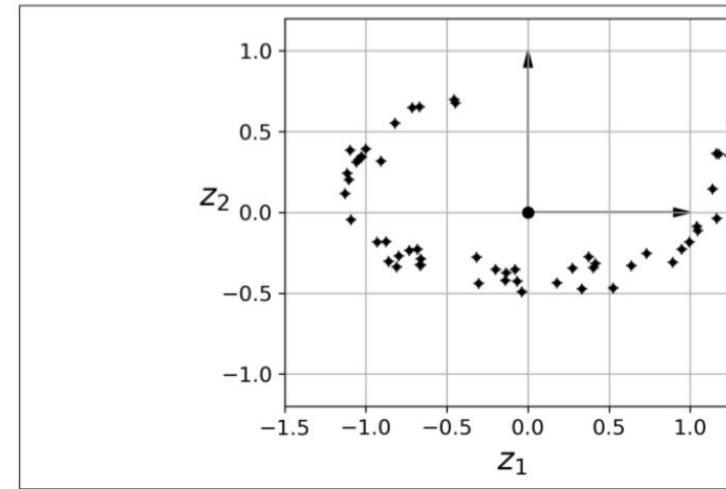
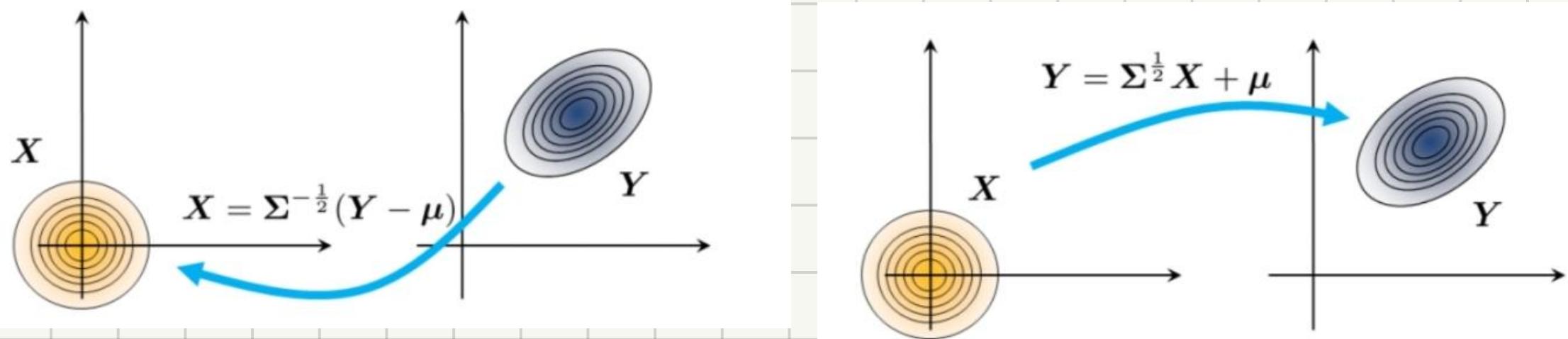


Figure 8-3. The new 2D dataset after projection

In the above picture, we have a 3D dataset, however most of the training instances lies on a hyperplane. One can project all points in this hyperplane and reduce the dimensionality to a 2D problem. Of course this is not the general case. Every dimensionality reduction comes with a loss of information. There are simple cases in which this loss is minimal because the information is easy to reconstruct or because there is a symmetry hidden in the data points.

A good example of information that can be recovered due to symmetry is the transformation property of a gaussian random variables.



Due to all informational content is stored in the properties of the covariance matrix Σ , we can only reconstruct the original data from a standard gaussian, using rotations and translations. Note that this is precisely the case of the hyperplane example.

However, a 3D plot twist and roll, like the Swiss roll dataset, may not have any clear symmetry and a projection may not be a good solution:

