



## ML 3rd unit - ml unit 3 notes

R18 B.Tech. Cse (Computer Networks) Iii & Iv Year Jntu Hyderabad (Jawaharlal Nehru  
Technological University, Hyderabad)

### UNIT-III

#### : BAYESIAN LEARNING :

It is based on the assumption that the quantities of interest are governed by Probability Distributions & the Optimal decisions can be made by reasoning about these probabilities together with observed data.

Bayesian Learning methods are similar to our Machine Learning for two different reasons:

1. These calculate explicit probabilities for hypothesis such as Naive Bayes Classifier.

2. Comparing the classifier with decision Tree and Neural Network Algorithms.

The Problem of learning to classify text documents such as electronic news articles., the Naive Bayes classifier is most effective algorithms.

2. They provide a useful perspective for understanding many learning algorithms that do not explicitly manipulate probabilities.

→ We analyze the algorithms such as FIND-S and CANDIDATE ELIMINATION ALGORITHM to determine conditions under which they output the most probable hypothesis given by Training Data.

## Features Of Bayesian Learning Methods :

1. Each Observed training example can incrementally decrease or increase the estimated probability that a hypothesis is correct.  
It also provides the complete elimination of hypothesis if it is inconsistent with any single example.
2. Prior knowledge can be combined with observed data to determine the final probability of hypothesis.
3. Bayesian methods can accommodate hypotheses that make probabilistic predictions.
4. New instances can be classified by combining the predictions of multiple hypotheses, weighted by their probabilities.
5. Bayesian methods prove computationally intractable, they can provide a standard of optimal decision making against which other practical methods can be measured.

11<sup>th</sup> May, Presubies

50, 55, 57, 61, 62, 63, 64  
71, 66, 72, 73, 75, 80, 82, 83, 81  
87, 89, 92  
510

10<sup>th</sup> May, presubies

50, 55, 57, 61, 62, 63, 66, 70, 73  
75, 82, 89, 92  
LE-502, 504, 505, 507, 512, 513  
72, 75

(30)

## BAYE'S THEOREM

In ML we are often interested in determining the best hypothesis from Space  $H$ , on observed Training data  $D$ . To specify the best hypothesis, we demand the most Probable hypothesis.

Bayes Theorem provides a way to calculate the probability of hypothesis directly based on its prior probability.

Based on previous knowledge of System

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

$\Rightarrow P(A|B)$  is a hypothesis, finding the probability of  $A$ , where  $B$  is occurred.

$\Rightarrow P(B|A)$  finding the probability of  $B$ , where  $A$  is occurred

$\hookrightarrow$  Likelihood  
 $\Rightarrow P(A)$  is Prior  $\Rightarrow P(B)$  Marginal.

$$P(A|B) = \frac{P(A \cap B)}{P(B)} - (*) - P(B|A) = \frac{P(B \cap A)}{P(A)}$$

$\Rightarrow P(A|B) \cdot P(B) = P(A \cap B) \quad \text{--- } ①$

$P(B|A) \cdot P(A) = P(B \cap A) \quad \text{--- } ②$

From eq. ① & ② RHS are equal i.e.,  $A \cap B = B \cap A$

then  $P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$

i.e., 
$$\boxed{P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}} \quad \text{Bayes Theorem}$$

Terms ! A is Hypotheses  
B is Given Data

$p(A|B)$  - finding Probability of Hypothesis, when prob of Training Examples given.

$p(B|A)$  - finding Probability of given data provided with prob of hypotheses that is True.

$p(A)$  - Prob of hypothesis before observing given data.

$p(B)$  - Prob of hypothesis of given data is True.

Example:  $P(\text{King} / \text{face})$

Probability of finding King that it is a face card.  
Face Cards  $\Rightarrow J, Q, K$  i.e.,  $4 \times 3 = 12$  face cards in single set of pack.

$$= \frac{P(\text{face/King}) \cdot P(\text{King})}{P(\text{face})}$$

$\frac{4 \text{ Kings}}{4 \text{ Colors}} = \frac{4}{4} = 1$

$$= \frac{1 \cdot \frac{4}{52}}{\frac{12}{52}}$$

[ 4 Kings among 52 cards  
12 faces among 52 cards ]

$$= \frac{1 \times \frac{1}{13}}{\frac{3}{13}} = \frac{1}{3} \Rightarrow P(K/f) = \frac{1}{3}$$

1 face among 3 cards  
is probability.

Baye: It will give Probability of an event based on previous Conditions.

## Concepts of Bayes Theorem :

### 1. Product Rule :

Probability  $P(A \cap B)$  of a conjunction of two events

$A \& B$

$$P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

### 2. Sum Rule :

Probability of disjunction of two events A and B.

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

### 3. Bayes Theorem :

The Prior Probability

$P(h|D)$  of h given D.

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

### 4. Theorem of Total Probability:

If events  $A_1, \dots, A_n$  are mutually exclusive with

$$\sum_{i=1}^n P(A_i) = 1, \quad P(B) = \sum_{i=1}^n P(B|A_i)P(A_i)$$

## MAP Hypothesis [Maximum A Posterior]

Method : Learner Considers some set of hypothesis  $H$  and is interested in finding the most probable Hypothesis  $h \in H$ .

