

The linear algebra of Principal Component Analysis

(with some Python)

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 a large data set is reduced to maintain only the directions of the largest variance.

These notes are in a way a tutorial on PCA with a deeper focus on linear algebra and statistics aspects that govern this method, compared to other PCA tutorials that I have so far found on the web. I have an aim that the reader will both gain a deep and intuitive understanding as well as an appreciation for PCA after going through these notes. This requires a deep and intuitive understanding of several linear algebra concepts such as the eigendecomposition or matrix operations.

Some initial level of linear algebra knowledge and intuition is hence required. 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	1
2	Data sets for PCA	1
3	Data pre-processing	2
4	Covariance matrix	2
5	PCA workflow	2
6	Why does it need to be eigenvectors?	3
7	PCA Python example	4
8	Another look, or how it all links together	6
9	Interesting questions	6
9.1	Is the covariance matrix of PC-scores a diagonal matrix?	6
A	Properties of interesting matrices	6
A.1	Symmetric matrices	6

B Highlights of statistics

6

1 Motivation for data reduction

To motivate why studying techniques such as PCA and learning to apply them on real data sets is important, there are several questions which stimulated the development of data reduction and data decomposition techniques:

1. Can we send less data but preserve maximum information contained by the data (data compression)?
2. Can we predict what the outcome of another observation will be?
3. Can we build a model in a purely data-driven way - based on data that we have collected (e.g. from experiments or numerics)?
4. Can we make sense of a large, multi-dimensional data set (which cannot be plotted graphically)?
5. Can we extract low-dimensional features that underlay the data?

2 Data sets for PCA

For the rest of this document, we assume that the data set for performing PCA is structured as follows: each column represents one variable V_j , where $j \in \langle 1, Q \rangle$, and the rows correspond to variable observations O_i , where $i \in \langle 1, n \rangle$, (e.g. at different positions in space, or in time). The raw data matrix \mathbf{X} is size $(n \times Q)$.

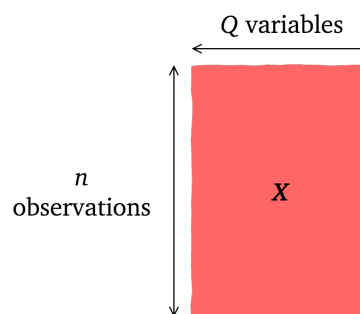


Figure 1: Data matrix for PCA.

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

`pca(X)` Rows of \mathbf{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.

Data centering allows to look at it as variations from some center. Typically we might center each variable by subtracting the mean of this variable's observations (realizations) - we would then look at the data as variations from the mean. This would substitute the original data set with:

$$\mathbf{X}' = \mathbf{X} - \text{mean}(\mathbf{X}) \quad (1)$$

Such centering will shift the "cloud" of data points (which in general is multi-dimensional) to the origin. However, other values for centering are possible.

Scaling, on the other hand, allow us to cancel the effect of unit that a variable might have and treat all variables with equal importance.

$$\mathbf{X}'' = \mathbf{X}' \mathbf{D}^{-1} \quad (2)$$

In the above equation, the matrix \mathbf{D} is a diagonal matrix whose entries are the corresponding scalings. Hence, every column of the matrix \mathbf{X}' gets divided by a corresponding scale from the diagonal of the matrix \mathbf{D} .

You might for instance think about a set of variables from a single experiment, representing temperature in the units of $[K]$ and range from 300-1500 K and associated pressures in the units of $[atm]$ which range from 1-1.1 atm . If we did not scale the data, the largest "spread" or the variance would be found in temperature, since on purely numerical grounds, the range 300-1500 is more significant than the range 1-1.1.

For simplicity, for the rest of this document we assume that \mathbf{X} represents pre-processed data (in place of \mathbf{X}' or \mathbf{X}'').

4 Covariance matrix

The starting point for performing PCA is to compute a *covariance matrix* from the data set. The covariance matrix is given by:

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}^T \mathbf{X} \quad (3)$$

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

We will start by exploring the meaning¹ of $\mathbf{X}^T \mathbf{X}$. Let's look at the graphical representation of this matrix multiplication in Fig. 2.

