

— Ist Internals.

05/01/2023

Issues in Decision Tree Learning.

- * Avoiding overfitting the data.
- * Incorporating continuous valued attributes
- * Alternate measures for selecting the attribute
- * Handling training examples with missing attribute value
- * " Attributes with differ in cost.

Shape +ve -ve

Triangle 0 1 $\Rightarrow 0$

Circle 2 0 $\Rightarrow 0$

$$\bar{x}(\text{Shape}) = 0$$

$$\text{Gain} = 0.918$$

Size

Big

Small

No

Shape

Circle

Triangle

Yes

No

Final Tree.

Continued.

Issues in Decision Tree learning.

* Incorporating continuous valued attributes.

\Rightarrow If we have " " s we can make it as Discrete value by using Discrete valued intervals.

Ex! - Age - 4, 10, 14, 24, 24, 34, 64, 84, 95

[child, Young, Adult, Teenage Senior]
 0-5 6-12 20-35 13-19 60-79 80-100

In some attributes we can find Discrete valued intervals we use other method where we find to calculate threshold value 'c'

Step 1: sort.

2:- Look for those value classification is changing

3:- And calculate gain as the previous sum.

Tennis	40	48	60	72	80	80
Play Tennis	No	No	Yes	Yes	Yes	No

$$\frac{48+60}{2} = 54$$

Threshold value 'c'

40 48

60 72 80 90

$$\frac{80+90}{2} = 85$$

85

K-NN Algorithm:-

→ K-Nearest Neighbour Algorithm

→ Multi classification

→ K stands for no. of neighbours.

→ The minimum value of K is always 3

Distance Measures.

→ Euclidean Distance

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

→ Manhattan Distance

$$|x_2 - x_1| + |y_2 - y_1|$$

Value of K will be changed based on the Accuracy

Confusion Matrix

$$\begin{matrix} & 0 & 1 \\ 0 & [00 \quad 01] \\ 1 & [10 \quad 11] \end{matrix}$$

0 → 0 → True +ve

0 → 1 → False Neg

1 → 0 → " +ve

1 → 1 → True -ve

No of Correct predictions

$$\text{Accuracy} = \frac{00 + 11}{00 + 01 + 10 + 11}$$

Issues continued.

3) Alternate measures for selecting the attribute.

Information Gain

Get favours for the attributes which has more value.

classmate

Date _____

Page _____

Disad:- Broader tree at one level.

Alternate measure is Gain Ratio.

↓
Split information

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

where S_1 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

$$\text{Gain}(S, A) \quad A(x) = \text{Temp}/\text{Humidity}/\text{Wind}/\text{Outlook}$$

Training example x , $C(x) > A(x) \Rightarrow A(x) = \text{Null}$

2 strategies

- whichever value is max in that attribute then fill it with that max value irrespective of $C(x)$
- $C(x) \neq$ that instance if it is true or -ve check the target concept whichever is true ~~take~~ or -ve then fill it correspondingly, acc to $C(x)$ value.

Handling attributes with diff costs.

Some times " are associated with costs/weight.

e.g.- To predict medical disease

Temp / Pulse / Blood test / Biopsy / Target

$$\frac{\text{Gain}(S, A)}{\text{Cost}(A)}$$

null
Based on ~~not~~ value.

v. imp Avoiding Overfitting in the Data.

~~less error in train data & less error in test data.~~

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

Given an hypothesis, it said to overfit the training examples if the hypothesis h is performing better over the training set than over the testing set. so Training error(h) $<$ Testing error(h).
 error less accuracy more.

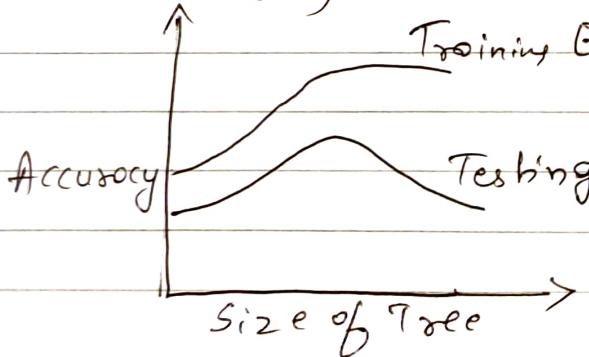
why is this overfitting causing in ID3
 → Noise.

→ If size of data set is small They are not in a position to fit.

eg: 1 (Noise)

• Outlook = sunny, Temp = hot, Humidity = Normal,
 wind = strong, Play tennis = No

(for the tree)



Avoiding overfitting in the Data.

Strategy:-

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

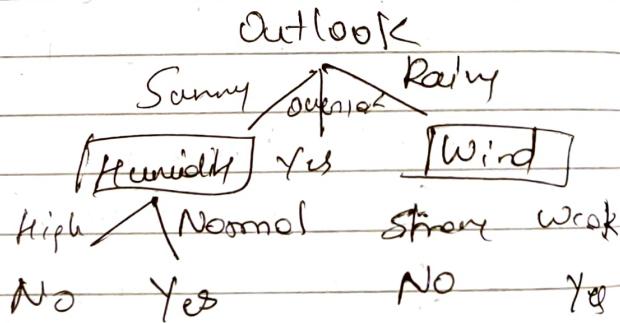
Disadv:- It is not easy to determine where we stop growing the tree. (practical difficult)

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 leaf 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

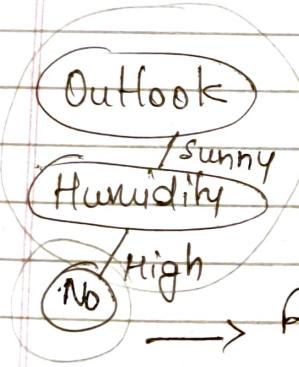
Pruning of tree continues until further pruning is harmful. The pruning is stopped once the accuracy ~~stop~~ starts decreasing.

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.

and

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.

- 1) * Build a Decision Tree from the Training set, grow the tree completely until training data is fit 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 in tree.
 - * The path from the root to leaf becomes Rule Antecedent or pre-condition.
 - * A classification at leaf node becomes the rules 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.

01/01/2024 3rd Module.

Bayesian learning.

Computing probabilities to find some inferences/prediction

Bayes Theorem

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

Prior Probability Marginal Probability

Likelihood

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

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

$$P(\text{king} | \text{facecard}) = \frac{P(\text{face} | \text{king}) \cdot P(\text{king})}{P(\text{face})}$$

$$= \frac{4/4 \times 4/52}{12/52} = \frac{1}{3}$$

$$P(K|F) = \frac{1}{3}$$

9. Bayesian Learning

- 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 in the estimated probability of 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 the predictions of multiple hypothesis weighted by their probabilities

Bayes Theorem:-

It provides a way to calculate the probability 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 most probable hypothesis is called maximum a posteriori hypothesis (MAP).

