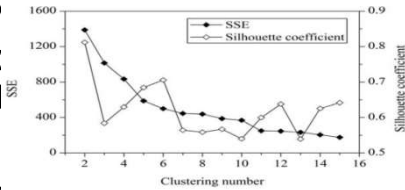


(iv) What should be the best choice for number of clusters based on the following results (as shown in Figure 1):



Time Allotted : 3 hrs

Figures out of the right margin inuicue juu marks.

Figure 1

Candidates are required to answer Group A and

any 5 (five) from Group B to E, taking at least one from each group.

Candidates are required to give following in their own words as far as practicable selection bias in a survey?

Group - A

(Multiple Choice Type Questions)

1. Choose the correct alternatives by the following:
 - (i) Using a random sample of students at a university to estimate the proportion of people who think the legal drinking age should be lowered.
 - (a) The mean is a measure of central tendency of the data
 - (b) The standard deviation is a measure of dispersion of the data
 - (c) The mode is a measure of central tendency of the data
 - (d) The range is a measure of dispersion of the data
 - (viii) In a fictitious state-wide database for road accidents, grouping of accidents involving a truck or accidents at night is an example of
 - (a) negatively low, negatively
 - (b) negatively high, positively
 - (c) positively low, negatively
 - (d) positively high, negatively
 - (ix) Which of the following machine learning algorithm is not used for inputting missing values of type categorical and continuous?
 - (a) Linear regression
 - (b) k-NN
 - (c) Logistic regression
 - (d) all of the mentioned.
 - (x) Consider the following confusion matrix (as shown in table below). What is the value of recall?

	Actual Yes	Actual No
Predicted Yes	100	10
Predicted No	5	50

 - (a) I and II
 - (b) I and III
 - (c) II and III
 - (d) I, II and III.
 - (iv) With respect to a typical 'Big Data Ecosystem', the technologies / tools that are not usually meant for 'Big Data Scientists' are _____, and security.
 - (a) benchmarking, deployment
 - (b) data integration, filesystem
 - (c) databases, scheduling
 - (d) programming, databases.

Group - B

2. (a) Compare Data Science (DS) with:
 - (i) Machine Learning (ML) - mention any two points; the former usually comes in two forms, 'Continuous' and _____
 - (ii) Databases (DB) - mention any two points; while Binary and _____ are two useful cases for the latter.
 - (iii) Business Intelligence (BI) - mention any one point.
 - (a) 'Discrete', 'Natural'
 - (b) 'Discrete', 'Ordinal'
 - (c) 'Discontinuous', 'Primal'
 - (d) 'Discontinuous', 'Secondary'.

