

→ Relationship Diagram

Linear Regression

Simple L.R (one-one)

$$y = \underbrace{\alpha_0 + \alpha_1 x_1}_{\text{regression coefficient}} + \text{error}$$

Actual value - Predicted value

If there is one independent variable that affect the dependant variable then the regression is called Simple L.R.

multiple L.R (many-one)

$$y = \underbrace{\alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_m x_m}_{\downarrow} + \text{error}$$

Many Independent variable
One Dependent variable.

When we have more than one independent variable that effect one dependent variable is called multiple linear Regression.

→ Linear Regression with least square criterion.

→ Best fit line equation :

$$\hat{y}^{(i)} = w_0 + w_1 x_1^{(i)} + \dots + w_n x_n^{(i)} ; i=1 \dots n$$

where w_0, w_1, \dots, w_n are the regression coefficient.

The method of linear regression is to choose the coefficients w_0 to w_n to minimize the residual sum of squares of these differences over all the n training instances.

The performance criteria is the sum of error squares.

$$\text{Residual error : } e^{(i)} = y^{(i)} - \hat{y}^{(i)} = y^{(i)} - \sum_{j=0}^n w_j x_j^{(i)} ; w_0^{(i)} = 1$$

Residual sum of error square:

$$E = \sum_{i=1}^N (e^{(i)})^2 = \sum_{i=1}^N \left(y^{(i)} - \sum_{j=0}^n w_j x_j^{(i)} \right)^2 ; w_0^{(i)} = 1$$

•> Minimal Sum of Error Square

$$X = \begin{vmatrix} 1 & 1 & \dots & 1 \\ x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(N)} \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(N)} \\ \vdots & \vdots & & \vdots \\ x_n^{(1)} & x_n^{(2)} & \dots & x_n^{(N)} \end{vmatrix} = [\bar{x}^{(1)} \bar{x}^{(2)} \dots \bar{x}^{(N)}]$$

$$y = [y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(N)}]^T$$

$$\bar{\omega} = [\omega_0 \ \omega_1 \ \omega_2 \ \dots \ \omega_N]^T$$

$$\bar{x} = [1 \ x_1 \ x_2 \ \dots \ x_n]^T$$

$$y^{(1)} = \bar{\omega}^T \bar{x}^{(1)}$$

$$y^{(N)} = \bar{\omega}^T \bar{x}^{(N)}$$

Therefore,

$$[y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(N)}] = \bar{\omega}^T [\bar{x}^{(1)} \bar{x}^{(2)} \ \dots \ \bar{x}^{(N)}] = \bar{\omega}^T x$$

or

$$\hat{y} = (\bar{\omega}^T x)^T$$

The vector of residual or becomes

$$\text{error} = y - (\bar{\omega}^T x)^T$$

Hence the error function can be written as

$$\begin{aligned} E(\bar{\omega}) &= [y - (\bar{\omega}^T x)^T]^T [y - (\bar{\omega}^T x)^T] \\ &= \bar{\omega}^T [x x^T] \bar{\omega} - 2 \bar{\omega}^T x y + y^T y \end{aligned}$$

→ In the least square estimation the objective is to find the optimum value of ω which is written as $\bar{\omega}^*$ which minimizes the error function.

$$\frac{\partial E(\omega)}{\partial \omega} = 2(x x^T) \bar{\omega} - 2x y = 0$$

This gives

$$\bar{\omega}^* = (x x^T)^{-1} x y$$

The fitted output values at the training data are

$$\hat{y} = x^T \bar{\omega}^* = x^T (x x^T)^{-1} x y$$

$$\bar{\omega}^* = x^T y$$

Q)

y	3	2	4	5	8
x_1	2	3	5	7	8
x_2	1	5	3	6	7

Find the estimated value of y .

