

Q.1 Consider the following dataset with 4 records.

[4+2 = 6 Marks]

Input X	Output Y
1	$\exp(2)$
2	$\exp(4)$
3	$\exp(6.3)$
4	$\exp(9.2)$

Assume output $y = e^{(\alpha * x)}$. Using linear regression,

- (a) Find the best value of α .
 (b) Find the optimal total sum of square error.

Solution:

1. Solution :-

a) Find the best value of α

$\ln(y_i) = \alpha x_i$

$J(\alpha) = (\alpha - 2)^2 + (2\alpha - 4)^2 + (3\alpha - 6.3)^2 + (4\alpha - 9.2)^2$

Taking derivative of $J(\alpha)$ w.r.t α and equating to 0.

$(\alpha - 2) + 2(2\alpha - 4) + 3(3\alpha - 6.3) + 4(4\alpha - 9.2) = 0$

or

$\alpha = 65.7/30$

$= 2.2$.

b) Find the optimal total sum of square error.

$(\exp(2.2) - \exp(2))^2 + (\exp(4.4) - \exp(4))^2 + (\exp(6.6) - \exp(6.3))^2 + (\exp(8.8) - \exp(9.2))^2$

Bias-Variance of the Squared Error

$$S = (y - \hat{y})^2$$

$$\begin{aligned}(y - \hat{y})^2 &= (y - E[y] + E[\hat{y}] - \hat{y})^2 \\ &= (y - E[y])^2 + (E[\hat{y}] - \hat{y})^2 + 2(y - E[y])(E[\hat{y}] - \hat{y})\end{aligned}$$

$$E[S] = E[(y - \hat{y})^2]$$

$$\begin{aligned}E[(y - \hat{y})^2] &= (y - E[y])^2 + E[(E[\hat{y}] - \hat{y})^2] \\ &= \text{Bias}^2 + \text{Var}\end{aligned}$$

Q.2 Consider inputs x_i which are real valued attributes and the outputs y_i which are real

valued of the form $y_i = f(x_i) + e_i$, where $f(x_i)$ is the true function and e_i is a random variable representing laplacian noise with PDF given by

$$f(y_i/\theta) = \frac{1}{2\theta} * e^{-\frac{|y_i - \mu|}{\theta}}$$

Implementing a linear regression model of the form $h(x_i) = \sum_{i=0}^n \theta_i x_i$, and $\mu = h(x_i)$,

find the maximum likelihood estimator of θ . Comment on the loss function. [4+1 Marks]

$$\begin{aligned}
f(y_i|\theta) &= \frac{1}{2\sigma} e^{-\frac{|y_i - \theta|}{\sigma}} \\
L(\theta|y) &\stackrel{\text{argmax}}{\rightarrow} \sum_{i=1}^n f(y_i|\theta) \\
&= \sum_{i=1}^n \frac{1}{2\sigma} e^{-\frac{|y_i - \theta|}{\sigma}} \\
\ln L(\theta|y) &= \ln \left(\prod_{i=1}^n \frac{1}{2\sigma} e^{-\frac{|y_i - \theta|}{\sigma}} \right) \\
&= \sum_{i=1}^n \ln \frac{1}{2\sigma} e^{-\frac{|y_i - \theta|}{\sigma}} \\
&= \sum_{i=1}^n -\ln [2\sigma \cdot e^{-\frac{|y_i - \theta|}{\sigma}}] \\
&= -n \ln 2\sigma - \sum_{i=1}^n \frac{|y_i - \theta|}{\sigma} \\
&= \underset{\theta}{\text{argmax}} \left[-n \ln 2\sigma - \sum_{i=1}^n \frac{|y_i - \theta|}{\sigma} \right] \\
&= \underset{\theta}{\text{argmin}} \left[n \ln 2\sigma + \sum_{i=1}^n \frac{|y_i - \theta|}{\sigma} \right] \\
\frac{d}{d\theta} \ln L(\theta|y) &= \frac{d}{d\theta} \left[n \ln 2\sigma + \sum_{i=1}^n \frac{|y_i - \theta|}{\sigma} \right] \\
&= \theta \frac{n}{\sigma} - \frac{1}{\sigma^2} \sum_{i=1}^n |y_i - \theta| = 0 \\
\Rightarrow \theta &= \frac{1}{\sigma^2} \sum_{i=1}^n |y_i - \theta|
\end{aligned}$$

$$\theta = \frac{1}{n} \sum_{i=1}^n |y_i - \theta|$$

Comment on Loss function: Instead of MSE, MAE is the maximum likelihood hypothesis. So MAE is appropriate for the loss function.

Q.3 Consider a result prediction system where student's efforts are encoded as percent of time a student has spent studying out of total available time.

- The input X is having just one feature representing the student's efforts having only four discrete values (25%, 50%, 75%, and 100%)
- The output Y is having 3 classes (First class, Second class, Fail)
- The priors for each class are: $P(Y = \text{First Class}) = 0.5$, $P(Y = \text{Second class}) = 0.3$, and $P(Y = \text{Fail}) = 0.2$.
- Based on the past data, the estimated the class-conditional probability $P(X|Y)$ are shown in the following table.

Student's efforts	$p(x y=\text{fail})$	$p(x y=\text{second class})$	$p(x y=\text{first class})$
25	0.7	0.4	0.1
50	0.2	0.3	0.1

75	0.1	0.2	0.3
100	0	0.1	0.7

Consider a following loss function $\ell(\hat{y}, y)$ where \hat{y} = predicted class label and y is true class label:

$$\ell(\hat{y}, y) = \begin{cases} 0 & \hat{y} = y \\ 1 & \hat{y} = \text{Fail and } \hat{y} \neq y \\ 2 & \hat{y} = \text{Second class and } \hat{y} \neq y \\ 4 & \hat{y} = \text{First class and } \hat{y} \neq y \end{cases}$$

Consider modified Naïve Bayes hypothesis function:

$$\hat{Y} \leftarrow \underset{y_k}{\operatorname{argmax}} l(y, \hat{y}) P(Y = y_k) \prod_i P(X_i | Y = y_k)$$

Use this modified hypothesis function to classify each of the examples in the given table. [5 Marks]

Solution:

	p(y x)			
	Fail	second class	first class	
25	0.35	0.12	0.02	
50	0.1	0.09	0.02	
75	0.05	0.06	0.06	
100	0	0.03	0.14	

	L(y-hat,y) * p(y/x)	L(y-hat,y) * p(y/x)	L(y-hat,y) * p(y/x)		
	Fail	second class	first class	Highest value	
25	0.35	0.24	0.08	0.35	fail
50	0.1	0.18	0.08	0.18	second class
75	0.05	0.12	0.24	0.24	first class
100	0	0.06	0.56	0.56	first class

Q.4 If we modify the loss function of the linear regression model as follows:

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n w^{(i)} \left(h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

Where $w^{(i)}$ is the weight assigned to each training example. Derive the equation to find the value of θ with this modified loss function. Suppose, we estimate the value of $w^{(i)}$ inversely proportional to the variance of the residuals, comment in **no more than 20 words** when you prefer to use this kind of modified loss function.

[3+2=5 Marks]

Solution:

4. Solution :-

Hint :- you might need to use a matrix w such that $\text{diag}(w) = [w_1, w_2, \dots, w_N]^T$

Define $Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$ and $X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix}$

Then $L(\beta) = (y - X\beta)^T w (y - X\beta)$.

