

# AI - ML IA3 QUESTION BANK (MOD 4)

1.a) DERIVE AN EQN. FOR MAP HYPOTHESIS USING BAYES THEOREM

Ans: → BAYES THEOREM

- \* Bayes theorem provides a way to calculate the probability of a hypothesis based on its prior probability, the probabilities of observing various data given in hypothesis and the observed data itself.
- \* In Bayesian learning, the best hypothesis from some hypothesis space H, given the observed training data D means the most probable hypothesis, given the data D plus any initial knowledge about the prior probabilities of the various hypotheses in H.

$$\text{BAYES THEOREM} : \quad P(h|D) = \frac{P(D|h) P(h)}{P(D)}$$

where  $P(h) / P(D) \rightarrow$  prior probability

$P(h|D) / P(D|h) \rightarrow$  posterior probability

→ MAP with BAYES THEOREM

- \* The learner considers some set of candidate hypotheses H and it is interested in finding the most probable hypothesis  $h \in H$  given the observed data D.

- \* Any such maximally probable hypothesis is called a MAXIMUM a POSTERIORI (MAP) HYPOTHESIS ( $h_{MAP}$ )

\* We determine the MAP hypotheses by using BAYES THEOREM to calculate the posterior probability of each candidate hypothesis.

\*  $h_{MAP}$  is a MAP hypothesis provided:

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

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

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

we dropped  $P(D)$   $\rightarrow$  it is a constant (independent of  $h$ )

$\therefore$  If every hypothesis in  $H$  is equally probable a priori  
i.e.,  $P(h_i) = P(h_j) \forall h_i, h_j \in H$

Any hypothesis that maximizes  $P(D|h)$  is called  
MAXIMUM LIKELIHOOD (ML) HYPOTHESIS ( $h_{ML}$ )

$$\Rightarrow \boxed{h_{ML} = \operatorname{argmax}_{h \in H} P(D|h)} \quad \boxed{\begin{array}{l} P(h) \text{ is constant} \\ \forall h \in H \end{array}}$$

→ 4b) Explain  $h_{MAP}$  using BAYES THEOREM.

1b) CONSIDER A FOOTBALL GAME B/W TWO RIVAL TEAMS, SAY TEAM A & TEAM B. SUPPOSE TA WINS 65% OF THE TIME & TB WINS THE REMAINING MATCHES. AMONG THE GAMES WON BY TA, ONLY 35% OF THEM COME FROM PLAYING AT TB'S FOOTBALL FIELD. ON THE OTHER HAND, 75% OF THE VICTORIES FOR TB ARE OBTAINED WHILE PLAYING AT HOME. IF TB IS TO HOST THE NEXT MATCH B/W THE TWO TEAMS, WHAT IS THE PROBABILITY THAT IT WILL EMERGE AS THE WINNER?

Ans:  $P(T_A)$  = Probability that TA wins : 0.65

$P(T_B)$  = Probability that TB wins ;  $1 - P(T_A) = 0.35$

Probability that TB hosted the match it had won

$$P(X_B | T_B) = 0.75$$

Probability that TB hosted the match won by TA

$$P(X_B | T_A) = 0.35$$

The given question can be solved by obtaining  $P(T_B | X_B)$

$\Rightarrow$  probability that TB wins the next match if host.

We know: BAYES THEOREM ;  $P(h|D) = \frac{P(D|h)P(h)}{P(D)}$

$\Rightarrow$  Using BAYES THEOREM,

$$P(T_B | X_B) = \frac{P(X_B | T_B) P(T_B)}{P(X_B)}$$

$$P(B) = \sum_{i=1}^n P(B|A_i) P(A_i) \quad \Rightarrow \quad P(X_B | T_B) P(T_B) + P(X_B | T_A) P(T_A)$$

$$= \frac{0.75 \times 0.35}{(0.75 \times 0.35) + (0.35 \times 0.65)} = 0.5357$$

$$P(T_B | X_B) = 0.5357 \approx 54\%$$

↗ 6a)

2a) EXPLAIN (1) BRUTE FORCE MAP LEARNING ALGORITHM  
 (2) NAIVE BAYES CLASSIFIER.

Ans (1) BRUTE FORCE MAP LEARNING ALGORITHM

- \* We use BAYES THEOREM in CONCEPT LEARNING as it provides a principled way to calculate the posterior probability of each hypothesis given the training data.
- \* It acts as a basis for a straight fwd. learning algorithm that calculates the probability for each possible hypothesis  $\Rightarrow$  outputs most PROBABLE.

