

SOCIETY OF ACTUARIES

**EXAM SRM - STATISTICS FOR RISK MODELING**

**EXAM SRM SAMPLE QUESTIONS AND SOLUTIONS**

These questions and solutions are representative of the types of questions that might be asked of candidates sitting for Exam SRM. These questions are intended to represent the depth of understanding required of candidates. The distribution of questions by topic is not intended to represent the distribution of questions on future exams.

April 2019 update: Question 2, answer III was changed. While the statement was true, it was not directly supported by the required readings.

July 2019 update: Questions 29-32 were added.

September 2019 update: Questions 33-44 were added.

December 2019 update: Question 40 item I was modified.

January 2020 update: Question 41 item I was modified.

November 2020 update: Questions 6 and 14 were modified and Questions 45-48 were added.

February 2021 update: Questions 49-53 were added.

July 2021 update: Questions 54-55 were added.

October 2021 update: Questions 56-57 added.

February 2022 update: Questions 58-60 were added. Question 59 corrected in July 2022.

June 2022 update: Questions 61-63 were added.

September 2022 update: Question 64 was added.

August 2023 update: Questions 66-67 were added, Questions 17, 28, 47, and 65 were deleted (no longer on the syllabus)

Copyright 2023 by the Society of Actuaries

## QUESTIONS

1. You are given the following four pairs of observations:

$$x_1 = (-1,0), \quad x_2 = (1,1), \quad x_3 = (2,-1), \quad \text{and} \quad x_4 = (5,10).$$

A hierarchical clustering algorithm is used with complete linkage and Euclidean distance.

Calculate the intercluster dissimilarity between  $\{x_1, x_2\}$  and  $\{x_4\}$ .

- (A) 2.2
- (B) 3.2
- (C) 9.9
- (D) 10.8
- (E) 11.7

2. Determine which of the following statements is/are true.

- I. The number of clusters must be pre-specified for both  $K$ -means and hierarchical clustering.
- II. The  $K$ -means clustering algorithm is less sensitive to the presence of outliers than the hierarchical clustering algorithm.
- III. The  $K$ -means clustering algorithm requires random assignments while the hierarchical clustering algorithm does not.

- (A) I only
- (B) II only
- (C) III only
- (D) I, II and III
- (E) The correct answer is not given by (A), (B), (C), or (D)

3. You are given:

i) The random walk model

$$y_t = y_0 + c_1 + c_2 + \dots + c_t$$

where  $c_t, t = 0, 1, 2, \dots, T$  denote observations from a white noise process.

ii) The following nine observed values of  $c_t$ :

|       |    |    |    |    |    |    |    |    |    |
|-------|----|----|----|----|----|----|----|----|----|
| $t$   | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
| $c_t$ | 2  | 3  | 5  | 3  | 4  | 2  | 4  | 1  | 2  |

iii) The average value of  $c_1, c_2, \dots, c_{10}$  is 2.

iv) The 9 step ahead forecast of  $y_{19}$ ,  $\hat{y}_{19}$ , is estimated based on the observed value of  $y_{10}$ .

Calculate the forecast error,  $y_{19} - \hat{y}_{19}$ .

- (A) 1
- (B) 2
- (C) 3
- (D) 8
- (E) 18

4. You are given:

i) The random walk model

$$y_t = y_0 + c_1 + c_2 + \dots + c_t$$

where  $c_t, t = 0, 1, 2, \dots, T$  denote observations from a white noise process.

ii) The following ten observed values of  $y_t$ .

|       |   |   |    |    |    |    |    |    |    |    |
|-------|---|---|----|----|----|----|----|----|----|----|
| $t$   | 1 | 2 | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 |
| $y_t$ | 2 | 5 | 10 | 13 | 18 | 20 | 24 | 25 | 27 | 30 |

iii)  $y_0 = 0$

Calculate the standard error of the 9 step-ahead forecast,  $\hat{y}_{19}$ .

(A) 4/3

(B) 4

(C) 9

(D) 12

(E) 16

5. Consider the following statements:
- I. Principal Component Analysis (PCA) provide low-dimensional linear surfaces that are closest to the observations.
  - II. The first principal component is the line in  $p$ -dimensional space that is closest to the observations.
  - III. PCA finds a low dimension representation of a dataset that contains as much variation as possible.
  - IV. PCA serves as a tool for data visualization.

Determine which of the statements are correct.

- (A) Statements I, II, and III only
- (B) Statements I, II, and IV only
- (C) Statements I, III, and IV only
- (D) Statements II, III, and IV only
- (E) Statements I, II, III, and IV are all correct

6. Consider the following statements:
- I. The proportion of variance explained by an additional principal component never decreases as more principal components are added.
  - II. The cumulative proportion of variance explained never decreases as more principal components are added.
  - III. Using all possible principal components provides the best understanding of the data.
  - IV. A scree plot provides a method for determining the number of principal components to use.

Determine which of the statements are correct.

- (A) Statements I and II only
- (B) Statements I and III only
- (C) Statements I and IV only
- (D) Statements II and III only
- (E) Statements II and IV only

7. Determine which of the following pairs of distribution and link function is the most appropriate to model if a person is hospitalized or not.
- (A) Normal distribution, identity link function
  - (B) Normal distribution, logit link function
  - (C) Binomial distribution, linear link function
  - (D) Binomial distribution, logit link function
  - (E) It cannot be determined from the information given.
8. Determine which of the following statements describe the advantages of using an alternative fitting procedure, such as subset selection and shrinkage, instead of least squares.
- I. Doing so will likely result in a simpler model
  - II. Doing so will likely improve prediction accuracy
  - III. The results are likely to be easier to interpret
- (A) I only
  - (B) II only
  - (C) III only
  - (D) I, II, and III
  - (E) The correct answer is not given by (A), (B), (C), or (D)

9. A classification tree is being constructed to predict if an insurance policy will lapse. A random sample of 100 policies contains 30 that lapsed. You are considering two splits:

Split 1: One node has 20 observations with 12 lapses and one node has 80 observations with 18 lapses.

Split 2: One node has 10 observations with 8 lapses and one node has 90 observations with 22 lapses.

The total Gini index after a split is the weighted average of the Gini index at each node, with the weights proportional to the number of observations in each node.

The total entropy after a split is the weighted average of the entropy at each node, with the weights proportional to the number of observations in each node.

Determine which of the following statements is/are true?

- I. Split 1 is preferred based on the total Gini index.
  - II. Split 1 is preferred based on the total entropy.
  - III. Split 1 is preferred based on having fewer classification errors.
- (A) I only
  - (B) II only
  - (C) III only
  - (D) I, II, and III
  - (E) The correct answer is not given by (A), (B), (C), or (D).



10. Determine which of the following statements about random forests is/are true?
- I. If the number of predictors used at each split is equal to the total number of available predictors, the result is the same as using bagging.
  - II. When building a specific tree, the same subset of predictor variables is used at each split.
  - III. Random forests are an improvement over bagging because the trees are decorrelated.
- (A) None
- (B) I and II only
- (C) I and III only
- (D) II and III only
- (E) The correct answer is not given by (A), (B), (C), or (D).

11. You are given the following results from a regression model.

| Observation number ( $i$ ) | $y_i$ | $\hat{f}(x_i)$ |
|----------------------------|-------|----------------|
| 1                          | 2     | 4              |
| 2                          | 5     | 3              |
| 3                          | 6     | 9              |
| 4                          | 8     | 3              |
| 5                          | 4     | 6              |

Calculate the sum of squared errors (SSE).

- (A) -35
- (B) -5
- (C) 5
- (D) 35
- (E) 46

12. Determine which of the following statements is true

- (A) Linear regression is a flexible approach
- (B) Lasso is more flexible than a linear regression approach
- (C) Bagging is a low flexibility approach
- (D) There are methods that have high flexibility and are also easy to interpret
- (E) None of (A), (B), (C), or (D) are true

13. Determine which of the following statements is/are true for a simple linear relationship,  $y = \beta_0 + \beta_1 x + \varepsilon$ .
- I. If  $\varepsilon = 0$ , the 95% confidence interval is equal to the 95% prediction interval.
  - II. The prediction interval is always at least as wide as the confidence interval.
  - III. The prediction interval quantifies the possible range for  $E(y | x)$ .
- (A) I only
- (B) II only
- (C) III only
- (D) I, II, and III
- (E) The correct answer is not given by (A), (B), (C), or (D).