Setting $\frac{dL(\beta)}{d\beta} = 0$, we get

$$\hat{\beta} = (X^T w X)^{-1} X^T w Y$$

Comment: Robust against outliers. Outliers will have higher variance of the residuals resulting into lower weight.

Q.5 Fit a logistic regression. Find the updated weights after the 3 iterations of modified Gradient Descent algorithm where gradient update happens after every training example using a learning rate of 0.5 and initial weights $(W_0, W_1, W_2) = (1, 1, 1)$ for the

following data with the logistic regression output given by

$$\frac{1}{1 + e^{(-W_0 + W_1 X_1^2 - W_2 X_2)}}$$

Assume the results obtained after 3 iterations is the final weights. Using this construct the confusion matrix for given below training data. [4+2=6 Marks]

Input X1	Input X2	Output Label
2	0	0
0	2	0
0	-2	0
-2	0	0
0	1	1
0	-1	1

Solution:

Wi-LR*[Y-Pred - Y]* Xi						
X1	X2	y	Y-Pred = h(X)	w0	w1	w2
2	0	0	0.05	0.975	0.95	1
0	2	0	0.95	0.5	0.95	0.05
0	-2	0	0.27	0.365	0.95	0.32
-2	0	0	0.05			

Handwritten calculations for logistic regression weights:

$$\hat{y} = \frac{1}{1 + e^{(-1+4)}} = \frac{1}{1 + e^3} \approx 0.05$$

$$w_0 = 0.5 \times (\hat{y} - y) \times x_1 = (-0.025) \approx 0.975$$

$$= 1 - 0.5(0.05 - 0)x_1$$

$$w_1 = 0.5 \times (\hat{y} - y) \times x_2 = 0.95$$

$$= 1 - 0.5(0.05 - 0)x_2$$

$$w_2 = 0.5 \times (\hat{y} - y) \times 1 = 1$$

$$= 1 - 0.5(0.05 - 0) \times 1$$

w0 0.365	w1 0.95	w2 0.32	Confusion Matrix				True Class			
X1	X2	y	Y-Pred = h(X)	Y	Predicted Class	Y=0	Y=1	Y=0	Y=1	
2	0	0	0.03	0	Y=0	3	0	Y=0	Y=1	
0	2	0	0.73	1	Y=1	1	2	Y=1	Y=0	
0	-2	0	0.43	0						
-2	0	0	0.03	0						
0	1	1	0.66	1						
0	-1	1	0.51	1						

- 1) calculate y_{pred} first
- 2) calculate w_0, w_1, w_2 acc. to formula given:
 $w_i = LR^* [y_{pred} - y] * x_i^*$

[consider $x_0 \rightarrow 1$] [LR = 0.5, given in question]

[calculate w_0, w_1, w_2 for 3 iterations & consider the 3rd one to calculate y_{pred} again]

x_1	x_2	Output	y_{pred}	$w_i - LR^* [y_{pred} - y] * x_i^*$
				w_0 w_1 w_2
2	0	0	0.05	0.975 0.95 0.1
0	2	0	0.9532	0.5 0.95 0.05
0	-2	0	0.2695	0.365 0.95 0.32
-2	0	0	0.05	
.

For iteration (1) :-

$$w_0 := w_0 - (0.5 \times (0.05 - 0) \times (x_0))$$

$$= 1 - 0.5 \times 0.05 \times 1$$

$$w_0' = 0.975$$

$$w_1 = w_1 - 0.5 \times 0.05 \times 2$$

$$= 1 - 0.05 = 0.95$$

$$w_2 = w_2 - 0 = 1$$

For iteration (2) :- $w_0 = 0.975 - (0.5 \times 0.95 \times 1)$

(use iteration
w₀ value
here)

$$w_0 = 0.5$$

$$w_1 = \underbrace{0.95}_{\text{iteration 1 value}} - (0.5 \times 0.95 \times 0)$$

$$= 0.95$$

$$w_2 = \underbrace{1 - (0.5 \times 0.95 \times 0)}_{\text{iteration 1 value}} = 0.05$$

⇒ Similarly, calculate for iteration-3 & those values are considered to further calculate y_{pred} .

⇒ Then take round-off of y_{pred} & consider it as the output prediction value \hat{Y} for confusion matrix.

$$\Rightarrow \begin{array}{c|ccc|cc} & w_0 & w_1 & w_2 & & \\ \hline & 0.365 & 0.95 & 0.32 & & \\ \begin{array}{c|c|c|c|c} x_1 & x_2 & y & y_{\text{pred.}} & y \\ \hline 2 & 0 & 0 & 0.03 & 0 \\ 0 & 2 & 0 & 0.73 & 1 \\ 0 & -2 & 0 & 0.43 & 0 \\ -2 & 0 & 0 & 0.03 & 0 \\ 0 & 1 & 1 & 0.66 & 1 \\ 0 & -1 & 1 & 0.51 & 1 \end{array} & & & & \end{array}$$

\therefore , confusion matrix :-

$$\begin{array}{c|cc|c} & & \text{True Class} & \\ \hline & & y=0 & y=1 \\ \hline \text{Predicted class} & y=0 & 3 & 0 \\ & y=1 & 1 & 2 \end{array}$$

Q.6 Consider the following set of training examples:

Instance	Classification	A ₁	A ₂
1	+	T	T
2	+	T	T
3	-	T	F
4	+	F	F
5	-	F	T
6	-	F	T

What is the information gain of A₂ relative to these training examples? Provide the equation for calculating the information gain as well as intermediate results.

[3 Marks]

6. Solution :-

$$\text{Entropy } E(S) = E([3+, 3-])$$

$$= -(3/6) \log_{10}(3/6) - (3/6) \log_{10}(3/6)$$

5
[Entropy $E(S) = 1$]

$$\text{Gain } (S, a_2) = E(S) - (4/6)E(T) - (2/6)E(F)$$

10
 $= 1 - (4/6) - (2/6) \approx 0$

$$E(T) = E([2+, 2-]) = 1$$

$$E(F) = E([1+, 1-]) = 1$$

Birla Institute of Technology & Science, Pilani
Work-Integrated Learning Programmes Division
MTech. Software Engineering at DSE (FC04, FA04_1-2021) Cluster

Second Semester 2021-2022
Mid-Semester Test
(EC-2 Regular)

Course No.	: DSECLZG565
Course Title	: Machine Learning
Nature of Exam	: Open Book
Weightage	: 30%
Duration	: 2 Hours
Date of Exam	: 10-07-2022(FN)

No. of Pages = 2
No. of Questions = 6

Note:

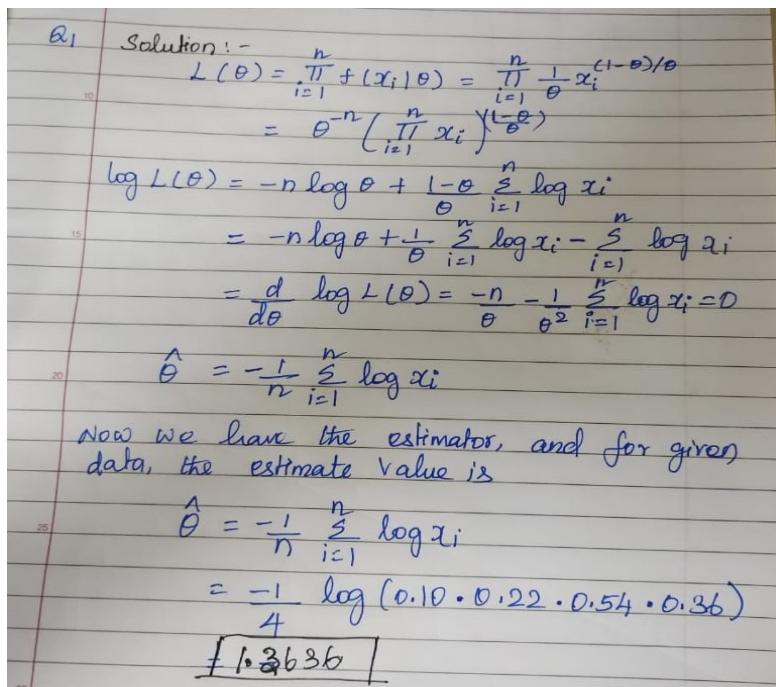
1. Please follow all the *Instructions to Candidates* given on the cover page of the answer book.
2. All parts of a question should be answered consecutively. Each answer should start from a fresh page.
3. Assumptions made if any, should be stated clearly at the beginning of your answer.

