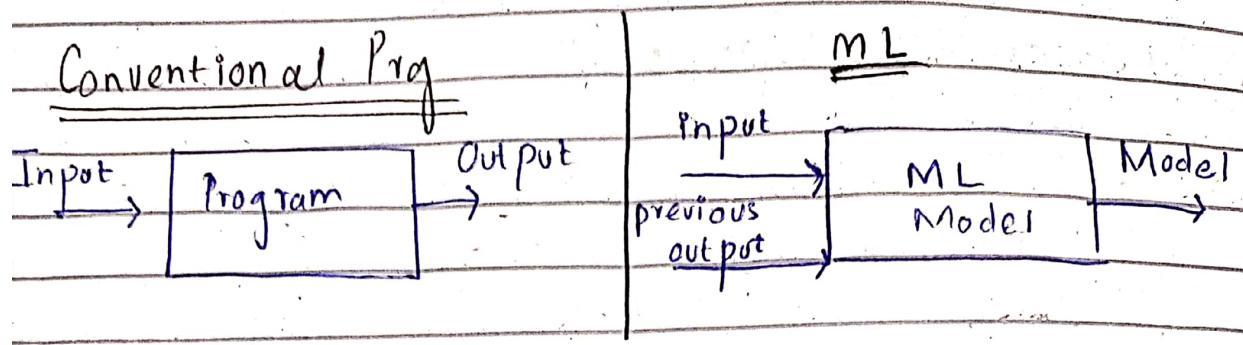


To transfer learning capability to machine without explicit program.

- Unseen data will be process to classify
- hard coded to replace kr dia.



Categories of ML

- 1) Supervised Learning
- 2) UnSupervised Learning
- 3) Reinforcement Learning

1. Supervised Learning:

class labels are available, you will train your model on the bases of previous data and predict the output of future data.

- Training data has class label
- We want to predict that label for future instance

Classification
Supervised Learning

Regression / Continuous value prediction

- Regression ~~Explain~~
Kita ha like predict the
price of a house, car

UnSupervised Learning

- class Label are not available
So we group up things that are
similar

- Class Label nahi hoty but feature ki base
py group kry like fruit ko color,
shape ki feature py check
kry.

- jitny zyda feature hain utna zyda
bhtr grouping.

Reinforcement :

Supervised py train kia or
test test kia ga it answers
aya humny us ko bataya galat
ha us ny apny knowledge base mn
us ko add kr lia, Learn ko lia us data ko

- Learning with experience
- Learning with passage of time.

Optimization : maximum benefit

Applications:

- To check whether loan should be awarded or not on the bases of credit card score.
 - Real Time Stock Value Prediction
 - Spam OR NOT SPAM / Fake News OR Real
 - Computer Vision to identify fertile Land and unfertile Land
 - Like Self Driving Car
- } Reinforcement
- } Supervised
- } Unsupervised

Decision Tree

Inductive

- Inductive means recursive.
- derive classes from I31.1
- ID3 is greedy (badi ki taraf jata ha)
but guarantee ni ha.

Concept Learning System

jiska name yahan skaien (e.g. Car, Animal etc) agar us k sub concept.

Algo

~~Ex 1)~~

Decision Tree is attribute-based description

- phly ek universal object ly lo
(Animal etc)
- then sub-concept is defined

Restaurant Problem

Should we wait for a table or not

~~Input Attribute~~

~~Alternate restaurant, etc~~

~~Output Will Wait~~

How to Construct Decision Tree

Task

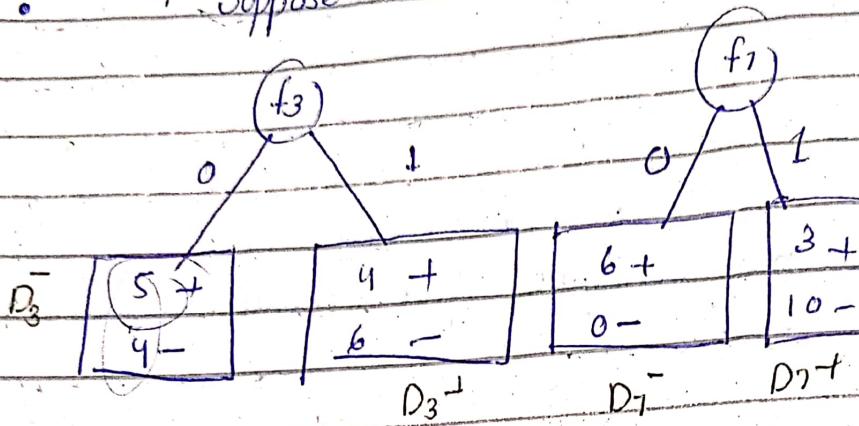
Measure of Node Impurity

Entropy

Total 19 instances.

