

# Machine Learning

- Machine Learning is a field of computer science that gives computers the ability to learn without being explicitly programmed.
- 1. Supervise learning
- 2. Regression & Applications
- 3. Classification & Applications
- 4. UnSupervised algorithm
- 5. Ensemble learning
- 6. Recommender Systems

## → Supervised learning:

Algorithms use the data developed in Supervised environments to deliver an output.

- \* It is a type of machine learning, where the model is trained on a labeled dataset.
- \* Here, Inputs & Output are known.

## Learning Methods

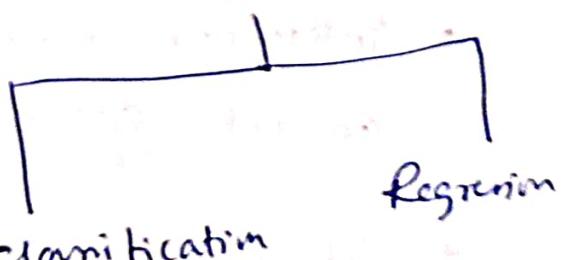
\* Linear Regression

\* Logistic Regression

\* Support Vector machines

\* Decision trees.

## Supervised



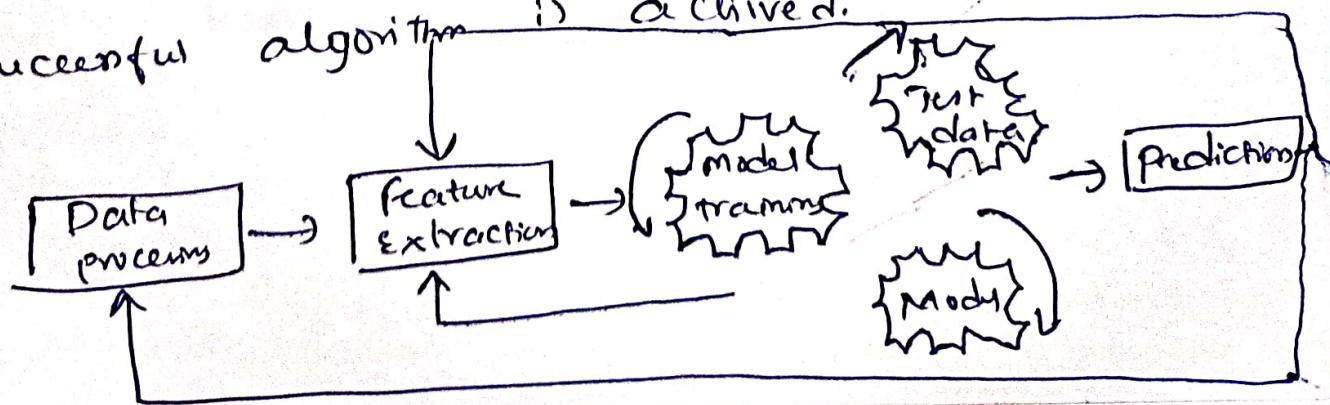
- Unsupervised learning :-
- \* No supervisor is required to prepare data as output may not be known for such algorithms.
  - \* It's trained on unlabeled data.
  - \* It tries to find hidden patterns & recognises their relation.

### Learning methods-

- \* K-means clustering
  - \* Hierarchical clustering
- Reinforcement Learning :-
- \* The program learns from its previous errors.
  - \* It's type of machine learning, & an agent learns to make decision by interacting with environment.
  - \* Interpreter rewards the algorithm when it finds the correct solution.
  - \* If output is not proper, the algorithm reiterates until it finds a better result.

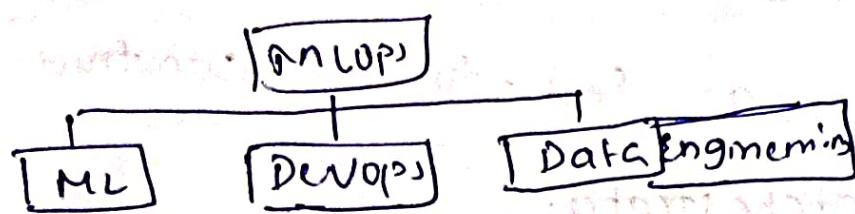
## Machine learning Pipeline :-

- \* It is a series of sequential steps used to codify and automate ML workflows to produce ML models.
- ML pipeline is an end-to-end construct that includes E. Orchestrator:
  - \* Data extraction
  - \* Raw data input
  - + Preprocessing
  - # Features
  - \* Outputs
  - + Model parameters
  - + Model Training
  - \* Deployment
  - + Predicting outputs
- \* The term pipeline implies a one-way flow of data.
- \* ML pipelines are cyclical & iterative sequence is repeated until a successful algorithm is achieved.
- \* Every step in the sequence is activated by the output of the previous step.



## → Machine learning & operations professionals (MLOPs)

- \* MLOPs is a set of practices that combines



- + Ensures reliable & efficient deployment & maintenance of ML models in production systems.

### Aims to:-

- \* Improve communication & collaboration between ML, DevOps, and Data Engineering teams.
- \* Shorten & manage complete development life cycle.
- + Ensure continuous delivery of high-quality predictive services.

### MLOPs consists of three phases:-

- \* Design
- \* Model development
- \* Operations

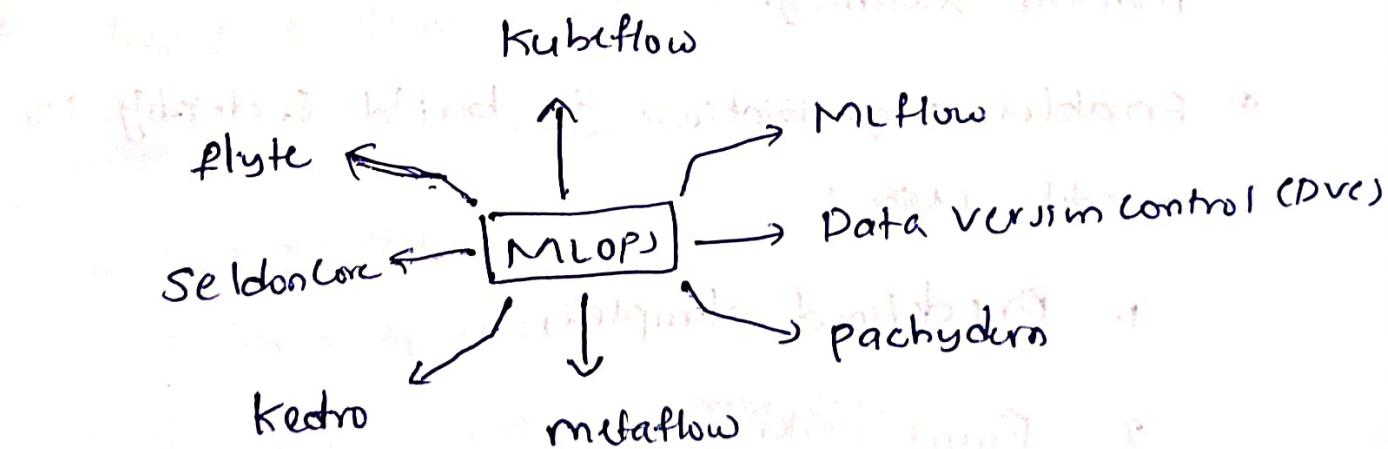
1. Design Phase: Understands that Business and data, & then designs the ML-powered Software.

2. Model development phase: Applicability of ML for the problem is verified by implementing proof of concept (PoC).

3. Operations Phase: Delivers the developed ML model in production.

All three MLOPs phase are interconnected & influence one other.

Tools are used for MLOPs:



## → CI/CD Pipeline Automation:

A continuous integration (CI) or continuous delivery (CD) system is used to:

- \* Test and deploy new pipeline implementations automatically, by providing rapid, reliable of the pipelines in production.
- \* A CI/CD pipeline automation helps with dynamic changes in the data & business environment.

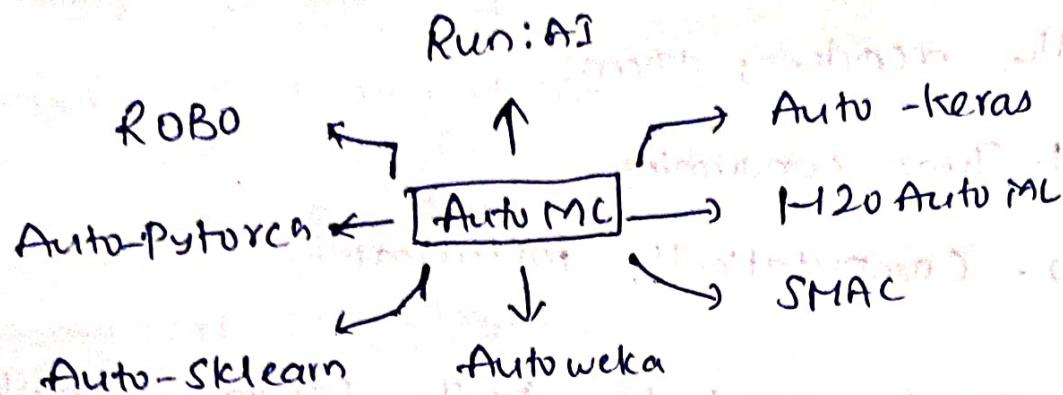
## → Automated Machine Learning (AutoML):

- \* Combines the best practices in automation & machine learning.

- \* Enables organizations to build & deploy ML models using:

1. Predefined templates
2. Frameworks
3. Processors to speed up time to completion.
4. Enhance functionality of ML models.

## Tools of Auto ML



Introduction to python packages used in ML

### & Python:-

→ Python is a high-level programming language that helps programmer to write code that is easy to read and write.

1. Clear
2. Organized
3. Logical

for small & large scale projects, released in 1991

→ Back in 1990, python gained immense popularity in recent times due to:

1. Ease of use
2. Object-Oriented approach
3. code readability owing to significant indentation

n. Evolving open source community

→ Advent of open-source python packages accelerated the field of data science

- Earlier Machine learning (ML), tasks were coded manually, rendering them:
  1. Time-consuming
  2. Computationally inefficient.
- Python may be used to easily perform complex ML tasks & build ML models.
- Python helps to rapidly build and test software prototypes!

Python packages: are folders & modules that form building blocks in python-based programming.

Python ~~base~~ libraries:  
libraries are a collection of packages or files containing prewritten code.

- Python specific
- 1. Numpy
- 2. Pandas
- 3. Tensorflow
- 4. Theano based on

5. matplotlib

6. Scipy

7. Scikit-learn

8. Keras

9. pyTorch for machine learning

1. Numpy & Pandas helps to manage preparation, loading, & manipulation of data.
2. Tensorflow & Theano can be used for fast numerical computing.
3. Matplotlib is used to plot the data.
4. Scipy helps solve the mathematical Equations & algorithms.
5. Scikit-learn provides efficient versions of common algorithms to develop ML models.
6. Keras makes the implementation of neural networks easy.
7. PyTorch specializes in deep learning applications & accelerates the path from prototyping to deployment.

## Supervised Learning :-

- \* machines are trained with labelled data as input.
- \* Machine learning model identify patterns & methods

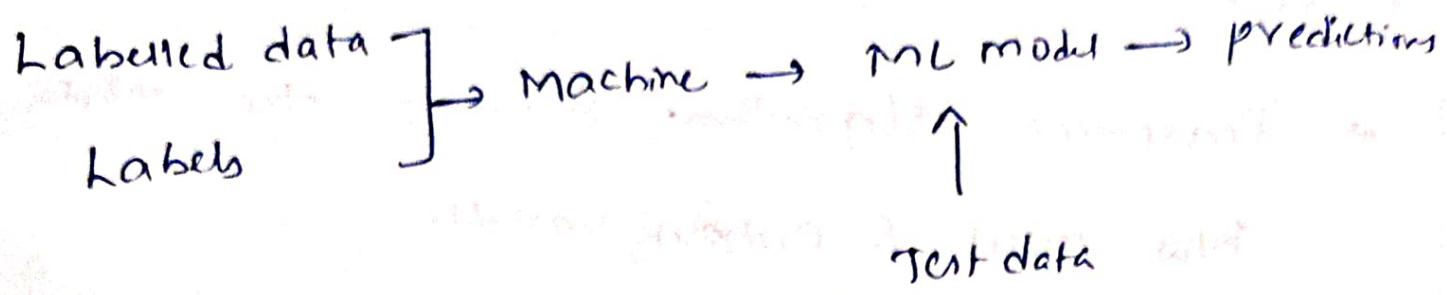
Learn from them

predict output

- \* If any prediction is incorrect, an Operator knows to correct it without consulting others
- \* process continues until the algorithm achieves the highest accuracy

It includes:-

- \* Linear & Logistic regression
- \* Multi-class Classification
- \* Decision trees
- \* Support vector machines (SVM)



Types:-

Supervised

Classification

Regression

Classification:-

\* Classification algorithm segregates data into two or more categories with one or more inputs & a classification module predicts the value of one or more outcomes.

\* A model designed to segregate the

Ex:- Classification

emails into Spam or Ham.

and can start with unlabelled/unstructured data.

## Regression:-

- \* Regression algorithm establishes the relationships b/w input & output variables.
- \* Suitable for situations where the output variable is a real or continuous value.
- \* Regression algorithm is used to forecast or predict the value of the stock market.

