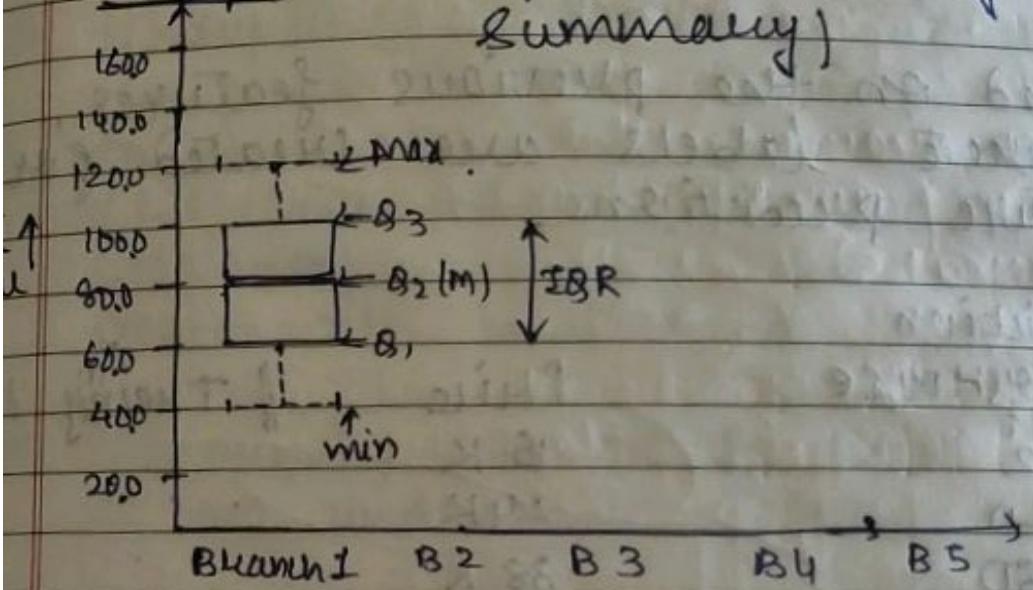


about the data residing at extremes.
Thus, 5 point summary is used.

steps:

- i) observe Q₁ Q₁
 ii) " (2)
 iii) " Median
 iv) " (3)
 v) " Q₃
 vi) " (4)
 minimum minimum
 maximum maximum (1)
 (5)

Box-plot (Representation of 5 number summary)



Classification and Prediction

Decision tree

Bayesian

ANN (Back prop.)

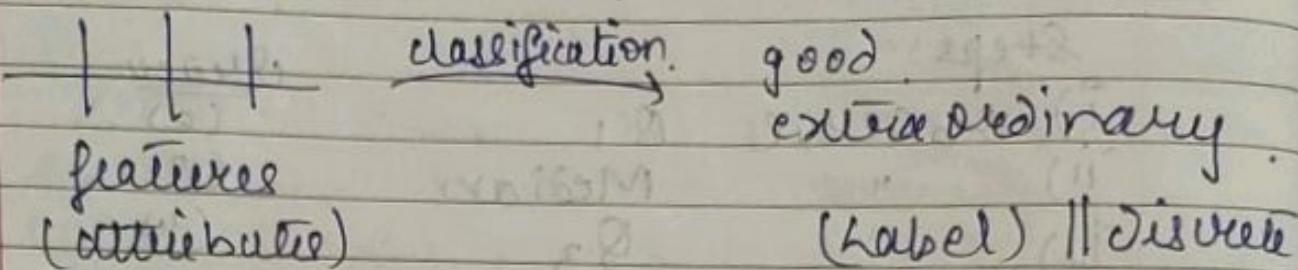
Others

- Fuzzy Set
- Rough Set
- Genetic Algo
- KNN (Clustering)
- Case Base Reasoning

Linear regression

Supervised learning

- class label given based on features



Ex - credit rating to predict whether person should be given loan or not.

Based on the previous features, discrete labels are created for future predictions.

Prediction

Size of house

12

20

50

40

Price

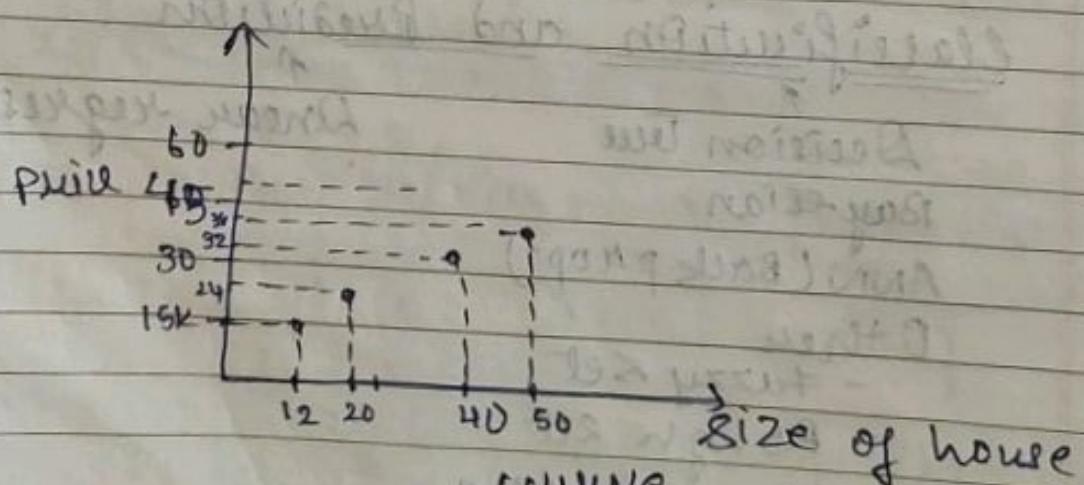
15 K

24 K

38 K

32 K

Training data



* best-fit line is evaluated and further analysis performed.

Prediction deals with continuous values
 Classification deals with discrete values

When dealing with curves, the degree of polynomial (sophistication level) needs to be carefully chosen.

If almost all the points fall on the curve, the model may fail to correctly classify if some outliers fed. (Analogue to memorization).

- ⇒ Unsupervised learning is where only the features are provided. We need to determine the labels as well.
 (No class label provided with feature vector).

Deals with clustering problems.

- ⇒ Determining the no. of clusters in advance is an undecidable problem.

1) Model construction

\uparrow learnt
 Rule, based on training data

Ex - how we calculate the values of β_0 & β_1 is the model for:

$$Y = \beta_0 + \beta_1 X$$

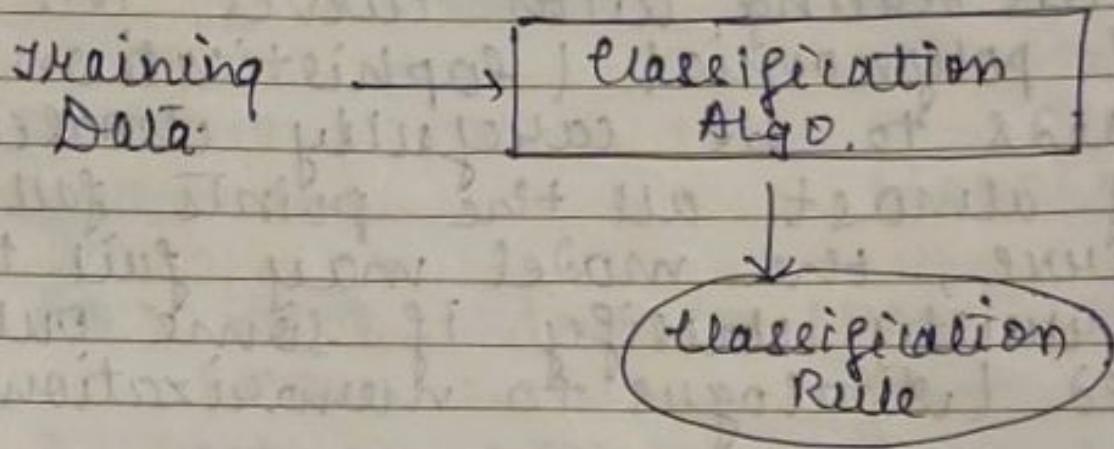
Method (Basic): Hold out

- Divide data into

should be random → Training & testing
 (80% of total) (20%)

