

Notes on Principal Component Analysis

Kamila Zdybał

Université libre de Bruxelles, kamila.zdybal@ulb.ac.be
camillejr.github.io/science-docs, 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®, Python

Contents

1	Motivation for data reduction
2	Datasets for PCA
3	Data pre-processing
4	Covariance matrix
5	PCA workflow
6	Why does it need to be eigenvectors?
7	PCA Python example
8	Interesting questions
A	APP1
B	APP2

1 Motivation for data reduction

There are several questions which stimulated the usage of data reduction and data decomposition techniques:

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).

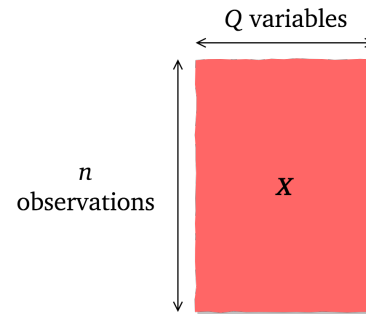


Figure 1: Data matrix for PCA.

This is also the data format that is needed for the MATLAB® command `pca`. From the documentation [1]:

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

3 Data pre-processing

In data science we are typically given a raw data set which is not centered and not scaled.

4 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:

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

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

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

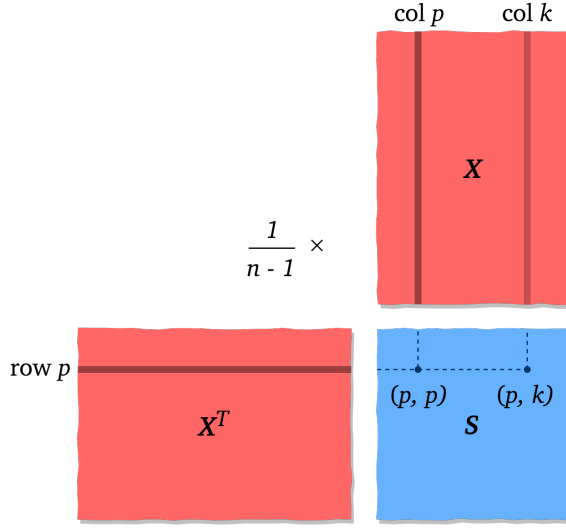


Figure 2: Covariance matrix \mathbf{S} graphical interpretation.

Any given column, say p or k , of the data matrix \mathbf{X} represents all measurements of a single variable Var_p (or Var_k) and can be viewed as a single n -dimensional vector. The same can be said about any row of a matrix \mathbf{X}^T .

Notice that an element at position (p, k) inside the matrix \mathbf{S} has the interpretation of a dot product between a vector formed by the p -th row of a matrix \mathbf{X}^T and a k -th column of a matrix \mathbf{X} . There is also the $\frac{1}{n-1}$ factor out front to which we will come back later.

$$s(p, k) = \frac{1}{n-1} \text{dot}(\mathbf{X}^T(p, :), \mathbf{X}(:, k)) \quad (2)$$

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

$$s(p, p) = \frac{1}{n-1} \text{dot}(\mathbf{X}^T(p, :), \mathbf{X}(:, p)) \quad (3)$$

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 \mathbf{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).

In other words, every element in the covariance matrix is:

$$s(i, j) = \frac{1}{n-1} \sum_{k=1}^n \text{Var}_i(k) \text{Var}_j(k) \quad (4)$$

The goal of PCA in terms of a covariance matrix

Principal Component Analysis aims to find a new, transformed data set \mathbf{Z} such that if we computed a new covariance matrix:

$$\mathbf{S}_Z = \mathbf{Z}^T \mathbf{Z} \quad (5)$$

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 \mathbf{Z} is carried by any other column of \mathbf{Z} .

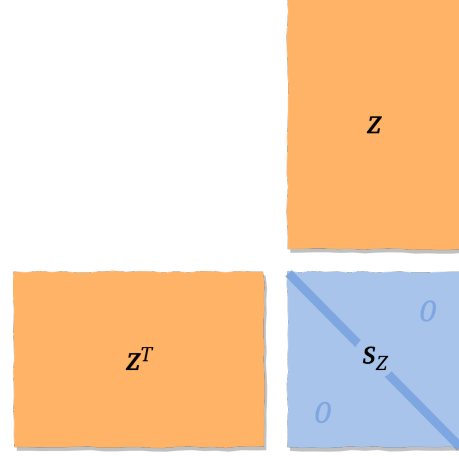


Figure 3: New transformed data set \mathbf{Z} and its covariance matrix \mathbf{S}_Z .

5 PCA workflow

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

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

The matrix \mathbf{A} is a matrix of eigenvectors and it is size $(Q \times Q)$. Each eigenvector is called a *principal component* (PC). The principal components are orthogonal to each other, and therefore the following important property holds: $\mathbf{A}^T = \mathbf{A}^{-1}$ (proof¹).