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

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

maximum
a posteriori probability hypothesis

classmate

Data

Page

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

Consider a medical diagnosis problem in which we have two 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 disease and also a lab 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. cancer is there
but it is +ve in 98%.

18/01/2024

Date _____
Page _____

Soluⁿ :- Consider a test for cancer.
 $P(\text{cancer}) = 0.008$
 $P(+/\text{cancer}) = 0.98$ $P(-/\sim \text{cancer}) = 0.97$
true already when they have cancer
 $P(\sim \text{cancer}) = 0.992$ $P(-/\sim \text{cancer}) = 0.008$

They have cancer & give -ve result, $P(-/\text{cancer}) = 0.02$

$h_1 = \text{cancer}$ $h_2 = \text{not cancer}$. $h_2 = \sim \text{cancer}$

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

$$h_{\text{mop}} = P(D|h) \times P(h)$$
$$= P(+/\text{cancer}) \cdot P(\text{cancer}) = 0.98 \times 0.008$$
$$= 0.0078.$$

$$h_{\text{mop}} = P(+/\sim \text{cancer}) * P(\sim \text{cancer})$$

$$= 0.02 * 0.992$$

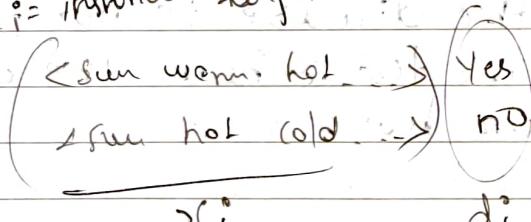
$$= 0.0298$$

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

~~inf~~ Bayes' Theorem and Concept learning.

$$\{x_i, d_i\} \rightarrow \langle x_1, d_1 \rangle, \langle x_2, d_2 \rangle, \dots, \langle x_m, d_m \rangle \}$$

$x_i = \text{instances}$, $d_i = \text{target} = d^*$



based on the instance
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 can be
made as set which
belongs to positive
zone

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 data set 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

imp.

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

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

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

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

CLASSMATE
Date _____
Page _____

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

Assumptions

Hence we must specify what values are to be used for $p(h)$ and $p(D|h)$. In order to choose the probability distribution $p(h)$ & $P(D|h)$ 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 that one hypothesis is more likely than ^{an} other. It is reasonable to assign the same prior probability to every hypothesis h in H . Also these prior probability sum to one. Hence we can write

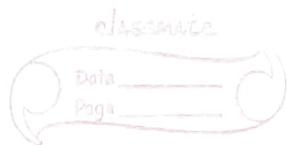
$$P(h) = \frac{1}{|H|} \quad \forall h \in H$$

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

$$d_i = h(x_i) \text{ and } 0 \text{ if } d_i \neq h(x_i) \therefore P(D|h) = \begin{cases} 1, & \text{if } d_i = h(x_i) \quad \forall d_i \in D \\ 0, & \text{otherwise.} \end{cases}$$

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

$$P(h|D) = \frac{P(D|h) \cdot P(h)}{P(D)} \quad \mathcal{H}$$



with h and 0 otherwise

Given these choices for $P(h)$ & $P(D|h)$

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

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

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

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

Case 2: h is consistent with D .

$$P(h|D) = \frac{1 * \frac{1}{|\mathcal{H}|}}{\frac{|\text{VS}_{H,D}|}{|\mathcal{H}|}} = \frac{1}{|\text{VS}_{H,D}|}$$

where $\text{VS}_{H,D}$ is a subset of Hypothesis from \mathcal{H} that are consistent with D , ie it is the version space of \mathcal{H} with respect to D .

To summarize we can write

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

With this we can say every consistent hypothesis is a map hypothesis

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

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

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

In this case further we simplify.
 $h_{\text{map}} = \underset{h \in H}{\text{argmax}} P(D|h) \neq P(h)$

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} .

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

22/01/2024

v. imp ML and LS Error Hypothesis

Maximum Likelihood & least squared Error
 $f: X \rightarrow R$ (Continuous valued target)

In the target by adding 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 funcⁿ & e_i is the random variable representing the same

int 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.

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

∴ we are now dealing with continuous values

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

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

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

$$\therefore h_m \approx LS$$

We show this by considering NL hypothesis and we start with the equation:

→ we write the same eqn with p to refer to the probability density

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

Probability of Density
 Normal Distribution = $p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$

classmate

Date _____

Page _____

Assuming that predicting examples are mutually 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 σ^2 and mean μ

$$\mu = f(x_i) \approx h(x_i)$$

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

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

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

constant we can discard.

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

$$\operatorname{argmax}_{h \in H} \sum_{i=1}^m -\frac{1}{2\sigma^2} (d_i - h(x_i))^2 \quad // \text{maximizing -ve term same minimizes +ve term.}$$

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

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

ML = least squared error

3/01/2024

classmate

Date _____

Page _____

imp. Maximum likelihood hypothesis for predicting Probabilities.

Learned \rightarrow Learn to predict probabilities.

Medical data set $x \rightarrow$ instance.

(it contains symptoms of Disease)

Target funcⁿ $f(x) = 1 \Rightarrow$ survives

$f(x) = 0 \Rightarrow$ Dies

It is a probabilistic funcⁿ

Set of patients 92% survives 8% dies.

$$f: x \rightarrow [0,1] \Rightarrow P(f(x)=1) = 0.92$$

$$f' = P(f(x)=1) = 0.92$$

What is the P that patient survives / $P(f(x)=1)$

$f' \Rightarrow$ Machine should learn this.

- * The objective is to learn to predict probability. Suppose we have a medical data set instance space X represents patient having in terms of symptoms of a disease. Target function $f(x) = 1$, if the patient survives the disease and $f(x) = 0$, if he doesn't.

Here f is a probabilistic function. Suppose we have a collection of patients exhibiting the same set of symptoms, we ^{wish} find that 92% survives & 8% do not. We want to find out what is the probability that $f(x) = 1$ for this we learn a target function $f': x \rightarrow [0,1]$ such that $f'(x) = P(f(x)=1)$

Based on our example $f'(x) = 0.92 = P(f(x)=1)$ To learn this function f' we desire a maximum likelihood hypothesis for f' .

because it is independent we can write in product form.

$$D = (x_1 d_1, x_2 d_2, \dots, x_n d_n)$$

CLASSTIME

Date _____
Page _____

For this, we must obtain an expression for $P(D|h)$

$$P(D|h) = \prod_{i=1}^m P(x_i; d_i | h) \rightarrow ①$$

Let us assume that training data D is of the form $D = \{<x_1; d_1>, <x_2; d_2>, \dots, <x_m; d_m>\}$ where $d_i = 0 \text{ or } 1$ for $f(x_i)$.

