# Principal Component Analysis (PCA) for Dimensionality Reduction
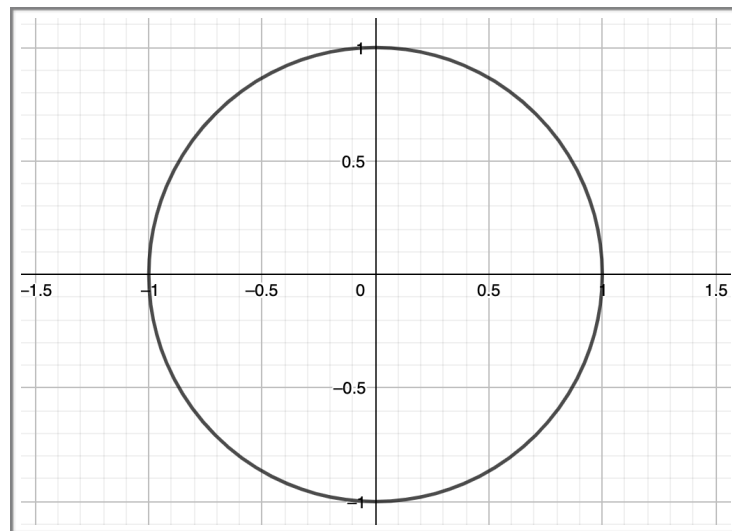
I would like to share with you how I learn topics that are or just seem to be difficult. Whenever I read an article, my first question is what motivated the author to write it and what are the ideas behind the formulas. This is because I am more interested in understanding the concepts than the equations; through concepts I can create simple situations, examples that help me much more than the theory that I usually do not comprehend in the first reading and easily forget it. I keep these examples in my notebooks.
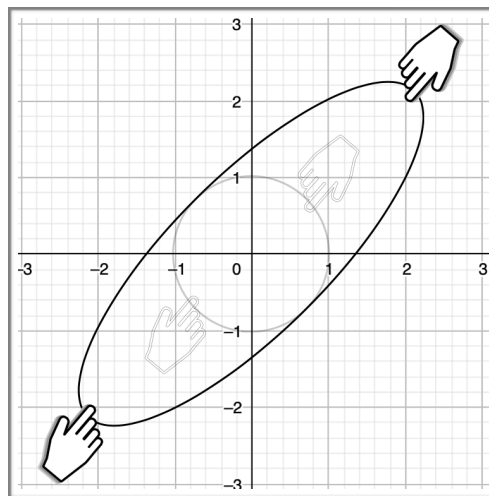


I recently saw a post about how difficult it is to understand the PCA technique. I agree that it will be challenging to get the picture without linear algebra, in particular eigenvalues and eigenvectors. With that in this first article in the Solid Foundation series, I would like to share with you an example that made me understand PCA technique to reduce dimensionality.

tutorsmik@gmail.com

Let's start using the set $C = \{(x, y) \in R^2 : x^2 + y^2 = 1\}$, circle of radius equal to 1.



Now, we are going to use a linear operator $T : R^2 \to R^2$ to stretch the set $C$. That is, transform the circle into an ellipse, as shown below (Any change in a circle is easy to analyze.).



The linear operator that did the transformation above was $T = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$. Now, the interesting part begins, if we find the eigenvalues and eigenvectors of the operator $T$, we can tell about the direction and the factor by which it was stretched.

So, let's first find the eigenvalues $\lambda$. By definition, the eigenvalues are the values that can be found when calculating the determinant of the following relation:

$$|T - \lambda I| = \begin{vmatrix} 2 - \lambda & 1 \\ 1 & 2 - \lambda \end{vmatrix} = 3 - 4\lambda + \lambda^2 = (\lambda - 3)(\lambda - 1)$$

Then, the eigenvalues are 1 and 3. On the other hand, to find the eigenvectors, we need to solve the linear systems below obtained through definition.