The diagonal matrix $\mathbf{\Lambda}$ of size $(Q \times Q)$ is a matrix of corresponding eigenvalues.

Given the eigendecomposition of matrix \mathbf{S} , we may state that: $\mathbf{S} = \mathbf{A} \mathbf{\Lambda} \mathbf{A}^T$.

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

$$\mathbf{Z} = \mathbf{X} \mathbf{A} \quad (7)$$

The new matrix \mathbf{Z} is still our dataset \mathbf{X} but represented in the basis associated with the matrix \mathbf{A} . It is also called the PC-scores matrix, since one may think of every element in this matrix as the "score" that the corresponding element in \mathbf{X} gets when represented in the new coordinate system after transformation.

In the matrix multiplication from eq.(7), every variable vector inside \mathbf{X} gets transformed by the transformation matrix \mathbf{A} and attains new scores in the basis associated with \mathbf{A} . The new representation of the old variable is now kept in the matrix \mathbf{Z} .

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

$$\mathbf{X} = \mathbf{Z} \mathbf{A}^T \quad (8)$$

¹Proof: for orthogonal columns of \mathbf{A} we have $\mathbf{A}^T \mathbf{A} = \mathbf{I}$ Multiplying both sides by \mathbf{A}^{-1} we get $(\mathbf{A}^T \mathbf{A}) \mathbf{A}^{-1} = \mathbf{I} \mathbf{A}^{-1}$. Since matrix multiplication is associative, we may also perform: $\mathbf{A}^T (\mathbf{A} \mathbf{A}^{-1}) = \mathbf{A}^T \mathbf{I}$. From definition of an inverse matrix, $\mathbf{A} \mathbf{A}^{-1} = \mathbf{I}$, hence: $\mathbf{A}^T = \mathbf{A}^{-1}$.

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 \mathbf{A} (recall that $\mathbf{A}^T = \mathbf{A}^{-1}$).

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

We shrink the transformation matrix \mathbf{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 \mathbf{X} . We will denote these truncated matrices \mathbf{A}_q and \mathbf{Z}_q respectively. Projecting \mathbf{Z}_q onto the basis \mathbf{A}_q^T will result in an approximation of the original data set:

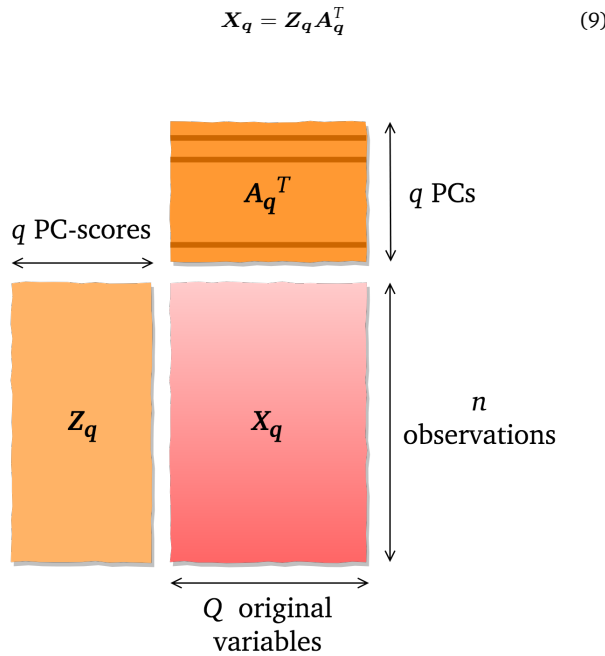


Figure 5: Data approximation with q -PCs.

6 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?

To begin the understanding, let's look back at Figure 3. What we want to achieve with PCA is to diagonalize the new covariance matrix. And we already have produced one diagonal matrix even before computing PC-scores - it was the matrix of eigenvalues $\mathbf{\Lambda}$.

7 PCA Python example

We will now go on to visualizing on a real 2D data set every step of PCA. We will use the PCA function from a Python library `sklearn.decompositions`. The full code can be accessed in the GitHub repository. Here, we will only recall elements of that code to explain what is happening.

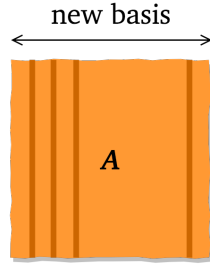


Figure 4: Eigenvalues of the covariance matrix form a new basis.