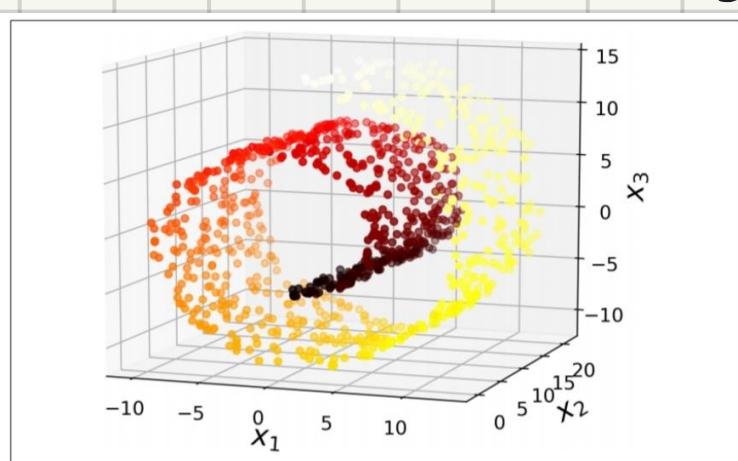


Figure 8-4. Swiss roll dataset

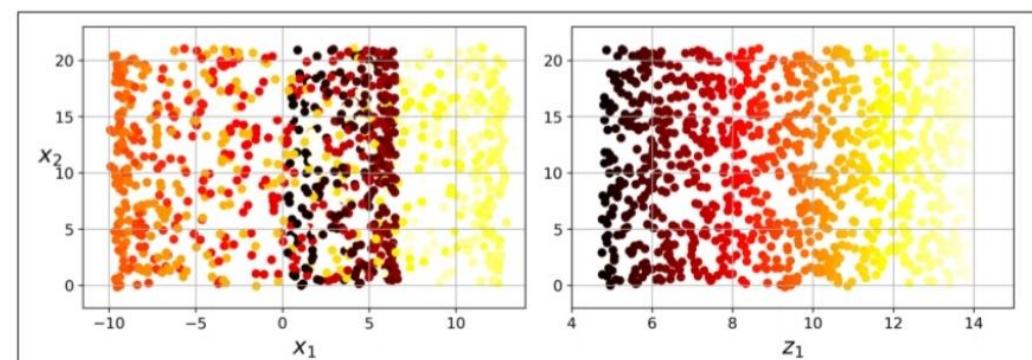


Figure 8-5. Squashing by projecting onto a plane (left) versus unrolling the Swiss roll (right)

Single projections in 2D will squish all different layers of the roll together and the information is lost.

B) Manifold Learning

The Swiss roll is a 2D manifold, i.e. a shape that can be bent and twisted in a higher-dimensional space. A d-dimensional manifold is an object that is part of a n-dimensional space, $d < n$, that locally resembles a d-dimensional object.

In this case, a dimensionality reduction algorithm works by modeling the manifold on which the training instances lies. This is called manifold learning. This is an elaboration over the so-called manifold hypothesis, in which every dataset in a high-dimensional space lies close to a low-dimensional manifold. This assumption can not be proven, but is empirically observed. Of course, we are also assuming that the task at hand will be simpler in this low-dimensional manifold. This does not always hold.