\* Consider the concept-learning problem :

- ① Assume  $L$  (LEARNER) considers some finite hypothesis  $H$  defined over the instance space  $X$ , in which TASK : to learn target concept  $C : X \rightarrow \{0, 1\}$

- ②  $L$  is given some sequence of training examples  $(x_1, d_1), \dots, (x_m, d_m)$  where :

$x_i \rightarrow$  some instance  $\in X$

$d_i \rightarrow$  target value of  $x_i \Rightarrow [d_i = c(x_i)]$

- ③ The sequence of target values ( $D$ ) :

$$D = (d_1, \dots, d_m)$$

- \* Design a st.fwd. CONCEPT LEARNING ALGORITHM to output MAP hypothesis, based on BAYES THEOREM.

BRUTE  
FORCE  
MAP  
LEARNING  
ALGORITHM

- i) For each hypothesis  $h$  in  $\epsilon H$  calculate POSTERIOR PROBABILITY

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

- ii) Output h<sub>MAP</sub> with highest POSTERIOR PROBABILITY

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

## \* ASSUMPTIONS FOR BRUTE-FORCE MAP LEARNING ALGORITHM

\* Let's choose  $P(h)$ ,  $P(D|h)$  to be consistent with the following assumptions

→ TRAINING DATA  $D \rightarrow$  noise free  $[d_i = c(x_i)]$

→ TARGET CONCEPT  $c$  is contained in HYPOTHESIS SPACE

→ Do not have a PRIORI REASON to believe that any hypothesis is more probable than any other.

$$* P(h) = \frac{1}{|H|} \text{ for all } h \in H \quad \left[ \begin{array}{l} \text{all prior probabilities} \\ \text{sum to 1} \end{array} \right]$$

↳ prior probability

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

↳ probability of observing the target values  $D$  for fixed set of instances  $X$ , given a world in which hypothesis  $h$  holds.

## BRUTE-FORCE MAP LEARNING ALGORITHM DERIVATION.

$$\textcircled{1} \text{ Bayes theorem : } P(h|D) = \frac{P(D|h) P(h)}{P(D)}$$

\textcircled{2}

when  $h$  is  
INCONSISTENT

$$P(h|D) = \frac{0 \cdot P(h)}{P(D)} = 0$$

$$\boxed{P(h|D) = 0}$$

$\boxed{VS_{H,D}}$   
Subset of hypotheses  
from  $H$  consistent  
with  $D$   
 $H' \subseteq D$

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

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

∴ Under assumed  $P(h)$  &  $P(D|h)$

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

## (2) NAIVE BAYES CLASSIFIER

- \* The Naive Bayes classifier applies to learning tasks where each instance  $x$  is described by a conjunction of attribute values and where the target function  $f(x)$  can take any value from some finite set  $V$ .
- \* A set of training examples of the target function is provided and a new instance is presented, described by :  $(a_1, a_2, \dots, a_m)$  (TUPLE OF ATTR. VALUES)
- \* The learner  $L$  is asked to predict the target value or classification, for this NEW INSTANCE.

\* Using BAYESIAN APPROACH,

$$\left[ \begin{array}{l} \text{Most probable} \\ \text{target value} \end{array} \right] V_{MAP} = \underset{V_j \in V}{\operatorname{argmax}} P(V_j | a_1, a_2, \dots, a_n) \rightarrow 1$$

\* Using BAYES THEOREM & EQN ①

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

$$\left[ \begin{array}{l} V_{MAP} = \underset{V_j \in V}{\operatorname{argmax}} P(a_1, a_2, \dots, a_n | V_j) P(V_j) \end{array} \right] \rightarrow 2$$

\* Assuming :  $P(a_1, a_2, \dots, a_n | V_j) = \prod_i P(a_i | V_j)$   $\rightarrow 3$   
 Given target value of instance, the probability of observing the conjunction  $(a_1, a_2, \dots, a_n) = \text{product of the probabilities for the individual attributes}$ .