hold 20%
 of the data

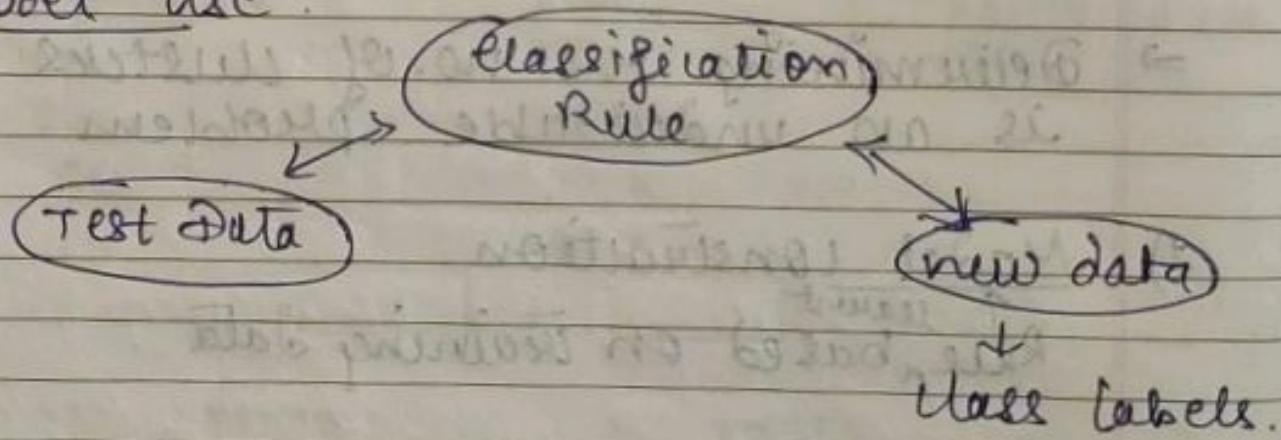
Train your model based on the 80% data & run error analysis on the held out 20% of data.



Model construction

After training the model, the model is 'repaired' based on the errors generated by the built model.

Model use



Hold out performed in this stage.

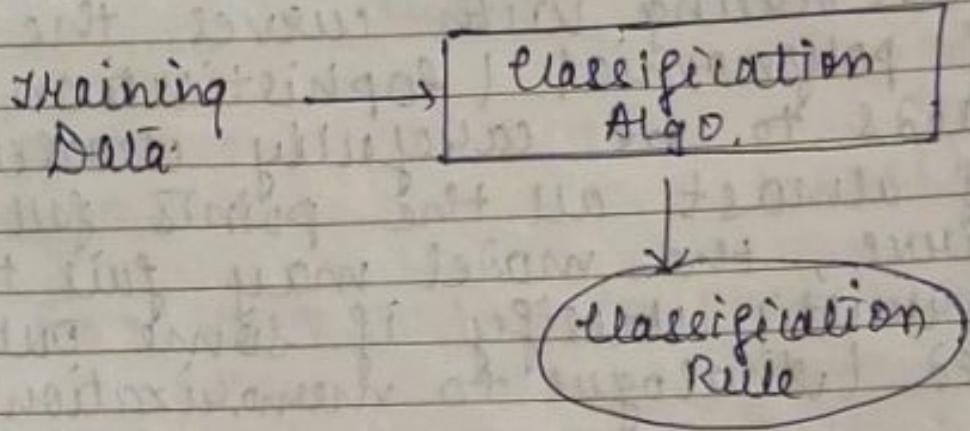
Before classification & prediction

1) Data cleaning

2) Relevance analysis

- Irrelevant data may lead the model to deviate from the intended path

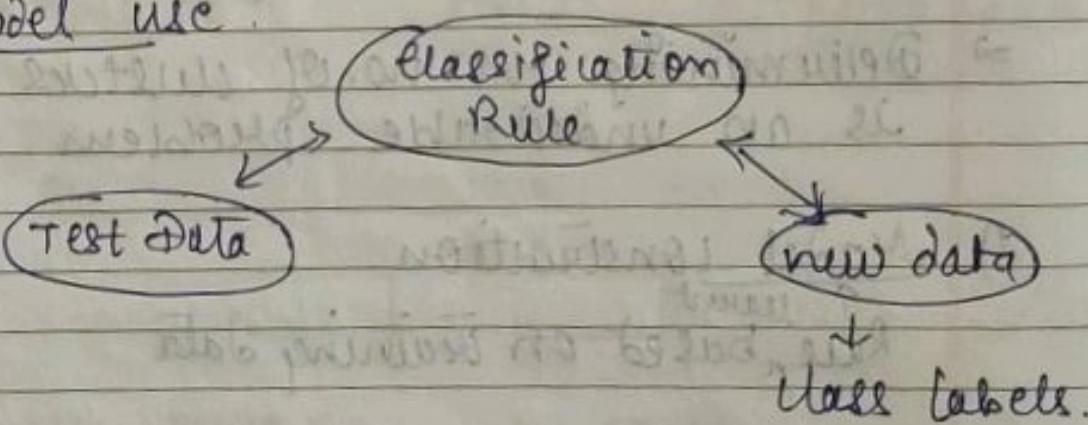
Train your model based on the 80% data & run error analysis on the held out 20% of data.



Model construction

After training the model, the model is 'repaired' based on the error generated by the built model.

2) Model use



Hold out performed in this stage

Before classification & prediction

1) Data cleaning

2) Relevance analysis

- Irrelevant data may lead the model to deviate from the intended path

3) Data Transformation

- especially when dealing with distance based algo. (for classification)

Income Age
 will dominate over other. 20,000 23
 age if normalization not performed. 18,000 24.

Comparing classification methods

Parameters:

i) Prediction accuracy

- Accuracy on the new (testing) data. How the classifier works on data other than the one on which it was trained.

ii) Speed

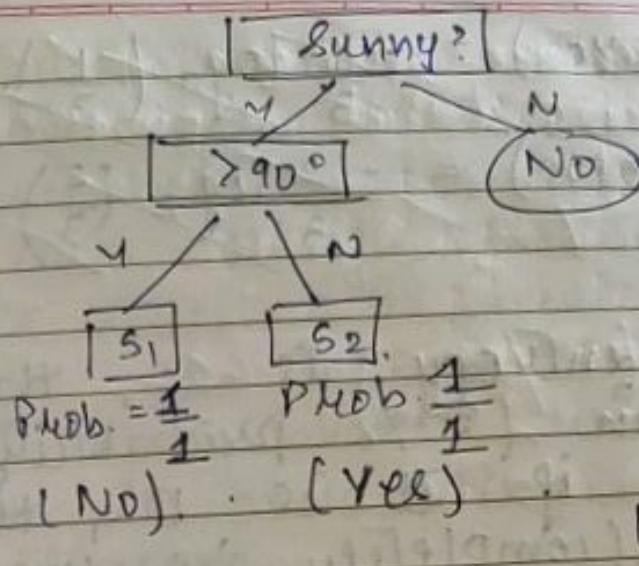
- The computational power of the classifier.

iii) Robustness

- How your classifier handles noisy data (or missing data). In the worked case, the classifier may give ambiguous results.