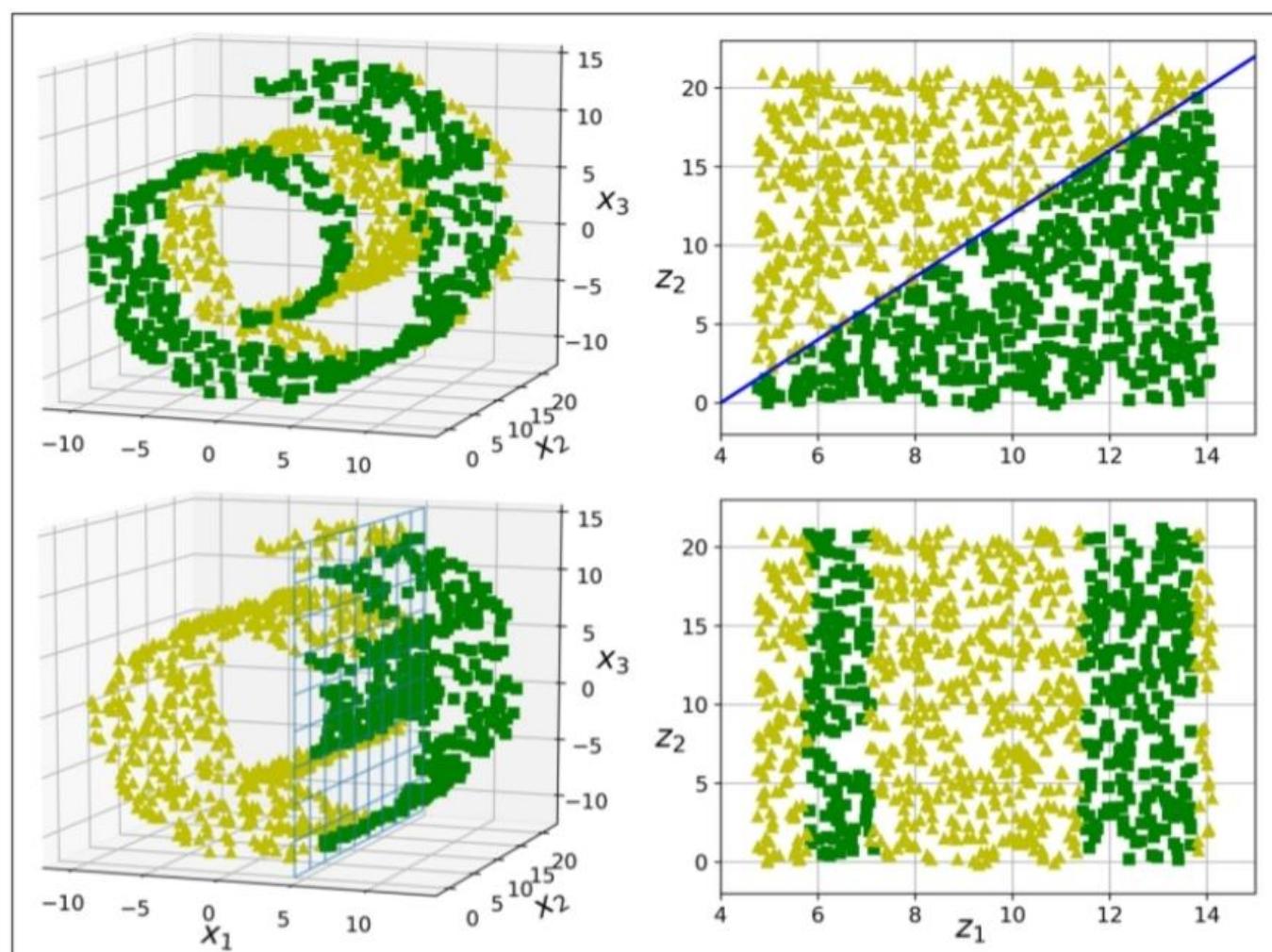


Figure 8-6. The decision boundary may not always be simpler with lower dimensions

Reducing dimensionality of your training set before training the model will usually speed up training but it may not always lead to better predictions.