Given the Observed data  $D$  [Atleast one of the marginally probable if there are many]

Such maximally probable hypothesis is called a maximum A posterior [MAP]

Determine : The MAP hypothesis by using Bayes Theorem.

$$h_{\text{map}} = \underset{h \in H}{\operatorname{argmax}} p(h/D)$$

$$= \underset{h \in H}{\operatorname{argmax}} \frac{p(D/h)(p(h))}{p(D)}$$

[ $\because p(D)$  is common for all hypothesis we ignore]

$$\boxed{\therefore h_{\text{map}} = \underset{h \in H}{\operatorname{argmax}} p(D/h)(p(h))}$$

## BRUTE FORCE MAP LEARNING :

The Simplifications does not alter the main conclusions of this method. We can design a straight-forward Concept Learning Algorithm to output the maximum Posterior Hypothesis.

### Brute-Force map Algorithm :

i. For each hypothesis  $h$  in  $H$ , calculate the Posterior Probability

$$P(h|D) = \frac{P(D|h) P(h)}{P(D)}$$

ii) Output the hypothesis  $h_{MAP}$  with the highest posterior Probability.

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

This may require Significant Computation because it applies Bayes Theorem to each hypothesis in  $H$  to calculate  $P(h|D)$ .

- \* It is impractical for large hypothesis Space

- \* Still the method is interested because it provides a standard against which we may judge the performance of other Concept learning algorithms.

31<sup>st</sup> May: Presenties  
549, 60, 64, 74, 75, 76, 78  
82, 83, 88, 89

→ Brute Force MAP Learning algorithm must specify values for  $p(h)$  and  $p(D|h)$

→  $p(h)$  and  $p(D|h)$  can be chosen to be consistent with the assumptions.

i. Training Data  $D$  is Noise free.

ii. The Target Concept  $C$  is contained in hypothesis space  $H$ .

iii. We have no prior reason to believe that any hypothesis is more probable than any other.

Given these assumptions, what values should we specify for  $p(h)$ .

→ Prior Probability can be given as :

$$p(h) = \frac{1}{|H|} \text{ for all } h \text{ in } H.$$

→ what choice shall we make for  $p(D|h)$

$$p(D|h) = \begin{cases} 1 & \text{if } d_i = h(x_i) \text{ for all } d_i \text{ in } D \\ 0 & \text{otherwise} \end{cases}$$

$\left[ \because d_i \text{ is Target } \& h(x_i) \text{ is Output} \right]$

The Probability of data  $D$  given hypothesis  $h$  is 1, if  $D$  is consistent with  $h$  otherwise 0.

Recalling Bayes Theorem :

$$P(h|D) = \frac{P(D|h) P(h)}{P(D)}$$

where  $h$  is inconsistent with Training Data  $D$ .

i.e.,  $P(D|h)$  to be 0.

$$\therefore P(h|D) = \frac{0 \cdot P(h)}{P(D)} = 0 \text{ if } h \text{ is inconsistent with } D.$$

where  $h$  is consistent with Training data  $D$

i.e.,  $P(D|h)$  to be 1.

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

$$= \frac{1 \cdot \frac{1}{|H|}}{\frac{|VS_{H,D}|}{|H|}} \quad [VS - \text{Version Space}]$$

$$P(h|D) = \frac{1}{|VS_{H,D}|} \text{ if } h \text{ is consistent with } D.$$

where  $VS_{H,D}$  is the subset of hypothesis from

$H$  that are consistent with  $D$ .

$$P(D) = \sum_{h_i \in H} P(D/h_i) P(h_i)$$

$$= \sum_{h_i \in VS_{H,D}} 1 \cdot \frac{1}{|H|} + \sum_{h_i \notin VS_{H,D}} 0 \cdot \frac{1}{|H|}$$

$$= \sum_{h_i \in VS_{H,D}} 1 \cdot \frac{1}{|H|} \quad [VS - Version Space]$$

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

To summarize, Bayes theorem implies that the posterior probability  $P(h/D)$  under our assumed  $p(h)$  and  $p(D/h)$  is

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

$$\frac{1}{|H|}$$

$$|VS_{H,D}|$$

$$|H|$$

$$|VS_{H,D}|$$

$$|H|$$

## Least Squared Error Hypothesis :

Bayesian analysis can sometimes be used to show that a particular learning algorithm outputs MAP hypotheses even though it may not explicitly use Bayes Rule.

→ This method considers the problem of learning a continuous valued target function.

→ Under certain assumptions any learning algorithm that minimizes the squared error hypothesis between the output hypothesis predictions and training data will output a maximum likelihood hypothesis.

Finding maximum likelihood hypothesis in Bayesian learning

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

[Maximum Likelihood]

[small  $p$  : Probability Density function]

Let us take training instances  
 $x_1, x_2, \dots, x_n$  and target values

$$D = \{d_1, d_2, \dots, d_m\}$$

∴ we write  $p(D|h)$  as product of  $p(d_i|h)$  as various probabilities

$$\text{i.e., } h_{MAP} = \underset{h}{\operatorname{argmax}} \prod_{i=1}^m p(d_i|h)$$

[Probabilities must be multiplied]

By assuming Normal Distribution for probability distribution input  $x \rightarrow$  variable.