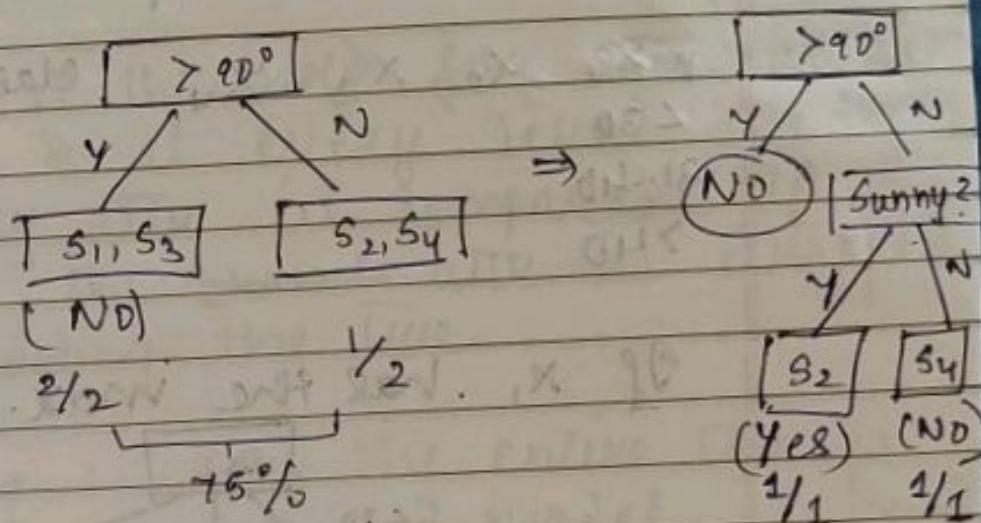
iv) Scalability

- Ability to handle a large set of data. This is interlinked with the speed parameter too (volume, variety, velocity of Big Data)



\therefore all paths give accurate results
(at least on training data)

consequently,



Since, training data here consists of all the possible combinations without any ambiguity, we get 100% results on the training data.

If we add the tuple (N, N, Y) , the model/tree will not give a 100% output or accuracy.

Most relevant attribute is generally considered at the root node.

$$I(S_1, S_2) = -\frac{1}{4} \log_2\left(\frac{1}{4}\right) - \frac{3}{4} \log_2\left(\frac{3}{4}\right)$$

$$I(S=\text{Yes}) = -\frac{1}{3} \log_2 \left(\frac{1}{3}\right)$$

$$I(S=\text{No}) = -\frac{2}{3} \log_2 \left(\frac{2}{3}\right)$$

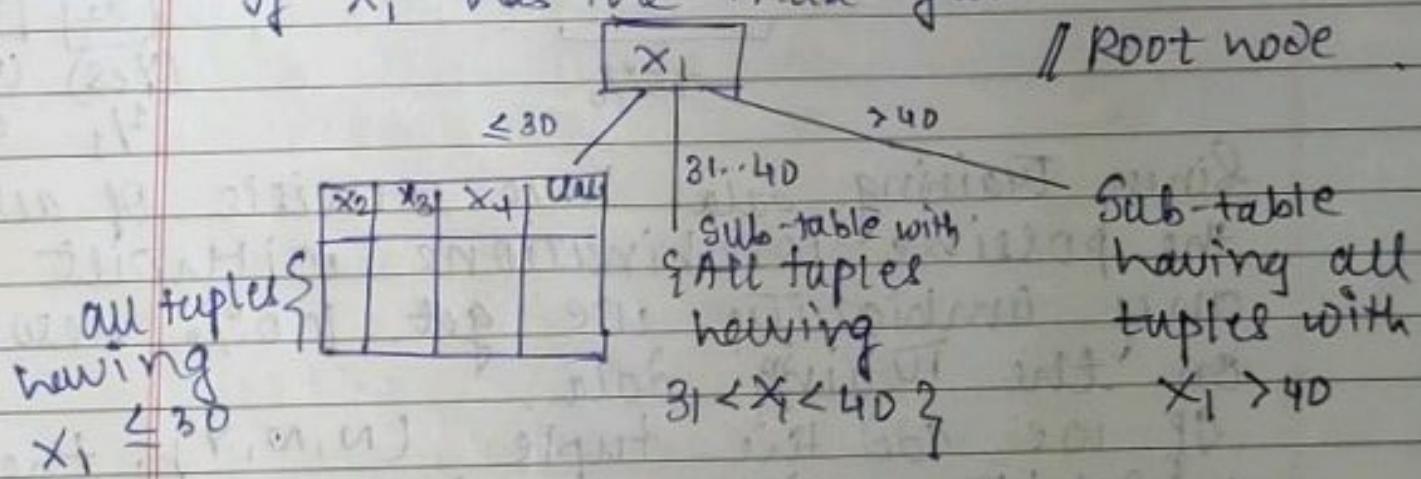
19/9/19

When entropy is 1, the result is uncertain (50% probability)

while if it is 0, the outcome is certain. (completely, i.e. 100% probability)

x_1	x_2	x_3	x_4	Class Label
≤ 30				
$31 \dots 40$				
> 40				

If x_1 has the max. gain:



If all the tuples in the sub-table have the same class label, a leaf node is created with the class label and the process for that branch terminated.

For all other sub-tables, repeat the entire procedure of calculating info. gain. (discard the previous table)

PAGE NO.	/ /
DATE	

→ Limitations:

- Outliers can create a serious problem
- in case of decision tree.
- Not scalable (in-memory processing there, mem. usage limitations)
- If let's say, > 40 branch has only one tuple and based on that we terminate the branch with leaf label 'NO', it is overfitting. May lead to misclassification.

Pre-pruning: A threshold for attr. i.e. set at every level, the info gain of every attr. is compared. If found less, the branch attr. is dropped. Done while building the tree.

Post-pruning: After the entire tree is built and model is found to misbehave, a branch is selected and cut out temporarily. The tree left is tested with the testing data again and if the error is compatible with the error obtained before, the branch is permanently cut (pruned).

$$\text{Bayes Rule: } P(Y|X) = \frac{P(X|Y) \cdot P(Y)}{P(X)}$$

↑ ↑ ↑
 Posteriori likelihood normalization
 prob. factor (prob. of
 obse. & sampler)
 [Evidence]

23/9/19

X - feature vector $\{X = (x_1, x_2, \dots, x_n)\}$
 Y - class label (target class)

Ex: $Y \rightarrow$ Yes
 No

$$\text{if } P(Y/x_{\text{new}}) > P(Y = \text{Yes}/x_{\text{new}}) \text{ - i} \\ \text{new sample} \quad \quad \quad P(Y = \text{No}/x_{\text{new}}) \text{ - ii}$$

If (i) $>$ (ii) then, outcome of $P(Y/x_{\text{new}})$ is predicted as 'Yes'

If (ii) $>$ (i) then outcome \rightarrow 'No'

$$P(Y/x) = \frac{P(x_1, x_2, x_3, \dots, x_n | Y) \cdot P(Y)}{P(x_1, x_2, \dots, x_n)} \\ \downarrow P(x_1 | x_2, x_3, \dots, x_n, Y) \cdot P(x_2 | x_3, x_4, \dots, x_n, Y) \cdot \\ \dots \cdot P(x_{n-1} | x_n, Y) \cdot P(x_n | Y)$$

What makes this theorem 'Naive' Bayes is the assumption that all the attr. (x_1, x_2, \dots, x_n) are independent of each other. Also, all the attr. are equally contributing to the target label.