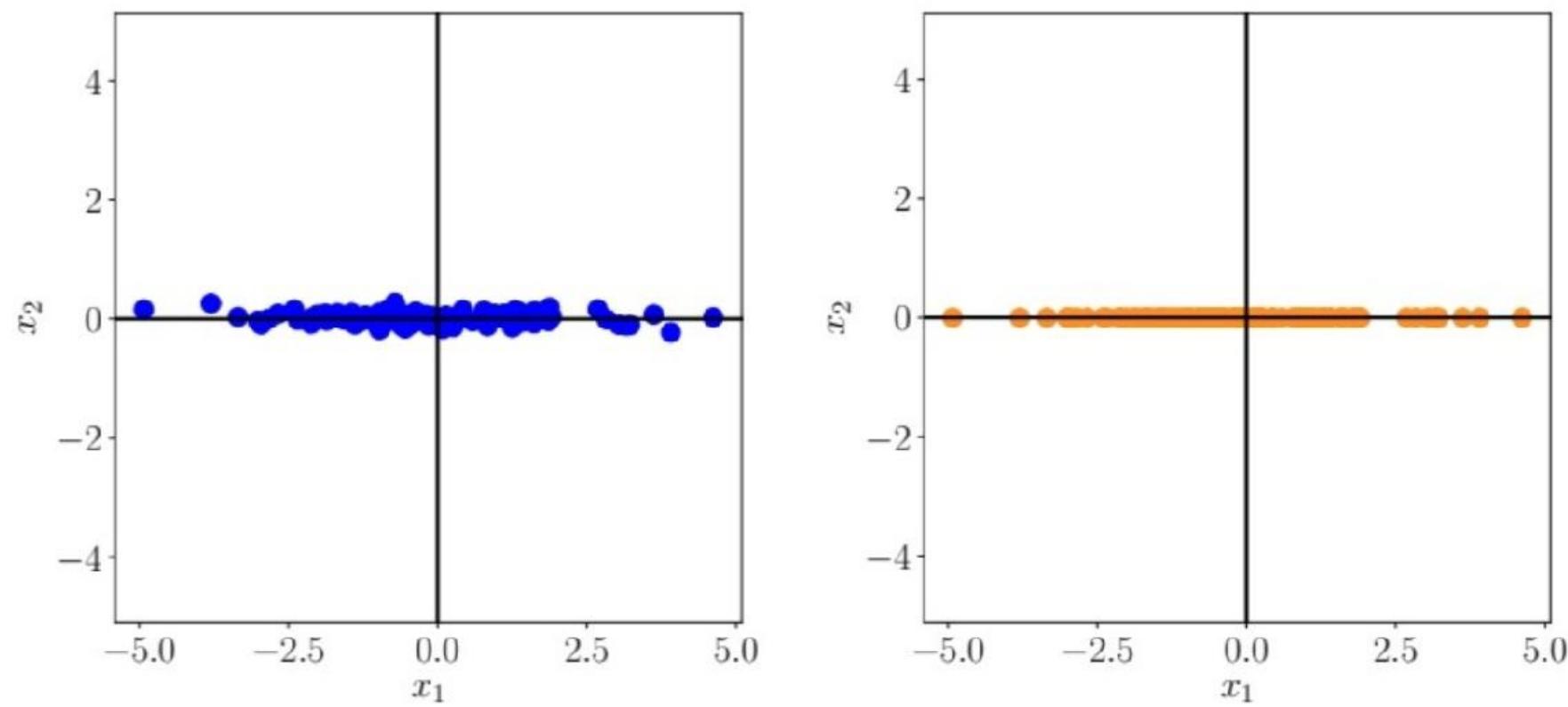
Principal Component Analysis

In what follows, we are going to deduce one of the most used methods for dimensionality reduction: the principal component analysis. This is a projection method and lies its roots in the algebraic properties of the covariance matrix or in the geometrical properties of Random Vectors over Euclidean Space.

§ The PCA Problem

Principal component Analysis looks for projections \tilde{x}_n of datapoints x_n , as similar as possible of the original, but with significantly lower intrinsic dimensionality.

This is in the heart of the so called feature extraction or feature selection



(a) Dataset with x_1 and x_2 coordinates.

(b) Compressed dataset where only the x_1 coordinate is relevant.

Suppose we have a vector \vec{x} , with $\dim(\vec{x}) = m$, and we want to transmit the encoded information with a vector $\tilde{\vec{x}}$ with $\dim(\tilde{\vec{x}}) = l$, and $l < m$. Of course, there will be a loss of information, however one might ask:

"There is a linear map $\tilde{\vec{x}} = \vec{T}\vec{x}$ so that a compression $\tilde{\vec{x}} = \vec{T}\vec{x}$ is optimal, given a loss function?".

This loss function must be related with (a) variational, (b) mean-squared-error. (c) geometry.

Figure 10.1
Illustration:
dimensionality
reduction. (a) The
original dataset
does not vary much
along the x_2
direction. (b) The
data from (a) can be
represented using
the x_1 -coordinate
alone with nearly no
loss.

This \hat{T} transformation has to possess the property that some of its components has low-variance. The transformation with this property is called Principal Component Analysis.

Consider a i.i.d dataset $\mathcal{X} = \{\vec{x}_i\}_{i=1}^N, \vec{x}_i \in \mathbb{R}^D$ and

$$(1) \quad \mathbb{E}[\vec{x}] = 0$$

and with a covariance matrix:

$$(2) \quad \overset{\leftrightarrow}{S} = \frac{1}{N} \sum_{i=1}^N \vec{x}_i \vec{x}_i^+$$

We claim that there is a transformation \hat{T} that give us a compressed representation:

$$(3) \quad \vec{z}_i = \hat{T} \vec{x}_i = \hat{B}^T \vec{x}_i \in \mathbb{R}^M$$

where we define the projection operator:

$$(4) \quad \hat{B} := [\vec{b}_1, \dots, \vec{b}_M] \in \mathbb{R}^{D \times M}$$

with orthonormal columns:

$$(5) \quad \vec{b}_i^T \cdot \vec{b}_j = \delta_{ij}$$

Our goal is to find a M -dimensional sub-space

$U \subseteq \mathbb{R}^D$, $\dim(U) = M < D$, onto which we project

the data. The projection is denoted as $\vec{\tilde{x}}_i \in U$, and

their coordinate with respect to the basis $\{\vec{b}_i\}$ of

U , by \vec{z}_i . We seek for projections $\vec{\tilde{x}}_i \in \mathbb{R}^P$, similar

to \vec{x}_i under a classification loss metric and can be
efficiently described in the $\{\vec{b}_i\}$ basis.