Treating both x_i and d_i as random variables and assuming that every training example is independent of other.

We can write $P(D|h)$ as ①

($\text{only } d_i \text{ depends on hypothesis not the patient details or } x_i$)

We assume that the probability of any particular instance x_i is independent of hypothesis h , when x is independent of h and survival of the patient alone i.e. d_i depends on h we can rewrite ① as below

$$P(D|h) = \prod_{i=1}^m P(d_i|h, x_i) \times P(x_i) \rightarrow ②$$

Now what is the prob $P(d_i|h, x_i)$ of observing $d_i=1$ for a single instance x_i , given that hypothesis h holds?

h is our hypothesis regarding this target func" which computes this probability.
Therefore $P(d_i=1|h, x_i) = h(x_i)$ and in general b:

$$P(d_i|h, x_i) = \begin{cases} h(x_i) & \text{if } d_i=1 \\ 1-h(x_i) & \text{if } d_i=0 \end{cases} \rightarrow 3$$

In order to substitute 3 in ~~to~~ 2 for $P(D|h)$

we express 3 as follows, that is found

$$P(d_i|h, x_i) = h(x_i)^{d_i} \cdot (1-h(x_i))^{1-d_i} \rightarrow 4$$

Substitute eq ~~4~~ in ~~eq 2~~

$$P(D|h) = \prod_{i=1}^m (1-h(x_i))^{1-d_i} \times P(x_i) \rightarrow 5$$

Now we write an expression for max likelihood hypothesis.

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m h(x_i)^{d_i} \cdot (1-h(x_i))^{1-d_i} \times P(x_i) \rightarrow 6$$

$P(x_i)$ is independent of h so we will drop it, for rest of the equation by applying ~~log~~ terms from both sides.

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m d_i \cdot h(x_i) + (1-d_i) \log(1-h(x_i)) \rightarrow 7$$

Equation 7 describes the quantity that must be maximised in order to obtain the maximum likelihood hypothesis for the given problem setting.

- * **Naive Bayes Classification:**
objective: predict class of the new sample / some ^{ex} ~~as~~
It is a highly practical bayesian learning ^{KNN} method. The objective of Naive Bayes classifier is to predict the class of the new sample. It applies to learning task where each instance x is described by conjunction of attribute values and where the target function $f(x)$ can take on any value from some finite set of V .

09/01/24

where we consider a conjunction of attributes.

classmate

Date _____

Page _____

Tasks $\rightarrow \langle a_1, a_2, a_3, a_4 \rangle$

if $a_1 \wedge a_2 \wedge a_3 \wedge a_4$ then ES = Yes / No.

(finite set $V = \{ \text{contains all target values yes/no} \}$
collection of target values)

"Naïve Bayes classifier always predicts the most probable target value." (V_{MAP})

$$V_{MAP} = \underset{v \in V}{\operatorname{argmax}} P(v_j | a_1, a_2, a_3, a_n)$$

A set of training examples of

and the new instances.
describes the

A set of training examples of target function is provided and a new instance is described by the tuple of attribute values $\langle a_1, a_2, a_3, \dots, a_n \rangle$ the learner is asked to predict the value of the target class for the new instance. the Bayesian approach for classifying the new approach is to assign the most probable target value denoted as V_{MAP} , given the attribute values $\langle a_1, a_2, a_n \rangle$ that describe the new instance. we can write V_{MAP} as:

$$V_{MAP} = \underset{v \in V}{\operatorname{argmax}} P(v_j | a_1, a_2, \dots, a_n)$$

We can use Bayes theorem to rewrite this expression as follows:

$$V_{MAP} = \underset{v_j \in V}{\operatorname{argmax}} \frac{P(a_1, a_2, \dots, a_n | v_j) \times P(v_j)}{P(a_1, a_2, \dots, a_n)}$$

$$V_{MAP} = \underset{v_j \in V}{\operatorname{argmax}} P(a_1, a_2, \dots, a_n | v_j) \times P(v_j) \rightarrow ①$$

We can attempt to calculate the second term in equation ① i.e. it is easy to estimate each of $V_{MAP} = \underset{v_j \in V}{\operatorname{arg}}$ the $P(v_j)$ by counting the frequency with which each target value v_j occurs in the training data, and we need to know how to estimate the first term.

The Naive Bayes classifier is based on the assumption that the attribute values are conditionally independent given the target value in other words the assumption is that given the target value of the instance the probability of the conjunction $\langle a_1, a_2, \dots, a_n \rangle$ is just the product of the probabilities of the individual attributes i.e. $P(a_1, a_2, \dots, a_n | v_j)$ can be written as $P(a_1, a_2, \dots, a_n | v_j) = \prod_i P(a_i | v_j)$ $\rightarrow ②$

Substitute ② in ①, we have the approach used by Naive Bayes classifier as follows.

$$V_{NB} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j) \times \prod_i P(a_i | v_j) \rightarrow ③$$

where V_{NB} denotes the target value output by the Naive Bayes classifier.

$$V_{NB} = \underset{v_j \in V}{\operatorname{argmax}} p(v_j) \times \prod_i (a_i | v_j)$$

classmate

Date _____

Page _____

Example - 1

* Play tennis dataset

new instance = (outlook = sunny, temp = cool, humidity = high, wind = strong)

$$V_{NB} = \underset{v_j \in \{Yes, No\}}{\operatorname{argmax}} p(v_j) \times P(\text{outlook} = \text{sunny} | v_j) \times$$

$$\times P(\text{temp} = \text{cool} | v_j) \times P(\text{humidity} = \text{high} | v_j) \times P(\text{wind} = \text{strong} | v_j)$$

$$V_{NB}(\text{Yes}) = \underset{v_j \in \text{Yes}}{\operatorname{argmax}} p(\text{Yes}) \times p(\text{sunny} | \text{Yes}) \times P(\text{cool} | \text{Yes}) \times$$

$$\times P(\text{high} | \text{Yes}) \times P(\text{strong} | \text{Yes})$$

$$= \frac{9}{14} \times \frac{2}{3} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9}$$

$$(0.642) (0.142) (0.214) (0.214) (0.214)$$

$$(0.642) (0.23) (0.34) (0.34) (0.34)$$

$$= 0.005$$

$$V_{NB}(\text{No}) = \underset{v_j \in \text{No}}{\operatorname{argmax}} p(\text{No}) \times p(\text{sunny} | \text{No}) \times P(\text{cool} | \text{No}) \times$$

$$\times P(\text{high} | \text{No}) \times P(\text{strong} | \text{No})$$

$$= \frac{5}{14} \times \frac{3}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5}$$

$$(0.357) (0.6) (0.2) (0.8) (0.6)$$

$$= 0.02$$

$$V_{NB}(\text{No}) > V_{NB}(\text{Yes})$$

hence new instance is classified as No.