Q.1 Let T_1, T_2, \dots, T_n be a random sample of a population describing the website loading time on a mobile browser with probability density function given as:

$$f(t/\theta) = \frac{1}{\theta} t^{\frac{(1-\theta)}{\theta}} \quad \text{where } 0 < t < 1 \text{ and } 0 < \theta < \infty$$

Find the maximum likelihood estimator of θ . What is the estimate of θ , if the website loading time from four samples are $t_1 = 0.10, t_2 = 0.22, t_3 = 0.54, t_4 = 0.36$. [5 Marks]

Solution:



Q1 Solution :-

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n f(x_i | \theta) = \prod_{i=1}^n \frac{1}{\theta} x_i^{\frac{(1-\theta)}{\theta}} \\ &= \theta^{-n} \left(\prod_{i=1}^n x_i \right)^{\frac{1-\theta}{\theta}} \end{aligned}$$

$$\begin{aligned} \log L(\theta) &= -n \log \theta + \frac{1-\theta}{\theta} \sum_{i=1}^n \log x_i \\ &= -n \log \theta + \frac{1}{\theta} \sum_{i=1}^n \log x_i - \frac{n}{\theta} \log \theta \\ &= \frac{d}{d\theta} \log L(\theta) = -n - \frac{1}{\theta} \sum_{i=1}^n \log x_i = 0 \end{aligned}$$

$$\hat{\theta} = -\frac{1}{n} \sum_{i=1}^n \log x_i$$

NOW we have the estimator, and for given data, the estimate value is

$$\begin{aligned} \hat{\theta} &= -\frac{1}{n} \sum_{i=1}^n \log x_i \\ &= -\frac{1}{4} \log (0.10 \cdot 0.22 \cdot 0.54 \cdot 0.36) \\ &= 1.3636 \end{aligned}$$

Marking Scheme: Derivation of $\theta = 3$ marks (step wise marks)

θ Computation = 2 marks (wrong value = 0 marks)

- Q.2 As a part of efforts to improve students' performance in the exams, you have been given the data showing number of study hours spent by students, their gender and their final results as pass or fail. Using this sample dataset, apply Naïve Bayes classification technique, to classify the test case {No of study hours = 3.5, Gender="male"} either as "Pass", or "Fail". [5 Marks]

No of study hours	Gender	Final result
4.5	Male	Pass
7	Female	Pass
2	Male	Fail
4	Female	Fail
2.5	Male	Fail
3	Female	Fail
8.3	Male	Fail
8	Female	Pass
9	Male	Pass

Solution:

1. Prior: [1M]

p(y=Pass)	p(y=Fail)
0.444444	0.555556

2. No of study hours –X1: continuous variable, applying class conditional PDF [1M]

3.

$$p(X_i = x | Y = y_k) = \frac{1}{\sqrt{2\pi\sigma_{ik}^2}} e^{-\frac{1}{2}\left(\frac{x-\mu_{ik}}{\sigma_{ik}}\right)^2}$$

	Variance	mean
Pass class	2.945	7.2
Fail class	4.64	3.9

4. X1=3.5, X2="male" [3M]

$$\hat{Y} \leftarrow \operatorname{argmax}_{y_k} P(Y = y_k) \prod_i P(X_i | Y = y_k)$$

p(X1/ y=Pass)	0.105614
p(X1/ y=Fail)	0.184564

p(X2/ y=Pass)	0.5
p(X2/ y=Fail)	0.6

P(y=Pass/X)	0.02346969
P(y=Fail/X)	0.061521395

Class : Fail

Q.3 The 2-input AND gate is implemented using logistic regression classifier with gradient descent optimization algorithm. The model parameters at time t are given by $\theta_0=0$, $\theta_1=0$, and $\theta_2=0$. Given binary input (x_1, x_2) , [2+3 = 5 Marks]

a) What will be value of the loss function at t ? [2M]

Solution:

Cross entropy loss:

$$= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)})) \right]$$

x1	x2	Target y	Actual Output-yhat	y.ln(yhat)+(1-y)ln(1-yhat)
0	0	0	0.5	0*ln0.5+(1-0)*ln(1-0.5)
0	1	0	0.5	0*ln0.5+(1-0)*ln(1-0.5)
1	0	0	0.5	0*ln0.5+(1-0)*ln(1-0.5)
1	1	1	0.5	1*ln(0.5)

total loss=	0.693147181
-------------	-------------

- Q.3 The 2-input AND gate is implemented using **logistic regression classifier with gradient descent optimization algorithm**. The model parameters at time t are given

by $\theta_0=0, \theta_1=0$, and $\theta_2=0$. Given binary input (x_1, x_2)

[2+3 = 5 Marks]

- a) What will be value of the loss function at t ? [2M]

Solution:

Cross entropy loss:

$$= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)})) \right]$$

x1	x2	Target y	Actual Output-yhat	$y \cdot \ln(yhat) + (1-y) \ln(1-yhat)$
0	0	0	0.5	$0 \cdot \ln(0.5) + (1-0) \cdot \ln(1-0.5)$
0	1	0	0.5	$0 \cdot \ln(0.5) + (1-0) \cdot \ln(1-0.5)$
1	0	0	0.5	$0 \cdot \ln(0.5) + (1-0) \cdot \ln(1-0.5)$
1	1	1	0.5	$1 \cdot \ln(0.5)$

$$\hat{y} = \frac{1}{1 + e^{-\theta^T x}}$$

$$-\frac{1}{m} \times y \cdot \ln(\hat{y}) + (1-y) \cdot \ln(1-\hat{y})$$

total loss = 0.693147181 ✓

- b) What will be the values of θ_0, θ_1 and θ_2 at $(t+1)$ with learning rate $\alpha=1$ and L2 regularization constant $\lambda=1$? [3M]

Solution:

Cost function

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Apply gradient descent update rule

y-hat	y	yhat-y	x0	(yhat-y)x0	w0-new
0.5	0	0.5	1	0.5	-0.25
0.5	0	0.5	1	0.5	
0.5	0	0.5	1	0.5	
0.5	1	-0.5	1	-0.5	

y-hat	y	yhat-y	x1	(yhat-y)x1	regularized w1-new
0.5	0	0.5	0	0	0
0.5	0	0.5	0	0	
0.5	0	0.5	1	0.5	
0.5	1	-0.5	1	-0.5	

y-hat	y	yhat-y	x2	(yhat-y)x1	regularized w2-new
0.5	0	0.5	0	0	0
0.5	0	0.5	1	0.5	
0.5	0	0.5	0	0	
0.5	1	-0.5	1	-0.5	