## Applications of Supervised learning:-

- Supervise learning can optimize & automate process.
- Used in HR operations to find the right candidate.
- \* It used in finance to segregate loans are:
  - \* good loans
  - \* Bad loans.
- \* Seafarers identifies sea level trends using supervised learning.

Supervised learning models are trained to Detect & prevent large scale fraud.

The model learns from:

- \* financial transaction with fraud pattern.
- \* Cyber attacks on servers
- \* fake social media profiles.
- \* waste management.

→ Overfitting & Underfitting:-

- \* Defines how machine learning model are learning & applying what they learned.

\* Let's understand the terms bias & Variance.

Bias :-

- ~ Bias is an error introduced in a model. (True Error)
- High bias means difference b/w the actual & predicted values. This is not good for model.
- Low bias indicates that the difference b/w the actual & predicted value is low.

Variance :- It indicates how scattered data is.

High Variance :- More scattered data

Low Variance :- less scattered data

<u>Underfit</u>	bias	Variance
over fit	low	high
under fit	High	High

→ Overfit: It indicates low bias & high variance in the data.

\* It happens when a model focuses on too many details in the training dataset.

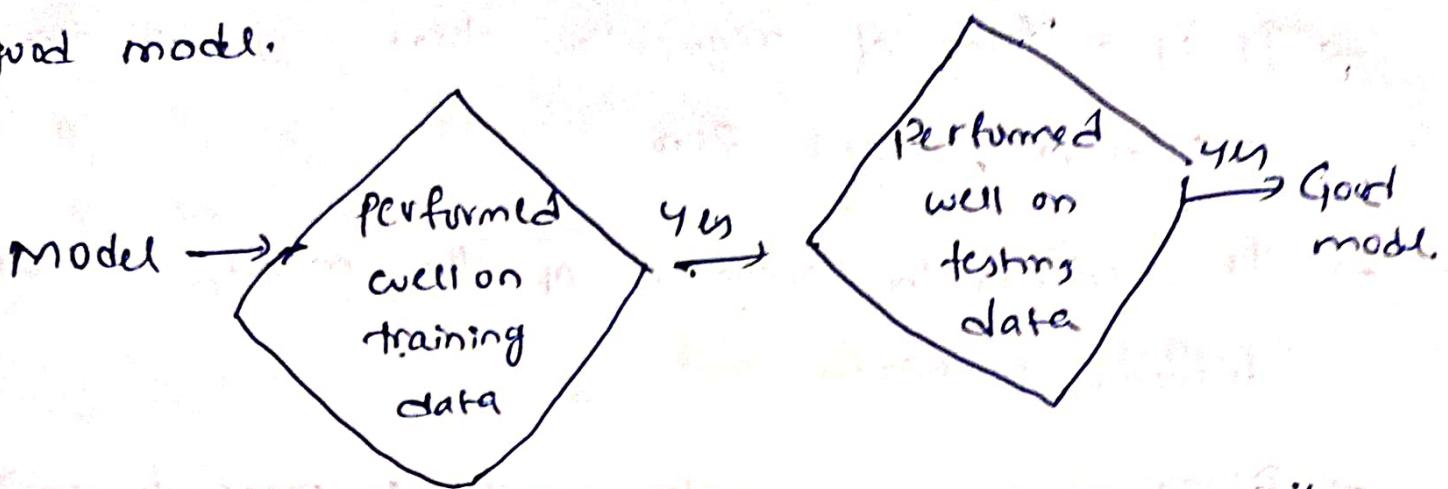
+ has a negative impact on the performance of the model on a new dataset.

→ Underfit: If it has high bias & high variance in the data.

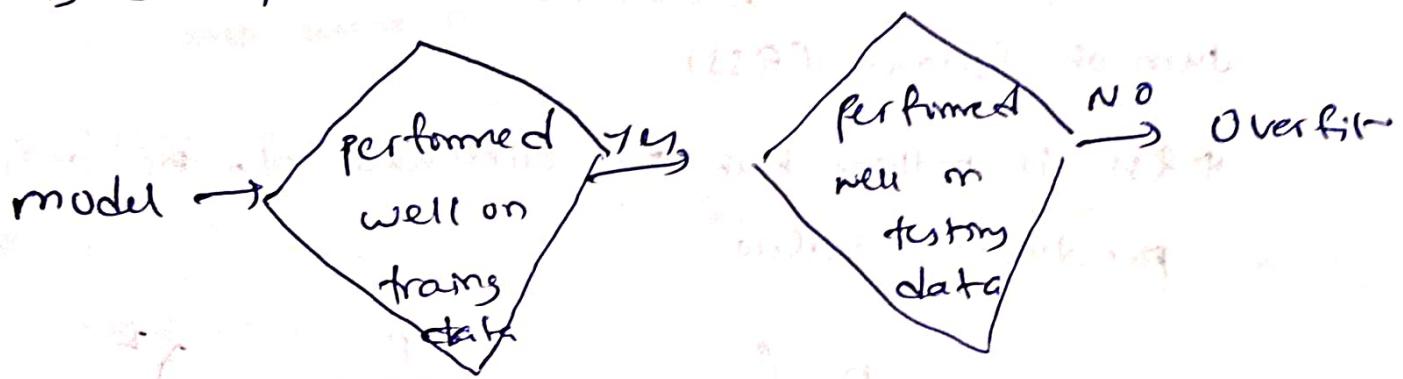
\* Underfitting is detectable as it exhibits poor performance on the training dataset.

+ A model is underfit if it is trained with limited features like temperature & wind speed.

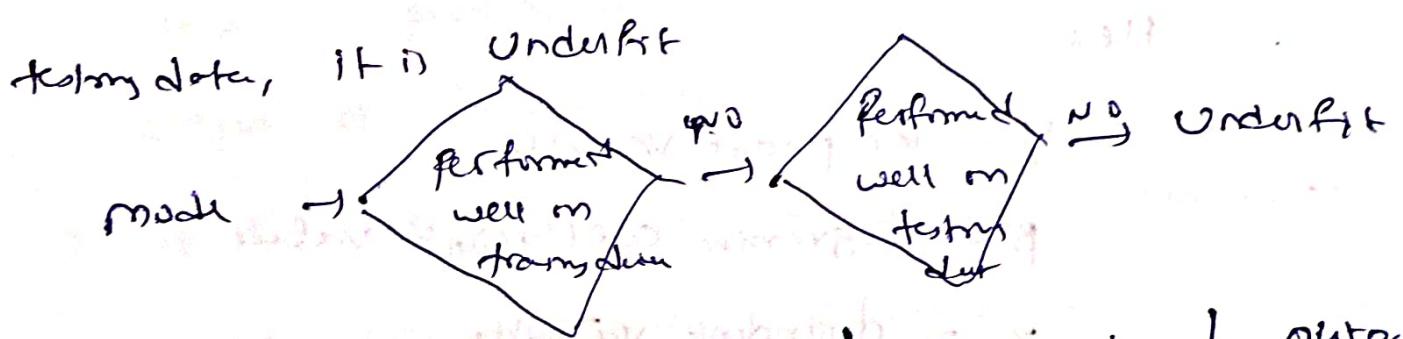
If model performs well on training data & testing data, it's good model.



If model performs well on training data but not with testing data, it is overfitted.



If model doesn't perform well on both training data & testing data, it is underfitted.



	Training data	Testing data	Output
model	✓	✓	Good
model	✓	✗	Over fit
model	✗	✗	Under fit

## Regularization

- It is a form of regression that shrinks the coefficient towards zero.
- It reduces the variance of the model, without increase in bias.
- Discourages more complex models & prevent overfitting.
- Fitting involves a loss function and known as residual sum of squares (RSS).
  - \* RSS is nothing but the difference b/w the actual & predicted value.

$$RSS = \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2$$

Here,

$y$  - independent variable

$\beta$  - regression coefficient value

$x$  - dependent variable

## Types of Regularization:-

### \* Dropout Regularization:

+ It works by removing a random selection.

+ The more units dropout, the stronger it is regularization.

\* This is good for training Neural Networks.

### \* Early Stopping:-

+ It uses large number of Epochs & plot the validation loss graph.

\* Stop training and save the model when the validation loss moves from decreasing to increasing.

### + Co-adaptation:-

+ Neurons predict patterns in the training data using output of specific neurons.

+ If the validation data doesn't have patterns that causes co-adaptation, it could cause overfitting.

\* Dropout regularization reduces co-adaption.

\* It ensures neurons cannot rely solely on other neurons.

## & Losses Regression Regularization - (L1 Regularization)

- + It ~~is~~ penalizes weights in proportion to the sum of the absolute values of the weights.
- + Drives the weights of irrelevant or barely relevant features to exactly zero.

## & Ridge Regression : (L2 Regularization)

- + Ridge Regression shrinks coefficients close to zero for unimportant predictors.
- + Never make them exactly zero.
- + In other words, the final model will include all predictors.

## Linear Regression Equation:-

$$\hat{y} = \beta_0 + \beta_1 x$$

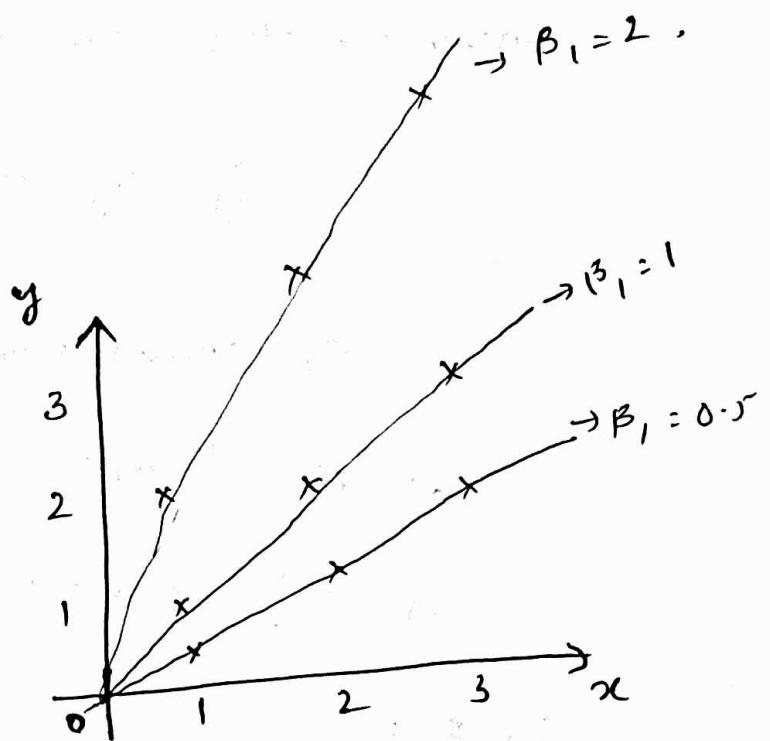
Here,

$$\rightarrow \beta_0 = 0, \beta_1 = 1, x = \{1, 2, 3\}$$

$$\Rightarrow x=1, \hat{y} = 0 + 1(1) = 1$$

$$\Rightarrow x=2, \hat{y} = 0 + 1(2) = 2$$

$$\Rightarrow x=3, \hat{y} = 0 + 1(3) = 3$$



$$\rightarrow \beta_0 = 0, \beta_1 = 0.5, x = \{1, 2, 3\}$$

$$\Rightarrow x=1, \hat{y} = 0 + 0.5(1) = 0.5$$

$$\Rightarrow x=2, \hat{y} = 0 + 0.5(2) = 1$$

$$\Rightarrow x=3, \hat{y} = 0 + 0.5(3) = 1.5$$

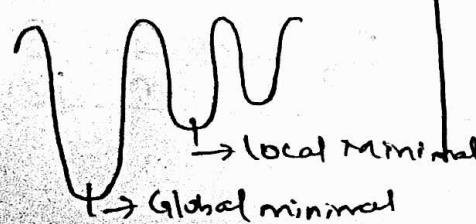
Mean Square Error (or) Cost function :-

$$J(\beta_1) = \frac{1}{2m} (\hat{y}_i - y_i)^2$$

$$J(1) = \frac{1}{2(3)} [(1-1)^2 + (2-2)^2 + (3-3)^2] = \frac{0}{6} = 0$$

$$J(0.5) = 0.58$$

$$J(2) = 2.5$$



$$\rightarrow \beta_0 = 0, \beta_1 = 2, x = \{1, 2, 3\}$$

$$\Rightarrow x=1, \hat{y} = 0 + 2(1) = 2$$

$$\Rightarrow x=2, \hat{y} = 0 + 2(2) = 4$$

$$\Rightarrow x=3, \hat{y} = 0 + 2(3) = 6$$



## Gradient Descent Optimization:-

Minimize the cost function  $J(\beta)$  (m.s.B).

$$\theta_j : \theta_j - \alpha \frac{\partial J}{\partial \theta_j}$$

→  $\theta_j$  = Model Parameters ( $\beta_0, \beta_1$ )

→  $\alpha$  = Learning rate

→  $\frac{\partial J}{\partial \theta_j}$  = partial derivative (Gradient).