Example-2

#	Color	Legs	Height	Smelly	Species
1	white	3	short	Yes	M
2	Green	2	Tall	No	M
3	Green	3	short	Yes	M
4	White	3	short	Yes	H
5	Green	2	short	No	H
6	White	2	Tall	No	H
7	white	2	short	Yes	H
8	white	2	Tall	No	H

New instance = Color = Green, Legs = 2, Height = Tall,
Smelly = No)

$$V_{NB} = \arg\max_{v_j \in \{M, H\}} P(v_j) \times P(\text{Color} = \text{Green}/v_j) \times P(\text{Legs} = 2/v_j) \times P(\text{Height} = \text{Tall}/v_j) \times P(\text{Smelly} = \text{No}/v_j)$$

$$V_{NB} = \arg\max_{v_j \in M} P(M) \times P(\text{Green}/M) \times P(2/M) \times P(\text{Tall}/M) \times P(\text{No}/M)$$

$$\frac{1}{8} \times \frac{2}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4}$$

$$(0.5)(0.5)(0.25)(0.25)(0.25) \\ = 0.0039$$

$$V_{NB} = \arg\max_{v_j \in H} P(H) \times P(\text{Green}/H) \times P(2/H) \times P(\text{Tall}/H) \times P(\text{No}/H)$$

$$= \frac{4}{8} \times \frac{1}{4} \times \frac{4}{4} \times \frac{2}{4} \times \frac{3}{4}$$

$$(0.5)(0.25)(1)(0.5)(0.75) = 0.0468$$

$$V_{NB}(H) > V_{NB}(M)$$

new instance is classified as H.

Example - 3

#	Color	Type	Origin	Stolen	
1	Red	Spoof	Domestic	Yes	
2	Red	Spoof	Domestic	No	
3	Red	Spoof	Domestic	Yes	
4	Yellow	Spoof	Domestic	No	
5	Yellow	Spoof	Imported	Yes	$P(\text{Yes}) = \frac{5}{10} = \frac{1}{2}$
6	Yellow	SUV	Imported	No	
7	Yellow	SUV	Imported	Yes	
8	Yellow	SUV	Domestic	No	
9	Red	SUV	Imported	No	
10	Red	Spoof	Imported	Yes	
	New	Red, SUV, Domestic			

$$V_{NB} = \max_{Y_j \in Yes} P(Y_j) \cdot P(\text{Red}/Y_j) \cdot P(\text{SUV}/Y_j) \cdot P(\text{Domestic}/Y_j)$$

$$= \frac{5}{10} \cdot \frac{3}{5} \cdot \frac{1}{5} \cdot \frac{2}{5}$$

$$(0.5)(0.6)(0.2)(0.4) = 0.024$$

$$V_{NB} = \max_{Y_j \in No} P(No) \cdot P(\text{Red}/No) \cdot P(\text{SUV}/No) \cdot P(\text{Domestic}/No)$$

$$= \frac{5}{10} \cdot \frac{2}{5} \cdot \frac{3}{5} \cdot \frac{3}{5}$$

$$(0.5)(0.4)(0.6)(0.6) = 0.072$$

$$V_{NB}(\text{No}) > V_{NB}(\text{Yes})$$

Hence new instance is classified as No.

Example - 4

A B C class.

1	0	0	0	+
2	0	0	1	-
3	0	1	1	-
4	0	1	1	-
5	0	0	1	+
6	1	0	1	+
7	1	0	1	-
8	1	0	1	-
9	1	1	1	+
10	1	0	1	+

new inst (0 1 0)

$$P(+)(0/+) \cdot P(1/+) \cdot P(0/+)$$

$$V_{NB} = \frac{5}{10} \cdot \frac{2}{5} \cdot \frac{1}{5} \cdot \frac{1}{5}$$

$$= (0.5) (0.4) (0.2) (0.2) = 0.008$$

$$V_{NB} = P(-) P(0/-) P(1/-) P(0/-)$$

$$(0.5) \cdot \frac{3}{5} \cdot \frac{2}{5} \cdot 0$$

$$V_{NB}(+) > V_{NB}(-)$$

Hence new instance can be classified as '+'

Example - 5

$$P(+)=\frac{4}{7}, P(-)=\frac{3}{7}$$

new instance (open, Audit, Red, 1990, spots)

$$V_{NB} = \frac{4}{7} \left(\frac{4}{4} \right) (0) (0) \left(\frac{1}{4} \right) (0) = 0$$

$$V_{NB} = \frac{3}{7} \left(\frac{1}{3} \right) \left(\frac{1}{3} \right) \left(\frac{2}{3} \right) (0) \left(\frac{1}{3} \right) = 0$$

$$V_{NB}(+) = V_{NB}(-) = 0$$

So we can choose arbitrarily either + or -

01/02/2024

most imp ★ Bayesian Belief Network.

Consider a set of random variables Y_1, Y_2, \dots, Y_n where each variable Y_i can take on a set of possible values $V(Y_i)$ we define the joint space of the set of variables Y , as the to be the cross product of $V(Y_1) \times V(Y_2) \times \dots \times V(Y_n)$. The probability distribution over this joint Space is joint " "

A Bayesian Belief network describes the joint probability distribution for the set of variables.

Conditional Independence.

Let x, y, z be three discrete valued random variables. we say that x is conditionally independent y given z , if the probability distribution learning x is independent of the value of y given a value for z , i.e if (for all) $\forall x_i, y_j, z_k \in V(x), V(y), V(z)$

$$P(X=x_i | Y=y_j, Z=z_k) = P(X=x_i | Z=z_k)$$

where $x_i \in V(x)$, $y_j \in V(y)$, $z_k \in V(z)$.