$\mu \rightarrow$  mean &  $\sigma \rightarrow$  Standard Deviation.

$$f(x|\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

According to Bayes,  $x$  with  $d_i$  &  $\mu$  with  $h(x_i)$

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

Use in function : when we apply log value, product  
Converts to summation i.e.,  $\prod \Rightarrow \sum$

$$\text{i.e. } h_{ML} = \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(d_i - h(x_i))^2$$

Maximizing this negative quantity is equivalent to  
minimizing the corresponding positive quantity.

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

we can discard constants, which are common to  
all data sets i.e.,  $\frac{1}{\sqrt{2\pi\sigma^2}}$

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

↓                      ↓  
Least                  Error minimised &  
                          Squared.

- The equation shows that the maximum likelihood hypothesis  $h_{ML}$ , which minimizes the sum of squared errors between observed training values & hypothesis prediction.
- It determines the most probable hypothesis.
- This approach is used in many neural networks like curve fitting.
- The analysis considers noise only in the target value of Training example.

Maximum Likelihood hypothesis for Predicting Probabilities

Collect Observed frequencies of 1's & 0's for each possible value of  $x$ , train neural network to update target frequency. Which Criterion Should we optimize in order to find a maximum likelihood hypothesis for  $f$ ? ( $f$ -dash)

- \* first obtain an expression for  $P(D|h)$
- \* Assume the training data  $D$  is of the form

$$D = \{(x_1, d_1), (x_2, d_2) \dots (x_m, d_m)\}, \text{ where } d_i \text{ is}$$

observed 0 or 1 for value for  $f(x_i)$ .

- \* both  $x_i$  &  $d_i$  as random variables and assuming that each training examples is drawn independently, we can

write  $P(D|h)$  as

$$P(D|h) = \prod_{i=1}^m P(x_i, d_i|h) \quad \text{--- (1)}$$

Applying Product Rule

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

The Probability  $p(d_i/h, x_i)$

$$p(d_i/h, x_i) = \begin{cases} h(x_i) & \text{if } d_i = 1 \\ (1-h(x_i)) & \text{if } d_i = 0 \end{cases} \quad -\textcircled{3}$$

Re-express it in a more mathematically manipulate form as:

$$p(d_i/h, x_i) = h(x_i)^{d_i} (1-h(x_i))^{1-d_i} \quad -\textcircled{4}$$

Substitute eq-\textcircled{4} in eq-\textcircled{1} to obtain eq-\textcircled{5}

$$p(D/h) = \prod_{i=1}^m h(x_i)^{d_i} (1-h(x_i))^{1-d_i} p(x_i) \quad -\textcircled{5}$$

$\therefore d_i, h$  are posterior examples.

we write an expression for maximum likelihood hypothesis.

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

The last term is a constant independent of  $h$ , so it can be dropped.

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

It is easier to work with log of likelihood yielding

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

$[\because \ln \text{ is length}]$

eq-\textcircled{7} describes the quantity that must be maximised in order to obtain the maximum likelihood hypotheses in our current problem setting.

## Minimum Description Length Principle:

Assumption :

Representing a Concept in minimum possible way, then it is said to be good concept.

Mathematically :

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

Applying logarithm:

$$\operatorname{argmax}_{h \in H} \log P(D|h) + \log P(h)$$

$$[\because \log ab = \log a + \log b]$$

$\therefore$  Convert it into min to get shortest

$$\operatorname{argmin}_{h \in H} -\log P(D|h) - \log(P(h))$$

[minimum length / short hypothesis is preferred]

Example :

\* Let us consider a probability of designing a code to transmit messages drawn at random from a set D. where probability of drawing an  $i^{th}$  msg =  $P_i$ ;

\* while transmitting, we need the code that minimize the expected number of bits.

To do this, we should assign shorter codes to the most probable.

We represent length of message  $i$  we write to see as  $L_C(i)$

$$\therefore h_{MAP} = \operatorname{argmin} L_{CH}(h) = L_{CDH}(D|h) \quad [i \text{ is message}]$$

$L_{CH}$  is optimal encoding for  $H$  given  $D$  given  $h$

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

$$\therefore h_{MDL} = \underset{h \in H}{\operatorname{argmin}} L_{C_1}(h) + L_{C_2}(D|h)$$

The above analysis shows that if we choose  $C_1$  to be the optimal encoding of hypotheses  $C_H$ , and if we choose  $C_2$  to be optimal encoding  $C_D|h$  then

$$h_{MDL} = h_{MAP}$$

### BAYES OPTIMAL CLASSIFIER :

Bayes Optimal Classifier is a probabilistic model that makes the most probable prediction for a new example:

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

for a dataset

$$x = \{x_1, x_2, x_3, \dots, x_n\} \quad \{y\}$$

$$p(y|x_1, x_2, \dots, x_n) = \frac{[p(x_1|y) \cdot p(x_2|y) \cdots p(x_n|y)] p(y)}{p(x_1) \cdot p(x_2) \cdots p(x_n)}$$

$$= \frac{p(y) \prod_{i=1}^n p(x_i|y)}{p(x_1) \cdot p(x_2) \cdots p(x_n)}$$

$$= p(y) \prod_{i=1}^n p(x_i|y)$$

[Denominator is Omitted  
it doesn't make change  
of probability]

CSE - III - II 6<sup>th</sup> Jun

Reserves: 50, 55, 57, 58, 60  
61, 63, 64, 66, 72, 74, 75, 82  
83, 87.  
LE: 504, 505, 506

(37)

Ex: Enjoy Sport with 14 days info  
We took only two attributes among all.