As our clustering, consider a dataset X ,
where we had subtracted the mean
of all instances, so that:

$$(6) \quad E[\vec{x}] = \vec{m}$$

$$(7) \quad \vec{x} \rightarrow \vec{x} - \vec{m} \Rightarrow E[\vec{x}] = \vec{0}$$

The projection over the orthonormal
 b -basis:

$$(8) \quad A = \vec{x}^T \cdot \vec{b} = \vec{b}^T \cdot \vec{x}$$

$$(9) \quad E[A] = \vec{b}^T \cdot E[\vec{x}] = 0$$

the variance of the projection:

$$\begin{aligned}(10) \quad \sigma^2 &= E[A^2] - (E[A])^2 \\&= E[\vec{b}^\top \vec{x} \vec{x}^\top \vec{b}] \\&= \vec{b}^\top E[\vec{x} \cdot \vec{x}^\top] \vec{b} \\&= \vec{b}^\top \overset{\leftrightarrow}{R} \vec{b}\end{aligned}$$

Where $\overset{\leftrightarrow}{R}$ is the covariance matrix. The covariance function is symmetric, and is a property of the data set. Thus, the variance of the projection is a function of the orthonormal basis:

$$\begin{aligned}(11) \quad \psi(\vec{b}) &= \sigma^2 \\&= \vec{b}^\top \overset{\leftrightarrow}{R} \vec{b}\end{aligned}$$

These vectors has to be chosen in a direction so that σ^2 is a extremum or least stationary. Thus, this condition means that:

$$(12) \quad \psi(\vec{b} + s\vec{b}) = \psi(\vec{b})$$

$$(13) \quad \phi(\vec{b} + \delta\vec{b}) = (\vec{b} + \delta\vec{b})^\top \overset{\leftrightarrow}{R} (\vec{b} + \delta\vec{b})$$

$$= \vec{b}^\top \overset{\leftrightarrow}{R} \vec{b} + \vec{b}^\top \overset{\leftrightarrow}{R} \delta\vec{b} + \delta\vec{b}^\top \overset{\leftrightarrow}{R} \vec{b} +$$

~~$\delta\vec{b}^\top \overset{\leftrightarrow}{R} \delta\vec{b}$~~

Up to 1st order
of perturbation

$$= \vec{b}^\top \overset{\leftrightarrow}{R} \vec{b} + 2\delta\vec{b}^\top \overset{\leftrightarrow}{R} \vec{b}$$

Using (12): $\delta\vec{b}^\top \overset{\leftrightarrow}{R} \vec{b} = 0$, and we recover (13).

This constraint gives us a clue about the direction of the projection over the basis:

$$(14) \quad \|\vec{b} + \delta\vec{b}\| = 1 \Rightarrow (\vec{b} + \delta\vec{b})^\top (\vec{b} + \delta\vec{b}) = 1$$

again up to 1st order in perturbation theory,
we have:

$$(15) \quad \delta\vec{b}^\top \vec{b} = 0$$

thus the variations has to be orthogonal to
the basis vectors! This means that only
transformations that the Euclidean norm
remains unitary are admissible, and those
ones are in the same direction of \vec{b} :

$$(16) \quad \max_{\vec{b}} \vec{b}^\top \overset{\leftrightarrow}{R} \vec{b}$$

subject to $\|\vec{b}^\top \overset{\leftrightarrow}{R} \vec{b}\| = 1$

Using the method of Lagrange multipliers:

$$(17) \quad \mathcal{L}(\hat{b}, \lambda) = \hat{b}^T \hat{R} \hat{b} + \lambda(1 - \hat{b}^T \hat{b})$$

$$(18) \quad \frac{\partial \mathcal{L}}{\partial \hat{b}} = \hat{b}^T \hat{R} - \lambda \hat{b} = 0 \Rightarrow \hat{b}^T \hat{R} = \hat{b}^T \lambda$$

$$(19) \quad \frac{\partial \mathcal{L}}{\partial \lambda} = 1 - \hat{b}^T \hat{b} = 0 \Rightarrow \hat{b}^T \hat{b} = 1$$

Resulting in a boundary value problem:

$$(20) \quad \hat{R} \hat{b} = \lambda \hat{b}$$

$$(21) \quad \hat{b}^T \hat{b} = 1$$

Let the $\lambda_1, \dots, \lambda_m$ eigenvalues ordered in decreasing order:

$$(22) \quad \lambda_1 > \dots > \lambda_m$$

and we have a set of basis vectors:

$$(23) \quad \hat{B} = [\hat{b}_1, \dots, \hat{b}_m]$$

we have:

$$(24) \quad \hat{R} \hat{B} = \hat{B} \hat{\Lambda}$$

where

$$(25) \quad \Lambda = \text{diag} [\lambda_1, \dots, \lambda_m]$$

the dataset can be represented as:

$$(26) \quad \overset{\leftrightarrow}{R} = \tilde{B} \tilde{\Lambda} \tilde{B}^+$$

and by means of the spectral theorem:

$$(27) \quad \overset{\leftrightarrow}{R} = \sum_i \lambda_i \tilde{b}_i \tilde{b}_i^+$$

