

## Principal Component Analysis:

→ What is PCA and why it is needed?

As we all know, including too many features during the training phase confuses the model which ultimately leads to overfitting of the model.

Now, to solve this overfitting problem, we generally do the dimensionality reduction. There are lots of ways of doing the dimensionality reduction to overcome the overfitting; 'PCA' is one of those.

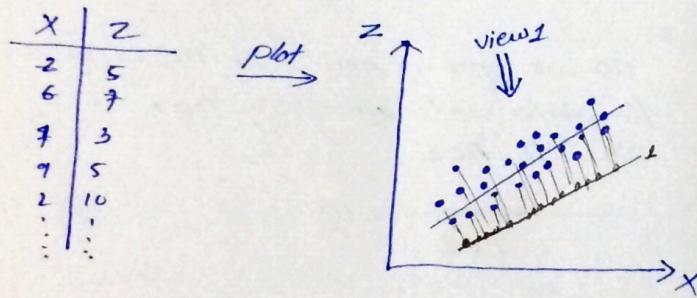
→ How do we reduce the dimensions in PCA?

It is all about the views.

Suppose we have a dataset, which contains 2 independent variables.

Note that we do not consider PCA for dependent variable.

Here for the sake of simplicity, we are taking only 2 features/variables.



Now, how to reduce the dimensions (X & Z) to one dimension.

# Note: We will see all the mathematical details shortly. For now, only concentrate only on concepts.

→ So if we look at the above graph from 'view 1', then we get a line which has few datapoints projected on it.

oooooooooooo

What we did here is basically, we projected our datapoints, which were plotted (distributed) across the dimensions (feature) and reduced the dimensions from two to one, which will eventually reduce overfitting as there is one less dimension now.

→ So by projecting these data points:

- ① we reduced the dimensions
  - ② we reduced the complexities
  - ③ we reduced the confusion for the model.

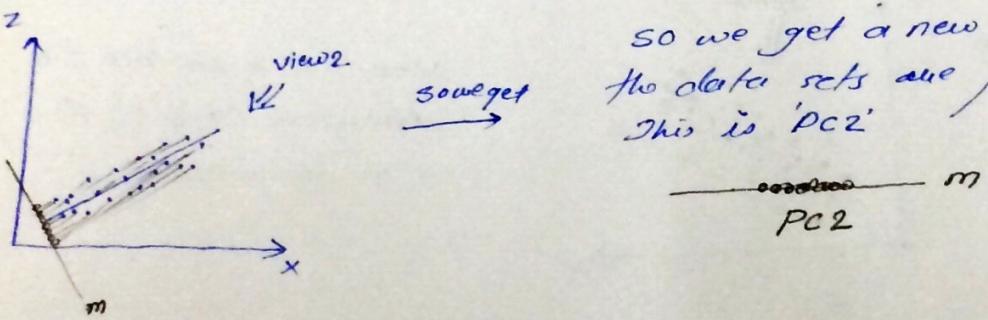
What is the actual motive of PCA.

- PC<sub>1</sub>: Now the component or new line (e) that we got after projecting the data points is called as a 'Principal Component'. This can be called as Principal Component 1 if it explains the most of the variability.

We will later see, how to choose a PC as  $PC_1$  or  $PC_2$  or  $PC_3$  or so on mathematically. For now assume this is ' $PC_1$ '.

~~60000 000 000 00~~ (l)  
PC1

Another PC can be generated. Let's look the same original graph from view 2.



so we get a new line ( $m$ ) on which  
the data sets are projected.  
This is 'PC2'

Hence we have generated two principal components here.  
From this we can conclude that:

From this we can conclude that:

No. of principal components  $\leq$  No. of features/variables/attributes  
↳ axiom 1.

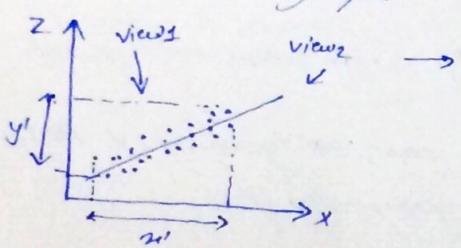
Note: Here in this example, we had only 2 attributes. However, if we have more than 2 attributes, we use the same approach and generate multiple principal components.