in  $\beta_j$ 's,

$$\beta_j = \beta_j - \alpha \frac{\partial J}{\partial \beta_j}$$

$$[\because J = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2]$$

$j=0$

$$\beta_0 = \beta_0 - \alpha \frac{\partial}{\partial \beta_0} \left[ \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \right]$$

$$[\because \hat{y}_i = \beta_0 + \beta_1 x_i]$$

$$\Rightarrow \beta_0 = \beta_0 - \alpha \frac{\partial}{\partial \beta_0} \left[ \frac{1}{2m} \sum_{i=1}^m [(\beta_0 + \beta_1 x_i) - y_i]^2 \right]$$

$$= \beta_0 - \alpha \frac{1}{2m} \sum_{i=1}^m [(\beta_0 + \beta_1 x_i) - y_i] \frac{\partial}{\partial \beta_0} (\beta_0)$$

$$\beta_0 = \beta_0 - \alpha \frac{1}{m} \sum_{i=1}^m [(\beta_0 + \beta_1 x_i) - y_i]$$

$$\boxed{\beta_0 = \beta_0 - \alpha \frac{1}{m} \sum_{i=1}^m [\beta_0 \hat{y}_i - y_i]}$$

$$\beta_1 = \beta_1 - \alpha \frac{\partial}{\partial \beta_1} \left[ \frac{1}{2m} \sum_{i=1}^m [(\beta_0 + \beta_1 x_i) - y_i]^2 \right]$$

$$= \beta_1 - \alpha \frac{1}{2m} \sum_{i=1}^m [(\beta_0 + \beta_1 x_i) - y_i] \frac{\partial}{\partial \beta_1} [\beta_1 x_i]$$

$$= \beta_1 - \alpha \frac{1}{2m} \sum_{i=1}^m [(\beta_0 + \beta_1 x_i) - y_i] x_i$$

$$\beta_1 = \beta_1 - \alpha \frac{1}{m} \sum_{i=1}^m [(\beta_0 + \beta_1 x_i) - y_i] x_i$$

$$\boxed{\beta_1 = \beta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) x_i}$$

$$\beta_j = \beta_j - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) x_j^{(i)} \text{ for all } j \in \{1, \dots, n\}$$

$$\beta_j = \beta_j - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) x_j$$

## Logistic Regression:

### Cost function:

$$J(\beta) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)]$$

$y_i \in \{0, 1\}$

- True class label

$$\hat{y}_i = h_{\beta}(x_i) = \frac{1}{1+e^{-x_i^T \beta}} \quad (\text{predicted probability})$$

Here,  $m$  = number of training Examples or Observations

$y_i$  = Actual class label (0 or 1)

$\hat{y}_i$  = predicted probability (from Sigmoid function)

$\beta$  = model parameters.

Likelihood Function:

$$P(y_i | \hat{y}_i) = (\hat{y}_i)^{y_i} (1 - \hat{y}_i)^{1-y_i}$$

\*  $y_i = 1 \rightarrow \hat{y}_i$

\*  $y_i = 0 \rightarrow 1 - \hat{y}_i$

Log-Likelihood Function:

$$L(\beta) = \prod_{i=1}^n P(y_i | \hat{y}_i) \quad [\because \log(ab) = b \log a]$$

Taking log on Both sides,

$$\begin{aligned} \Rightarrow \log L(\beta) &= \log \left[ \prod_{i=1}^n (\hat{y}_i)^{y_i} (1 - \hat{y}_i)^{1-y_i} \right] \\ \Rightarrow \log L(\beta) &= \sum_{i=1}^n \left[ \log(\hat{y}_i)^{y_i} + \log(1 - \hat{y}_i)^{1-y_i} \right] \\ &= \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \end{aligned}$$

Negative Log-Likelihood,

$$J(\beta) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

## Gradient Descent Optimization :-

Minimize the cost function  $J(\beta) - \text{M.S.E.}$

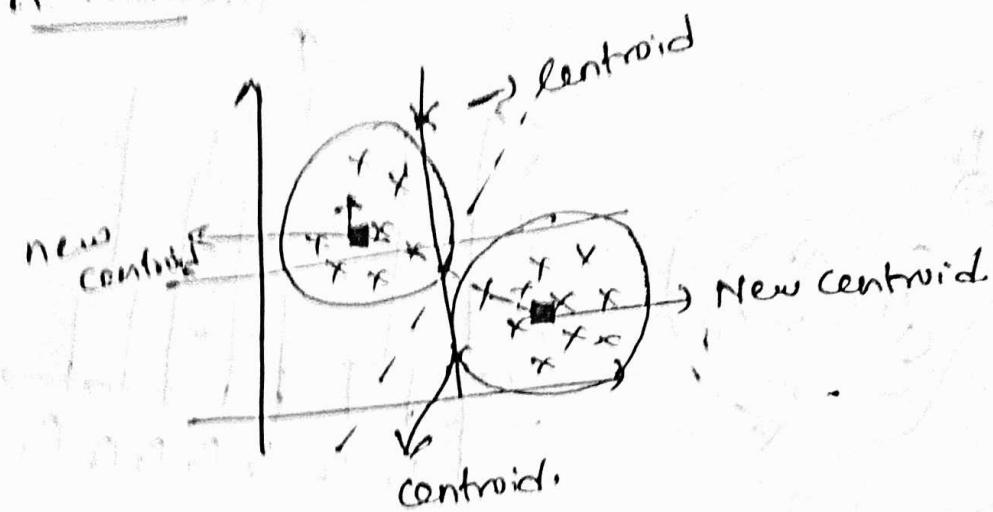
$$\theta_j : \theta_j - \alpha \frac{\partial J}{\partial \theta_j} \quad [ \because \hat{y}_i = h_{\beta}(x_i) \\ \beta_j = \beta_j - \alpha \frac{\partial J}{\partial \beta_j} \qquad \qquad \qquad = \frac{1}{1 + e^{-x_i \beta}} ]$$

$$J = 0, \quad \beta_j = \beta_0 \Rightarrow \beta_0 = \beta_0 - \alpha \frac{\partial}{\partial \beta_0} \left[ -\frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)] \right]$$

$$\beta_0 = \beta_0 - \alpha \frac{\partial}{\partial \beta_0} \left[ -\frac{1}{m} \sum_{i=1}^m [y_i \log(h_{\beta_0} x_i) + (1-y_i) \log(1-h_{\beta_0} x_i)] \right]$$

## K-means

(Large Dataset)



1. We find k values  $\Rightarrow$  Suitable  $k=2$

2. Initialise k-number of centroids

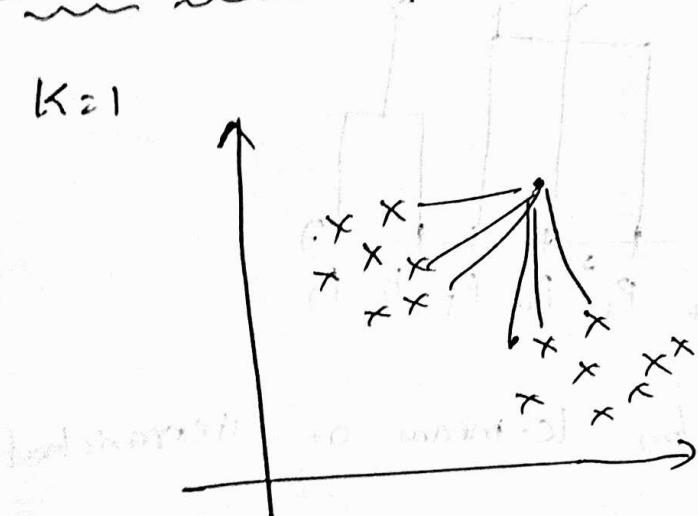
3. Compare the AVG to update centroids

Deciding k-value for validating purpose

Elbow method (k-value)

WCSS - Within Cluster Sum of Square

$K=1$

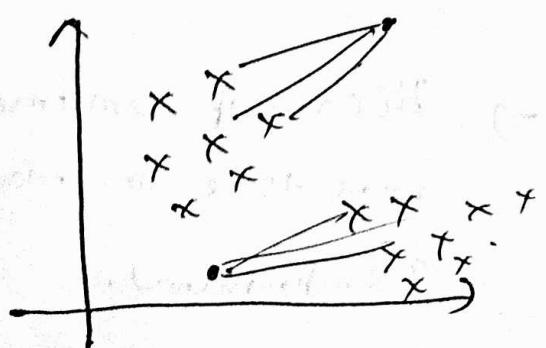


WCSS

Elbow curve

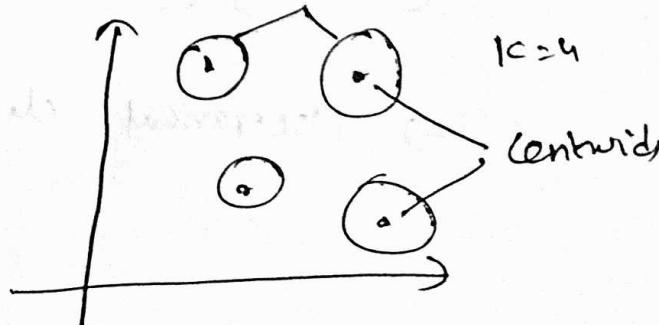
$K=4$  abrupt change

$K=2$



clusters-

$K=4$



## K-means:-

Point	X	Y	Nearest Centroid	Cluster	Centroid
A	1	2	Centroid 1	1	1.5
B	2	1	Centroid 1	1	1.5
C	4	5	Centroid 2	2	1, 4.5
D	5	4	Centroid 2	2	1, 4.5

Step-1 :-  $k=2$ , (Assume)

Centroid-1  $\rightarrow (1, 2)$  (Same as point A)

Centroid-2  $\rightarrow (5, 4)$  (Same as point B)

Step-2 :- Calculate Euclidean Distance from Each Point to Each centroid.

Each centroid.

$$d(p, q) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Point

A(1, 2) :

$$d = \sqrt{(1-1)^2 + (2-2)^2} = 0$$

Centroid 1 (1, 2) :

$$d = \sqrt{(1-5)^2 + (2-4)^2} = \sqrt{20} = 4.47$$

Centroid 2 (5, 4) :

$\therefore$  point A closer to Centroid-1

Finally

cluster 1, (1, 2)  
(2, 1)

cluster 2, (4, 5)  
(5, 4)

Centroids1 (1.5, 1.5)  
2 (4.5, 4.5)

Point B(2, 1)

$$\text{to Centroid 1 - (1, 2)}: d = \sqrt{(2-1)^2 + (1-2)^2} = \sqrt{1+1} = \sqrt{2} = 1.41$$

$$\text{to Centroid 2 - (5, 4)}: d = \sqrt{(2-5)^2 + (1-4)^2} = \sqrt{9+9} = \sqrt{18} = 4.24$$

$\therefore$  Point B is closer to Centroid - 1

Point C(4, 5)

$$\text{to Centroid 1 - (1, 2)}: d = \sqrt{(4-1)^2 + (5-2)^2} = \sqrt{9+9} = \sqrt{18} = 4.24$$

$$\text{to Centroid 2 - (5, 4)}: d = \sqrt{(4-5)^2 + (5-4)^2} = \sqrt{1+1} = \sqrt{2} = 1.41$$

$\therefore$  Point C is closer to Centroid - 2

Point D(5, 4)

$$\text{to Centroid 1 - (1, 2)}: d = \sqrt{(5-1)^2 + (4-2)^2} = \sqrt{16+4} = \sqrt{20} = 4.47$$

$$\text{to Centroid 2 - (5, 4)}: d = \sqrt{(5-5)^2 + (4-4)^2} = 0$$

$\therefore$  Point D is closer to Centroid - 2

Step-3:

Point	Nearest centroid	Cluster
A	Centroid 1	Cluster - 1
B	Centroid 1	Cluster - 1
C	Centroid 2	Cluster - 2
D	Centroid 2	Cluster - 2

#### Step-4 : Update Centroids:

Now, we calculate the new position of each centroid by finding the mean (average) of points in its cluster.