\* With eqn ② & eqn ③, we get NAIVE BAYES CLASSIFIER:

target value  
output  
by NB classifier

$$V_{NB} = \underset{V_j \in V}{\operatorname{argmax}} P(V_j) \left[ \prod_i P(a_i | V_j) \right]$$

2b) The following table gives data set about stolen vehicles  
Using NAIVE BAYES CLASSIFIER classify the new data  
(Red, SUV, Domestic)

N	COLOR	TYPE	ORIGIN	STOLEN?
1	red	sports	domestic	YES 1
2	red	sports	domestic	NO
3	red	sports	domestic	YES 2
4	yellow	sports	domestic	NO
5	yellow	sports	imported	YES 3
6	yellow	SUV	imported	NO
7	yellow	SUV	imported	YES 4
8	yellow	SUV	domestic	NO
9	red	SUV	imported	NO
10	red	sports	domestic	YES 5

Ans:  $\rightarrow$  NEW INSTANCE = {Red, SUV, Domestic}

$$\rightarrow P(\text{STOLEN} = \text{YES}) = \frac{5}{10} = 0.5$$

$$\rightarrow P(\text{STOLEN} = \text{NO}) = \frac{5}{10} = 0.5$$

COLOR	YES		NO		TYPE	YES		NO		ORIGIN	YES NO	
	Red	3/5	2/5			Sports	4/5	2/5			Domestic	2/5
Yellow	2/5	3/5			SUV	1/5	3/5			Imported	3/5	2/5

$$\rightarrow P(\text{YES} | \text{NEW INSTANCE}) = P(\text{YES}) * P(\text{Color} = \text{Red} | \text{YES}) * P(\text{Type} = \text{SUV} | \text{YES}) * P(\text{Origin} = \text{Domestic} | \text{YES})$$

$$= \frac{5}{10} * \frac{3}{5} * \frac{1}{5} * \frac{2}{5} = \underline{\underline{0.024}}$$

$$\rightarrow P(\text{NO} | \text{NEW INSTANCE}) = P(\text{NO}) * P(\text{Color} = \text{Red} | \text{NO}) * P(\text{Type} = \text{SUV} | \text{NO}) * P(\text{Origin} = \text{Domestic} | \text{NO})$$

$$[0.024 > 0.072] \quad \begin{bmatrix} \text{not} \\ \text{stolen} \end{bmatrix} = \frac{5}{10} * \frac{2}{5} * \frac{3}{5} * \frac{3}{5} = \underline{\underline{0.072}}$$

$\Rightarrow$  NEW INSTANCE classified as 'NO'

3a) Discuss Minimum Description Length Principle in brief:

Ans: \* Inductive bias are assumptions that are made by the learning algorithm to form a hypothesis or a generalization beyond the set of training instances in order to classify unobserved data.

\* Occam's razor is a simple inductive bias that involves preference for a simpler hypothesis that best fits the data.

BAYESIAN interpretation

\* MINIMUM DESCRIPTION LENGTH PRINCIPLE (MDL) is inspired from Occam's RAZOR that recommends choosing the hypothesis that captures the regularities in the data and this regularity is in turn used to compress/minimize it.

COMPRESS ↑ LEARNING ↑

\* The MDL is motivated by interpreting the definition of hMAP in the light of basic concepts from information theory.

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

which can be equivalently expressed in terms of maximizing the  $\log_2$ :

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} \log_2 P(D|h) + \log_2 P(h)$$

or  $h_{MAP} = \underset{h \in H}{\operatorname{argmin}} -\log_2 P(D|h) - \log_2 P(h) \rightarrow (1)$

$$\Rightarrow -\log_2 P(h)$$

the description length of  $h$  under the optimal encoding for the hypothesis space  $H$

$$[L_{C_H}(h) = -\log_2 P(h)]$$

$C_H \rightarrow$  optimal code for hypothesis space  $H$ .

$$\Rightarrow -\log_2 P(D|h)$$

the description length of the training data  $D$  given hypothesis  $h$ , under the optimal encoding from the hypothesis space  $H$ .

$$[L_{C_{D|h}}(D|h) = -\log_2 P(D|h)]$$

$C_{D|h} \rightarrow$  optimal code for describing data  $D$

(assuming Sender + Receiver know hypothesis  $h$ )

$\Rightarrow$  Rewrite Eqn ①, to show that  $h_{MAP}$  is the hypothesis  $h$  that minimizes the sum given by the description length of the hypothesis plus the description length of the data given the hypothesis

$$h_{MAP} = \underset{h \in H}{\operatorname{argmin}} L_{C_H}(h) + L_{C_{D|h}}(D|h) \quad ②$$

$\Rightarrow$  MDL principle recommends choosing the hypothesis that minimizes the sum of these two description lengths

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

where  $C_1$  &  $C_2$  are codes to represent  $H, D$ .