9 pos & 10 neg

Suppose it has 7 features



$p :=$ proportion of +ive

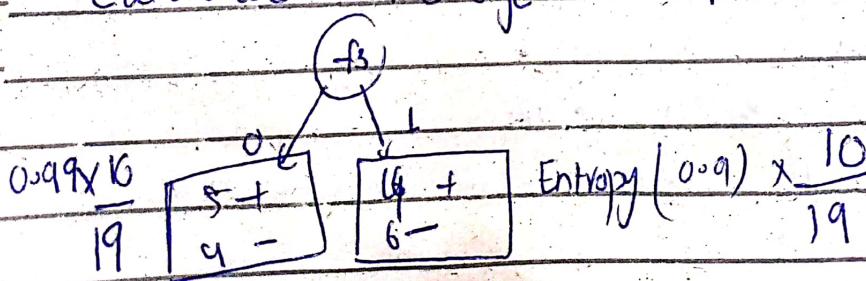
$$H = -p \log_2 p (1-p) \log_2 (1-p)$$

$$\log_2 0 = 0 \quad \log_2 1 = \text{undefined}$$

$$D_3^- = -\frac{5}{9} \log_2 \frac{5}{9} - \left(\frac{4}{9}\right) \log_2 \left(\frac{4}{9}\right) = 0.6551 - 0.37$$

$$D_3^+ = -\frac{4}{10} \log_2 \frac{4}{10} - \left(\frac{6}{10}\right) \log_2 \left(\frac{6}{10}\right)$$

Calculate Average Entropy



$$\frac{9}{19} \times 0.99 + \frac{10}{19} \times 0.97$$

$$= AE$$

jiski AE woh better split hogा

Class Entropy.

info. Gain ko maxi

⇒ Stopping ?

1) Stop if entropy is below some threshold.

2) Number of element in data set is below threshold

3) Any split is entropy decrease nہیں hogے

underfit kisi ko classify he nہیں کر skya

Example

: Step 1 : Class Level ka +1 (Entropy)

$$\text{Total} = 6$$

$$+3, -3$$

Class Level Entropy :

$$H = -P \log_2(P) - (1-P) \log_2(1-P)$$

$$= -\frac{3}{6} \log_2\left(\frac{3}{6}\right) - \left(\frac{3}{6}\right) \log_2\left(\frac{3}{6}\right)$$

f_1 do same for all

f_1	0	1
2+		2+
1-		1-

Calculate AE

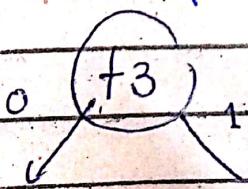
$$AE_1 = 0.92$$

$$AE_2 = 0.92$$

$$AE_3 = 0.8$$

$$AE_4 = 1$$

Choose f_3 & split



0 K instance

1 K Samy

instance

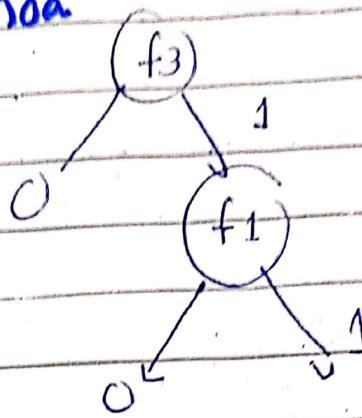
ap f₃ ko chor dia f₁, f₂, f₄ ko mila
kore

$$AE_1 = 0.58$$

$$AE_2 = 0.55$$

$$AE_3 = 0.95$$

F₁ Choose hoa



Stopping Criteria ajaye toh stop kr do
ni toh batir ko f₂, f₄
split kry jao

Pruning

Pre Pruning jaisy stopping criteria
ka zikr kia.

Post Pruning para kr k phir prune
kr lo

Leaf Node Ka content o

What is over fitting

performing good ↑ on some training
data but

It can't be generalized so won't work
in real env.

