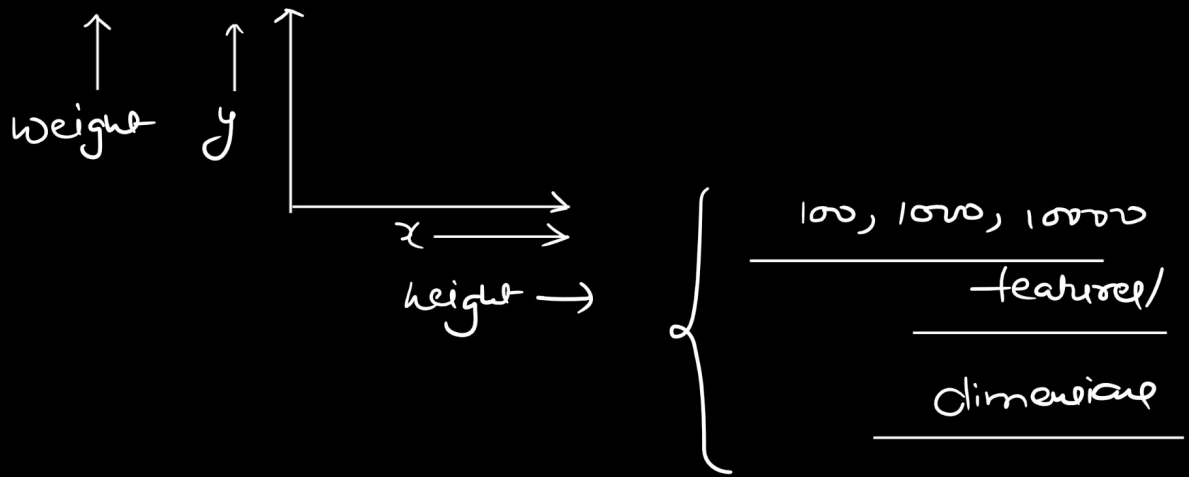


# Dimensionality Reduction

1 Dimensional → feature /  
Column



## ① feature selection

20 features → ① feature importance  
{ DT, Random forest }  
10 features → highest feature  
importance value  
(feature)

② Lasso Regularization  
technique

③ Correlation  
heatmap

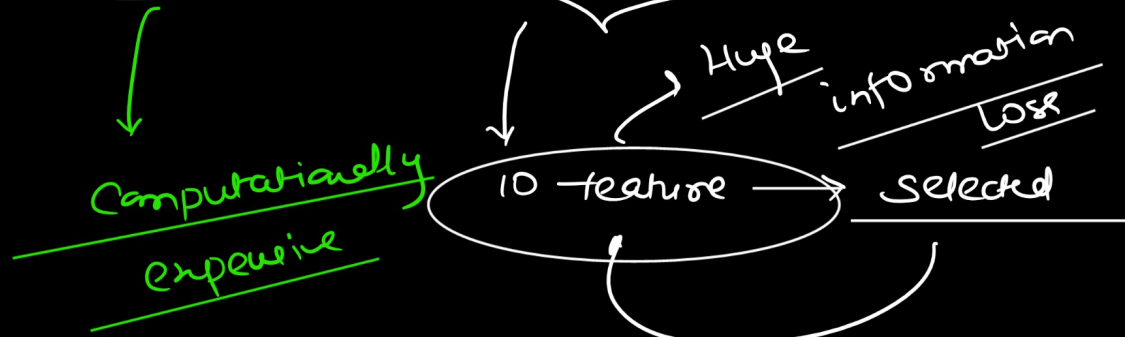
## ② feature Extraction / Projection

Redesign new features 99%

or

Dimensionality Reduction

10 → 90%  
95%  
80%  
{  $f_1, f_2, f_3, f_4, f_5, \dots, f_{100}$  }



EDA → Data Preprocessing

Unsupervised

ML Technique

PCA

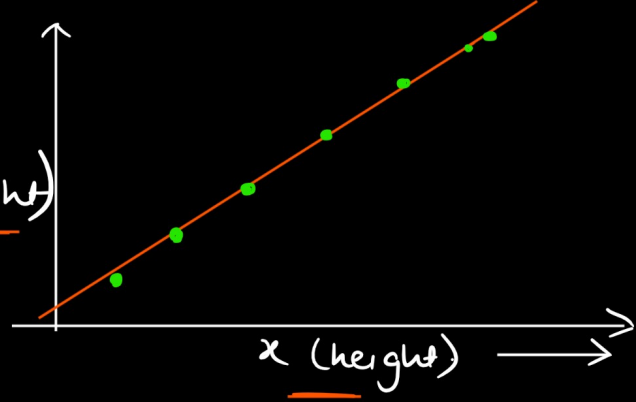
→ Principal Component

Analysis

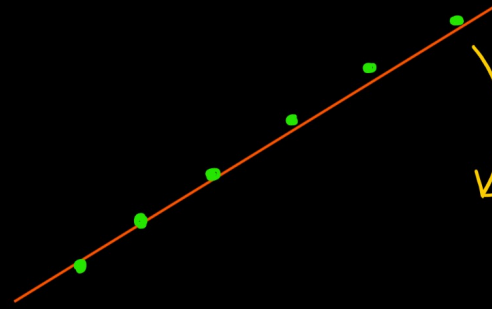
↓  
Eigen values & Eigen vectors

2D data

y (weight)