We can rewrite this expression as

$$P(X|Y, Z) = P(X|Z)$$

This

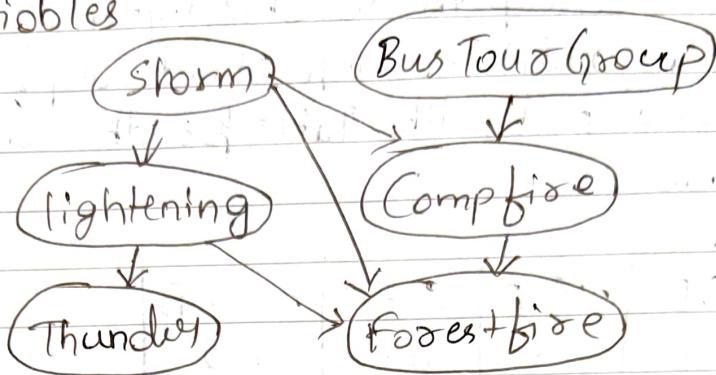
can also be extended to a set of variables

$$P(X_1, X_2, \dots, X_m | Y_1, Y_2, \dots, Y_n, Z_1, Z_2, \dots, Z_l) =$$

$$P(X_1, \dots, X_m | Z_1, \dots, Z_l)$$

Representation of Bayesian Belief Network.

IBBN represents the joint PD for a set of variables



The Bayesian n/w in the figure represents the joint probability Dist over the Boolean Variables, storm, lightning, thunder, Bustourgroup, Compfire and forestfire.

Each variable in the joint space bu is represented by a node in the network. For each variable 2 types of information are specified.

- ① → The network arcs represent that the variable is conditionally independent of its non-descendent in the n/w, Given its immediate predecessor in the n/w. we say that x is a descendent of y , if there is a directed path or arc from y to x .
- ② → A Condin Prob table is given for each variable describing the PD for that variable Given the values of its immediate predecessor

In this fig consider the node compfire the n/w nodes and arcs represent the campfire is conditionally independent of its non-descendents lightning and thunder, given its immediate parents

C = campfire

P = Probability

S = storm

B = BusTown Group

classmate

Date _____

Page _____

Bus Town Group

or predecessors storm and BTG. this means that once we know the values of variables storm and BusTown Group, the variable lightning and thunder provide no additional info about campfire. The conditional probability Table for the campfire node is as follows.

Campfire

	S, B	~S, B	S, ~B	~S, ~B
C	0.6	0.5	0.1	0.2
~C	0.4	0.5	0.9	0.8

- The top left entry expresses the assertion that the prob of ($P(C = \text{True} / S = \text{True}; B = \text{True})$)
- The set of conditional prob tables for all the variables together with the set of conditional independence assumptions described by the network describes the full joint probability distribution for the network.

Example-2 There is a burglar alarm installed at home. It is fairly reliable at detecting burglary but also sometimes it responds to minor earthquakes. You have 2 neighbours J & M who have promised to call you at work when they hear the alarm. J always calls when he hears the alarm but sometimes confuses telephone ringing with the alarm and calls too. M likes loud music and sometimes misses the alarm.

a) what is the prob that the alarm has sounded but neither burglary nor earthquake has occurred by both J and M calls.

$$P(B) = 0.01$$

Burglary

Earthquake

classmate

Date _____

Page _____

$$P(E) = 0.002$$

Alarm

J calls

M calls

CPT for M calls

$$A \quad P(M|A) \quad P(\sim M|\bar{A})$$

$$\text{True} \quad 0.70 \quad 0.30$$

$$\text{False} \quad 0.01 \quad 0.99$$

(CPT) Conditional prob table for alarm.

$$B \quad E \quad P(A|B,E) \quad P(\sim A|B,E) \quad \text{CPT for J calls}$$

$$T \quad T \quad 0.95 \quad 0.05 \quad A \quad P(J|A) \quad P(\sim J|\bar{A})$$

$$T \quad F \quad 0.94 \quad 0.06 \quad \text{True} \quad 0.90 \quad 0.1$$

$$F \quad T \quad 0.29 \quad 0.71 \quad \text{false} \quad 0.05 \quad 0.95$$

$$F \quad F \quad 0.001 \quad 0.999$$

Joint Prob Distribution

$$P(J \wedge M \wedge A \wedge \sim B \wedge \sim E)$$

$$= P(J|A) \cdot P(M|A) \cdot P(A| \sim B, \sim E) \cdot P(\sim B) \cdot P(\sim E)$$

$$= 0.90 \times 0.70 \times 0.001 \times 0.998$$

$$= 0.0006224$$

es

05/2/2024

$$P(J) = P(J|A) * P(A) + P(J|\sim A) * P(\sim A)$$

$$= 0.90 * P(A) + 0.05 * P(\bar{A})$$

$$P(A) = P(A|B,E) * P(B \wedge E) + P(A|\sim B,E) * P(\sim B \wedge E)$$

$$+ P(A|B,\sim E) * P(B \wedge \sim E) + P(A|\bar{B},\bar{E}) * P(\bar{B} \wedge \bar{E})$$

$$P(A) = 0.95 * P(B) * P(E) + 0.29 * P(\bar{B}) * P(E) + 0.94 * P(B) * P(\bar{E}) + 0.001 * P(\bar{B}) * P(\bar{E})$$

$$= 0.95(0.01)(0.002) + 0.29(0.99)(0.002) + 0.94(0.01)(0.998) + 0.001(0.99)(0.998)$$

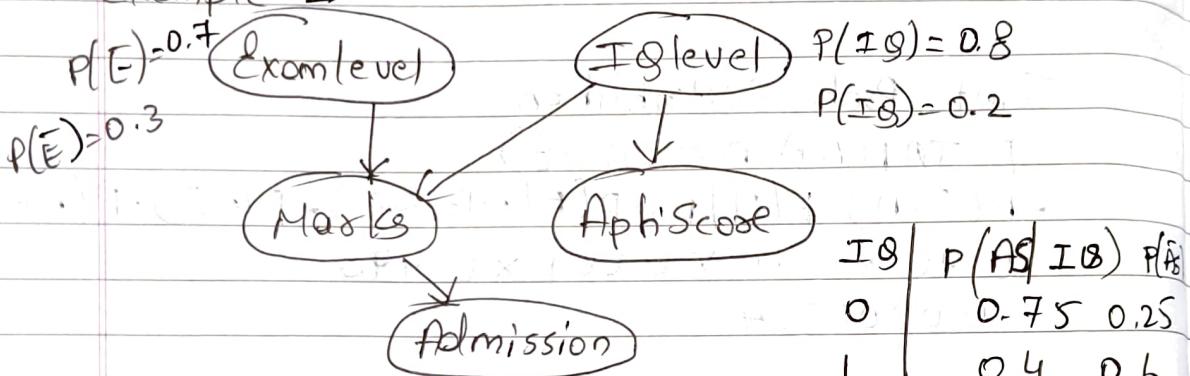
$$P(A) = 0.0109$$

Because $P(B \wedge E)$ is independent $P(B) \cdot P(E)$

$$\begin{aligned}
 P(\bar{A}) &= P(\bar{A} | B, E) * P(B \cap E) + P(\bar{A} | \bar{B}, E) * P(\bar{B} \cap E) \\
 &\quad + P(\bar{A} | B, \bar{E}) * P(B \cap \bar{E}) + P(\bar{A} | \bar{B}, \bar{E}) * P(\bar{B} \cap \bar{E}) \\
 &= (0.05)(0.01)(0.002) + 0.000050.71(0.99)(0.002) \\
 &\quad + (0.06)(0.01)(0.998) + 0.999(0.99)(0.998) \\
 &= 0.9890 \quad 0.000050.71 - 0.0109
 \end{aligned}$$

$$\begin{aligned}
 P(J) &= 0.90(0.0109) + 0.05(0.9890) \\
 &= 0.0592
 \end{aligned}$$

