

The linear algebra of Principal Component Analysis

(with Python examples)

This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) license.

Kamila Zdybał

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

Preface

Principal Component Analysis (PCA) is a multivariate dimensionality reduction technique in which a high-dimensional data set is projected on the directions of the largest variance. The technique exploits the fact that the original basis to represent the data set might not be an *optimal* one and there might be a redundancy of dimensions. Once a new, orthogonal basis is found, the data can be transformed to that new basis which give certain advantages. The projected data has lower rank and is thus easier to analyze. This in turn can be helpful in a variety of problems such as extracting information, data compression or analysis of structures hidden in the data [10].

These notes are in a way a tutorial on PCA with a deeper focus on linear algebra and statistics aspects governing this method. The aim is to present both a deep and intuitive approach of several linear algebra concepts such as the eigendecomposition or matrix operations that underlay the PCA technique.

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	Nomenclature	1
2	Data sets for PCA	1
3	Data pre-processing	1
4	Covariance matrix	2
4.1	Construction	2
4.2	Properties	2
5	PCA workflow	3
6	Why eigenvectors?	3
7	Eigenvalues as energy	4
8	Python visual examples	4
8.1	Plotting the workflow	4
8.2	Low-rank approximations	6
8.3	Local PCA	7

1 Nomenclature

\mathbf{A}	is a matrix
\mathbf{A}^T	is a matrix transpose
\mathbf{a}	is a vector
a	is a scalar
\mathbf{a}_j	is the j^{th} column of a matrix \mathbf{A} , it is equivalent to $\mathbf{A}(:, j)$
\mathbf{a}_i	is the i^{th} row of a matrix \mathbf{A} , it is equivalent to $\mathbf{A}(i, :)$
$a_{i,j}$	is an element from i^{th} row and j^{th} column of a matrix \mathbf{A} , it is equivalent to $\mathbf{A}(i, j)$

2 Data sets for PCA

For the rest of this document, we assume that the data set for performing PCA is a matrix \mathbf{X} . Each column of \mathbf{X} represents all observations of one variable (feature). Each row of \mathbf{X} corresponds to a single observation (sample) of all variables. Let $i \in \langle 1, n \rangle$ numerate the observations and $j \in \langle 1, Q \rangle$ numerate the variables. The data matrix \mathbf{X} is hence size $(n \times Q)$.¹ This structure can be seen in Fig.1.

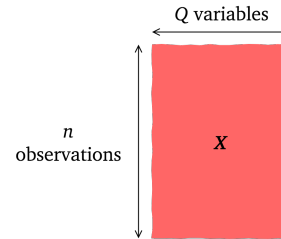


Figure 1: Data matrix for PCA.

3 Data pre-processing

- 1 In data science we are typically given a raw data set which is not centered and not scaled. However, standardizing the data set might be a good idea before applying a data science technique [quote].
- 2 In particular, *centering* allows to look at data as variations from some center. Graphically, centering shifts the center of the *cloud* of data points (which in general is multi-dimensional) to a new, selected center. One of the popular choices is to center each variable by subtracting the mean of this variable's observations - the center will be shifted to the origin². Other centering that could be encountered³ is subtracting the minimum of the variable's observations.
- 3 Centering thus substitutes the original data set with:

¹This is also the data format that is needed for the MATLAB® function `pca` and for the Python function `sklearn.decomposition.PCA`.

²See Fig.7-8.

³For instance in the *Min-Max* standardizing.

$$\mathbf{X}_c = \mathbf{X} - \mathbf{C} \quad (1)$$

Here, a matrix of centers \mathbf{C} is created, for instance by computing the mean of each column of \mathbf{X} . A specific center c_j is then subtracted from a corresponding column \mathbf{X}_j .

