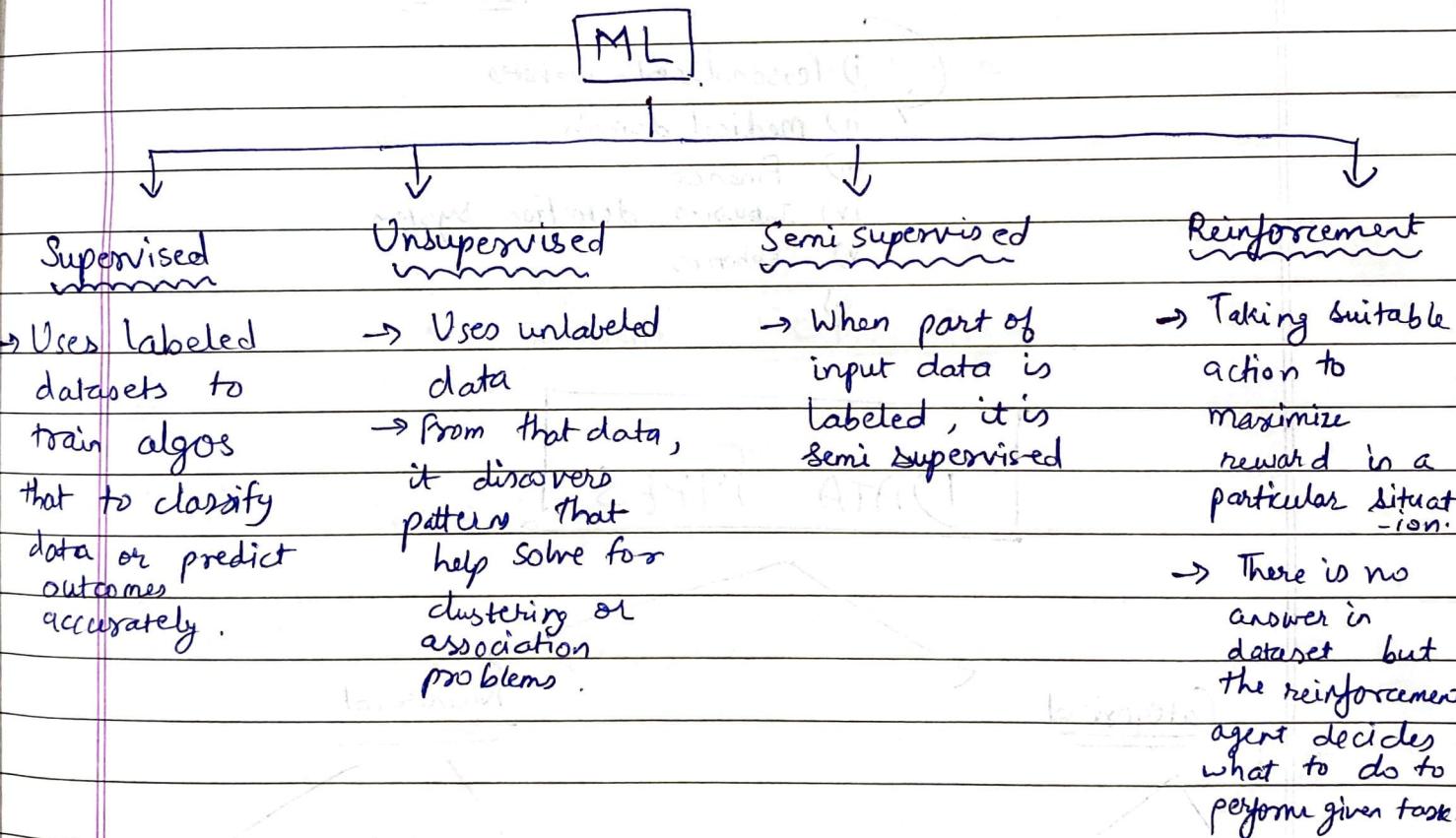


What is ML?

→ Branch of AI and CS which focuses on use of data and algos to imitate the way that humans learn, gradually improving its accuracy.



