

Welcome back everyone.

⇒ PCA

⇒ Projects



↓  
end to end

⇒ DL / NLP

Cit

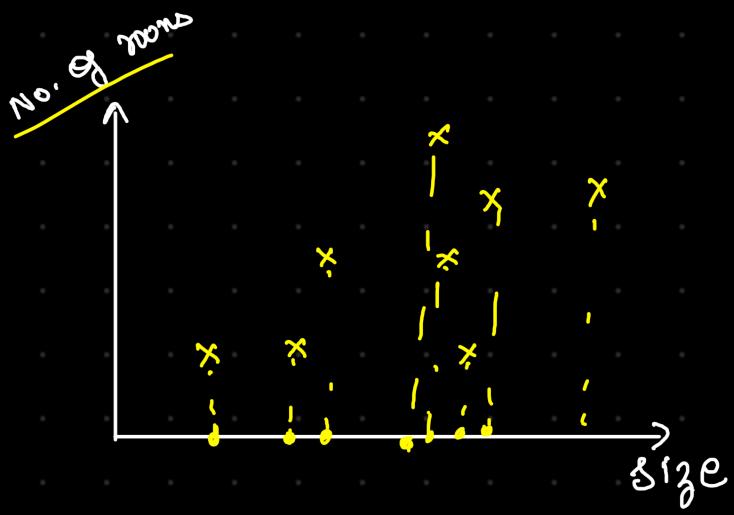
Web Scrapping

# Principle Component analysis (PCA)

size of rooms | No. of rooms | Price

reduce the dimensionality

2 independent  $\rightarrow$  1 independent

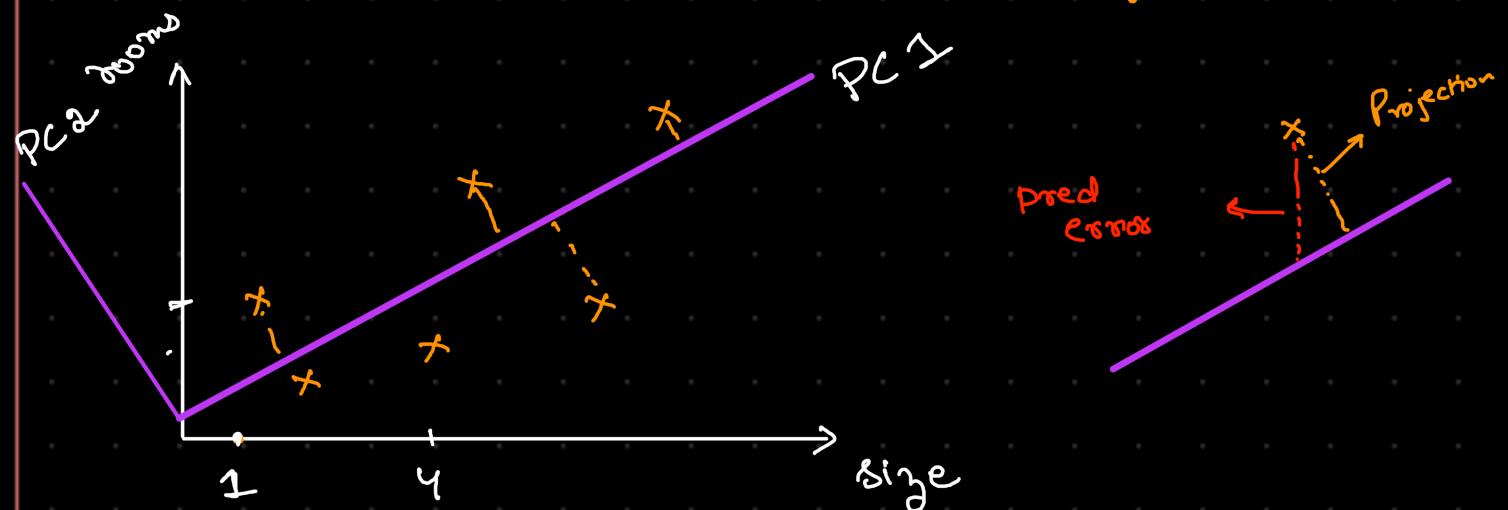


① Feature Selection

$$2D \Rightarrow 1 D$$

Problem is that all information about no. of rooms is lost

All info lost



pred error

Projection

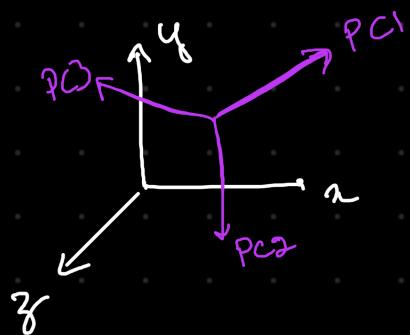
Now, will be able to capture a lot more variance.