Scaling, on the other hand, allows us to cancel the effect of units that variables in our data set might have, and treat all variables with equal importance. It can also remove the effect of various ranges of variables. A centered and scaled data set can then be written as follows:

$$\mathbf{X}_{cs} = \mathbf{X}_c \mathbf{D}^{-1} \quad (2)$$

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}_c gets divided by a corresponding scale from the diagonal of the matrix \mathbf{D} .

To motivate data scaling with an illustrative example, you might 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.

One of the popular scaling methods used in literature is the *auto scaling*⁴ in which each column is divided by the standard deviation of that column. After auto scaling all columns have a standard deviation equal to 1.

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

4 Covariance matrix

The covariance between two random vectors \mathbf{x} and \mathbf{y} is defined as:

$$\text{cov}(\mathbf{x}, \mathbf{y}) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y}) \quad (3)$$

where N is the number of samples (weights of a vector). If we look at any data matrix as a composition of vectors formed by its columns, we may compute the covariances of these vectors and store the result in another matrix, called a *covariance matrix*. This matrix is symmetric due to symmetry: $\text{cov}(\mathbf{x}, \mathbf{y}) = \text{cov}(\mathbf{y}, \mathbf{x})$. The off-diagonal elements have the meaning of covariance of two random vectors and the elements on the diagonal represent variance of each corresponding column, since $\text{cov}(\mathbf{x}, \mathbf{x}) = \text{var}(\mathbf{x})$ where:

$$\text{var}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (4)$$

4.1 Construction

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

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}^T \mathbf{X} \quad (5)$$

and is therefore size $(Q \times Q)$. Notice the similarity with the eq.(3). We will start by exploring the meaning of $\mathbf{X}^T \mathbf{X}$. Let's look at the graphical representation of this matrix multiplication in Figure 2.

⁴also known as *standard scaling* or *z-score* (in that case refers to a combination of centering by the mean value and scaling by the standard deviation).

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 .

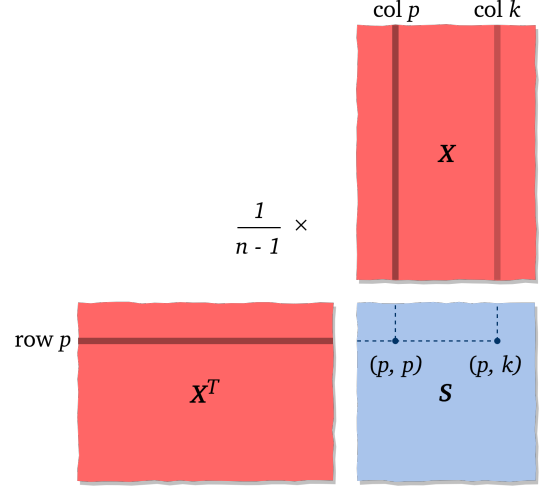


Figure 2: Covariance matrix \mathbf{S} graphical interpretation.

Notice that an element at position (p, k) inside the covariance 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 the k -th column of a matrix \mathbf{X} . There is also the factor $\frac{1}{n-1}$ out front, which we discuss in a closer detail in the Appendix ??.

$$s_{p,k} = \frac{1}{n-1} (\mathbf{x}_p^T \circ \mathbf{x}_k) \quad (6)$$

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

$$s_{p,p} = \frac{1}{n-1} (\mathbf{x}_p^T \circ \mathbf{x}_p) \quad (7)$$

In general, the dot product between two vectors \mathbf{x} and \mathbf{y} represents how much vector \mathbf{x} lays in the direction of vector \mathbf{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).

4.2 Properties

The covariance matrix is a very special matrix and it is worth pointing out some of its interesting properties that the Principal Component Analysis makes use of.

First of all, a matrix constructed as: $\mathbf{S} = \mathbf{C}^T \mathbf{C}$ (where \mathbf{C} is any real matrix) is square and symmetric (the reason for this is easy to see from Fig.2). Apart from that, the eigenvalues of such matrix are real and all are at least non-negative - such matrix is called 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.