But all are dependent on target

Thus, the above formula becomes:
 $P(Y/x) = P(x_1 | Y) \cdot P(x_2 | Y) \dots P(x_n | Y) \cdot P(Y)$
 $P(x_1) \cdot P(x_2) \cdot P(x_3) \dots P(x_n)$

OR

$$P(Y|X) = P(Y) \cdot \prod_{i=1}^n P(X_i|Y)$$

$$\prod_{i=1}^n P(X_i)$$

evidence
(constant)
irrespective of
target class

Ex: $x_{\text{new}} = \text{Age} = " < 30 ", \text{income} = \text{"medium"},$
 $\text{Student} = \text{"Yes"}, \text{Credit rating} = \text{"fair"}$
 $y_{\text{new}} = \text{Buy-computer}$

not present in given data. Thus, we need to predict

Now, we don't need to calculate the prob. of age > 30 or $40-50$ etc. because we don't have to predict any value for it.

$\therefore \prod_{i=1}^n P(X_i)$ is going to be same for all values (distinct) of "Buy-computer", we may skip its calculation.

$$\therefore P(Y|X) \propto P(Y) \cdot \prod_{i=1}^n P(X_i|Y)$$

$$y = \text{sigmoid}(P(Y) \cdot \prod_{i=1}^n P(X_i|Y))$$

Here, ^{from} data given,

$$P(\text{Buy-computer} = \text{Yes}) = \frac{9}{14}$$

$$P(\text{Buy-computer} = \text{No}) = \frac{5}{14}$$

Now calculate:

$\frac{2}{9}$	$P(\text{age} = " \leq 30 " \text{Buy-computer} = " \text{Yes" })$
$\frac{3}{15}$	$P(\text{"} \text{Buy-computer} = " \text{No" })$
$\frac{4}{9}$	$P(\text{income} = " \text{medium" } \text{Buy-computer} = " \text{Yes" })$
$\frac{2}{15}$	$P(\text{"} \text{Buy-computer} = " \text{No" })$
$\frac{6}{9}$	$P(\text{student} = " \text{Yes" } \text{Buy-computer} = " \text{Yes" })$
$\frac{1}{15}$	$P(\text{"} \text{Buy-computer} = " \text{No" })$
$\frac{6}{9}$	$P(\text{credit-rating} = " \text{fair" } \text{Buy-computer} = " \text{Yes" })$
$\frac{2}{15}$	$P(\text{"} \text{Buy-computer} = " \text{No" })$

Now,

$$P(\text{Buy-computer} = " \text{Yes" } | X) = \frac{9}{14} \cdot \frac{2}{9} \cdot \frac{4}{9} \cdot \frac{6}{9} \cdot \frac{6}{9} \\ = 0.028$$

$$P(\text{Buy-computer} = " \text{No" } | X) = \frac{5}{14} \cdot \frac{3}{5} \cdot \frac{2}{5} \cdot \frac{1}{5} \cdot \frac{2}{5} \\ = 0.007$$

\therefore Predicted value = Yes.

In case of continuous values:

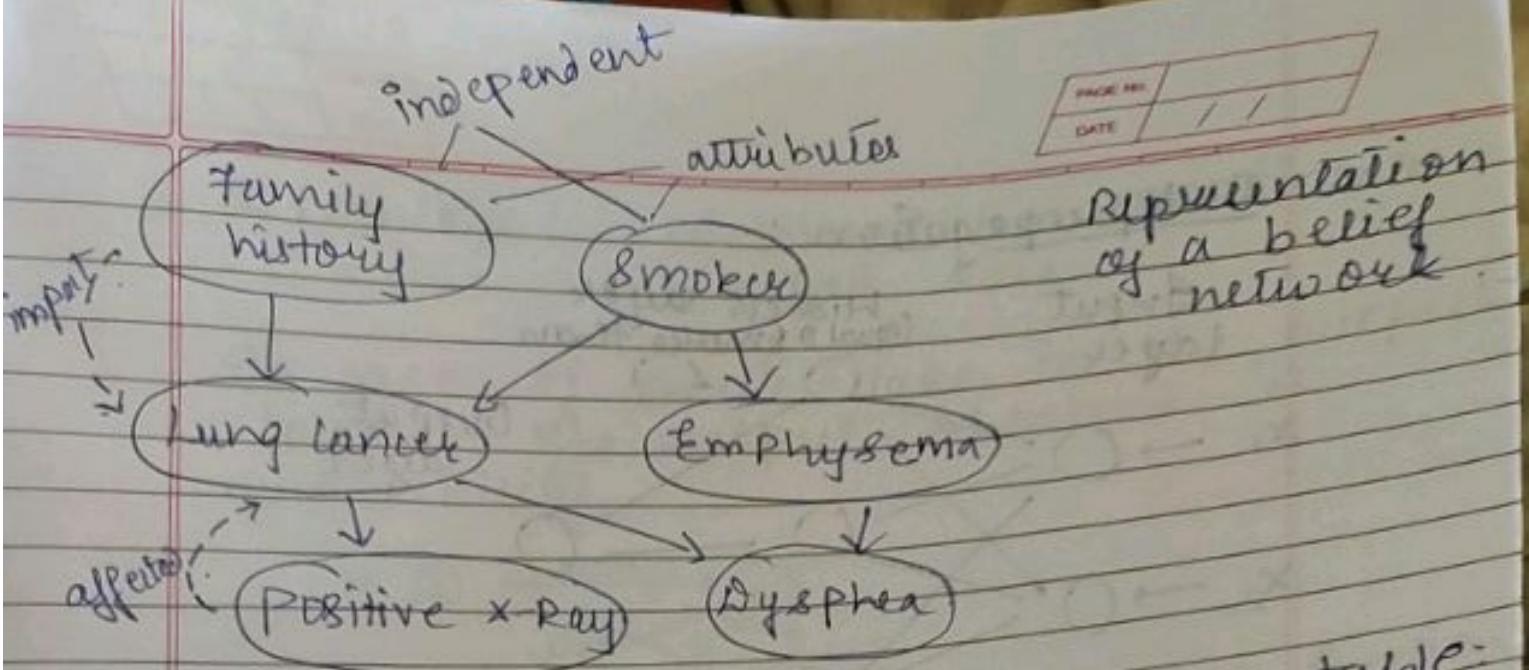
Assume variable as Gaussian

$$P(X_i | Y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} e^{-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}}$$

where, μ_y = mean of values of Y

σ_y = standard deviation

If we remove the assumption of attr being independent, Bayesian Belief Network comes into picture.