$y \quad x_1 \quad x_2 \quad \text{equation}$

$$\begin{array}{lll} 3 & 2 & 1 \\ 2 & 3 & 5 \\ 4 & 5 & 3 \\ 5 & 7 & 6 \\ 8 & 8 & 7 \end{array} \quad \left. \begin{array}{l} 3 = 1\omega_0 + 2\omega_1 + 1\omega_2 + \epsilon_1 \\ 2 = 1\omega_0 + 3\omega_1 + 5\omega_2 + \epsilon_2 \\ 4 = 1\omega_0 + 5\omega_1 + 3\omega_2 + \epsilon_3 \\ 5 = 1\omega_0 + 7\omega_1 + 6\omega_2 + \epsilon_4 \\ 8 = 1\omega_0 + 8\omega_1 + 7\omega_2 + \epsilon_5 \end{array} \right\} \begin{array}{l} \text{error} \\ \text{It is to be calculated} \\ \text{for determining the} \\ \text{matrix property} \end{array}$$

$$\underbrace{\begin{pmatrix} 3 \\ 2 \\ 4 \\ 5 \\ 8 \end{pmatrix}}_y = \underbrace{\begin{pmatrix} 1 & 2 & 1 \\ 1 & 3 & 5 \\ 1 & 5 & 3 \\ 1 & 7 & 6 \\ 1 & 8 & 7 \end{pmatrix}}_X \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix}}_{\beta} + \underbrace{\begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \epsilon_4 \\ \epsilon_5 \end{pmatrix}}_{\epsilon}$$

$$X = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 2 & 3 & 5 & 7 & 8 \\ 1 & 5 & 3 & 6 & 7 \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 \\ 1 & 3 & 5 \\ 1 & 5 & 3 \\ 1 & 7 & 6 \\ 1 & 8 & 7 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 5 & 25 & 22 \\ 25 & 151 & 130 \\ 22 & 130 & 120 \end{bmatrix}_{3 \times 3}$$

$$\Rightarrow (X^T X)^{-1} = \frac{1}{|X^T X|} \text{Adj}(X^T X)$$

$$= \begin{bmatrix} 1.201 & -0.138 & -0.071 \\ -0.138 & 0.114 & -0.098 \\ -0.071 & -0.098 & 0.128 \end{bmatrix}$$

$$= (X^T X)^{-1} X Y =$$

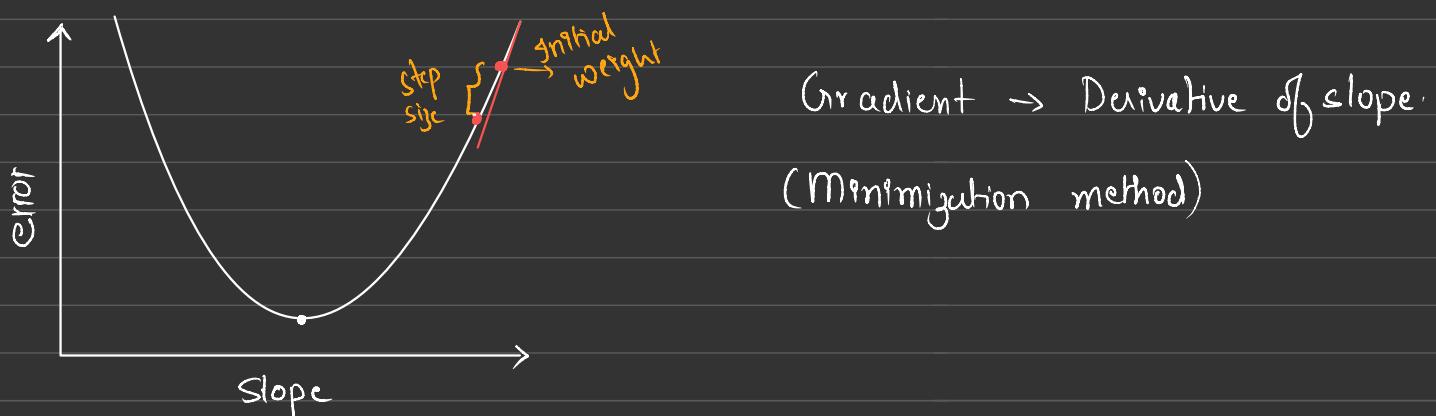
$$X^T Y = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ 2 & 3 & 5 & 7 & 8 \\ 1 & 5 & 3 & 6 & 7 \end{pmatrix}_{3 \times 5} \begin{pmatrix} 3 \\ 2 \\ 4 \\ 5 \\ 6 \end{pmatrix}$$