Another important property is that the eigenvectors of the covariance matrix are orthogonal, which is a very special thing indeed. You may already see some interesting uses for such eigenvectors. One of which could be: they can form a new, Q -dimensional coordinate system - a

basis. The question thus remains: can such basis be any interesting basis?

5 PCA workflow

As you may have already anticipated, the next step in PCA is to perform the eigendecomposition of the covariance matrix:

$$\text{eig}(S) = [A, \Lambda] \quad (8)$$

The matrix 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: $A^T = A^{-1}$ (proof⁵).

The diagonal matrix Λ of size $(Q \times Q)$ is a matrix of the corresponding eigenvalues.

Given the eigendecomposition of matrix S , we may state that:

$$S = A\Lambda A^T \quad (9)$$

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

$$Z = XA \quad (10)$$

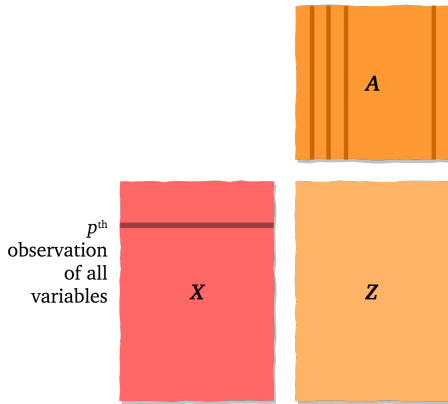


Figure 3: Data transformation to a new basis.

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* matrix, since one may think of every element in this matrix as a "score" that the corresponding element in X gets when represented in the new coordinate system after transformation.

In the matrix multiplication from eq.(10), every variable vector inside X gets transformed by the transformation matrix A and attains new scores in the basis associated with A . The new representation of the old variable is now kept in the matrix 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:

$$X = ZA^T \quad (11)$$

⁵Proof: for orthogonal columns of A we have $A^T A = I$. Multiplying both sides by A^{-1} we get $(A^T A)A^{-1} = IA^{-1}$. Since matrix multiplication is associative, we may also perform: $A^T(AA^{-1}) = A^{-1}$. From definition of an inverse matrix, $AA^{-1} = I$, hence: $A^T = A^{-1}$.

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

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:

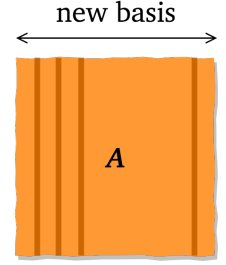


Figure 4: Eigenvalues of the covariance matrix form a new basis.

$$X_q = Z_q A_q^T \quad (12)$$

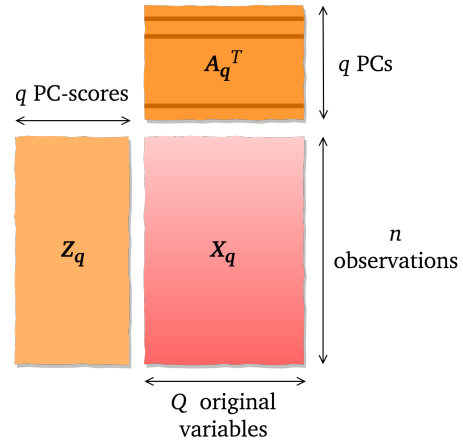


Figure 5: Data approximation with q PCs.

6 Why eigenvectors?

In this section we come back to the eq.(8) and answer the question: why are principal components the eigenvectors of a covariance matrix?

The goal of PCA in terms of a covariance matrix

To begin the understanding, let's look at Fig.6. Principal Component Analysis aims to find a new, transformed data set Z such that if we computed a new covariance matrix in such a way:

$$S_Z = \frac{1}{n-1} Z^T Z \quad (13)$$

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 . This removes the redundancy of information that could have been present in the original data set X ; each column of Z now contributes to a "unique" piece of information that cannot be found in any other column of Z .

