

Get feature for the attributes which has more value.

Decision - Broaden tree at base level.

Alternate measure is Gain Ratio.

$$\text{Split info} = \sum_{i=1}^c \frac{|S_i|}{|S|} \log \frac{|S_i|}{|S|}$$

Split information

where  $S$  to  $S_c$  are 'c' subsets of examples resulting from partitioning 's' by the 'c' valued attribute  $A$ .

$$\text{Gain Ratio}(S, A) = \frac{\text{Gain}(S, A)}{\text{Split Info}(S, A)}$$

08/01/2024

) Handle the training examples with missing attribute

$$\text{Gain}(S, A) = \text{G}(x) / \text{Humidity} / \text{Wind} / \text{Dust} / \text{Temp} / \text{Example}(x, C(x)) > A(x) > \text{Null}$$

3 strategies

- 1) whichever value is max in that attribute then fill it with that max value irrespective of  $C(x)$
- 2) that instance if it for old or new check the target concept whichever is true take old -ve then fill it correspondingly acc to  $c(x)$  value eg - 1 (Noise)

Handling attributes with diff. costs  
Some gives +ve and associated with cost/weight.  
eg - To predict medical disease

$$\text{Temp} | \text{Pulse} | \text{Blood test} | \text{Biopsy} | \text{Tongue} \\ \text{Gain}(S, A) \\ \text{Cost}(A)$$

Based on ent. value.

Avoiding Overfitting in the Data.

more errors in train data  $\Rightarrow$  less errors in test data

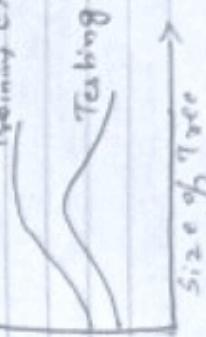
when we apply hypothesis on training sample it shows less error  
" " " some " "  
it shows more errors because it has learnt more than cut is required.

Given an hypothesis  $h$  said to overfit the training examples if the hypothesis  $h$  is performing better on the training set than over the testing set. So Training error(h)  $<$  Testing error(h).  
Cross len Accuracy more

why is this overfitting causing in ID3

$\rightarrow$  Noise.  
 $\rightarrow$  If size of above set is small they are not in a position to test.  
eg - 1 (Noise)  
outlook=sunny, Temp=hot, Humidity=Normal,  
wind=snowy, Play Tennis = No

(for the Tree) Training Example



97

## Avoiding overfitting in the Data.

Strategy:-

- ① ↳ Stop ~~the~~ growing the tree after certain level
- ② ↳ Grow the tree completely, allow overfitting to occur, then we go for pruning the tree

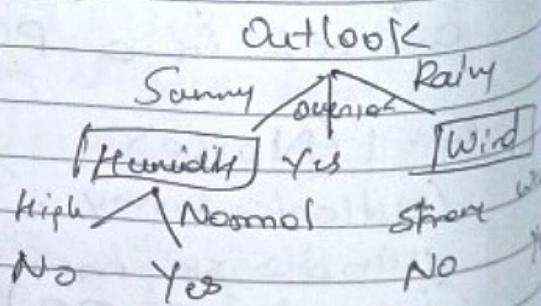
① Disadv: It is not easy to determine when we stop growing the tree. (practical difficult)

② It is mostly used strategy.

- ② ↳ i) reduced error pruning.  
ii) Rule-post pruning.

; In the tree that is already built, every Decision node becomes candidate to prune, pruning is made until we get high accuracy or while pruning accuracy is checked if it is high it is chopped off and replaced with deaf node.

Humidity & wind



i) It considers each of the decision nodes in the tree to be candidate for pruning. Pruning a decision node consists of removing the subtree rooted at that node, making it a leaf node and assigning it to the most ~~complication~~ common

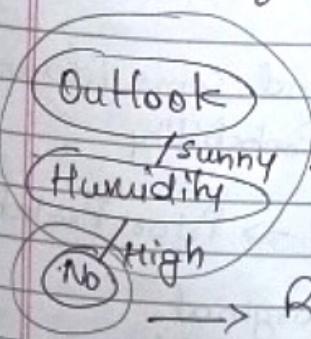
Nodes are removed only if the resulting pruned tree performs better than the original tree for

Pruning of tree continues until further pruning is beneficial. The pruning is stopped once the accuracy stops starts decreasing.

ii) Rule - post - pruning.

After the tree is built it is converted into set of rules. For every leaf node there is a rule. Rule is in the form of path from the root node to leaf node.

Example: if ( $\text{outlook} = \text{sunny}$ )  $\wedge$  ( $\text{Humidity} = \text{high}$ ) then ( $\text{play tennis} = \text{No}$ )



This path is known as  
Rule Antecedent (pre-condition)

Rule Consequent (post-condition)

- Only pre-condition is pruned.