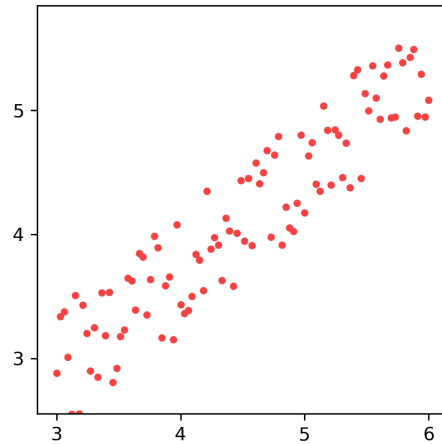


Figure 6: Raw dataset \mathbf{X} .

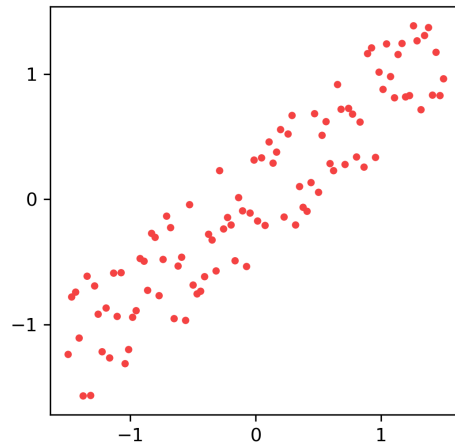


Figure 7: Raw dataset \mathbf{X} , centered.

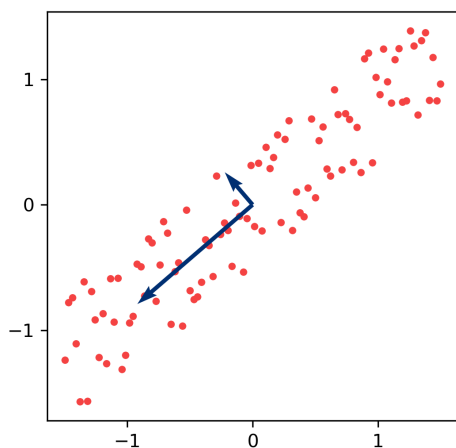


Figure 8: Data set with principal components.

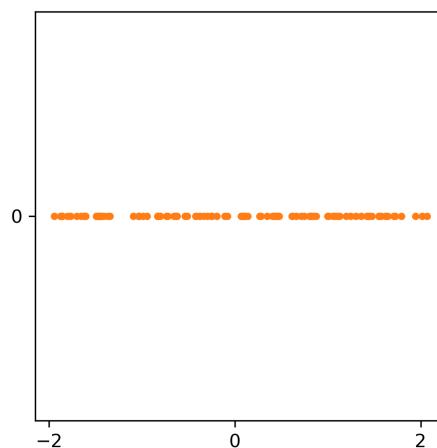


Figure 10: Data projection.

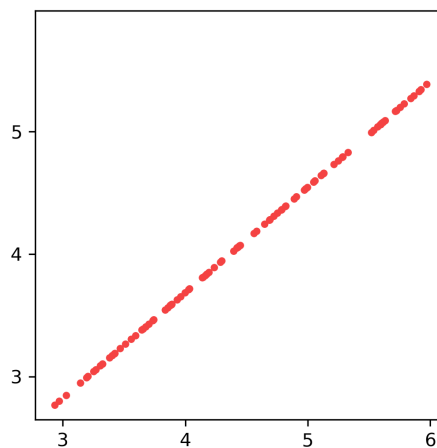


Figure 11: Data approximation with $q = 1$.

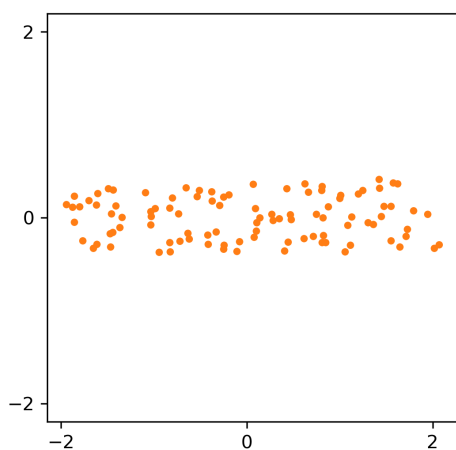


Figure 9: PC-scores.

8 Interesting questions

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

A APP1

B APP2

References

- [1] <https://nl.mathworks.com/help/stats/pca.html>
- [2] Ian T. Jolliffe, *Principal Component Analysis*, Second Edition, 1986
- [3] Gilbert Strang, *Introduction to Linear Algebra*, Fifth Edition, 2016

- [4] Jonathon Shlens, *A Tutorial on Principal Component Analysis*, 2016, <https://arxiv.org/abs/1404.1100>
- [5] <http://people.sju.edu/~pklingsb/dot.cov.pdf>
- [6] J. Edward Jackson, *A User's Guide To Principal Components*, 1991
- [7] Lindsay I. Smith, *A tutorial on Principal Component Analysis*, 2002