In other words, what we want to achieve with PCA is to diagonalize this new covariance matrix S_Z . As a "template" PCA uses a diagonal matrix that is easily available (and which we have already produced), and which is guaranteed to be diagonal - the matrix of eigenvalues Λ . We will now ask ourselves: how do we need to construct Z , so that the matrix $S_Z = \Lambda$?

We find Z through a series of transformations in which we combine eq.(9) with eq.(5):

$$A\Lambda A^T = \frac{1}{n-1} X^T X / A^T \times \quad (14)$$

$$A^T A\Lambda A^T = \frac{1}{n-1} A^T X^T X / \times A \quad (15)$$

$$A^T A\Lambda A^T A = \frac{1}{n-1} A^T X^T X A \quad (16)$$

$$I\Lambda I = \frac{1}{n-1} A^T X^T X A \quad (17)$$

$$S_Z = \Lambda = \frac{1}{n-1} A^T X^T X A \quad (18)$$

It is thus visible, that if we chose $Z = XA$, the product $Z^T Z$ give a diagonal matrix.⁶

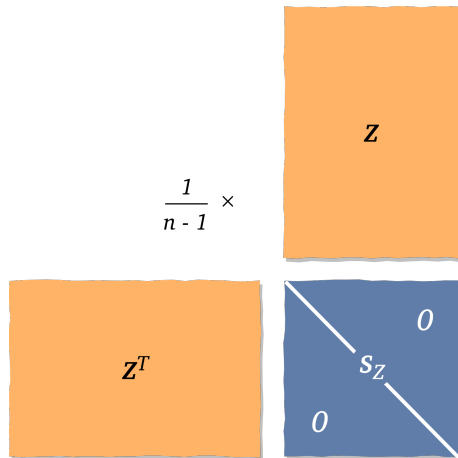


Figure 6: New transformed data set Z and its diagonal covariance matrix S_Z .

There is one more thing owed explaining. We mentioned before that the elements on the diagonal are maximized, so is this also achieved by taking the eigenvalues matrix as a template? It turns out that the answer is yes. PCA achieves both things at the same time: it diagonalizes the new covariance matrix and it makes the diagonal elements maximum possible. Pretty neat, no?

We can now answer the previously posed question: the basis associated with the eigenvectors of the covariance matrix is a very interesting basis! It lets us diagonalize the product $Z^T Z$ and it orders the PC-scores from the most to the least *informative* ones.

⁶Note here that if $Z = XA$, then $Z^T = (XA)^T = A^T X^T$.

7 Eigenvalues as energy

You may sometimes encounter the following statement in the literature: the eigenvalues in PCA represent the energy contributions of each of the Principal Components. The energy here is considered in a physical sense, so that since the PCs are ordered, the first ones are said to carry most of the energy, which in mathematical terms was previously treated in terms of variance. Variance and energy can be encountered as terms used interchangeably. How do we know that the eigenvalues associated to PCs can be viewed that way?

In order to understand this, let's introduce a measure called *inertia*, that can be computed for any vector. We would like that measure to capture the amount of carried by the elements of that vector, on purely numerical grounds.

8 Python visual examples

8.1 Plotting the workflow

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 that are for performing PCA.

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 Fig.7. 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.

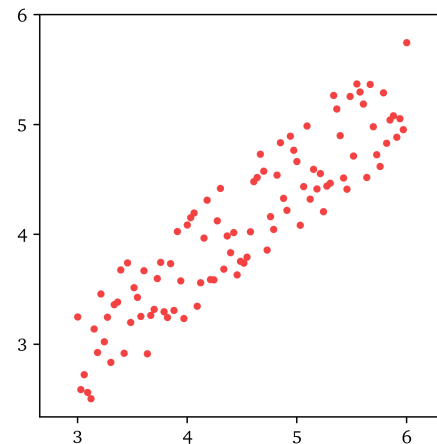


Figure 7: Raw data set.