### Applications

#### NLP

- Autonomous tagging of stackoverflow ques
- Automated essay grading
- Sentence to sentence semantic similarity
- Fight online abuse
- Open domain question answering
- Social chat / conversational bots

#### Dataset

- |  |
|--|
| Stacklite<br>Essays with human graded scores<br>Quora ques pairs with similar ques marked<br>Toxic comments on Kaggle<br>NCERT books for K-12 school students in India<br>Reddit dataset |
|--|

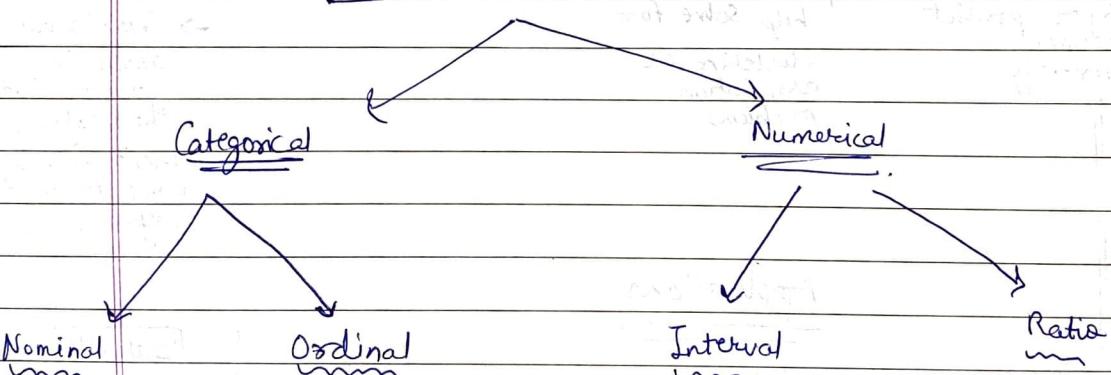
\* Applications that can't be programmed by hand :-

- Handwriting Recognition
- NLP : (i) Chatbot (ii) Speech Recognition (iii) Text <sup>summarization</sup>
- Part-of-Speech Tagging
- Sentiment Analysis

Many more applications in PPT (i)

- i) Personalised websites
- ii) Medical domain
- iii) Finance
- iv) Intrusion detection system
- v) Robotics
- vi) ...

## DATA TYPES



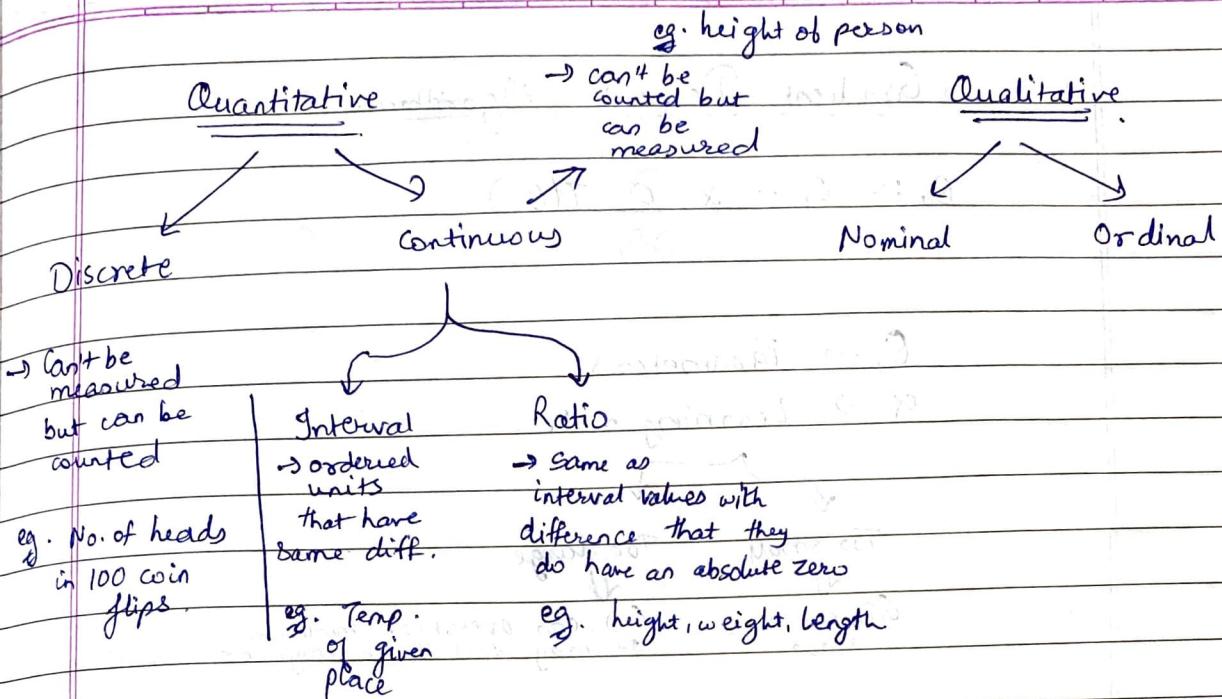
- No order
- Non quantitative value
- Ordering matters
- Non numeric feature