Pruning is the technique to avoid
over fitting

ID3

Pruning:

Statistical method is used to prune the
least reliable

Two common stra

- Pre Prune (by halting)
- Statistical Significance (between attribute significant)

Sometime practitioner combine both tech

pre + post (hybrid) approach

e.g lagai hoi pre ha but agr error

ajaye toh post prune kar do

simple

Post Pruning is divided into

1) Training Sample

2) 10 → 20% for valid (Training Test)

3) 5 → 10% for testing

Regression tree

- target column could be a number
- Use standard deviation

$$X = 25, 30, 46, 45, 52, 23, 43, 35, 38, \\ 46, 48, 52, 44, 30$$

$$\text{Average} = 39.78$$

Standard deviation

$$\sqrt{\frac{(25 - 39.78)^2 + (30 - 39.78)^2 + \dots}{N}}$$

$$\underline{\text{Ex 2}} \quad 25, 30, 35, 38, 48$$

$$\text{Avg} = \frac{25 + 30 + 35 + 38 + 48}{5} \\ = 35.2$$

$$\text{Stand} = \sqrt{\frac{(25 - 35.2)^2 + (30 - 35.2)^2 + (35 - 35.2)^2 + (38 - 35.2)^2 + (48 - 35.2)^2}{5}}$$

$$= \sqrt{(10.2)^2 + (-5.2)^2 + (0.2)^2 + (2.8)^2 + (12.8)^2} / \sqrt{5}$$

$$\sqrt{104.04 + 27.04 + 0.04 + 7.84 + 163.84} = 15$$

$$\sqrt{302.8} = \sqrt{60.56}$$

sunny = 7.78

overcast = 3.49

rainy = 10.87

Summarized

	Std	Instan.
Overcast	3.49	4
Rain	10.87	5
Sunny	7.78	5
		14

Weighted Standard

$$\frac{4}{14} \times 3.49 + \frac{5}{14} \times 10.87 + \frac{5}{14} \times 7.78$$

$$= 7.68$$

Naive Bayes

spam classification Example

First we built a histogram

$$P(\text{Dear} \mid \text{Normal}) = \frac{8}{17}$$

8 times dear occurred in normal message.

Class Prior Probability

$$P(S) \times P(\text{Dear} \mid S) \times P(\text{Friend} \mid N)$$

KNN

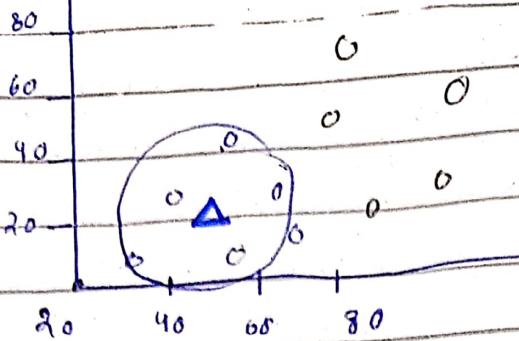
- Supervised Learning
- non-parametric classification
- lazy learning algo:
phly sy train ni krtai,
jab query ky jaye tab train krtai.
- active learning
already prepared

KNN is Lazy

Example:

- Sensitive to imbalanced class.

bacterial



near by mn \in viral hair
to be viral classify

Step 1 Measure the distance, for the
new data point that's why it
is considered Lazy.

Step 2

Step 3 Sort distance

Step 4 Identify closest depend on k
determine majority class

k shows the number of closest
point to select.

Leave One Out Cross Validation

How good is the classifier?

= 1 method is Leave One out Cross Validation

- One by one check algo kis row ko
algo kya predict kita phir us
ki actual or prediction match
kr lo.

- Each tuple is checked across all training data

$$\text{Accuracy} = \frac{\text{Correct Instances}}{\text{Total instances}} = 0.833$$

Total instance

\$ Training test -

Problem of leave one out
 when we have 1000+ datapoints
 than calculating distance
 is time consuming.

Solution

Split the data and into 2 parts

- Large training data
- Small testing data

80, 90% training

20, 10% testing

Problem with imbalanced Dataset

agr strength bhar jaye toh
 probability bhar jaye ge

How to find optimal value of k

- Try k even na ho-
- Should not be too high & low
- k should be odd and > 1

$$k = \sqrt{n}$$

More than 2 groups standardization

- jiskai gap kam ho uska impact o ho jay
- highest gap will influence more

perform standardization (assign weight)

$$Z_i = \frac{x_i - \bar{x}}{s_p}$$

equal impact & result will be less biased

Advantage

simple & easy

Disadvantage

- all training data is used
- sensitive to imbalanced data