→ Cluster 1: Points A(1,2) & B(2,1)

New centroid 1:

$$x = (1+2)/2 = 1.5$$

$$y = (2+1)/2 = 1.5$$

$$\boxed{\text{New centroid } 1 = (1.5, 1.5)}$$

→ Cluster 2: Points C(4,5) & D(5,4)

New centroid 2:

$$x = \frac{(4+5)}{2} = 4.5$$

$$y = \frac{(5+4)}{2} = 4.5$$

$$\boxed{\text{New centroid } 2 = (4.5, 4.5)}$$

## Step-5 - Recalculate Distance (Repeat Step-2)

$\Rightarrow$  • Distance from A(1,2) & <sup>new</sup>Centroid-1 (1.5, 1.5)

$$d = \sqrt{(1-1.5)^2 + (2-1.5)^2} = \sqrt{0.25 + 0.25} = \sqrt{0.5} \approx 0.71$$

Centroid 2 (4.5, 4.5)

$$d = \sqrt{(1-4.5)^2 + (2-4.5)^2} = \sqrt{12.25 + 6.25} = \sqrt{18.5} \approx 4.30116$$

$\therefore$  Point A is closer to new centroid -1

$$\Rightarrow \text{Point B } C(2,1) \quad c_1 = \sqrt{(2-1.5)^2 + (1-1.5)^2} = \sqrt{0.25} = 0.5$$

$$c_2 = \sqrt{(2-4.5)^2 + (1-4.5)^2} = \sqrt{18.5} \approx 4.30$$

$\therefore$  point B is closer to new centroid -2

$$\Rightarrow \text{Point C } C(4,5) \quad c_1 = \sqrt{(4-1.5)^2 + (5-1.5)^2} = \sqrt{6.25 + 12.25} = \sqrt{18.5} \approx 4.30116$$

$$c_2 = \sqrt{(4-4.5)^2 + (5-4.5)^2} = \sqrt{0.25} = 0.5$$

$\therefore$  point C is closer to new centroid -2

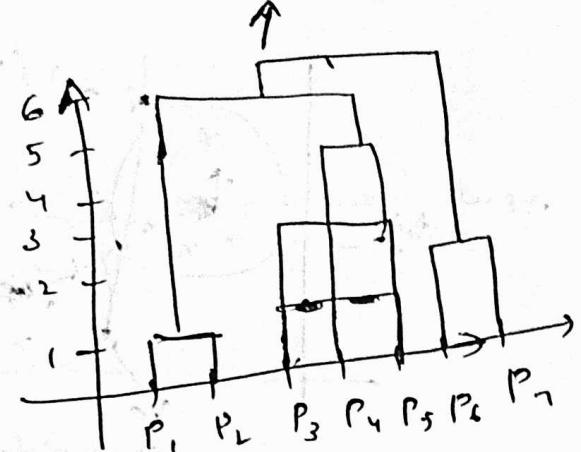
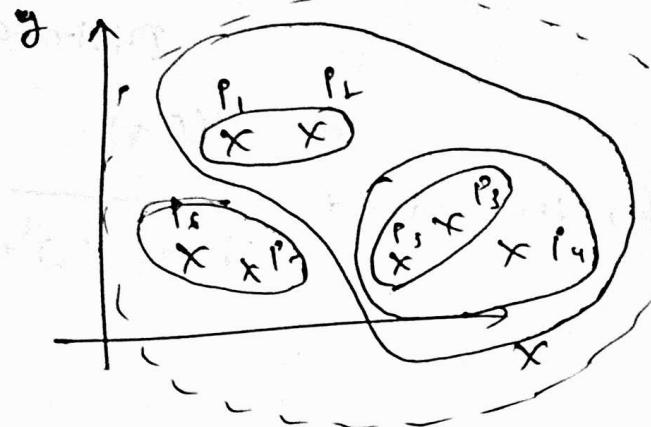
$$\Rightarrow \text{Point D } C(5,4) \quad c_1 = \sqrt{(5-1.5)^2 + (4-1.5)^2} = \sqrt{18.5} \approx 4.30116$$

$$c_2 = \sqrt{(5-4.5)^2 + (4-4.5)^2} = \sqrt{0.25} = 0.5$$

$\therefore$  point D is closer to new centroid -2.

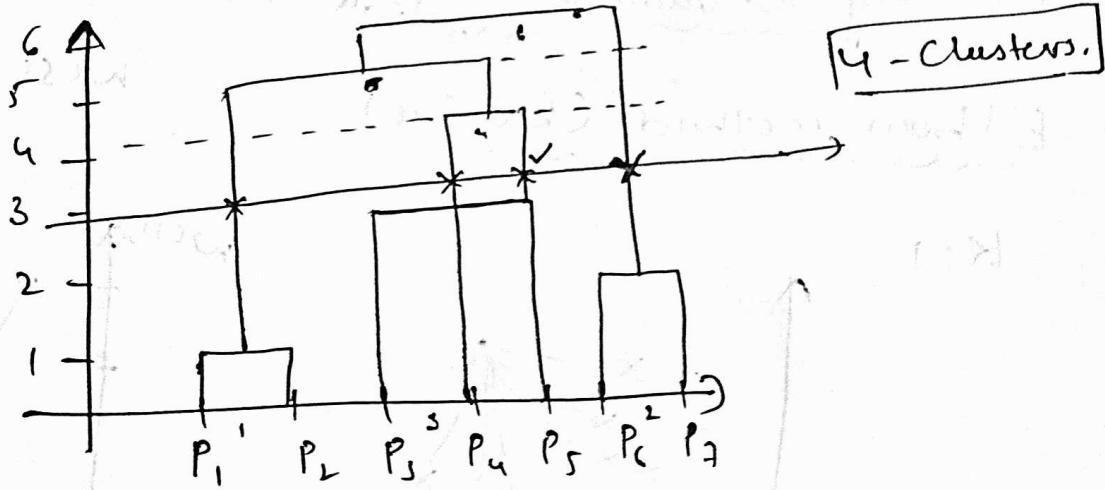
$\therefore$  Cluster assignments remain the same. So the algorithm stops.

## 2. Hierarchical clustering: (Small dataset) Dendrogram



Dendrogram - Bottom root to top.

+ You need to find out the longest vertical line that has no horizontal line passed through it.



\* Max time is taken by k-means or hierarchical clustering !!

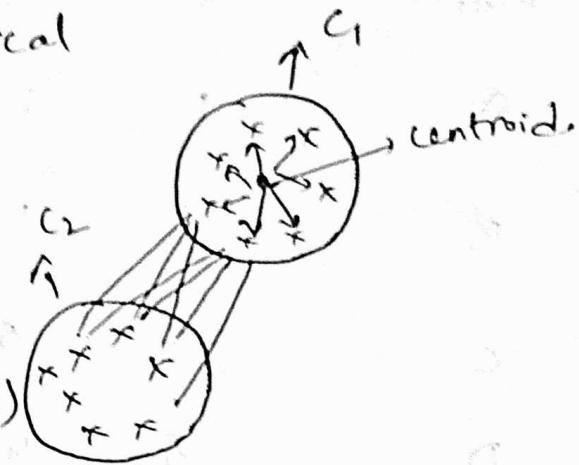
→ Hierarchical clustering → Because of constructing more time to make Dendrogram.

## Validate Clustering models:

\* Silhouette Score  $\rightarrow$  k-means  
 \*  $\rightarrow$  Hierarchical

$$\rightarrow a(i) = \frac{1}{|C_i|-1} \sum_{j \in C_i, i \neq j} d(i, j)$$

$$b(i) = \min_{j \neq i} \frac{1}{|C_j|} \sum_{j \in C_j} d(i, j)$$



Good clustering model:  $a(i) > b(i) \gg a(i)$   
 $a(i) > b(i)$  or  $b(i)$  is small (good separation)

$\rightarrow$  value b/w -1 to +1

+1  $\rightarrow$  good model

-1  $\rightarrow$  bad model

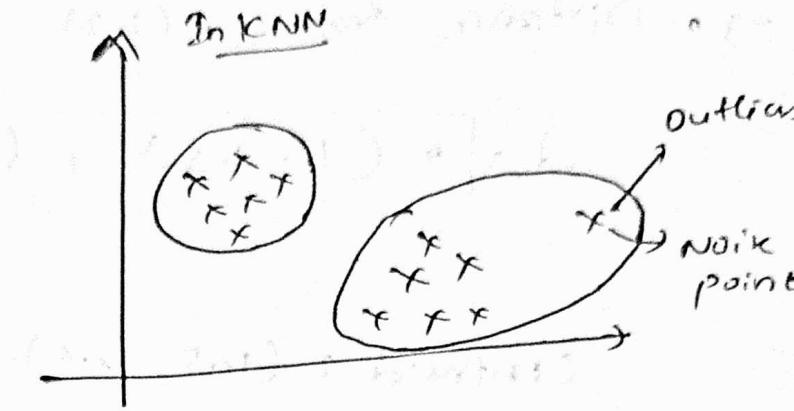
0  $\rightarrow$  clustering need to be improve

$$\text{or } \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

$$S(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \Rightarrow -1 \rightarrow +1$$

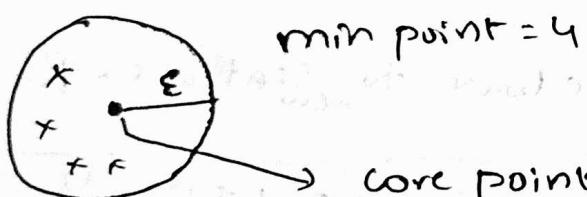
# DBScan Clustering (Density Based Spatial Clustering Application with Noise (DBSCAN))

1. Min point
2. Core point
3. Border point
4. Noise point



\* Min point  $\rightarrow$  kind of hyperparameter

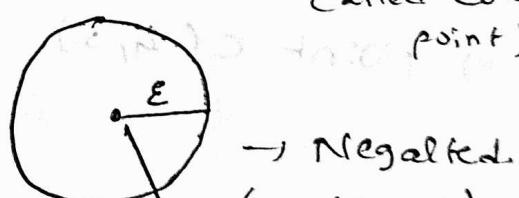
L) called as Epsilon.  $\rightarrow$  radius of circle.



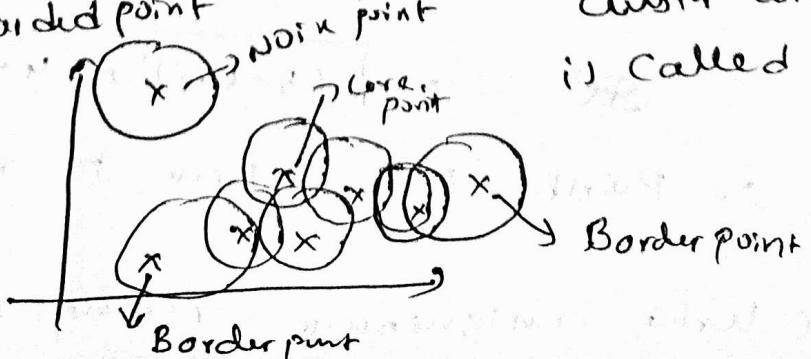
(Assume min points are present in cluster)  $\rightarrow$  called core point



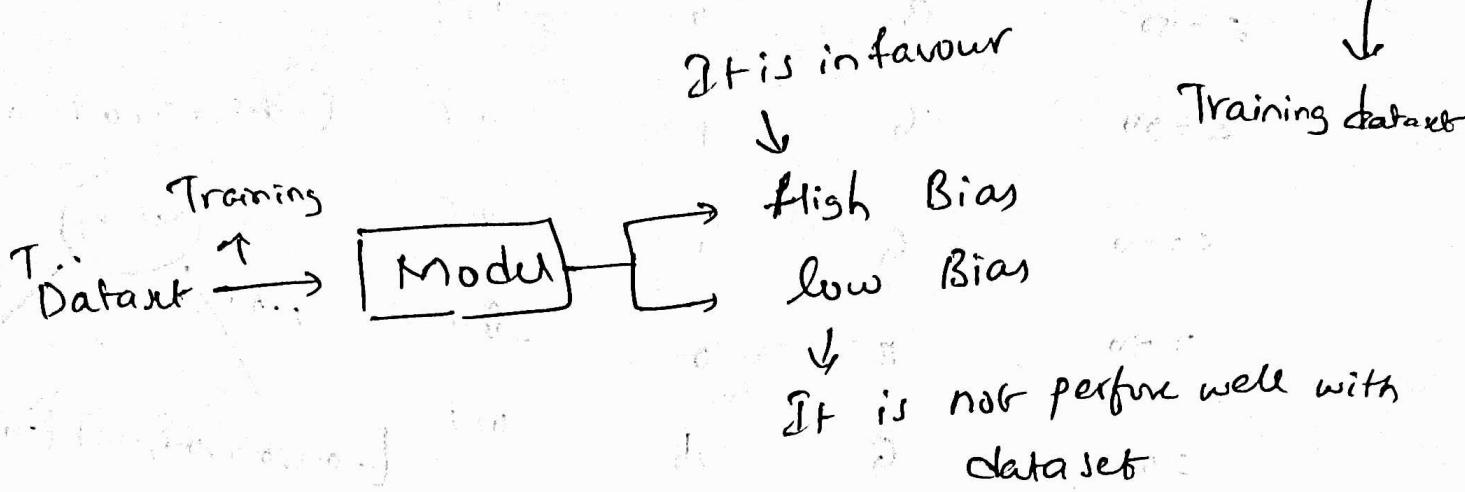
Cluster with one core point (at least)  
i) called Border point



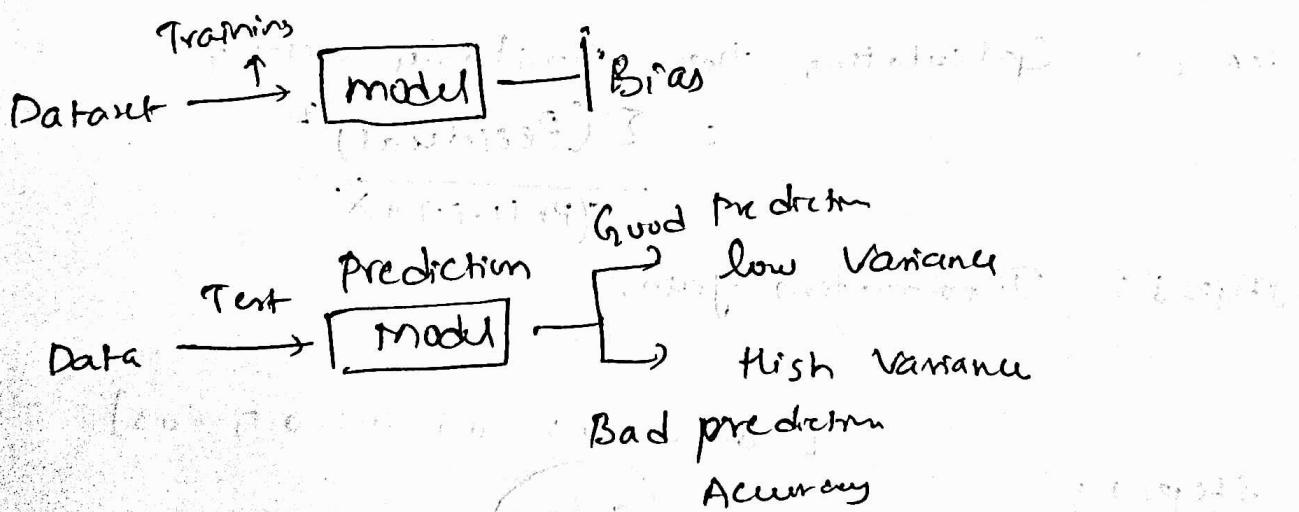
Cluster with no point  
i) called Noise point



Bias :- It is a phenomenon that skews the results of an algorithm in favour or against an idea.