- \* Build a Decision Tree from the Training set. Grow the tree completely until training data is fit to allow the overfitting to occur.
  - \* Convert the tree into set of rules by creating a rule for each path from the root node to leaf node.
  - \* Prune each rule by removing any tree conditions so that the accuracy is improved.
  - \* We generate one rule for each leaf node.
  - \* The path from the root to leaf becomes rule Antecedent or pre-condition.
  - \* A classification at Leaf node becomes the rule consequent or the post-condition.
- ↓
- \* Each rule is pruned by removing the Antecedent or the precondition.
  - \* This process is continued until there is improvement in the estimated accuracy.
  - \* No pruning is performed if there is no improvement in the accuracy.

16/01/2024 3<sup>rd</sup> Module.

Bayesian learning.

Computing probabilities to find some inference prediction

Bayes Theorem

Posterior  
Probability

Prior Probability

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

Likelihood

Prior Probability

Marginal  
Probability.

Ex:- What is the probability of finding the king given that it is a facecard in a pack of card.

$$\text{facecards} = J \& K \quad 4 \times 3 = 12$$

$$P(\text{King} | \text{Facecard}) = \frac{P(\text{Face}) \cdot P(\text{King})}{P(\text{Face})}$$

$$= \frac{4/12 \times 4/52}{12/52} = \frac{1 \times \frac{4}{12}}{\frac{52}{12}} = \frac{4}{52} = \frac{1}{13}$$

$$P(K/F) = \frac{1}{13}$$

- Bayesian learning provides a probabilistic approach for inferences.
- Bayesian learning Algorithms calculate explicit probabilities for hypothesis, and are among the most practical approaches for certain types of learning problems.

#### Features of Bayesian learning

- \* Each training example can incrementally increase or decrease the estimated probability that the hypothesis is correct. This provides more flexibility than algorithms that completely eliminate hypothesis if it is found to be inconsistent.
- \* Prior information can be combined with observed data to determine the final probability of the hypothesis.
- \* Bayesian method can accommodate hypothesis that makes probabilistic predictions.

→ New instances can be classified by combining predictions of multiple hypothesis weighted by their probabilities

Bayes Theorem:-

It provides a way to calculate the prob of hypothesis based on its prior probability, The probability of observing various data given the hypothesis, and the observed data itself

$$P(h|D) = \frac{P(D|h) * P(h)}{P(D)} \rightarrow \text{Prior } P \text{ of } h \\ P(D) \rightarrow \text{marginal P}$$

$P(D)$  = prior probability that training Data D will be observed, i.e. the probability of  $D$  given no knowledge about which hypothesis holds.

$P(D|h)$  = probability of observing data  $D$  in which the hypothesis  $h$ , holds. We are interested in  $P(h|D)$  i.e. the prob that  $h$  holds given the observed training data  $D$ . It is also called posterior prob of  $h$

Many a times we have set of hypothesis  $h$  and we are interesting in finding the most probable hypothesis Data  $D$ . Any such maximally <sup>probable</sup> hypothesis is called max(A) min. maximum a posteriori hypothesis (MAP)

$$h_{MAP} = \arg \max_{h \in H} P(h|D)$$

$$\Rightarrow \arg \max_{h \in H} \frac{P(D|h) * P(h)}{P(D)}$$

maximum  
a posteriori probability  
Hypothesis

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

$$h_{\text{map}} = \text{argmax } P(D|h) \times P(h)$$

Consider a medical diagnosis problem in which there are 2 alternative hypothesis

- That the patient has particular form of cancer
- That the ~~doesn't~~ doesn't have.

- The available data is from a particular test with 2 possible outcomes +ve or -ve. There is prior knowledge about the test.

The test returns a possible result in only 98%.

- The available data is from a particular lab test with 2 possible outcomes i.e +ve and -ve. There is prior knowledge that over entire population of people only 0.008 have this Disease also lab test is only an imperfect indicator of this disease. The test returns a correct +ve result in only 98% of the cases in which the disease is actually present and a correct -ve result is only 97% of cases in which the disease is not present. In other cases the test returns the opposite result. Suppose we have a new patient for whom the lab test returns a +ve result. Should be diagnosed the patient has having cancer or not.

conditional probability. result is there  
but it is +ve in 98%.