$\Rightarrow$  If  $C_1 = C_H$  &  $C_2 = C_{D|h}$ , then  $h_{MDL} = h_{MAP}$

3b) Explain BAYESIAN BELIEF NETWORK + CONDITIONAL INDEPENDENCE with example.

Ans: ① BAYESIAN BELIEF NETWORK (BBN)

- \* A BBN describes the probability distribution governing a set of variables by specifying a set of CONDITIONAL INDEPENDENCE assumptions along with a set of CONDITIONAL PROBABILITIES.
- \* It allows stating CONDITIONAL INDEPENDENCE assumptions that apply to a subset of variables represented by DIRECTED ACYCLIC GRAPHS.

→ NOTATION

- \* Consider an arbitrary set of random variables  $\{Y_1, Y_2, \dots, Y_n\}$  where  $Y_i \rightarrow$  can take on the set of possible values  $V(Y_i)$
- \* The joint space of the set of variables  $Y$  = cross product  $V(Y_1) \times V(Y_2) \times \dots \times V(Y_n)$
- \* Each item in JOINT SPACE corresponds to one of the possible assignments of values to the tuple of variables  $(Y_1, Y_2, \dots, Y_n)$ . The probability distribution over this JOINT SPACE is called JOINT PROBABILITY DISTRIBUTION. (J.P.D)
- \* J.P.D specifies the probability for each of the possible variable bindings for tuple  $(Y_1, \dots, Y_n)$
- \* BBN describes the J.P.D for a set of variables

## ② CONDITIONAL INDEPENDENCE

\* Let  $X, Y$  and  $Z$  be three discrete valued random variables.  $X$  is CONDITIONALLY INDEPENDENT of  $Y$  given  $Z$ , if the PROBABILITY DISTRIBUTION governing  $X$  is independent of the value of  $Y$ , given a value for  $Z$  that is, if:

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

where  $x_i \in V(X)$   
 $y_j \in V(Y)$   
 $z_k \in V(Z)$

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

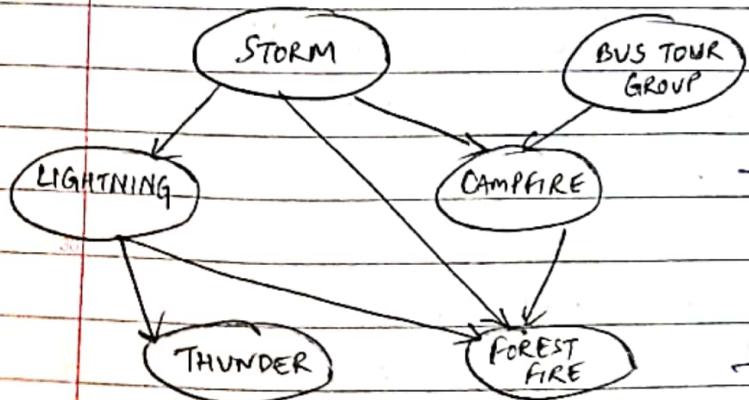
$$\Rightarrow \text{for set of variables } P(X_1, \dots, X_l | Y_1, \dots, Y_m, Z_1, \dots, Z_n) = P(X_1, \dots, X_l | Z_1, \dots, Z_n)$$

\* Using NB CLASSIFIER for  $P(A_1, A_2 | V)$  where  $A_1$  is conditionally independent of  $A_2$  given target value  $V$ :

$$P(A_1, A_2 | V) = P(A_1 | A_2, V) P(A_2 | V)$$

$$= \underline{P(A_1 | V)} P(A_2 | V)$$

$\Rightarrow$  EXAMPLE



$\Rightarrow$  ASSERTION: CAMPFIRE conditionally independent of LIGHTNING, THUNDER  
 Immediate PARENTS: STORM, BUS TOUR GROUP, CAMPFIRE