We center the data set, 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.

```
Dataset_mean = np.mean(Dataset, axis=0)
Dataset_proc = Dataset - Dataset_mean
```

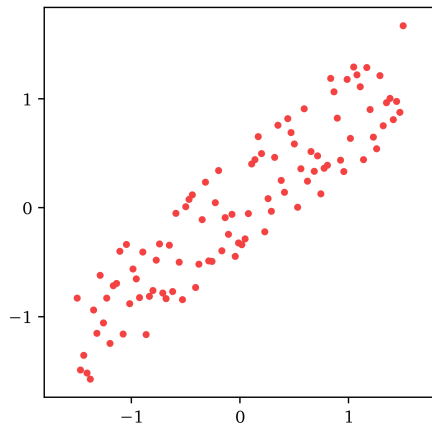


Figure 8: Data set centered.

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()
pca.fit(Dataset)

eigvals = pca.explained_variance_ratio_
PCs = pca.components_
PCscores = pca.transform(Dataset)
```

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 data set in Fig.9. 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 perpendicular to the first one.

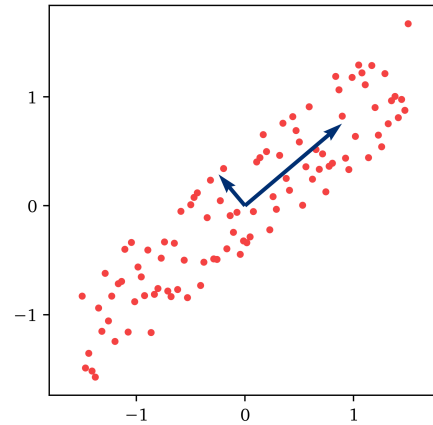


Figure 9: Data set with principal components.

We also compute the transformed data set `PCscores` represented in the basis associated with the obtained PCs.

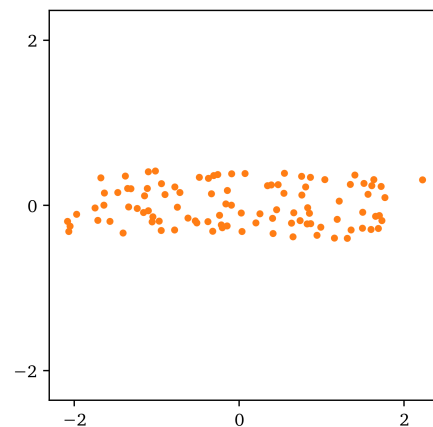


Figure 10: PC-scores.

Next, the PC-scores can be projected on the first PC, to reduce the dimensionality - from two dimensions to one. The data representation from Fig.11 can be viewed as the "scores" each data point would attain if represented on 1-dimensional structure associated with the first Principal Component.

```
q = 1
Dataset_projected = np.dot(Dataset_proc,
np.transpose(pca.components_[:q, :]))
```

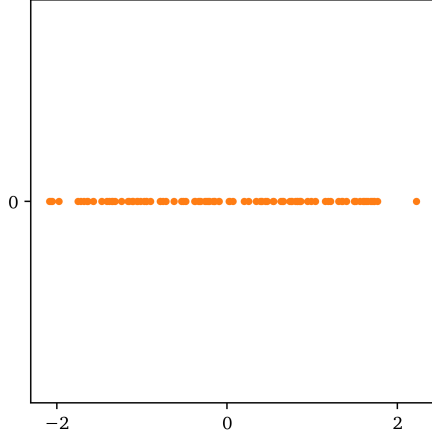


Figure 11: Data projection on lower dimension.

We reconstruct the original data set from the reduced space. 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.