¹It is worth noticing here, that a matrix constructed as: $\mathbf{S} = \mathbf{A}^T \mathbf{A}$ is square and symmetric and has got at all eigenvalues at least non-negative. It is therefore at least a *positive semidefinite* matrix. In a practical case that we are most interested in, this matrix has got only positive eigenvalues and we call it a *positive definite* matrix.

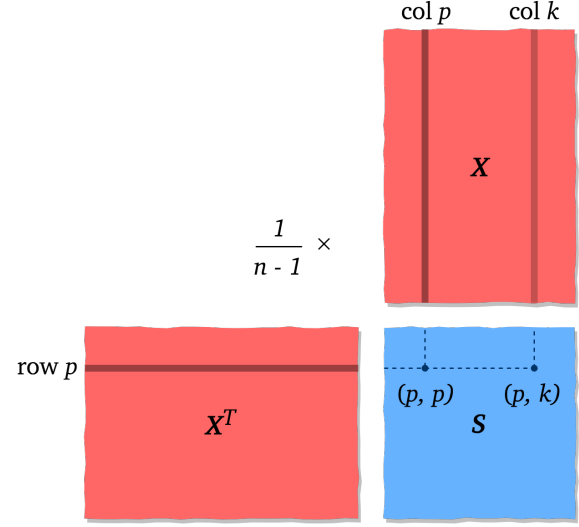


Figure 2: Covariance matrix \mathbf{S} graphical interpretation.

Any given column, say p or k , of the data matrix \mathbf{X} represents all measurements of a single variable V_p (or V_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 come back in the Appendix B.

$$s(p, k) = \frac{1}{n-1} \text{dot}(\mathbf{X}^T(p, :), \mathbf{X}(:, k)) \quad (4)$$

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

$$s(p, p) = \frac{1}{n-1} \text{dot}(\mathbf{X}^T(p, :), \mathbf{X}(:, p)) \quad (5)$$

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 V_p is carried by a variable V_k (and vice versa).

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

$$s(i, j) = \frac{1}{n-1} \sum_{k=1}^n V_i(k) V_j(k) \quad (6)$$

5 PCA workflow

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

$$\text{eig}(\mathbf{S}) = [\mathbf{A}, \mathbf{\Lambda}] \quad (7)$$

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 compo-

nents are orthogonal to each other, and therefore the following important property holds: $\mathbf{A}^T = \mathbf{A}^{-1}$ (proof²). The diagonal matrix $\mathbf{\Lambda}$ of size $(Q \times Q)$ is a matrix of the 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 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:

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

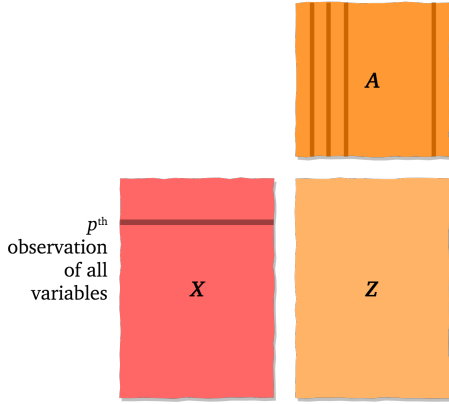


Figure 3: Data transformation to a new basis.

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.(8), 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:

$$\mathbf{X} = \mathbf{Z}\mathbf{A}^T \quad (9)$$

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:

$$\mathbf{X}_q = \mathbf{Z}_q\mathbf{A}_q^T \quad (10)$$

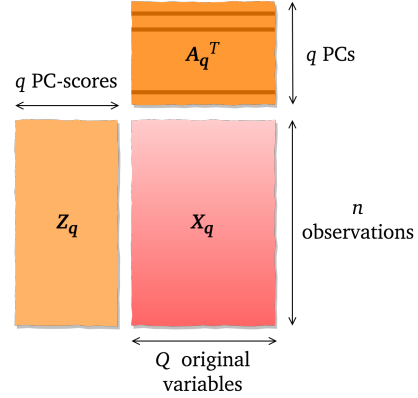


Figure 5: Data approximation with q -PCs.

6 Why does it need to be eigenvectors?

In this section we come back to equation (7) 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 14. 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 $\mathbf{\Lambda}$.

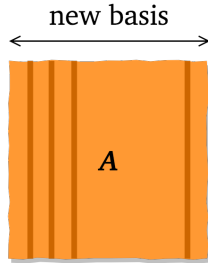


Figure 4: Eigenvalues of the covariance matrix form a new basis.

