

# Quiz 1, ME-781, October 16, 2020

Max Marks: 50

Name:

Roll No:

2x5

1. Which data scale do the following variables belong to:

- Name of a person (Nominal)
- Cell number of a person (Ratio)
- Marks obtained in an exam (Ratio)
- Ranking of a student in his/her class (Ordinal)
- Rate of change of fever of a hospitalized patient (in deg C per minute) (Ratio)

2x3

2. Let set  $A=\{1,2\}$  and set  $B=\{-1,0,1\}$ . Then write the following sets:

- $C = A \times B \{(1,-1), (2,-1), (1,0), (2,0), (1,1), (2,1)\}$
- $D = B \times A\{(-1,1), (-1,2), (0,1), (0,2), (1,1), (1,2)\}$
- $D \cap C \{(1,1)\}$

3

3. Probability of heart disease (Y) is seen to be related to age (X1), gender(X2) and body mass index (X3) of a person. How would you write this relation, if Y has a quadratic relation with X (with no term of the relation having both X1 and X3 variables)

Without loss of generality the variable X2 can be taken as 0 or 1 for Male and Female respectively. Thus, the relation would be:

$$y = \beta_0 + \beta_2 X_2 + (\alpha_{11}) X_1^2 + (\beta_{11} + \beta_{12} X_2) X_1 + (\alpha_{31}) X_3^2 + (\beta_{31} + \beta_{32} X_2) X_3$$

3x2

4. In a random health checkup, a person is tested positive for a rare disease. It is found that on an average 1 person in million in the world have this disease. However, the medical test which was conducted is very accurate with 99% accuracy (i.e. Sensitivity = Specificity = 99%).

- What is the probability that the person, who has been tested positive by this medical test, has this disease. (PPV =  $0.99 \times 1e-6 / (0.99 \times 1e-6 + (1-0.99) \times (1-1e-6)) = 9.899e-05$ )
- However, if the same disease is prevalent in 10% of the population then what would be the probability that the person, who has been tested positive by this medical test, has the disease. (PPV =  $0.99 \times 1e-1 / (0.99 \times 1e-1 + (1-0.99) \times (1-1e-1)) = 0.91667$ )

2x3

5. Show (very briefly) whether the following are dissimilarity measures (or not) for two point (x1, y1) and (x2, y2) in a 2D space:

- $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$  a.) Yes, it is a dissimilarity measure
- $\sqrt{|(x_1 - x_2) \cdot (y_1 - y_2)|}$  b.) No.  $\because d(x, y) = 0 \nRightarrow x = y$
- $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} + \sqrt{|(x_1 - x_2) \cdot (y_1 - y_2)|}$  c.) No.  $\because d(x, z) \nless d(x, y) + d(y, z)$

2

6. Let  $y'_i$  be the value predicted by a linear regression model at point  $x_i$ . If  $y_i$  is the true value and  $y_{mean}$  be the mean of values (over the entire data set of n values), then find  $\sum_1^n (y_i - y'_i)^2$  if  $\sum_1^n (y_{mean} - y'_i)^2 = 45$  and  $\sum_1^n (y_i - y_{mean})^2 = 49$ .  
(answer=4)

$$\because y_{mean} = y'_{mean}, \therefore \sum_1^n (y_{mean} - y'_i)^2 = nVar(Y')$$

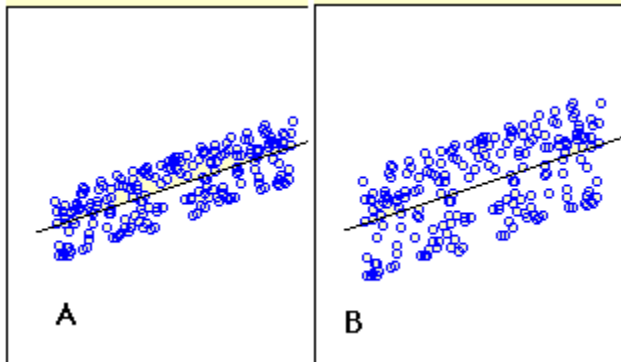
And  $(y_i - \hat{y}_i)$  and  $\hat{y}_i$  are orthogonal.)

2

- 3

Since the number of predictors  $p$  is large, we expect on an average  $p$ -value associated with 5% (i.e. about 2 to 3 predictors) would be below 0.05 by random chance, even if there is no relation between predictor variables and response.

9. Below graphs show two fitted linear regression lines (A & B) on a randomly generated data. Note: 1) Scale is same in both graphs for both axis. 2) X axis is independent variable and Y-axis is dependent variable.



- 2

- 2

11. Which of the following are true? Cosine similarity is always preferred over Euclidean distance similarity, if . (a is true)
- a. the scale of the magnitude of each datapoint is irrelevant.
  - b. the datapoints are not collinear.
  - c. the datapoints are collinear.
  - d. the datapoints are sparsely scattered in the space. (scattering need not be random)

Note: In all the above cases, assume that we can define both the similarity measures.

2x2

12. The relationship between number of beers consumed ( $x$ ) and blood alcohol content ( $y$ ) was studied in chimpanzee using least squares regression. The following regression equation was obtained from this study:  $Y = 0.0037 + 0.0180 x$

- a. Which of the following is implied by the above equation?
  - i. each beer consumed increases blood alcohol by 0.37%
  - ii. on average it takes 1.8 beers to increase blood alcohol content by 1%
  - iii. each beer consumed increases blood alcohol by an average of amount of 1.8%  
(True)
  - iv. each beer consumed increases blood alcohol by exactly 0.018
- b. The equation seems to suggest that even if no beer is consumed, there is alcohol present in the blood. What could be the reason for this? (suggest a reason based on data analytics)  
Biologically it is expected that the value of blood alcohol should be zero if no alcohol is consumed i.e.  $\beta_0 = 0$ . However, 95% of the times the estimated value of  $\beta_0$ , i.e.  $\hat{\beta}_0$  would be between  $\pm 2 SE(\hat{\beta}_0)$  which depends on the standard deviation of random error.