- b) What will be the values of θ_0 , θ_1 and θ_2 at (t+1) with learning rate $\alpha=1$ and L2 regularization constant $\lambda=1$? [3M]

Solution:

Cost function

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

$\lambda \sum |\theta_j|$ ✓

Apply gradient descent update rule

y-hat ✓	y ✓	yhat-y ✓	x0 ✓	(yhat-y)x0	w0-new
0.5	0	0.5	1	0.5	-0.25
0.5	0	0.5	1	0.5	
0.5	0	0.5	1	0.5	
0.5	1	-0.5	1	-0.5	

$$\theta_0 = \theta_0 - \frac{\alpha}{m} [y - \hat{y}] \cdot x_0$$

$$\theta_0 = 0 - \frac{1}{4} \times 1$$

$$\theta_0 = -0.25$$

y-hat	y	yhat-y	x1	(yhat-y)x1	regularized w1-new
0.5	0	0.5	0	0	0
0.5	0	0.5	0	0	
0.5	0	0.5	1	0.5	
0.5	1	-0.5	1	-0.5	

$$\theta_1 = 0 - \frac{1}{4} \times 0 = 0$$

y-hat	y	yhat-y	x2	(yhat-y)x2	regularized w2-new
0.5	0	0.5	0	0	0
0.5	0	0.5	1	0.5	
0.5	0	0.5	0	0	
0.5	1	-0.5	1	-0.5	

$$\theta_2 = 0 - \frac{1}{4} \times 0 = 0$$

Q.4 We claim that there exists a value for α in the following data : (1.0, 4.0), (2.0, 9.0), (3.0, α) such that the line $y = 2 + 3x$ is the best least-square fit for the data. Is this claim true? If the claim is true, find the value of α . Otherwise, explain why the claim is false. Give **detailed** mathematical justification for your answer. [5 Marks]

For the line $y = a + bx$, the MSE turns out to be

$$(4.0 - (a + b))^2 + (9.0 - (a + 2b))^2 + (\alpha - (a + 3b))^2$$

$$\begin{aligned} \frac{\partial \text{MSE}}{\partial a} &= 0 \\ &\Rightarrow 2(4.0 - a - b)(-1) + 2(9.0 - a - 2b)(-1) \\ &\quad + 2(\alpha - a - 3b)(-1) = 0 \\ \Rightarrow 4.0 - a - b + 9.0 - a - 2b + \alpha - a - 3b &= 0 \\ \alpha + 13 - 3a - 6b &= 0 \\ \Rightarrow \alpha &= 3a + 6b - 13 \quad -\textcircled{1} \end{aligned}$$

$$\frac{\partial \text{MSE}}{\partial b} = 0$$

From ① and ②, the parameters a and b for the best fit must solve

$$\frac{6a + 14b - 22}{3} = 3a + 6b - 13$$

$$\Rightarrow 6a + 14b - 22 = 9a + 18b - 39$$

$$3a + 4b = 17$$

Substituting $a = 2$ and $b = 3$, we see that the above requirement is not met \Rightarrow the given claim is false

Marking Scheme: calculation of 1 and 2 – 3M

Equation of a and b = 1M

Final answer = 1M

Q.5 Consider a basis function $\phi_j(x) = x^j$, which is used to model nonlinear function of the input variables of the form $y(x, \theta) = \sum_{j=0}^2 \theta_j \phi_j(x)$. Determine θ_0, θ_1 and θ_2 for the table given below. [6 Marks]

x	y
0	1
1	3
2	7
5	31

Solution:

$$\text{Polynomial Regression: } y = \theta_0 + \theta_1 x + \theta_2 x^2 \quad [2M]$$

Solution: Method 1 [4M]

Q.5 Consider a basis function $\varphi_j(x) = x^j$, which is used to model nonlinear function of the input variables of the form $y(x, \theta) = \sum_{j=0}^2 \theta_j \varphi_j(x)$. Determine θ_0 , θ_1 , and θ_2 for the table given below.

x	y
0	1
1	3
2	7
5	31

Solution 1: Using Polynomial Regression:

$$y(x, \theta) = \sum_{j=0}^2 \theta_j \varphi_j(x) \quad [\varphi_j(x) = x^j]$$

$$y \Rightarrow \theta_0 x^0 + \theta_1 x^1 + \theta_2 x^2$$

$$y \Rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$$

Use the following equation to determine θ_0 , θ_1 , and θ_2

$$n a_0 + (\sum x_i) a_1 + (\sum x_i^2) a_2 = \sum y_i \rightarrow ①$$

$$(\sum x_i) a_0 + (\sum x_i^2) a_1 + (\sum x_i^3) a_2 = \sum x_i y_i \rightarrow ②$$

$$(\sum x_i^2) a_0 + (\sum x_i^3) a_1 + (\sum x_i^4) a_2 = \sum x_i^2 y_i \rightarrow ③$$

x	y	x^2	x^3	x^4	$x * y$	$x^2 * y$
0	1	0	0	0	0	0
1	3	1	1	1	3	3
2	7	4	8	16	14	28
5	31	25	125	625	155	775
Sum	42	30	134	642	172	806

$n = 4$

$\sum x_i = 8$ $\sum x_i y_i = 172$
 $\sum y_i = 42$ $\sum x_i^2 y_i = 806$
 $\sum x_i^2 = 30$
 $\sum x_i^3 = 134$
 $\sum x_i^4 = 642$

$n = 4$

Substitute the above values in all 3 equation

$$4a_0 + 8a_1 + 30a_2 = 42$$

$$8a_0 + 30a_1 + 134a_2 = 172$$

$$30a_0 + 134a_1 + 642a_2 = 806.$$

Coefficient Matrix

$$\begin{matrix} R_1 & \left[\begin{array}{ccc} 4 & 8 & 30 \end{array} \right] \\ R_2 & \left[\begin{array}{ccc} 8 & 30 & 134 \end{array} \right] \\ R_3 & \left[\begin{array}{ccc} 30 & 134 & 642 \end{array} \right] \end{matrix} \left[\begin{array}{c} a_0 \\ a_1 \\ a_2 \end{array} \right] = \left[\begin{array}{c} 42 \\ 172 \\ 806 \end{array} \right]$$

constant Matrix

Step 1 :- $R_1 \rightarrow R_1 / 4$ $R_{11} (4/4 \Rightarrow 1) \quad R_{12} (8/4 \Rightarrow 2)$

$$\begin{matrix} R_1 & \left[\begin{array}{ccc|c} 1 & 2 & 7.5 & 10.5 \end{array} \right] \\ R_2 & \left[\begin{array}{ccc|c} 0 & 8 & 30 & 172 \end{array} \right] \\ R_3 & \left[\begin{array}{ccc|c} 0 & 30 & 134 & 806 \end{array} \right] \end{matrix} \quad \begin{matrix} R_{13} (30/4 \Rightarrow 7.5) \\ C_1 \rightarrow 42/4 \Rightarrow 10.5 \end{matrix}$$

