

Notes on Principal Component Analysis

Kamila Zdybał

camllejr.github.io/science-docs

kamila.zdybal@ulb.ac.be, kamila.zdybal@gmail.com

Preface

These are notes on **Principal Component Analysis** (PCA), a dimensionality reduction technique in which the data set is reduced to maintain only the directions of the largest variance.

The deep and intuitive understanding of PCA requires the deep and intuitive understanding of basic linear algebra. I highly recommend a very pleasant to watch course by 3Blue1Brown YouTube channel, which, to my belief, will be everything you need to master about linear algebra to gain a great intuition behind PCA.

This document is still in preparation. Please feel free to contact me with any suggestions, corrections or comments.

Keywords

principal component analysis, data reduction, dimensionality reduction, linear algebra, MATLAB®

Contents

1	Motivation for data reduction	1
2	Datasets for PCA	1
3	Covariance matrix	2
3.1	The goal of PCA in terms of a covariance matrix	2
4	PCA workflow	2
5	Why does it need to be eigenvectors?	3
6	PCA example	3
7	Linear vs. non-linear PCA	3
8	Interesting questions	4
A	APP1	4
B	APP2	4

1 Motivation for data reduction

There are several questions which stimulated the development of data reduction techniques and reduced-order modeling:

1. Can we send less data but preserve maximum information contained by the data (data compression)?
2. If the data were measurements from an experiment, can we predict what the outcome of another measurement will be?
3. Can we build a model that will make predictions about a system?
4. Can we make sense of a large data set with high-dimensionality (which cannot be plotted graphically)?

2 Datasets for PCA

For the rest of this document, we assume that the dataset for performing PCA is structured as follows: each column represents one variable Var_j and the rows correspond to variable observations Obs_i (e.g. at different positions in space, or in time).

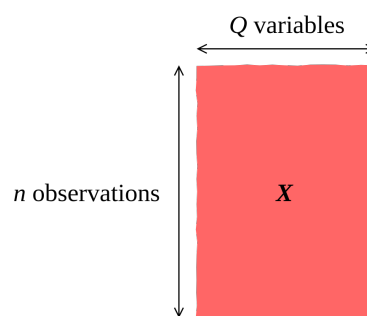


Figure 1: Data matrix for PCA.

This is also the data format that is needed for the MATLAB® command `pca`. From the documentation:

`pca(X)`
Rows of X correspond to observations and columns correspond to variables.

3 Covariance matrix

We start our discussion of PCA with explaining the concept of a *covariance matrix*, since computing this matrix is a starting point for performing PCA. The covariance matrix is given by:

$$S = \frac{1}{Q-1} X^T X \quad (1)$$

The covariance matrix S is therefore size $(Q \times Q)$.

We will start with exploring the meaning of $X^T X$. Let's look at the graphical representation of this matrix multiplication in Fig. 2.

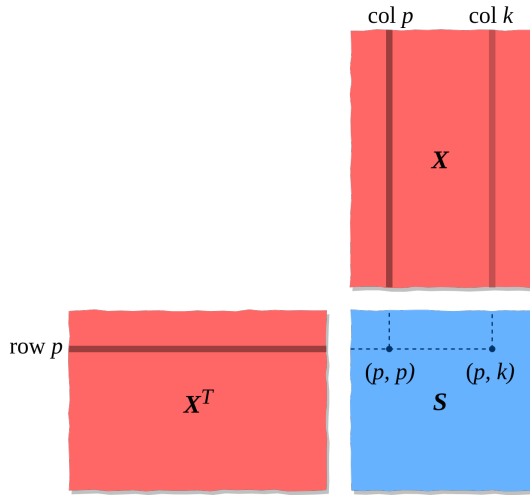


Figure 2: Covariance matrix S graphical interpretation.

Any given column, say p or k , of the data matrix X represents all measurements of a single variable Var_p and can be viewed as a single n -dimensional vector. The same can be said about any row of a matrix X^T , since this is simply a transpose of the data matrix X .

Notice that an element at position (p, k) inside the matrix S has the interpretation of a dot product between a vector formed by the p -th row of a matrix X^T and a k -th column of a matrix X .

$$s_{(p,k)} = \text{dot}(X^T(p,:), X(:,k)) \quad (2)$$

In a special case, where we multiply the row p with the column p , we get a dot product of a vector with itself.

$$s_{(p,p)} = \text{dot}(X^T(p,:), X(:,p)) \quad (3)$$

In general, the dot product between two vectors x and y represents how much vector x lays in the direction of vector y (and vice versa) - and it is zero when two vectors are perpendicular to each other. This intuition can be carried to our covariance matrix S . If any off-diagonal element is non-zero, say element at position (p, k) this means that some information about variable Var_p is carried by a variable Var_k (and vice versa).

A covariance of two vectors:

$$\text{cov}(x, y) = \quad (4)$$

3.1 The goal of PCA in terms of a covariance matrix

What Principal Component Analysis aims to do is to find a new, transformed data set Z such that if we computed a new covariance matrix:

$$S_Z = Z^T Z \quad (5)$$

the variances (the elements on the diagonal) are maximized and the covariances (the off-diagonal elements) are zero. This means that no more information about any column of Z is carried by any other column of Z .

4 PCA workflow

The next step in PCA is to perform the eigendecomposition of the covariance matrix:

$$\text{eig}(S) = [A, \Lambda] \quad (6)$$

The matrix A is a matrix of eigenvectors and it is size $(Q \times Q)$. Each eigenvector is called a *principal component* in the PCA world. The diagonal matrix Λ of size $(Q \times Q)$ is a matrix of corresponding eigenvalues.

