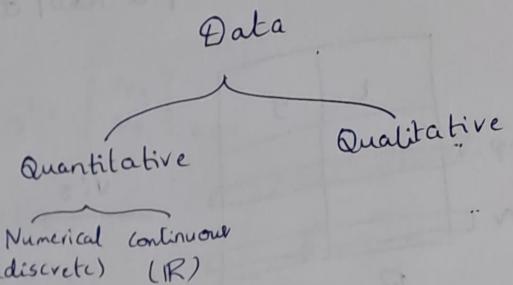


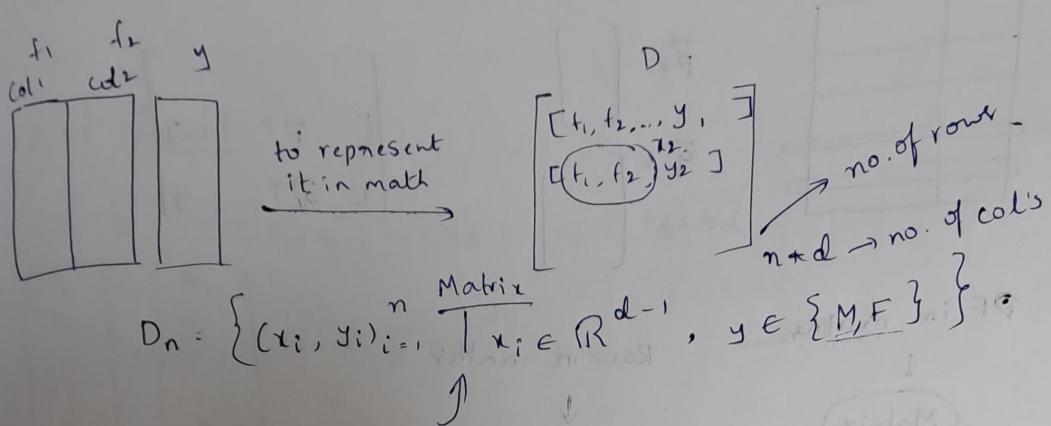
3/12/2022

# Maths

- dimensions / features / Column
- datapoint / query point / Row



It is impossible to visualize more dimensions. This is the reason we go with Mathematical functions



$D_n = \left\{ (x_i, y_i)_{i=1}^n \mid x_i \in \mathbb{R}^{d-1}, y \in \mathbb{R} \right\} \rightarrow \text{Regression.}$

In maths,  $y = mx + c$  → line Equation

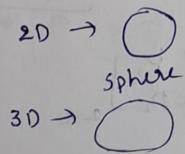
$$y = mx + c \rightarrow \text{line Equation}$$

3D → plane

4D

⋮ ⌊ hyperplane ⌋

D

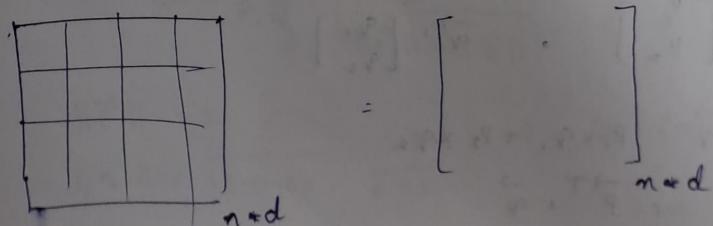


≥ 4D → Hyper sphere

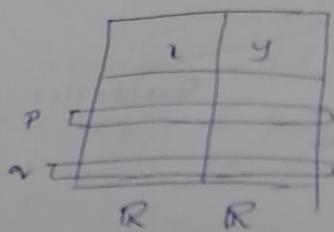
Hypercube

3/12/22

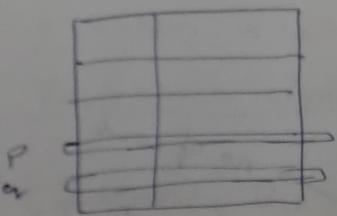
Representation of a Dataframe (Using matrix)



• Representation of a Row/ data point (Using column vector)



$\rightarrow$  Vis



$$\vec{P} = \begin{bmatrix} \quad \\ \quad \end{bmatrix}$$

$$\vec{q} = \begin{bmatrix} \quad \\ \quad \end{bmatrix}$$

$$\begin{bmatrix} \quad & \vec{P}^T \\ \quad & \vec{q}^T \end{bmatrix}$$

$$\begin{bmatrix} \quad \\ \quad \end{bmatrix}$$

$$\begin{bmatrix} \quad \\ \quad \end{bmatrix}$$

→ DF in mathe

↓  
Matrix

Row in mathe

↓

column Matrix + (Vectors)

↓

Then converted into row vectors in  
matrix.

## Mathematical Operations on Vectors (i.e., Column Vectors)

$$\vec{P} = \begin{bmatrix} P_1 \\ P_2 \end{bmatrix}_{2 \times 1}$$

$$\vec{q} = \begin{bmatrix} q_1 \\ q_2 \end{bmatrix}_{2 \times 1}$$

⇒ Addition

$$\vec{P} + \vec{q} = \begin{bmatrix} P_1 + q_1 \\ P_2 + q_2 \end{bmatrix}$$

$$\vec{P} = \begin{bmatrix} 150 \\ 50 \end{bmatrix}$$

$$\vec{q} = \begin{bmatrix} 180 \\ 60 \end{bmatrix}$$

$$\vec{P} + \vec{q} = \begin{bmatrix} 330 \\ 130 \end{bmatrix}_{2 \times 1}$$

⇒ Multiplication (dot product only)

$$\vec{P} \cdot \begin{bmatrix} P_1 \\ P_2 \end{bmatrix} . \vec{q} \cdot \begin{bmatrix} q_1 \\ q_2 \end{bmatrix}$$

$$\vec{P} \cdot \vec{q} = P_1 \times q_1 + P_2 \times q_2$$

# Vector Algebra

$$\vec{P} = \begin{bmatrix} P_1 \\ Q_1 \end{bmatrix}_{2 \times 1}, \quad \vec{Q} = \begin{bmatrix} Q_1 \\ Q_2 \end{bmatrix}_{2 \times 1}$$

Addition

$$\vec{P} + \vec{Q}$$

$$= \begin{bmatrix} P_1 + Q_1 \\ P_2 + Q_2 \end{bmatrix}$$

Vector

Subtraction

$$\vec{P} - \vec{Q}$$

$$= \begin{bmatrix} P_1 - Q_1 \\ P_2 - Q_2 \end{bmatrix}$$

Vector

Multiplication

$$\vec{P} \cdot \vec{Q}$$

$$= P_1 \cdot Q_1 + P_2 \cdot Q_2$$

Scalar.

## Dot Product

$$\vec{P} \cdot \vec{Q}$$

$$\begin{bmatrix} P_1 \\ P_2 \\ \vdots \\ P_d \end{bmatrix} \cdot \begin{bmatrix} Q_1 \\ Q_2 \\ \vdots \\ Q_d \end{bmatrix}$$

$$P_1 \times Q_1 + P_2 \times Q_2 + \dots + P_d \times Q_d$$

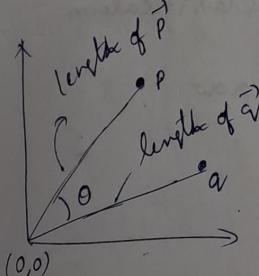
$$\vec{P} \cdot \vec{Q} = \sum_{i=1}^d P_i \times Q_i$$

$$\rightarrow 0$$

$$[P_1, P_2, \dots, P_d] \cdot \begin{bmatrix} Q_1 \\ Q_2 \\ \vdots \\ Q_d \end{bmatrix}$$

$$\vec{P} \cdot \vec{Q} = \vec{P}^T \times \vec{Q}$$

$$\rightarrow ②$$



$$\vec{P} \cdot \vec{Q} = \| \vec{P} \| \times \| \vec{Q} \| \times \cos \theta$$

length of vector  
i.e., magnitude of  
vector/norm of  
vector

Angle b/w  $\vec{P}$  &  $\vec{Q}$

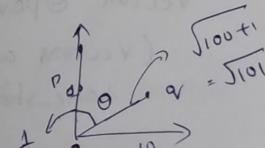
Question :

$$\vec{P} \cdot \vec{Q} = 2$$

$$\vec{P} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \vec{Q} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\therefore \vec{P} \cdot \vec{Q} = \| \vec{P} \| \times \| \vec{Q} \| \times \cos \theta$$

$$2 = \| \vec{P} \| \times \| \vec{Q} \| \times \cos \theta$$



$$\cos \theta = \frac{1}{\| \vec{P} \| \times \| \vec{Q} \|}$$

$$\| \vec{P} \| = 1$$

$$\| \vec{Q} \| = \sqrt{1+100} = \sqrt{101}$$

$$\Rightarrow \theta = \cos^{-1} \left[ \frac{1}{\sqrt{101}} \right]$$

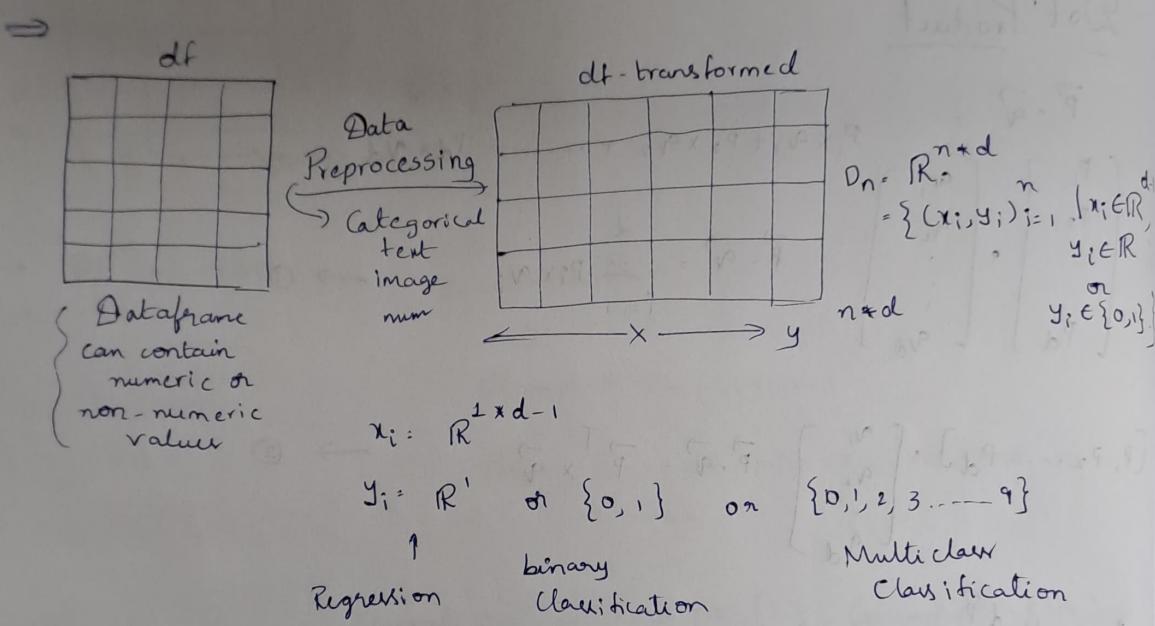
$$= \cos^{-1}(0.099) = 1.471^\circ$$

$$\|\vec{P}\| = \sqrt{\sum_{i=1}^d p_i^2}$$

$$\|\vec{q}\| = \sqrt{\sum_{i=1}^d q_i^2}$$

To get inverse of cos from numpy it is  
`np.arccos()`

$$\Theta_{\vec{P}, \vec{q}} = \cos^{-1} \left( \frac{\sum_{i=1}^d p_i q_i}{\sqrt{\sum p_i^2} \cdot \sqrt{\sum q_i^2}} \right)$$