$$w^* = (X^T X)^{-1} X^T Y = X^T Y = \begin{pmatrix} 22 \\ 31 \\ 111 \end{pmatrix}_{3 \times 1}$$

$$\begin{pmatrix} 0.50 \\ 1 \\ -0.25 \end{pmatrix} \quad Y = \beta_0 + \beta_1 u_1 + \beta_2 u_2$$

$$Y = 0.5 + 1. u_1 - 0.25 u_2$$

•) Gradient Descent method (Error Detection)



\rightarrow It is a first order iterative optimization algorithm for finding the minimum of the error function. To achieve this goal it performs two steps iteratively.

Step 1: Compute the slope or the 1st order derivative of the current point

Step 2: Move in opposite direction of the slope starting from the current point by the computed amount

- └ Step size
- └ Learning rate (η)

Two methods to check.

- 1) Batch Gradient Descent
- 2) Stochastic Gradient Descent

\rightarrow In Batch gradient D. all the training data are taken into account to take a single step, then we take the average of the gradient to update the weight parameter.

\rightarrow Stochastic : Here we update the parameter after every single observation & every time the weights are updated which is known as iteration.

Batch

:

Stochastic

→ Used for Smaller Data set → Used for Larger Dataset.

* Derivation of Gradient Descent method.

∇^* to minimize the error function

$$\overline{w}_{k+1} = \overline{w}_k - \eta \frac{\partial E}{\partial \overline{w}}|_k$$

\downarrow new value \downarrow old value

$$w_{new} = w_{old} - \eta * \underbrace{\frac{\partial E}{\partial w}}_{learning\ rate}|_k$$

step size

} Training rule
of the
gradient Descent.

In this method we start with an arbitrary initial weight vector then the weight vector parameters are updated after each iteration till the global minimum error is attained.

The gradient with respect to weight ω_j

$$\frac{\partial E}{\partial \omega_j} = \frac{\partial}{\partial \omega_j} \left[\frac{1}{2} \sum_{i=1}^N \left(y^{(i)} - \left(\sum_{j=1}^n \omega_j x_j^{(i)} + \omega_0 \right) \right)^2 \right]$$

Error $e^{(i)}$ for the i^{th} sample is given by

$$e^{(i)} = y^{(i)} - \left(\sum_{j=1}^n \omega_j x_j^{(i)} + \omega_0 \right)$$

$$\frac{1}{2} \sum_{i=1}^N \left(\frac{\partial}{\partial \omega} (e^{(i)})^2 \right) = \sum_{i=1}^N e^{(i)} \frac{\partial e^{(i)}}{\partial \omega_j}$$

$$\delta \frac{\partial e^{(i)}}{\partial \omega_j} = -x_j^{(i)}$$

$$\Rightarrow \frac{\partial E}{\partial \omega_j} = \frac{\partial}{\partial \omega_j} \left[\frac{1}{2} \sum_{i=1}^N (e^{(i)})^2 \right] \\ = - \sum_{i=1}^N \frac{1}{2} \times 2 e_i \times \frac{\partial e_i}{\partial \omega_j}$$

$$\frac{\partial E}{\partial \omega_j} = - \sum_{i=1}^N e_i x_j \\ = - \sum_{i=1}^N \left(y^{(i)} - \left(\sum_{j=1}^n \omega_j x_j^{(i)} + \omega_0 \right) \right) x_j$$