## \* Outlook

	Yes	No	$P(Y)$	$P(N)$
Sunny	2	3	2/9	3/5
Outcast	4	0	4/9	0/5
Rainy	3	2	3/9	2/5
Total	9	5	100%	100%

$p(Y) = \text{Probability of Yes}$   
 $p(N) = \text{Probability of No}$

## \* Temperature

	Yes	No	$P(Y)$	$P(N)$
Hot	2	2	2/9	2/5
Mild	4	2	4/9	2/5
Cold	3	1	3/9	1/5
Total	9	5	100%	100%

## \* Play

Yes	9	9/14
No	5	5/14
Total	14	100%

Total (Sunny, Hot)

$$p(\text{Yes} | \text{Sunny, Hot}) = p(\text{Sunny} | \text{Yes}) \times p(\text{hot} | \text{Yes}) \times p(\text{Yes})$$

$$= \frac{2}{9} \times \frac{2}{9} \times \frac{9}{14} = 0.031$$

$$p(\text{No} | \text{Sunny, hot}) = p(\text{Sunny} | \text{No}) \times p(\text{hot} | \text{No}) \times p(\text{No})$$

$$= \frac{3}{5} \times \frac{2}{5} \times \frac{5}{14} = 0.08571$$

$$\text{Total} = 0.031 + 0.08571$$

$$= 0.27$$

$$P(\text{Yes}) = \frac{0.031}{0.27} = 0.114$$

$$P(\text{No}) = \frac{0.0857}{0.27} = 0.317.$$

Probability is No is more  
∴ player will not enjoy sport.

### : GIBBS ALGORITHM :

- \* chooses one hypothesis at Random, according to  $p(h/D)$
- \* Use this classify new instance.
- \* ~~Uses~~ Error frequency of Gibbs is two times less than Bayesian classifier.
- \* Simply applies a hypothesis drawn at random according to current posterior probability Distribution.

(38)

## BAYESIAN BELIEF NETWORKS :

Before going to learn about Bayesian Belief Networks we come across about:

- \* Directed Acyclic Graph. [DAG]
- \* Conditional Probability Table. [CPT]

DAG :

Rain

when rainy  
dog barks

Dog bark

when dog barks  
cat hides

Cat hide

A directed graph with no cycles.

As Rain falls, Dog barks if dog barks cat hides.

Two probabilities for above example for all nodes i.e., may or may not.

The Conditional Probability Table :- Ex: 1

	R	$\sim R$
B	9/48	18/48
$\sim B$	3/48	18/48

R - Rains &  $\sim R$  - Not Rains  
B - Barks &  $\sim B$  - Not Barks

Rain is a node with two hypothesis true or false  $\Rightarrow 2$   
Total we consider 48 events

$$B=T \& R=T \Rightarrow 9/48 = 0.19$$

$$B=T \& R=F \Rightarrow 18/48 = 0.375$$

$$B=F \& R=T \Rightarrow 3/48 = 0.06$$

$$B=F \& R=F \Rightarrow 18/48 = 0.375$$

$\Rightarrow$  Rain is Parent of Dog bark  
Dog bark is parent of Cat hide

$\Rightarrow$  Calculating the probability, can be done by parent i.e., Using Parent node not child. So, we haven't used the concept of Cat-hide.

(81)

64      9th: 50, 51, 54, 61, 70  
66      72, 73, 74, 75, 77, 78  
58      82, 83, 84, 86, 87, 89, 90

②  
49, 60, 61, 64  
71, 82, 87, 89  
91, 502, 505, 506

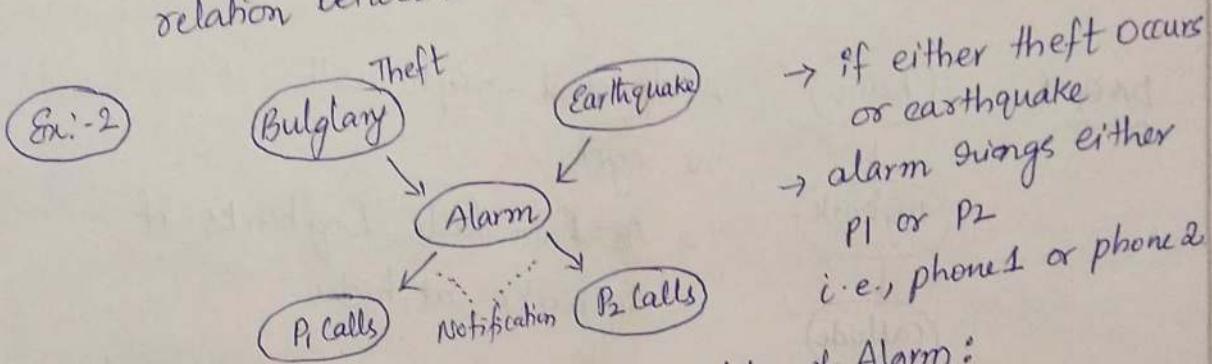
(3rd)

7th June: 49, 50, 55, 60, 61, 63  
64, 70, 71, 72, 75, 78, 81, 83, 87  
89, 91, LE-502, 504, 505, 506  
510, 512

## Bayesian Belief Networks :

→ It is a probabilistic Graphical model (PGM) that represents conditional dependencies between random variables through DAG.