e.g. Gender

\* Male  
\* Female

e.g. Good

Better  
Best



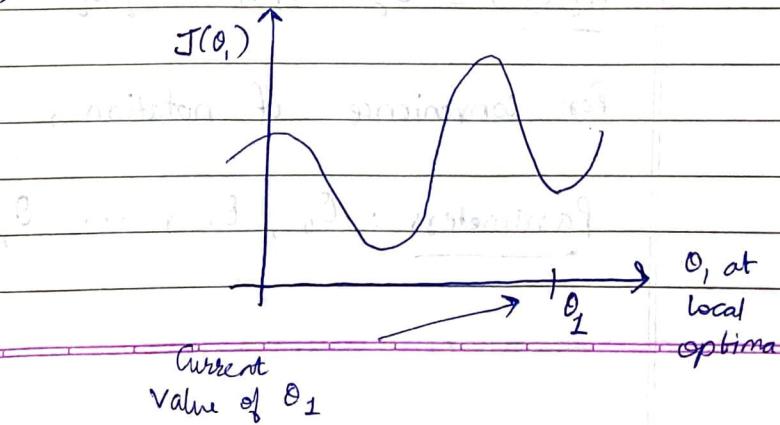
## Simple linear Regression

$$\hat{y} = h_0(x) = \theta_0 + \theta_1 x$$

minimize  $\sum_{i=1}^m (y_i - \hat{y}_i)^2$

$$\Rightarrow \text{Cost function}, J(\theta_0, \theta_1) = \left(\frac{1}{2m}\right) \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

$\Rightarrow$  Minimize  $J(\theta_0, \theta_1)$   
 $\theta_0, \theta_1$



## Gradient Descent Algorithm:

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

$\theta_j \rightarrow$  Parameters

$\alpha \rightarrow$  Learning rate.

Too small  
↓

GD can be slow

Too large  
↓

GD can overshoot minimum.  
It may fail to converge or even diverge

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \quad \left. \begin{array}{l} \text{update} \\ \theta_0 \text{ and } \theta_1 \end{array} \right\}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \quad \left. \begin{array}{l} \text{simultaneously} \\ \text{update} \end{array} \right\}$$

}

## Multiple Linear Regression

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n = \theta^T x$$

For convenience of notation, define  $x_0 = 1$ .

Parameters :  $\theta_0, \theta_1, \dots, \theta_n$ .

Cost function :

$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

Gradient descent :

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$$

}

(simultaneously update for every  $j=0, \dots, n$ )

New algo ( $n \geq 1$ ) { $i \rightarrow 1$ }

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

}

(simultaneously update  $\theta_j$  for  $j=0, \dots, n$ )

## Feature Scaling

Get every feature into approximately  $-1 \leq x_i \leq 1$  range

### \* Min Max normalization

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{newMax}_A - \text{newMin}_A) + \text{newMin}_A$$

→ For mapping to  $[0, 1]$  range, formula becomes

$$v' = \frac{v - \min_A}{\max_A - \min_A} (1 - 0) + 0$$

$$v' = \frac{v - \min_A}{\max_A - \min_A}$$

### \* Z-score normalization

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

$$v' = \frac{v - \bar{x}}{\sigma_A}$$

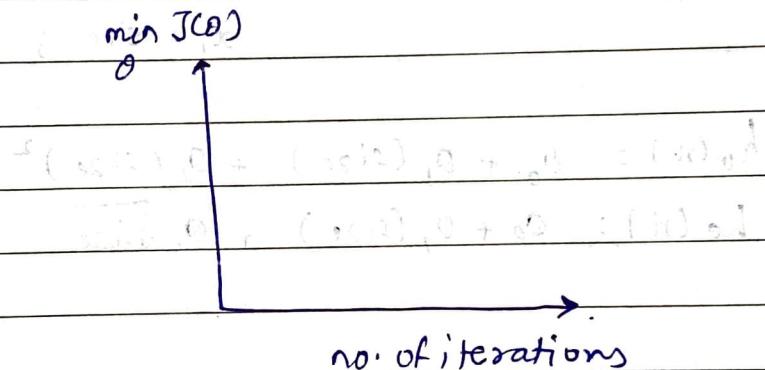
Learning rate

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

"Debugging" : How to make sure gradient descent is working correctly

How to choose learning rate  $\alpha$ .