Column Vector Representation of a Datapoint or row

$$x_i = \begin{bmatrix} ] \\ ] \\ ] \end{bmatrix}$$

Column Vector (By Default)

Vector Operation  
 (Vectors are of same shape)

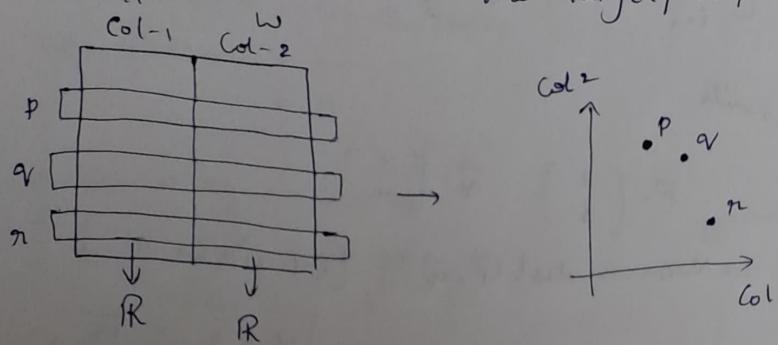
Addition  $\rightarrow$  component wise operation  $\rightarrow$  Vector

Subtraction  $\rightarrow$  "  $\rightarrow$  Vector

Dot product  $\rightarrow$  "  $\rightarrow$  Scalar

To get dot product from numpy  
`np.dot(a, b)`

Understand the similarities b/w datapoints are important to determine the target/output variable.

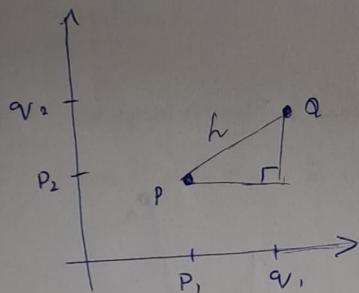


$$\text{Sim}(\vec{P}, \vec{Q}) > \text{Sim}(\vec{P}, \vec{R})$$

Similarity b/w two points/vectors can be computed by looking into the distance between the two points

$$\uparrow \downarrow \text{Similarity} \propto \frac{1}{\text{Distance}} \downarrow \uparrow$$

Distance between the two vectors:



$$h = \sqrt{(side-1)^2 + (side-2)^2}$$

$$\text{dist}(\vec{P}, \vec{Q}) = \sqrt{(v_1 - p_1)^2 + (v_2 - p_2)^2}$$

$$\vec{P} = \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_d \end{bmatrix} \quad \vec{Q} = \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_d \end{bmatrix}$$

$$\text{dist}(\vec{P}, \vec{Q}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2}$$

$$= \sqrt{\sum_{i=1}^d (q_i - p_i)^2}$$

$$\boxed{\text{dis}(\vec{P}, \vec{Q}) = \left[ \sum_{i=1}^d (q_i - p_i)^2 \right]^{1/2}}$$

↑ Euclidean Distance (i.e., Shortest distance b/w two vectors)

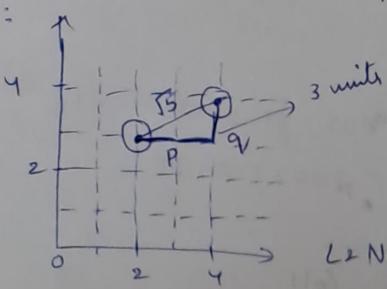
or

L2 Norm.

$$\text{dist}(\vec{p}, \vec{q}) = \left[ \sum_{i=1}^d \text{abs}(q_i - p_i)^2 \right]^{1/2}$$

absolute ||

Eg:



$$\vec{p} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad \vec{q} = \begin{bmatrix} 4 \\ 3 \end{bmatrix}$$

$$L_2 \text{ Norm} \rightarrow \text{dist}(\vec{p}, \vec{q}) = \left[ (4-2)^2 + (4-3)^2 \right]^{1/2} = \sqrt{4+1} = \sqrt{5}$$

$$L_1 \text{ Norm} \rightarrow \text{dist}(\vec{p}, \vec{q}) = \left[ \text{abs}(4-2) + \text{abs}(4-3) \right]^{1/2}$$

Manhattan distance

$L_p$  Norm:

If  $p=1$ ,  $L_1$  norm / Manhattan distance

If  $p=2$ ,  $L_2$  norm / Euclidean distance

$$(L_p \text{ Norm}) = \left( \sum_{i=1}^d \text{abs}(q_i - p_i)^p \right)^{1/p}$$

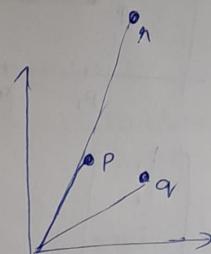
Minkowski  
distance

Angular distance:

Distance  $\uparrow \downarrow \Rightarrow$  Sim  $\downarrow \uparrow$

$$\text{similarity}(\vec{p}, \vec{q}) < \text{sim}(\vec{p}, \vec{r})$$

$$\text{Angular distance}(\vec{p}, \vec{q}) > \text{Angular distance}(\vec{p}, \vec{r})$$



Cosine Similarity:

$$\vec{p} \cdot \vec{q} = \|\vec{p}\| \|\vec{q}\| \cdot \cos \theta_{\vec{p}, \vec{q}}$$

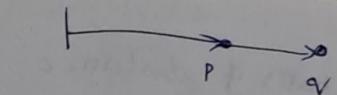
$$\boxed{\cos \theta_{\vec{p}, \vec{q}} = \frac{\vec{p} \cdot \vec{q}}{\|\vec{p}\| \|\vec{q}\|}} \quad \leftarrow \text{Cosine Similarity}$$

$-1 \leq \cos \theta \leq +1$

$$\text{Cosine distance} = 1 - \text{Cosine Similarity}$$

$$\cos \theta(\vec{p}, \vec{q}) = 1$$

$$\cos \theta \left\{ \begin{array}{l} \max = 1 \\ \min = -1 \\ 0 \end{array} \right.$$



angular-distance ( $\vec{p}, \vec{q}$ )  $\downarrow \Rightarrow \text{sim}(\vec{p}, \vec{q}) \uparrow$

Cos distance = 1 - Cos Similarity

$$\begin{aligned}\text{Cos-distance} &= 1 - 1 \\ &= 0\end{aligned}$$

Cos-Similarity = Cos  $\theta$

$$\begin{aligned}&= \cos 180^\circ \\ &= -1\end{aligned}$$



$$\text{Cos distance} = 1 - (-1) = 2$$

Cos-dist

0

2

1

Cos-sim

1

-1

0

$\rightarrow 0^\circ$  angle or  $360^\circ$  angle

$\rightarrow 180^\circ$

$\rightarrow 90^\circ$  &  $270^\circ$  angle

No. of nearest neighbours

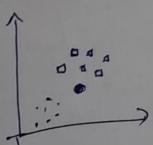
## RNN Algorithm:

Type - Supervised Learning

Task - Classification & Regression.

① Given  $\rightarrow D_n = \{(x_i, y_i); i=1, \dots, n | x_i \in \mathbb{R}^{d-1}, y_i \in \{0, 1\}\}$

$x_i \in \mathbb{R}^{d-1}$  - input vector



② Decide - Value of  $k$

- Distance Metrics

Euclidean

Manhattan

Minkowski

Angular

distances = [ ]

③ for datapoint in  $D_n$ :

- calculate the distance b/w  $x_q$  and data point
- distances.append ((datapoint, distance))

distances = [(d<sub>p1</sub>, d<sub>1</sub>), (d<sub>p2</sub>, d<sub>2</sub>), ..., (d<sub>pn</sub>, d<sub>n</sub>)]

- ④ Sort the distances list on the basis of distance  $d_i$  i.e.,  $i$  from 1 to  $n$   
 distances =  $[(d_{P1}, d_1), (d_{P2}, d_2), \dots]$

- ⑤ KNNpts = distances [:k]

$$\text{count\_pos} = 0$$

$$\text{count\_neg} = 0$$

for pt in KNNpts:

$$\text{if pt}['y'] == 1:$$

$$\text{count\_pos} += 1$$

else:

$$\text{count\_neg} += 1$$

Issue with KNN Algorithm:

Time &

Space complexity for KNN  
is very high.

$\Rightarrow$  Mode.

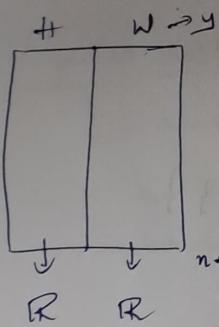
- ⑥ if count\_pos > count\_neg:

$$\text{print } ['y_q = 1']$$

else:

$$\text{print } ['y_q = 0']$$

### KNN for Regression:



① Given

$$D_n = \mathbb{R}^{n+d}$$

$$D_n = \{(x_i, y_i)\}_{i=1}^n \mid x_i \in \mathbb{R}^{d-1}, y_i \in \mathbb{R}\}$$

Know the query point  $x_q \in \mathbb{R}^{d-1}$

- ② Decide  $\rightarrow k$  (by default - 5)

Distance (by default = Euclidean)

- ③ distances = []

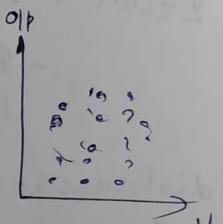
- ④ for  $x_i$  in  $D_n$ :

- calculate the distance b/w  $x_i$  &  $x_q$

- distances.append(( $d(x_i, x_q)$ ,  $y_i$ ))

- ⑤ Sort the distances list based on  $d(x_i, x_q)$

- ⑥ KNNpts = distances [:k]



② Mean/Median

$$\text{sum} = 0$$

for pt in KNNpts :

$$\text{sum} = \text{sum} + \text{pt}[y_i]$$

$$\text{avg} = \frac{\text{sum}}{K}$$

return avg  $\leftarrow y_v$

⇒ Hyperparameter in KNN?