Step 2 :- $R_2 \rightarrow R_2 - R_{21} * R_1$ $R_{21} \rightarrow (8 - 8 * 1 \Rightarrow 8 - 8 \Rightarrow 0)$

$$\begin{matrix} R_1 & \left[\begin{array}{ccc|c} 1 & 2 & 7.5 & 10.5 \end{array} \right] \\ R_2 & \left[\begin{array}{ccc|c} 0 & 0 & 14 & 88 \end{array} \right] \\ R_3 & \left[\begin{array}{ccc|c} 0 & 30 & 134 & 806 \end{array} \right] \end{matrix} \quad \begin{matrix} R_{22} \rightarrow (30 - 8 * 14 \Rightarrow 30 - 112 \Rightarrow 14) \\ R_{23} \rightarrow (134 - 8 * 14 \Rightarrow 134 - 112 \Rightarrow 22 \Rightarrow 74) \\ C_2 \Rightarrow (172 - 8 * 10.5 \Rightarrow 172 - 84 \Rightarrow 88) \end{matrix}$$

$$\begin{array}{l}
 R_1 \left[\begin{array}{ccc|c} 1 & 2 & 7.5 & 10.5 \end{array} \right] \\
 R_2 \left[\begin{array}{ccc|c} 0 & 1 & 7.4 & 8.8 \end{array} \right] \\
 R_3 \left[\begin{array}{ccc|c} 0 & 7.4 & 41.7 & 49.1 \end{array} \right]
 \end{array}$$

Step 4: $R_2 \rightarrow R_2 | R_{22}$ $R_{21} \Rightarrow 0/14 \Rightarrow 0$
 $R_1 \left[\begin{array}{ccc|c} 1 & 2 & 7.5 & 10.5 \end{array} \right]$ $R_{22} \Rightarrow 14/14 \Rightarrow 1$
 $R_2 \left[\begin{array}{ccc|c} 0 & 1 & 5.285714 & 6.285714 \end{array} \right]$ $R_{23} \Rightarrow 74/14 \Rightarrow 5.285714$
 $R_3 \left[\begin{array}{ccc|c} 0 & 7.4 & 41.7 & 49.1 \end{array} \right]$ $C_2 \Rightarrow 88/14 \Rightarrow 6.285714$

Step 5: $R_3 \rightarrow R_3 - R_{32} * R_2$ $R_{31} \Rightarrow 0 - 74 * 0 \Rightarrow 0$
 $R_1 \left[\begin{array}{ccc|c} 1 & 2 & 7.5 & 10.5 \end{array} \right]$ $R_{32} \Rightarrow 74 - 74 * 1 \Rightarrow 0$
 $R_2 \left[\begin{array}{ccc|c} 0 & 1 & 5.285714 & 6.285714 \end{array} \right]$ $R_{33} \Rightarrow 41.7 - 74 * 5.285714$
 $R_3 \left[\begin{array}{ccc|c} 0 & 0 & 25.85714 & 25.85714 \end{array} \right]$ $\Rightarrow 25.85714$
 \downarrow $C_3 \Rightarrow 49.1 - 74 * 6.285714$
 $1 * a_0 + 2 * a_1 + 7.5 * a_2 = 10.5 \Rightarrow ①$
 $1 * a_1 + 5.285714 * a_2 = 6.285714 \Rightarrow ②$
 $25.85714 * a_2 = 25.85714 \Rightarrow ③$

Solve the above equation ③
 $25.85714(a_2) = 25.85714$
 $a_2 = 1$

Substitute $a_2 = 1$ in equation ②
 $a_1 + 5.285714 * 1 = 6.285714$
 $a_1 = 6.285714 - 5.285714$
 $a_1 = 1$

Substitute $a_1 = 1$ & $a_2 = 1$ in equation ①
 $a_0 + 2 * 1 + 7.5 * 1 = 10.5$
 $a_0 + 9.5 = 10.5$
 $a_0 = 10.5 - 9.5$
 $a_0 = 1$

$\boxed{a_0 = 1 ; a_1 = 1 ; a_2 = 1}$

Solution: Method 2

Using closed form solution: [4M]

$$\theta = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\theta_0 = 1, \theta_1 = 1 \text{ and } \theta_2 = 1$$

- Q.6 Consider the dataset of binary values in terms of attribute-value pairs where F is the value, and A,B, C are attributes. What is the entropy of the dataset? Fill in the columns for A and B, if it is known that A has maximum information gain and B has minimum information gain. Give mathematical justification for your answer. [4 Marks]

A	B	C	F
		0	0
		1	1
		0	1
		1	0
		0	1
		1	1
		0	1
		1	1

Solution:

For the entropy problem, the column F is the output attribute.

2 Marks:

Let column A = column F, so that the i th entries of the two column match each other. The information gain can be written as

$$InformationGain(S, A) = Entropy(S) - \sum \frac{|S_A|}{|S|} Entropy(S_A) \text{ where the sum is over the attribute values of A.}$$

Since the column entries of A match with those of F we see that the set $S_{A=0}$ is full of 0s and $S_{A=1}$ is full of 1s, so that $Entropy(S_{A=0})$ is 0 and so is $Entropy(S_{A=1})$. From the equation on Information Gain we can see that we get the maximum information gain possible in this case.

2 Marks:

The information gain with respect to column B can be written as $InformationGain(S, B) = Entropy(S) - \sum \frac{|S_B|}{|S|} Entropy(S_B)$.

For minimum information gain we see that if let the column B be the column of all 1s, then we have $S_{B=1} = S$ and $S_{B=0} = \emptyset$. Once again plugging this into the information gain equation shows that the information gain with respect to B is 0.

The arguments above work for maximum information gain when A is taken to be complement of F, and B is taken to be all zeroes rather than all 1s.

- (1) It is known that a natural law obeys the quadratic relationship $y = ax^2$. What is the best linear curve that can be used to model this data if all of the data points are drawn uniformly at random in the interval $(0,1)$?

Let $y = px + q$ be the equation of the linear curve

$$\text{We need to minimize } L = \int (ax^2 - px - q)^2 dx$$

$$L = \frac{a^2}{3} + f_3^2 + q^2 - \frac{2af}{4} - \frac{2aq}{3} + \frac{2pq}{2}$$

$$\frac{\partial L}{\partial p} = 0, \frac{\partial L}{\partial q} = 0 \Rightarrow \frac{2f}{3} + q = \frac{a}{2}; 2q + p = \frac{2a}{3}$$

$$\Rightarrow p = a, q = -\frac{a}{6}$$

Marking Scheme

Setting up the integral $\rightarrow 1$ mark

Getting this expression $\rightarrow 2$ marks

Final solution $\rightarrow 2$ marks

- (2) Consider the following dataset for text classification where three training instances are given with corresponding classifications into the '+' or '-' category:

Hindi India India	+
India Kannada Hindi	+
Chinese Hindi India	-

Showing all intermediate calculations, find the appropriate classification for the test instance: Chinese Kannada Chinese using the Naïve Bayes text classification algorithm.

