

PCA : Principal Component Analysis

Intuitive explanation

Réda Arab

1 Goal

The goal of PCA is to transform our feature space of dimension p to a new feature space of dimension $k < p$.

It can be used for different purposes : **visualization** ($k = 2, 3$), **compression**, **computational** (lower complexity) or to avoid **overfitting** (however, it is not a "good practice"; it is better to regularize for example, one of the reason is because of the 'fitting' part).

Important use : having new *decorrelated features*.

Principle : We project our data $x_1, x_2, \dots, x_n \in \mathbb{R}^p$ into a new space of dimension $k < p$ such that the *total variance* (i.e the *spread* of our data) is maximized.

Notation : We write $X \in \mathbb{R}^{n \times p}$ with n the size of the sample, p the number of features.

$$X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix} = (X_1 \quad X_2 \quad \dots \quad X_p)$$

In practice, we *center and reduce* our data before applying PCA (it simplifies the calculation and it prevents from having one feature which contains all the variance; the importance of doing this transformation will be clear later).

$$x \rightarrow \frac{x - \hat{x}}{\sigma_x}$$

2 Preliminary Mathematics

2.1 SVD : Singular Value Decomposition

The reader may refer to the Wikipedia article (the schemes are interesting)
https://en.wikipedia.org/wiki/Singular_value_decomposition

Statement : For X a real matrix of $\mathbb{R}^{n \times p}$, it can be written as $X = UDV^T$ with U and V two orthogonal matrices of $\mathbb{R}^{n \times n}$ and $\mathbb{R}^{p \times p}$, and D a matrix with diagonal terms $D_{11} \geq D_{22} \geq \dots \geq D_{mm}$ with $m = \min(n, p)$ and the other elements being null.

For example, if $n > p$, D will be of the form :

$$D = \begin{pmatrix} D_{11} & 0 & \dots & \dots \\ 0 & D_{22} & & \\ \vdots & & \ddots & \\ \vdots & & & D_{mm} \\ 0 & \dots & \dots & 0 \end{pmatrix}$$

Therefore :

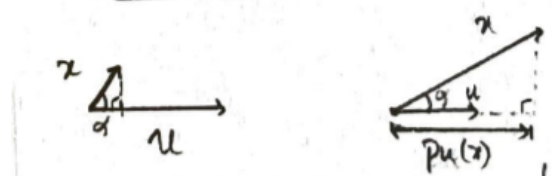
$$X^T X = (UDV^T)^T (UDV^T) = V D^T D V^T$$

$$X^T X = V \Lambda V^T$$

$\Lambda = D^T D$ contains **the eigenvalues** of $X^T X$ (singular values squared) which are real and positive or null.

V contains the **eigenvectors** which constitute an orthonormal basis of \mathbb{R}^p .

2.2 Scalar product and projection on \mathbb{R}^2



We have : $\cos(\alpha) = \frac{p_u(x)}{\|x\|_2}$ where $p_u(x)$ the orthogonal projection of x on the line with u as direction vector.

Therefore $\langle x, u \rangle = \|x\|_2 \cdot \|u\|_2 \cdot \cos(\alpha) = \|u\|_2 \cdot p_u(x)$.

If $\|u\|_2 = 1$, then $\langle x, u \rangle = p_u(x)$

2.3 Characterization of the orthogonal projection on a plane in $\mathbb{R}^n, n \geq 3$

Reminder : Given P a plane, $\Pi_P(x) = \operatorname{argmin}_{y \in P} \|y - x\|_2^2$ where $\Pi_P(x)$ is the orthogonal projection of x on P .

Let us consider an orthonormal basis of $P : (v_1, v_2)$. $\forall y \in P, y = \lambda_1 v_1 + \lambda_2 v_2$.
So :

$$\|y - x\|_2^2 = \|x\|_2^2 - 2 \langle x, y \rangle + \|y\|_2^2$$

$$\|y - x\|_2^2 = \|x\|_2^2 - 2\lambda_1 \langle x, v_1 \rangle - 2\lambda_2 \langle x, v_2 \rangle + \lambda_1^2 + \lambda_2^2$$

Let us consider the function $\phi(\lambda_1, \lambda_2) = -2\lambda_1 \langle x, v_1 \rangle - 2\lambda_2 \langle x, v_2 \rangle + \lambda_1^2 + \lambda_2^2$.