K & distance

↑

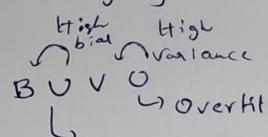
# of nearest neighbours

⇒ Explain Overfitting & Underfitting in KNN?

$K \downarrow \rightarrow$  Overfitting / high variance / low bias

$K \uparrow \rightarrow$  Underfitting / high bias / low variance

In KNN underfit & Overfit can be controlled by changing the value of K.



⇒ How do you detect Overfitting / Underfitting?

⇒ Can you explain the impact of outliers on KNN models?

⇒ Can you explain the impact of K on KNN models?

Same solution i.e., Overfit & Underfit

$K \downarrow$   $K \uparrow$

Min value of K = 1

Max value of K = n / total no. of datapoints in given dataset.

⇒ What is the space & complexity of KNN?

⇒ Why do we call KNN as lazy learner?

Since there is no learning in KNN algorithm. It simply remembers the data.

At prediction time it performs the distance calculation, computes the nearest neighbours & do the prediction for a query point.

Data → Algo → Model  
 $x_q \rightarrow$  Model →  $\hat{y}_q$

KNN

Data → Algorithm → Nothing but remembers the data  
 $x_q \rightarrow$  Algo →  $\hat{y}_q$   
↳ distance Calculation  
↳ Sorting

⇒ Is K-NN a parametric Algo or a Non-parametric Algorithm

Non-Parametric

⇒ What is the difference b/w parameters & hyperparameters?

DT / ID3 Algorithm 10%  
Algo → C4.5 Algo

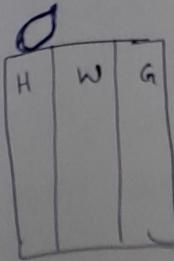
threshold - Acceptable  
↓ value  
(10-20% larger than  
max value)

Naive Bayes' Algo → MAP (Max A Posteriori)

KNN Algo → Lazy learner.

## Naive Bayes Algorithm: (Maximum A Posteriori)

(Mainly for Classification task)



Supervised learning  
Target Variable  
Classification

- Different approaches
  - Distance Based Approach
  - Boundary Based Approach
  - Probabilistic " "
  - Rule Based Approach
  - Ensemble " "
  - Deep learning " "

Task: Given  $x_q$  find  $y_q$

$$P(y_q = M | x_q(h, w)), P(y_q = F | x_q(h, w))$$

① Identify target variable

② Task

③ Evaluation Matrix

④ Algorithms

\* Given  $D_n = \{(x_i, y_i)\}_{i=1}^n \mid x_i \in \mathbb{R}^{d-1}, y_i \in \{c_1, c_2, \dots, c_k\}\}$

$x_q =$

Task → find  $y_q$

How? → Given  $x_q$ , compute the prob of  $x_q$  belonging to each of the classes whichever is maximum, assign that class to  $x_q$   
i.e.,  $P(c_1|x_q), P(c_2|x_q) \dots P(c_k|x_q)$

$$y_q = \arg \max \{ P(c_i|x_q) \}$$

$\hookrightarrow$  Posterior probability  $= \frac{P(c_i)}{P(x_q)}$

$c_1, c_2, \dots, c_k$  mark

From Bayes Theorem  $\frac{\text{Likelihood}}{\text{Prior}} = \frac{P(x_q|c_i) * P(c_i)}{P(x_q)}$

$$P(c_i|x_q) = \frac{P(x_q|c_i) * P(c_i)}{P(x_q)}$$

$\hookrightarrow$  Marginal

Probability →  
Study of  
Uncertainty

Probability:

① Random Experiment →

an experiment for which exact outcome is not known

for Eg: Tossing a coin, Rolling a dice

$$\{H, T\}$$

$$\{1, 2, 3, 4, 5, 6\}$$

- $\Rightarrow$  ⑥ {Sample Space}  $\rightarrow$  It is a set of all the possible outcomes for a random experiment.
- $\Rightarrow$  ⑦ {Event}  $\rightarrow$  Subset of a sample space

For Eg:

R.E  $\rightarrow$  Uncertain experiment

\* Rolling a dice

$$S.S \rightarrow \{1, 2, 3, 4, 5, 6\}$$

\* set of all the possible outcome.

Event  $\rightarrow$  Subset of a Sample Space

$$E_1 - \text{a prime no.} \in \{2, 3, 5\}$$

$$P(E) = \frac{3}{6} = 50\% \text{ or } 0.5$$

Conditional Probability:

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

Compulsory Event is the subset of Sample Space that why probability is  $0 \leq P(E) \leq 1$

Axioms of probability

Frequency table:

$$P(SS) = 1$$

Temperature: Yes No

Marginal Probability:

	Hot	Mild	Cool	
P(A)	2	4	3	9
P(B)	2	2	1	5
	4	6	4	

Joint Probability:

$$P(A \cap B)$$

Wind:

Yes No

outlook

yer No

Weak 2 6 8

Strong 3 3 6

Sums

(3) (2) 5

5 9

Rain

$P(\text{Playing Tennis if wind strong})$

$$\frac{3}{14} = \frac{3}{14}$$

$P(\text{Play tennis} = \text{'No'} \mid \text{outlook} = \text{Rainy})$

$$= \frac{P(\text{'No'} \wedge \text{Rainy})}{P(\text{Rainy})} = \frac{\frac{2}{14}}{\frac{5}{14}} = \frac{2}{5}$$

R.E: tossing a coin & Rolling a dice

S.S:  $\{(H, 1), (H, 2), (H, 3), (H, 4), (H, 5), (H, 6)$

$(T, 1), (T, 2), (T, 3), (T, 4), (T, 5), (T, 6)\}$

Event: What is probability of getting an even number on the dice.

$$= \{(H, 2), (H, 4), (H, 6), (T, 2), (T, 4), (T, 6)\}$$

Probability:  $\frac{6}{12} = \frac{1}{2} = 0.5 \text{ or } 50\% \text{ (chance)}$

$$P(A|B) = \frac{P(A \wedge B)}{P(B)} = \frac{\frac{3}{12}}{\frac{6}{12}} = \frac{3}{6} = 50\% \text{ chance (or) } 0.5$$

Given that B has already happen

Baye's theorem:

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

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

$$P(A \wedge B) = P(B|A) P(A)$$

Independent Events:

$$P(A|B) = P(A)$$

$$P(B|A) = P(B)$$

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

$$P(A \wedge B) = P(B|A) \cdot P(A)$$

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

Similarly,

$$P(A \wedge B) = P(A|B) * P(B)$$

$$= P(A) * P(B)$$

$P(A \wedge B)$   
 $= P(A|B) * P(B)$   
 $= P(A) * P(B)$

More than event (A & B)

Marginal

$$\begin{matrix} P(A) \\ P(B) \end{matrix}$$

$$\left. \begin{matrix} P(A \wedge B) \\ P(B|A) \end{matrix} \right\} \text{Both are same}$$

$$\left. \begin{matrix} P(B \wedge A) \\ P(A|B) \end{matrix} \right\} \text{Both are same}$$

Conditional

$$P(A|B)$$

$$P(B|A)$$

$$\Rightarrow \rightarrow P(C_1 = 0 | x_q) \quad P(C_2 = 1 | x_q)$$

man

$$\Rightarrow \max \{ P(C_1 | x_q), P(C_2 | x_q) \}$$

$$\max \{ P(x_q | C_1) * P(C_1), P(x_q | C_2) * P(C_2) \}$$

$$\max \{ P(f_1, \dots, f_{d-1} \cap C_1), P(f_1 \cap f_2, \dots, f_{d-1} \cap C_2) \}$$

$y_q = a$

For Tennis dataset:

$$P(C_1 = \text{Yes} | x_q) = \frac{s/c}{s/a} = \frac{(8/14)}{(8/12)} = (8/14)$$

$\Rightarrow$

$$(8/14) / (6/14) = (4/3)$$

$\Rightarrow$

Color = White

Texture = Smooth

Sound = Woof

$$P(C_1 = \text{Dog} | x_q)$$

$\downarrow$

$$P(x_q | \text{Dog}) + P(\text{Dog})$$

$\downarrow$

$$P(x_q \cap \text{Dog})$$

$\downarrow$

$$P(\text{White} \cap \text{Smooth} \cap \text{Woof} \cap \text{Dog})$$

$$P(\text{White} | \text{Smooth} \cap \text{Woof} \cap \text{Dog}) *$$

$$P(\text{Smooth} | \text{Woof} \cap \text{Dog}) +$$

$$P(\text{Woof} | \text{Dog}) + P(\text{Dog})$$

$$P(C_2 = \text{Cat} | x_q)$$

$$(8/14) \downarrow$$

$$P(x_q | \text{Cat}) + P(\text{Cat})$$

$\downarrow$

$$P(x_q \cap \text{Cat})$$

$\downarrow$

$\Sigma_f$  contains Independent features

$$P(\text{white}|\text{Dog}) * P(\text{smooth}|\text{Dog}) * P(\text{woof}|\text{Dog}) * P(\text{Dog})$$

$$P(\text{white}|\text{cat}) * P(\text{smooth}|\text{cat}) * P(\text{woof}|\text{cat}) * P(\text{cat})$$