²Proof: 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: $\mathbf{A}^T(\mathbf{A}\mathbf{A}^{-1}) = \mathbf{A}^{-1}$. From definition of an inverse matrix, $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$, hence: $\mathbf{A}^T = \mathbf{A}^{-1}$.

7 PCA Python example

We will now go on to visualizing on an artificial 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.

We create the **Dataset** (the equivalent of X) as follows:

```
import numpy as np
Np = 100
x = np.linspace(3, 6, Np)
y = 0.8*x + 1*np.random.rand(Np)
Dataset = np.column_stack((x, y))
```

Let's assume that the first column of this data set are realizations of the first variable x and the second column are realizations of the second variable y .

The two variables x and y span two dimensional space but the data set exhibits a low-dimensional structure which is easily visible to the eye just by looking at the Figure 6. We already see that the data seems to be spread along some linear function. Perhaps changing the basis to a basis associated with this linear function will be a more effective representation by our data set? Note here also, that for multidimensional data sets "seeing" such data structures is no longer possible (as it still is in 2-D or 3D). We need to rely on the dimensionality reduction technique that we chose to find this low-dimensional structure for us.

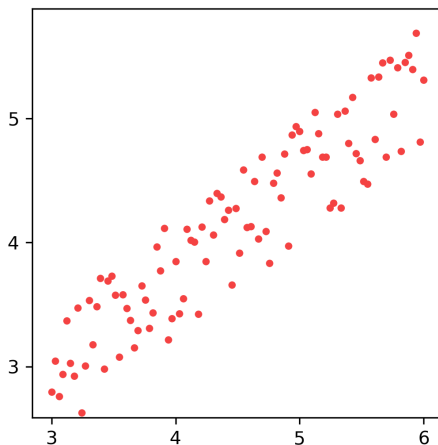


Figure 6: Raw data set.

We center the dataset, which simply moves the center of the cloud of points to the origin of the coordinate system. If necessary, data set would also be scaled to allow for even comparison of the two variables.

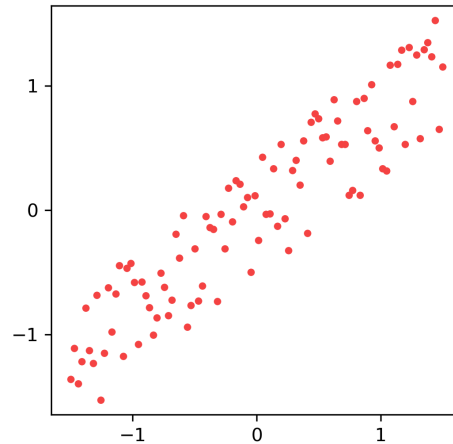


Figure 7: Data set centered.

Next, PCA is performed on the dataset and the eigenvectors (the Principal Components) PCs are found, with the corresponding eigenvalues `eigvals`.

```
from sklearn.decomposition import PCA
pca = PCA(n_components=2)
pca.fit(Dataset)
PCs = pca.components_
eigvals = pca.explained_variance_ratio_
```

In the above code, we create an object `pca` of class `PCA`. We train the model with our **Dataset** using the `fit` function.

The eigenvectors are plotted on the dataset in Figure 8. Their lengths are proportional to their corresponding eigenvalue. Notice that PCA was able to find the direction of the largest variance in the data marked by the direction of the first, longest Principal Component. The second PC is simply perpendicular to the first one.

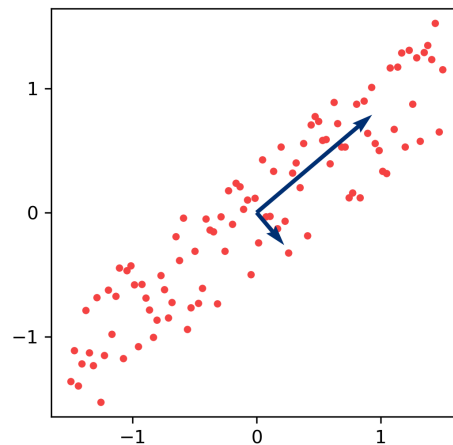


Figure 8: Data set with principal components.