	S, B	S, TB	TS, B	TS, TB
C	0.4	0.1	0.8	0.2
TC	0.6	0.1	0.2	0.8

$$\rightarrow P(y_1, \dots, y_n) = \prod_{i=1}^n P(y_i | \text{PARENTS}(y_i))$$

$$\rightarrow P(\text{Campfire} = T | \text{Storm} = T, \text{BG} = T) = 0.4$$

4a) EXPLAIN  
① EM ALGORITHM  
② GIBBS ALGORITHM

Ans: ② GIBBS ALGORITHM

- \* Although the Bayes optimal classifier obtains the best performance that can be achieved from the given training data, it can be quite costly to apply.
- \* An alternative, less optimal method is GIBBS ALGORITHM defined as follows:
  - ① Choose a hypothesis  $h$  from  $H$  at random, according to the POSTERIOR PROBABILITY DISTRIBUTION over  $H$ .
  - ② Use  $h$  to predict the classification of the next instance  $x$ .
- \* Given a new instance to classify, the GIBBS ALGORITHM simply applies a hypothesis drawn at random according to current posterior probability distribution.
- \* Suppose  $p(x, y)$  is a p.d.f / p.m.f. that is difficult to sample from directly. Using conditional distributions, GIBBS SAMPLING can be done as;
  - \* Initialize  $y^0, x^0$
  - \* for  $j=1, 2, 3 \dots N$  do
    - sample  $x^{j,i} \sim p(x|y^{j-1})$
    - sample  $y^{j,i} \sim p(y|x^j)$
  - end for

\* On comparison of the errors;

$$\mathbb{E} [\text{error}_{\text{GIBBS}}] \leq 2 \mathbb{E} [\text{error}_{\text{BAYES OPTIMAL}}]$$

## → EM ALGORITHM

- \* Applying EM algorithm to the existing problem; it initializes the hypothesis to  $h = \langle \mu_1, \mu_2 \rangle$ .
- \* Then, iteratively re-estimates  $h$  by repeating the following two steps until the procedure converges to a stationary value for  $h$ .

STEP 1: Calculate the expected value  $E[z_{ij}]$  of each hidden variable  $z_{ij}$ , assuming  $h = \langle \mu_1, \mu_2 \rangle$

STEP 2: i) Calculate a new max. likelihood hypothesis

$$h' = \langle \mu'_1, \mu'_2 \rangle$$

assuming the value taken on by each hidden variable  $z_{ij}$  is its expected value  $E[z_{ij}]$  calculated in STEP ①.

ii) Then, replace  $h = \langle \mu_1, \mu_2 \rangle$  with  $h' = \langle \mu'_1, \mu'_2 \rangle$  and iterate.

$$\begin{aligned} ① E[z_{ij}] &= \frac{p(x=x_i | \mu=\mu_j)}{\sum_{n=1}^2 p(x=x_i | \mu=\mu_n)} \\ &= \frac{e^{-\frac{1}{2\sigma^2}(x_i - \mu_j)^2}}{\sum_{n=1}^2 e^{-\frac{1}{2\sigma^2}(x_i - \mu_n)^2}} \end{aligned}$$

$$\begin{aligned} ② \quad \mu_j &\leftarrow \frac{\sum_{i=1}^m E[z_{ij}] x_i}{\sum_{i=1}^m E[z_{ij}]} \end{aligned}$$

5a) Prove that MAXIMUM LIKELIHOOD (BAYESIAN LEARNING) can be used in any learning algorithms that are used to minimize the squared error between actual output hypothesis & predicted output hypothesis.

Ans: Consider the following problem setting:

\* Learner L considers an instance space X and a hypothesis space H consisting of some REAL-VALUED FUNCTIONS defined over X,

The task of the learner is to output MAXIMUM LIKELIHOOD HYPOTHESIS given by:

$$\Rightarrow h_{ML} = \underset{h \in H}{\operatorname{argmax}} p(d|h), \quad - (1)$$

PROBABILITY DENSITY FN,

\* We assume a fixed set of training instances  $\langle x_1, \dots, x_m \rangle$  with corresponding target values:

$D = \langle d_1, \dots, d_m \rangle$  such that:

$$\Rightarrow \boxed{d_i = \underbrace{f(x_i)}_{\text{NOISE-FREE TARGET VALUE}} + \underbrace{e_i}_{\text{NOISE}} \rightarrow \text{NOISE}} \quad - (2)$$

\* Assuming the training examples are mutually independent given h, we get:

From (1), (2)  $\Rightarrow \boxed{h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m p(d_i|h)} \quad - (3)$