11<sup>y</sup>, we can have PC<sub>3</sub>, PC<sub>4</sub>, PC<sub>5</sub>...etc. It all depends on number of features as per attrmnt.

### Priority of Principal Components:

- PC<sub>1</sub> has the highest priority as it explains/contributes maximum to the model.
- It is followed up by PC<sub>2</sub>, PC<sub>3</sub>, PC<sub>4</sub> etc..
- Which component to assign as PC<sub>1</sub>?

Conceptually: We check which feature explains the maximum variance, which component/attribute has values the most. (This is why standardization/normalization of your data is important).  
e.g. see the below graph:



If we see here, the 'x' variable/feature is more widely distributed. It varies more than 'z' feature.  
 $x > y > z$

Hence, whatever PC we will get by projecting datapoints from 'views' will be our 'PC<sub>1</sub>'. Hence, we are compromising with our 'z' attribute as it is varying less when compared to 'x' attribute.

11<sup>y</sup>, whatever component we get from 'views' will be our PC<sub>2</sub> as it is less varying.

Mathematically: We calculate the values of  $\lambda$  (lambda). It is also called as 'Eigen value'. Which ever value of  $\lambda$  among all lambdas ( $\lambda_1, \lambda_2, \lambda_3$  etc...) is greatest, the corresponding 'eigen vectors' are assigned as 'PC<sub>1</sub>', then the second largest as PC<sub>2</sub> and so on....

How many PCs should we have?

- well, whatever number of PCs explains the 75-80% of the variance, those many PCs are taken as optimised number of Principal components.
- You can check the variance by plotting 'scree plot'.
- **Orthogonal property:** The PCs that we get after PCA, there should be an orthogonal property between them.

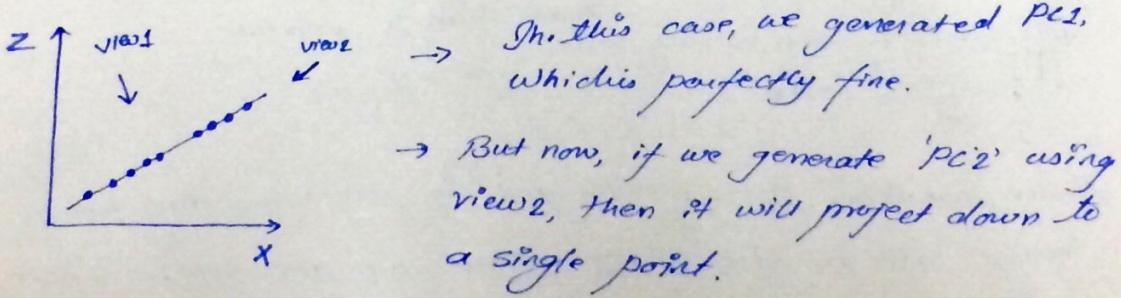
Basically, each 'PC' should be orthogonal to the other.

In simple words, PC1 should be totally independent of PC2, and vice versa.

Each PC should be independent of each other.

See, there is a logical reasoning behind it, which we shall see, using a graph.

Let's say, we have  $x$  &  $z$  perfectly correlated with each other.



Hence, the sum of squared distance, on which basically the PCA works will be zero and hence our PC2 will be of no use.

Basically, because of this reason only, we say the Principal Components should be independent of each other.