We may now find the PC-scores matrix which represents the "scores"

each point from the data set attains when the new coordinate system was associated with the directions of the Principal Components.

```
scores = pca.transform(Dataset)
```

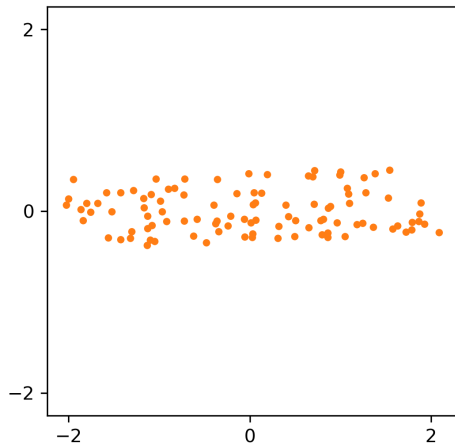


Figure 9: PC-scores.

Next, the PC-scores can be projected on the first PC, to reduce the dimensionality - from two dimensions to one. The data representation from Figure 10 can be viewed as the "scores" each data point would attain if it represented on 1-dimensional structure associated with the first Principal Component.

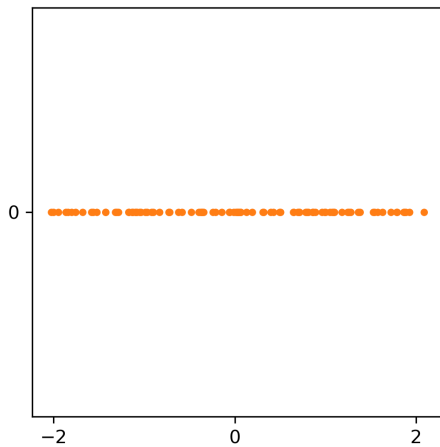


Figure 10: Data projection on low dimension.

The data set can be reconstructed. This represents going back from the 1-dimensional space to the original dimensions. The mean of the data set is added back to undo the data centering.

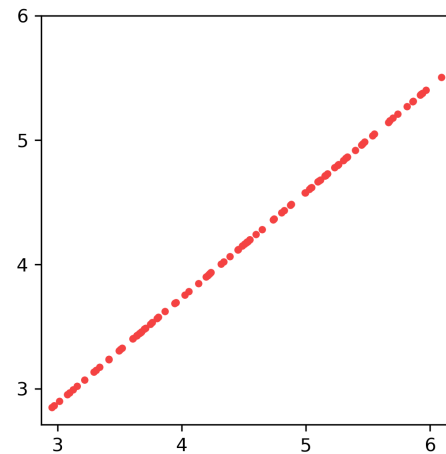


Figure 11: Data approximation with $q = 1$.

8 Another look, or how it all links together

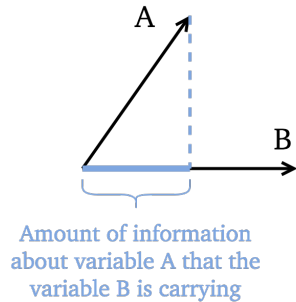


Figure 12: Variable basis might not be optimal.

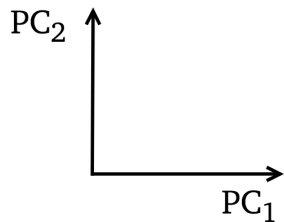


Figure 13: PC-basis is orthogonal.

9 Interesting questions

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

The goal of PCA in terms of a covariance matrix

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

$$S_Z = Z^T Z \quad (11)$$

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 .

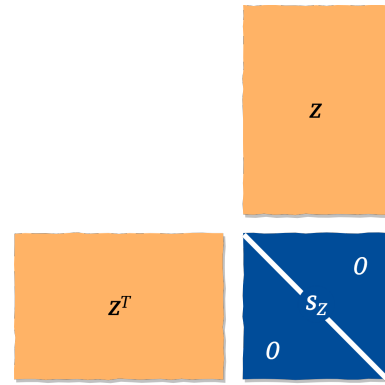


Figure 14: New transformed data set Z and its covariance matrix S_Z .

A Properties of interesting matrices

A.1 Symmetric matrices

The properties of symmetric matrices:

1. Their eigenvectors are orthogonal.
2. Their eigenvalues are real.

B Highlights of statistics

The intuition behind $\frac{1}{1-n}$.

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