Date \_\_\_\_\_  
Page \_\_\_\_\_

18/01/2020

Consider a

Solu<sup>n</sup> :-  $p(\text{cancer}) = 0.008$   $p(\Theta/\sim \text{cancer}) = 0.9$   
 $p(+/\text{cancer}) = 0.98$   $p(\Theta/\text{cancer}) = 0.02$   
 true already when they have cancer

$P(\sim \text{cancer}) = 0.992$   $p(+/\sim \text{cancer}) = 0.98$   
 They have cancer & give -ve result,  $P(\Theta/\text{cancer}) = 0.02$

$h_1 = \text{cancer}$   $h_2 = \sim \text{cancer}$

$h_{\text{map}} = \arg \max_{h \in H} p(h|D)$

$$h_{\text{map}} = P(D|h) \times P(h)$$

$$= P(+/\text{cancer}) \cdot P(\text{cancer}) = 0.98 \times 0.008$$

$$= 0.0078$$

$h_2^{\text{map}} = P(+/\sim \text{cancer}) * P(\sim \text{cancer})$

$= 0.03 * 0.992$

$= 0.0298$

$\boxed{h_{\text{map}} = \sim \text{cancer}}$

~~IC~~ \* Bayes Theorem and Concept learning.

$$\langle x_1, d_1 \rangle, \langle x_2, d_2 \rangle, \dots, \langle x_m, d_m \rangle$$

$x_i$  = instances target =  $d_i$

$$d_i = c(x_i)$$

$\langle \text{sun warm, hot} \rangle$   
 $\langle \text{sun hot cold} \rangle$

Yes  
no  
 $c_i$   
 $d_i$

based on the inst  
concept is learnt  
target is found

$\langle x_1, x_2, \dots, x_m \rangle$  - set of instances fixed.

$D = \{d_1, d_2, \dots, d_m\}$  left column (on b  
mode as set of categories  
belongs to)

## Brute force map learning Algorithm.

### Bayes Theorem and Concept Learning

Concept Consider a finite Hypothesis Space  $H$  defined over an instance space  $X$ . Task is to learn some target concept  $c: X \rightarrow \{0, 1\}$

Let the sequence of training examples be  $\{(x_1, d_1), (x_2, d_2), \dots, (x_m, d_m)\}$  where  $x_i$  is some instance from  $X$  and  $d_i$  is the target value of  $x_i$ , i.e.  $d_i = c(x_i)$

We assume that the sequence of instances  $(x_1, x_2, \dots, x_m)$  is held fixed so that the training dataset  $D$  can be written just as a sequence of target values  $D = \{d_1, d_2, \dots, d_m\}$

We can design a concept learning algorithm to output the map hypothesis based on Bayes Theorem which we call it as Brute force MAP learning Algorithm.

### \* Brute force MAP learning Algorithm

Step 1: - For each hypothesis  $h$  in  $H$  calculate the posterior probability  $P(h|D) = \frac{P(D|h) * P(h)}{P(D)}$

Step 2: - Output the hypothesis  $h_{MAP}$  with the highest posterior probability i.e.  $h_{MAP} = \operatorname{argmax}_{h \in H} P(h|D)$

$$h_{MAP} = \operatorname{argmax}_{h \in H} \frac{P(D|h) * P(h)}{P(D)}$$

$$P(h) = \frac{1}{|H|}$$

$$h_1 = \dots \quad h_2 = \dots \quad h_3 = \dots$$

classmate

when  $d_i = h(x_i) \Rightarrow 1$  (consistent)  
 $\neq n \Rightarrow 0$  (in " "

## Assumptions.

Here we must specify what values are used for  $p(h)$  and  $p(D|h)$ . In order to choose the probability distribution  $p(h)$  we make the following assumptions.

- 1) The training data is noise free i.e  $d_i = c(x_i)$
  - 2) The target concept  $c$  is contained in the hypothesis space  $H$
  - 3) We have no prior reason to believe that any hypothesis is more probable than the other

Given these assumptions what values do we have to

Given that no prior knowledge is available to us, one hypothesis is more likely than <sup>an</sup> other. It is reasonable to assign the same prior probability to every hypothesis  $h$  in  $H$ . Also these prior probabilities sum to one. Hence we can write

$$p(h) = \frac{1}{|H|} + h \in H$$

$P(D|h)$  is the probability of observing the target values  $D = \{d_1, d_2, \dots, d_m\}$  for the fixed set of instances  $\langle x_1, x_2, \dots, x_m \rangle$  given that hypothesis  $h$  holds. Since we assume noise free data the following probability of observing the classifier

$d_i = h(x_i)$  and 0 if  $d_i \neq h(x_i) \therefore P(D|h)$

In other words the probability of data given hypothesis  $h$  is 1 if  $d$  is in consistent.

with  $h$  and 0 otherwise

Given these choices for  $P(h) \in P(D|h)$

we now have a fully defined brute force map learning Algorithm.

Compute the posterior probability  $P(h|D)$  of every hypothesis using Baye's Theorem

Case 1:-  $h$  is inconsistent with the Training data  $D$ ,  $P(h|D) = \frac{0 \times P(h)}{P(D)} = 0$ .  $P(D|h) = 0$ .

The Posterior probability of hypothesis inconsistent with  $D$  is always 0.