\* Given noise  $e_i$  obeys NORMAL DISTRIBUTION

with ZERO MEAN & SOME VARIANCE  $\sigma^2$

$\Rightarrow d_i$  must also obey a NORMAL DISTRIBUTION with variance  $\sigma^2$  centered around true  $f(x_i)$

rather than ZERO. Then:

$$\boxed{p(d_i|h) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} (d_i - \mu)^2}} \quad - (4)$$

\* From (3), (4)  $\Rightarrow h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - h(x_i))^2}$

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

\* Maximizing  $\ln p$  also maximizes  $p$  (MONOTONIC FN.)

$$\Rightarrow h_{ML} = \underset{h \in H}{\operatorname{argmax}} \left\{ \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{1}{2\sigma^2} (d_i - h(x_i))^2 \right\}$$

constant

$$\Rightarrow h_{ML} = \underset{h \in H}{\operatorname{argmax}} \left\{ - \frac{1}{2\sigma^2} (d_i - h(x_i))^2 \right\}$$

\* Maximizing -ve quantity = Minimizing +ve quantity

$$\Rightarrow h_{ML} = \underset{h \in H}{\operatorname{argmin}} \left\{ \frac{1}{2\sigma^2} (d_i - h(x_i))^2 \right\}$$

\* After discarding constants (independent of  $h$ ):

$$\Rightarrow h_{ML} = \underset{h \in H}{\operatorname{argmin}} \left\{ (d_i - h(x_i))^2 \right\}$$

$\Rightarrow$  This shows that MAXIMUM LIKELIHOOD HYPOTHESIS  $h_{ML}$  minimizes the sum of the squared errors between observed training values  $d_i$  & hypothesis predictions  $h(x_i)$ .

5b) Estimate Conditional Probabilities of each attribute

$\Rightarrow \{\text{colour, legs, height, smelly}\}$  for species classes  $\{M, H\}$

Using given data in table and provide probability values for NEW INSTANCE : (Colour = GREEN, Legs = 2, Height = Tall, Smelly = No)

COLOUR	LEGS	HEIGHT	SMELLY	SPECIES
White	3	Short	Yes	M
Green	2	Tall	No	M
Green	3	Short	Yes	M
White	3	Short	Yes	M
Green	2	Short	No	H
White	2	Tall	No	H
White	2	Tall	No	H
White	2	Short	Yes	H

Ans:

$$P(M) = \frac{4}{8} = 0.5$$

$$P(H) = \frac{4}{8} = 0.5$$

COLOUR	M	H	LEGS	M	H	HEIGHT	M	H	SMELLY	M	H
White	$\frac{2}{4}$	$\frac{3}{4}$	2	$\frac{1}{4}$	$\frac{4}{4}$	Short	$\frac{3}{4}$	$\frac{2}{4}$	Short	$\frac{3}{4}$	$\frac{1}{4}$
Green	$\frac{2}{4}$	$\frac{1}{4}$	3	$\frac{3}{4}$	$\frac{0}{4}$	Tall	$\frac{1}{4}$	$\frac{2}{4}$	Tall	$\frac{1}{4}$	$\frac{3}{4}$

$$\begin{aligned} \rightarrow P(M | \text{NEW INSTANCE}) &= P(M) * P(\text{COLOUR} = \text{Green} | M) * P(\text{LEGS} = 2 | M) \\ &\quad * P(\text{HEIGHT} = \text{Tall} | M) * P(\text{SMELLY} = \text{No} | M) \\ &= 0.5 * \frac{2}{4} * \frac{1}{4} * \frac{1}{4} * \frac{1}{4} = \frac{1}{256} = \underline{0.0039} \end{aligned}$$

$$\begin{aligned} \rightarrow P(H | \text{NEW INSTANCE}) &= P(H) * P(\text{COLOUR} = \text{Green} | H) * P(\text{LEGS} = 2 | H) \\ &\quad * P(\text{HEIGHT} = \text{Tall} | H) * P(\text{SMELLY} = \text{No} | H) \\ &= 0.5 * \frac{1}{4} * \frac{1}{4} * \frac{2}{4} * \frac{3}{4} = \underline{0.0117} \end{aligned}$$

$\therefore 0.0039 < 0.0117 \Rightarrow \text{NEW INSTANCE } \in \text{Species H}$