→ Declare convergence if  $J(\theta)$  decreases by less than  $10^{-3}$  in one iteration.



GD not working  $\Rightarrow$  use smaller  $\alpha$

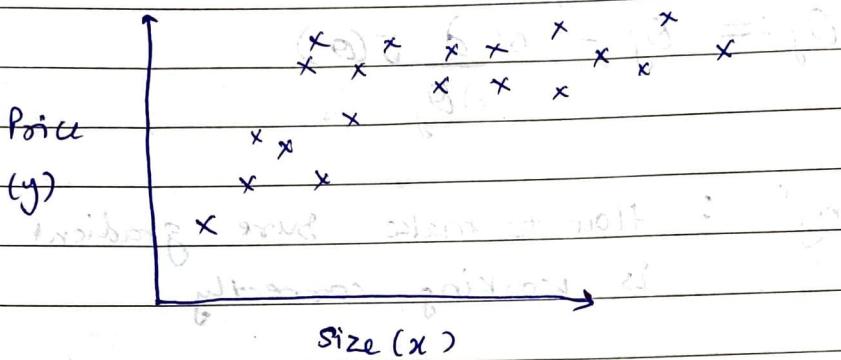


For sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration.

\* But if  $\alpha$  is too small, G.D can be slow to converge.

If  $\alpha$  is too large,  $J(\theta)$  may not decrease on every iteration; may not converge.

## Features & Polynomial regression



$$h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$\text{and Parameters } (\theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3)$$

$$x_1 = \text{size} ; x_2 = (\text{size})^2 ; x_3 = (\text{size})^3$$

$$h_0(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2$$

$$h_0(x) = \theta_0 + \theta_1(\text{size}) + \theta_2 \sqrt{\text{size}}$$

contribution

### Normal Equation

	$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	2104	(5)	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 1 & 95 \\ 1 & 1416 & 3 & 2 & 4 & 0 \\ 1 & 1534 & 3 & 2 & 3 & 0 \\ 1 & 852 & 2 & 1 & 3 & 6 \end{bmatrix} \quad y = \begin{bmatrix} 900 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

$$\theta = (X^T X)^{-1} X^T y$$

training Examples  $\rightarrow m$

features  $\rightarrow n$

G. D

N.E.

- Need to choose  $\alpha$
- No need to choose  $\alpha$
- Needs many iterations
- Don't need to iterate
- Works well even when  $n$  is large
- Need to compute  $(X^T X)^{-1}$
- Slow if  $n$  is very large

# What if  $X^T X$  is non-invertible?

→ Redundant features (linearly dependent)

e.g.  $x_1 = \text{size in feet}^2$

$x_2 = \text{size in m}^2$

→ Too many features (e.g.  $m \leq n$ )

\* Delete some features, or use regularization



Eq<sup>n</sup> of least squares line

$$y = w_0 + w_1 x$$

Here,  $w_0 = \bar{y} - w_1 \bar{x}$

$$w_1 = \frac{\sum_{i=1}^d (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^d (x_i - \bar{x})^2}$$

### Error measures

Absolute error :  $|y_i - y'_i|$

Squared error :  $(y_i - y'_i)^2$

Mean absolute error :  $\frac{1}{d} \sum_{i=1}^d |y_i - y'_i|$

Mean absolute % error :  $\frac{1}{d} \sum_{i=1}^d \left| \frac{y_i - y'_i}{y_i} \right| \times 100$

RMSE :  $\sqrt{\frac{1}{d} \sum_{i=1}^d (y_i - y'_i)^2}$

NRMSE =  $\frac{\text{RMSE}}{y_{\max} - y_{\min}}$

For Normal eqn :-

$$\theta = (X^T X + \lambda L)^{-1} X^T y$$

$$L = \begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & & 0 \\ 0 & 0 & \ddots & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

is square matrix of size  $(n+1) \times (n+1)$

## Overfitting

→ If we have too many features, learned hypothesis may fit training set very well (i.e.  $J(\theta) \approx 0$ ), but fail to generalize to new examples

Options :- (i) Reduce no. of features

- Manually select features to be kept
- Model selection algo

(ii) Regularization

- Keep all features, but reduce magnitude of parameters  $\theta_j$ .
- Works well when we have a lot of features, each of which contributes a bit to predicting  $y$ .