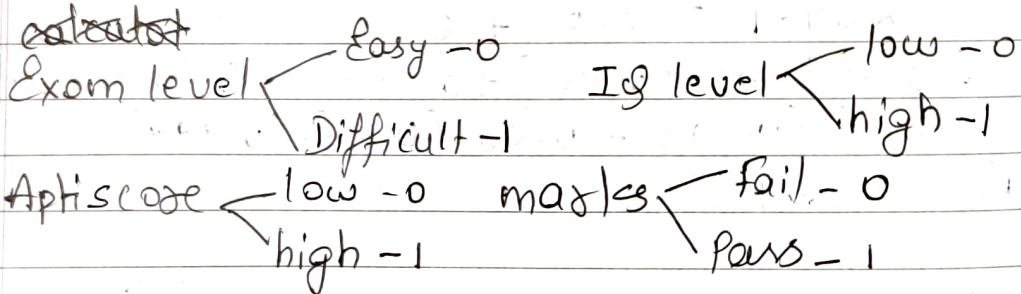
Example - 2



IQ	E	$P(M IQ, E)$
0	0	0.6
0	1	0.9
1	0	0.5
1	1	0.8

M	$P(A M)$	$P(\bar{A} M)$
0	0.6	0.4
1	0.9	0.1

→ Calculate the probability that inspite of the exam level being difficult the student having a low IQ level and a low aphiscose manages to pass the exam and secure admission in the university.



$P(E) = \text{Difficult Exam level}$

$P(IQ) = \text{low IQ level}$

Admission ↗ Secure/No = 0
Yes = 1

$$P(A=1, m=1, i=0, e=1, s=0)$$

$$P(A=1/m=1) \approx P(m=1/e=1, i=0) * P(i=0) *$$

$$P(e=1) * P(s=0/i=0)$$

$$= (0.9)(0.9)(0.8)(0.7)(0.75)$$

$$= 0.3402$$

→ Calculate the probability that the student has a high IQ level & API score, the exam being easy yet fails to pass and doesn't secure admission.

$$P(A=0, m=0, i=1, e=0, s=1)$$

$$P(A=0/m=0) \approx P(m=0/e=0, i=1) * P(i=1) * P(e=0)$$

$$P(s=1/i=1)$$

$$= 0.6 * 0.5 * 0.2 * 0.3 * 0.4$$

$$= 0.0072$$

Example - 3:

$P(E) = 0.7$

$P(D) = 0.25$

Exercise

Diet

E	D	$P(HD/E, D)$
Y	Y	0.25
Y	N	0.45
N	Y	0.55
N	N	0.75

Heart Disease

Heart Burn

Blood Poisoning

Chest pain

D	$P(HB/D)$
Y	0.2
N	0.85

HD	$P(BP/HD)$
Y	0.85
N	0.2

HD	$P(CP)$
Y	0.8
Y	0.6
N	0.4
N	0.1

what is the Prob that the Person is having high BP
not solved.

06/02/2029

Association Analysis.

Methodology to identify interesting pattern or relationships.

Market Basket Transaction.

collection of items

Tid	Bread	Butter	Jam	Milk	eg.
0	1	1	0	1	
1	0	0	1	0	

2 issues:-

→ computational expensive

→ result in spurious relationships.
(by chance)

$$T = \{t_1, t_2, \dots, t_n\} \quad I = \{i_1, i_2, \dots, i_m\}$$

Itemset: Collection of items.

Let $I = \{i_1, i_2, i_3, \dots, i_d\}$ be the set of all items in the market basket data. Let $T = \{t_1, t_2, t_3, \dots, t_n\}$ be the set of all transaction. Each Transaction t_j contains subset of items chosen from I .

as on.

A collection of two or more items is called itemset.

If an itemset contains ' k ' items it is called k itemset. $\{Bread, Butter, Jam\}$ is an example of 3 itemset.

The null or empty itemset doesn't contain any items.

Transaction width

τ is the number of items present in the transaction. A Transaction t_i is said to contain an itemset x , if x is a subset of t_i .

An important property of an itemset is called support count.

\times Support count is the no. of transaction that contain a particular itemset x . It is represented as $\sigma(x)$, support is called as "support count" of x . Mathematical Representation.

$$\sigma(x) = |\{t_i \in T : t_i \subseteq x\}|$$

~~↳ A root~~ ~~↳~~ Association Rule.

It is an implication expression of form $x \rightarrow y$ where x & y are disjoint item set i.e $x \cap y = \emptyset$.

Two important notes
 \rightarrow Support - A measure of trust.

Support Count

No of transaction that contain a particular itemset x , it is represented as $\sigma(x)$

ex: $x = \{\text{Bread, Butter}\}$

Tid Bread Butter Jam Support count {t₁}

t ₁	1	1	0
----------------	---	---	---

t ₂	1	0	0
----------------	---	---	---

Represented by σ_x where x is itemset.

$$\sigma_x = |\{t_i | x \subseteq t_i, t_i \in T\}|$$

Association Rule

we write in the form of implication ($X \rightarrow Y$)
 Ex:- $\{Bread, Butter\} \rightarrow \{Rice\}$
 If I buy these, then I will have

If it is a implication expression of form $X \rightarrow Y$ where X and Y are disjoint itemset i.e. $X \cap Y = \emptyset$

The strength of an association rule can be measured in terms of support(s) and confidence(c)

support: $s(X \rightarrow Y) = \frac{|T(X \cup Y)|}{N}$ where N is no of transaction.

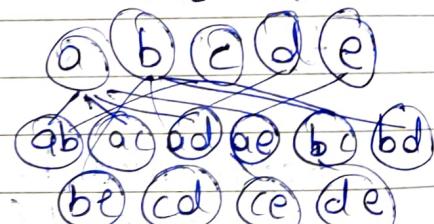
confidence: $c(X \rightarrow Y) = \frac{|T(X \cup Y)|}{|T(X)|}$

formulation of Association Rule Mining Problem or strength of rule is found by SEC. Association Rule Discovery.