Class

Sound Dog Cat

	Dog	Cat		Texture	Dog	Cat	
Woof	3	1	4	Furry	2	3	5
Purr	1	3	4	Smooth	2	1	3

4 4

Colour Dog Cat

	Dog	Cat
White	2	1
Black	2	3
	4	4

4 4

$$\Rightarrow \frac{P(\text{white} \cap \text{Dog})}{P(\text{Dog})} * \frac{P(\text{white} \cap \text{Cat})}{P(\text{Cat})}$$

$$\frac{P(\text{white} \cap \text{Dog})}{P(\text{Dog})} * P(\text{Dog})$$

$$\Rightarrow \frac{218}{418} * \frac{218}{418} + \frac{318}{418} + \frac{4}{8}$$

$$\Rightarrow \frac{1}{2} + \frac{1}{2} + \frac{3}{4} + \frac{1}{2}$$

$$\Rightarrow \frac{3}{32} = 0.0937$$

$$\Rightarrow \frac{118}{418} + \frac{118}{418} + \frac{118}{418} + \frac{4}{8}$$

$$\Rightarrow \frac{1}{4} + \frac{1}{4} + \frac{1}{4} + \frac{1}{2}$$

$$\Rightarrow \frac{1}{128} = 0.0078$$

The predicted class is Dog

Issue  $\rightarrow$  Zero probability Issue

If a feature in query point contain an unseen value {in historical data}

In this scenario NB fails provide the output

\* float point precision  $\rightarrow$  Numerical Instability

To bring out from Numerical Instability,

log transformation

$$\log(a+b) = \log a + \log b$$

$$\hat{y}_q = \underset{i=1, \dots, k}{\operatorname{argmax}} \left\{ \log \left( P(c_i) + \sum_{j=1}^d P(f_j | c_i) \right) \right\}$$

$$\hat{y}_{av} = \underset{i \in \{1, \dots, k\}}{\text{arg max}} \left\{ \log(p(c_i)) + \log \left[ \prod_{j=1}^d p(f_j | c_i) \right] \right\}$$

$$\log \sum_{j=1}^d b_j = \log(b_1 + b_2 + b_3 + \dots + b_d)$$

$x, y \rightarrow$  [NB Algo]  $\rightarrow$  Model  
 freq table of  $f_i$  vs  $y$   
 & if  $f_i$  is categorical

$$\begin{cases} f_i \xrightarrow{My} y \xrightarrow{\text{if } f_i \text{ is num}} \text{from Normal dist} \end{cases}$$

$$x \rightarrow y_{av} = \text{argmax}\{\cdot\} \rightarrow \hat{y}_{av}$$

$\rightarrow$  Zero Probability

Solution: Apply Regularization  $\rightarrow$  Laplace Smoothing

for Eg:  $P(\text{Brown} | \text{Dog}) = 0 + \text{some very small value}$

$P(\text{Brown} | \text{Cat}) = 0 + \text{"}$

$$\Rightarrow P(B|D) = \frac{p(f_j | c_i)}{n + K + d} \quad n = P(\text{Dog}) \quad K = \# \text{ of unique values in } A \text{ (now)} \\ \text{i.e., } K = 2$$

$$\Rightarrow P(B|C) = \frac{0 + \alpha}{m + K + d} \quad m = P(\text{Cat}) \quad \alpha = \text{hyperparameter} \quad K = 2$$

$\alpha \rightarrow$  hyperparameter (we have to decide what to take)

Questions:

1) Is there any hyperparameter in Vanilla N.B?

Ans: No. It will be having zero probability value

2) How do you resolve zero prob issue?

Ans: By applying Regularization

3) Can you explain Overfitting & Underfitting in N.B?

$\alpha \rightarrow \uparrow \rightarrow$  Underfitting

$\alpha \rightarrow \downarrow \rightarrow$  Overfit

4) Is there any impact of outliers on N.B?

Not affected by outliers that much

No impact

5) Why we called N.B as MAP (Maximum Posteriori)

Because  $x_N$

We try to find

$$P(c_i|x_N) = P(x_N|c_i) * P(c_i) / P(x_N)$$

$$P(c_2|x_N)$$

$$P(c_3|x_N)$$

Max of these probability is taken as Class

From Sklearn.

Naive Bayes

Gaussian NB

Data → Numerical  
features ↑  
few categorical  
features

Multinomial NB

Data → Text data

Bernoulli NB

lot of Categorical  
features.

6) Assumption in NB Algorithm?

+ features should be conditionally independent

+ Another assumption for Gaussian NB

- All num features should be gaussian distributed

### Decision Tree Algorithm

Highly used Algorithms. → commercialized (not available freely)

→ Decision Tree Model (ID3, C4.5, C5.0, SPRINT) → commercialized (not available freely) where data is distributed on different machines

Algorithms to create DT Models.

\* ID3 → Iterative Dichotomiser 3 (Old Algorithm only works for the data which contains categorical features → works only for classification)

→ Rule Based Algorithm

\* C4.5 → Solves both Classification as well as Regression

Data → Categorical / Numerical

Rules: for example Iris data

1. for Setosa

- PL < 2.5

- PW < 0.8

For Versicolor

- PL < 5.0

- PW < 1.6

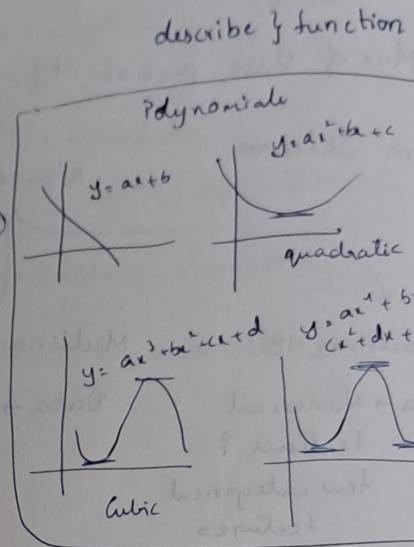
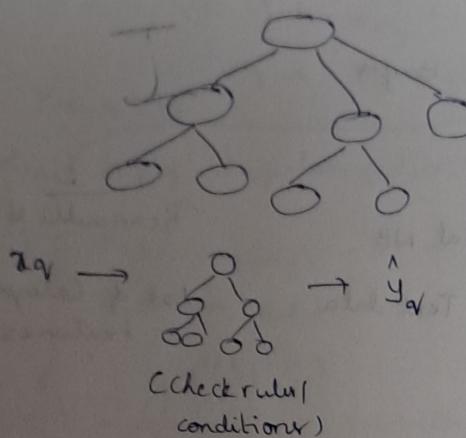
For Virginica

- PL ≥ 5.0

- PW ≥ 1.6

shape  
column  
index } properties

Tree Structure :



\* For DT Algorithm we don't need to apply Standardization / Normalization technique on categorical column only apply label encoder.

→ Max depth refer to length of longest path from root to a leaf. I limit to stop further splitting of nodes.

To picturize DT model

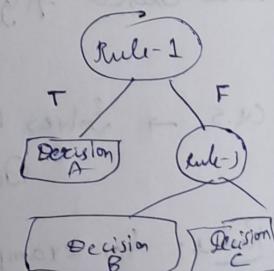
- » from sklearn.tree import plot\_tree
- » filled = True (Colour filling)
- » rounded = True (for edge roundness)

Learnings :

x-train, y-train →

ID3/  
c4.5

→ Model  
List of rules  
arranged in  
tree structure

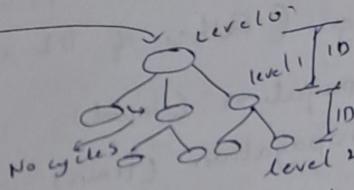
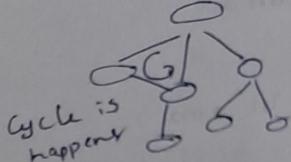


Tree Structure norms:

\* Nodes (Root, Intermediate, leaf nodes)

\* Branches (No cycles)

difference b/w graph & trees:



SQL → indexing for easy & quick retrieval of data  
↳ indexing is stored using B trees & B+ trees former.

ID3 Algorithm : (Entropy, Information Gain, Gini Impurity)

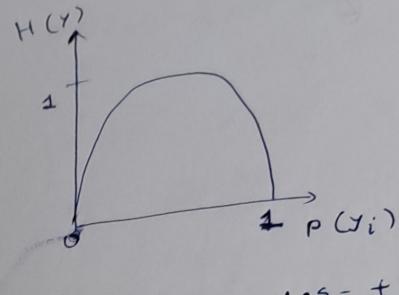
1. Entropy  $\rightarrow$  Information theory  
 $\rightarrow$  Bivariate

$\rightarrow H(\text{Target}) \rightarrow$  no. of classes.

$$\rightarrow H(Y) = - \sum_{i=1}^k p(y_i) \cdot \log_2 p(y_i)$$

$$0 \leq H(Y) \leq 1 \leftarrow \text{Max Entropy}$$

No imp.

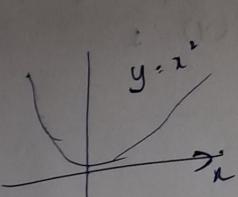
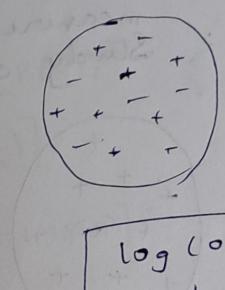


$$\Rightarrow H(Y) = - \{ p(\text{Yes}) \cdot \log_2(p(\text{Yes})) + p(\text{No}) \cdot \log_2(p(\text{No})) \}$$

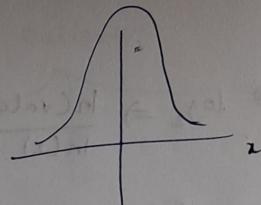
$$= - \left\{ \frac{9}{14} \cdot \log_2 \left( \frac{9}{14} \right) + \frac{5}{14} \cdot \log_2 \left( \frac{5}{14} \right) \right\}$$

$$= 0.606.$$

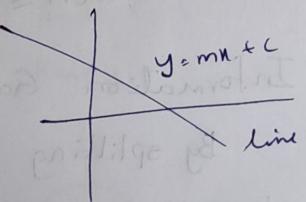
$\log(0 \text{ to } 1)$   
↓  
-ve value



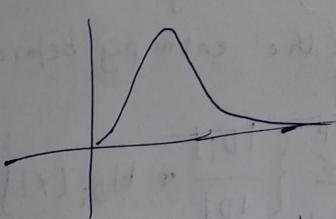
Parabola  
function



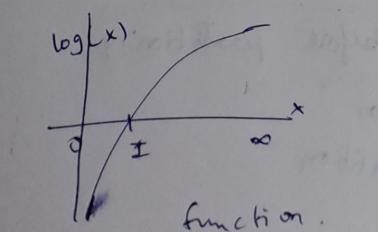
Normal distri



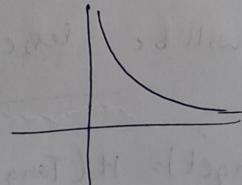
function



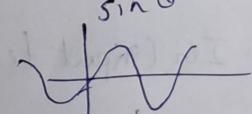
Right skewed dist



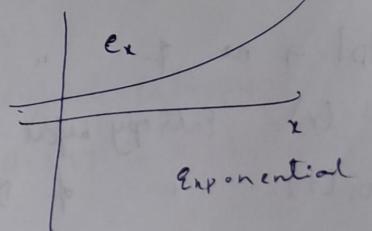
function.



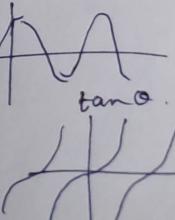
Power law distribution



sin θ



Exponential graph.



Geoffrey

In deep learning (Father of DL - Hinton)

- forward propagation (input)
- backward propagation (train) } used in 2012

- \* Make a split of impure data such that the new partitions should have lesser impurity  $\rightarrow$  Entropy ↓

Ways to create Decision Trees?