spread and variance of both room & size.

no. of features : no. of Principal Component

variance } information of data  
spread }

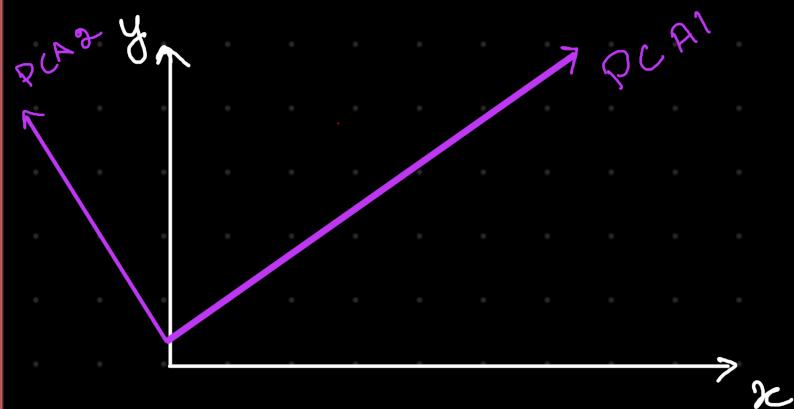
3 dim<sup>n</sup> → PC1   PC2   PC3   always  $\perp \propto$   
will take first 2



reducing dimensionality of my data

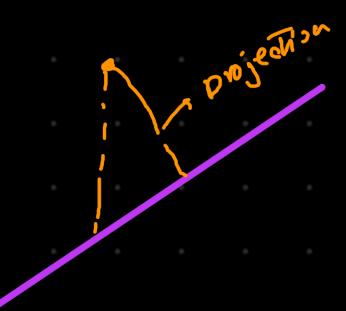
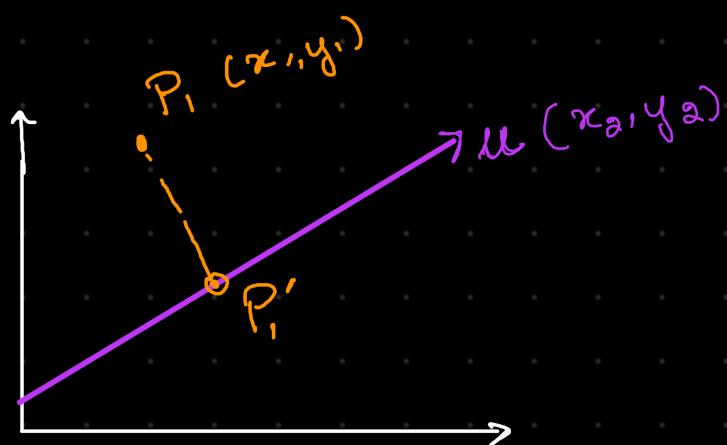
1000dim → 10 dim  
→  $\frac{1}{10}$  sec

# Mathematical Intuition of PCA



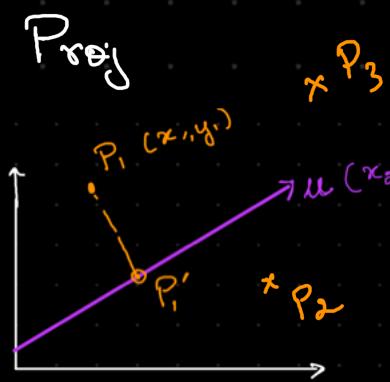
2 features  
⇒ 2 PCA  
vectors

- ① Projections? we are projecting our data point
- ② Optimization? which vector/PCA gives us max<sup>m</sup> variance.