Variance :- Variance refers to changes in the model when using different portions of the training or test data.



Model 1  
Trainng ACC = 90%  
Test ACC = 70%

{  
↓  
low Bias }  
High Variance

Model 2  
Train ACC = 60%  
Test ACC = 55%

{  
↓  
high Bias }  
high variance

Model 3  
Train ACC = 90%  
Test ACC = 92%

{  
↓  
low Bias }  
low variance

↓  
Generalized model

# XGBoost Classifier - Extreme Gradient Boosting

Salary	Credit	Approval	Residual	Base model	Pr = 0.5
< 250	B	0	-0.5		
< 250	G	1	0.5	$[-0.5, 0.5, 0.5, -0.5, 0.5]$	$[-0.5, 0.5]$
< 250	G	1	0.5	$[-0.5, 0.5, 0.5, -0.5, 0.5]$	$[-0.5, 0.5]$
> 250	B	0	-0.5		
> 250	G	1	0.5	$[-0.5, 0.5, 0.5, -0.5]$	$[-0.5, 0.5]$
> 250	N	1	0.5	$\downarrow$ $s.w = 0$	$s.w = 0.33$
< 250	N	0	-0.5		

Step-1 :- Create a Binary Decision Tree Using the feature.

Step-2 :- Calculating the similarity weight

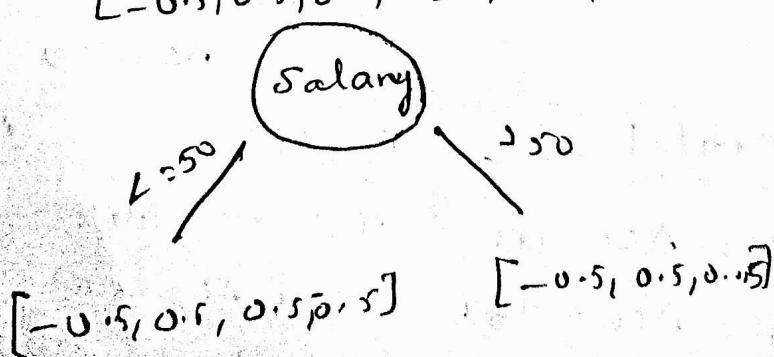
$$= \sum (\text{Residual})^2$$

$$\frac{\sum (\text{Residual})^2}{\sum (Pr(1-Pr) + \lambda)}$$

Step-3 :- Information gain.

$$[-0.5, 0.5, 0.5, -0.5, 0.5, 0.5, -0.5]$$

Step-1 :-



Similarity weight  $\lambda \geq 0 \rightarrow$  kind of Hyperparameter

Step-2 :- Residuals

$$\text{Similarity weight} = \frac{[-0.5 + 0.5 - 0.5]^2}{0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5) + 0}$$

$$= 0$$

$$S.W. 2 = \frac{[-0.5 + 0.5 + 0.5]^2}{0.5(1-0.5) + 0.5(1-0.5) + 0.5(1-0.5)} \\ = \frac{0.25}{0.25 + 0.25 + 0.25} = \frac{0.25}{0.75} = 0.33$$

$$S.W. \text{ root class} = \frac{[-0.5 + 0.5 + 0.5 + 0.5 + 0.5 + 0.5 - 0.5]^2}{0.25 + 0.25 + 0.25 + 0.25 + 0.25 + 0.25} \\ = \frac{0.25}{1.75} = 0.142$$

$$\text{Step-3 :- Informative Gain} = S.W. 1 + S.W. 2 - S.W. \text{root}$$

$$= 0 + 0.33 - 0.14$$

$$= 0.19$$

$$[-0.5, 0.5, 0.5, -0.5, 0.5, 0.5, -0.5] \rightarrow S \cdot w = 0.14$$

Informative

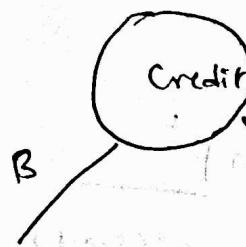
$$\text{Gain} = 0 + 0.33 - 0.14 \\ = 0.19 \quad L = 50$$

$$S \cdot w = 0 \leftarrow [-0.5, 0.5, 0.5, -0.5]$$

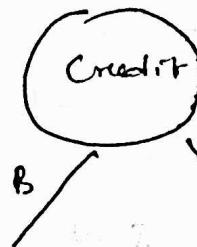
$$[0.5, 0.5, 0.5] \rightarrow S \cdot w = 0.33$$

Salary

\$50



$$S \cdot w = 1 \leftarrow [-0.5]$$



$$[0.5, 0.5, -0.5]$$

$$\rightarrow S \cdot w = 0.33$$

Information gain

$$= 1 + 0.33 - 0$$

$$= 1.33$$

for first value in the table

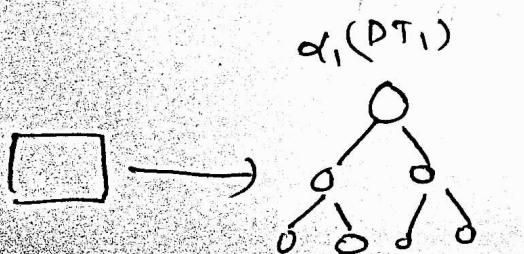
$\alpha$  - Learning rate.

$$\sigma[0 + \alpha(1)]$$

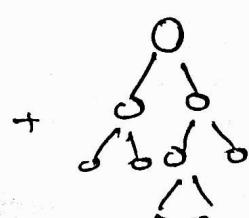
$$\sigma[\alpha_0 + \alpha_1(DT_1) + \alpha_2(DT_2) + \alpha_3(DT_3) + \dots + \alpha_n(DT_n)]$$

XG Boost - Black Box Model

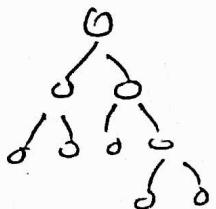
→ Pre pruning



$\alpha_2(DT_2)$



$\alpha_n(DT_n)$



$w_{\text{weakLearn}} + w_{\text{weakLearn}} + \dots + w_{\text{weakLearn}}$   
= Strong learner,

# XG Boost Regressor! - Bias-Variance

Avg of Salaries.

Exp Gap Salary

Exp	Gap	Salary	$\hat{R}_1$
2	Yes	40k	-11k
2.5	Yes	42k	-9
3	No	52k	1
4	No	60k	9
4.5	Yes	62k	"

①

[Base model]  $\rightarrow 51k$

$$[-11, 9, 1, 9, 11] \rightarrow s.w = 6\%$$

Exp

52

$$s.w = 60.5 \left[ \begin{matrix} -11 \\ 9 \end{matrix} \right] \rightarrow [9, 1, 9, 11] \rightarrow s.w = 28.5$$

$$2 \cdot 4 = 89.13$$

$$\text{Similarity Weight} = \frac{\sum (\text{Residuals})^2}{\text{No. of residuals} + \lambda}$$

Here,  $\lambda = 0 \rightarrow \text{Hyperparameter Tuning}$

$$= \frac{[-11]^2}{1+0} = \frac{121}{2} = 60.5 \quad \left| \frac{[-9+1+9+11]^2}{4+1} = \frac{144}{5} = 28.5 \right.$$

$$\text{root} = \frac{[-11+9+1+9+11]^2}{5+1} = \frac{1}{6} =$$

$$\text{Information Gain} = 60.5 + 28.5 - \frac{1}{6} = \underline{\underline{89.13}}$$

$$[-11, -9, 1, 9, 11] \xrightarrow{S.W.} 1/6$$

②

EXP

$\leftarrow$  L.R.W.  $\rightarrow$  R.R.W.

$$\begin{aligned} S.W. &= [ -11, 9 ] \\ (33.33) &\Leftrightarrow [ 1, 9, 11 ] \xrightarrow{S.W.} 1/0.25 \end{aligned}$$

$$I.G = 133.33 + 1/0.25 - 1/6$$

$$I.G = 245.38$$

1st Q.G << 2nd I.G

$$51 + \alpha_1 (\text{Avg}(-11, -9))$$

$$51 + \alpha_2 [\text{Avg}(1, 9, 11)]$$

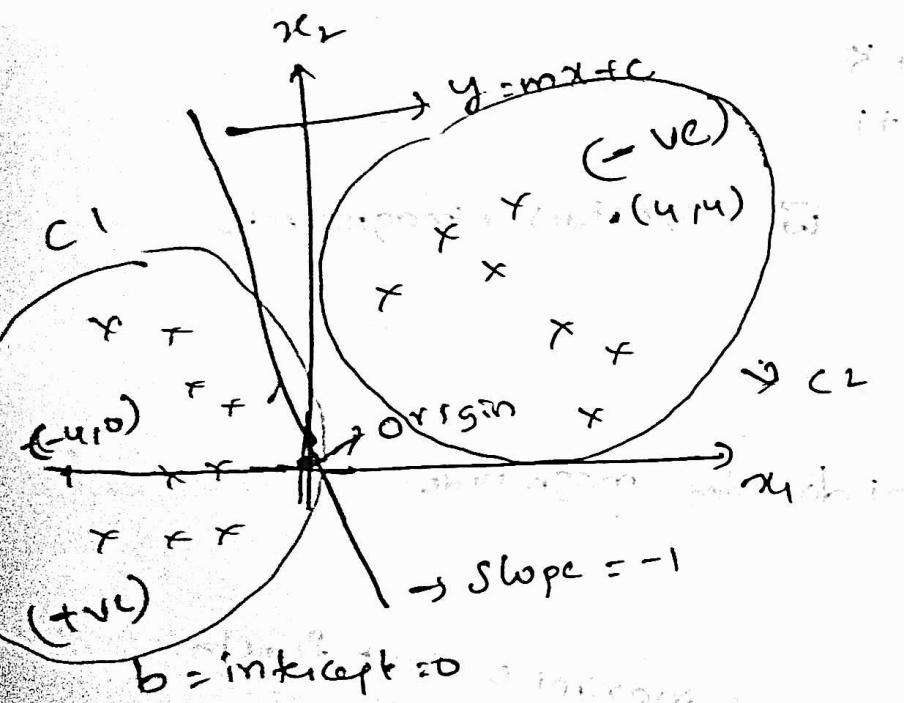
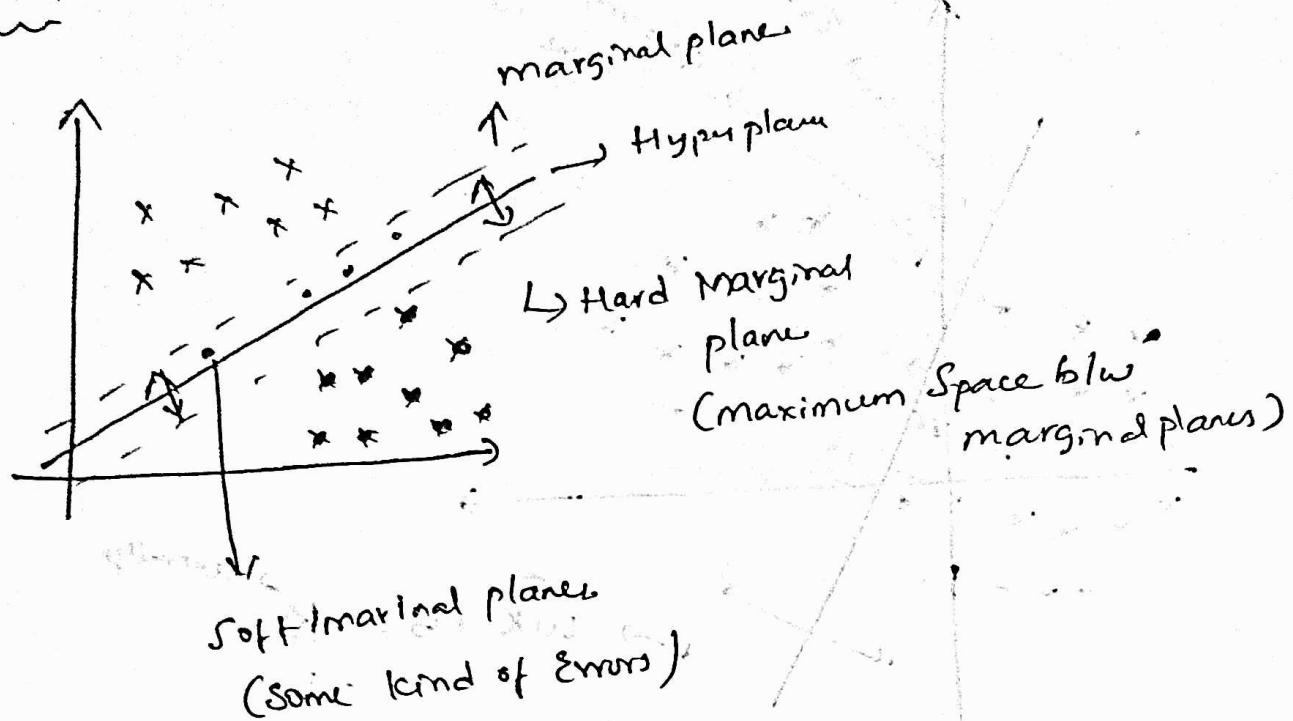
$$O/P = \text{Base} + \alpha_1 (P\tau_1) + \alpha_2 (P\tau_2) + \dots + \alpha_n (P\tau_n)$$

model

Avg.