First we observe $P(+)=\frac{2}{3}$ and $P(-)=\frac{1}{3}$

$\text{docs}_+ = \text{Hindi India India India Kannada Hindi}$

$\text{docs}_- = \text{Chinese Hindi India}$

Vocabulary = {Hindi, India, Kannada, Chinese}

$$P(\text{Hindi}/+) = \frac{2+1}{6+4} = \frac{3}{10} \quad P(\text{Hindi}/-) = \frac{1+1}{3+4} = \frac{2}{7}$$

$$P(\text{India}/+) = \frac{3+1}{6+4} = \frac{4}{10} \quad P(\text{India}/-) = \frac{1+1}{3+4} = \frac{2}{7}$$

$$P(\text{Kannada}/+) = \frac{1+1}{6+4} = \frac{2}{10} \quad P(\text{Kannada}/-) = \frac{0+1}{3+4} = \frac{1}{7}$$

$$P(\text{Chinese}/+) = \frac{0+1}{6+4} = \frac{1}{10} \quad P(\text{Chinese}/-) = \frac{1+1}{3+4} = \frac{2}{7}$$

$$\text{Decision } P(+)/P(\text{Chinese}/+)\ P(\text{Kannada}/+)\ P(\text{Chinese}/+)$$

$$P(-)/P(\text{Chinese}/-) P(\text{Kannada}/-) P(\text{Chinese}/-)$$

$$\frac{2}{3} \times \frac{1}{10} \times \frac{2}{10} \times \frac{1}{10} \text{ vs } \frac{1}{3} \times \left(\frac{2}{7}\right) \times \frac{1}{7}$$

$$0.00132 \text{ vs } 0.0038$$

Marking Scheme:

Positive Conditional Probabilities → 1.5 mark

Negative Conditional Probs → 1.5 mark

Decision - 2 marks

Q.3. There are two varieties of cucumbers – C_1 and C_2 which have different distributions of length. The joint probability density function of the length of the cucumber and category 1 is denoted by $p(x, C_1)$, and is a uniform distribution over the range (10cm, 30cm). Similarly $p(x, C_2)$ is a uniform distribution over the range (20cm, 50cm). What is the error of classification we will make if we assert that all cucumbers of length less than 25cm are of Variety 1 and all cucumbers of length greater than 25cm are of Variety 2?

[5]

Marks]

The information in this question is
incompletely specified

The functions given for $p(x, c_1)$ and $p(x, c_2)$ in the question are actually $p(x/c_1)$ and $p(x/c_2)$ respectively. To get $p(x, c_1)$ and $p(x, c_2)$ we need to multiply $p(x/c_1)$ and $p(x/c_2)$ by $p(c_1)$ and $p(c_2)$ respectively.
 $p(c_1)$ and $p(c_2)$ are not specified in the question, so $p(\text{mistakes})$ should be calculated in terms of $p(c_1)$ and $p(c_2)$.

We have

$$\begin{aligned} p(\text{mistakes}) &= \int_{R_1} p(x, c_2) dx + \int_{R_2} p(x, c_1) dx \\ &= p(c_2) \int_{R_1} p(x/c_2) dx + p(c_1) \int_{R_2} p(x/c_1) dx \\ &= p(c_2) \int_{20}^{25} p(x/c_2) dx + p(c_1) \int_{25}^{30} p(x/c_1) dx \end{aligned}$$

$$= p(c_2) \frac{5}{30} + p(c_1) \frac{5}{20}$$

$$= p(c_2) \frac{1}{6} + p(c_1) \frac{1}{4}$$

Marking Scheme

Formula for $p(\text{mistake}) = 3 \text{ marks}$

Final Calculation = 2 marks

- (3) Consider the standard set of Gaussian Naïve Bayes assumptions used in the derivation of the logistic regression expression, but with one modification – the class conditional density for each class has unique values for both the mean and variance, rather than a common value for the variance, i.e $P(X_i|Y = y_k) = N(\mu_{ik}, \sigma_{ik}^2)$. Find the expression for $P(Y = 1|X_1, X_2, \dots, X_n)$ in this case and find the decision boundary.

$$P(Y = 1|X) = \frac{P(Y = 1)P(X|Y = 1)}{P(Y = 1)P(X|Y = 1) + P(Y = 0)P(X|Y = 0)}$$

$$= \frac{1}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}}.$$

$$= \frac{1}{1 + \exp(\ln \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)})}$$

$$= \frac{1}{1 + \exp((\ln \frac{1-\pi}{\pi}) + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)})}$$

$$\text{Now } \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)} = -\frac{1}{2} \left(\frac{(X_i - \mu_{i0})^2}{\sigma_{i0}^2} \right) + \frac{1}{2} \left(\frac{(X_i - \mu_{i1})^2}{\sigma_{i1}^2} \right)$$

$$= \frac{1}{2} \frac{X_i^2}{\sigma_{i1}^2} - \frac{1}{2} \frac{X_i^2}{\sigma_{i0}^2} - \frac{\mu_{i1} X_i}{\sigma_{i1}^2} + \frac{\mu_{i0} X_i}{\sigma_{i0}^2} + \frac{\mu_{i0}^2}{\sigma_{i0}^2} - \frac{\mu_{i1}^2}{\sigma_{i1}^2}$$

This will give rise to an expression of the form $P(Y = 1|X) = \frac{1}{1 + \exp(\omega_0 + \sum_i \omega_i x_i + \alpha x_i^2)}$

$$\text{where } \alpha = \frac{1}{2} \left(\frac{1}{\sigma_{ii}^2} - \frac{1}{\sigma_{i,i}^2} \right)$$

The decision boundary is quadratic instead
of linear

Marking Scheme

Derivation for the new expression
= 4 marks

Decision boundary quadratic \rightarrow 1 mark

(4) Consider the following dataset

price	maintenance	capacity	Safety measures	Beneficial
lowpriced	cheap	5	yes	yes
lowpriced	average	5	yes	yes
lowpriced	cheap	5	yes	no
lowpriced	excessive	3	no	no
fair	average	5	no	no
fair	average	5	no	yes
fair	excessive	3	yes	no
overpriced	average	5	yes	yes
overpriced	excessive	3	yes	no
overpriced	excessive	6	yes	no

Classify the new instance given: "price = fair, maintenance = cheap, capacity = 5, safety measures = yes". Use Laplace smoothing only when needed for an attribute.

$$\text{We have } P(\text{Beneficial} = \text{Yes}) = \frac{5}{11}$$

$$P(\text{Beneficial} = \text{No}) = \frac{6}{11}$$

$$P\left(\frac{\text{Price} = \text{fair}}{\text{Yes}}\right) = \frac{2}{5} \quad \left| \begin{array}{l} P\left(\frac{\text{Price} = \text{fair}}{\text{No}}\right) = \frac{2}{6} \\ P\left(\frac{\text{Price} = \text{excessive}}{\text{No}}\right) = \frac{4}{6} \end{array} \right.$$

$$P\left(\frac{\text{maint} = \text{cheap}}{\text{Yes}}\right) = \frac{1}{5} \quad \left| \begin{array}{l} P\left(\frac{\text{maint} = \text{cheap}}{\text{No}}\right) = \frac{1}{6} \\ P\left(\frac{\text{maint} = \text{average}}{\text{No}}\right) = \frac{5}{6} \end{array} \right.$$

$$P\left(\frac{\text{capacity} = 5}{\text{Yes}}\right) = \frac{4}{5} \quad \left| \begin{array}{l} P\left(\frac{\text{capacity} = 5}{\text{No}}\right) = \frac{2}{6} \\ P\left(\frac{\text{capacity} = 3}{\text{No}}\right) = \frac{4}{6} \end{array} \right.$$

$$P\left(\frac{\text{Safety} = \text{Yes}}{\text{Yes}}\right) = \frac{4}{5} \quad \left| \begin{array}{l} P\left(\frac{\text{Safety} = \text{Yes}}{\text{No}}\right) = \frac{4}{5} \\ P\left(\frac{\text{Safety} = \text{No}}{\text{No}}\right) = \frac{1}{5} \end{array} \right.$$

Compute $P(\text{Yes}) P\left(\frac{\text{fair}}{\text{Yes}}\right) P\left(\frac{\text{cheap}}{\text{Yes}}\right) P\left(\frac{5}{\text{Yes}}\right) P\left(\frac{\text{yes}}{\text{Yes}}\right)$
 with $P(\text{No}) P\left(\frac{\text{fair}}{\text{No}}\right) P\left(\frac{\text{cheap}}{\text{No}}\right) P\left(\frac{5}{\text{No}}\right) P\left(\frac{\text{yes}}{\text{No}}\right)$

$$\frac{5}{11} \times \frac{2}{5} \times \frac{1}{3} \times \frac{4}{7} \times \frac{4}{5} \textcircled{vs} \frac{5}{11} \times \frac{2}{7} \times \frac{1}{6} \times \frac{2}{7} \times \frac{4}{6}$$

$$\underline{0.023} \text{ vs } 0.0067$$

The inference should be classified as Y_0

Marking Scheme

$$P(Y_i), P\left(\frac{X_i}{Y_i}\right) \rightarrow 2 \text{ Marks}$$

$$I(N_0), I\left(\frac{X_i}{N_0}\right) \rightarrow 2 \text{ Marks}$$

final decision $\rightarrow 1 \text{ Mark}$

- (5) Consider the dataset in terms of attribute-value pairs where F is the value, and A,B, C are attributes. What is the entropy of the dataset? Compute the information gain with respect to the attributes A, B and C. Which attribute would be used at the root of the decision tree constructed by the ID3 algorithm?