Principal component Analysis and the eigen decomposition of the covariance matrix are basically the same:

$$(28) \quad \phi(\tilde{b}_i) = \lambda_i$$

From this development, we see that every single vector in the feature space can be represented as a linear combination of the eigenvectors of the covariance matrix

$$(29) \quad \begin{aligned} \vec{x} &= \sum_i \alpha_i \tilde{b}_i \\ &= \begin{bmatrix} \tilde{b}_1 & \cdots & \tilde{b}_m \end{bmatrix} \begin{bmatrix} \alpha_1 & \cdots & \alpha_m \end{bmatrix} \end{aligned}$$

where the a_i 's one called principal component of the feature vector.

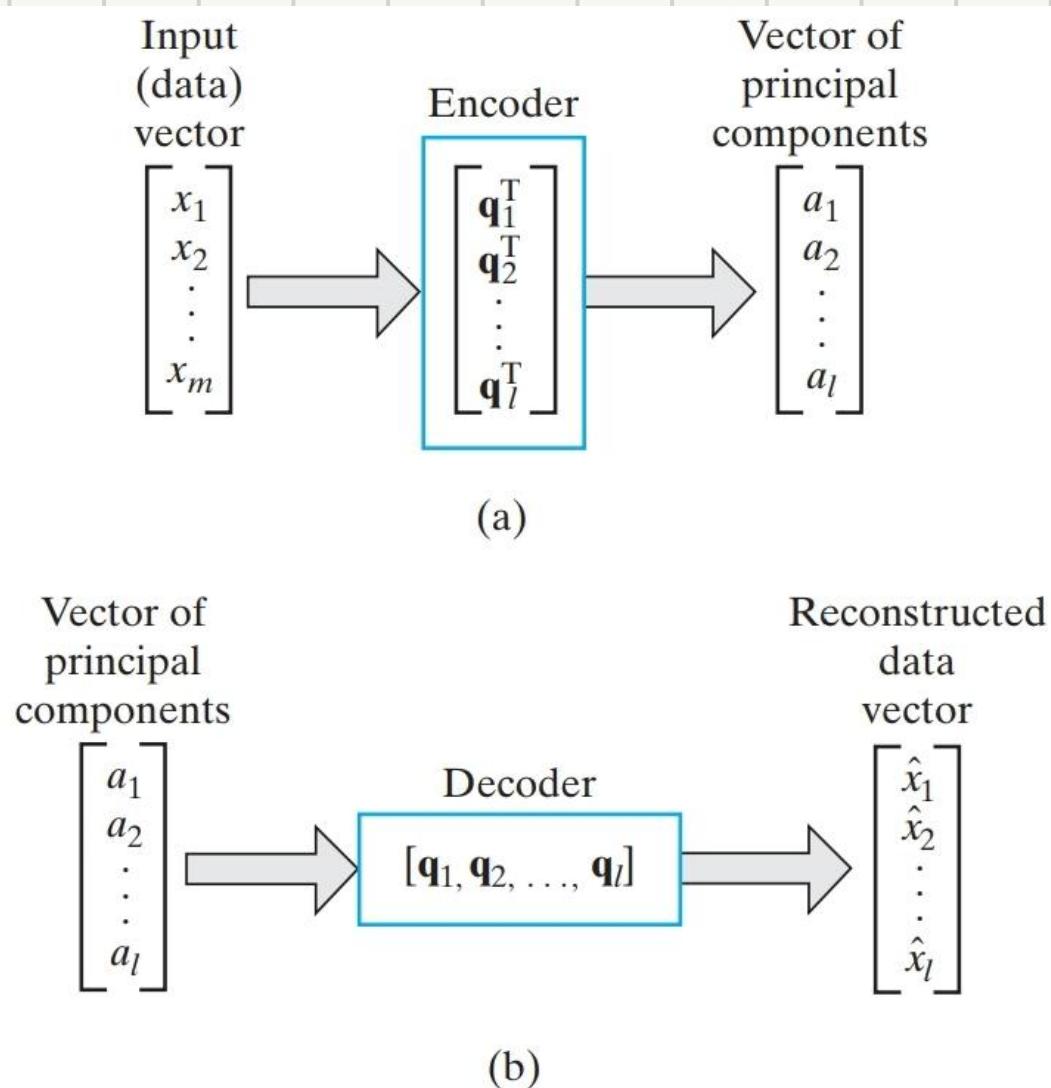


FIGURE 8.2 Illustration of two phases of principal-components analysis: (a) Encoding, (b) Decoding.