Brute force T // Identify Given a set of Transaction
 $T \Rightarrow$ Support \geq min support \Rightarrow Strong rules
 $T \Rightarrow$ Confidence \geq min threshold of rules confidence \Rightarrow confidence which obeys both 2 tasks

- ① \rightarrow Frequent Itemset Generation \Rightarrow "In frequent minus"
- ② \rightarrow Rule Generation

① Itemset $I = \{a, b, c, d, e\}$
 calculate only support Lattice structure
 {objective:- To identify strong rules}



abc abd ace

abcd abcde

abcde

$2^k - 1$ itemsets.

classmate

Date _____

Page _____

T _i	Condition	* Reduce the Complexity of ① 2 types
1. a,b,c	a → 3 ab → 2	* Reduce the no. of itemset
2. ab	b → 3 ac → 2	* Reduce the no. of itemset
3. ade	c → 3 ;	comparison.
4. de	d → 3 ;	
5. acd	e → 3 ;	We use Apriori Principle
6. b,e		To Reduce the no. of itemset to calculate the support.

Apriori Principle

states that if a particular itemset is found to be frequent then all its subset are also said to be frequent. Conversely if an itemset is found to be infrequent all its superset items are also said to be infrequent.

- ① Given set of T, identify the itemset that satisfy the min supp are frequent itemsets.

9/02/24

② Rule Generation

Once frequent itemset are identified, next is to generate Rule generation. (Only strong rule are) Here itemset should satisfy minconfidence > minconf. Finally only strong rule are selected.

- * find the frequent itemset and generate Association Rule for the given Problem.

T:

$$T = \{A, B, C\}$$

$$D \rightarrow BC$$

$$B \rightarrow AC$$

$$C \rightarrow AB$$

$$AB \rightarrow C$$

$$AC \rightarrow B$$

$$BC \rightarrow A$$

$$I = \{ \underline{A}, \underline{B}, \underline{C} \}$$

$$A, B, C, \underline{AB}, \underline{AC}, \underline{BC}$$

$$\text{minSupp} = 50\% = \frac{50}{100} \times 5 = 2.5$$

Date _____
Page _____

$$T = \{T_1, T_2, T_3, T_4, T_5\} \quad \text{minSupp} = 2$$

$$I = \{i_1, i_2, i_3, i_4, i_5, i_6\} \quad \text{minConf} = 75\%$$

Tid | Transactions

T ₁	i ₁ , i ₃ , i ₄
T ₂	i ₂ , i ₃ , i ₅ , i ₆
T ₃	i ₁ , i ₂ , i ₃ , i ₅
T ₄	i ₂ , i ₅
T ₅	i ₁ , i ₃ , i ₅

To find strong rules we first take
 → Find out frequent itemsets and support count.

I-itemsets | Support count

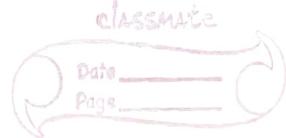
i ₁	P ₁ = 3
i ₂	P ₂ = 3
i ₃	P ₃ = 4
i ₄	P ₄ = 1 X
i ₅	P ₅ = 4
i ₆	P ₆ = 1 X

→ Consider 2 itemsets and find out frequent itemset but using apriori principle where it is says some itemset & its infrequent if its superset are also infrequent. but i₄ & i₆ are infrequent cause support count < minsupport.

I2-itemsets | Support count

i ₁ , i ₂	1 X
i ₁ , i ₃	3
i ₁ , i ₅	2 X
i ₂ , i ₃	2 X
i ₂ , i ₅	3
i ₃ , i ₅	3

~~o(x)~~



3-itemsets	Supp count
i_1, i_3, i_5	2
i_2, i_3, i_5	2

4-itemset	Supp count
i_1, i_2, i_3, i_5	1 \times

→ satisfy min supp condition.

Frequent itemsets	Supp count
i_1	3
i_2	3
i_3	4
i_5	4

i_1, i_3	3
i_1, i_5	2
i_2, i_3	2
i_2, i_5	3
i_3, i_5	3
i_1, i_3, i_5	2
i_2, i_3, i_5	2

Rule Generation

is based on frequent itemset generation.

1-itemset we will leave if cause we cannot write $X \rightarrow Y$ format.

Association Rules

$i_1 \rightarrow i_3$

confidence, $conf(X \rightarrow Y) = \frac{sup(X \rightarrow Y)}{sup(X)}$

$$3/3 = 1 = 100\%$$

$\neg(X)$

$i_3 \rightarrow i_1$

$$3/4 = 75\%$$

$i_1 \rightarrow i_5$

$$2/3 = 66.7\% \times$$

$i_5 \rightarrow i_1$

$$2/4 = 50\% \times$$

$i_2 \rightarrow i_3$

$$2/3 = 66.7\% \times$$

$i_3 \rightarrow i_2$

$$2/4 = 50\% \times$$

$i_2 \rightarrow i_5$

$$3/3 = 100\%$$

$i_5 \rightarrow i_2$

$$3/4 = 75\%$$

$i_3 \rightarrow i_5$

$$3/4 = 75\%$$

$i_5 \rightarrow i_3$

$$3/4 = 75\%$$

$i_1 \rightarrow i_3, i_5$

$$2/3 = 66.7\% \times$$

$i_3 \rightarrow i_1, i_5$

$$2/4 = 50\% \times$$

$i_5 \rightarrow i_1, i_3$
 $i_1, i_3 \rightarrow i_5$
 $i_1, i_5 \rightarrow i_3$
 $i_3, i_5 \rightarrow i_1$
 $i_2 \rightarrow i_3, i_5$
 $i_3 \rightarrow i_2, i_5$
 $i_5 \rightarrow i_2, i_3$
 $i_2, i_3 \rightarrow i_5$
 $i_2, i_5 \rightarrow i_3$
 $i_3, i_5 \rightarrow i_2$

$2/4 = 50\% \quad$ finally we
 $2/3 = 66\% \quad X \quad$ have to
 $2/2 = 100\% \quad$ write down
 $2/3 = 66\% \quad X \quad$ all strong
 $2/4 = 50\% \quad X \quad$ rules.
 $2/4 = 50\% \quad X \quad$ which is not
 $2/2 = 100\% \quad$ crossed.

III IA.