Case 2:  $h$  is consistent with  $D$ .

$$P(h|D) = \frac{1 * \frac{1}{|H|}}{|VS_{H,D}|} = \frac{1}{|VS_{H,D}|}$$

where  $VS_{H,D}$  is a subset of Hypothesis from  $H$  that are consistent with  $D$ , ie it is the union space of with respect to  $D$ .

To summarize we can write

$$P(h|D) = \begin{cases} \frac{1}{|VS_{H,D}|}, & \text{if } h \text{ is consistent} \\ 0, & \text{otherwise} \end{cases}$$

\* with this we can say every consistent hypothesis is a map hypothesis

Note:-  $h_{map} = \operatorname{argmax}_{h \in H} P(D|h) * P(h)$  when Prior  $P$  is not given we take some value & hence ignore it and write it as

$$h_{ML} = \operatorname{argmax}_{h \in H} P(D|h)$$

In some cases we assume that every hypothesis in this equally probable that is  $p(h_i) = p(h_j)$  for all  $h_i \in h_j$  belonging to  $H$ .

$$p(h_i) = p(h_j) \forall h_i, h_j \in H.$$

In this case further we simplify.

$$h_{\text{mop}} = \underset{h \in H}{\operatorname{argmax}} P(D|h) * P(h)$$

and need to consider only the term  $P(D|h)$  to find the most probable hypothesis.

$P(D|h)$  is often called as the likelihood of the  $P(D|h)$  and any hypothesis that maximises  $P(D|h)$  is called maximum likelihood hypothesis denoted by  $h_{ML}$ .

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} P(D|h)$$

22/01/2024

### v-imp ML and LS Error Hypothesis

Maximum Likelihood & least squared Error

$$f: X \rightarrow R \quad (\text{Continuous valued target})$$

In the target by associating noise we can make it continuous valued.

We consider a problem of learning a continuous value target function. Consider an instance space  $X$  and the hypothesis space  $H$  consisting of some class of Real valued functions over  $X$  i.e. each  $h \in H$  is of the form of function  $h: X \rightarrow R$  where  $R$  is the set of Real numbers.

The problem faced by the learner  $L$  is to learn an unknown target function  $f: X \rightarrow R$ . A set of ' $m$ ' training examples are provided

$p$  = probability density function

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

where the target value of each example is corrupted by a random noise drawn according to a normal probability distribution.

In other words each training ex is of the form  $\langle x_i, d_i \rangle$  where  $d_i = f(x_i) + e_i$ , where  $f(x_i)$  is noise free value of target func<sup>n</sup> &  $e_i$  is the random variable representing the same

~~Note~~ Bayesian Analysis will show that under certain assumptions any learning algorithm that minimizes the squared error between the output hypothesis predictions and the training data will output a maximum likelihood hypothesis.

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} P(D|h)$$

$\therefore$  we are now dealing with continuous values

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} p(D|h)$$

$$\text{Error}^2 = (\text{Actual} - \text{Predicted})^2 \text{ or } (\text{Predicted} - \text{Actual})^2$$

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} P(D|h) \Rightarrow LS$$

$$\therefore h_{ML} = LS$$

We show this by considering ML hypothesis and we start with the equation.

→ we write the some equ<sup>n</sup> with  $p$  to refer to the probability density

we are taking every instance as independent hence we are taking product . ex:-  $P(A) \cdot P(B)$

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

Probability Density Function  
 Normal Distribution =  $b(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$

$\log e^x = x$

Assuming that training examples are mutual independent, given  $h$  we can write  $p(D|h)$  as the product of various  $p(d_i|h)$  that can be shown as

$$h_{ML} = \operatorname{argmax}_{h \in H} \prod_{i=1}^m p(d_i|h)$$

$p(d_i|h)$  can be written as a normal distribution with variance  $\sigma^2$  and mean  $\mu$ .

$$\mu = f(x_i) \text{ or } h(x_i)$$

$$h_{ML} = \operatorname{argmax}_{h \in H} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - \mu)^2}$$

$$= \operatorname{argmax}_{h \in H} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - h(x_i))^2}$$

$$= \operatorname{argmax}_{h \in H} \sum_{i=1}^m \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \ln \left( e^{-\frac{1}{2\sigma^2}(d_i - h(x_i))^2} \right)$$

(equivent w/ cov.  
discord.)

$$= \operatorname{argmax}_{h \in H} \sum_{i=1}^m \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) \xrightarrow{\text{independent of } h} -\frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

$$\operatorname{argmax}_{h \in H} \sum_{i=1}^m -\frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

/ maximizing -ve func  
some minimi +ve func

$$h_{ML} = \operatorname{argmin}_{h \in H} \sum_{i=1}^m \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

$$h_{ML} = \operatorname{argmin}_{h \in H} \sum_{i=1}^m (d_i - h(x_i))^2$$

ML = least squared error