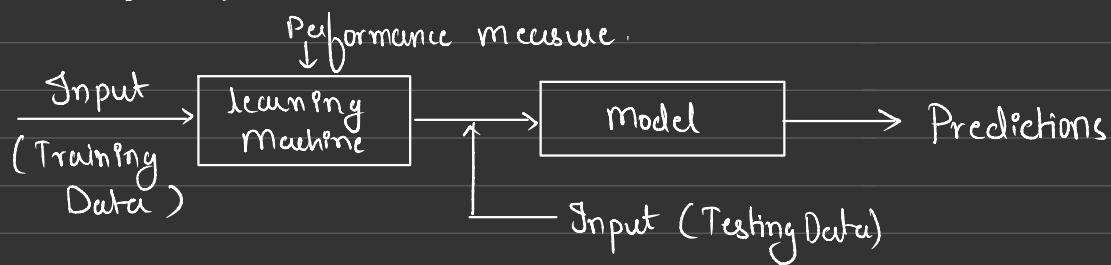
$$\frac{\partial E}{\partial \omega_0} = - \sum_{i=1}^N e^{(i)} = - \sum_{i=1}^N \left[y^{(i)} - \left(\sum_{j=1}^n \omega_j x_j^{(i)} + \omega_0 \right) \right]$$

$$\omega_j \leftarrow \omega_j + \eta \sum_{i=1}^N \left(y^{(i)} - \left(\sum_{j=1}^n \omega_j x_j^{(i)} + \omega_0 \right) \right) x_j$$

$$\omega_0 \leftarrow \omega_0 + \eta \sum_{i=1}^N \left(y^{(i)} - \left(\sum_{j=1}^n \omega_j x_j^{(i)} + \omega_0 \right) \right)$$

Unit - 4 (Generalisation)

•> Generalisation:



→ What is Generalisation / Concept of Generalisation:

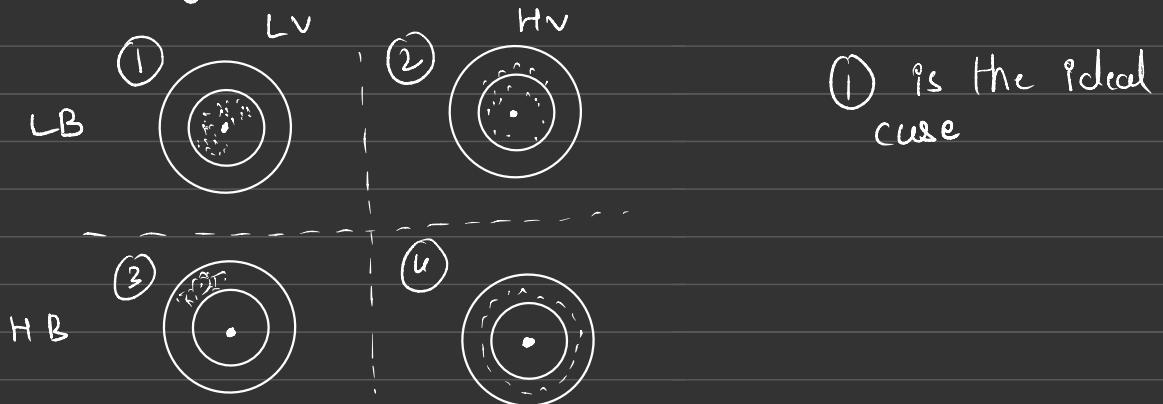
- ⇒ Generalisation refers to the model's ability to adopt properly to new data or previously new data drawn from the same distribution.
- ⇒ The generalization means how good our model is at learning from our given data and applying the learned information elsewhere.

•> Poor Generalisation

→ Factors

- 1) Bias & Variance
- 2) Underfitting & Overfitting.

1) Bias & Variance



Bias: Technically called as error.

Variance: How much data points are scattered between each other.

• function used for representing Bias & Variance

⇒ $h(x)$ is called as hypothesis function. It is the function that best maps input to output. It is also called approximation & mapping function.

$h(x) = \text{Mapping function} / \text{approximation function}$
 $f(x) = \text{True function.}$

$$\text{error}_D[h] = E_D \{ [h(x; D_j) - f(x)]^2 \}$$

E_D represents the expected value or the arithmetic mean of the large no. of independent values