12/2/24.

d) $I = \{A, B, C, D, E\}$ $T = \{T_1, T_2, T_3, T_4, T_5, T_6, T_7, T_8, T_9\}$

T_1	A, B
T_2	B, D
T_3	B, C
T_4	A, B, D
T_5	A, C
T_6	B, C
T_7	A, C
T_8	A, B, C, E
T_9	A, B, C

min sup = 2 9% min support
 min conf = 50% is given in
 percentage

ex:- 60%

$$\text{min support} = \frac{60}{100} \times 9$$

$$\frac{54}{10} = 5.4$$

no. of transactions

1-itemset | support count

T_1	A	6
T_2	B	7
T_3	C	6
T_4	D	2
T_5	E	1 X
T_6		
T_7		
T_8		
T_9		

2-itemset	SC	3-itemset	SUPP
A B	4	ABC	2
A C	4		
A D	1 X	4-itemset	
B C	4		
B D	2		
C D	0 X		

Frequent itemset Supp

A	6
B	7
C	6
D	2
AB	4
AC	4
BC	4
BD	2
ABC	2

Association Rules. Confidence.

$A \rightarrow B$	$4/6 = 66.67\%$	$D \rightarrow B$	$2/2 = 100\%$
$B \rightarrow A$	$4/7 = 57.14\%$		
$A \rightarrow C$	$4/6 = 66.67$		
$C \rightarrow A$	$4/6 = 66.67$		
$B \rightarrow C$	$4/7 = 57.14\%$		
$C \rightarrow B$	$4/6 = 66.67$		
$B \rightarrow D$	$2/7 = 28.57\%$		
$A \rightarrow B, C$	$2/6 = 33.33\%$		
$B \rightarrow A, C$	$2/7 = 28\%$		
$C \rightarrow A, B$	$2/6 = 33\%$		
$A, B \rightarrow C$	$2/4 = 50\%$		
$B, C \rightarrow A$	$2/4 = 50\%$		
$A, C \rightarrow B$	$2/4 = 50\%$		

minconf = 60%

T₁ Bread, Buns, Sauce min sup = 33.33%T₂ Bread, Buns $\frac{33.33}{100} \times 6 = 1.99$ T₃ Bread, Coke, ChipsT₄ Chips, CokeT₅ Chips, SauceT₆ Bread, Coke, Chips $= \underline{\underline{2}}$

1-itemset Supp C

Bread 4

Buns 2

Sauce 2

Coke 3

Chips 4

3-itemset Supp C

Bo, Co, Ch 2

frequent items Supp C

Bo 4

Bu 2

Sa 2

2-itemset Supp C

Br, Bu 2

Br, Sa 1 X

Br, Co 2

Br, Ch 2

Bu, Sa 1 X

Bu, Co 0 X

Bu, Ch 0 X

Sa, Co 0 X

Sa, Ch 0 X

Co, Ch 2

Bo 3

Ch 4

Br, Bu 2

Br, Co 2

Br, Ch 2

Co, Ch 3

Br, Co, Ch 2

Association Rules

Br \rightarrow Bu $2/4 = 50\% \times 60\% = 30\%$ Bu \rightarrow Br $2/2 = 100\% \times 60\% = 60\%$ Br \rightarrow Co $2/4 = 50\% \times 75\% = 37.5\%$ Co \rightarrow Br $2/3 = 66.6\% \times 50\% = 33\%$ Ch \rightarrow Br $2/4 = 50\% \times 60\% = 30\%$ Br \rightarrow Ch $2/4 = 50\% \times 50\% = 25\%$

$$\begin{array}{l} \text{O Ch} \rightarrow \text{Br} \quad 2/3 = 66.7 \\ \text{Br} \rightarrow \text{Ch} \quad 2/2 = 100\% \\ \text{Br Ch} \rightarrow \text{O} \quad 2/3 = 66.7 \end{array}$$

13/2/24

Apriori Algorithm.

N = No. of transactions.
 F = Frequent itemsets.
 K = no. of items in itemsets.

1. $K=1$
 2. $F_K = \{i \mid i \in I \wedge \sigma(\{i\}) > N \times \text{minisupport}\}$ // Generate all itemsets, $K=1, 2, 3$
 3. Repeat.
 4. $|C|=K+1$
 5. C_K = Candidate generation (F_{K-1})
 6. C_K = Candidate Prune (C_K, F_{K-1})
 7. for each transaction $t \in T$ do
 8. C_K = subset (C_K, t)
 9. for each candidate itemset $c \in C_K$ do
 10. $\sigma(c) = \sigma(c) + 1$
 11. end for
 12. end for
 13. $F_K = \{c \mid c \in C_K \wedge \sigma(c) > N \times \text{minisupport}\}$
 14. until $F_K = \emptyset$
 15. Result UF_K
- $U = \text{Union}$.

Generation of freq. itemsets.

two imp. process: \rightarrow ① Candidate Generation

to generate frequent itemsets $\xrightarrow{\text{methods}} \rightarrow$ ② Candidate Pruning -
 ① putting every possibility of itemsets are Candidate Itemset.
 All Cond. I. ~~is~~ may be or not frequent Itemset.

1 - Itemset

$$\begin{cases} A & 3 - FI \\ B & 0 \\ C & 1 \end{cases}$$

Candidate Itemsets. $\begin{cases} B & 0 \\ C & 1 \end{cases}$ this C_K of $K=2$ is given F_1 .

$$\begin{cases} P & 4 - FI \\ G & 3 - FI \end{cases}$$

② pruning is done based on support pruning.

① Candidate Generation

Candidate Generation Procedure
it should follow both types to generate CG, I

Complete \Rightarrow it should not omit any possibilities.

Non-R \Rightarrow Any itemset should not occur more than once.

Different ways we can generate CL Itemsets.

1) Brute force method.

Candidate C

1-itemset - 2-IS.

A AB, A.C

B AD, A.E

C BC, BD, BE

D CD, CE

E DE

* Disadvantage! - Pruning is expensive which takes lot of time

K = level of itemsets.

2) $f_{k-1} \times f_1$ Itemsets.

Take Bread, Butter example (Refer last page.)
consider we have to generate 3 itemset. we can do this in $f_2 \times f_1$

f_2 f_1 3-itemset

Bread Bun

Bread Coke

Bread Chips

Chips Coke

Bread

Buns

sauce

coke

chips

Redundancy

Bread, Bun, Sauce

Bread, Bun, Coke

Bread, Buns, Chips

Bread, Coke, Sauce

Bread, Chips, Buns

Dis :- Candidate itemset generated are not non-Redundant in 2)

* Items should be arranged in lexicographical order

Bread chips Coke Bun

hence we can remove Redundant itemsets.

e.g. Bread chips combining with Bun we cannot proceed cause bun < chips in alphabetical order.

3) $F_{k-1} \times F_{k-1}$ to form 3-itemsets.

e.g.: - { Bread, Sauce }

{ Bread Jam }

{ Bread Milk }

{ Jam, Milk, Sugar }

{ Jam, Milk }

1-item

should

match

when
 $F_{k-1} \times$

whenever we are considering this at least F_{k-1} , which 1st item should match.