→ It is also a Black Box Model.

SVM



$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} -4 & 0 \end{bmatrix}$$

$= -4 + \text{tvc Value}$

$$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} \begin{bmatrix} 4 & 0 \end{bmatrix}$$

$= -4 + 0 = -4 \therefore -\text{vive Value}$

$$0 = a_1 x_1 + a_2 x_2 + b$$

Slope

$$y = mx + c \rightarrow \text{Intercept}$$

$$a_1 x_1 + a_2 x_2 + b = 0 \rightarrow \text{Eq of straight line}$$

$$by = -a_1 x_1 - a_2 x_2$$

$$y = -\frac{a_1}{b} x_1 - \frac{a_2}{b} x_2$$

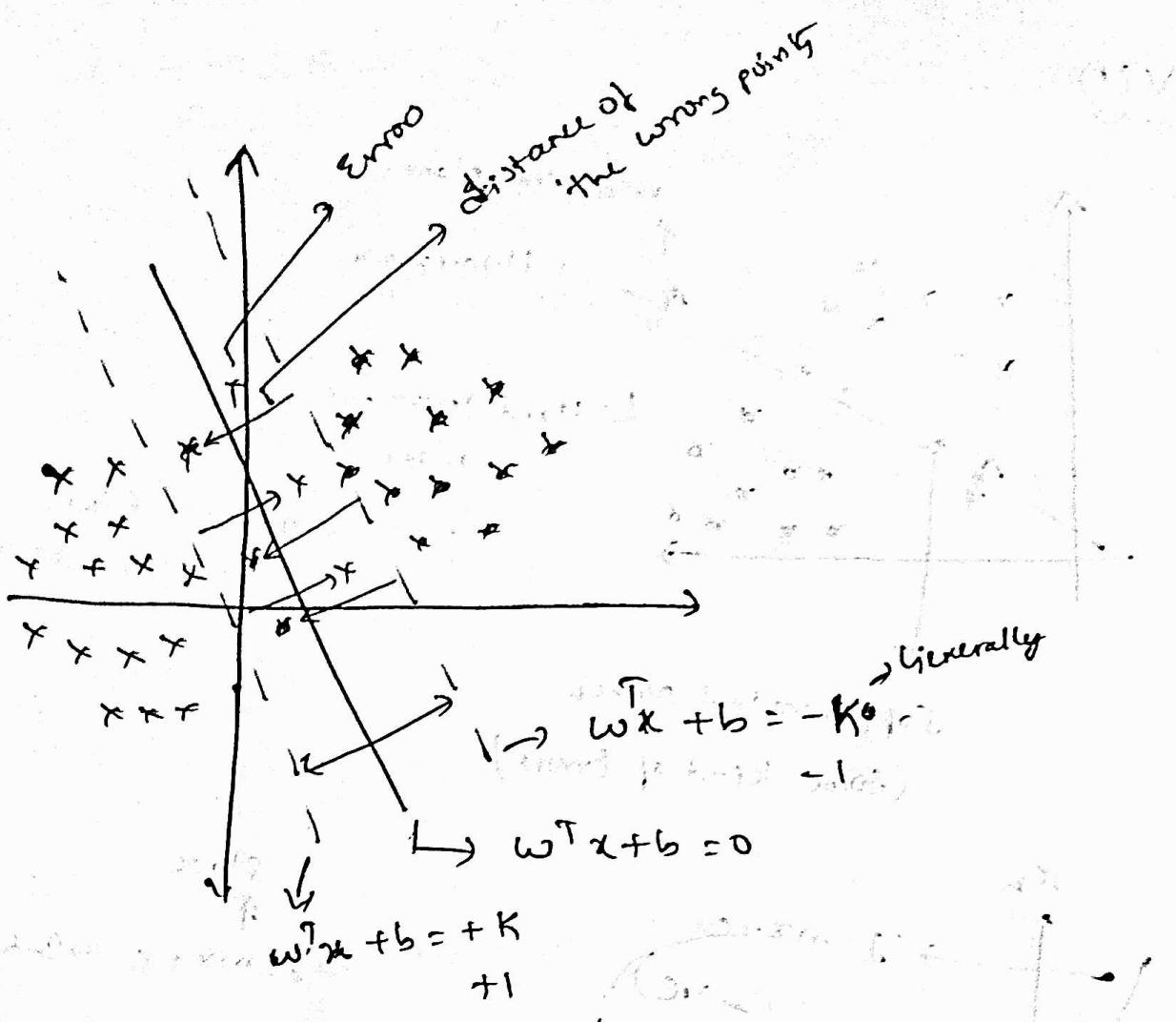
Slope

$$m = -\frac{a_1}{b} \quad c = -\frac{c}{b}$$

$$y = mx + c$$

$$y = w_1 x_1 + w_2 x_2 + \dots + b$$

$$\boxed{\therefore y = w^T x + b}$$



$\vec{w}$  vectors + magnitude

$$w^T x_1 + b = 1$$

$$\underline{\underline{(+)}} \quad w^T x_2 + b = -1 \quad \underline{\underline{(-)}}$$

$$\frac{w^T(x_1 - x_2)}{\|w\|} = 2 \rightarrow \text{dividing magnitude}$$

$$\Rightarrow \text{maximize } (\omega, b) \quad \frac{2}{\|w\|} \Rightarrow \text{maximize cost function}$$

$$y_i \begin{cases} +1 & , w^T x_i + b \geq 1 \\ -1 & , w^T x_i + b \leq -1 \end{cases}$$

major Aim  $y_i * (w^T x_i + b) \geq 1$

for correct points

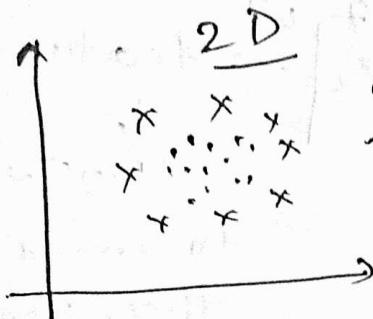
$$\text{maximum}_{(w,b)} \frac{2}{\|w\|} \Rightarrow \min_{(w,b)} \frac{\|w\|}{2}$$

Regularization

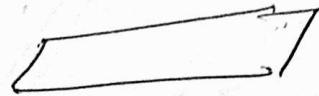
$\Rightarrow \min_{(w,b)} \frac{\|w\|}{2} + C \sum_{i=1}^n \xi_i$

value of the Error  
Sum of distance of the wrong data points  
How many Errors we can have

### SVM kernel



convert  
→ 3 dimensional

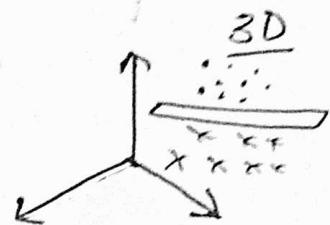


### kernel

converts

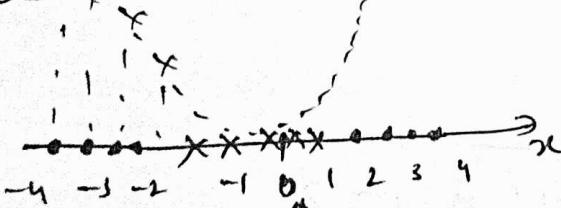
low dimension

to high dimensionality

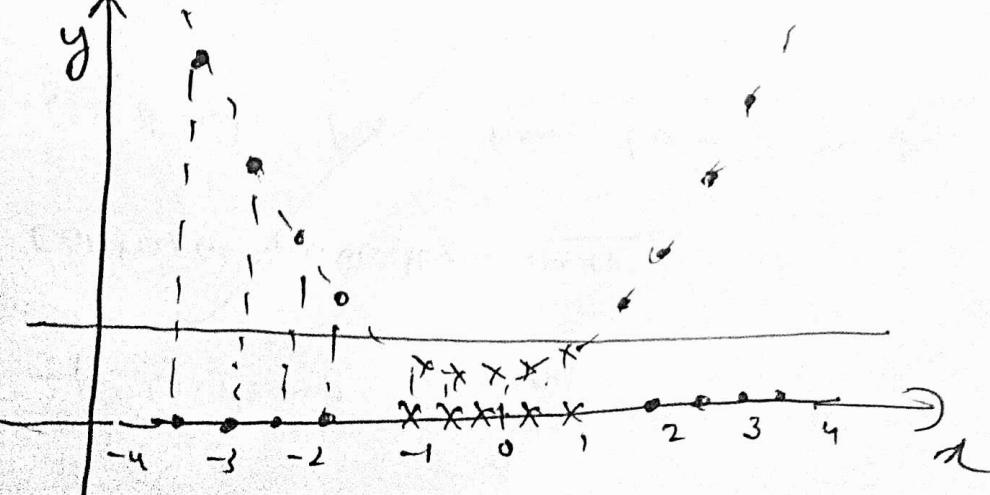


1-Dimension:

$$y = f(x) = x^2$$



$1D \rightarrow 2D$



## Polynomial Kernel:

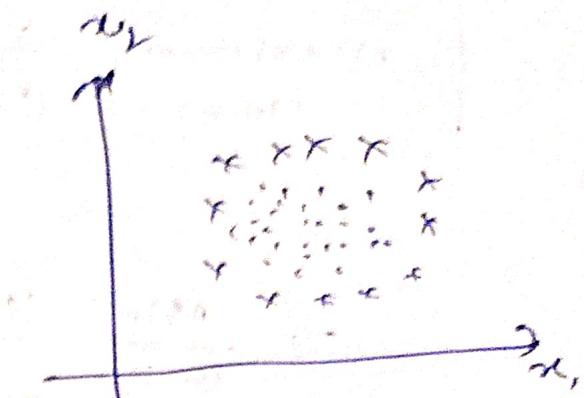
data  $x_1 \ x_2 \ y_0 \ y_1 \ f(x_1, x_2)$

$$f(x_1, x_2) = (x_1^T \cdot x_2 + 1)^d$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} x_1 & x_2 \end{bmatrix}$$

$$= \begin{bmatrix} x_1^2 & x_1 \cdot x_2 \\ x_2 \cdot x_1 & x_2^2 \end{bmatrix}$$

$$\rightarrow x_1 \ | \ x_2 \ | \ y_0 \ | \ x_1^2 \ | \ x_2^2 \ | \ x_1 \cdot x_2^2 \ | \ x_1^2 \cdot x_2^2$$



