

Machine Learning Algorithms in short

① Linear Regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Cost function \rightarrow Squared Error problem

$$\Rightarrow \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x) - y)^2$$

minimize it by adjusting θ_0, θ_1

Convergence Algorithm

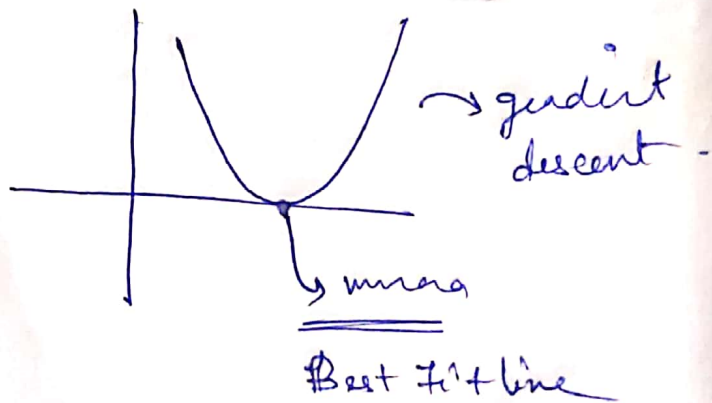
repeat until {
converge

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

}

learning rate

$\alpha \rightarrow$ generally
0.01



Reduced convergence algorithm

repeat until converge

$$\left\{ \begin{aligned} \theta_0 &= \theta_0 - \alpha \frac{1}{m} (h_0(x) - y) \\ \theta_1 &= \theta_1 - \alpha \frac{1}{m} (h_0(x) - y)x \end{aligned} \right\}$$

Performance Metrics

$$R^2 = 1 - \frac{SS_{\text{residual}}}{SS_{\text{total}}} = 1 - \frac{\sum (y_i - h_0(x_i))^2}{\sum (y_i - \bar{y})^2}$$

→ The problem with R^2 is that it works in such a way that if features are added even if they are no way correlated then R^2 increases

$$\text{Adjusted } R^2 = \frac{1 - (1 - R^2)(N - 1)}{N - P - 1}$$

$N - P - 1$
↙ ↘
no. of data points no. of features

Assumptions of Linear Regression

① If features are in Normal/Gaussian Distribution, model will get trained well

② Standardisation { scaling data } with z score

$$\mu = 0, \sigma = 1$$

③ Linearity

Note:-

Overfitting:- When model performs well with training data but fails to perform well with testing data

→ Low Bias and High variance

→ Model is too complex and size of data set is small

Underfitting:- Model performs bad with training and testing data set.

→ High Bias and low variance

Ridge and Lasso

Ridge (L2 Regularisation)

$$\sum (\hat{y}_i - y)^2 + \lambda (\text{slope})^2$$

→ To prevent overfitting

→ Hypersensitive

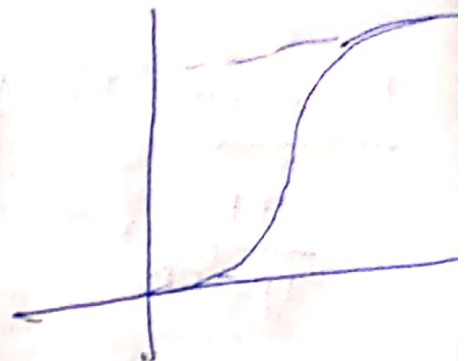
Lasso (L1 Regularisation)

$$\sum (\hat{y}_i - y)^2 + \lambda \|\text{slope}\|$$

→ helps in feature selection

Logistic Regression

$$h_\theta(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}} \rightarrow \frac{1}{1 + e^{-z}}$$



Cost function

$$J(\theta_i) = \begin{cases} -\log(h_\theta(x)) & y=1 \\ -\log(1-h_\theta(x)) & y=0 \end{cases}$$

$$\text{overall} \Rightarrow \frac{1}{n} [-y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))]$$

repeat until convergence

$$\theta_j = \theta_j - \alpha \frac{d}{d\theta_j} (J(\theta_i))$$

}

Performance metrics {classification problem}

	Actual	
Pred	1	0
	TP	FP
	FN	TN

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

<u>Precision (P)</u>	<u>Recall (Sensitivity)</u>
$\frac{TP}{TP + FP}$	$\frac{TP}{TP + FN}$

$$\text{F-score} \rightarrow (1 + \beta^2) \cdot \frac{P \times R}{\beta^2 P + R}$$

$$\beta = 2 \quad 12N \neq PP$$

$$\beta = 0.5 \quad FP \neq PR$$

$$\beta = 1 \quad FN = FP$$

Naive Bayes' Algorithm

$$P(A \text{ and } B) = P(A) \times P(B/A)$$

$$P(A \text{ and } B) = P(B \text{ and } A)$$

$$\boxed{P(A) \times P(B/A) = P(B) \times P(A/B)} \rightarrow \text{Bayes' theorem}$$

$$P(B/A) = \frac{P(B) \times P(A/B)}{P(A)}$$

Let B be output and A be input

$$P(y/x_1, x_2, \dots, x_n) = \frac{P(y) \times P(x_1, x_2, x_3, \dots, x_n/y)}{P(x_1) \times P(x_2) \times \dots \times P(x_n)}$$

$$P(y/x_1, x_2, \dots, x_n) \Rightarrow \frac{P(y) \times P(x_1/y) \times \dots \times P(x_n/y)}{P(x_1) \times P(x_2) \times P(x_3) \times P(x_4) \times \dots \times P(x_n)}$$