14. From an investigation of the residuals of fitting a linear regression by ordinary least squares it is clear that the spread of the residuals increases as the predicted values increase. Observed values of the dependent variable range from 0 to 100.

Determine which of the following statements is/are true with regard to transforming the dependent variable to make the variance of the residuals more constant.

- I. Taking the logarithm of one plus the value of the dependent variable may make the variance of the residuals more constant.
  - II. A square root transformation may make the variance of the residuals more constant.
  - III. A logit transformation may make the variance of the residuals more constant.
- 
- (A) None
  - (B) I and II only
  - (C) I and III only
  - (D) II and III only
  - (E) The correct answer is not given by (A), (B), (C), or (D).

15. You are performing a  $K$ -means clustering algorithm on a set of data. The data has been initialized randomly with 3 clusters as follows:

| Cluster | Data Point |
|---------|------------|
| A       | (2, -1)    |
| A       | (-1, 2)    |
| A       | (-2, 1)    |
| A       | (1, 2)     |
| B       | (4, 0)     |
| B       | (4, -1)    |
| B       | (0, -2)    |
| B       | (0, -5)    |
| C       | (-1, 0)    |
| C       | (3, 8)     |
| C       | (-2, 0)    |
| C       | (0, 0)     |

A single iteration of the algorithm is performed using the Euclidian distance between points and the cluster containing the fewest number of data points is identified.

Calculate the number of data points in this cluster.

- (A) 0
- (B) 1
- (C) 2
- (D) 3
- (E) 4

16. Determine which of the following statements is applicable to  $K$ -means clustering and is not applicable to hierarchical clustering.
- (A) If two different people are given the same data and perform one iteration of the algorithm, their results at that point will be the same.
  - (B) At each iteration of the algorithm, the number of clusters will be greater than the number of clusters in the previous iteration of the algorithm.
  - (C) The algorithm needs to be run only once, regardless of how many clusters are ultimately decided to use.
  - (D) The algorithm must be initialized with an assignment of the data points to a cluster.
  - (E) None of (A), (B), (C), or (D) meet the stated criterion.

17. DELETED

18. For a simple linear regression model the sum of squares of the residuals is

$$\sum_{i=1}^{25} e_i^2 = 230 \text{ and the } R^2 \text{ statistic is } 0.64.$$

Calculate the total sum of squares (TSS) for this model.

- (A) 605.94
- (B) 638.89
- (C) 690.77
- (D) 701.59
- (E) 750.87

19. The regression model  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \varepsilon$  is being investigated.

The following maximized log-likelihoods are obtained:

- Using only the intercept term:  $-1126.91$
- Using only the intercept term,  $X_1$ , and  $X_2$ :  $-1122.41$
- Using all four terms:  $-1121.91$

The null hypothesis  $H_0 : \beta_1 = \beta_2 = \beta_3 = 0$  is being tested at the 5% significance level using the likelihood ratio test.

Determine which of the following is true.

- (A) The test statistic is equal to 1 and the hypothesis cannot be rejected.
- (B) The test statistic is equal to 9 and the hypothesis cannot be rejected
- (C) The test statistic is equal to 10 and the hypothesis cannot be rejected.
- (D) The test statistic is equal to 9 and the hypothesis should be rejected.
- (E) The test statistic is equal to 10 and the hypothesis should be rejected.

20. An analyst is modeling the probability of a certain phenomenon occurring. The analyst has observed that the simple linear model currently in use results in predicted values less than zero and greater than one.

Determine which of the following is the most appropriate way to address this issue.

- (A) Limit the data to observations that are expected to result in predicted values between 0 and 1.
- (B) Consider predicted values below 0 as 0 and values above 1 as 1.
- (C) Use a logit function to transform the linear model into only predicting values between 0 and 1.
- (D) Use the canonical link function for the Poisson distribution to transform the linear model into only predicting values between 0 and 1.
- (E) None of the above.

21. A random walk is expressed as

$$y_t = y_{t-1} + c_t \text{ for } t = 1, 2, \dots$$

where

$$E(c_t) = \mu_c \text{ and } \text{Var}(c_t) = \sigma_c^2, \quad t = 1, 2, \dots$$

Determine which statements is/are true with respect to a random walk model.

- I. If  $\mu_c \neq 0$ , then the random walk is nonstationary in the mean.
- II. If  $\sigma_c^2 = 0$ , then the random walk is nonstationary in the variance.
- III. If  $\sigma_c^2 > 0$ , then the random walk is nonstationary in the variance.

- (A) None
- (B) I and II only
- (C) I and III only
- (D) II and III only
- (E) The correct answer is not given by (A), (B), (C), or (D).

22. A stationary autoregressive model of order one can be written as

$$y_t = \beta_0 + \beta_1 y_{t-1} + \varepsilon_t, \quad t = 1, 2, \dots$$

Determine which of the following statements about this model is false

- (A) The parameter  $\beta_0$  must not equal 1.
- (B) The absolute value of the parameter  $\beta_1$  must be less than 1.
- (C) If the parameter  $\beta_1 = 0$ , then the model reduces to a white noise process.
- (D) If the parameter  $\beta_1 = 1$ , then the model is a random walk.
- (E) Only the immediate past value,  $y_{t-1}$ , is used as a predictor for  $y_t$ .



23. Toby observes the following coffee prices in his company cafeteria:

- 12 ounces for 1.00
- 16 ounces for 1.20
- 20 ounces for 1.40

The cafeteria announces that they will begin to sell any amount of coffee for a price that is the value predicted by a simple linear regression using least squares of the current prices on size.

Toby and his co-worker Karen want to determine how much they would save each day, using the new pricing, if, instead of each buying a 24-ounce coffee, they bought a 48-ounce coffee and shared it.

Calculate the amount they would save.

- (A) It would cost them 0.40 more.
- (B) It would cost the same.
- (C) They would save 0.40.
- (D) They would save 0.80.
- (E) They would save 1.20.

24. Sarah performs a regression of the return on a mutual fund ( $y$ ) on four predictors plus an intercept. She uses monthly returns over 105 months. Her software calculates the  $F$  statistic for the regression as  $F = 20.0$ , but then it quits working before it calculates the value of  $R^2$ . While she waits on hold with the help desk, she tries to calculate  $R^2$  from the  $F$ -statistic.

Determine which of the following statements about the attempted calculation is true.

- (A) There is insufficient information, but it could be calculated if she had the value of the residual sum of squares (RSS).
- (B) There is insufficient information, but it could be calculated if she had the value of the total sum of squares (TSS) and RSS.
- (C)  $R^2 = 0.44$
- (D)  $R^2 = 0.56$
- (E)  $R^2 = 0.80$

25. Determine which of the following statements concerning decision tree pruning is/are true.

- I. The recursive binary splitting method can lead to overfitting the data.
  - II. A tree with more splits tends to have lower variance.
  - III. When using the cost complexity pruning method,  $\alpha = 0$  results in a very large tree.
- (A) None
  - (B) I and II only
  - (C) I and III only
  - (D) II and III only
  - (E) The correct answer is not given by (A), (B), (C), or (D).

26. Each picture below represents a two-dimensional space where observations are classified into two categories. The categories are representing by light and dark shading. A classification tree is to be constructed for each space.

Determine which space can be modeled with no error by a classification tree.

I.



II.



III.



- (A) I only
- (B) II only
- (C) III only
- (D) I, II, and III
- (E) The correct answer is not given by (A), (B), (C), or (D).

27. Trevor is modeling monthly incurred dental claims. Trevor has 48 monthly claims observations and three potential predictors:

- Number of weekdays in the month
- Number of weekend days in the month
- Average number of insured members during the month

Trevor obtained the following results from a linear regression:

|                           | <b>Coefficient</b> | <b>Standard Error</b> | <b><i>t</i> Stat</b> | <b><i>p</i>-value</b> |
|---------------------------|--------------------|-----------------------|----------------------|-----------------------|
| Intercept                 | -45,765,767.76     | 20,441,816.55         | -2.24                | 0.0303                |
| Number of weekdays        | 513,280.76         | 233,143.23            | 2.20                 | 0.0330                |
| Number of weekend days    | 280,148.46         | 483,001.55            | 0.58                 | 0.5649                |
| Average number of members | 38.64              | 6.42                  | 6.01                 | 0.0000                |

Determine which of the following variables should be dropped, using a 5% significance level.