$P_1'$  will be my final point which I will be using in my model training.





$$\text{Proj}_{P_1} u = \frac{P_1 \cdot u}{\|u\|} \rightarrow 1 \text{ for unit vector}$$

$$= P_1 \cdot u$$

$$= \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} \cdot [x_2, y_2]$$

$$= x_1 x_2 + y_1 y_2 \Rightarrow \text{scalar value.}$$

$$\begin{array}{cccc} P_1 & P_2 & P_3 & P_4 \\ \downarrow & \downarrow & \downarrow & \downarrow \\ P_1' & P_2' & P_3' & P_4' \end{array}$$

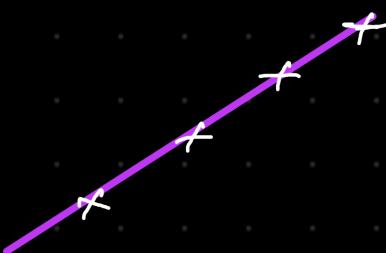
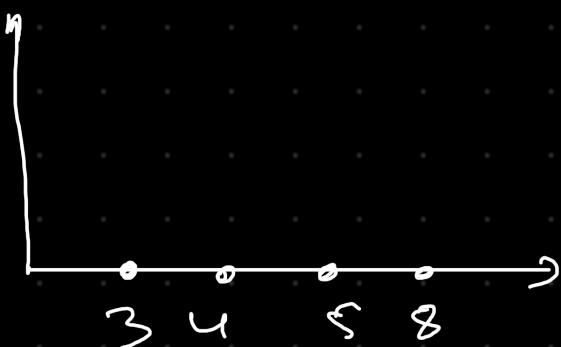
$\Rightarrow$  scalar value

and can then calculate  
the variance

$$x_1' \quad x_2' \quad x_3' \quad x_4'$$

## ② Max<sup>n</sup> variance

$$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n} \Rightarrow \text{cost f^n}$$



How do we find  $\mu$  ?

Eigen values decomposition

Eigen value  $\xleftarrow{\quad}$  Eigen vector

Eigen values tells how much a specific eigen vector is capturing the variance

high eigen value  $\rightarrow$  high variance factor.

Steps.

① Covariance matrix w/o features

$$\text{Cov} \{ f_1, f_2 \} \Rightarrow A$$

② Eigen value & eigen vector will be found out

using covariance matrix

$$A v = \xrightarrow{\quad} \xrightarrow{\quad} \text{vector}$$

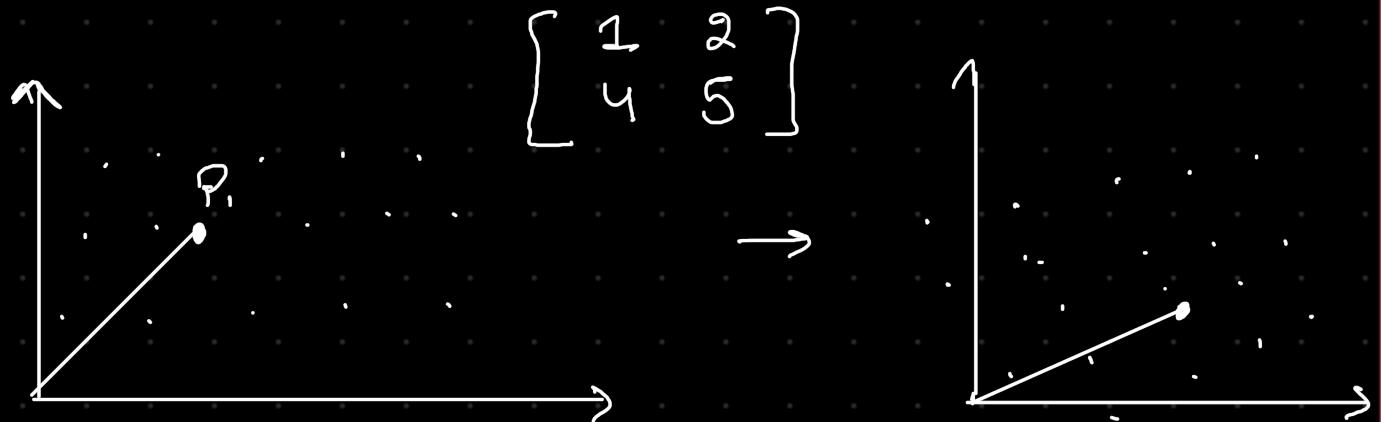
$\downarrow$

Cov matrix eigen value.