→ It is also suitable for representing probabilistic relation between multiple events. [more than Two]



Given Probabilities are:

$$P(B=T) = 0.001$$

$$P(B=F) = 0.999$$

$$P(E=T) = 0.002$$

$$P(E=F) = 0.998$$

: Probability of Alarm:

	Burglary	Earthquake	P(A=T)	P(A=F)
T	T	0.95	0.05	
T	F	0.99	0.06	
F	T	0.29	0.71	
F	F	0.001	0.999	

: Probability of P<sub>2</sub>

A	P <sub>2</sub> =T	P <sub>2</sub> =F
T	0.70	0.30
F	0.01	0.99

Probability of P<sub>1</sub>

Alarm (A)	P(P <sub>1</sub> =T)	P(P <sub>1</sub> =F)
T	0.90	0.10
F	0.05	0.95

Ex:- Find the Probability of P<sub>1</sub> is T, P<sub>2</sub> is T & A is T, B is F & E is F

$$\text{i.e., } P(P_1, P_2, A, \neg B, \neg E)$$

$$= P(P_1|A) P(P_2|A) P(A) P(\neg B, \neg E) \\ = P(A|P_1, P_2, \neg B, \neg E) \cdot P(P_1) \cdot P(P_2) \cdot P(\neg B) \cdot P(\neg E)$$

$$= 0.90 \times 0.70 \times 0.001 \times 0.999 \times 0.998$$

$$= 0.00062$$

No Burglary, No Earthquake then ringing of alarm is very little possibility.

## EM-ALGORITHM

### Expectation - Maximization

→ Used to find Latent variable.

→ This is Basic for many unsupervised clustering Algorithm.

The variables that directly <sup>cannot</sup> derived or observed are called as Latent Variables. These are observed with old or past derived variables.

### STEPS:

1. Initially, set of values initially considered.

A set of incomplete data is given to system.

2. Estimation or Expecting value from

incomplete set of data [fill the incomplete data]

3. Maximization or M-Step.

Here, we use the complete generated data in preceding e-step to update the value.

4. We check if values are converging or not i.e., we are getting expected value.

If not repeat the Steps ③ & ④ i.e., E & M.

### USAGE:

1. Used to fill missing data.

2. Used for unsupervised Clustering.

3. Used to discover values of latent variables.

### Advantages:

1. with each iteration, likelihood increases.

2. E-step and m-step are easy to implement.

### Disadvantages:

1. slow convergence.
2. makes convergence to Local optimal only but Cannot be done by global optimal.

### (CHAP 7) COMPUTATIONAL LEARNING THEORY:

The Theoretical characterization of difficulties of several types of machine learning problems and the capabilities of several ML Algorithms.

Answers the questions like :

→ "Under what Conditions is a particular Algorithm is possible & Impossible"

→ "Under what Conditions a particular Algorithm is assured of learning successfully"

Two Specific frameworks for analyzing learning algorithms are

Considered by :-

⇒ PAC : Probably Approximately Correct.

Depending on attributes of learning problem such as :

→ the size or complexity of hypothesis space considered

by Learner.

→ the accuracy to which the target concept must be approximated.

→ the probability that the learner will output a successful hypothesis.

5F 9th June:  
50, 55, 58, 60, 63, 64, 69  
70, 72, 73, 74, 75, 76, 77, 78  
84, 86, 87, 88, 89  
15, 502, 505, 506, 507, 510

Our Goal is to answer such as :

Sample Complexity : How many training examples are needed for a learner to converge for a successful hypothesis.

Computational Complexity : How much computational effort is needed for a learner for a successful hypothesis.

Mistake Bound : How many training examples will the learner misclassify before converging a successful Hypothesis.

Probably Learning An Approximately Correct Hypothesis

Consider a particular setting for learning problem called PAC [Probably Approximately Correct] model.

- Specify the Problem Setting that define PAC.
- Consider how many training examples.
- How much computation is required.

1. Problem Setting:

- Let  $X$  is set of all possible instances over which target function may be defined.
- Let  $C$  be set of target Concepts that our Learner might be called upon to learn.
- Each target concept  $c$  in  $C$  corresponds to some subset of  $X$

i.e.,  $c : X \rightarrow \{0, 1\}$

if  $x$  is positive example then  $c(x) = 1$ .