- I. Intercept
  - II. Number of weekdays
  - III. Number of weekend days
  - IV. Number of members.
- (A) I only
  - (B) II only
  - (C) III only
  - (D) IV only
  - (E) None should be dropped from the model

28. DELETED
29. Determine which of the following considerations may make decision trees preferable to other statistical learning methods.
- I. Decision trees are easily interpretable.
  - II. Decision trees can be displayed graphically.
  - III. Decision trees are easier to explain than linear regression methods.
- (A) None
- (B) I and II only
- (C) I and III only
- (D) II and III only
- (E) The correct answer is not given by (A), (B), (C), or (D).
30. Principal component analysis is applied to a large data set with four variables. Loadings for the first four principal components are estimated. Determine which of the following statements is/are true with respect the loadings.
- I. The loadings are unique.
  - II. For a given principal component, the sum of the squares of the loadings across the four variables is one.
  - III. Together, the four principal components explain 100% of the variance.
- (A) None
- (B) I and II only
- (C) I and III only
- (D) II and III only
- (E) The correct answer is not given by (A), (B), (C), or (D).

31. Determine which of the following indicates that a nonstationary time series can be represented as a random walk

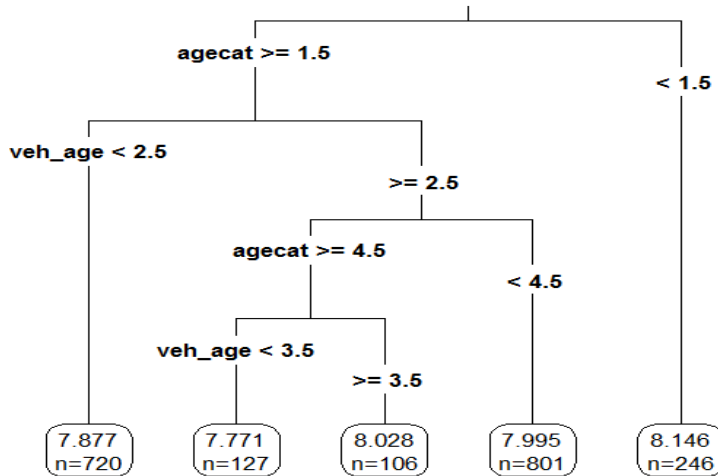
- I. A control chart of the series detects a linear trend in time and increasing variability.
  - II. The differenced series follows a white noise model.
  - III. The standard deviation of the original series is greater than the standard deviation of the differenced series.
- (A) I only  
(B) II only  
(C) III only  
(D) I, II and III  
(E) The correct answer is not given by (A), (B), (C), or (D).

32. You are given a set of  $n$  observations, each with  $p$  features.

Determine which of the following statements is/are true with respect to clustering methods.

- I. The  $n$  observations can be clustered on the basis of the  $p$  features to identify subgroups among the observations.
  - II. The  $p$  features can be clustered on the basis of the  $n$  observations to identify subgroups among the features.
  - III. Clustering is an unsupervised learning method and is often performed as part of an exploratory data analysis.
- (A) None  
(B) I and II only  
(C) I and III only  
(D) II and III only  
(E) The correct answer is not given by (A), (B), (C), or (D).

33. The regression tree shown below was produced from a dataset of auto claim payments. Age Category (agecat: 1, 2, 3, 4, 5, 6) and Vehicle Age (veh\_age: 1, 2, 3, 4, ...) are both predictor variables, and log of claim amount (LCA) is the dependent variable.



Consider three autos I, II, III:

- I: An Auto in Age Category 1 and Vehicle Age 4
- II: An Auto in Age Category 5 and Vehicle Age 5
- III: An Auto in Age Category 5 and Vehicle Age 3

Rank the estimated LCA of Autos I, II, and III.

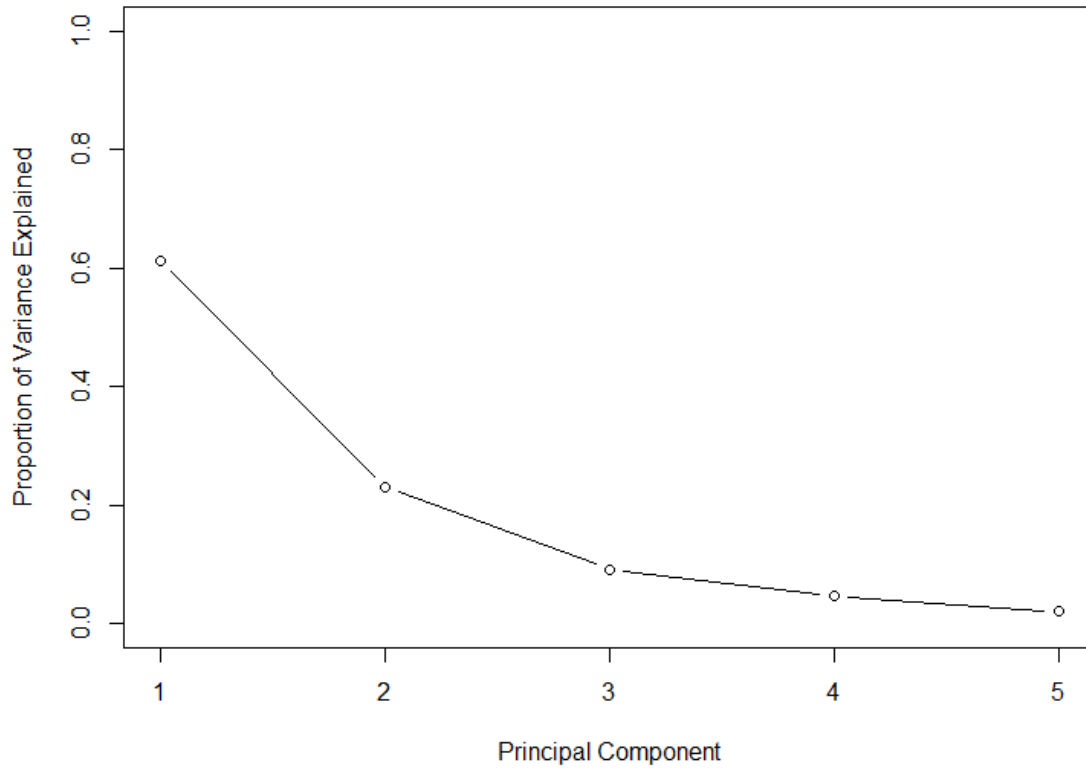
- (A)  $LCA(I) < LCA(II) < LCA(III)$
- (B)  $LCA(I) < LCA(III) < LCA(II)$
- (C)  $LCA(II) < LCA(I) < LCA(III)$
- (D)  $LCA(II) < LCA(III) < LCA(I)$
- (E)  $LCA(III) < LCA(II) < LCA(I)$

34. Determine which of the following statements is/are true about clustering methods:

- I. If  $K$  is held constant,  $K$ -means clustering will always produce the same cluster assignments.
  - II. Given a linkage and a dissimilarity measure, hierarchical clustering will always produce the same cluster assignments for a specific number of clusters.
  - III. Given identical data sets, cutting a dendrogram to obtain five clusters produces the same cluster assignments as  $K$ -means clustering with  $K = 5$ .
- (A) I only
  - (B) II only
  - (C) III only
  - (D) I, II and III
  - (E) The correct answer is not given by (A), (B), (C), or (D).



35. Using the following scree plot, determine the minimum number of principal components that are needed to explain at least 80% of the variance of the original dataset.



- (A) One
- (B) Two
- (C) Three
- (D) Four
- (E) It cannot be determined from the information given.

36. Determine which of the following statements about hierarchical clustering is/are true.

- I. The method may not assign extreme outliers to any cluster.
  - II. The resulting dendrogram can be used to obtain different numbers of clusters.
  - III. The method is not robust to small changes in the data.
- 
- (A) None
  - (B) I and II only
  - (C) I and III only
  - (D) II and III only
  - (E) The correct answer is not given by (A), (B), (C), or (D).