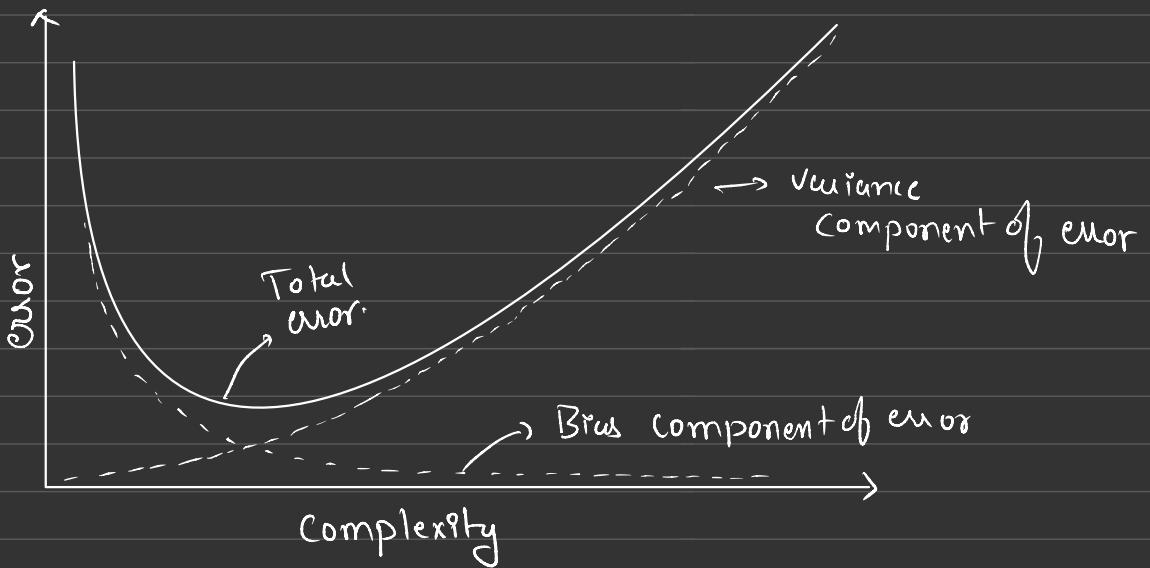
can be split into two funcn.

$$\Rightarrow \underbrace{[E_D \{ h(x; D_j) \} - f(x)]^2}_{(\text{bias})^2} + \underbrace{E_D \{ [h(x; D_j) - E_D(x; D_j)]^2 \}}_{\text{Variance}}$$

⇒ In the figure A in linear hypothesis function slope & intercept are two adjustable function if the experiment is repeated multiple times the estimates are scattered around 10.81, while the theoretical value is scattered around 8 for $x = 3$

⇒ For second case for a 5th order polynomial as there are many adjustable parameters the distribution of estimates is scattered while the theoretical value y is scattered around g . So it is understand that for a simple model bias is high & variance is low but for complex model bias is low & variance is high. This is known as Bias - variance tradeoff.

•> Graph of Bias - Variance Tradeoff.

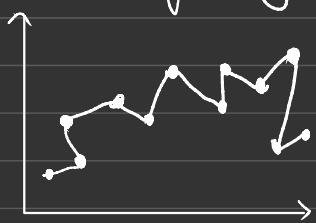


•> Occam's Razor Principle : The simpler explanations are more reasonable, and any unnecessary complexity should be shaved off.

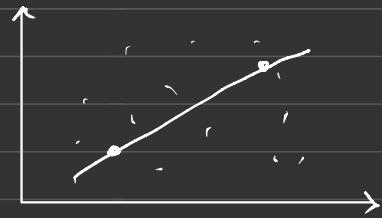
→ It suggests hypothesis function that avoid overfitting of the training data.

→ 'Simple' may imply needing lesser parameters, lesser training time, fewer attributes for data representation & less computational complexity

Overfitting



Underfitting



→ High variance
→ Overcomplicated
→ Covers all / maxm of the actual data

→ High bias
→ Oversimplifies
→ Line is not fitting properly.

•) Regularization

→ Regularization Technique optimizes hypothesis complexity for a given training data set.