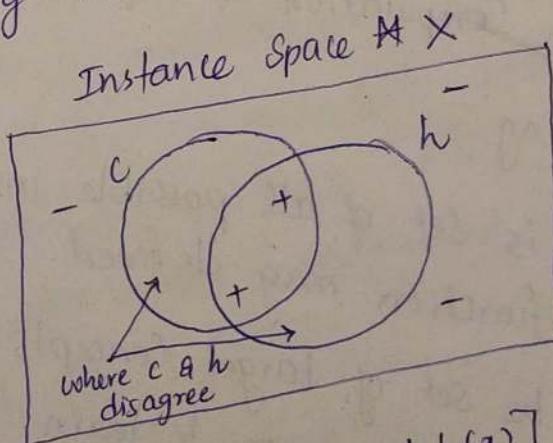
if  $x$  is negative example then  $c(x) = 0$ .

The instances we assume are generated at random from  $X$  according to some probability distribution  $D$ , it may be any distribution, generally not known to learner. Training examples are generated by drawing an instance  $x_i$ , along with its target value  $c(x_i)$ .

This setting is interested in characterizing the performance of various learners  $L$  using various hypothesis spaces  $H$ .

Error Of Hypothesis: We define true error of a hypothesis  $h$  with respect to target concept  $c$  and instance distribution  $D$ .

True Error: denoted by  $\text{error}_D(h)$ , the error which will misclassify an instance drawn at random according to  $D$ .



$$\text{error}_D(h) = \Pr_{x \in D} [c(x) \neq h(x)]$$

$\therefore \Pr$  indicates probability taken over instance distribution  $D$ .

The randomly drawn instance will fall into the region where  $h$  and  $c$  disagree on its classification.

The  $+/-$  indicates positive and negative examples.

We use Training Error to refer to fraction of training examples misclassified by  $h$ .

### PAC Learnability:

PAC: Probably Approximately Correct.

It characterizes the classes of target concepts that can be reliably learned from randomly drawn training examples and reasonable amount of computation.

→ characterizing number of training examples needed to learn a hypothesis  $h$  for which  $\text{error}_D(h) = 0$ .

Definition: Consider a concept class  $C$  defined over a set of instances  $X$  of length  $n$  and a learner  $L$  using hypothesis space  $H$ .

$C$  is PAC-learnable by  $L$  using  $H$ , if for all  $c \in C$ , distributions  $D$  over  $X$ , such that  $0 < \epsilon < 1/2$  and  $\delta$  such that  $0 < \delta < 1/2$ , Learner  $L$  with probability atleast  $(1 - \delta)$  output a hypothesis  $h \in H$  such that  $\text{error}_D(h) \leq \epsilon$ , in time that is polynomial in  $1/\epsilon, 1/\delta$ .

14<sup>th</sup> presentation: (30)

⑩ 54, 55, 58, 60, ⑪ 63, 64, 66, 69  
 ⑫ 71, 72, 73, 74, 75, ⑬ 76, 77, 78, 82, 81  
 83, 84, 87, ⑭ 89,

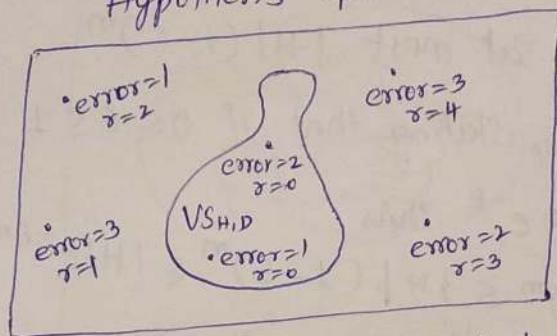
## SAMPLE COMPLEXITY FOR FINITE HYPOTHESIS SPACES

PAC-learnability is largely determined by the number of training examples required by the learner.

- The increase in number of required number of training examples with problem size, called Sample Complexity.
- The success of learner is limited availability of training examples.
- A general bound on sample complexity for a very broad class of learners called Consistent learners.
- A learner is consistent if output hypothesis perfectly fits the training data.
- Derive a bound on Training Examples required by any consistent learner.
- To accomplish we recall Version Space  $V_{S,H,D}$   
i.e.,  $V_{S,H,D} = \{h \in H \mid (\forall (x, c(x)) \in D) (h(x) = c(x))\}$

The significance of version space is that every consistent learner outputs a hypothesis belonging to version space.  
Instance Space  $X$ , hypothesis  $H$ , Training data  $D$ .

Def: Consider a hypothesis space  $H$ , target concept  $c$ , instance distribution  $D$ , training examples  $D$  of  $c$ . The Version Space is said to be  $\epsilon$ -exhausted with respect to  $c$  &  $D$ .  
 $V_{S,H,D}$  has an error less than  $\epsilon$  with respect to  $c$  and  $D$ .  
 $(\forall h \in V_{S,H,D}) \text{error}_D(h) < \epsilon$

Hypothesis Space  $H$ 

$$\text{error}(h) = \Pr_{x \in D} [c(x) \neq h(x)]$$

The Version Space is  $\epsilon$ -exhausted just in the case that all the hypothesis consistent with the observed Training examples have zero training error  $r=0$

Theorem:  $\epsilon$ -exhausting the version Space.