37. Analysts W, X, Y, and Z are each performing Principal Components Analysis on the same data set with three variables. They use different programs with their default settings and discover that they have different factor loadings for the first principal component. Their loadings are:

|   | Variable 1 | Variable 2 | Variable 3 |
|---|------------|------------|------------|
| W | -0.549     | -0.594     | 0.587      |
| X | -0.549     | 0.594      | 0.587      |
| Y | 0.549      | -0.594     | -0.587     |
| Z | 0.140      | -0.570     | -0.809     |

Determine which of the following is/are plausible explanations for the different loadings.

- I. Loadings are unique up to a sign flip and hence X's and Y's programs could make different arbitrary sign choices.
  - II. Z's program defaults to not scaling the variables while Y's program defaults to scaling them.
  - III. Loadings are unique up to a sign flip and hence W's and X's programs could make different arbitrary sign choices.
- (A) None
  - (B) I and II only
  - (C) I and III only
  - (D) II and III only
  - (E) The correct answer is not given by (A), (B), (C), or (D).

38. You are given two models:

Model L:  $y_t = \beta_0 + \beta_1 t + \varepsilon_t$

where  $\{\varepsilon_t\}$  is a white noise process, for  $t = 0, 1, 2, \dots$

Model M:  $y_t = y_0 + \mu_c t + u_t$

$$c_t = y_t - y_{t-1}$$

$$u_t = \sum_{j=1}^t \varepsilon_j$$

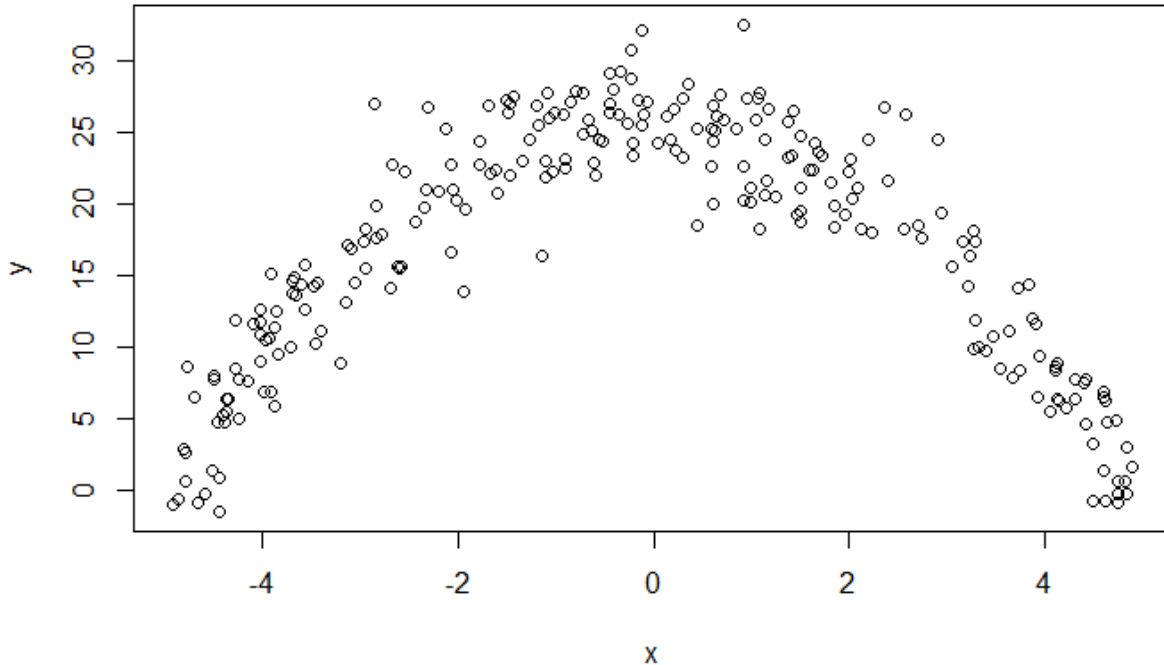
where  $\{\varepsilon_t\}$  is a white noise process, for  $t = 0, 1, 2, \dots$

Determine which of the following statements is/are true.

- I. Model L is a linear trend in time model where the error component is not a random walk.
  - II. Model M is a random walk model where the error component of the model is also a random walk.
  - III. The comparison between Model L and Model M is not clear when the parameter  $\mu_c = 0$ .
- 
- (A) I only
  - (B) II only
  - (C) III only
  - (D) I, II and III
  - (E) The correct answer is not given by (A), (B), (C), or (D).

39. You are given a dataset with two variables, which is graphed below. You want to predict  $y$  using  $x$ .

Determine which statement regarding using a generalized linear model (GLM) or a random forest is true.



- (A) A random forest is appropriate because the dataset contains only quantitative variables.
- (B) A random forest is appropriate because the data does not follow a straight line.
- (C) A GLM is not appropriate because the variance of  $y$  given  $x$  is not constant.
- (D) A random forest is appropriate because there is a clear relationship between  $y$  and  $x$ .
- (E) A GLM is appropriate because it can accommodate polynomial relationships.

40. Determine which of the following statements about clustering is/are true.
- I. Cutting a dendrogram at a lower height will not decrease the number of clusters.
  - II.  $K$ -means clustering requires plotting the data before determining the number of clusters.
  - III. For a given number of clusters, hierarchical clustering can sometimes yield less accurate results than  $K$ -means clustering.
- (A) None
  - (B) I and II only
  - (C) I and III only
  - (D) II and III only
  - (E) The correct answer is not given by (A), (B), (C), or (D).
41. For a random forest, let  $p$  be the total number of features and  $m$  be the number of features selected at each split.
- Determine which of the following statements is/are true.
- I. When  $m = p$ , random forest and bagging are the same procedure.
  - II.  $\frac{p-m}{p}$  is the probability a split will not consider the strongest predictor.
  - III. The typical choice of  $m$  is  $\frac{p}{2}$ .
- (A) None
  - (B) I and II only
  - (C) I and III only
  - (D) II and III only
  - (E) The correct answer is not given by (A), (B), (C), or (D).

42. Determine which of the following statements is NOT true about the linear probability, logistic, and probit regression models for binary dependent variables.
- (A) The three major drawbacks of the linear probability model are poor fitted values, heteroscedasticity, and meaningless residual analysis.
  - (B) The logistic and probit regression models aim to circumvent the drawbacks of linear probability models.
  - (C) The logit function is given by  $\pi(z) = e^z / (1 + e^z)$ .
  - (D) The probit function is given by  $\pi(z) = \Phi(z)$  where  $\Phi$  is the standard normal distribution function.
  - (E) The logit and probit functions are substantially different.
43. Determine which of the following statements about clustering methods is NOT true.
- (A) Clustering is used to discover structure within a data set.
  - (B) Clustering is used to find homogeneous subgroups among the observations within a data set.
  - (C) Clustering is an unsupervised learning method.
  - (D) Clustering is used to reduce the dimensionality of a dataset while retaining explanation for a good fraction of the variance.
  - (E) In  $K$ -means clustering, it is necessary to pre-specify the number of clusters.

44. Two actuaries are analyzing dental claims for a group of  $n = 100$  participants. The predictor variable is sex, with 0 and 1 as possible values.

Actuary 1 uses the following regression model:

$$Y = \beta + \varepsilon .$$

Actuary 2 uses the following regression model:

$$Y = \beta_0 + \beta_1 \times \text{Sex} + \varepsilon .$$

The residual sum of squares for the regression of Actuary 2 is 250,000 and the total sum of squares is 490,000.

Calculate the  $F$ -statistic to test whether the model of Actuary 2 is a significant improvement over the model of Actuary 1.