→ Here we look for the weights that both decrease the error & shrink words zero.

$$Ex = y = 0.9 + 2u_1 + 20u_2 + \underbrace{39u_3}_{\text{Has predominant effect on } y}$$

* \rightarrow Ridge regularization (L_2 norm)

\rightarrow Lasso regularization (L_1 norm)

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

$$y = \beta_0 + \beta_1 u_1 + \beta_2 u_2 + \dots + \beta_p u_p$$

•) Ridge Regression:

→ Equation for ridge regression

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

λ = shrinkage penalty Ordinary least square

When $\lambda = 0$, RSS will produce OLS estimate

When $\lambda \rightarrow \infty$, the impact of shrinkage penalty grows & the ridge regression coefficient estimates will approach towards zero.

Ridge Vs OLS

$$L_{\text{ridge}}(\hat{\beta}) = \sum_{i=1}^n (y_i - x_i^\top \beta)^2 + \lambda \sum_{j=1}^m \hat{\beta}_j^2 = \|y - x\hat{\beta}\|^2 + \lambda \|\hat{\beta}\|^2.$$

Solving

$$\hat{\beta} = (x^\top x + \lambda I)^{-1} (x^\top y)$$

$$\hat{\beta}^* = \underbrace{(x x^\top)^{-1} x^\top y}_{\text{pseudomatrix}}$$

•) Cross Validation (λ)

$1 2 3 4 5 6 7 8 9$
<u>testing data</u> <u>training data</u>

$1 2 3 4 5 6 7 8 9$
training testing training

$1 2 3 4 5 6 7 8 9$
training testing

- ⇒ A commonly used technique for forecasting the success rate of a learning method taking into account a fixed data sample is called the K-Fold cross validation.
- ⇒ Here training & testing is done K times. The given is divided into K subsets or folds in iteration K one partition is set aside for testing and other partitions are employed to train the model.
- ⇒ Ultimately the K error estimate received from K iterations, we averaged to give rise to an overall error estimation.

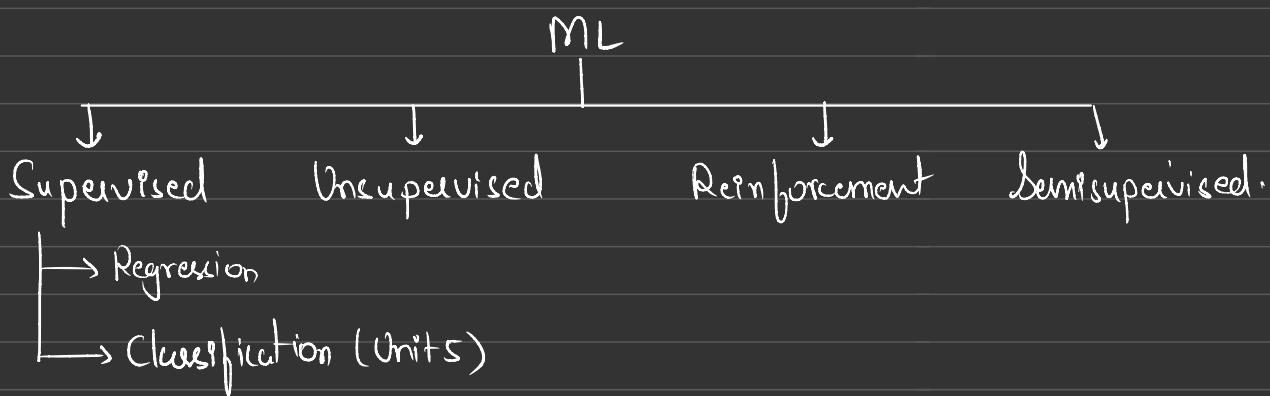
•) Lasso Regression:

↳ Least absolute Shrinkage & Selection Operator

→ Has feature selection capability.

•> Unit - 5 (Classification)

→



$P(A|B)$: Conditional Probability.