This approach is called minimal variance procedure.

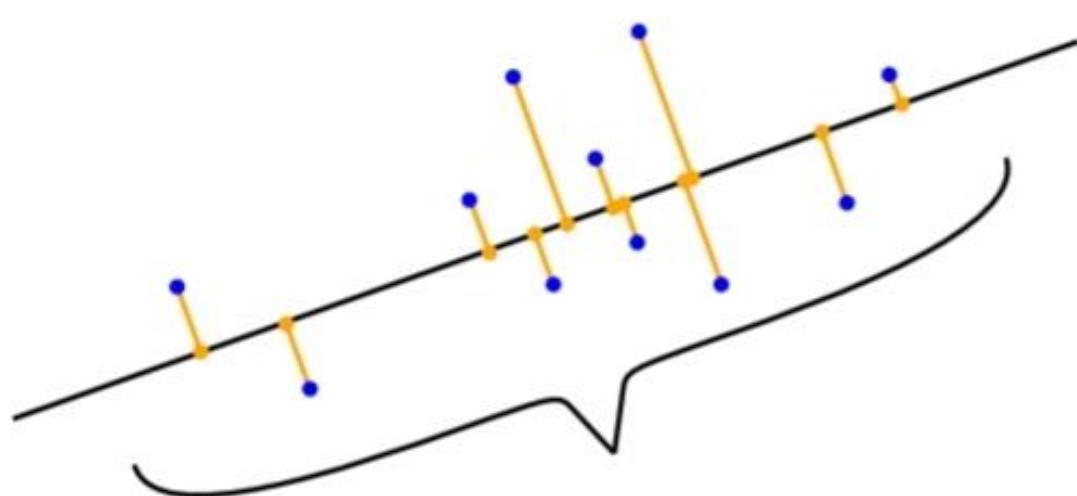


Figure 10.4 PCA finds a lower-dimensional subspace (line) that maintains as much variance (spread of the data) as possible when the data (blue) is projected onto this subspace (orange).

When we looks for a new basis when w , can span the data set, truncating in a dimension $k < m$, retaining as most variance as possible.

The linear map that goes from the original data space \mathbb{R}^m to \mathbb{R}^l is called "encoder". And the way back transformation is called "decoder".

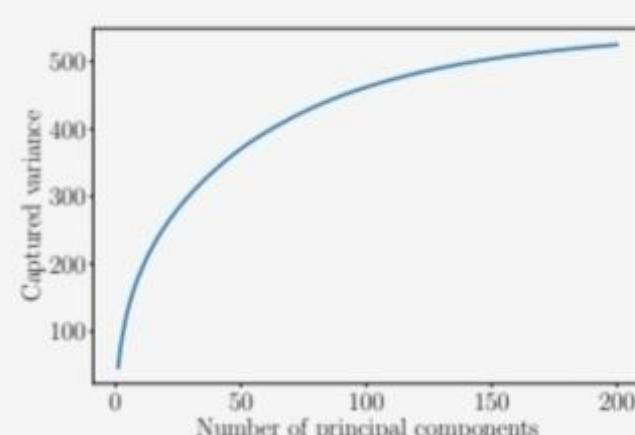
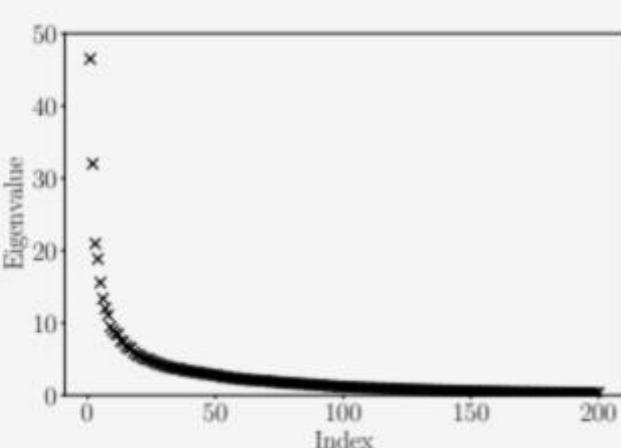
The general approach to dimensionality reduction is to select the first l biggest eigenvalues of the ordered eigenvalues and encode all information into these l basis vector.