A	B	C	F
0	0	0	0
0	0	1	1
0	1	0	1
0	1	1	0
1	0	0	1
1	0	1	1
1	1	0	1
1	1	1	1

$$\text{Entropy}(S) = -p_1 \log_2 p_1 - p_2 \log_2 p_2 -$$

$$= -\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8}$$

$$\text{Gain}(S, A) = \text{Entropy}(S) - \frac{|S_{A=0}|}{|S|} \text{Entropy}(S_{A=0})$$

$$- \frac{|S_{A=1}|}{|S|} \text{Entropy}(S_{A=1})$$

$$= \left(-\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8} \right) - \frac{4}{8} \left(\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4} \right)$$

$$= 0.81 - 0.5 = \underline{\underline{0.31}}$$

$$\text{Gain}(S, B) = \text{Entropy}(S) - \frac{|S_{B=0}|}{|S_B|} \text{Entropy}(S_{B=0})$$

$$\frac{S_B=1}{|S|} \text{ Entropy}(S_B=1)$$

$$= \left(-\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8} \right) - \frac{4}{8} \left(-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \right)$$

$$- \frac{4}{8} \left(-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \right) = 0.81 - 0.81 = 0$$

$$\text{Gain}(S, C) = \left(-\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8} \right) - \frac{4}{8} \left(-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \right)$$

$$- \frac{4}{8} \left(-\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} \right) = 0.81 - 0.81 = 0$$

$$\text{Gain}(S, A) > \text{Gain}(S, B) = \text{Gain}(S, C)$$

\Rightarrow A is the root attribute

Marking Scheme

Entropy of data set \rightarrow 2 Marks
 Gains of A, B, C \rightarrow 2 Marks
 Root attribute \rightarrow 1 Mark

Format of Question paper for Mid-Semester Test

**Birla Institute of Technology & Science, Pilani
Work Integrated Learning Programmes Division
Second Semester 2020-2021**

Mid-Semester Test (EC-2 Regular/ EC-2 Make-up)

Course No.	:	
Course Title	:	
Nature of Exam	:	Open Book
Weightage	:	30% or 35% (As per Course Handout)
Duration	:	2 Hours
Date of Exam	:	06/03/2021 or 19/03/2021 (FN/AN)
		No. of Pages = No. of Questions =

Note to Students:

1. Please follow all the *Instructions to Candidates* given on the cover page of the answer book.
2. All parts of a question should be answered consecutively. Each answer should start from a fresh page.
3. Assumptions made if any, should be stated clearly at the beginning of your answer.

Please Note:

A maximum of 12 main questions with Q.Nos. 1 to N. Please do not include any objective type questions. Please do not include any reference material like data tables in the question paper, as all Mid-Semester Test are completely Open Book examinations.

- Q.1. We have five data points $x_1 = 1, x_2 = 3, x_3 = -1, x_4 = 4, x_5 = -3$ which are obtained from sampling a Gaussian distribution of zero mean. What is the Maximum Likelihood Estimate of the variance of the Gaussian distribution? Show all the steps in the calculation. [5 Marks]

We are sampling from a normal distribution $N(0, \sigma^2)$ when $\sigma \approx 5$

The likelihood is $L = \prod_{i=1}^5 \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \frac{x_i^2}{\sigma^2}}$

$$L = \left(\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \frac{x_1^2}{\sigma^2}} \right) \left(\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \frac{x_2^2}{\sigma^2}} \right) \dots \left(\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \frac{x_5^2}{\sigma^2}} \right)$$

$$L = \frac{1}{(2\pi)^{5/2} \sigma^5} e^{-\frac{1}{2} \frac{(x_1^2 + x_2^2 + \dots + x_5^2)}{\sigma^2}}$$

To find the value of σ that maximizes likelihood we have to take the derivative with respect to σ of L and equate it to

6.
It is more convenient to work with $\log L$

$$\begin{aligned} \frac{\partial}{\partial \sigma} (\log L) &= 0 \\ \Rightarrow \frac{1}{\sigma} \left(\log \frac{1}{(2\pi)^{1/2}} + \log \frac{1}{\sigma^5} + -\frac{1}{2} \frac{(x_1^2 + x_2^2 + \dots + x_5^2)}{\sigma^2} \right) \\ &= 0 \\ \Rightarrow \frac{1}{\sigma} \left(-5 \log \sigma - \frac{x_1^2 + x_2^2 + \dots + x_5^2}{2\sigma^2} \right) &= 0 \\ \Rightarrow -\frac{5}{\sigma} + \frac{1}{\sigma^3} (x_1^2 + x_2^2 + \dots + x_5^2) &= 0 \\ \Rightarrow 5\sigma^2 &= (x_1^2 + x_2^2 + \dots + x_5^2) \\ \Rightarrow \sigma^2 &= \frac{x_1^2 + x_2^2 + \dots + x_5^2}{5} \\ \sigma^2 &= \frac{1^2 + 3^2 + (-1)^2 + (-3)^2 + (4)^2}{5} \end{aligned}$$

$$\sigma^2 = \frac{36}{5} = 7.2$$

Thus the maximum likelihood estimate
of the variance of the Gaussian
distribution is 7.2

Marking Scheme

2 Marks → Setting up max likelihood expression

3 Marks - differentiation and final calculation

- Q.2. Consider a table with a single attribute "wind" and category "rain", where "wind" can take two attribute values – high and low, and "rain" has two classes – yes and no. There are 10 entries in the table, and it is known that 8 entries in the table have wind=high. It is also known that 8 entries in the table also have rain=yes. What is the highest and lowest possible information gain if we split the table on the attribute "wind"? [5 Marks]

For the highest information gain the table looks like this:

Wind	Rain
High	Yes
Low	No
Low	No

$$\text{Gain}(S, \text{Wind}) = \text{Entropy}(S) - \frac{|S_{\text{W=H}}| \text{Entropy}(S_{\text{W=H}})}{|S|} - \frac{|S_{\text{W=L}}| \text{Entropy}(S_{\text{W=L}})}{|S|}$$

$$\text{Entropy}(S) = -P_1 \log P_1 - P_0 \log P_0 =$$

$$= -\frac{8}{10} \log \frac{8}{10} - \frac{2}{10} \log \frac{2}{10} = 0.721$$

$$\text{Entropy}(S_{W=H}) = -\frac{8}{8} \log \frac{8}{8} - \frac{0}{8} \log 0 = 0$$

$$\text{Entropy}(S_{W=L}) = \frac{0}{2} \log \frac{0}{2} - \frac{2}{2} \log 1 = 0$$