# PCA - How it works? - step by step.

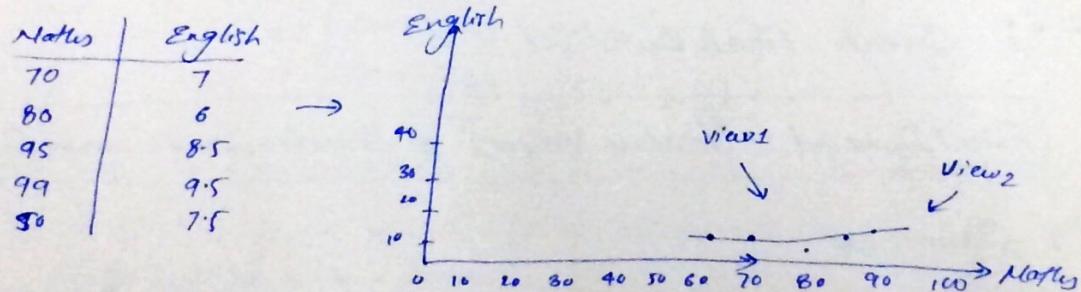
## Step 1: Standardization of your features (inputs)

PCA is quite sensitive regarding the variances of the initial variables. If there are large differences between the range of initial variables, those variables with larger ranges will dominate over those with small ranges, which will lead to bias results.

Hence data transformation to common scale is necessary.

Let me give you an example: let say we have maths subject marks out of 100 and English subject marks out of 10. Both subjects are important and contributes almost equally to the modes.

### Without standardization:



If we see from 'view 1', we can clearly see that maths is contributing a lot towards variance, while English is almost negligible (see on English axis, it shows very less variation, view 2). However it is not the case in reality.

Both the subjects are equally contributing to the variation. Hence standardization of data is important.

## Step 2: Calculate

Step 1: Standardization of features (inputs)

Step 2: Covariance Matrix Computation

Step 3: Compute Eigenvalues & Eigen vectors of the co-variance matrix

Assign the PCs accordingly

Step 4: Keep Important PCs and discard the less significant ones

Step 5: Create the final feature vector by taking eigen vectors kept in above steps (according to PCs)

Step 6: Create Final Data set.

$$\boxed{\text{Final Data Set} = (\text{Feature Vector})^T * (\text{Standardized Original Dataset})^T}$$

T → transpose.

Note: Standardised Original Dataset does not contain target variable.

PCA - Example: For the sake of simplicity, we are taking only 2-d features ( $X$  &  $Y$ ). However the same steps and concepts are applied for more than 2 features as well:

① Data we have is :

	X	Y	
2.5	2.4		→ total 10 observations
0.5	0.7		→ attributes are $X$ & $Y$
2.2	2.9		→ <u>Standardization</u> , we are skipping here, but if you do it programmatically, make sure you do the standardization.
1.9	2.2		
3.1	3.0		
2.3	2.7		
2	1.6		
1	1.1		mean ( $X$ ) i.e. $\bar{X} = 1.81$ & $\bar{Y} = 1.91$
1.5	1.6		
1.1	0.9		
	$\bar{X}$	$\bar{Y}$	
1.81	1.91		

Step 2: Calculating co-variance matrix.

$$C = \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix}$$

Note: We had only 2 dimensions of data here, hence we have  $2 \times 2$  covariance matrix, If we have higher dimension data, then we have higher dimension matrix.

e.g. If we have 3 dimensions data, then,  $(x, y, z)$ , the covariance matrix would have been:

$$C = \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) & \text{cov}(x, z) \\ \text{cov}(y, x) & \text{cov}(y, y) & \text{cov}(y, z) \\ \text{cov}(z, x) & \text{cov}(z, y) & \text{cov}(z, z) \end{bmatrix}$$

$$\boxed{\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1)}}$$

Applying the formula, we get:

$$\begin{aligned} \text{cov}(x, x) &= 0.6165 & \text{cov}(y, x) &= 0.6154 \\ \text{cov}(x, y) &= 0.6154 & \text{cov}(y, y) &= 0.7165 \end{aligned}$$

∴ Covariance matrix is:

$$C = \begin{bmatrix} 0.6165 & 0.6154 \\ 0.6154 & 0.7165 \end{bmatrix}$$

Step 3: Compute Eigen value & Eigen vectors.