Note:- Suppose we do binary classification then any probability > 0.5 will be 1 and < 0.5 will be equal to 0.

→ Normalisation

$$P(\text{yes}/x_i)_{\text{normal}} = \frac{P(\text{yes}/x_i)}{P(\text{yes}/x_i) + P(\text{no}/x_i)}$$

Q) Dataset

Day	Outlook	Temp	Humidity	wind	Playtime
				windy	No
1	Sunny	Hot	High	slow	No
2	Sunny	Hot	High	windy	yes
3	Overcast	Hot	High	slow windy	yes
4	Rain	mild	High	windy	yes
5	Rain	cool	Normal	slow	no
6	Rain	cool	normal	slow	yes
7	Overcast	cool	normal	windy	no
8	Sunny	mild	High	windy	yes
9	Sunny	cool	Normal	windy	yes
10	Rain	mild	Normal	slow	yes
11	Sunny	mild	High	slow	yes
12	Overcast	Hot	normal	windy	yes
13	Overcast	mild	High	slow	no
14	Rain				

Find $P(\text{yes} | \text{sunny, hot})$

$$\Rightarrow \frac{P(\text{yes}) \times P(\text{sunny/yes}) \times P(\text{hot/yes})}{P(\text{sunny}) \times P(\text{hot})}$$

Outlook (x_i)

	Y	N
Sunny	2	3

Overcast	4	0
----------	---	---

Rain	3	2
	9	5

$P(\frac{x_i}{y_1})$ $P(\frac{x_i}{no})$

$\frac{2}{9}$ $\frac{3}{5}$

$\frac{4}{9}$ $\frac{0}{5}$

$\frac{3}{9}$ $\frac{2}{5}$

$P(\frac{x_i}{yes})$ $P(\frac{x_i}{no})$

$\frac{2}{9}$

$\frac{2}{5}$

$\frac{4}{9}$

$\frac{2}{5}$

$\frac{3}{9}$

$\frac{1}{5}$

Temp

	Y	N
Hot	2	2

Hot

mild

cold

	3	1
--	---	---

	9	5
--	---	---

\Rightarrow

$$\frac{1}{14} \times \frac{2}{9} \times \frac{2}{9} \Rightarrow \frac{2}{63}$$

$$P(no/sunny, hot) = \frac{1}{14} \times \frac{3}{5} \times \frac{2}{5} \Rightarrow \frac{3}{35}$$

\Rightarrow

$P(yes/u_1, u_2)$

$$\frac{2}{63}$$

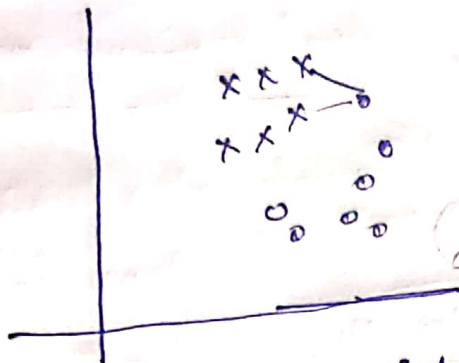
$$\frac{3}{35} + \frac{2}{63}$$

normalisation

K nearest Neighbours

↓
Classification

↓
Regression

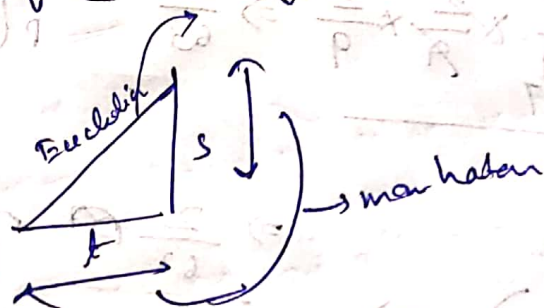


- Find nearest K points
- If the ~~most~~ no. of points $\geq \frac{K}{2}$ in any of the subsets (types) then it belongs to that group

Eucledian distance = $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

Manhattan distance

→ $t + s$



- For regression the predicted value will be avg. of K nearest neighbours

Problems faced with K nearest neighbours

① Outlier

② Imbalanced dataset.