If the hypothesis space  $H$  is finite and is  $D$  is sequence of  $m \geq 1$  Independent Randomly drawn examples of some target Concept  $C$ , for any  $0 \leq \epsilon \leq 1$ .

$$K \leq |H|e^{-\epsilon m}$$

Proof: Let  $h_1, h_2, h_3, \dots, h_K$  all hypothesis in  $H$ . we have  
 $\Rightarrow$  we fail to  $\epsilon$ -exhaust the Version Space if and only if atleast one of these  $K$ -hypotheses happens to be considered consistent with  $m$  independent training examples.  
 $\Rightarrow$  the probability any single hypothesis having true error greaterthan  $\epsilon$  is consistent with  $(1-\epsilon)$ .  
 $\therefore (1-\epsilon)^m$  probability of drawing  $m$  independent examples.

$$K(1-\epsilon)^m$$

$\therefore K \leq |H|$  is at most  $|H|(1-\epsilon)^m$ .

General Inequality stating that if  $0 \leq \epsilon \leq 1$  then

$$(1-\epsilon) \leq e^{-\epsilon} \text{ Thus}$$

$$K(1-\epsilon)^m \leq |H|(1-\epsilon)^m \leq |H|e^{-\epsilon m}.$$

proves the theorem.

### Sample Complexity for Infinite Hypothesis Space

As the sample complexity for PAC learning grows as the algorithm of the size of hypothesis space.

There are two drawbacks for finite hypothesis space

1. It can lead to quite weak bounds
2. Cannot apply the equation for infinite bounds.

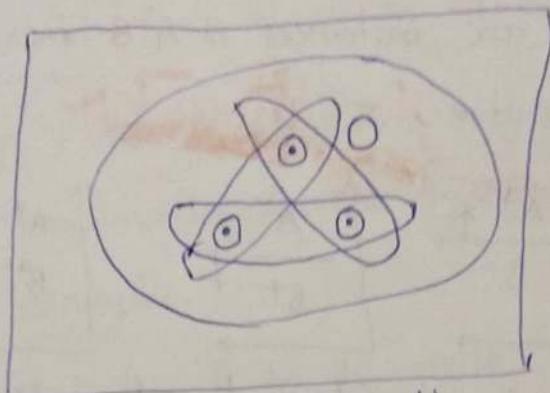
### Shattering a Set of Instances:

We measure the complexity of  $H$  called VC - Vapnik-Chervonenkis dimension of  $H$

i.e.,  $VC(H)$ .

VC dimension measures the complexity of hypothesis space  $H$ , not by number of distinct hypothesis  $|H|$ , but instead by number of distinct instances from  $X$ .

Def: A set of instances  $S$  is shattered by Hypothesis space  $H$ , if and only if for every dichotomy of  $S$  there exists some hypothesis in  $H$  consistent with this dichotomy.

Instance Space  $H$ 

A set of three Instances shattered by eight hypothesis,  
for every possible dichotomy of the instances, there exists a  
Corresponding hypothesis.

NOTE : If a set of instances is not shattered by a hypothesis space, then there must be some concept that can be defined over the instances, but cannot be represented by hypothesis space.

VAPNIK CHERVONENKIS DIMENSION :

The ability to shatter a set of instances is closely related to inductive bias of a hypothesis space.

Def: The  $VC(H)$ , of hypothesis space  $H$  defined over instance space  $X$  is the size of largest finite subset of  $X$  shattered by  $H$ .

If large finite sets of  $X$  can be shattered by  $H$ , then  $VC(H) = \infty$ .

Note that for any finite  $H$ ,  $VC(H)$  consider few example

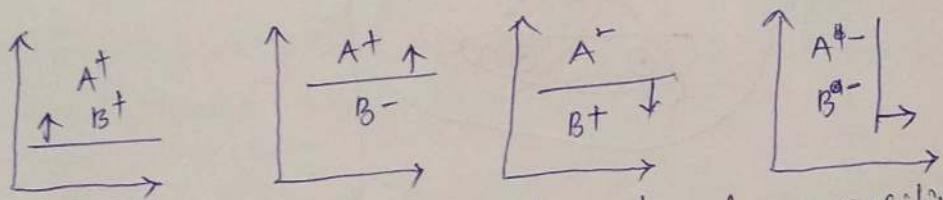
hypothesis.  $VC(H) \leq \log_2 |H|$

i.e.,  $VC(H) = d$ , then  $H$  will require

$2^d$  distinct hypothesis to shatter  $d$  instances.  $2^d \leq |H| \Leftrightarrow$

$$d = VC(H) \leq \log_2 |H|$$

Ex:- if we consider two instances A & B i.e.,  $2^2$  distinct hypothesis can be done i.e.,



Shattering can be done based on positivity & negative i.e., probability of an instance.

Ex:- if three instances then  $2^3$  distinct hypothesis are taken i.e., 8 hypothesis and shattered according to their probability.

### MISTAKE BOUND MODEL FOR LEARNING :

- We consider the number of mistakes the learner make before it learns the correct hypothesis.
- Learner is evaluated by the number of mistakes it makes.
- Set of training examples are provided to the learner. The learner must predict the target value  $c(x)$  before it is shown the correct value by the trainer.
- we check how many mistakes it makes in its predictions, before it learns the target concept.
- Learning the target concept exactly means converging to a hypothesis.