- (A) 92
- (B) 93
- (C) 94
- (D) 95
- (E) 96



45. The actuarial student committee of a large firm has collected data on exam scores. A generalized linear model where the target is the exam score on a 0-10 scale is constructed using a log link, resulting in the following estimated coefficients

| Predictor Variables                        | Coefficient |
|--|-------------|
| Intercept                                  | - 0.1       |
| Study Time (in units of 100 hours)         | 0.5         |
| Attempt (1 for first attempt, else 0)      | 0.5         |
| Master's degree (1 for Yes, 0 for No)      | - 0.1       |
| Interaction of Attempt and Master's degree | 0.2         |

The company is about to offer a job to an applicant who has a Master's degree and for whom the exam would be a first attempt. It would like to offer half of the study time that will result in an expected exam score of 6.0.

Calculate the amount of study time that the company should offer.

- (A) 123 hours
- (B) 126 hours
- (C) 129 hours
- (D) 132 hours
- (E) 135 hours

46. A time series was observed at times 0, 1, ..., 100. The last four observations along with estimates based on exponential and double exponential smoothing with  $w = 0.8$  are:

|                                 |      |      |      |       |
|---------------------------------|------|------|------|-------|
| Time ( $t$ )                    | 97   | 98   | 99   | 100   |
| Observation ( $y_t$ )           | 96.9 | 98.1 | 99.0 | 100.2 |
| Estimates ( $\hat{s}_t^{(1)}$ ) | 93.1 | 94.1 | 95.1 |       |
| Estimates ( $\hat{s}_t^{(2)}$ ) | 88.9 | 89.9 |      |       |

All forecasts should be rounded to one decimal place and the trend should be rounded to three decimal places.

Let  $F$  be the predicted value of  $y_{102}$  using exponential smoothing with  $w = 0.8$ .

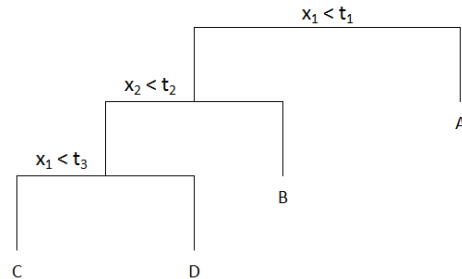
Let  $G$  be the predicted value of  $y_{102}$  using double exponential smoothing with  $w = 0.8$ .

Calculate the absolute difference between  $F$  and  $G$ ,  $|F - G|$ .

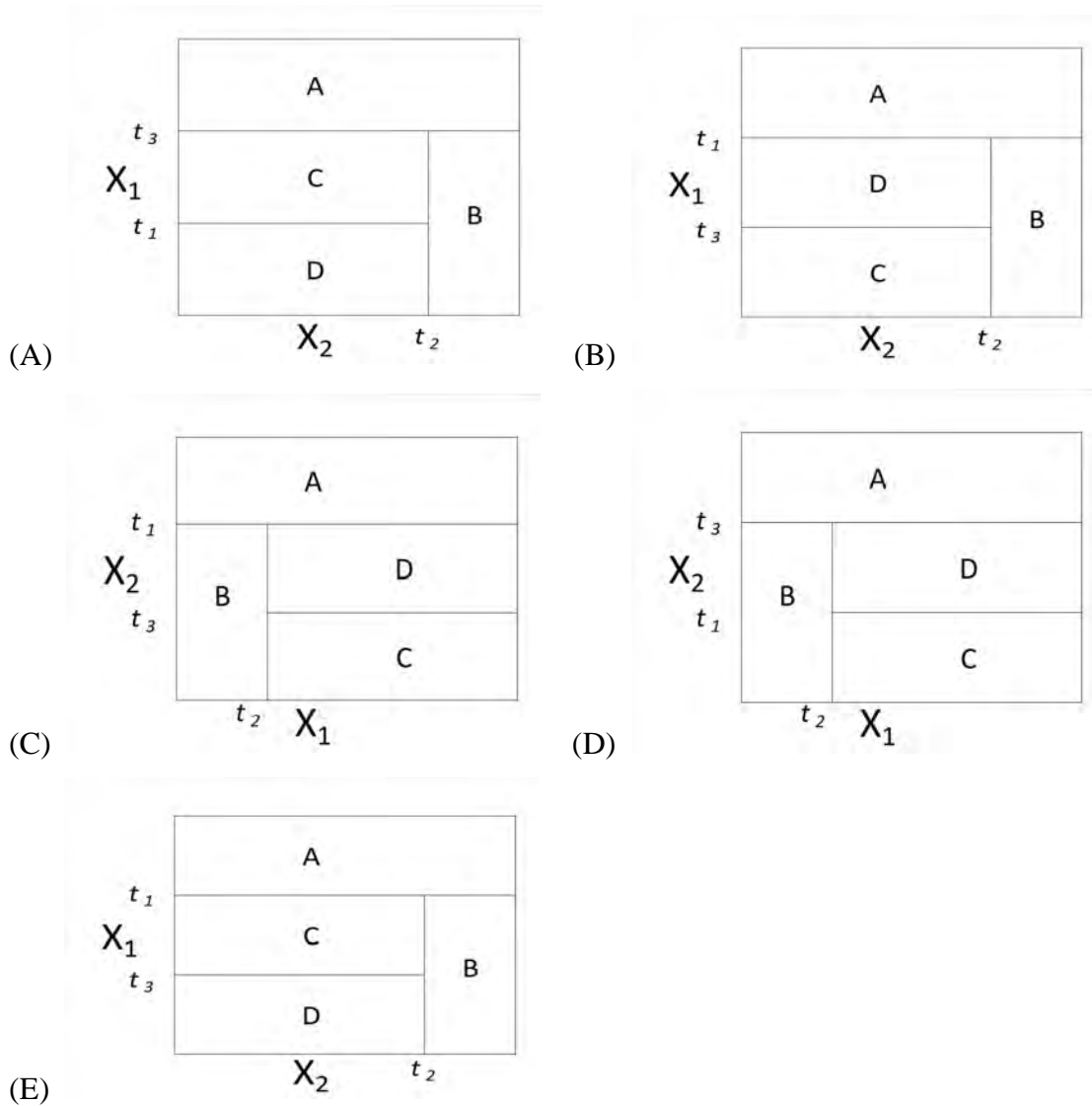
- (A) 0.0
- (B) 2.1
- (C) 4.2
- (D) 6.3
- (E) 8.4

47. DELETED

48. The following tree was constructed using recursive binary splitting with the left branch indicating that the inequality is true.



Determine which of the following plots represents this tree.



49. Trish runs a regression on a data set of  $n$  observations. She then calculates a 95% confidence interval  $(t, u)$  on  $y$  for a given set of predictors. She also calculates a 95% prediction interval  $(v, w)$  on  $y$  for the same set of predictors.

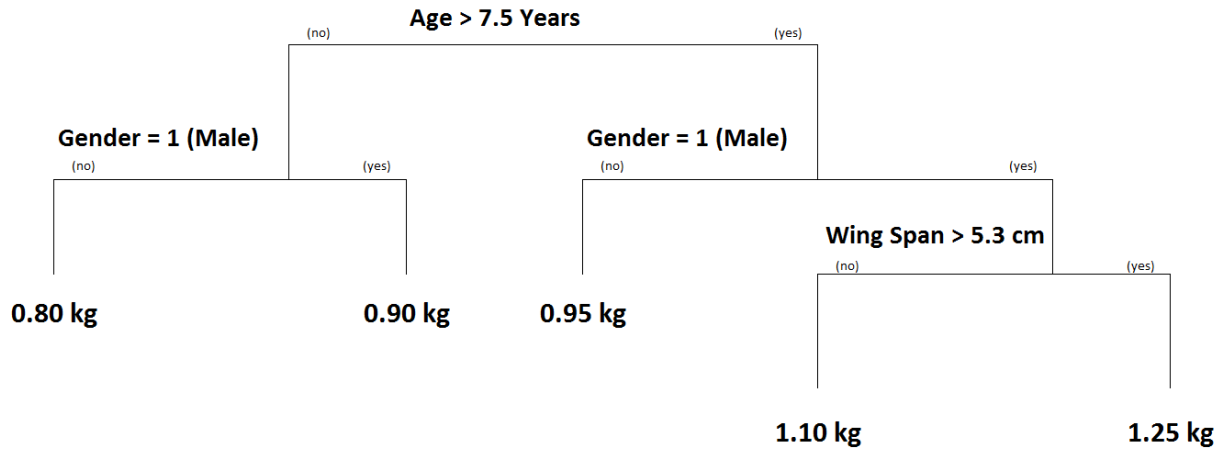
Determine which of the following must be true.