Decision

Tree

- Regression
- Classification → Nested if else

Classification of a particular leaf node should be done until

pure split

→ pure split is when there is purely one class.

→ D-Tree is made by 2 diffrent

Parameters

- Entropy
- Gini Index



Entropy

$$H(s) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$GI = 1 - \sum_{i=1}^n (P_i)^2$$

Features are selected using Information gain.

$$\text{Gain}(S, A) = H(S) - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S_v)$$

Annotations for the equation above:

- $H(S)$: Entropy of root node
- $|S_v|/|S|$: no. of samples in leaf nodes
- $H(S_v)$: entropy of leaf node
- $|S|$: no. of samples in root node

Hyperparameter → Problem with D-Tree is overfitting

→ Post pruning (cutting after seeing the tree)

→ Pre pruning

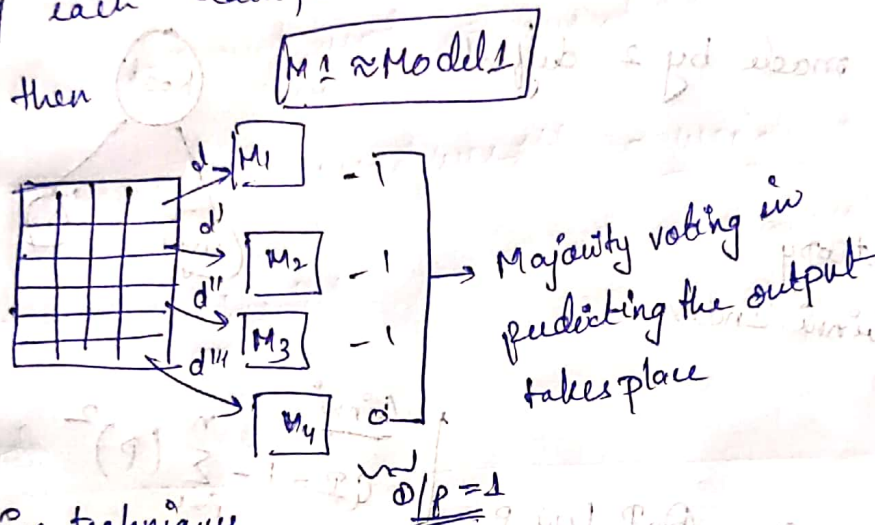
Note:- More the Info gain of a feature, splitting takes place first

Ensemble Learning

- Bagging
- Boosting

Custom Bagging technique

- Dividing the datapoints to various sets
 - Training each datapoint with various ML algorithms
- and then



Bagging techniques

- ① Random Forest classifiers
- ② Random Forest Regression [Mean of M_1, M_2, M_3, M_4]

In Random Forest M_1, M_2, M_3, M_4 are all D-Trees

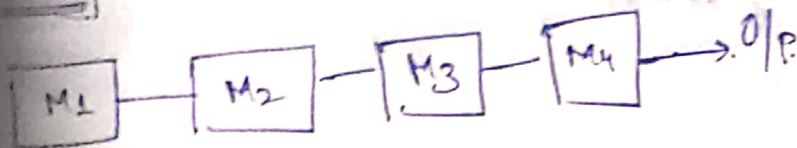
- No need of normalisation
- Not affected by outliers

Outliers

- An outlier is an observation that lies in a abnormal distance from other values in a random sample from population

Boosting

Boosting



Weak
Learner

→ Multiple ML model in sequence are used to predict the output

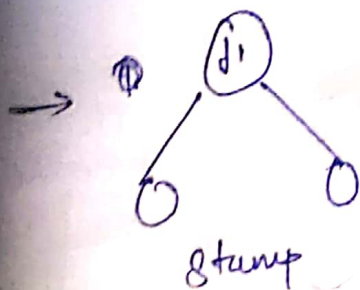
→ O/P is got by voting again

→ These data points and features are selected.

Adaboost

	d1	d2	d3	d4	O/P	Weight
1						$\frac{1}{4}$
2						$\frac{1}{4}$
3						$\frac{1}{4}$
4						$\frac{1}{4}$
5						$\frac{1}{4}$
6						$\frac{1}{4}$
7						$\frac{1}{4}$

→ Assigning equal weights for all nodes



Stump

1 level D-Tree is called a stump
 the error is then found out
 Let the no. of errors be n
 and $\frac{1}{n}$

If ~~no. of~~ weight of 1 word is w and no. of words are

→ n then $T_E = w \times n$
(Total error)

→ Performance of stump $(P_s) = \frac{1}{2} \log_e \left(\frac{1 - T_E}{T_E} \right)$

→ new sample weight $\underset{\text{for count data}}{=} \text{weight} \times e^{-P_s}$

Incorrect word weight $= \text{weight} \times e^{P_s}$

→ normalised word weight $= \frac{\text{weight of it on}}{\text{Total (which is } n)}$

then create bucket

~~example~~