15 - Percentiles  
50, 55, 58, 63, 64, 66, 68, 70, 71  
72, 73, 75, 76, 77, 78, 81, 82, 83  
84, 87, 88, 89, LE: 502, 504, 505  
506, 507,

## Mistake Bound for FIND-S Algorithm

Recall Find-S algorithm, which incrementally computes the maximally specific hypothesis consistent with the training examples.

FIND-S :  $\Rightarrow$  Initialize  $h$  to the most specific hypothesis

$$l_1 \perp \neg l_1 \wedge l_2 \perp \neg l_2 \dots \perp \neg l_n$$

[ literals (+ve & -ve) ]

$\Rightarrow$  for each positive training instance  $x$   
 \* Remove from  $h$  any literal that is not satisfied by  $x$ .

$\Rightarrow$  Output hypothesis  $h$ .

FIND-S converges in the limit to a hypothesis that makes no errors on  $C \subseteq H$  & provide noise-free Training Data.

$\Rightarrow$  For the hypothesis that makes no errors, this generalization step consists of deleting unsatisfied literals.  
 $\Rightarrow$  Never classifies a negative example as positive.  
 $\Rightarrow$  Current hypothesis  $h$  is always at least as specific as target concept  $C$ .

$\Rightarrow$  Learn how many mistakes occur before FIND-S learns  $C$  & consider first positive example encountered by FIND-S  
 $\Rightarrow$  Learner makes a mistake in classifying this example, because its initial hypothesis labels every instance negative.  
 $\Rightarrow$  The number of mistakes will be 1 more than  $n$  hypothesis i.e.,  $1+n$  or  $n+1$ . This number of mistakes will require in worst case.

## Weighted-Majority Algorithm:

This algorithm make predictions by taking a weighted vote among a pool of prediction algorithms, and learns altering the weights.

It begins by assigning a weight of 1 to each prediction algorithm, then considers training examples.

Algorithm:  $a_i$  denotes the  $i$ th prediction algorithm in a pool of  $A$  of algorithms.  $w_i$  denotes the weight associated with  $a_i$ :

- For all  $i$  initialize  $w_i \leftarrow 1$
- For each training example  $\langle x, c(x) \rangle$ 
  - Initialize  $q_0$  and  $q_1$  to 0
    - For each prediction algorithm  $a_i$ 
      - if  $a_i(x) = 0$  then  $q_0 \leftarrow q_0 + w_i$
      - if  $a_i(x) = 1$  then  $q_1 \leftarrow q_1 + w_i$
  - If  $q_1 > q_0$  then predict  $c(x) = 1$
  - If  $q_0 > q_1$ , then predict  $c(x) = 0$
  - If  $q_0 = q_1$ , then predict 0 or 1 at random for  $c(x)$ .
- For each prediction algorithm  $a_i$  in  $A$  do
  - If  $a_i(x) \neq c(x)$  then  $w_i \leftarrow \beta w_i$

## Instance Based Learning

In general explicit description of the target function when training examples are provided.

Instance Based Learning methods simply stores the Training Examples.

It uses Nearest Neighbour & Locally weighted regression methods.

These are sometimes called "Lazy" learning methods because they delay processing until a new instance must be classified.

The methods used by this learning are straight forward for approximating real-valued or discrete-valued.

These use more complex, symbolic representations for instances.

"neighbouring" instances is elaborated accordingly.

### K - Nearest Neighbour Learning:

The most basic instance-based method is

K-NEAREST NEIGHBOUR algorithm.

The nearest neighbour of an instance are defined in terms of Examples.

It is an Lazy Learning, because it function will not learn from Training Data, it just Memorize.

Ex:-

Example: Given Data & very

$x = [\text{maths} = 6, \text{cs} = 8]$  and  $k = 3$

classification : pass/fail.

constant gives  
inf. about nn

	Maths	CS	Result	we use Euclidean's Approach for distance
1)	4	3	F	
2)	6	7	P	
3)	7	8	P	
4)	5	5	F	
5)	8	8	P	

$$d_1 = \sqrt{(6-4)^2 + (8-3)^2} \quad [6 \text{ is observed, } 4 \text{ is Actual}] \quad \text{Maths}$$

$$= \sqrt{2^2 + 5^2} = \sqrt{29} = 5.38$$

$$d_2 = \sqrt{(6-6)^2 + (8-7)^2} = \sqrt{0+1} = 1$$

$$d_3 = \sqrt{(6-7)^2 + (8-8)^2} = \sqrt{1+0} = 1$$

$$d_4 = \sqrt{(6-5)^2 + (8-5)^2} = \sqrt{1+9} = \sqrt{10} = 3.16$$

$$d_5 = \sqrt{(6-8)^2 + (8-8)^2} = \sqrt{4+0} = \sqrt{4} = 2$$

∴ Among all distances we have to find nearest three  
neighbours i.e.,  $d_2, d_3, d_5$  are nearest with their

values 1, 1, 2

check the result of  $d_2, d_3, d_5$  i.e., 2, 3, 5 all  
of them are pass. so, what we selected are nearest  
neighbours