Name: Rutvik Kolhe

Student ID: 200258232

Unity ID: rkolhe

	Question 1:						
	step1: We find the entropy of training sel- say T, using the following formula:						
	H[T]=-P+10g=P+-P-10g=P-						
	Here positive class: Buy negative class: Not Buy.						
	1. P+ = = P- = 4 9 9						
	H[T]= -5 log 5 - 4 log 4 9						
	H[T] = 0,991.						
	Step 2: Now, for each attribute trai- divides Timo subsets Ti, we calculate						
	-4. Outopil Apr each cubock						
->	Consider Age: (n: no goccuences in T) Age, n P+ P- entropy						
377	Yery young 2 1/2 /2 -1/2 10921/2 - 1/2 10921/2 = 1						
	Young 1 0 1 0						
	Middle 3 1/3 = 1/3 1092 1/3 - 2/3 1092 2/3 = 0.918						
	Very old 2 2/2 0 0						
197	Very old 2 2/2 0 0						

Now, calculate arecage entropy Avceage entropy [H(T,age)]= Formula for calculating average entropy:-H(T, a) = = P/ × H(T)) :. H[T, age] = 2(1)+0+3+0+0 H[T, age] = 0.528 Consider Ticker type: Ticker type 2/4 Long Lo cal Short H[T, Ticker +ype] = A(1) + 0 + 4(1) H [T, Tilker-type] = 0.888

	Consider La	ngu	age,	feegu	ency	
	Language	h	p+	P-	2ntropy	
	Fluent Not fluent	<u>s</u>	3/3	0 2/3	0.918	
R.J.	Accent-	2		1/2		
	Foreign		1/2		0	
	foreign		<u> </u>	Ò	. 0	
	: H[T, language] = 0+ 3 (0.912) Frequency 9					
9 -	Fre	quen	cy		9	
	$ \frac{1}{9} + \frac{2}{9} + 0 $ $ \frac{1}{9} + \frac{2}{9} + \frac{2}{9} + 0 $ $ \frac{1}{9} + \frac{2}{9} + \frac$					
	Ψ	V		3,230		
->	-> consider Type of call					
	U Company of the comp					
_	Local Local	n	PH	P-	Entopy	
	Local	4	2/4	214		
	Long distance	2	2/2	0	O	
in Lings	Intern	3		43	0-918	
25.52					Sapera)	
	H[T, Type of call]: 4(1)+0+3(0.918)					
10200	- dime	V				
	1. 5				V (**)	
	H[T, type of cau] = 0.450					
	_	V				
		7			317	
				. M		

Step 3:- Now that we have calculated average entropy for each attribute, we calculate information gain using 1-I(T,a) = H[T] - H[T,a] I(T; age) = 0.991 - 0.528 = 0.463 I(T; Ticket type) = 0.991 - 0.888 = 0.103 I(T, Language) = 0.991-0.528 = 0.463 frequency) ILT, type of call) = 0.991 - 0.750 = 0.2411 Ans!- The ganting (lowest to highest)

of attributes based on information

l gain 191-Ticket Type > Type of call

> Age, Dlanguage Frequency

Question 2:

stepl: center the data by subtracting mean of the attribute from each value

. The new table wated is:

χ_2
2-
-3
-4
-6
- 2

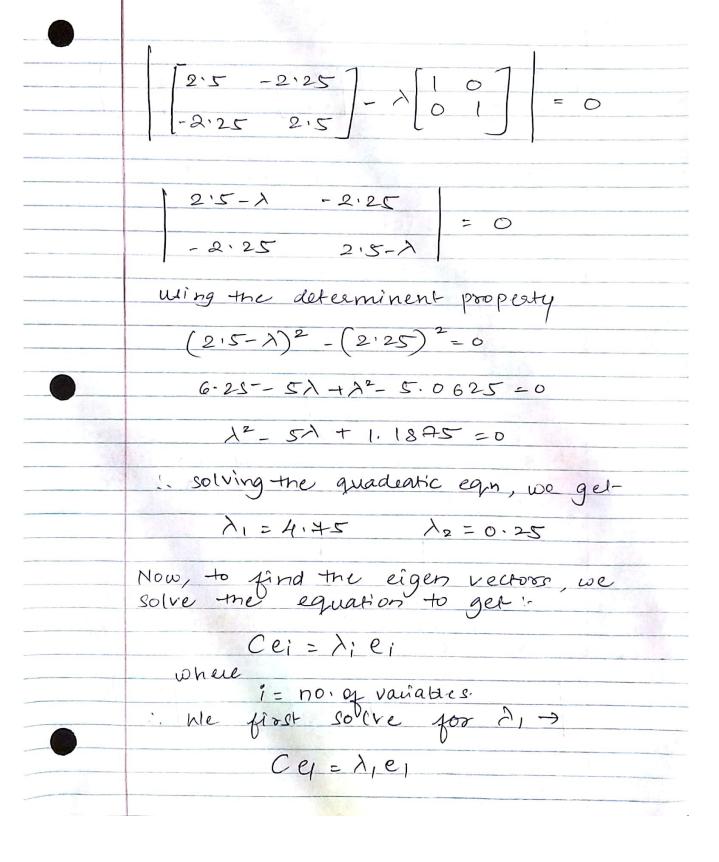
n= 5

Step II, we compute the covariance using:

Cij = 1 \(\int (\times im - \times i) \) (\times im - \times j)

$$cov(x_1, x_1) = \frac{1}{4}(1^2 + 0 + 1^2 + 2^2 + 2^2)$$

 $Cov(2_2, x_2) = \frac{1}{4} \left[(-1)^2 + (1)^2 + 0 + (-2)^2 + (2)^2 \right]$ COV(72, 22): 2.5 : $cov(\alpha_1, 22) = 1[(n(-1) + 0 + 0 + (2)(-2) + (-2)(2)]$ (or (21, 22) = -2.25 Hence, we get the covaciance matrix c: $C = \begin{bmatrix} 2.5 & -2.25 \\ -2.25 & 2.5 \end{bmatrix}$ step III :- calculate eigen vectors: Eigen values of C are solutions of A to (C-)I/=0



$$\begin{bmatrix} 2.5 & -2.25 \\ -2.25 & 2.5 \end{bmatrix} \begin{bmatrix} e_{11} \\ e_{12} \end{bmatrix} = \frac{2.36}{1.75} \begin{bmatrix} e_{11} \\ e_{12} \end{bmatrix}$$

$$\therefore (2.5)(e_{11}) + (2.25)(e_{12}) = 4.75(e_{11})$$

$$2.25(e_{11}) = -2.25(e_{12})$$

$$\vdots e_{11} = -e_{12}$$

$$e_{11} = -e_{12}$$

$$\vdots e_{11} = -1$$

$$\vdots e_{1} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$\vdots To avoid multiple solutions, we converted the aunit vector by dividing the Eucledian norm q the vector $\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$

$$\vdots v_{11} = \sqrt{2}$$

$$\vdots e_{1} = \begin{bmatrix} -0.71 \\ 0.71 \end{bmatrix}$$

$$\vdots e_{1} = \begin{bmatrix} -0.71 \\ 0.71 \end{bmatrix}$$$$

Question 3 a):

Feature selection algorithms are used to obtain the best features from the given set, thereby decreasing the complexity of the dataset. There are different approaches to select these features of which one is the **Wrapper** based approach. This approach consists of three different algorithms:

- Best Subset Selection
- Forward Stepwise Selection
- Backward Stepwise Selection

As given in the question, we consider that the dataset consists of "p" covariates of which we want to select "k" predictors where k<<p. These algorithms are briefly described as follows:

Best Subset Selection:

- At first, we consider the value of k= 1,2,3...p
- For k=1, we find the relation of each feature with the output by finding the value of RSS for that particular feature. In the end we will get "p" number of models for k=1.
- ° Now, we find the model with the **least** value for RSS and name it as $\mu(1)$.
- \circ Similarly, we repeat the step till k=p. As a result, we will get "p" models- μ from (1 to p).
- $^{\circ}~$ From these models, we choose the best model using different techniques such as C_p , AIC, BIC, adjusted R^2

• Forward Stepwise Selection:

- Similar to the above algorithm, we find the bet model when k=1
- \circ Say the best model for k=1 is X_1 . Now when k=2, we find the model with X_1 as one of the feature and find another feature from the remaining set of features in the dataset.
- New features get added to the best model as the value of k increases, whereas the old features are not dropped.
- ° From the different models that we obtain in $\mu(1 \text{ to p})$ we choose the best model using different techniques such as C_p , AIC, BIC, adjusted R^2

Backward Stepwise Selection:

- In this methods, we consider the value of k=p and then start the computation. At this step we get a single model with p features as the best model
- Now at k=p-1, we choose the best "p-1" features to be part of the best model for k=1 (i.e. $\mu(p-1)$)
- \circ For k=p-2, we choose the best p-3 features from μ (p-1) which contains p-1 features only. The feature that was discarded in the previous step is not considered again.
- $^{\circ}$ This process continues till k=1 as we get "p" models (µ(p-1 to 1))
- From these models, we choose the best model using different techniques such as C_p, AIC, BIC, adjusted R²

From the above description, we see that the Best Subset Selection checks every possible combination to find the model with the least value of RSS. However, it is not the case with Forward/Backward stepwise selection techniques. In Forward Stepwise Selection method, the feature once not selected in the previous set of best models $\mu(k)$ (where value of k is smaller) cannot be a part of the model for consequent values of k. Similarly for Backward Stepwise Selection, the feature once discarded from being a part of the best model cannot occur in the consequent models. Hence, there may arise a case where these 2 methods ignore a feature that could potentially be a part of the best model. In contrast to this, the Best Subset Selection follows an exhaustive approach and checks the value of RSS for each

and every possible model from the given set of p covariates. Due to this the value of RSS differs, where **Best Subset Selection** method provides us with the smallest training RSS.

Question 3 b):

In the above question, we discussed about the Wrapper based approach to feature selection. Another approach is the Embedded approach, where the following two methods are widely used:

- 1. Ridge regression
- 2. LASSO

The **similarity** between these two is that they both are examples of Regularization method where an additional constraint is added to coefficient estimates so that they shift toward zero. This results in decreasing variance and hence improves the prediction

• Ridge Regression:

• It is a regularization method which minimizes the value of RSS by adding a constraint to the coefficient as follows:

$$\Rightarrow \sum_{i=1}^{N} \left(y_{i}^{2} - \beta_{o} - \sum_{j=1}^{N} \beta_{j}^{2} \chi_{ij}^{2}\right) + \lambda \leq \beta_{o}^{2}$$

$$\Rightarrow RSS + \lambda \leq \beta_{j}^{2}$$

 $^{\circ}$ We can see that ridge regression adds a penalty to the value of β . This, reduces the values of β towards zero. It's significance is that it minimizes the effect of certain features that are unimportant but have a high impact on the class variable/ target variable.

LASSO

• It is similar to the above method, but differs with respect to the penalty added to the coefficients. LASSO minimizes the RSS by the given formula:

$$\Rightarrow \sum_{i=1}^{N} \left(y_i - \beta_i - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + A \sum_{j=1}^{p} |\beta_j|$$

- \checkmark As mentioned above, the major **difference**, between the two methods is the constraint on β . In LASSO, the coefficient can be equal to zero in some cases, which is not possible in the the case of ridge regression
- \checkmark The occurrence of $\beta=0$, signifies that the feature with which it is associated need not be fitted in the model and hence is dropped. In this way, LASSO acts as a variable selection algorithm in some cases.
- ✔ This reduces the curse of dimensionality and provides us with better predicted values.
- ✔ Therefore, I would prefer LASSO because of the reasons stated above.