Highest Information
gain: $0.721 - \frac{8}{10} \times 0 - \frac{2}{10} \times 0 = \underline{0.721}$

For the lowest information gain the table looks like this:

Wind	Rain
High	No
High	No
High	Yes
Low	Yes
Low	Yes

In this case

$$\text{Gain}(S, \text{Wind}) = \text{Entropy}(S) - \frac{|S_{W=H}|}{|S|} \text{Entropy}(S_{W=H})$$

$$- \frac{|S_{W=L}|}{|S|} \text{Entropy}(S_{W=L})$$

$$\text{Entropy}(S_{\omega=H}) = -\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8} = 0.8075$$

$$\text{Entropy}(S_{\omega=L}) = -\frac{2}{2} \log \frac{2}{2} - \frac{0}{2} \log \frac{0}{2} = 0$$

$$\text{Gain}(S, \text{Wind}) = 0.721 - \frac{8}{10} \times 0.8075 - \frac{2}{10} \times 0 \\ = \underline{0.075}$$

Marking Scheme

Highest Information Gain = 2.5 marks

Lowest Information Gain = 2.5 marks

- Q.3. What is the best curve of the form $y = a + bx + cx^2$ in terms of minimizing square error that fits the following data of the form (x, y) : $(-1, 0)$, $(1, 10)$, $(2, 24)$, $(-2, 4)$? [5 Marks]

Computing the square loss function gives

$$L = (a - b + c - 0)^2 + (a + b + c - 10)^2 \\ + (a + 2b + 4c - 24)^2 + (a - 2b + 4c - 4)^2$$

Setting $\frac{\partial L}{\partial a} = 0$, $\frac{\partial L}{\partial b} = 0$ and $\frac{\partial L}{\partial c} = 0$

We get

$$4a + 16c = 38$$

$$6a + 16b + 4c = 100$$

$$2a + 68c = 244$$

Final answer: $a = 2, b = 5, c = 3$

Marking Scheme

Setting up loss function = 2 marks

Final calculation = 3 marks

- Q.4. There exists a training set consisting of 100 documents for text classification consisting of two types of document '+' and '-'. 75 of the 100 documents are '+' and the remaining are '-'. The total number of words including duplicates in the '+' documents is 150 and the total number of words including duplicates in the '-' documents is 100. The number of words in the vocabulary is 1000. What classification is given to a test text with 5 words consisting of words belonging to the vocabulary but which have not occurred in the training set at all using the Naïve Bayes algorithm for text classification?

[5 Marks]

$$\text{We have } P(+)=\frac{75}{100}=\frac{3}{4}$$

$$P(-)=\frac{25}{100}=\frac{1}{4}$$

Since the words ω_k in the test document do not occur in the training set, we

$$\begin{aligned} \text{have } P\left(\frac{\omega_k}{+}\right) &= \frac{n_k + 1}{n + |V|} = \frac{0 + 1}{150 + 100} \\ &= \frac{1}{1150} \end{aligned}$$

$$P\left(\frac{\omega_k}{-}\right) = \frac{n_k + 1}{n + |V|} = \frac{1}{100 + 100} = \frac{1}{100}$$

$$P(+)^{\prod_{k=1}^5 P\left(\frac{\omega_k}{+}\right)} = \frac{3}{4} \left(\frac{1}{1150}\right)^5$$

$$P(-)^{\prod_{k=1}^5 P\left(\frac{\omega_k}{-}\right)} = \frac{1}{4} \left(\frac{1}{100}\right)^5$$

$$\frac{3}{4} \left(\frac{1}{1150}\right)^5 > \frac{1}{4} \left(\frac{1}{1100}\right)^5$$

∴ we give a true classification to the new text.

Marking Scheme
3 marks → setting up $P\left(\frac{\omega_K}{+}\right)$ and $I\left(\frac{\omega_K}{-}\right)$
2 marks → final calculation

Q.5. One percent of women over 50 have breast cancer. Ninety percent of women who have breast cancer test positive on mammograms. Eight percent of women will have false positives. What is the probability that a woman has cancer if she has a positive mammogram result?

[5 Marks]

Let C denote the event of cancer for women over 50.
 $P(C) = 0.01$
 $I(M/C) = 0.9$

$P(M/c)$ = probability of showing a positive mammogram but given no cancer = false positive = 0.08

$$P(c/M) = ?$$

$$P(c/M) = \frac{P(M/c) P(c)}{P(M/c) P(c) + P(M/\bar{c}) P(\bar{c})}$$

$$= \frac{(0.1)(0.01)}{(0.08)(0.99) + (0.9)(0.01)}$$

$$= \frac{0.001}{0.0882} = 0.102$$

Marking scheme :

2 Marks \rightarrow calculating all the probabilities
 $P(c), P(\bar{c}), P(M/c), P(M/\bar{c})$

3 Marks \rightarrow Bayes equation

Q.6. Let X_1, X_2 be two real-valued features and Y be a Boolean-valued function of the given features such that the Gaussian Naïve-Bayes assumptions are satisfied. Suppose

$P(Y/X_1, X_2) = \frac{1}{1 + \exp(-0.1 - 0.2X_2 - 0.3X_3)}$. Assume that $P(X_1/Y=0) = N(1.0, \sigma_1)$ and

$P(X_1/Y=1) = N(2.0, \sigma_1)$. Similarly $P(X_2/Y=0) = N(1.0, \sigma_2)$ and $P(X_2/Y=1) = N(2.0, \sigma_2)$. Calculate the standard deviations σ_1 and σ_2 and the probability $P(Y=1)$. [5 Marks]

From the formula for logistic regression we see that

$\frac{\mu_{j0} - \mu_{i0}}{\sigma_j}$ = coefficient of X_i in the formula $P(Y=1/X) = \frac{1}{1 + \exp(\theta_0 + \sum \theta_i X_i)}$

$$\therefore \frac{1.0 - 2.0}{\sigma_1^2} = -0.2 \Rightarrow \sigma_1^2 = 5$$

$$\frac{1.0 - 2.0}{\sigma_2^2} = -0.3 \Rightarrow \sigma_2^2 = \frac{10}{3}$$

Also $\ln \left(\frac{1-\pi}{\pi} \right) + \frac{\bar{\mu}_{11} - \bar{\mu}_{10}}{2\sigma_1^2} + \frac{\bar{\mu}_{21} - \bar{\mu}_{20}}{2\sigma_2^2}$

= constant term where $\pi = \text{prior probability}$
 $= P(Y=1)$

$$\ln \frac{1-\pi}{\pi} + \frac{4-1}{2 \times 5} + \frac{4-1}{2 \times \frac{10}{3}} = -0.1$$

$$\ln \frac{1-\pi}{\pi} = -0.1 - \frac{3.0}{10} - \frac{9}{20} = -0.85$$

$$\frac{1-\pi}{\pi} = e^{-0.85} \quad \frac{1}{\pi} = 1 + e^{-0.85} \quad \Rightarrow \pi = \frac{1}{1 + e^{-0.85}}$$

$$\frac{\mu_{10} - \mu_{11}}{\sigma_i} = \text{Coefficient of } x_i \text{ in the formula } P(Y=1/x_1, x_2) = \frac{1}{1 + \exp(\theta_0 + \sum \theta_i x_i)}$$

Marking Scheme

2 Marks \rightarrow getting the formulae right

3 Marks \rightarrow final calculation