3.1 → Compute Eigen values: To get the eigen values, use the below form  $[C - \lambda I = 0]$ , then solve for the determinants.

where,

$C$  = covariance matrix

$\lambda$  = constant / eigen value

$I$  = Identity matrix. (Its dimensions depends on ' $C$ ')

If ' $C$ ' is  $3 \times 3$  matrix, ' $I$ ' is also  $3 \times 3$  matrix.

$$\therefore \begin{bmatrix} 0.6165 & 0.6154 \\ 0.6154 & 0.7165 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

$$\therefore \begin{bmatrix} 0.6165 - \lambda & 0.6154 - 0 \\ 0.6154 - 0 & 0.7165 - \lambda \end{bmatrix} = 0$$

Solve for determinants, this will form we get quadratic equation in terms of lambda:

$$(0.6165 - \lambda) * (0.7165 - \lambda) - 0.6154 * 0.6154 = 0$$

Here the above equation is quadratic, hence it has 2 roots.  
after solving we get the roots as:

$$\boxed{\lambda_1 = 0.4908}$$

$$\boxed{\lambda_2 = 1.2840}$$

### 3.2 → Compute Eigen vectors:

⑨ Substitute  $\lambda_1$  &  $\lambda_2$  in: ~~(Ax=0)~~ (Ax=0)

$$\begin{bmatrix} 0.6165 - \lambda & 0.6154 \\ 0.6154 & 0.7165 - \lambda \end{bmatrix} \begin{bmatrix} x_i \\ y_i \end{bmatrix} = 0$$

& solve for  $x_i$ .

Hence  $\lambda_1 = 0.4908$

$$\begin{bmatrix} 0.6165 - 0.4908 & 0.6154 \\ 0.6154 & 0.7165 - 0.4908 \end{bmatrix} \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = 0$$

$$\therefore 0.1257 x_1 + 0.6154 y_1 = 0$$

$$0.6154 x_1 + 0.2257 y_1 = 0$$

∴ we get  $\boxed{x_1 = -0.735 \quad y_1 = 0.677}$  ---- vector ①

↳ Eigen vector for  $\lambda_1$  ∴  $v_1 = \begin{bmatrix} -0.735 \\ 0.677 \end{bmatrix}$

11) for  $\lambda_2 = 1.2840$

$$\begin{bmatrix} -0.6675 & 0.6154 \\ 0.6154 & -0.5675 \end{bmatrix} \begin{bmatrix} x_2 \\ y_2 \end{bmatrix} = 0$$

$$\boxed{x_2 = -0.677 \quad y_2 = -0.735} \rightarrow \text{Eigen vectors for } \lambda_2 \\ \dots \text{vector ②} \quad v_2 = \begin{bmatrix} -0.677 \\ -0.735 \end{bmatrix}$$

Note: Here we had 2 eigen values. hence we got 2 eigen vectors  
if we had 3 eigen value we should have got 3 eigen vectors  
and so on..

### Assigning PC1 & PC2 :

Here  $\lambda_2 = 1.2840$  &  $\lambda_1 = 0.4908$   $\boxed{x_2 > x_1}$

∴ corresponding eigen vector of  $\lambda_2$  will be our PC1.

$$\therefore \boxed{\begin{bmatrix} -0.678 & -0.73 \\ -0.735 & 0.677 \end{bmatrix} \rightarrow \text{PC1}}$$

$$\boxed{\begin{bmatrix} -0.735 & 0.677 \end{bmatrix} \rightarrow \text{PC2}}$$

Step 4&5: keeping both here. However we can discard PC2 as well.

If we keep both:

$$\text{Feature vector} = \begin{bmatrix} v_{11} \\ -0.735 \\ 0.677 \end{bmatrix} \quad \begin{bmatrix} v_{12} \\ -0.677 \\ -0.735 \end{bmatrix}$$

If we discard PC2:  $(v_1)$

$$\begin{bmatrix} -0.735 \\ 0.677 \end{bmatrix}$$

However we keep both here.

Step 6:

$$\boxed{\text{FinalDataset} = (\text{Feature vector})^T * [\text{Standardization(Data)}]^T}$$