Find out

$$x_1^2 \ x_2^2 \ x_1 \cdot x_2$$

$$x_2$$

$$x_1 \cdot x_2$$

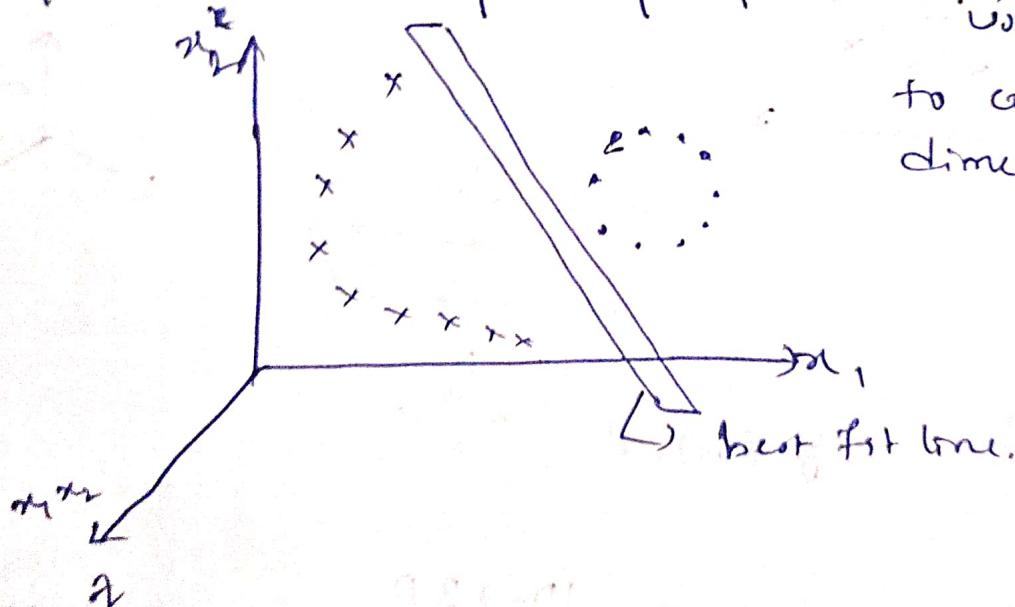
$$x_2^2$$

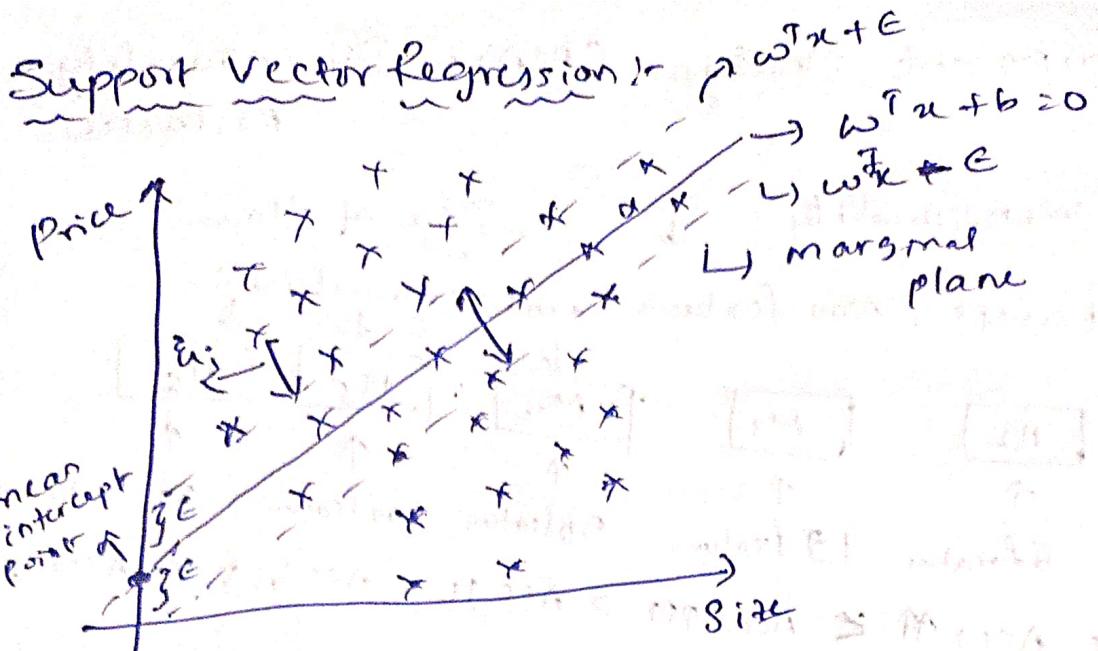
$$x_1^2 \cdot x_2^2$$

features

Best fit  
line

Using these features  
to convert low  
dimensional to high





Cost function:

$$\min_{(\omega, b)} \frac{\|\omega\|}{2}$$

Constraint

$$|y_i - \omega^T x_i| \leq \epsilon \rightarrow \text{It is very good}$$

Cost function:  $\min_{(\omega, b)} \frac{\|\omega\|}{2} + C_i \sum_{j=1}^n |\xi_{ij}| \rightarrow \text{Hinge loss.}$

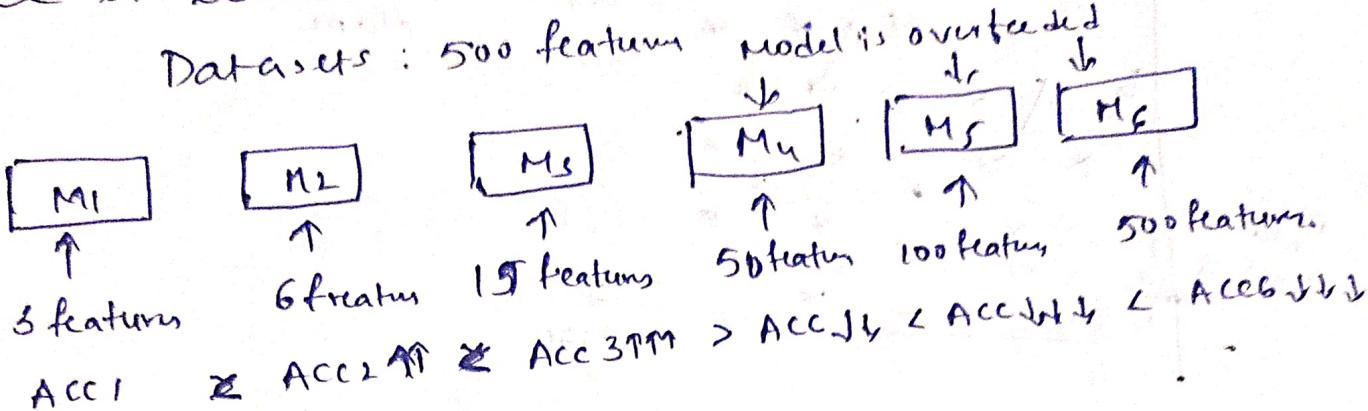
$\hookrightarrow$  Hyperparameter

constraint:

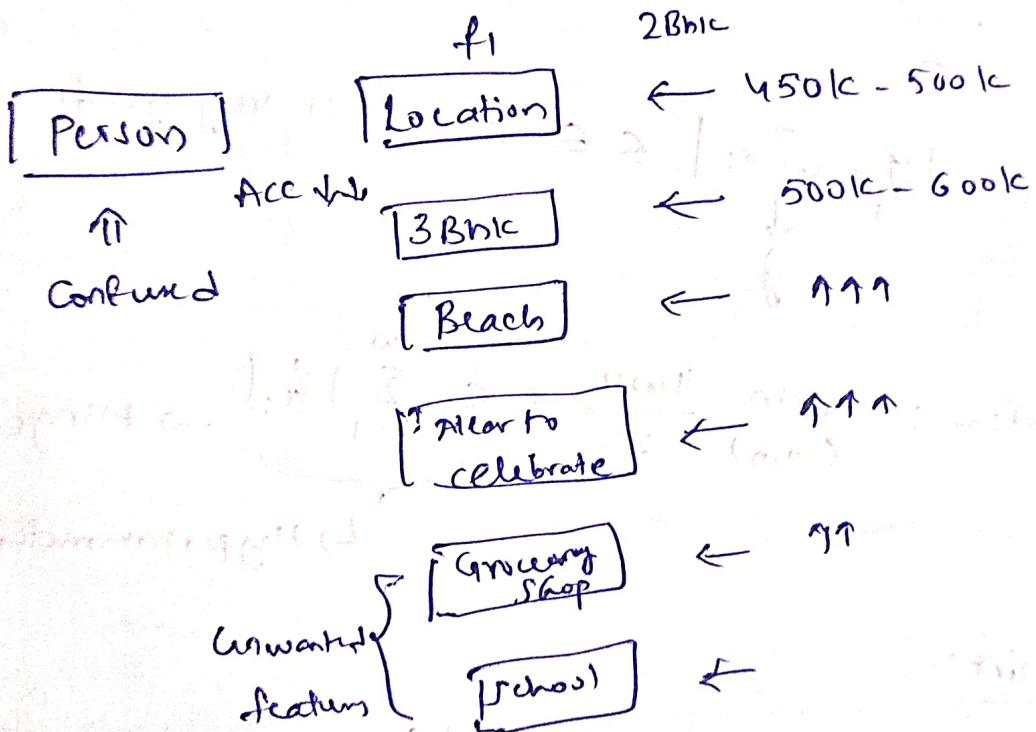
$$|y_i - \omega^T x_i| \leq \epsilon + \xi_i$$

# Principal Component Analysis (PCA) :- Dimensionality Reduction

## 1. Curse of Dimensionality :-



## 2. Model performance Degradation - when we use more features (Dimensions)



Two different ways to remove curse of Dimensionality

1. Feature Selection

↓  
takes imp features

2. Dimensionality Reduction

↓  
Feature Extraction

# feature Selection Vs feature Extraction:-

Dimensionality

Reduction.

## 1. why Dimensionality Reduction?

- \* Prevent  $\rightarrow$  Curse of Dimensionality
- \* Improve the Performance of the model
- \* Visualize the data  $\rightarrow$  Understand the data.

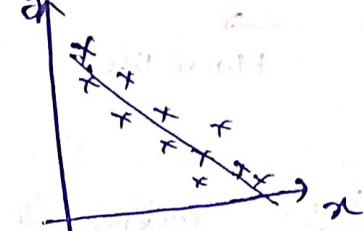
### feature Selection:-

1lp olp

	+ve	-ve
x	$x \uparrow y \uparrow$	$x \uparrow y \downarrow$
-	$x \downarrow y \downarrow$	$x \downarrow y \uparrow$
y		

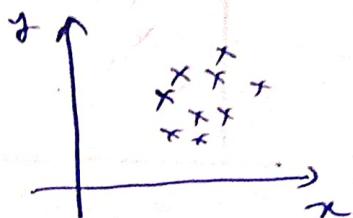


	+ve	-ve
x	$x \uparrow y \downarrow$	$x \downarrow y \uparrow$
-	$x \downarrow y \downarrow$	$x \uparrow y \uparrow$
y		



No relationship

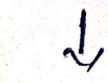
,  $x \nparallel y$



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

$$\left. \begin{array}{l} = 0, \text{No relation} \\ = +\text{ve}, \text{positive cov} \\ = -\text{ve}, \text{negative cov} \end{array} \right\}$$

\* Correlation (Pearson) =  $\frac{\text{cov}(x, y)}{\sigma_x \sigma_y} = -1 \text{ to } +1$



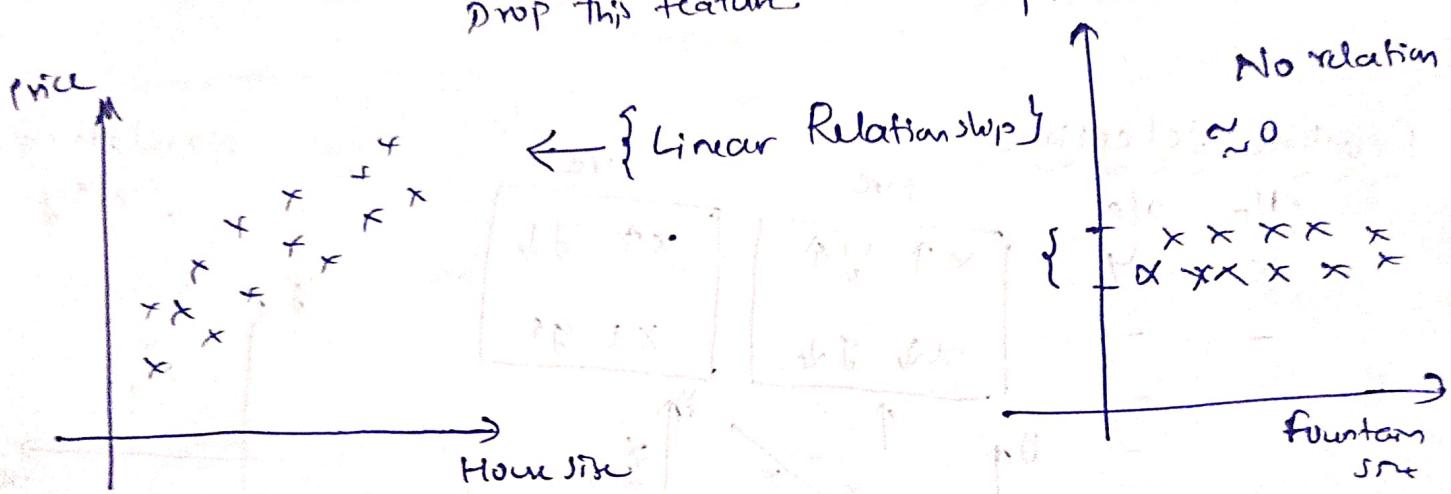
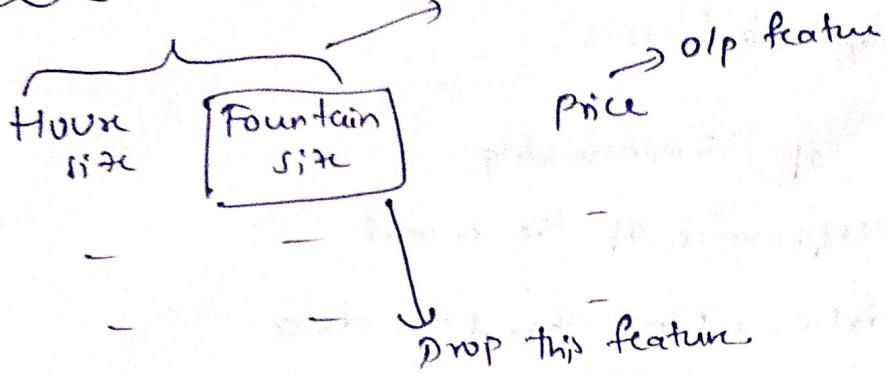
\* The more toward the value of +1 then more pos correlated

\* The more toward the value of -1 then more neg correlated

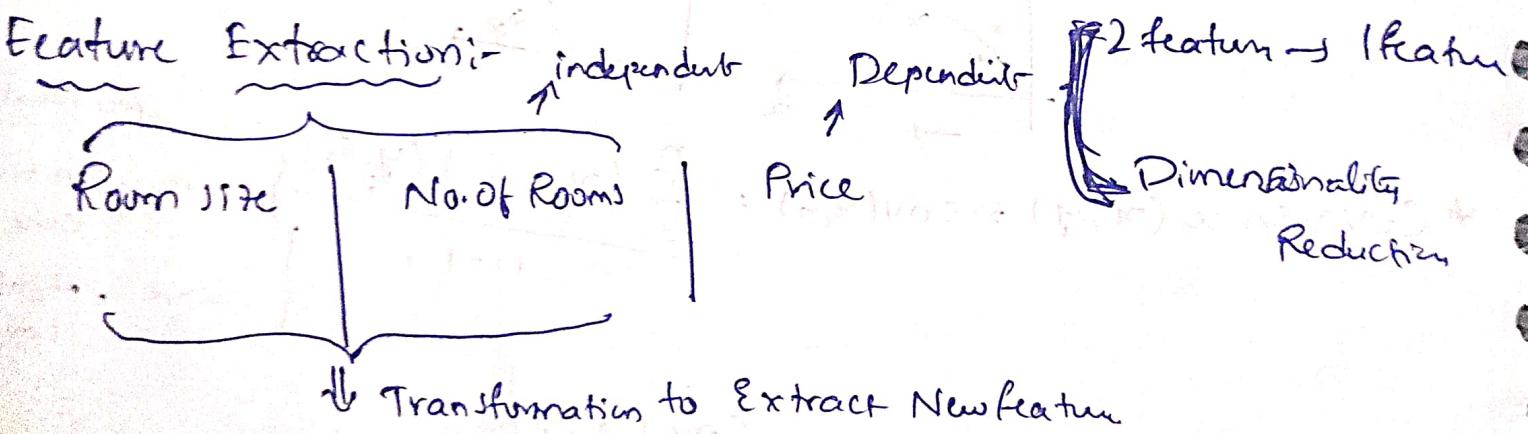
\* The more toward the value of 0 is zero then no correlation.