$$\frac{\partial \phi}{\partial \lambda_i} = -2 \langle x, v_i \rangle + 2\lambda_i \quad i = 1, 2$$

$$\frac{\partial^2 \phi}{\partial \lambda_i^2} = 2 \quad i = 1, 2$$

$$\frac{\partial^2 \phi}{\partial \lambda_1 \partial \lambda_2} = 0$$

Thereby, $\frac{\partial^2 \phi}{\partial \lambda^2} \succ 0$ (Hessian matrix positive definite) so the function is strictly convex. The minimum is reached for :

$$\frac{\partial \phi}{\partial \lambda} = 0 \quad \text{i.e}$$

$$\lambda_i = \langle x, v_i \rangle \quad i = 1, 2$$

$$\text{And } \Pi_P(x) = \langle x, v_1 \rangle v_1 + \langle x, v_2 \rangle v_2$$

Generalization

The previous formula can be generalized easily to a projection on a space E_k of dimension k by considering (v_1, \dots, v_k) an orthonormal basis of E_k . We get :

$$\Pi_{E_k}(x) = \sum_{i=1}^k \langle x, v_i \rangle v_i = \sum_{i=1}^k (x^T v_i) v_i$$

3 PCA - Introduction

The **total variance** is defined as :

$$\frac{1}{n} \sum_{i=1}^n \|x_i - \hat{x}\|_2^2 = \sum_{i=1}^n \|x_i\|_2^2 = \sum_{i=1}^p \|X_i\|_2^2$$

as the datapoints are centered (i.e features). The term $\frac{1}{n}$ (or $\frac{1}{n-1}$) is not important for the problem.

→ We are looking for a space of dimension k such that the projection of the x_i on this space will give us new datapoints (with new features associated) which keeps the total variance as large as possible.

As a result, we will have a new matrix Y such that :

$$Y = \begin{pmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_n^T \end{pmatrix} \text{ with } y_i \in \mathbb{R}^k \text{ such that } \sum_{i=1}^n \|y_i\|_2^2 \text{ is "maximal" in the sense defined above.}$$

Example

Let us consider an example when we project our datapoints from a p dimensional feature space to a space of dimension 1 (line). On the schemes, you can see from dimension 2 to 1.

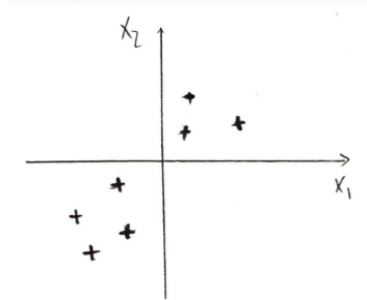


FIGURE 1 – Datapoints in 2D

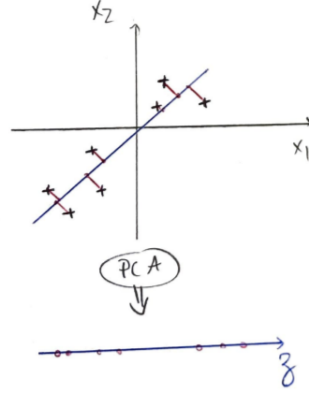


FIGURE 2 – New feature created with a large variance

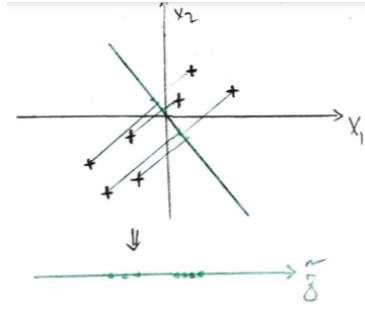


FIGURE 3 – New feature created with a low variance

Consider the case when we have datapoints with 2 features that we want to compress into one feature. So we are looking for a vector u (we can restrict our search to $\|u\|_2 = 1$ as we are looking for a direction) such that the projected datapoints on u ($x_i^T u$, $i = 1, \dots, n$) keeps a maximal total variance. .

i.e $Xu = \begin{pmatrix} x_1^T u \\ x_2^T u \\ \vdots \\ x_n^T u \end{pmatrix}$ and we know that $x_i^T u = p_u(x_i)$ for u with norm equals to 1.

Mathematically, we are looking for $u \in \mathbb{R}^2$ such that :
$$u = \operatorname{argmax}_{v \in \mathbb{R}^2, \|v\|_2=1} \|Xv\|_2^2$$

(total variance maximized).

More generally, if we have a feature space of dimension p :
$$u = \operatorname{argmax}_{v \in \mathbb{R}^p, \|v\|_2=1} \|Xv\|_2^2$$

We have : $\|Xv\|_2^2 = v^T X^T X v = v^T V D^T D V^T v$ with $\|v\|_2 = 1$.