Other way is conditional prob. table-

	FH, S	$\bar{F}H, \bar{S}$	$F\bar{H}, S$	$\bar{F}\bar{H}, \bar{S}$
(lung cancer) LC	0.8	0.5	0.7	0.1
$\bar{L}C$	0.2	0.5	0.3	0.9

$$P(z_1, z_2, \dots, z_n) = \prod_{i=1}^n P(z_i | \text{Parent}(z_i))$$

26/9/19 Q. A murder takes place
Witness 1: 90% sure that X is murderer
Total 20 people ~~were~~ present at the place

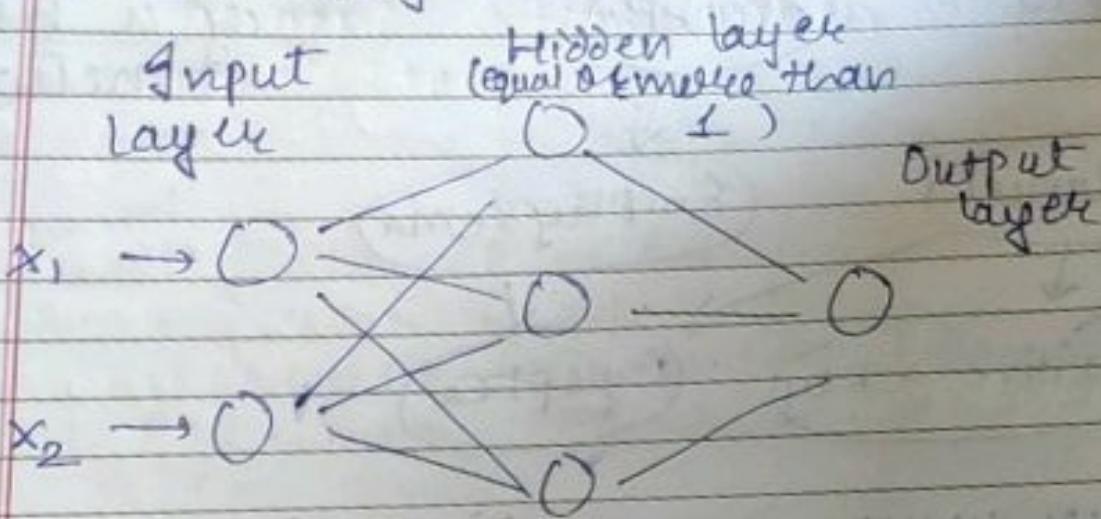
$$\rightarrow P(X = \text{guilty}) = \frac{1}{20} = 0.05$$

$$P(Y = \text{eyewitness-identifies} | X = \text{guilty}) = 0.9$$

$$P(Y = \text{eyewitness-identifies} | X = \text{not-guilty}) = 0.1$$

$$\begin{aligned} \therefore P(X = \text{guilty} | Y = \text{eye-witness-identifies}) &= \frac{0.9 \times 0.05}{(0.9 \times 0.05) + (0.1 \times 0.95)} \\ &= 0.3213 \\ &= 32\% \end{aligned}$$

Back propagation:



$$x = (x_1, x_2 \dots x_n)$$

x_1 : height x_2 : weight
 y : Adult or child.

Every like link is a weighted edge to the next node/neuron. The weight determines the attribute importance/relevance towards the particular node.

Neural network \approx function approximator

We go for NN only if we have a large amount of training data.

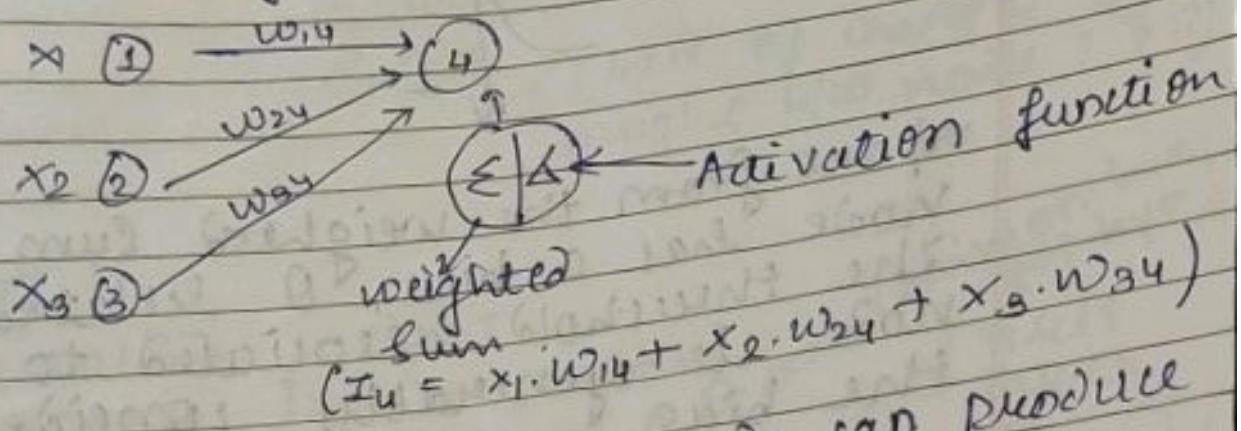
Disadvantage: Interpretability, because it involves learning weight matrix which may not always be straightforward.

Output layer - 1 node \Rightarrow Yes/No type of O/P.

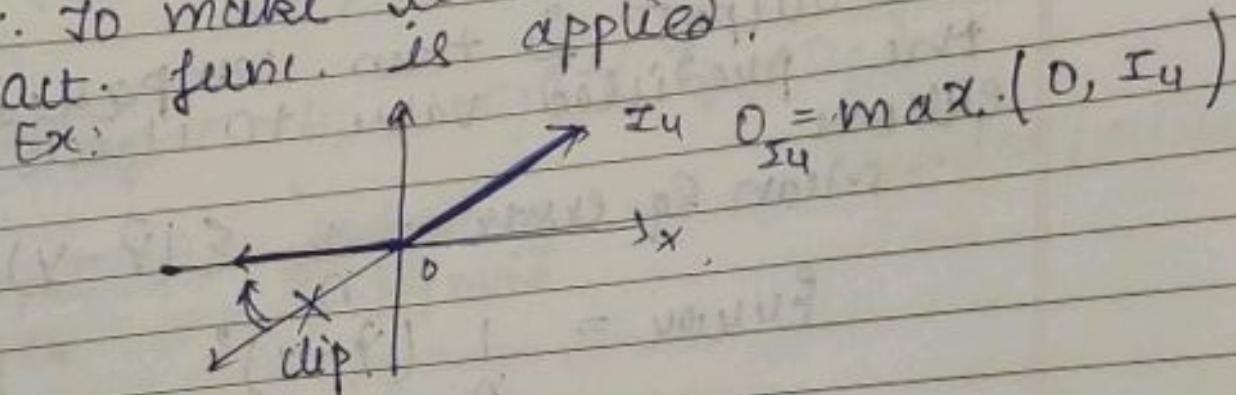
PAGE NO. / / /
DATE / / /

Advantage : Accommodates non-linear shaped data.

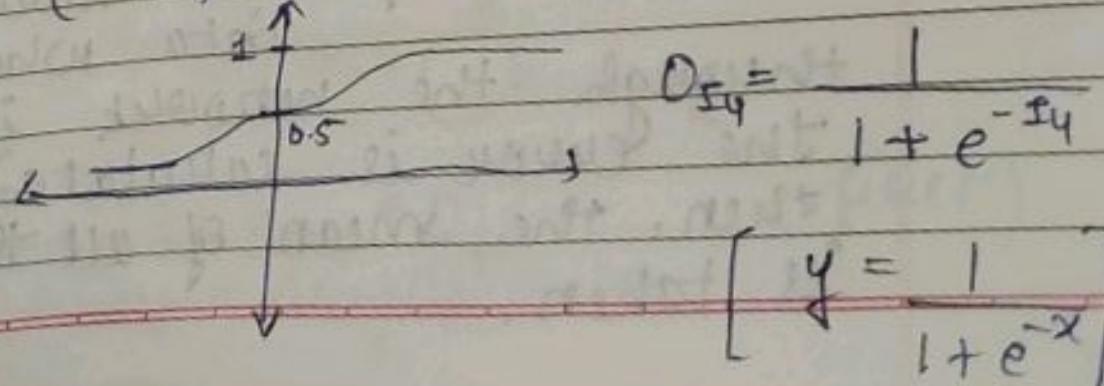
In case of multi-class in the output layer, the output obtained is probability.