- I.  $\lim_{n \rightarrow \infty} (u - t) = 0$
  - II.  $\lim_{n \rightarrow \infty} (w - v) = 0$
  - III.  $w - v > u - t$
- (A) None
  - (B) I and II only
  - (C) I and III only
  - (D) II and III only
  - (E) The correct answer is not given by (A), (B), (C), or (D).

50. Determine which of the following statements regarding statistical learning methods is/are true.

- I. Methods that are highly interpretable are more likely to be highly flexible.
  - II. When inference is the goal, there are clear advantages to using a lasso method versus a bagging method.
  - III. Using a more flexible method will produce a more accurate prediction against unseen data.
- (A) I only
  - (B) II only
  - (C) III only
  - (D) I, II and III
  - (E) The correct answer is not given by (A), (B), (C), or (D).

51. You are given the following regression tree predicting the weight of ducks in kilograms (kg):



You predict the weight of the following three ducks:

X: Wing Span = 5.5 cm, Male, Age = 7 years

Y: Wing Span = 5.8 cm, Female, Age = 5 years

Z: Wing Span = 5.7 cm, Male, Age = 8 years

Determine the order of the predicted weights of the three ducks.

- (A)  $X < Y < Z$
- (B)  $X < Z < Y$
- (C)  $Y < X < Z$
- (D)  $Y < Z < X$
- (E)  $Z < X < Y$

52. Determine which of the following statements is/are true about Pearson residuals.
- I. They can be used to calculate a goodness-of-fit statistic.
  - II. They can be used to detect if additional variables of interest can be used to improve the model specification.
  - III. They can be used to identify unusual observations.
- (A) I only  
(B) II only  
(C) III only  
(D) I, II, and III  
(E) The correct answer is not given by (A), (B), (C), or (D).
53. Determine which of the following statements is NOT true about the equation

$$Y = \beta_0 + \beta_1 X + \varepsilon.$$

- (A)  $\beta_0$  is the expected value of  $Y$ .  
(B)  $\beta_1$  is the average increase in  $Y$  associated with a one-unit increase in  $X$ .  
(C) The error term,  $\varepsilon$ , is typically assumed to be independent of  $X$ .  
(D) The equation defines the population regression line.  
(E) The method of least squares is commonly used to estimate the coefficients  $\beta_0$  and  $\beta_1$ .

54. For a regression model of executive compensation, you are given:
- i) The following statistics:

| Executive Compensation |           |            |                     |                 |
|------------------------|-----------|------------|---------------------|-----------------|
| Coefficients:          | Estimate  | Std. Error | <i>t</i> -statistic | <i>p</i> -value |
| (INTERCEPT)            | -28,595.5 | 220.5      | -129.7              | <0.001          |
| AGEMINUS35             | 7,366.3   | 12.5       | 588.1               | <0.001          |
| TOPSCHOOL              | 50.0      | 119.7      | 0.4                 | 0.676           |
| LARGECITY              | 147.9     | 119.7      | 1.2                 | 0.217           |
| MBA                    | 2,490.9   | 119.7      | 20.8                | <0.001          |
| YEARSEXP               | 15,286.6  | 7.2        | 2132.8              | <0.001          |

- ii) The acceptable significance level is  $\alpha = 0.10$ .

Determine which variable or variables should be removed first prior to rerunning the model.

- (A) (INTERCEPT)
- (B) AGEMINUS35, MBA, and YEARSEXP
- (C) TOPSCHOOL
- (D) TOPSCHOOL and LARGECITY
- (E) YEARSEXP

55. You are given the following eight observations from a time series that follows a random walk model:

|                       |   |   |   |   |    |    |    |    |
|-----------------------|---|---|---|---|----|----|----|----|
| Time ( $t$ )          | 0 | 1 | 2 | 3 | 4  | 5  | 6  | 7  |
| Observation ( $y_t$ ) | 3 | 5 | 7 | 8 | 12 | 15 | 21 | 22 |

You plan to fit this model to the first five observations and then evaluate it against the last three observations using one-step forecast residuals. The estimated mean of the white noise process is 2.25.

Let  $F$  be the mean error (ME) of the three predicted observations.

Let  $G$  be the mean square error (MSE) of the three predicted observations.

Calculate the absolute difference between  $F$  and  $G$ ,  $|F - G|$ .

- (A) 3.48
- (B) 4.31
- (C) 5.54
- (D) 6.47
- (E) 7.63



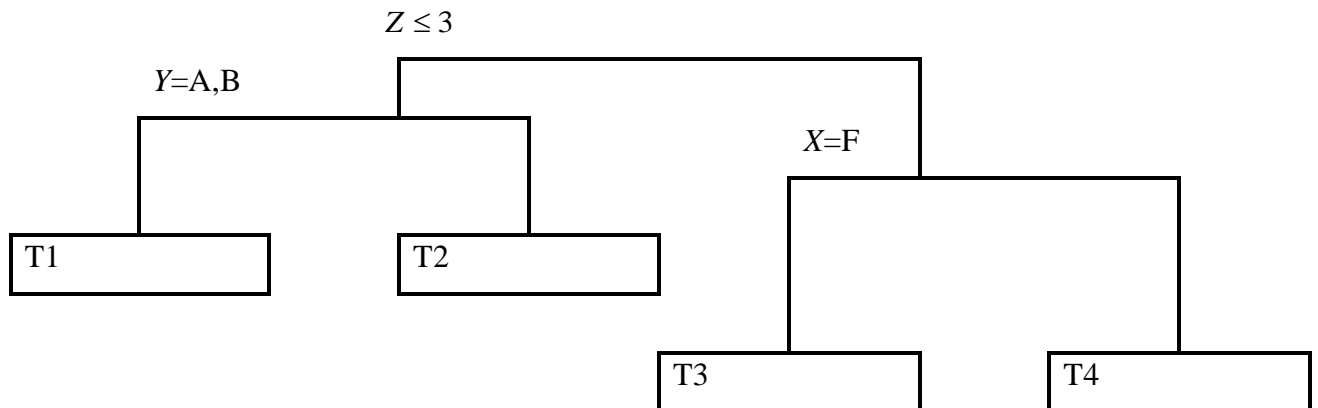
56. Determine which of the following statements about prediction is true.
- (A) Each of several candidate regression models must produce the same prediction.
  - (B) When making predictions, it is assumed that the new observation follows the same model as the one used in the sample.
  - (C) A point prediction is more reliable than an interval prediction.
  - (D) A wider prediction interval is more informative than a narrower prediction interval.
  - (E) A prediction interval should not contain the single point prediction.

57. You are given:

- i) The following observed values of the response variable,  $R$ , and predictor variables  $X$ ,  $Y$ ,  $Z$ :

|     |      |      |      |      |      |      |      |      |      |
|-----|------|------|------|------|------|------|------|------|------|
| $R$ | 4.75 | 4.67 | 4.67 | 4.56 | 4.53 | 3.91 | 3.90 | 3.90 | 3.89 |
| $X$ | M    | F    | M    | F    | M    | F    | F    | M    | M    |
| $Y$ | A    | A    | D    | D    | B    | C    | B    | D    | B    |
| $Z$ | 2    | 4    | 1    | 3    | 2    | 2    | 5    | 5    | 1    |

- ii) The following plot of the corresponding regression tree:



Calculate the Mean Response (MR) for each of the end nodes.

- (A)  $MR(T1) = 4.39$ ,  $MR(T2) = 4.38$ ,  $MR(T3) = 4.29$ ,  $MR(T4) = 3.90$   
 (B)  $MR(T1) = 4.26$ ,  $MR(T2) = 4.38$ ,  $MR(T3) = 4.62$ ,  $MR(T4) = 3.90$   
 (C)  $MR(T1) = 4.26$ ,  $MR(T2) = 4.39$ ,  $MR(T3) = 3.90$ ,  $MR(T4) = 4.29$   
 (D)  $MR(T1) = 4.64$ ,  $MR(T2) = 4.29$ ,  $MR(T3) = 4.38$ ,  $MR(T4) = 3.90$   
 (E)  $MR(T1) = 4.64$ ,  $MR(T2) = 4.38$ ,  $MR(T3) = 4.39$ ,  $MR(T4) = 3.90$

58. You are given the following six observed values of the autoregressive model of order one time series

$$y_t = \beta_0 + \beta_1 y_{t-1} + \varepsilon_t \text{ with } \text{Var}(\varepsilon_t) = \sigma^2.$$

|       |    |    |    |    |    |    |
|-------|----|----|----|----|----|----|
| $t$   | 1  | 2  | 3  | 4  | 5  | 6  |
| $y_t$ | 31 | 35 | 37 | 41 | 45 | 51 |

The approximation to the conditional least squares method is used to estimate  $\beta_0$  and  $\beta_1$ .