- ① ID3 (Iterative Dichotomiser 3) - Only classification task for categorical variables

Entropy: measure

Steady of randomness

or impurity



$$\text{Entropy} = 0$$

∴ No impurity value



$$H(Y) = \text{Entropy} = 0$$

$$\# \text{ of -ve} = \# \text{ of +ve}$$

$$\rightarrow H(Y) = 1$$

$$\rightarrow H(Y) = 1$$

$$\Rightarrow 0 \leq H(Y) \leq 1$$

$$\Rightarrow \log_2 \Rightarrow \frac{\ln(\text{value})}{\ln(2)}$$

2. Information Gain (helps in determining the Best split)

By splitting into multiple partitions the weighted entropy of partitions will be lesser than the entropy before split.

$$IG_i(\text{input feature, target}) = H(\text{Target}) - \sum_{i=1}^p \left\{ \frac{|D_i|}{|D|} * H_{D_i}(Y) \right\}$$

here,  $p \rightarrow$  Total no of partitions

$|D_i| \rightarrow$  no. of datapoints in each partition

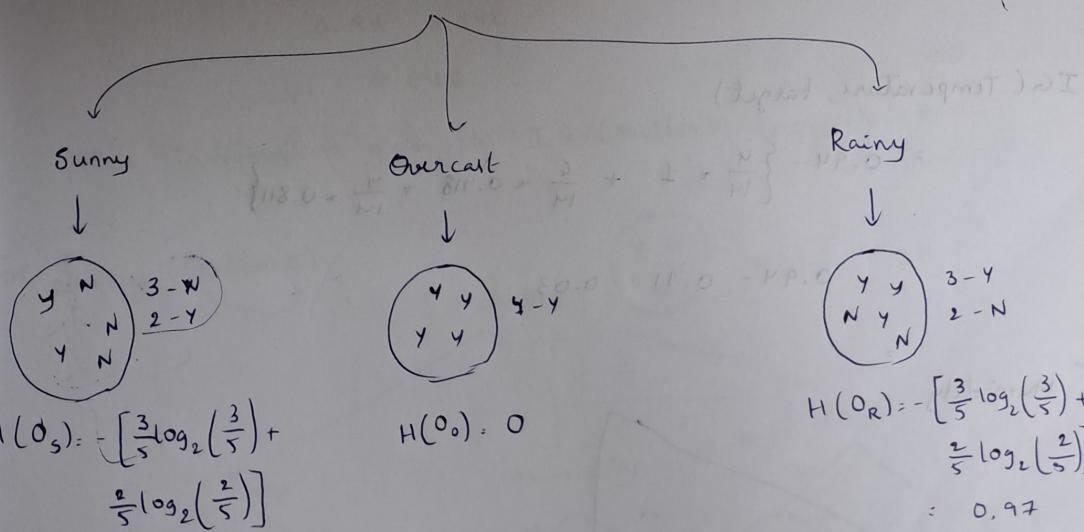
$|D| \rightarrow$  no. of datapoints before partitioning.

$H(Y) =$  Entropy before partition

Data:	Outlook	Temp	Humidity	Wind	PlayTennis
S	Hot	High	Weak	N	
S	H.	H	Strong	N	
O	H	H	W	Y	
R	Mild	H	W	Y	
R	Cool	Normal	W	Y	
R	C	N	S	N	
O	C	N	S	N	
S	M	H	W	Y	
S	C	H	W	N	
R	M	N	W	Y	
S	M	N	W	Y	
O	M	N	S	Y	
O	H	H	S	Y	
R	M	N	W	Y	
		H	S	N	

$$\begin{aligned}
 H(\text{Target}) &= - \sum_{i=1}^k P(x_i) \log_2 (P(x_i)) \\
 &= - [(P(Y) \cdot \log_2 P(Y)) + (P(N) \cdot \log_2 P(N))] \\
 &= - \left[ \frac{9}{14} \cdot \log_2 \left( \frac{9}{14} \right) + \frac{5}{14} \cdot \log_2 \left( \frac{5}{14} \right) \right] \\
 &= - [-0.6374 * (9/14) + (5/14) * (-1.4854)] \\
 &= -[-0.409 + (-0.5305)] \\
 &= -(-0.9394) = 0.94
 \end{aligned}$$

For Outlook feature



Avg = Weighed Entropy  
 $= 0.647 \approx 0.65$

IG (outlook, target)

$$= H(\text{target}) - \sum_{i=1}^P \left\{ \frac{|D_i|}{|D|} \cdot H(D_i | Y) \right\}$$

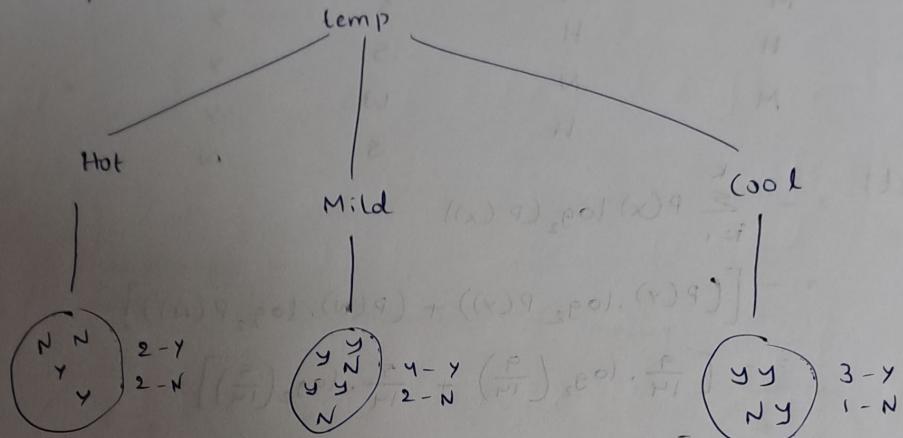
$$= 0.94 - \left\{ \frac{5}{14} \cdot (0.97) + \left( \frac{4}{14} \cdot 0 \right) + \left( \frac{5}{14} \cdot 0.97 \right) \right\}$$

$$= 0.94 - \left\{ 2 \left[ \frac{5}{14} \cdot 0.97 \right] \right\}$$

$$= 0.94 - 0.6714$$

$$= 0.2686$$

/ For temp feature



$$H(T_h) = - \left[ \frac{2}{4} \log_2 \left( \frac{2}{4} \right) + \dots \right] \quad H(T_m) = - \left[ \frac{4}{6} \log_2 \left( \frac{4}{6} \right) + \dots \right] \quad H(T_c) = - \left[ \frac{3}{4} \log_2 \left( \frac{3}{4} \right) + \dots \right]$$

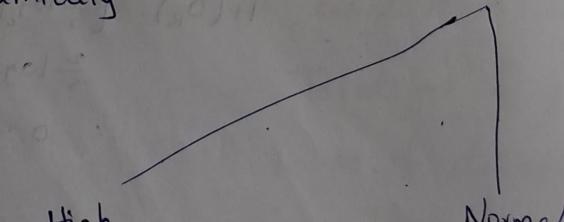
$$= 1 \quad = 0.918 \quad = 0.811$$

IG (Temperature, target)

$$= 0.94 - \left\{ \frac{4}{14} \cdot 1 + \frac{6}{14} \cdot 0.918 + \frac{4}{14} \cdot 0.811 \right\}$$

$$= 0.94 - 0.91 = 0.03$$

Humidity



$$\begin{array}{c} N \ N \ Y \\ Y \ N \ N \ Y \end{array} \quad Y-3 \quad N-4$$

$$\begin{array}{c} Y \ Y \\ N \ * Y \ Y \\ Y \end{array} \quad Y-6 \quad N-1$$

$$H(H_n) = -\left[\frac{4}{7} \log_2\left(\frac{4}{7}\right) + \frac{3}{7} \log_2\left(\frac{3}{7}\right)\right] \quad H(H_N) = -\left[\frac{6}{7} \log_2\left(\frac{6}{7}\right) + \frac{1}{7} \log_2\left(\frac{1}{7}\right)\right]$$

= 0.785

= 0.591

$IG(\text{Humidity}, \text{Target})$

$$= 0.94 - \left\{ \frac{7}{14} * 0.785 + \frac{7}{14} * 0.591 \right\}$$

$$= 0.94 - 0.788$$

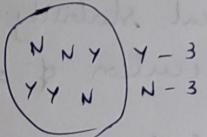
$$= 0.152$$

Wind

Wind

Weak

Strong



$$H(W_w) = -\left[\frac{6}{8} \log_2\left(\frac{6}{8}\right) + \frac{2}{8} \log_2\left(\frac{2}{8}\right)\right] \quad H(W_s) = -\left[\frac{3}{7} \log_2\left(\frac{3}{7}\right) + \frac{4}{7} \log_2\left(\frac{4}{7}\right)\right]$$

$$= 0.811$$

$$= 1$$

$IG(\text{Wind}, \text{Target})$

$$= 0.94 - \left\{ \frac{8}{14} * 0.811 + \frac{6}{14} * 1 \right\}$$

$$= 0.94 - 0.892$$

$$= 0.048$$

for Sunny (Outlook)

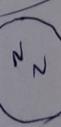
$$\text{H}(S) = 0.97$$

$$\frac{10}{14} \rightarrow \frac{10}{14} \times H(S) = 0.97$$

Temp

H

C



O



1

H

N



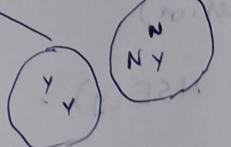
0

Wind

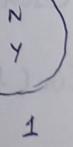
W

S

1



1



1

Take a decision on which feature to choose for splitting

	Classification	Regression
* KNN	✓ (model)	✓ (mean)
Supervised Learning	* Naive Bayes (Max a posterior)	✓ X
	* Decision Tree	brill

- \* KNN cannot be used if Latency is very important
- \* Numerical stability concerns how errors introduced during the execution of an algorithm affect the result.
- ⇒ CART (Classification and Regression Trees) is very similar to C4.5 and this algo is used by sklearn library
- ⇒ C4.5 (Entropy + Information Gain Ratio)
- ⇒ ID3 ( " + Information Gain)
- ⇒ CART (Gini Impurity / Gini Index + Weighed Sum of Entropy)

### Information Gain Ratio:

Information gain favours high branching features  
lets look at IG Ratio

$$\text{Gain Ratio} = \frac{\text{IG}(Y)}{\text{Intrinsic Information} / \text{Split Info}}$$
$$= H(Y) - \sum_i \left\{ \frac{|D_i|}{|D|} H_{D_i}(Y) \right\}$$
$$= \sum_i \left\{ \frac{|D_i|}{|D|} * \log_2 \frac{|D_i|}{|D|} \right\}$$

### IG<sub>r</sub> (for Regression)

$$IG_r(f_r, Y) = MSE(Y) - \sum_i \left\{ \frac{|D_i|}{|D|} \bar{MSE}_{D_i}(Y) \right\}$$
$$\frac{\sum}{n} (y - \frac{\sum y_i}{n})^2$$

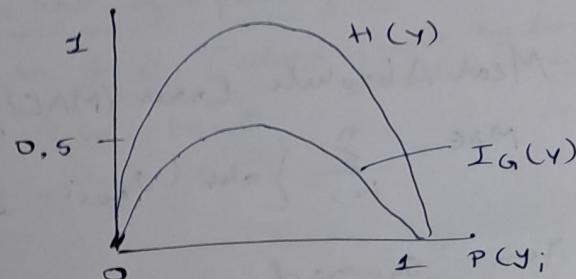
Gini Index (Another form of measure of impurity)

### Gini Impurity

$$\rightarrow I_G(\text{Target}) = 1 - \sum_{i=1}^c (p(y_i))^2$$

$$0 \leq I_G(y) \leq 1 - \frac{1}{c}$$

$$\rightarrow 1 - \left\{ p(y_{+ve})^2 + p(y_{-ve})^2 \right\}$$



### Measure of Randomness

$\rightarrow H(y)$

$\leftarrow I_G(y)$

Cuts with Gini Impurity instead of Entropy.