1D data



eigen values & eigen vectors

Rotate axis

{



Step-by-step Guide of PCA

to get/extract new features

from the original features

sklearn.preprocessing

- 1) Standardise the data → to maintain the range of all the features

Age	Standardise	→ $z = \frac{x_i - \mu}{\sigma}$
23		
25		
47		
48		
72		
64		

- 2) Covariance matrix ←

A covariance matrix is a square, symmetric matrix in statistics that describes the variance and covariance between different variables in a dataset. The diagonal elements of the matrix contain the variances of the individual variables, indicating their individual spread or variability, while the off-diagonal elements represent the covariances between pairs of variables, showing how they vary together.

Covariance matrix ⇒ matrix @ matrix <sup>matrix multiplication</sup>  
num of - 1  
records

Covariance provides the measure of strength of correlation between two variable or more set of variables. The covariance matrix element  $C_{ij}$  is the covariance of  $x_i$  and  $x_j$ . The element  $C_{ii}$  is the variance of  $x_i$ .

If  $\text{COV}(x_i, x_j) = 0$  then variables are uncorrelated

If  $\text{COV}(x_i, x_j) > 0$  then variables positively correlated

If  $\text{COV}(x_i, x_j) < 0$  then variables negatively correlated

*Handwritten notes:*  
 $-\infty$  to  $+\infty$   
 $f_1 \uparrow$  (Project)  $f_2 \uparrow$  (better expertise)  
 $f_1 \uparrow$  (TV series)  $f_2 \downarrow$  (lower expertise)

## Difference b/w Correlation Coefficient vs Covariance ??

Covariance	Correlation
Covariance is a measure of how much two random variables vary together	Correlation is a statistical measure that indicates how strongly two variables are related.
Involves the relationship between two variables or data sets	Involves the relationship between multiple variables as well
Lie between -infinity and +infinity	Lie between -1 and +1
Measure of correlation	Scaled version of covariance
Provides direction of relationship	Provides direction and strength of relationship
Dependent on scale of variable	Independent on scale of variable
Have dimensions	Dimensionless

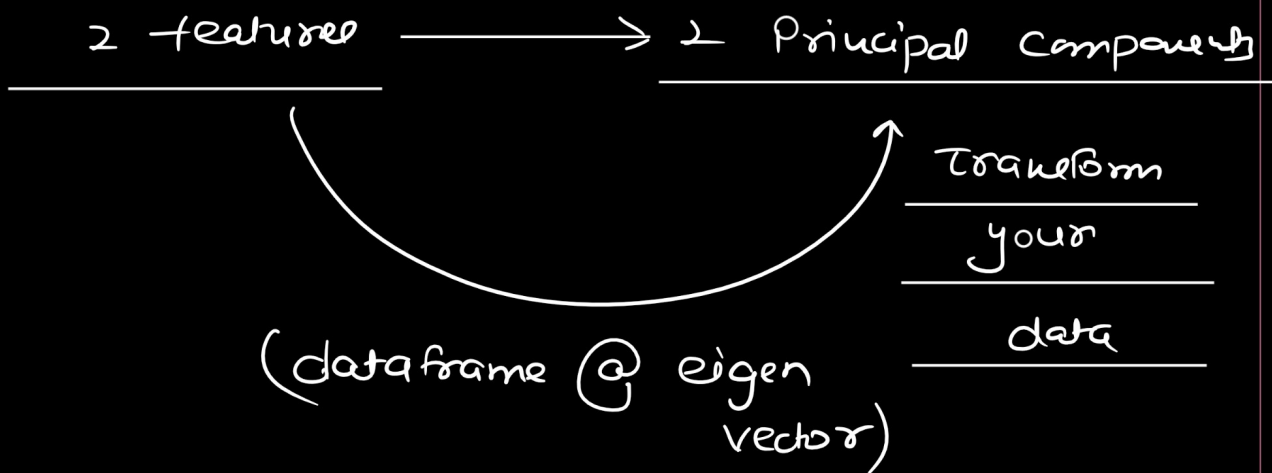
They key difference is that Covariance shows the direction of the relationship between variables, while

## ③ Eigen Value & Eigenvectors using Covariance matrix

2 feature  $\rightarrow$  2 eigen value  
 $\downarrow$   
2 eigen vector

④ To evaluate the new principal components

$$PC1 > PC2 > PC3 \dots$$



### Properties of Principal components

- ① → uncorrelated Each principal component is orthogonal to all others.
- ② a) PC1 captures the most variance/information, subsequent components will have lesser

# Real time examples of higher dimensional data

① Healthcare data

② NLP (Natural Language

Processing

↳ Textual  
data)

③ Image & video data

④ financial data

43%

23%

array([0.72962445, 0.22850762, 0.03668922, 0.00517871])

96%

PC1 & PC2

Explained

Variance

Ratio

{