Let us write $a = V^T v$. So $\|a\|_2 = \|V^T v\|_2 = \|v\|_2 = 1$ (as V orthogonal).

$$\rightarrow \|Xv\|_2^2 = a^T D^T D a = \sum_{i=1}^m a_i^2 D_{ii}^2 \leq D_{11}^2 \sum_{i=1}^m a_i^2 = D_{11}^2$$

Therefore $a = (1, 0, \dots, 0)^T$ maximizes $\|Xv\|_2^2$ i.e. $\boxed{v = Va = V_1}$ which is **the eigenvector associated to D_{11}^2** . The direction which maximizes $\|Xv\|_2^2$ is V_1 and in this case $\|XV_1\|_2^2 = \|D_{11}U_1\|_2^2 = D_{11}^2$

Remarks

- The projected datapoints $x_i^T v$ are centered ($\sum_{i=1}^n x_i^T v = (\sum_{i=1}^n x_i^T)v = 0$), so the total variance computed is $\sum_{i=1}^n (x_i^T v)^2 = \|Xv\|_2^2$
- PCA and linear regression are different.

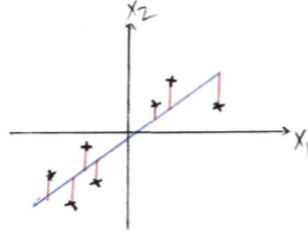


FIGURE 4 – Linear regression

The quantities we want to minimize in each case (the red lines in figure 2 4) are different.

4 Formulation of the problem

The problem can be summarized as :

$$\boxed{\operatorname{argmax}_{(v_1, v_2, \dots, v_k) \in \mathbb{R}^p} \sum_{i=1}^k (\|Xv_i\|_2^2) \quad \text{such that} \quad v_i^T v_j = \delta_{i,j}}$$

Indeed, for a space E_k of dimension k and an orthonormal basis (v_1, \dots, v_k) :

$$\Pi_{E_k}(x_i) = \sum_{l=1}^k (x_i^T v_l) v_l \quad \text{and} \quad \sum_{i=1}^n \|\Pi_{E_k}(x_i)\|_2^2 = \sum_{i=1}^n \sum_{l=1}^k (x_i^T v_l)^2$$

In this new space, we can write **the coordinates of the projected** x_i according to the orthonormal basis as :

$$Y = \begin{pmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_n^T \end{pmatrix} = \begin{pmatrix} x_1^T v_1 & x_1^T v_2 & \dots & x_1^T v_k \\ \vdots & \vdots & \dots & \vdots \\ x_n^T v_1 & x_n^T v_2 & \dots & x_n^T v_k \end{pmatrix} = (Xv_1 \quad Xv_2 \quad \dots \quad Xv_k)$$

The total variance in this case is : $\sum_{i=1}^n \sum_{l=1}^k (x_i^T v_l)^2 = \sum_{i=1}^k \|Xv_i\|_2^2$

5 Problem resolution

We can solve the problem step by step, *direction by direction* (the problem is separable, it is a sum).

We can start with V_1 as a first new feature (*cf PCA -Introduction*) then, we are looking for a vector $v^{(2)}$ such that $\|Xv^{(2)}\|_2^2$ is maximized and $\|v^{(2)}\|_2 = 1$ and $(v^{(2)})^T V_1 = 0$.

We can easily show that $\boxed{v^{(2)} = V_2}$.

Recursively, we are looking a vector $v^{(j)}$ for $j = 2, \dots, k$, defined as below :

$$\boxed{\begin{array}{l} v^{(j)} \text{ maximise } \|Xv\|_2 \text{ over } v \in \mathbb{R}^p \text{ with the constraints} \\ \|v\|_2 = 1 \text{ and } (v^{(l)})^T v = 0 \text{ for all } l < j \end{array}}$$

We get $\boxed{(V_1, \dots, V_k)}$, **the eigenvectors** $X^T X$ associated to the eigenvalues $D_{11}^2 \geq D_{22}^2 \geq \dots \geq D_{kk}^2 \geq 0$.

Therefore, the new datapoints are :

$$Y = \begin{pmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_n^T \end{pmatrix} = (XV_1 \quad XV_2 \quad \dots \quad XV_k) = (D_{11}U_1 \quad D_{22}U_2 \quad \dots \quad D_{kk}U_k)$$

The total variance of our new projected datapoints is : $\sum_{l=1}^k ||D_{ll}U_l||_2^2 = \sum_{l=1}^k D_{ll}^2$

The initial total variance, before the projection, is : $\sum_{l=1}^m D_{ll}^2$

Choice of k

We can define a way of choosing k (number of new features) as below.

$$\min\{k \mid \frac{\sum_{l=1}^k D_{ll}^2}{\sum_{l=1}^m D_{ll}^2} \geq a\} \quad \text{with } a = 0.80, 0.90 \text{ or } 0.95 \text{ for example}$$

Final remarks

- Our new datapoints are **centered** as shown in a previous section.
- The new features are **decorrelated** :

$$i \neq j, (XV_i)^T (XV_j) = V_i^T V \Lambda V^T V_j = 0.$$

- The new features are less interpretable, explainable.
- The general idea of this method of dimensionality reduction : **keeping the spread of our data as large as possible** .