$P(A|B)$ → Prob of event A given that Prob of B is already given.

$$\Rightarrow \frac{P(A \cap B)}{P(B)}$$

Ex: 40 students belong to set A who likes CI. Then 30 student belongs to set B who likes ML 20 students like both CI & ML Calculate $P(A|B)$

$$\frac{20}{30} \Rightarrow \frac{2}{3} = 0.67$$

Bayes Theorem.

$$:- P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$\therefore P(A|B) \cdot P(B) = P(A \cap B) \quad \text{--- (1)}$$

$$\text{Similarly } P(B|A) = \frac{P(B \cap A)}{P(A)}$$

$$\therefore P(B|A) \cdot P(A) = P(B \cap A) \quad \text{--- (2)}$$

$$\therefore P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

posterior = Prior × likelihood } Here A is hypothesis
 Evidence B is Evidence (data)

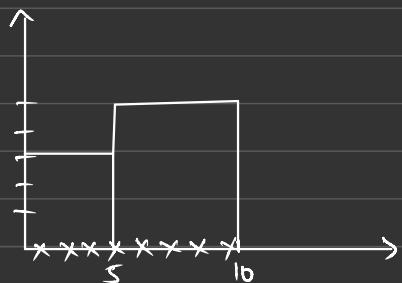
Q) $\text{Prob}(\text{King Face})$: find out that my given card is a King and it was given that it is a face card.

$$\Rightarrow \frac{4}{12} = \frac{1}{3} \Rightarrow \frac{\text{Prob(Face King)}}{\text{P(Face)}} \times \frac{P(K)}{P(Face)} = \frac{1 \times 4/52}{12/52} = 0.33$$

•) Probability Density function

- A function that represents a continuous probability distribution is called as Probability Density function.
- KNN methods estimates the density value at point x based on the distance between x and its k^{th} nearest neighbour.

•) Histogram: density estimation.



→ Origin = 0
→ Bin width = 5

consider a dataset
[2, 3, 5, 7, 9, 4, 6, 8]

- The Bin represents the input space is divided into equal size intervals named as Bins.

from	0 - 5	,	2	3	4
from	5 - 10	,	5	6	7, 8, 9

$$\hat{p}(x) = \frac{x^t \text{ in th same bin}}{N \times h} = \frac{\text{how many data points are in the same bin}}{\text{Total dataset} \times \text{bin width}}$$

→ The general formulation for density estimation is

$$\hat{P}(n) = \frac{k}{N V_k(n)}$$

→ k represents the number of neighbors.

→ V_k represents the volume of the neighborhood around n .

The class conditional density function

$$P(x|y_q) = \frac{k_q}{N_q V_k(n)}$$

⇒ k_q represents the number of neighbors out of k nearest that belongs to class y_p

N_q represents no of samples that belongs to class y_q

V_k represents volume of hyperspace

$$P(y_q | n) = \frac{P(n|y_q) \cdot P(y_q)}{P(n)}$$

$$\therefore \frac{\frac{k_q}{N_q \times V_k(n)}}{\frac{1}{N V_k(n)}} = \frac{k_q}{k}$$

$$P(y_q | n) = \frac{k_q}{k}$$

The posterior probability of an unknown sample n belonging to a specific class q can be calculated as the ratio of k_q to k

Q) Given a dataset with two classes where 1st class [0, 2, 3] and 2nd class are [6, 8, 10] estimate the class label for n = 5.

use kNN method with k = 3.

$$\Rightarrow P(y_1 | n=5) = \frac{k_1}{k} = \frac{2}{3} = 0.67$$

$$P(y_2 | n=5) = \frac{k_2}{k} = \frac{1}{3} = 0.33$$

∴ n = 5 belongs to 1st class for k = 3.

n = 5 belongs to 2nd class for k = 1.

This is because classification in kNN method depends on the samples in the feature space.

If the samples are scattered or if the features are high dimensional the density estimation with kNN becomes problematic.

If K varies prob will vary

•) kNN classification:

Problem: $k = 3$

	Age	loan	Default	Distance	x_5 Dist
1	25	40000	N	162000	0.7652
2	35	60000	N	82000	0.5200
3	45	80000	N	62000	0.3160
4	20	20000	N	122000	0.9245
5	35	120000	N	22000	2
6	52	18000	N	124000	0.6220
7	23	95000	Y	67000	0.6669
8	40	62000	Y	80000	0.4437
9	60	10000	Y	62000	3
10	48	220000	Y	78000	0.3861
11	33	150000	Y	8000	1
	48	142000	?		

$$D = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2} = \sqrt{(48 - 33)^2 + (142000 - 150000)^2} = 8000$$

$k=3$ we have 2 Y & 1 N

[∴ prediction = Y]

•) While calculating the distance Directly from the training set where variables have different measurement scale there is a influence on the calculated distance.