$$\begin{pmatrix} 2x & y \\ x & 2y \end{pmatrix} = 3 \begin{pmatrix} x \\ y \end{pmatrix} \text{ and } \begin{pmatrix} 2x & y \\ x & 2y \end{pmatrix} = \begin{pmatrix} x \\ y \end{pmatrix}.$$
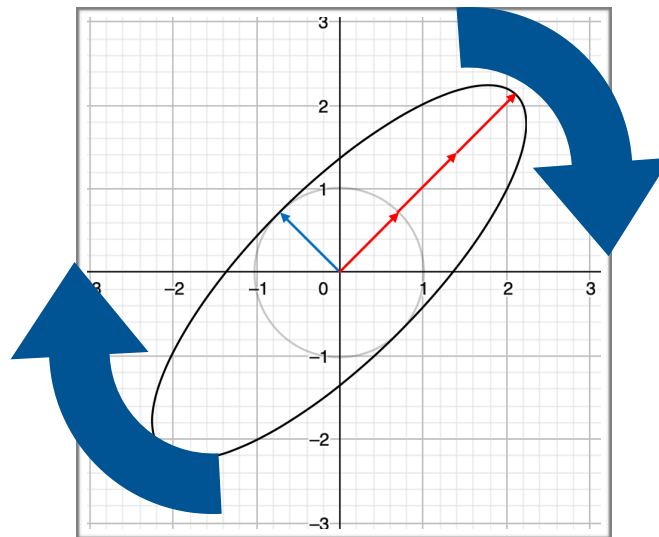
We get that the eigenvectors are $(1,1)$ and $(-1,1)$, normalizing them we get

$$\left( \frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2} \right) \text{ and } \left( -\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2} \right).$$
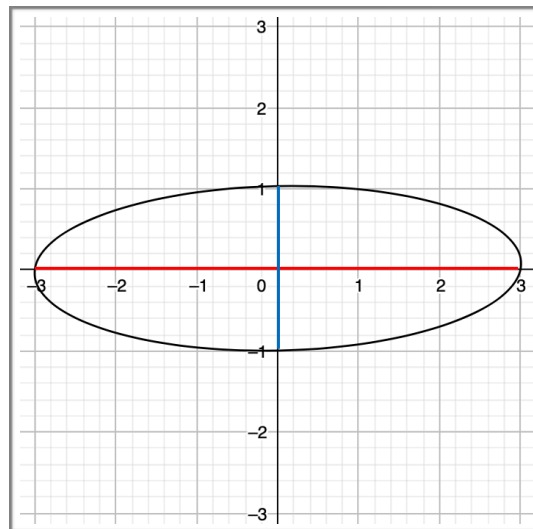
If we create a matrix where the columns are the autonomous vectors, we will see that it is a rotation matrix.

$$V = \begin{pmatrix} \dfrac{\sqrt{2}}{2} & -\dfrac{\sqrt{2}}{2} \\ \dfrac{\sqrt{2}}{2} & \dfrac{\sqrt{2}}{2} \end{pmatrix} = \begin{pmatrix} \cos 45° & -\sin 45° \\ \sin 45° & \cos 45° \end{pmatrix}.$$

Through the matrix above, it will be possible to rotate the figure in order to turn the eigenvectors to the new space axes. And, the eigenvalues can inform us which are the main axes to project our set, reducing the dimension, with less loss of information from the original set. So, if we do $V(T(X))$, we are going to rotate the previous figure 45° clockwise.

Moreover, we can see that variance in the x-axes is bigger than the y-axes. Then, it makes more sense to project the informations into the x-axis, where the range is $[-3,3]$ , 3 times bigger than the range $[-1,1]$.



Well, believe it or not. Now you are familiar with PCA technique for Dimensionality Reduction.

Because:

**PCA** is defined as an orthogonal linear transformation that transforms the data to a new coordinate system such that the greatest variance by some scalar projection of the data comes to lie on the first coordinate (called the first principal component), the second greatest variance on the second coordinate, and so on.(source: Wikipedia)

Okay, I know that the data is generally not as regular as the example shown, much less there is a linear transformation behind it to guide us. It is at this point that the variance-covariance matrix take place to perform the linear operator *T* purpose.

A **variance-covariance matrix** is a square matrix that contains the variances and covariances associated with several variables. The diagonal elements of the matrix contain the variances of the variables and the off-diagonal elements contain the covariances between all possible pairs of variables.(source: Wikipedia)

I will not go deeper in order not to extend our conversation. But I will list comment some commonly used steps:

1.  Mean normalize the data (Balance the data / Create an origin.)
2.  Compute the covariance matrix of your data (Analyzing the stretches).
3.  Compute the eigenvectors and the eigenvalues of your covariance matrix (Find the new axes and consequently the necessary degree to turn the eigenvectors with greater eigenvalues the principal components).
4.  Multiply the first *n* eigenvectors by your normalized data (projecting the dates on the first *n* axes.)

To check further information about the topic (Please visit: https://github.com/Tributino/MIK).

# SOLID FOUNDATION

MASTER THE BASIC CONCEPTS 2 UNDERSTAND COMPLEX IDEAS

tutorsmik@gmail.com