The principal components now form a new basis in which we can represent our data set. We therefore perform a transformation of the original data matrix X from the original space to the new space represented by the eigenvectors of the covariance matrix. This transformation is achieved by the following multiplication:

$$Z = XA \quad (7)$$

The new matrix Z is still our dataset X but represented in the basis associated with the matrix A . It is also called the *PC-scores* since one may think of every element in this matrix as the "score" that the corresponding element in the original data matrix would get when represented in the new coordinate system after transformation.

We now approach the dimensionality reduction but first let's obtain the original data set back given that we know the PC-scores and the transformation matrix:

$$X = ZA^T \quad (8)$$

The above equation is our route back to obtain the original data set in which the PC-scores are projected on the basis associated with a transposed eigenvectors matrix A .

Suppose that we would like to find the approximation of the data matrix X with only q principal components (we project the PC-scores onto only q out of Q principal components).

We shrink the transformation matrix A to be of size $(Q \times q)$ (we only keep q principal components). To match the matrix sizes we also need to shrink in size the PC-scores matrix which originally is size $(n \times Q)$ - the same size as the data matrix X . We will denote these truncated matrices A_q and Z_q respectively.

Projecting Z_q onto the basis A_q^T will result in an approximation of the original data set:

$$X_q = Z_q A_q^T \quad (9)$$

5 Why does it need to be eigenvectors?

In this section we come back to equation (6) and answer the question: why are principal components the eigenvectors of a covariance matrix? If you are new to PCA, this indeed might seem like a not-related-to-anything thing to do. So why is it eigenvectors? And why does it need to be eigenvectors of this special matrix called covariance matrix?

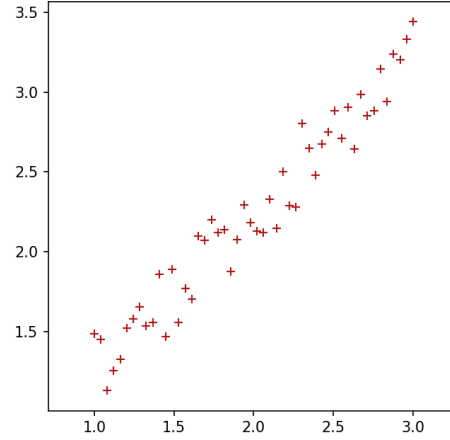


Figure 3: Data set for linear-PCA.

We may approximate the above data set simply by:

$$\text{Data} \approx A_1 P\vec{C}1 \quad (11)$$

6 PCA example

7 Linear vs. non-linear PCA

The linear-PCA is adequate for data sets that can be approximated by a linear combination:

$$\text{Data} \approx A_1 P\vec{C}1 + A_2 P\vec{C}2 + \dots + A_n P\vec{C}n \quad (10)$$

where $P\vec{C}1, P\vec{C}2, \dots, P\vec{C}n$ are consecutive principal components (PCs).

The example of a data set that is suited for linear-PCA is presented in Figure 3. What is unique about such data set is that it has got a preferred direction in a lower dimension (in this case 1D) along which the data points are aligned. The first PC will be associated with this direction. There will also be a second PC which is perpendicular to the first one, however the weight of this second PC is much lower - there is much less variance of data along this other direction.

The data reduction that we can perform in our heads is that this data set almost aligns with a linear function $f(x) = x$ for $x \in \langle 1; 3 \rangle$.

When we take a look at a data set in Figure 4, we can see that the data is spread almost equally along two dimensions. The two PCs of such data set will have almost the same weight (the length of the two eigenvectors will be almost the same). It seems that there is no preferred direction in this data. So, does that mean that it cannot be reduced?

The underlying pattern in this data set is very clearly visible. Intuitively, we can say that the data points are almost aligned with a circle centered at $(2; 2)$ with radius $r = 1$. If we knew a function $f(x)$ describing that circle, perhaps we might be able to write it in terms of just the first PC and approximate:

$$\text{Data} \approx f(P\vec{C}1) \quad (12)$$

where $f(P\vec{C}1)$ is essentially non-linear.

One idea to find the function f might be to associate the first PC with the radius of the circle (that is find the relation $r(PC1)$) and write:

$$f(x, P\vec{C}1)^2 = x^2 + r(P\vec{C}1)^2 \quad (13)$$

This will give us a centered approximation to the original data.

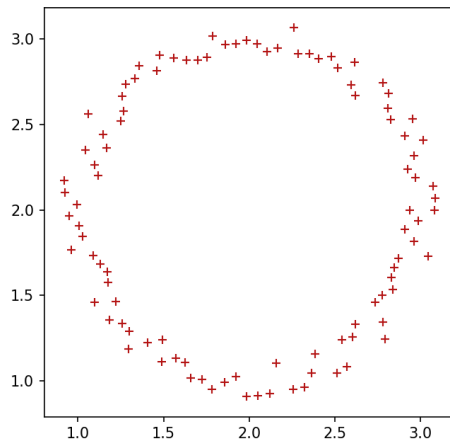


Figure 4: Data set for non-linear-PCA.

To conclude, the non-linear-PCA is helpful at times when there is no preferred direction in a data set, however there is an underlying non-linear relation that describes the data. The linear-PCA would in such case not give a satisfactory approximation after truncating the number of PCs.

8 Interesting questions

Is the covariance matrix of PC-scores a diagonal matrix?

A APP1

B APP2

References

- [1] Ian T. Jolliffe *Principal Component Analysis*, Second Edition, 1986
- [2] Gilbert Strang *Introduction to Linear Algebra*, Fifth Edition, 2016
- [3] Jonathon Shlens, *A Tutorial on Principal Component Analysis*, 2016, <https://arxiv.org/abs/1404.1100>
- [4] <http://people.sju.edu/~pklingsb/dot.cov.pdf>