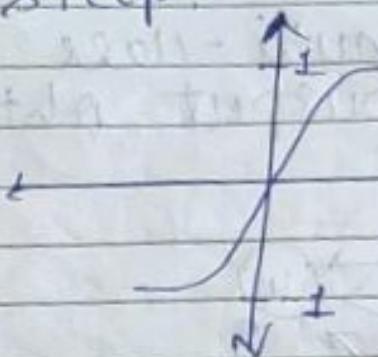
without Act. func., the node can produce a result between $(-\infty, \infty)$.
 \therefore To make it more interpretable, act. func. is applied.



Other ex. is sigmoid func.
 $(-\infty, \infty) \rightarrow [0, 1]$.



Yet another func. is hyperbolic tangent. Used when derivation is very steep.



apart from the weighted sum, every node has a bias, θ related to it. The threshold associated to every node is generally considered as the bias.

$$\text{ex: } w_{11} \cdot x_1 + w_{22} \cdot x_2 - T > 0$$

The output is then compared with the predicted value to check for errors.

$$\text{Mean Sq. error} = \frac{1}{n} \sum (\hat{y} - y)^2$$

$$\text{Error} = \frac{1}{n} (\hat{y} - y)^2$$

Epoch updating - The entire training data when iterated through the network is 1 epoch. The error is calculated. Then, the mean of all the epoch's is taken.

PAGE NO. / / /
DATE / / /

Back updating - Error is calculated passed through the network and the attributes adjusted.

Weights are adjusted based on:

$$\frac{\partial \text{Error}}{\partial w_{ij}} \quad (\text{Rate of change of error})$$

$$\hookrightarrow (\text{Rate of change of weight b/w node } i \text{ & } j)$$

Steps:

1) Initialize weight

2) Propagate input forward

$$I_j = \sum_i w_{ij} \cdot O_i + D_i$$

common bias
for layer

[first
layer]
 $O_j = x_j$
 $O/P = I/P$

$$3) O_j = \frac{1}{1 + e^{-I_j}}$$

true O/P
Actual O/P.

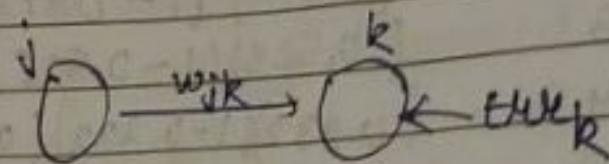
$$4) \text{Error}_j = O_j(1-O_j)(T_j - O_j)$$

Predicted / Actual

for Output layer:
derivative
of Sigmoid

5) For hidden layer:

$$\text{Error}_j = O_j(1-O_j) \sum_k E_{hk} \cdot w_{jk}$$



$$6) \Delta w_{ij} = (l) \text{Error}_j * O_i \quad (\text{first})$$

$$w_{ij}^{new} = w_{ij} + \Delta w_{ij} O_i = x_i$$

learning rate (η /a) parameter

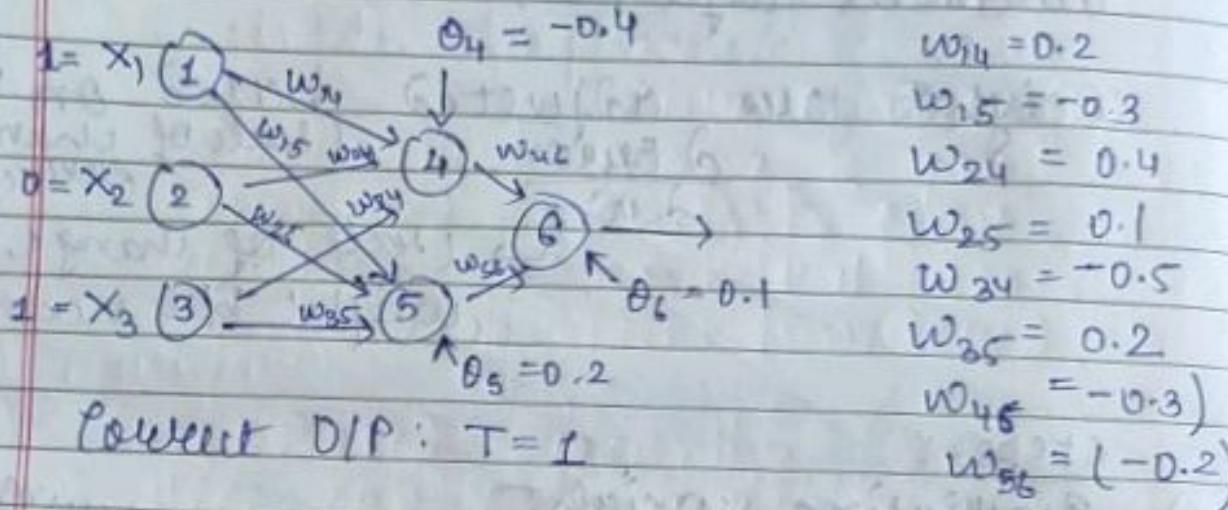
Page No.	
DATE	///

$$\Delta \theta_j = (\eta) E_{\text{NLL}} j$$

$$\theta_j = \theta_j + \Delta \theta_j$$

80 | 9 | 19

Ex:



Current D/P: $T = 1$

unit / j	Net I/P, I_j	Output, O_j
4	$[0.2 + 0 + (-0.5)] + (-0.4)$ $= -0.7$	$\frac{1}{(1 + e^{(-0.7)})} = 0.33$
5	0.1	0.525
6	$[(0.332)(-0.7) + (-0.2)(0.525)] + 0.1 = -0.105$	0.474

$$I_j = \sum_i w_{ij} D_i + \theta_j$$

There is a huge difference between obtained D/P & expected output. Thus we adjust the values by back propagation.

	E_{NLL}	$E_{\text{NLL}} j$	$\times T_j \times \theta_j$
6	$(0.474)(1 - 0.474)$	$(1 - 0.474)$	$= 0.1311$
5	$(0.525)(1 - 0.525)$	$(0.1311 * (-0.2))$	$= -0.065$
4	$(0.332)(1 - 0.332)$	$(0.1311 * (-0.3))$	$= -0.0087$

Now, update the weights & biases:

Given, learning rate, $\lambda = 0.9$

$$\text{weight \& bias} \quad \text{New value}$$

$$w_{46} \quad (-0.3) + [(0.9)(0.1311 * 0.932)] \\ = -0.261$$

$$w_{56} \quad (-0.2) + [(0.9)(0.1311 * 0.525)] = -0.138$$

$w_{14}, w_{24}, w_{34}, w_5, w_6, w_{35}$

$$o_6 \quad (0.1) + [(0.9)(0.1311)] = 0.218$$

o_5

o_4

- Neural network is not interpretable.
- Disadvantage: Representation is not easy.

→ Drop-out (ignoring or removing some neurons) is used to increase the accuracy of neural network. It also avoids over-fitting.

Other Classification methods: (Practically not adaptable)

1) K-Nearest Neighbour

Euclidean distance is used to find the k' nearest nodes to a point and are classified as one.