The truncation error is given by:

$$(30) \quad \hat{e} = \vec{x} - \vec{\tilde{x}} \\ = \sum_{i=l+1}^m a_i \vec{b}_i$$

and the retained variance of the original data is given by:

$$(31) \quad \hat{\sigma}^2 = \sum_{j=1}^m \hat{\sigma}_j^2 = \sum_{j=1}^m \lambda_j$$



(a) Eigenvalues (sorted in descending order) of the data covariance matrix of all digits "8" in the MNIST training set.

(b) Variance captured by the principal components.

The dimension of the subspace is a tradeoff between complexity on 1 retained amount of variance.

Eigenvector Computation

The space of principal components is a vector space formed by the eigenvectors of the covariance matrix:

$$(32) \quad \overleftrightarrow{R} = \frac{1}{N} \sum_{i=1}^N \vec{x}_i \vec{x}_i^T = \vec{X} \cdot \vec{X}^T$$

Please, note that $\vec{X} \in \mathbb{R}^{N \times D}$, being N the number of observations and D the number of features, this is a huge rectangular matrix.

The eigendecomposition can be performed in two ways:

i) Directly from the matrix \overleftrightarrow{R} ,

ii) Using the singular value decomposition.

The second approach is more reliable to our purposes:

$$(33) \quad \overleftrightarrow{X} = \underbrace{\overleftrightarrow{U}}_{D \times N} \underbrace{\overleftrightarrow{\Sigma}}_{D \times D} \underbrace{\overleftrightarrow{V}^T}_{D \times N} \quad \overleftrightarrow{V}^T \in \mathbb{R}^{N \times N}$$

where $\overleftrightarrow{U} \in \mathbb{R}^{D \times D}$ and $\overleftrightarrow{V} \in \mathbb{R}^{N \times N}$ one orthogonal matrices and $\overleftrightarrow{\Sigma} \in \mathbb{R}^{D \times N}$ is a matrix where

whose non-zero entries are the singular values $\sigma_{ii} \geq 0$. It follows

$$\begin{aligned}
 (34) \quad \overset{\leftrightarrow}{R} &= \frac{1}{N} \overset{\leftrightarrow}{X} \overset{\leftrightarrow}{X}^T \\
 &= \frac{1}{N} \overset{\leftrightarrow}{U} \sum \underbrace{\overset{\leftrightarrow}{V} \overset{\leftrightarrow}{V}^T}_{\mathbb{I}_N} \sum \overset{\leftrightarrow}{U}^T \\
 &= \frac{1}{N} \overset{\leftrightarrow}{U} \sum \sum \overset{\leftrightarrow}{U}^T
 \end{aligned}$$

the columns of $\overset{\leftrightarrow}{U}$ are the eigenvectors of $\overset{\leftrightarrow}{R}$

and the eigenvalues of λ_d of $\overset{\leftrightarrow}{R}$ are related with the singular values of $\overset{\leftrightarrow}{X}$ via:

$$(35) \quad \lambda_d = \frac{\sigma_d}{N}$$

thus, the maximum variance principal component analysis is just the singular value decomposition of the $\overset{\leftrightarrow}{X}$ matrix.

In order to maximize the variance of the projected data, PCA chooses the columns of $\overset{\leftrightarrow}{U}$ as the eigenvectors of a M -dimensional space of the M -longest eigenvalues of $\overset{\leftrightarrow}{R}$.

This means that, we identify $\overset{\leftrightarrow}{B}$ as the projection of $\overset{\leftrightarrow}{U}$ into this M -dimensional space. The optimal low-rank approximation is then obtained by the Eckart-Young theorem:

$$(36) \quad \tilde{X}_N \stackrel{\sim}{=} \underset{\text{rank } \tilde{A} \leq M}{\operatorname{argmin}} \| \tilde{X} - \tilde{A} \|_F^2 \in \mathbb{R}^{D \times N}$$

where $\| \cdot \|_F$ is the spectral norm:

Thus, we obtain \tilde{X}_N by extracting the SVD representation of the top- M singular values:

$$(37) \quad \tilde{X}_N = \underbrace{U_M}_{D \times N} \sum_M \underbrace{V_M^T}_{M \times M} \in \mathbb{R}^{D \times N}$$

§ PCA Dimensionality Reduction Program

Given an arbitrary dataset, the PCA procedure runs the following steps:

1- Mean subtraction

2- Standardization

3- Spectral factorization of the Covariance matrix.

4- Projection.



" "