## Evaluation Matrix

### For Regression :-

① Mean Absolute Error (MAE)

$$MAE = \frac{1}{n} \sum_{i=1}^n \{ \text{abs}(y_{\text{acti}} - y_{\text{predi}}) \}$$

② Mean Squared Error (MSE)

Algorithm  $\rightarrow$  Decision Tree

How ? (CART)

- Compute impurity in your data using entropy / Gini impurity
- Choose one feature & split your data into multiple partitions.

Representation of datapoint or row using maths

↳ Vector  $\rightarrow$  By default representation  $\rightarrow$  Column

$x_1$ ,  $x_2$

## Equation of a line:

$$y = mx + c$$

$\hookrightarrow$  intercept  
 w.r.t y

Slope of  
 the line w.r.t x.

## Parameter of a plane / line:

$m \& c$  [where y cuts by line i.  
 at which point]

$$\hookrightarrow m = \tan \theta$$

2D dim  $\rightarrow a_1x_1 + b_1x_2 + c = 0$  (Not better approach) } Eq of line

$$w_1x_1 + w_2x_2 + \cancel{w_3x_3} + w_0 = 0$$

3D dim

$$w_1x_1 + w_2x_2 + w_3x_3 + w_0 = 0 \rightarrow \text{Eq of plane}$$

4D dim

$$w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + w_0 = 0 \rightarrow \text{Eq of hyperplane}$$

100 dim  $\rightarrow w_1x_1 + \dots + w_{100}x_{100} + w_0 = 0 \rightarrow$  It is also Eq of hyperplane.

For 2D  $\rightarrow w_1x_1 + w_2x_2 + w_0 = 0$ .

$$w_1x_1 = -w_2x_2 - w_0$$

$$x_1 = -\frac{w_2}{w_1}x_2 - \frac{w_0}{w_1}$$

Slope:  $-\frac{w_2}{w_1}$  w.r.t  $x_2$ , w.r.t  $x_1$  # of slopes = d-1

Intercept:  $-\frac{w_0}{w_1}$  w.r.t  $x_1$  # of intercepts = 1

$$w_1x_1 + w_2x_2 + \dots + w_{100}x_{100} + w_0 = 0$$

$$\vec{w} = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$\rightarrow \pi_{100} = \sum_{i=1}^{100} w_i x_i + w_0 = 0$$

$\vec{w} \in \mathbb{R}^2$   
 Real value with 2 components

Represent the Eq of Hyperplane as a dot Product

$$\pi_d: \vec{w} \cdot \vec{x} + w_0 = 0$$

$$\pi_d: \sum_{i=1}^d w_i x_i + w_0 = 0$$

$$\pi_d: \vec{w}^T \vec{x} + w_0 = 0$$

$$\therefore \vec{w}^T \vec{x} + w_0 = 0$$

$$\pi_d:$$

$$\|\vec{w}\| \|\vec{x}\| \cos \theta + w_0 = 0$$

20/12/22

→ Equation of Circle :

$$x^2 + y^2 = r^2$$

→ Equation of Sphere :

$$x^2 + y^2 + z^2 = r^2$$

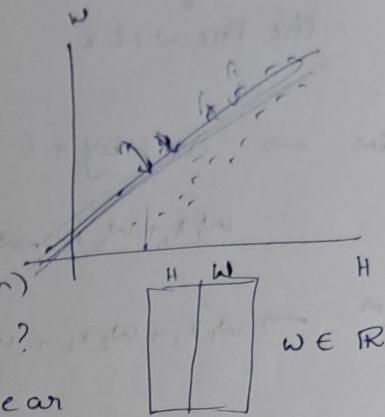
Regression task

Algorithm: Linear Regression

How? Find a line that best fits the data (Role of linear Regression)

Why it is named as Linear Regression?

Coefficient terms should be linear



Classification task

Algorithm: Logistic Regression

How? Find a line the best separation

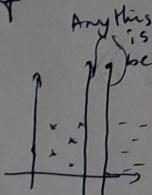
we from -ve.

Logistic Regression only find the linear pattern in case of non-linear distribution it simply fails

\* Gives only straight boundary.

Logistic (Vs) Support Vector

It takes any of the best fit separator for analyzing.



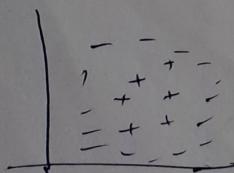
It only gives the best fit separator among all the best fit separators

\* It is only for classification tasks

Feature Engineering:

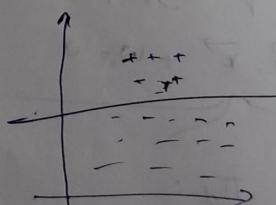
learns new features from the data.

By SVM

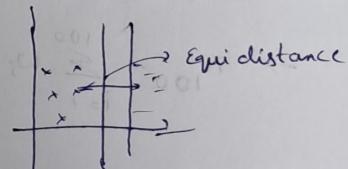


Automatic feature engg.

Low dimensionality



Higher dimensionality



This is only best separator

How? Try to find HYPERPLANE that BEST FITS the given data

### Type of Linear Regression:

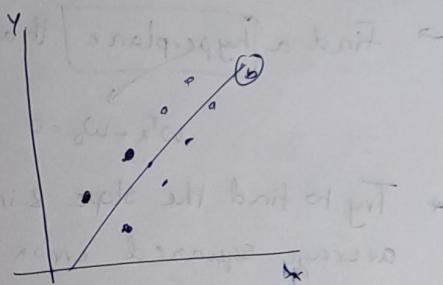
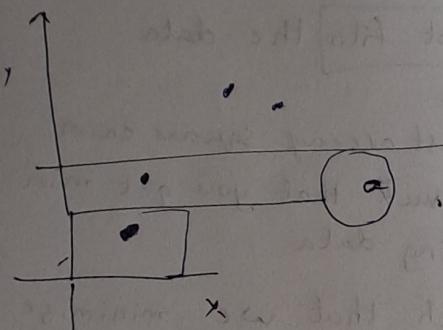
1. Simple Linear Regression (If your input contains one feature)  $y = mx + c$

2. Multiple linear Regression (If your input contains more than one feature)

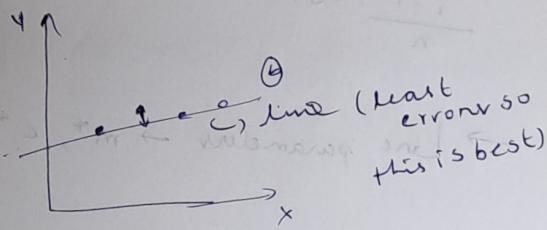
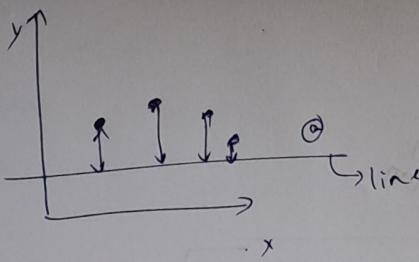
$$y = m_1x_1 + m_2x_2 + c$$

3. Linear Regression with features Engg / Polynomial Regression

What do I mean by Predict BEST FIT?



The line which fits the best is having the least error.



Compute the error done by line  $\hat{a}$  by line  $\hat{b}$  Compute

$$\text{Point } 1 \rightarrow \text{error}_1: y_1 - \hat{a}_{y_1} = y_1 - (mx_1 + c)$$

$$2 \rightarrow \text{error}_2: y_2 - \hat{a}_{y_2}$$

3

$$\vdots \rightarrow y_5 - \hat{a}_{y_5}$$

$$\frac{\sum \text{error}_i}{n} >$$

$$\text{Point } 1 \rightarrow \text{error}_1: y_1 - \hat{b}_{y_1}$$

Actual Predict

$$\frac{\sum \text{error}_i}{n}$$

(b) line best fits

## Issues with Linear Regression:

- 1. Mean is Corrupted
- 2. +ve & -ve values

↳ To solve this problem.

$$\text{for } Q \rightarrow \frac{\sum_{i=1}^n (y_i - (mx_i + c))^2}{n}$$

Actual error predicted

The best fit line is a line that has the least average squared error

→ find a [hyperplane] that has [best fits] the data

$$w^T x + w_0 = 0$$

least average square error.

\* Try to find the slope & intercept such that you get min average squared error on training data

\* Try to find the  $\vec{w}, w_0$  such that we minimize

$$\frac{\sum_{i=1}^n (y_i - (w^T x_i + w_0))^2}{n}$$

on the training data

→ Line parameters  $\rightarrow m^*, c^*$

x-train, y-train

$$\vec{w}^*, w_0^* = \arg \min_{\vec{w}, w_0} \left\{ \sum_{i=1}^n (y_i - \{w^T x_i + w_0\})^2 \right\}$$

minimizing the error

↳ loss function

Error function

Cost function

Optimization means minimizing a function

Ordinary Least Square Equation. ↴ i.e., minimizer loss function.

\* How do you solve the Optimization Eqn's?

By Using Optimizers

1. Gradient Descent

2. Stochastic Gradient

3. Momentum Based Gradient Descent

4. RMS Prop

\* The dot product b/w the vectors is zero whenever these vectors are perpendicular or orthogonal to each other.

$$\vec{w} \perp \vec{x}$$

$$\vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \\ w_0 \end{bmatrix}$$

$\Rightarrow \vec{w}$  is called a norm of a hyperplane.

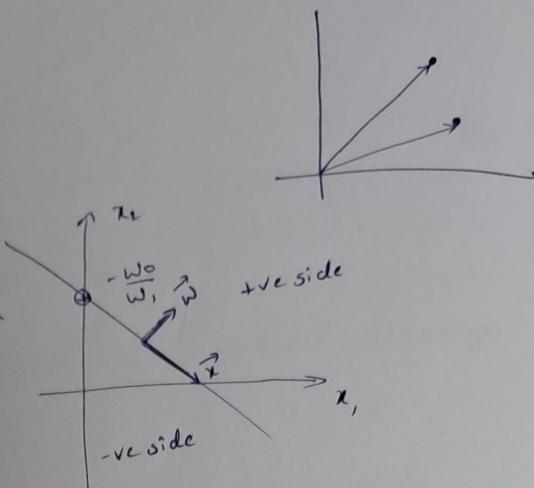
$\vec{w} \perp$  to the hyperplane

\* It helps us identify the direction in which the hyperplane is facing.