③ Eigen vector with most eigen value will be selected as that will capture highest variance.

How to calculate eigen vector & eigen value

{Linear Transformation} of vector



eigen vector is one where our dirr remains same

We have to find

eigen value  $\Rightarrow$  value of vector

$$\begin{bmatrix} & \\ & \end{bmatrix}$$

Cov Matrix

$$V = \lambda v$$

eigen value

Eigen vector  $\rightarrow$  max magnitude

↓

max eigen value

↓

PC1

Steps to calculate eigen value and eigen vector

$2d \rightarrow 1d$

- ⑤ Always standardize before applying PCA
- ⑥ Covariance of features

$$\underbrace{x_1 \ x_2}_{\Downarrow} \quad t_1$$

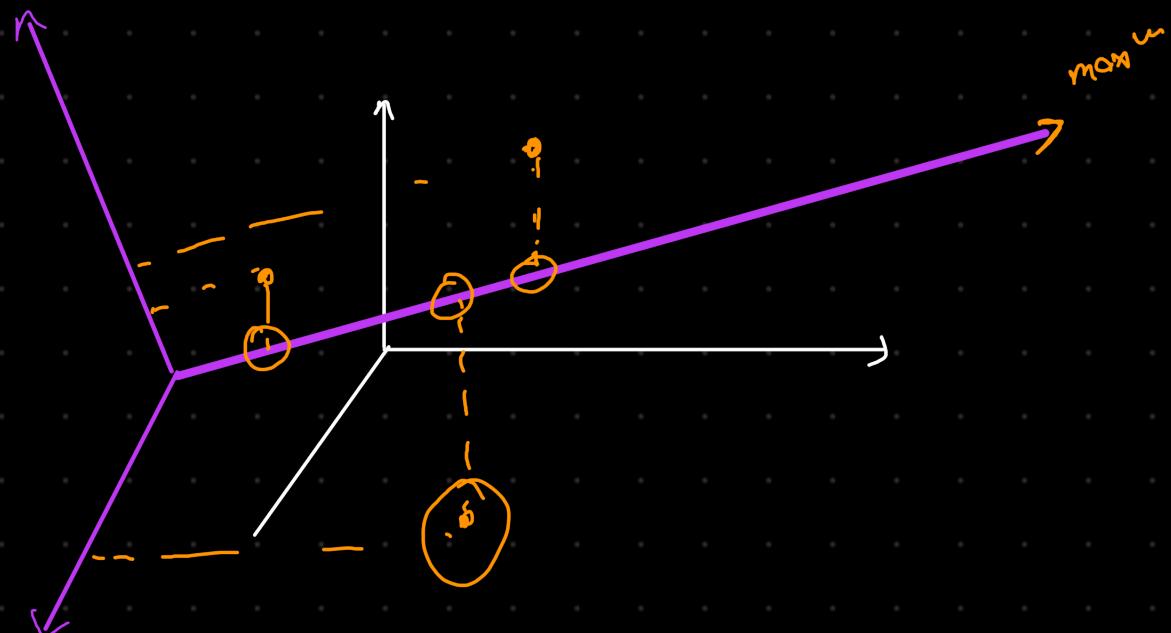
$$\text{Cov}(x_1 x_2) = \sum_{i=1}^n \frac{(x_i - \bar{x})(y_i - \bar{y})}{n}$$

$$\text{Cov} \begin{bmatrix} x_1 & x_2 \\ x_1 & \text{Var}(x_1) \text{ Cov}(x_1, x_2) \\ x_2 & \text{Cov}(x_2, x_1) \text{ Var}(x_2) \end{bmatrix} \Rightarrow A$$

$$A \cdot v = \lambda \cdot v$$

$3d \rightarrow 2d$

$$\begin{array}{cc} \lambda_1 & \lambda_2 \\ \Downarrow & \Downarrow \\ PC_1 & PC_2 \end{array}$$



# Projects

Python  
Statistics  
Data Science  
ML

DB API, file  
Cir  
Web Scrolling

DL/nlp  $\Rightarrow$  Info

Projects

.ipynb

80%

Kaggle

Codex

First we will do projects on ".ipynb" to have data science concepts clear.