```
Dataset_approx =
np.dot(pca.transform(Dataset)[:,:q],
pca.components_[:q,:]) + Dataset_mean
```

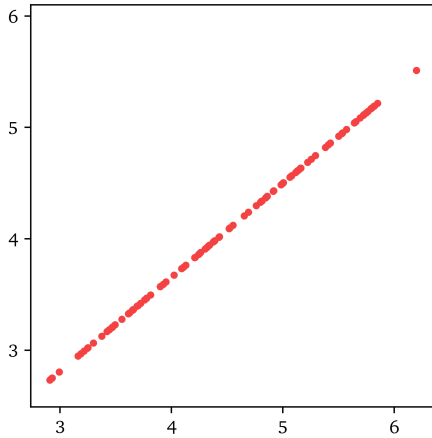


Figure 12: Data approximation with $q = 1$.

8.2 Low-rank approximations

We perform PCA on three artificially generated matrices of size (10×6) : a **random** matrix which is populated by random floats in the range 0-1 and a **semi-structured** and **structured** matrices whose elements are also in the range 0-1 but were populated by the user and have an increasing level of structure that was judged visually. These matrices are presented in Fig.13.

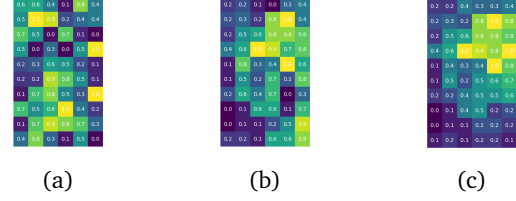


Figure 13: Original data matrices: (a) random matrix, (b) semi-structured matrix, (c) structured matrix.

The visual judgment of the level of imposed "structure" is quite objective in this exercise but the aim was to group the elements of high numerical value (most yellow) in a single region of the matrix and elements of low numerical value (most purple) in other regions of the matrix.

The level of the imposed structure can also be observed quantitatively after performing PCA from the eigenvalue distribution, presented in Fig.14. The structured matrix has got the strongest decaying behaviour which suggest that the matrix can be approximated by relatively low number of modes and hence exhibits the strongest low-rank structure. The first PC is expected to carry 80% of the total variance in the structured data matrix.

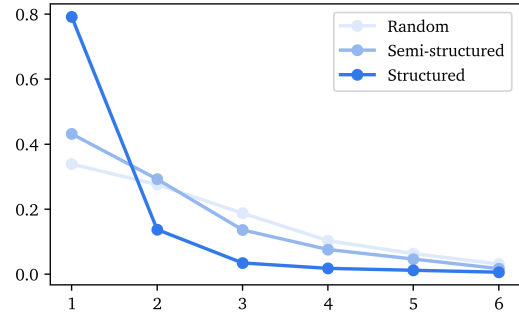


Figure 14: Eigenvalue distribution after performing PCA on original data matrices.

We reconstruct the original data matrices using a certain number q of PC-scores and corresponding PCs. Using the Matlab notation we may write the approximation as:

$$\mathbf{D}_{app} = \text{PC-scores}(:, 1:q) \cdot \text{PCs}^T(1:q, :) + \mathbf{D}_{mean} \quad (19)$$

which is equivalent to eq.(12) and to the Python approximation presented in Sec.8.1.

The above multiplication is presented in three cases in Fig.16. We can see a rank-1 approximation of the original matrices using the 1st Principal Component found by PCA. The vectors (10×1) represent the PC-scores and the vectors (1×6) represent the PCs.

What can be seen visually is that for the semi-structured and structured matrix the low numerical value region is clearly separated. For the random matrix, only few regions of lowest and highest numerical values are recovered in the rank-1 approximation.

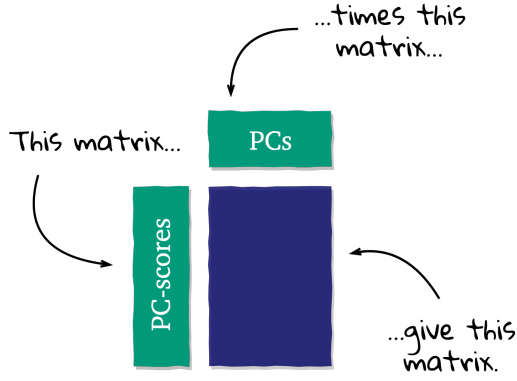