Type of lazy learning. (Until a new point is not added, the model does not learn about any new knowledge)

While, decision tree is an eager

deanne classifier as it can build
its model based on training
data itself. It does not need to
wait for testing data.

If it is a discrete classifier, the
max. value among the points of
considered will be taken if it is
deployed as a dual classifier,
the average is considered.

2) Case Based Reasoning

Use CAR (complete symbolic
descriptions).

For any new input, it searches
for the exact same case in the
previous gained knowledge. If found,
that exact output is given.

If not then, the most identical
case is considered.

Ultimately the output depends on
the importance or weightage of
an attribute.

If cases are represented as graphs,
the new sample graph should
resemble a subgraph.

3) Genetic Algorithm

Based on neural re-evaluation.

- Rules are represented as strings of bits

- Initially generate population randomly

Attr.

$A_1 \& A_2$

Target class

$C_1 \& C_2$

Rule: If $\frac{A_1}{1} \& \frac{\text{not } A_2}{0}$ then $\frac{C_2}{0}$

if it \leftarrow every attr. is represented by 0/1
not only 2 values Target class: $C_1: 1$
ex-M/F $C_2: 0$.

Rule: If $\frac{\text{not } A_1}{0} \text{ and not } A_2$ then $\frac{C_1}{1}$

If attr. has k values, every attr. is
represented by k bits.

- Fitness of rule / crossover and mutation
- New population is generated on the basis of fitness of rule.

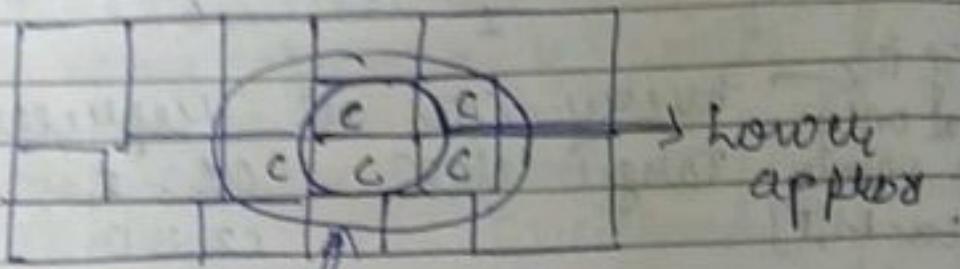
4) Rough set approach

- Useful in handling imprecise and noisy data
 - Discrete values
- ↓
- If continuous values \rightarrow discretization required.

- Based on equivalence class
- Data tuples in one class should be indistinguishable w.r.t. certain attr.
- In real world, however, complete

distinguishability may not be possible
therefore, approximation may be used

- ↳ Lower approx.
- ↳ Upper approx.



Upper approx.

(cannot guarantee whether tuple does not belong to the class)

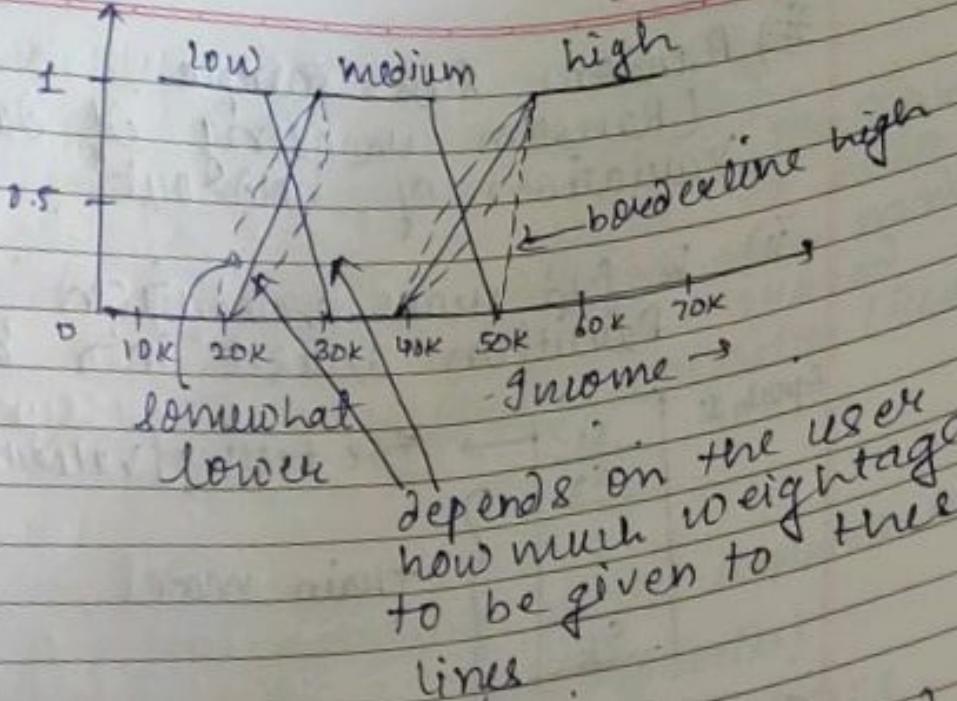
) Fuzzy Set approach

- Rule based system have a sharp cut-off.
- To avoid it, some fuzziness is introduced.

If $\text{year completed} \geq 2 \wedge (\text{income} \geq 50k)$
Then credit = "Approved"

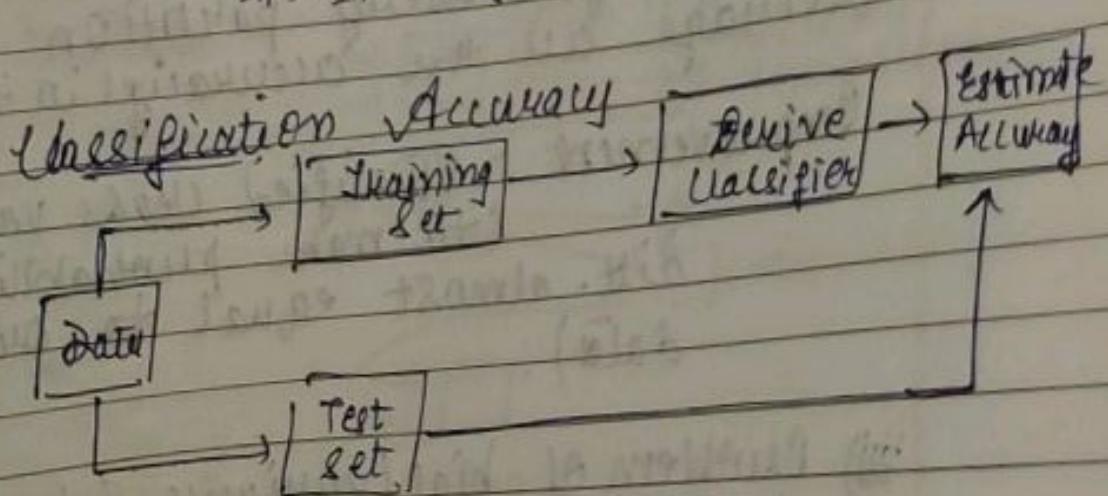
Here if income = 49999 then also his/her credit is not approved which is very harsh.

Thus, instead of having 0/1 discrete values, we consider all the dual values between 0 & 1.