Ex:-

Data set Housing :- Independent feature



Feature Extraction:-



House size

Price.

# PCA Geometric Intuition:

House Dataset

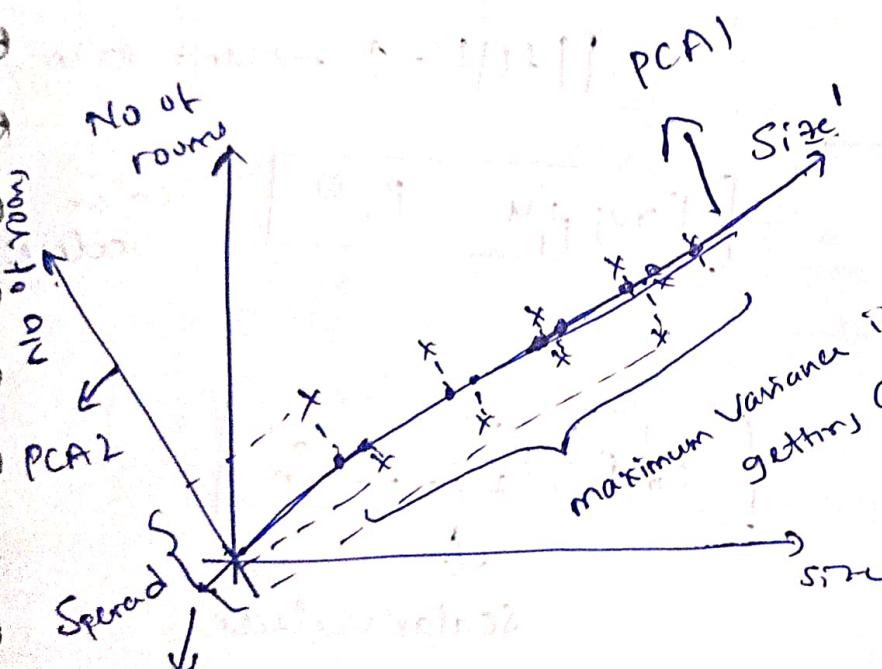
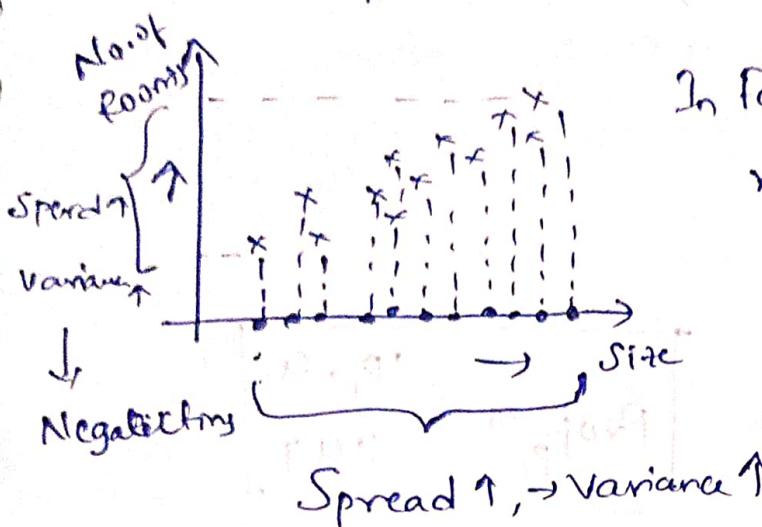
Size (Haus)	No. of Rooms	Price
10	1	100
15	2	150
20	3	200
25	4	250
30	5	300
35	6	350
40	7	400
45	8	450
50	9	500
55	10	550

PCA

2 Dimension  $\rightarrow$  1 Dimension

In Feature Extraction

$\Rightarrow$  Loss of information (No. of Rooms)



Eigen Decomposition on matrix

Transformation

2D  $\rightarrow$  1D

much information is lost

\* 2D  $\xrightarrow{\text{axis are}} \text{PC1}, \text{PC2}$

$\text{Var}(\text{PC1}) > \text{Var}(\text{PC2})$

To get the best principal

component which captures maximum variance.

3D  $\rightarrow$  1D  $\rightarrow$  3D  $\rightarrow$  1D.

$\boxed{\text{PC1}}$   $\text{PC2}, \text{PC3}$

3D  $\rightarrow$  2D

2D

$\text{Var}(\text{PC1}) > \text{Var}(\text{PC2}) > \text{Var}(\text{PC3})$

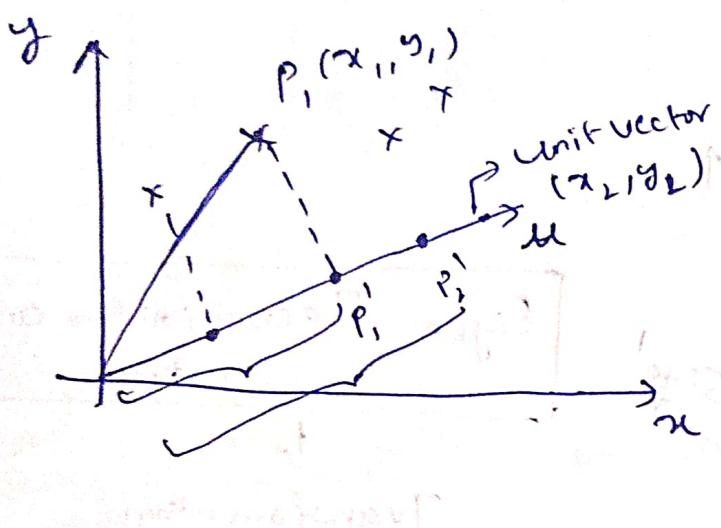
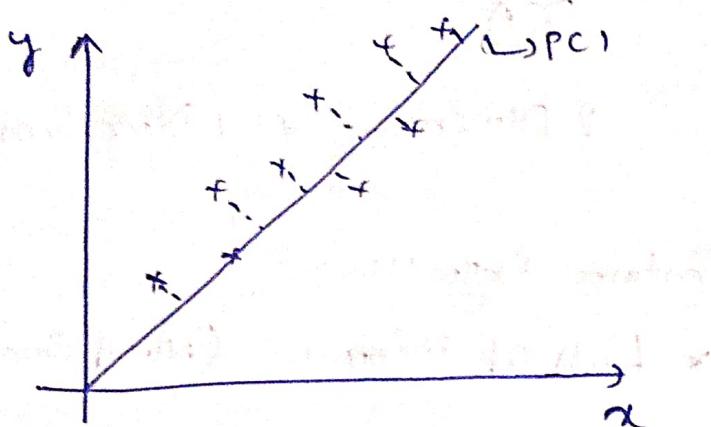
$\text{Var}(\text{PC1}) > \text{Var}(\text{PC2}) > \text{Var}(\text{PC3})$

# Maths intuition behind PCA Algorithm:-

$2D \rightarrow 1D$

1. Projection ✓

2. Cost function  $\rightarrow$  Variance ✓

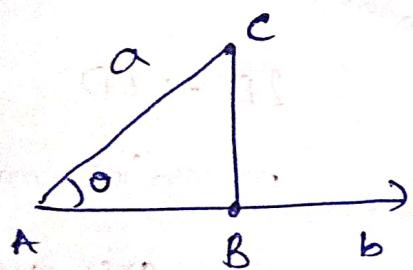


$$\text{Proj}_{P_1} u = \frac{P_1 \cdot u}{\|u\|}$$

$\Rightarrow \|u\| = 1 \rightarrow \text{unit vector}$

$$\text{Proj}_{P_1} u = P_1 u$$

$\rightarrow$  Scalar  
vector



$$P_0^1, P_1^1, P_2^1, \dots, P_n^1$$

$\downarrow$   
scalar values.

$$x_0^1, x_1^1, x_2^1, \dots, x_n^1$$

$$\therefore \text{Cost function} \rightarrow \text{Variance} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

Goal - find the best unit vector which captures the maximum variance

# Eigen Vector & Eigen Values:

$$A\mathbf{v} = \lambda \mathbf{v}$$

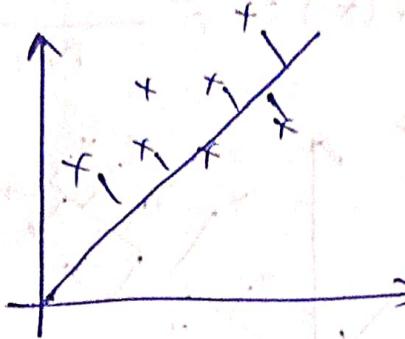
$\hookrightarrow$  Linear Transformation of matrix.

1. Covariance matrix b/w features

2. Eigen Vector & Eigen Values will be found out from this covariance matrix.

3. Eigen Vector  $\rightarrow$  larger Eigen Value  $\rightarrow$  magnitude of the Eigen vector

$\downarrow$   
Capturing the maximum variance



[Eigen Decomposition of Covariance matrix]

Eigen vector & Eigen Values.

$$\begin{bmatrix} A & v \end{bmatrix} * \begin{bmatrix} v \end{bmatrix} = \lambda \cdot \begin{bmatrix} v \end{bmatrix}$$

Eigen

Value

matrix.

$$A * v = \lambda * v$$

$\downarrow$   
Eigen Vector  $\rightarrow$  Maximum Magnitude

$\downarrow$   
use as principal component

$\downarrow$   
maximum Variance

# Steps to calculate Eigen Value & Vectors!

## 1. Covariance of features

$$\begin{bmatrix} x, y \\ \downarrow \\ x' \end{bmatrix} \quad z$$

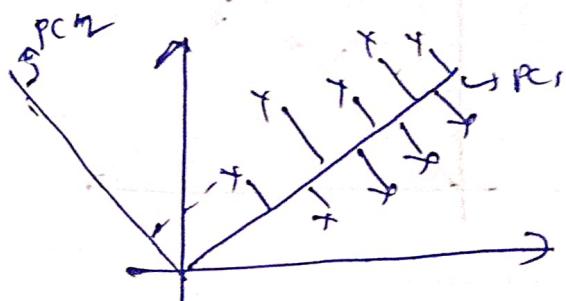
$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{N-1}$$

$2 \times 2$

	x	y
x	$\text{Var}(x)$	$\text{Cov}(x, y)$
y	$\text{Cov}(y, x)$	$\text{Var}(y)$

$$\text{Cov}(x, x) = \text{Var}(x) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N-1}$$

	x	y	z
x	$\text{Var}(x)$	$\text{Cov}(x, y)$	$\text{Cov}(x, z)$
y	$\text{Cov}(y, x)$	$\text{Var}(y)$	$\text{Cov}(y, z)$
z	$\text{Cov}(z, x)$	$\text{Cov}(z, y)$	$\text{Var}(z)$



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

$$\begin{bmatrix} \lambda_1 \\ \downarrow \\ \text{PC}_1 \end{bmatrix} \quad \begin{bmatrix} \lambda_2 \\ \downarrow \\ \text{PC}_2 \end{bmatrix}$$

$$\lambda_1, \lambda_2, \lambda_3$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$\text{PC}_1 \quad \text{PC}_2 \quad \text{PC}_3$$

$$\checkmark \quad \checkmark$$

$$\frac{1D}{\downarrow} \quad 1D$$

$$3D - 2D$$

$$3D - 1D$$

$$\lambda_1, \lambda_2, \lambda_3$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$\text{PC}_1 \quad \text{PC}_2 \quad \text{PC}_3$$

Projection:

$$\begin{bmatrix} x_1 & x_2 \end{bmatrix} \rightarrow \text{Eigen Value}$$

$$\downarrow$$

$$\begin{bmatrix} \text{PC}_1, \text{PC}_2 \end{bmatrix}$$