## Regularization

Small values of parameters  $\theta_0, \theta_1, \dots, \theta_n$

- "Simpler" hypothesis
- Less prone to overfitting

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

For regularized Linear Regression,  $\min_{\theta} J(\theta)$

G.D :-

repeat until convergence {

$$\theta_0 := \theta_0 - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)} \right]$$

$$\theta_j := \theta_j \left( 1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \left[ \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right]$$

where  $j \in \{1, 2, \dots, n\}$

# Logistic Regression

an algorithm  $h_0(x) = g(\theta^T x)$ , where  $g(z) = \frac{1}{1+e^{-z}}$

$$\therefore h_0(x) = \boxed{h_0(x) = \frac{1}{1+e^{-\theta^T x}}}$$

$h_0(x)$  = estimated probability that  $y=1$  on input  $x$

## Accuracy Measures

### Validation techniques

- Resubstituting
- Hold-out
- K fold cross validation.
- Leave one out cross validation
- Random subsampling.
- Bootstrapping.

### Confusion matrix

Actual class / Predicted class	$C_1$	$\sim C_1$
$C_1$	TP	FN
$\sim C_1$	FP	TN

$$\text{# Accuracy} = \frac{(TP + TN)}{AU}$$

$$\begin{aligned}\text{# Error rate} &= 1 - \text{accuracy} \\ &= \frac{(FP + FN)}{AU}\end{aligned}$$

$$\text{# Sensitivity} = \frac{TP}{P}$$

$$\text{# Specificity} = \frac{TN}{N} \quad (TN \text{ rate})$$

$$\text{# } 1 - \text{Specificity} = \frac{FP}{N} \quad (FP \text{ rate})$$

# Precision =  $\frac{TP}{TP + FP}$

# Recall =  $\frac{TP}{TP + FN}$

# F-score,  $F = \frac{2(\text{Precision})(\text{Recall})}{\text{Precision} + \text{Recall}}$

#  $F_\beta = \frac{(1+\beta^2)(\text{Precision})(\text{Recall})}{\beta^2 \times \text{Precision} + \text{Recall}}$

Additional aspects for model selection

- Accuracy classifier accuracy : predicting class label
- Speed Training time ; prediction time
- Robustness handling noise & missing values
- Scalability efficiency in disk-resident databases
- Interpretability understanding & insight provided by model

## Total Probability theorem :-

$$P(B) = \sum_{i=1}^M P(B|A_i) P(A_i)$$

## Bayes' Theorem :-

$$P(H|X) = \frac{P(X|H) \cdot P(H)}{P(X)}$$

↑ Prior probability

(Numerator) (Conditioned) (denominator)  
↓ posterior probability

with Gaussian      Naive Bayes

$$P(C_i|X) = \frac{P(X|C_i) P(C_i)}{P(X)}$$

Now since  $P(X)$  is constant for all classes,

$$\boxed{P(C_i|X) = P(X|C_i) P(C_i)} \text{ needs to be maximised.}$$

$$\text{Now } P(X|C_i) = \prod_{k=1}^n P(x_k|C_i)$$

$$= P(x_1|C_i) \times P(x_2|C_i) \times \dots \times P(x_n|C_i)$$

Now if  $A_k$  is categorical  $\Rightarrow P(x_k | c_i)$

$$= \frac{\text{# of tuples in } c_i \text{ having value } x_k \text{ for } A_k}{\text{# of tuples of } c_i \text{ in } D}$$

If  $A_k$  is continuous  $\Rightarrow P(x | c_i) = g(x_k, \mu_{c_i}, \sigma_{c_i})$

$$g(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

④ In order to avoid zero probability problem, add Laplace correction i.e. add 1 to each case

⑤ For a remedy of zero conditional probability problem, conditional probabilities estimated with ( $m$ -estimates) are used.