Prediction

- 1) Linear regression (example solved ahead)
- 2) Non-linear
 $y = \alpha_1 x + \alpha_2 x^2 + \dots$ (based on shape)
 x_1, x_2, \dots (if more dimensions (≥ 3))



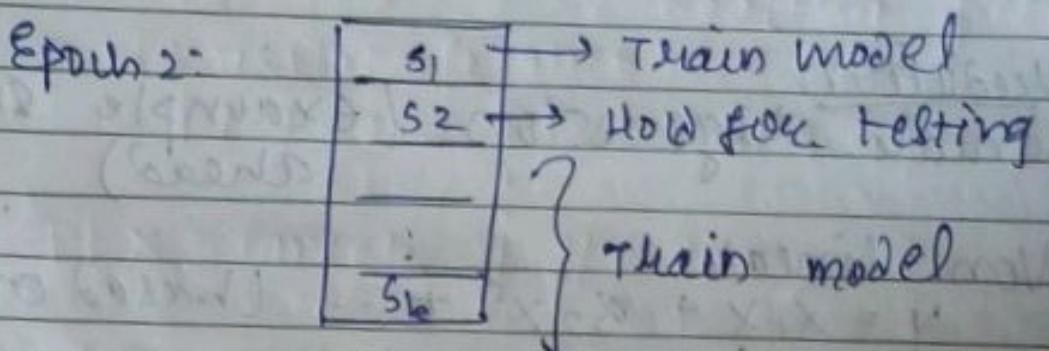
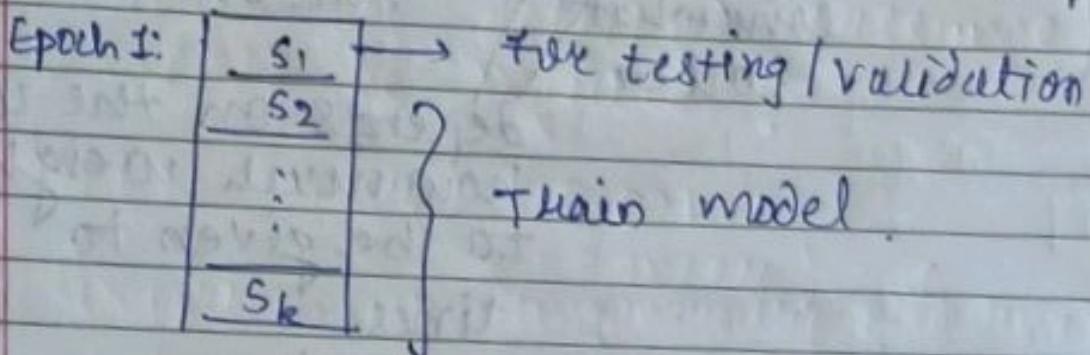
How to estimate improve accuracy?

- i) Hold out
 - Train - Build classifier
 - Test - Test accuracy on randomly

iii) Random sub-sample
(Random choosing of data)
variation of hold-out

ii) k-fold cross validation

- Partition dataset into k partitions



Repeated for every partition
Average all the accuracies in the end

Improvement - Stratified cross validation
(to make probability
dist. almost equal to original
data).

iv) Problem of bias/variance solved upto
some extent.

Increase classification accuracy:
⇒ Pruning is a common solution

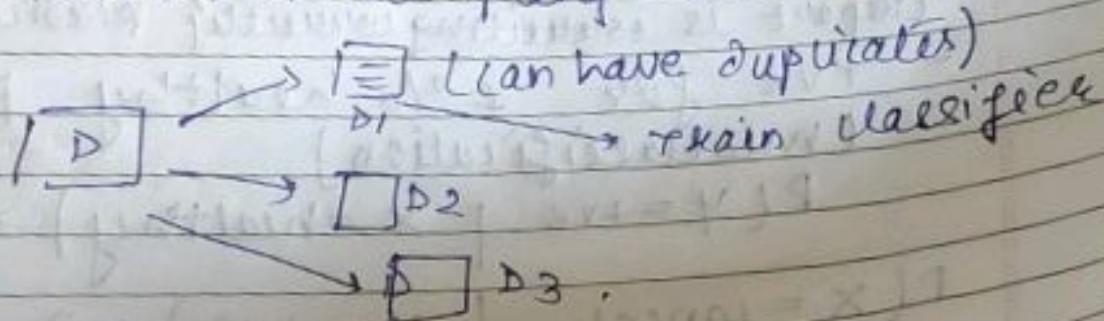
Other methods:-

Bagging (Bootstrap Aggregator)

- Train your dataset onto multiple classifiers and a model is built.
- Testing data is fed to each model.
- After this, voting is carried out and a majority opinion (result) is taken.

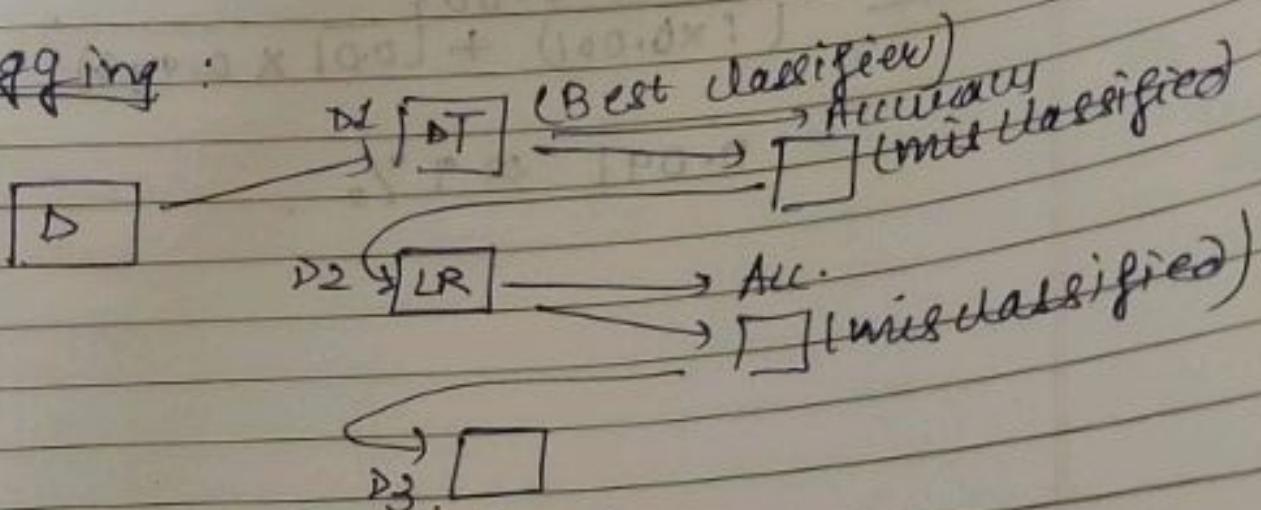
Boosting

- Random Sampling with replacement



Late, voting and majority

Bagging:



(weighted voting)

$$\Rightarrow \text{Sensitivity} = \frac{\text{total - positive}}{\text{positive}} \rightarrow \begin{array}{l} \text{identified} \\ \text{(actual)} \end{array}$$

Page No.	
Date	/ /

$$\text{Specificity} = \frac{t-\text{neg}}{\text{neg}}$$

$$\text{Precision} = \frac{t-\text{pos}}{(t-\text{pos}) + f-\text{pos})}$$

⇒ Cancer is detected (assume)
0.1% is infected.

$$P(X = \text{cancer}) = 0.001$$

$$P(Y = +ve | X = \text{cancer}) = 1$$

[Report is correctly predicting]

- +ve for 1% healthy people
(misclassification)

$$P(Y = +ve | X = \text{healthy}) = 0.01$$

$$P(X = \text{cancer} | Y = +ve) = ?$$

$$\text{Ans. } \frac{1 \times 0.001}{(1 \times 0.001) + (0.01 \times 0.999)}$$

$$= 0.091 \approx 9\%$$