Calculate the mean squared error  $s^2$  that estimates  $\sigma^2$ .

- (A) 13
- (B) 21
- (C) 22
- (D) 26
- (E) 35

59. You apply 2-means clustering to a set of five observations with two features. You are given the following initial cluster assignments:

| Observation | X1 | X2 | Initial cluster |
|-------------|----|----|-----------------|
| 1           | 1  | 3  | 1               |
| 2           | 0  | 4  | 1               |
| 3           | 6  | 2  | 1               |
| 4           | 5  | 2  | 2               |
| 5           | 1  | 6  | 2               |

Calculate the total within-cluster variation of the initial cluster assignments, based on Euclidean distance measure.

- (A) 32.0
- (B) 70.3
- (C) 77.3
- (D) 118.3
- (E) 141.0

60. Determine which of the following statements about selecting the optimal number of clusters in  $K$ -means clustering is/are true.

- I.  $K$  should be set equal to  $n$ , the number of observations.
- II. Choose  $K$  such that the total within-cluster variation is minimized.
- III. The determination of  $K$  is subjective and there does not exist one method to determine the optimal number of clusters.

- (A) I only
- (B) II only
- (C) III only
- (D) I, II and III
- (E) The correct answer is not given by (A), (B), (C), or (D).

61. A linear model has been fit to a dataset containing six predictor variables, F, G, H, I, J, and K. Determine which of the following statements regarding using Akaike information criterion (AIC) or Bayesian information criterion (BIC) to select an optimal set of predictor variables for this linear model is/are true.

- I. AIC and BIC provide a direct estimate of the test error.
- II. When choosing between the subsets {F, G, H} and {I, J, K}, AIC and BIC will always select the same subset.
- III. For large sample sizes ( $n > 7$ ), the number of variables selected by BIC will be less than or equal to the number selected by AIC.

- (A) None
- (B) I and II only
- (C) I and III only
- (D) II and III only
- (E) The correct answer is not given by (A), (B), (C), or (D).

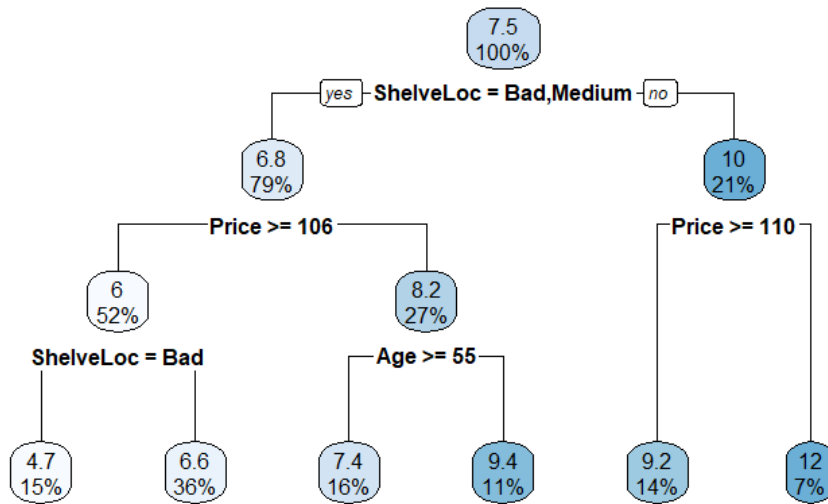
62. In a simple linear regression model based on over 100 observations, you are given the following estimates.

- i) The estimated slope is  $-1.03$ .
- ii) The standard error of the estimated slope is  $0.06$ .

Calculate the 95% confidence interval for the slope.

- (A)  $(-1.15, -0.91)$
- (B)  $(-1.13, -0.93)$
- (C)  $(-1.11, -0.95)$
- (D)  $(-1.09, -0.97)$
- (E)  $(-1.07, -0.99)$

63. You have constructed the following regression tree predicting unit sales (in thousands) of car seats. The variable ShelfLoc has possible values Good, Medium, and Bad.



| Variable    | Observed Value |
|-------------|----------------|
| ShelveLoc   | Good           |
| Price       | 120            |
| Age         | 57             |
| Advertising | 12             |

Determine the predicted unit sales (in thousands) for the above observation based on the regression tree.

- (A) 4.7
- (B) 6.6
- (C) 7.4
- (D) 9.2
- (E) 9.4

64. You are given a stationary AR(1) model,  $y_t = \beta_0 + \beta_1 y_{t-1} + \varepsilon_t$ ,  $t = 2, \dots, T$ .

Determine which of the following is always true.

(A)  $\beta_0 \neq 0$

(B)  $\beta_0 = 1$

(C)  $\beta_1 = 0$

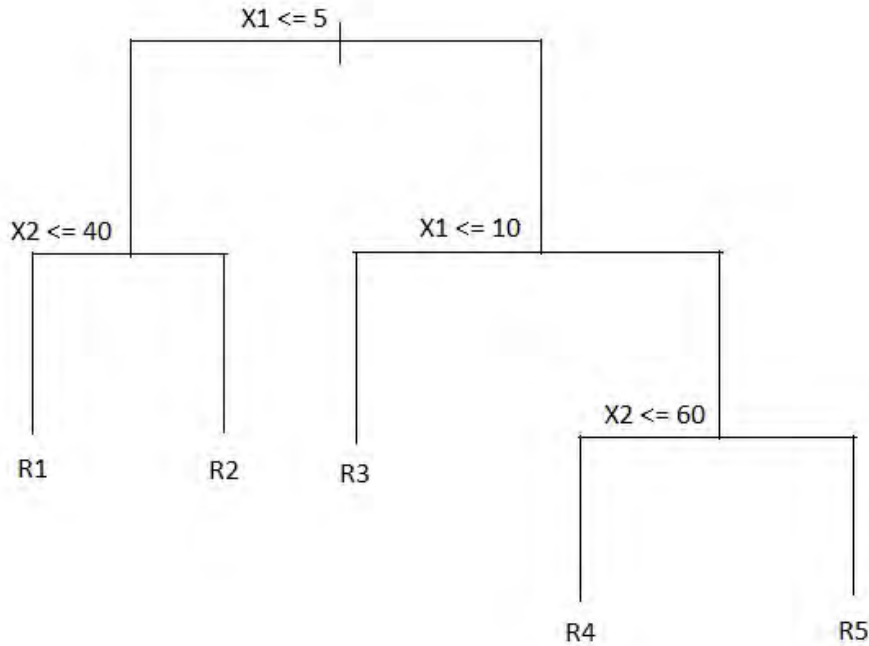
(D)  $\beta_1 = 1$

(E)  $|\beta_1| < 1$

65. DELETED



66. You are given the following regression tree (where the inequality represents the values in the left hand branch of a split).



Determine the set of regions that represents this regression tree.