Figure 15: Matrix multiplication from eq.(19) shown graphically.

It is worth noticing here that indeed the obtained matrices are necessarily rank-1, since they are formed as a linear combination of a single vector. This can be seen in two ways: either you may take the vector of PC-scores (10×1) and assume that it forms every column of the (10×6) matrix through multiplying it by the corresponding element from the PC vector. Or, you may assume that the PC is a vector that forms every row of the (10×6) matrix when multiplied by the corresponding element from the PC-scores vector. In either case, the full matrix (10×6) becomes a linear combination of a single vector - hence, it is rank-1.

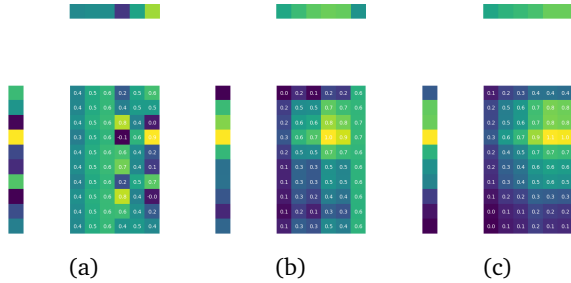


Figure 16: Reconstruction with 1st Principal Component of: (a) random matrix, (b) semi-structured matrix, (c) structured matrix.

In Fig.17 we present a rank-2 approximation, where we maintained two first PCs. Again, the matrices (10×2) represent the PC-scores and the vectors (2×6) represent the PCs.

In the semi-structured and structured matrix, the single matrix elements with highest numerical values (yellow) are already recovered in their actual positions. In the random matrix this is still not the case with rank 2-approximation.

Following the analogous reasoning as for Fig.16 we may notice that the matrix reconstructed with two PCs is a rank-2 matrix.

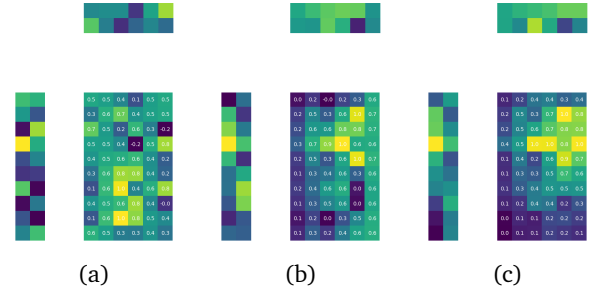


Figure 17: Reconstruction with 2 Principal Components of: (a) random matrix, (b) semi-structured matrix, (c) structured matrix.

The original data matrices are not completely retrieved until all 6 Principal Components and 6 PC-scores are taken into account. In Fig.18 we obtain the final data matrices of rank-6.

The PC-scores are low-dimensional representations of the original data matrix and we return to the original dimensions by the transformation from eq.(12). In the case of taking all 6 PCs, the PC-scores $Z_q = Z$ and the eq.(19) becomes the eq.(11), rather than the approximation from eq.(12).

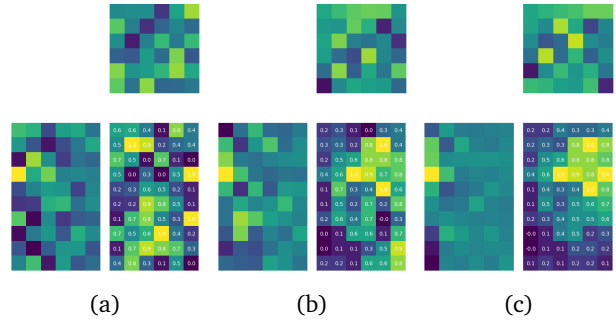


Figure 18: Reconstruction with 6 Principal Components of: (a) random matrix, (b) semi-structured matrix, (c) structured matrix.

Indeed, the multiplications from Fig.18 show full PC-scores matrices transformed to the original 6 dimensions using the inverse (transposition) of the basis matrix made from PCs. The resulting matrix has to be the original data matrix.

8.3 Local PCA