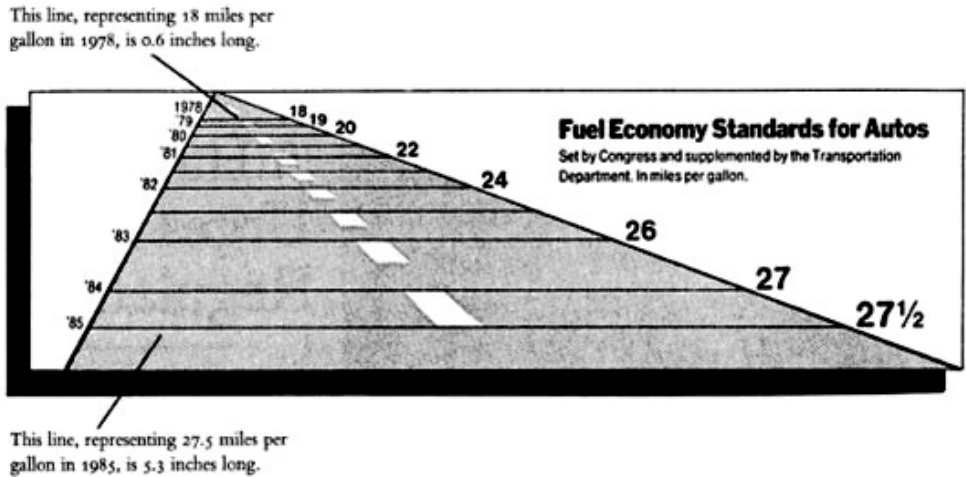
- (b) You are given a set of data **Group-C** of teenage boys and girls, aged between 16 and 19 years, residing in diverse localities within Kolkata, coming from various ethnic communities, studying in different schools under respective boards, and their respective overall percent marks in class-Board Examinations. You are expected to do some analysis on the data set to find, explain the kind, how the marks will marks might be equal or not by gender, ethnic community, residential locality, and school board. While carrying out the necessary 'data preparation' steps.
4. (a) Define 'quantile' and 'inter-quartile range (IQR)'. How are '50th percentile' and 'median' related?
 (b) What is 'expected value (EV)'? Work out one single EV of a Seminar attendee for the following case:
 (i) mention any two typical issues you might face.
 (ii) outline your possible approach(es) towards handling such issues with suitable examples.
 (iii) highlight the relative advantages and disadvantages thereof.
 A marketing company promotes a new cloud technology that offers two levels of service: one priced at USD 300 per month and another at USD 50 per month. It offers free seminars to generate leads, and it figures that 5% of the attendees will sign up for the USD 300 service, 15% for the USD 50 service, and 80% will not sign up for anything.
3. (a) What roles does exploratory data analysis play in a data mining project?
 (d) What is 'random sampling'? Explain, in brief, the difference between 'sampling with replacement' and 'sampling without replacement'.
 (b) Does the experiment below suggest any bias? Justify your answer.
Experiment: Mall shoppers are asked to fill out and return a form rating their shopping experience at each of the 26 stores to identify the most popular stores in each of four categories.
5. (a) A student counted the number of words in an essay she had written, recording the total every 10 lines (as shown in table below). How many words has she written if she writes 65 lines? Are you happy with the estimate. If not, what would you do to make yourself satisfied.
 (c) Assume that the number of words in an essay is proportional to the number of lines. If the relationship is that one variable decreases when the other increases, then will the values of Pearson and Spearman correlation coefficients be the same?
 (d) Consider the data given in table below. Show the contingency table and compute the distance between Jack and Mary using Jaccard's Coefficient.
- | No. of Lines (x) | 10 | 20 | 30 | 40 | 50 | 60 | 70 | 80 |
|------------------|-----|-----|-----|-----|-----|-----|-----|-----|
| No. of Words (y) | 150 | 186 | 210 | 234 | 258 | 282 | 306 | 330 |
-
- | Name | Fever | Cough | Medical Test 1 | Medical Test 2 | Medical Test 3 | Medical Test 4 |
|------|-------|----------|----------------|----------------|----------------|----------------|
| Jack | Yes | Negative | Positive | Negative | Negative | Negative |
| Mary | Yes | Negative | Positive | Negative | Positive | Negative |
| Jim | Yes | Positive | Negative | Negative | Negative | Negative |

$(5 + 1 + 2) + (2 + 2) = 12$

- (b) Which ~~intuitive~~ ~~of the~~ ~~descriptive~~ ~~statistical~~ clustering algorithms of the single linkage method is more appropriate for hierarchical clustering algorithms is robust to outliers? Justify your answer.
- (b) Draw a schematic diagram for a typical visualization process.
- (c) The table below contains the pairwise distances for the five objects (A, B, C, D and E). Use the 'Ansambe's Quartet' [Francis Anscombe, 1973] example to illustrate the importance of looking at a set of data graphically before starting to analyse it according to a particular type of relationship, and the inadequacy of basic statistical properties for describing real-life datasets.

Object	A	B	C	D	E
A	0	1	2	2	3
B	1	0	2	4	3
C	2	2	0	5	4
D	2	4	1	0	3
E	3	3	4	3	0

9. (a) Explain expressiveness and effectiveness in the context of visual encoding.
- (b) Give two examples of each of the following:
 (i) Quantitative Mackinlay's retinal variables, (ii) Ordinal Mackinlay's retinal variables and (iii) Nominal Mackinlay's retinal variables
- (d) Calculate the lie factor for the following figure below. Is this always a good approach? Justify your answer.



- (d) By definition, Bayesian classification works with categorical variables. Explain, in brief, the two typical approaches for applying naive Bayes technique to numerical variables.

Group - E

8. (a) One formal definition of 'visualization' says "... is the *process* of extracting salient *features* from *sets of data* and *displaying* the features