$$\hat{P}(x_j = a_{jk} | C = c_i) = \frac{n_c + mp}{n + m}$$

$n_c \rightarrow$  no. of tuples for which  $x_j = a_{jk}$  &  $C = c_i$

$n \rightarrow$  no. of tuples for which  $C = c_i$

$p \rightarrow$  prior estimate [usually,  $p = 1/t$  for  $t$  possible values]

$m \rightarrow$  weight to prior (no. of "virtual" examples,  $m \geq 1$ ) of  $x_j$

# Multinomial Naive Bayes

$$\hat{P}(x_i|w_j) = \frac{\sum t f(x_i, d \in w_j) + \alpha}{\sum N_{d \in w_j} + \alpha V}$$

where,

$x_i \rightarrow$  Word from feature vector  $x$  of particular sample

$\sum t f(x_i, d \in w_j) \rightarrow$  sum of raw term freq. of word  $x_i$  from all documents belonging to class  $w_j$ .

$\sum N_{d \in w_j} \rightarrow$  sum of all term frequencies in training dataset for class  $w_j$ .

$\alpha \rightarrow$  An additive smoothing parameter ( $\alpha = 1$  for Laplace smoothing)

$V \rightarrow$  no. of different words in training set?

# Multivariate Bernoulli Naive Bayes

$$P(X|w_j) = \prod_{i=1}^m P(x_i|w_j)^b \cdot (1 - P(x_i|w_j))^{(1-b)} \quad (b \in \{0, 1\})$$

Now,  $\hat{P}(x_i|w_j) = \frac{df_{x_i,y} + 1}{df_y + 2}$

where,

$df_{x_i,y} \rightarrow$  no. of documents that contain  $x_i$  and belong to class  $w_j$ .

$df_y \rightarrow$  no. of documents belong to class  $w_j$ .

+1 and +2  $\rightarrow$  Laplace smoothing.

Doubt Why formula in theory and sum solving are diff?

Formula in problem solving :-

$$P(S/b_i) \propto P(S) \prod_{i=1}^n [b_{ii} P(w_i | S) + (1-b_{ii})(1-P(w_i | S))]$$

## DECISION TREES

### # ID3 Algorithm

→ Select attribute that has highest information gain

$$\rightarrow \text{Entropy}(S) = \sum_{i=1}^c -p_i \log_2 p_i \quad \text{where } p_i = \text{proportion of } S \text{ belonging to class } i$$

→ for collection  $S$  having +ve and -ve examples, entropy is given by,

$$E(S) = -p^+ \log_2 p^+ - p^- \log_2 p^-$$

where  $p^+$  is proportion of +ve examples,  $p^-$  is proportion of -ve examples

→  $E(S) = 0$  if all members belong to same class  
 $E(S) = 1$  if all members are split equally

Information gain for feature column A is calculated as :

$$IG(S, A) = \text{Entropy}(S) - \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} \log_2 \text{Entropy}(S_v)$$

$S_v \rightarrow$  set of rows in S for which feature column A has value v

$|S_v| \rightarrow$  no. of rows in  $S_v$

$|S| \rightarrow$  no. of rows in S

## # C4.5 algorithm

$$\text{Splitinfo}_A(D) = - \sum_{j=1}^k \frac{|D_j|}{|D|} \log_2 \frac{|D_j|}{|D|}$$

$$\text{Gain Ratio}(A) = \frac{\text{Gain}(A)}{\text{Splitinfo}(A)}$$

→ Attribute with maximum gain ratio is selected as splitting attribute.

## # CART Algorithm

If dataset  $D$  contains examples from  $n$  classes, gini index,  $\text{gini}(D)$  is defined as

$$\text{gini}(D) = 1 - \sum_{j=1}^n p_j^2 \quad ; \quad p_j \rightarrow \text{relative frequency of class } j \text{ in } D$$

If dataset  $D$  is split on  $A$  into two subsets  $D_1$  and  $D_2$ , the gini index  $\text{gini}_A(D)$  is defined as

$$\text{gini}_A(D) = \frac{|D_1|}{|D|} \text{gini}(D_1) + \frac{|D_2|}{|D|} \text{gini}(D_2)$$

Reduction in impurity :

$$\Delta \text{gini}(A) = \text{gini}(D) - \text{gini}_A(D)$$

→ Attribute provides smallest  $\text{gini}_{\text{split}}(D)$  (or largest reduction in impurity) is chosen to split the node.