→ The solution is to standardize the training dataset.

•> Method of Standardization

$$x_s = \frac{x - \text{min}}{\text{max} - \text{min}} \quad \left. \right\} \text{standardize}$$

$$= \frac{150000 - 18000}{220000 - 18000} = 0.65$$

Unit 6: Clustering:

Unsupervised Learning

- ⇒ Unsupervised machine learning learns from the unlabeled data. It automatically finds the pattern between the training dataset and put the similar dataset in one group and another dataset in another group.
- ⇒ There is no specific methodology.

Clustering

i) It is an unsupervised task which involves the grouping of the similar data points.

Association

It is an unsupervised task that is used to find the implicit relationship b/w data points

•> Clustering By K-means Approach.

Sno.	Height	Weight
1	185	72
2	170	56
3	168	60
4	179	68
5	182	72
6	188	77
7	180	71
8	180	70
9	183	84
10	180	88
11	180	67
12	177	76

Q4 Divide the training dataset into two clusters where cluster C_1 is given $(185, 72)$ 1st row & $C_2 (170, 56)$ and row by using k-means approach by taking Euclidean distance.

→ Step 1 :- Check that the remaining data points are belonging to cluster 1 or cluster 2 by using euclidean distance.

$$\therefore \sqrt{(x_o - x_c)^2 + (y_o - y_c)^2}$$

x_o = Observed value of height.

x_c = Centroid value of particular cluster $(185, 72)$ & $(170, 56)$

y_o = Observed value of weight

y_c = Centroid value of weight

ED for 3rd row.

$$C_1 = \sqrt{(168 - 185)^2 + (60 - 72)^2} \\ = 20.8$$

$$C_2 = \sqrt{(168 - 170)^2 + (60 - 56)^2} \\ = 4.47$$

The shortest distance decides the cluster.

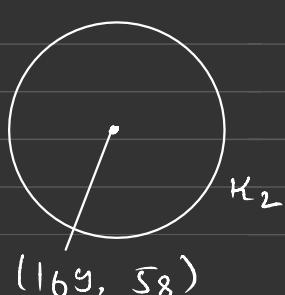
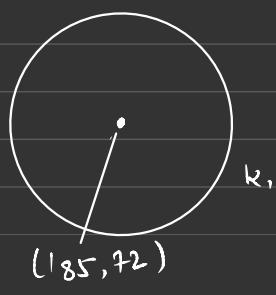
∴ 3rd row belongs to C_2

Once a row goes to a particular cluster, always we have to compute new cluster.

Step 2 :- Re computing of centroid

$$\text{New centroid for } C_2 = \left(\frac{170 + 168}{2}, \frac{56 + 60}{2} \right)$$

$$C_2 = (169, 58)$$



Euclidean Distance for 4th row :-

$$K_1 = \sqrt{(179-185)^2 + (68-72)^2} \\ = 7.21$$

$$K_2 = \sqrt{(179-169)^2 + (68-58)^2} \\ = 14.14$$

4th row belongs to cluster 1 \therefore New centroid for cluster 1

$$\left(\frac{185+179}{2}, \frac{72+68}{2} \right)$$

$$\text{Centroid for } C_1 = (182, 70) \quad \{ K_1 = [1, 4] \}$$

$$\text{Centroid for } C_2 = (169, 58) \quad \{ K_2 = [2, 3] \}$$

ED for 5th row,

$$K_1 = \sqrt{(182-182)^2 + (72-70)^2} = 2 \\ K_2 = \sqrt{(182-169)^2 + (72-58)^2} = 19.10$$

5th row belongs to cluster 1

1

2

3

Repeat these steps

Summary:

- 1) Decide the cluster.
- 2) Initialize the centroid
- 3) Calculate other datasets distance w.r.t the clusters.
- 4) Recomputing the centroid.
- 5) Stopping Criteria: Here k-means clustering stops when there is no change in centroid.