PCA can also be applied on portions of the entire data set, in *local clusters*. This can have certain advantages. Firstly, the reconstruction of the data set from low-rank approximations in the local clusters can allow for further reduction in dimensionality of a non-linear data set. This is due to the fact that clustering creates portions that are locally linear, or that are close to being locally linear. Secondly, the interpretation of Principal Components in local clusters can have more physical meaning, since they become better suited to represent that specific cluster. Below is an example of the general idea of performing Local PCA. A data set composed of two distinct clouds of data can be first partitioned (for instance by techniques such as K-Means clustering) and then PCA is performed separately on both portions of the entire data set. What happens in practice is that the same PCA workflow is applied but on a subset of the data. Notice that the found PCs have different directions.

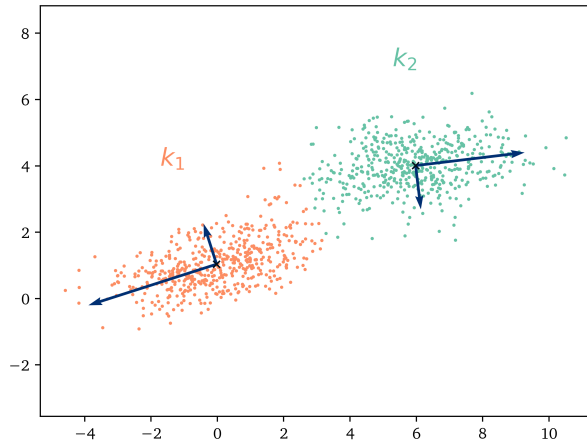


Figure 19: Local PCA.

In the second example presented in the figure below we have a data set with a non-linear behaviour. Performing PCA on the entire data set will result in two PCs that will not represent very accurately the direction of variance of the most "bent" region at the top of the figure. You can observe that after dividing the data into three clusters the local PCs adjust to the direction of variance in each cluster.

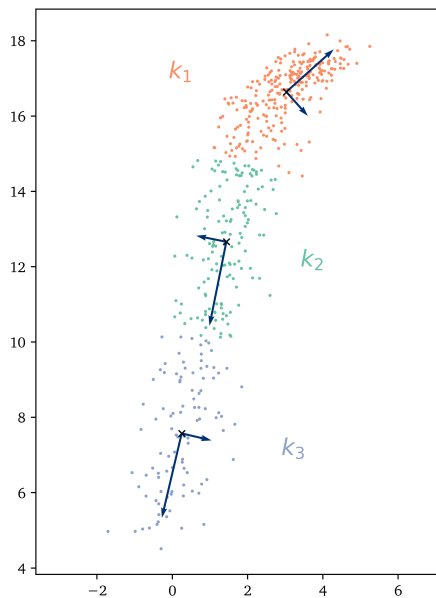


Figure 20: Local PCA.

The reconstruction from the Local PCA is slightly more involved since now we have to account for contributions from each cluster separately.

References

- [1] 3Blue1Brown, *Essence of linear algebra*
- [2] <https://nl.mathworks.com/help/stats/pca.html>
- [3] Ian T. Jolliffe, *Principal Component Analysis*, Second Edition, 1986

- [4] Gilbert Strang, *Introduction to Linear Algebra*, Fifth Edition, 2016
- [5] Jonathon Shlens, *A Tutorial on Principal Component Analysis*, 2016, <https://arxiv.org/abs/1404.1100>
- [6] <http://people.sju.edu/~pklingsb/dot.cov.pdf>
- [7] J. Edward Jackson, *A User's Guide To Principal Components*, 1991
- [8] Lindsay I. Smith, *A tutorial on Principal Component Analysis*, 2002
- [9] Cosma Shalizi, *Course on statistics: The Truth about Principal Components and Factor Analysis*, 2009 <https://www.stat.cmu.edu/~cshalizi/350/>
- [10] Herve Abdi, Lynne J. Williams, *Principal component analysis*, 2010