If  $\vec{w} \cdot \vec{p} > 0 \Rightarrow \vec{p}$  is on the +ve side

$\vec{w} \cdot \vec{p} < 0 \Rightarrow \vec{p}$  is on the -ve side

$\vec{w} \cdot \vec{p} = 0 \Rightarrow \vec{p}$  is on the line



$$\Rightarrow \sum_{i=1}^n [m_{f_1} - (\text{obs}_i)_{f_1}] + [m_{f_2} - (\text{obs}_i)_{f_2}] \\ = \text{Covariance}(f_1, f_2)$$

$$f = \frac{\text{Covariance}(f_1, f_2)}{c(f_1) + c(f_2)}$$

\* Slope determines the rate of change

### Logistic Regression:

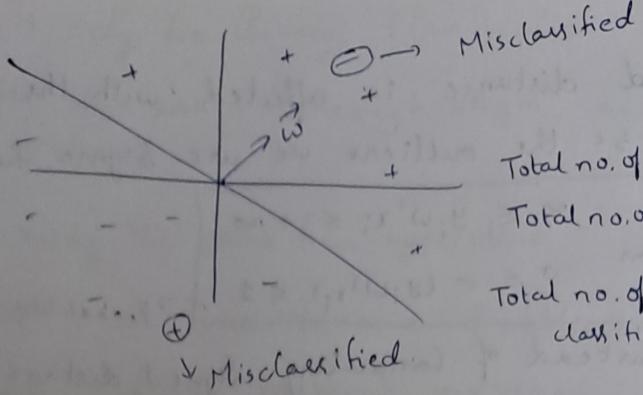
find a hyperplane that ~~Best Separator~~  $\rightarrow$  +ve is from -ve's.

max the correctly classified points

Signed distance  $\rightarrow$  act; \* predicted;

$\hookrightarrow$  -ve or +ve

Best model to be linear and correct rate of change is represented.



Total no. of Dpts = 13  
 Total no. of misclassified = 2  
 Total no. of correctly classified = 11

Case - 1:

if actual = +1 & predicted = +ve  
 actual \* predicted = +ve      Correctly classification

Case - 2:

if actual = +1 & predicted = -ve  
 actual \* predicted = -ve      Misclassification

Case - 3:

if actual = -1 & predicted = -ve  
 actual \* predicted = +ve      Correct Classification

Case - 4: if actual = -1 & predicted = +ve

actual \* predicted = -ve      Misclassification

Distance of a point from the hyperplane:

$$d = \frac{\vec{w} \cdot \vec{p} + w_0}{\sqrt{\sum w_i^2}}$$

Let assume the  $w_0$  is 0 { If it is passing through Origin }

Assume the  $\vec{w}$  is unit vector

$$\|\vec{w}\| = 1$$

$$w_0 = 0$$

$$d = \vec{w} \cdot \vec{p}$$

Signed distance is affected with the outliers. To minimize the outliers we use Sigma function then

$$-\infty \leq y_i w^T x_i \leq +\infty$$

After Sigma function  $0 \leq e^{-y_i w^T x_i} \leq 1$  → probability

Instead of computing signed distance for now, we will compute Sigmoid (Signed Distance)

$$w^* = \arg \max \left\{ \prod_{i=1}^n e^{-y_i w^T x_i} \right\}$$

$$= \arg \max \left\{ \prod_{i=1}^n \left( \frac{1}{1 + \exp(-y_i w^T x_i)} \right) \right\}$$

Problem → Numerical Instability Because the product below no. of values is zero.

Solution → log transformation.

$$w^* = \arg \max \left\{ \log \left( \prod_{i=1}^n \left( \frac{1}{1 + \exp(-y_i w^T x_i)} \right) \right) \right\}$$

$$= \arg \max \left\{ \sum_{i=1}^n \log \left( \frac{1}{1 + \exp(-y_i w^T x_i)} \right) \right\}$$

$$= \arg \max \left\{ \sum_{i=1}^n (\log(1) - \log(1 + \exp(-y_i w^T x_i))) \right\}$$

$$\log \left( \prod_{i=1}^n \frac{1}{1 + \exp(-y_i w^T x_i)} \right) = \sum_{i=1}^n \log \left( \frac{1}{1 + \exp(-y_i w^T x_i)} \right)$$

$$w^* = \arg \max \left\{ - \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i)) \right\}$$

### Optimization of Logistic Regression

Q) How do we solve the above Eq?

Any Optimizer → Gradient Ascent.

According to Optimization Theory

$$\max(-f(x)) = \min(f(x))$$

$$w^* = \arg \min \left\{ \sum_{i=1}^n \underbrace{\log(1 + \exp(-y_i w^T x_i))}_{\text{Loss function}} \right\}$$

final Optimization Equation.

→ Error function / Cost function  
Logistic loss function.

⇒ How do we solve?

Gradient Descent.

Logistic Regression is only for Binary Classification

- \* We assume the line passes through origin i.e., there is no intercept value.

Magnitude: how far away the point from hyperplane

In Logistic regression we have the issue of outliers  $\Theta \cdot P_7$

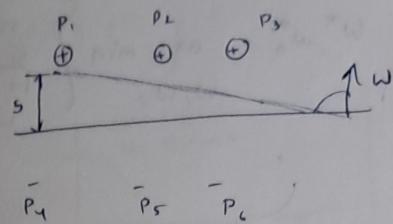
Solution: Sigmoid function to treat outliers.

for Eg :: First hyperplane

# of Data points = 7

# of Correctly classified = 6

# of Misclassified = 1



Signed distance:

acti \* predi

↓

$$y_i + w^T z_i$$

$$y_i + (\vec{w} \cdot \vec{z}_i)$$

$$\sum (acti * predi)$$

$$= [s + s + s + s + s + s - (-50)]$$

= -20

$$P_1 = +1 * (+ve\ value) = +ve\ mag = +5$$

$$P_2 = +1 * (+ve\ value) = +ve\ mag = +5$$

$$P_3 = - - - - - = +5$$

$$P_4 = -1 * (-ve\ value) = +ve\ mag = +5$$

$$P_5 = - - - - - = +5$$

$$P_6 = - - - - - = +5$$

$$P_7 = -1 * (+ve\ value) = -ve\ mag = -50$$

Second hyperplane

$$P_1 = +1 * (+ve\ value) = +1$$

$$P_2 = +1 * (+ve\ value) = +5$$

$$P_3 = +1 * (+ve\ value) = +10$$

$$P_4 = -1 * (+ve\ value) = -1$$

$$P_5 = -1 * (+ve\ value) = -5$$

$$P_6 = -1 * (+ve\ value) = -10$$

$$P_7 = -1 * (-ve\ value) = +1$$

P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>
+	+	+
P <sub>4</sub>	P <sub>5</sub>	P <sub>6</sub>

P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>
-	-	-
P <sub>4</sub>	P <sub>5</sub>	P <sub>6</sub>

$$\sum (acti * predi) = [s + y_0 + v + (-1) + (-5) + (-10) + (-1)] = -1$$

Outlier

## Support Vector Machine:

$\Rightarrow$   $D_{n \times d} \xrightarrow[\text{trick}]{\text{Kernel}} D_{n \times d'}$  i.e.,  $d' \gg d$   
 (exposed to larger dimensions)

It projects the data into higher dimensional space

## Linear Regression:

$$w^*, w_0^* = \arg \min_{w, w_0} \left\{ \sum_{i=1}^n (y_i - (w^T x_i + w_0))^2 \right\}$$

SSE Loss function

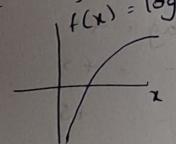
Optimization fn for linear Regression.

$$= \arg \min_{w, w_0} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - (w^T x_i + w_0))^2 \right\}$$

MSE loss function

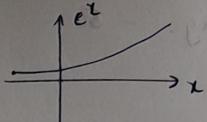
## Gradient Descent:

$$x^* = \arg \min_x \{f(x)\}$$

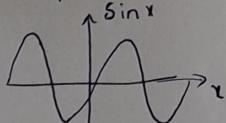


$$\star f(x) = \log x$$

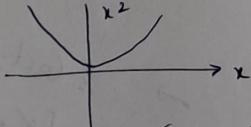
$$\star f(x) = e^x$$



$$\star f(x) = \sin x$$



$$\star f(x) = x^2$$



$$x^* = \arg \min_x \{x^2\} \rightarrow \left\{ \begin{array}{l} \text{find that value of } x \\ \text{which minimizes } x^2 \text{ function} \end{array} \right\}$$

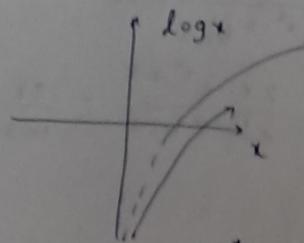
## Linear Reg:

$$w^*, w_0^* = \arg \min_{w, w_0} \left\{ \frac{1}{n} \sum (y_i - (w^T x_i + w_0))^2 \right\} \rightarrow \left\{ \begin{array}{l} \text{find the value of } w \\ w_0 \text{ which minimizes} \\ \text{the MSE loss} \\ \text{function} \end{array} \right\}$$

$$f(x) =$$

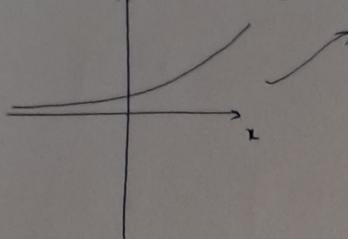


$\Rightarrow$



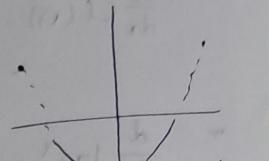
No Maxima  
No Minima

It is a continuous &  
always increasing fn  
 $\rightarrow$  Monotonic fn



One Maxima  
No Minima

Concave fn



Convex fn



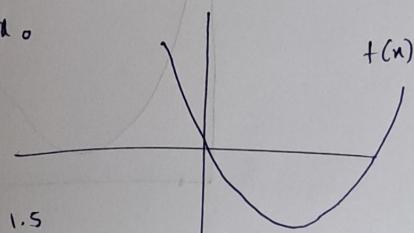
\* Linear & Logistic Regression are Convex functions.

Problem:

$f(x) = x^2 - 3x + 2$ . Is it is Convex or Concave fn

differentiation of a function @ a point  $x_0$

$\Rightarrow$  Slope of tangent @  
point  $x_0$



$$\frac{d}{dx} f(x) = \frac{d}{dx} (x^2 - 3x + 2) \Rightarrow 2x - 3 = 0 \Rightarrow x = 1.5$$

$$@ x = 1.5$$

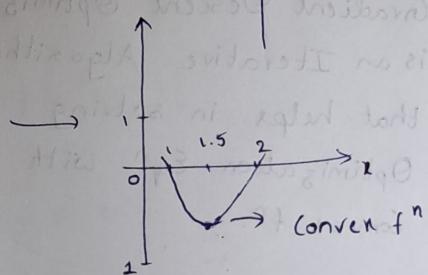
$$f(x) = (1.5)^2 - 3(1.5) + 2 = -0.25$$

Then we check for prior and after pts

$$@ x = 1 \Rightarrow f(x) = 0$$

$$@ x = 2 \Rightarrow f(x) = 0$$

So, here it is minima



Optimizer is a function / also that helps you solve Optimization function (with Convex fn's).

Gradient Descent should be used for Convex fn's then only it gives proper results.

27/12/2022

$$f(x) = e^{x^2+3}$$

$$= e^{x^2+3} \frac{d}{dx}(x^2+3)$$

$$= e^{x^2+3} \cdot (2x)$$

$$= 2x \cdot e^{x^2+3}$$

$$f(x) = \log_3 e^{x^2+3}$$

$$= \frac{1}{e^{x^2+3} \ln 3} \cdot \frac{d}{dx}(e^{x^2+3})$$

$$= \frac{1}{e^{x^2+3} \ln 3} \cdot e^{x^2+3} \cdot 2x$$

$$= \frac{2x \cdot e^{x^2+3}}{e^{x^2+3} \ln 3} = \frac{2x}{\ln 3}$$

$$\frac{d}{dx} f(x) = \frac{d}{dx} ax = a \frac{dx}{dx} = a$$

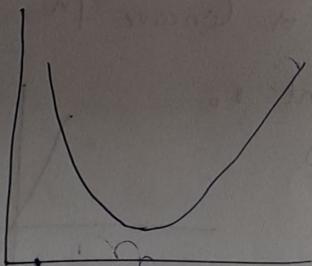
$$\frac{d}{dx}(f(x)) = \frac{d}{dx} e^x = e^x \frac{d}{dx}(x) = e^x$$

$$\frac{d}{dx} \log_a x = \frac{1}{x \ln a} \cdot \frac{dx}{dx}$$

↳ natural logarithm

$$\frac{d}{dx} x^n = n \cdot x^{n-1}$$

## Gradient Descent:



$$x^* = \underset{x}{\operatorname{argmin}} \{f(x)\}$$

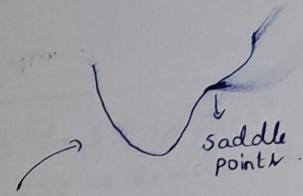
Assuming  $f(x)$  to be a convex function

### Algorithm:

1. Randomly pick the value of  $x_0$  (minima if slope of tangent  $= 0$ )
2. Move forward minimize & in order to do that we have to move either towards left or right
3. Reiterate step 2 until you reach minima.

Gradient Descent Optimizer is an Iterative Algorithm that helps in solving Optimization Eqn with Convex fn.

- ① Randomly initialize  $\rightarrow x_i$
  - ② Move towards minima  $\rightarrow$  hyperparameter of Gradient Descent
- $$x_{i+1} = x_i - \eta \left[ \frac{d}{dx} f(x) \right]$$
- ↑ learning rate
- [Update Eqn of Gradient Descent]



Issue  $\rightarrow$  Oscillation Problem  
 $\rightarrow$  Gradient Descent can get stuck at Saddle points

### Hyperparameter of GD

①  $\eta$  = Learning rate

② # of iteration.

$$\begin{aligned}
 m_{i+1} &= m_i - \eta \left[ -\frac{\sum_{i=1}^n}{n} \right] \\
 m_{i+1} &= m_i - \eta \left[ \frac{\partial}{\partial m} f(m) \right] m_i \\
 &= m_i - \eta \left[ \frac{\partial}{\partial m} \left[ + \sum_{i=1}^n \log \left( 1 + \exp \left\{ -y_i w^\top x_i \right\} \right) \right] \right] m_i \\
 &= m_i - \eta \left[ \frac{\partial}{\partial m} \left[ + \sum_{i=1}^n \log \left( 1 + \exp \left\{ -y_i m_i x_i \right\} \right) \right] \right] m_i \\
 &= m_i - \eta \left[ + \sum_{i=1}^n \frac{\partial}{\partial m} \left[ \log \left( 1 + \exp \left\{ -y_i m_i x_i \right\} \right) \right] \right] m_i \\
 &= m_i - \eta \left[ + \sum_{i=1}^n \left[ \left( \frac{\partial}{\partial m} (\log) \right) + \frac{\partial}{\partial m} \left( \log \left( \exp \left\{ -y_i m_i x_i \right\} \right) \right) \right] m_i \right] \\
 &= m_i - \eta \left[ + \sum_{i=1}^n \left[ [0] + \frac{1}{\exp \left\{ -y_i m_i x_i \right\}} \ln a \cdot \frac{\partial}{\partial m} \left( \exp \left\{ -y_i m_i x_i \right\} \right) \right] m_i \right]
 \end{aligned}$$

$$\begin{aligned}
 &= m_i - n \left[ \sum_{i=1}^n \frac{\partial}{\partial m_i} \left[ \log \left( \overbrace{1 + \exp \{-y_i m_i x_i\}}^{} \right) \right] \right]_{m_i} \\
 &= m_i - n \left[ \frac{1}{(1 + \exp \{-y_i m_i x_i\}) \log 2} \cdot \frac{\partial}{\partial m_i} (\exp \{-y_i m_i x_i\}) \right] \\
 &= m_i - n \left[ \frac{1}{1 + \exp \{-y_i m_i x_i\} \log 2} \left[ \frac{\partial}{\partial m_i}(1) + \frac{\partial}{\partial m_i}(\exp \{-y_i m_i x_i\}) \right] \right] \\
 &= m_i - n \left[ \frac{1}{1 + \exp \{-y_i m_i x_i\} \log 2} \left[ 0 + \exp \{-y_i m_i x_i\} \frac{\partial}{\partial m_i}(-y_i m_i x_i) \right] \right] \\
 &= m_i - n \left[ \frac{1}{1 + \exp \{-y_i m_i x_i\} \log 2} \left[ \exp \{-y_i m_i x_i\} \cdot (-y_i x_i) \right] \right] \\
 &\quad \text{Note: } \log 1 = 0 \\
 &= m_i - n \left[ \frac{1}{1 + 0} \left[ \exp \{y_i m_i x_i\} \cdot (y_i x_i) \right] \right] \\
 &= m_i - n \left[ \exp \{y_i m_i x_i\} \cdot (y_i x_i) \right]
 \end{aligned}$$

Derivative of Sigmoid fn?

$$\begin{aligned}
 \frac{d}{dx} \sigma(x) &= \frac{d}{dx} \left( \frac{1}{1 + e^{-x}} \right) \\
 &= (1 + e^{-x})^{-1} \\
 &= -1 \cdot (1 + e^{-x})^{-2} \cdot \frac{d}{dx} (1 + e^{-x}) \\
 &= -1 \cdot (1 + e^{-x})^{-2} \cdot (0 + e^{-x}) \\
 &= -\frac{1}{(1 + e^{-x})^2} \cdot e^{-x} \\
 &= \frac{-e^{-x}}{(1 + e^{-x})^2} \\
 &= \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{e^{-x}}{(1 + e^{-x})(1 + e^{-x})} \\
 &= \frac{e^{-x}}{1 + 2e^{-x} + (e^{-x})^2} = \frac{e^{-x}}{1 + 2e^{-x} + e^{-2x}}
 \end{aligned}$$

$$\frac{d}{dx} \sigma(x) = \frac{e^{-x}}{(1+e^{-x})} + \frac{1}{(1+e^{-x})}$$

$$= (1 - \sigma(x)) * \sigma(x)$$

Linear Regression takes some time for training whereas for predicting it takes less time.

For Logistic Regression,

$$D_n = \{(x_i, y_i)\}_{i=1}^n \mid x_i \in \mathbb{R}^{d-1}, y_i \in \{-1, +1\}\}$$

why because when we take  $0_+$ , then signed distance becomes '0'.

Evaluation matrix: (for Regression)

$$MAE \rightarrow \frac{\sum}{n} \underbrace{\text{abs}(y_{\text{acti}} - y_{\text{predi}})}_{\text{error}}$$

$$* MSE \rightarrow \frac{\sum}{n} (y_{\text{acti}} - y_{\text{predi}})^2$$

$$* RMSE \rightarrow \sqrt{\frac{\sum}{n} (y_{\text{acti}} - y_{\text{predi}})^2}$$

$\sum (y_{\text{acti}} - y_{\text{predi}})^2$   
RSS  $\rightarrow$  Residual Sum of squares

$$* R^2 \rightarrow \frac{1}{1 - \frac{RSS}{TSS}}$$

$$[-\infty \leq R^2 \leq 1]$$

$$* \text{Adjusted } R^2 = 1 - \left[ \frac{(1-R^2)(n-1)}{n-p-1} \right]$$

$$RSS = \sum_{i=1}^n (\text{residual})^2$$

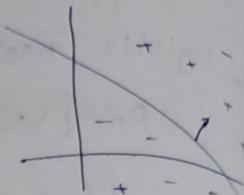
$$= \sum_{i=1}^n (y_{\text{acti}} - y_{\text{predi}})^2$$

$$TSS = \sum (y_{\text{acti}} - \bar{y})^2$$

## Evaluation Matrix (for Classification) :

Accuracy  $\rightarrow \frac{\# \text{ of Datapoints which were correctly classified}}{\text{Total } \# \text{ of Datapoints}}$

$$\rightarrow \frac{8}{10} = 0.8$$



S.N.	Actual	Predicted	So called
1	A	Not A	True Positive
2	Not A	Not A	False True Negative
3.	A	Not A	False Negative
4.	Not A	A	False Positive

← Confusion Matrix  $\rightarrow$  By this matrix we determine the accuracy of predicted & actual values

$$\rightarrow \text{Accuracy} = \frac{TP + TN}{(TP + TN + FP + FN)}$$

$$\rightarrow \text{Precision} = \frac{TP}{(TP + FP)}$$

$$\rightarrow \text{Recall} = \frac{TP}{(TP + FN)}$$

$$\rightarrow \text{F1 Score} = \frac{2(p * r)}{(p + r)} \quad \therefore p - \text{precision} \quad r - \text{recall}$$

from picture:

		Predicted	
		Positive	Negative
Actual	Positive	TP	FN
	Negative	FP	TN