- (A) Region 1:  $\{X \mid X_1 \leq 5, X_2 \leq 40\}$   
 Region 2:  $\{X \mid X_1 \leq 5, X_2 > 40\}$   
 Region 3:  $\{X \mid X_1 > 5\}$   
 Region 4:  $\{X \mid X_1 > 10, X_2 \leq 60\}$   
 Region 5:  $\{X \mid X_1 > 10, X_2 > 60\}$
- (B) Region 1:  $\{X \mid X_1 \leq 5, X_2 \leq 40\}$   
 Region 2:  $\{X \mid X_1 \leq 5, X_2 > 40\}$   
 Region 3:  $\{X \mid X_1 > 5, X_2 \leq 60\}$   
 Region 4:  $\{X \mid X_1 > 10, X_2 \leq 60\}$   
 Region 5:  $\{X \mid X_1 > 10, X_2 > 60\}$
- (C) Region 1:  $\{X \mid X_1 \leq 5, X_2 > 40\}$   
 Region 2:  $\{X \mid X_1 \leq 5, X_2 \leq 40\}$   
 Region 3:  $\{X \mid 5 < X_1 \leq 10\}$   
 Region 4:  $\{X \mid X_1 > 10, X_2 \leq 60\}$   
 Region 5:  $\{X \mid X_1 > 10, X_2 > 60\}$

- (D) Region 1:  $\{X \mid X_1 \leq 5, X_2 \leq 40\}$   
Region 2:  $\{X \mid X_1 \leq 5, X_2 > 40\}$   
Region 3:  $\{X \mid 5 < X_1 \leq 10\}$   
Region 4:  $\{X \mid X_1 > 10, X_2 \leq 60\}$   
Region 5:  $\{X \mid X_1 > 10, X_2 > 60\}$

- (E) Region 1:  $\{X \mid X_1 \leq 5, X_2 \leq 40\}$   
Region 2:  $\{X \mid X_1 \leq 5, X_2 > 40\}$   
Region 3:  $\{X \mid X_1 > 5, X_2 \leq 60\}$   
Region 4:  $\{X \mid X_1 > 10, X_2 > 60\}$   
Region 5:  $\{X \mid X_1 > 10, X_2 \leq 60\}$

67. A and B each attempt a task 10 times. A is successful 8 out of 10 times while B is successful 2 out of 10 times. They fit a logistic regression model to predict an individual's success with linear term  $\beta_0 + \beta_1 X$ , where  $X$  is an indicator variable for individual A.

The full model estimates that A will be successful 80% of the time and B will be successful 20% of the time. Under the null hypothesis that  $\beta_1 = 0$ , the estimate is that both A and B will be successful 50% of the time. The likelihood ratio test is applied to the logistic regression model.

Determine the lowest available  $p$ -value at which the null hypothesis is rejected.

- (A) 0.100  
(B) 0.050  
(C) 0.025  
(D) 0.010  
(E) 0.005

## SOLUTIONS

1. **Key: E**

First, calculate the distance between pairs of elements in each set. There are two pairs here:

$$x_1, x_4 : \sqrt{(-1-5)^2 + (0-10)^2} = \sqrt{136} = 11.66$$

$$x_2, x_4 : \sqrt{(1-5)^2 + (1-10)^2} = \sqrt{97} = 9.85.$$

For complete linkage, the dissimilarity measure used is the maximum, which is 11.66.

2. **Key: C**

I is false because the number of clusters is pre-specified in the  $K$ -means algorithm but not for the hierarchical algorithm.

II is also false because both algorithms force each observation to a cluster so that both may be heavily distorted by the presence of outliers.

III is true.

3. **Key: D**

$$y_{10} = y_0 + c_1 + \cdots + c_{10} = y_0 + 20$$

$$y_{19} = y_{10} + c_{11} + \cdots + c_{19} = y_0 + 20 + c_{11} + \cdots + c_{19} = y_0 + 20 + 26 = y_0 + 46$$

$$\hat{y}_{19} = y_{10} + \hat{c}_{11} + \cdots + \hat{c}_{19} = y_0 + 20 + 9(2) = y_0 + 38$$

$$y_{19} - \hat{y}_{19} = (y_0 + 46) - (y_0 + 38) = 8.$$

4. **Key: B**

$$c_t = y_t - y_{t-1} \text{ and hence } c_1, c_2, \dots, c_{10} = 2, 3, 5, 3, 5, 2, 4, 1, 2, 3.$$

The mean of the  $c$  values is 3, the variance is  $(1+0+4+0+4+1+1+4+1+0)/9 = 16/9$ . The standard deviation is  $4/3$ . The standard error of the forecast is  $(4/3)\sqrt{9} = 4$ .

5. **Key: E**

Statement I is correct – Principal components provide low-dimensional linear surfaces that are closest to the observations.

Statement II is correct – The first principal component is the line in p-dimensional space that is closest to the observations.

Statement III is correct – PCA finds a low dimension representation of a dataset that contains as much variation as possible.

Statement IV is correct – PCA serves as a tool for data visualization.

6. **Key: E**

Statement I is incorrect – The proportion of variance explained by an additional principal component decreases or stays the same as more principal components are added.

Statement II is correct – The cumulative proportion of variance explained increases or stays the same as more principal components are added.

Statement III is incorrect – We want to use the least number of principal components required to get the best understanding of the data.

Statement IV is correct – Typically, the number of principal components is chosen based on a scree plot.

7. The intent is to model a binary outcome, thus a classification model is desired. In GLM, this is equivalent to binomial distribution. The link function should be one that restricts values to the range zero to one. Of linear and logit, only logit has this property.

8. **Key: D**

Alternative fitting procedures will tend to remove the irrelevant variables from the predictors, thus resulting in a simpler and easier to interpret model. Accuracy will likely be improved due to reduction in variance.

9. **Key: E**

The total Gini index for Split 1 is  $2[20(12/20)(8/20) + 80(18/80)(62/80)]/100 = 0.375$  and for Split 2 is  $2[10(8/10)(2/10) + 90(22/90)(68/90)]/100 = 0.3644$ . Smaller is better, so Split 2 is preferred. The factor of 2 is due to summing two identical terms (which occurs when there are only two classes).

The total entropy for Split 1 is  $-[20(12/20)\ln(12/20) + 20(8/20)\ln(8/20) + 80(18/80)\ln(18/80) + 80(62/80)\ln(62/80)]/100 = 0.5611$  and for Split 2 is  $-[10(8/10)\ln(8/10) + 10(2/10)\ln(2/10) + 90(22/90)\ln(22/90) + 90(68/90)\ln(68/90)]/100 = 0.5506$ . Smaller is better, so Split 2 is preferred.

For Split 1, there are  $8 + 18 = 26$  errors and for Split 2 there are  $2 + 22 = 24$  errors. With fewer errors, Split 2 is preferred.

10. **Key: C**

It is false because with random forest a new subset of predictors is selected for each split.

11. **Key: E**

Solution: SSE is sum of the squared differences between the observed and predicted values. That is,  $[(2-4)^2 + (5-3)^2 + (6-9)^2 + (8-3)^2 + (4-6)^2] = 46$ .

12. **Key: E**

A is false, linear regression is considered inflexible because the number of possible models is restricted to a certain form.

B is false, the lasso determines the subset of variables to use while linear regression allows the analyst discretion regarding adding or moving variables.

C is false, bagging provides additional flexibility.

D is false, there is a tradeoff between being flexible and easy to interpret.

13. **Key E**

I is true. The prediction interval includes the irreducible error, but in this case it is zero.

II is true. Because it includes the irreducible error, the prediction interval is at least as wide as the confidence interval.

III. is false. It is the confidence interval that quantifies this range.

14. **Key: B**

Adding a constant to the dependent variable avoids the problem of the logarithm of zero being negative infinity. In general, a log transformation may make the variance constant.

Hence I is true. Power transformations with the power less than one, such as the square root transformation, may make the variance constant. Hence II is true. A logit

transformation requires that the variable take on values between 0 and 1 and hence cannot be used here.

15. **Key: D**

The cluster centers are A: (0, 1), B: (2, -2), and C: (0, 2). The new assignments are:

| Cluster | Data Point | New Cluster |
|---------|------------|-------------|
| A       | (2, -1)    | B           |
| A       | (-1, 2)    | C           |
| A       | (-2, 1)    | A           |
| A       | (1, 2)     | C           |
| B       | (4, 0)     | B           |
| B       | (4, -1)    | B           |
| B       | (0, -2)    | B           |
| B       | (0, -5)    | B           |
| C       | (-1, 0)    | A           |
| C       | (3, 8)     | C           |
| C       | (-2, 0)    | A           |
| C       | (0, 0)     | A           |

Cluster C has the fewest points with three.

16. **Key: D**

(A) For  $K$ -means the initial cluster assignments are random. Thus different people can have different clusters, so the statement is not true for  $K$ -means clustering.

(B) For  $K$ -means the number of clusters is set in advance and does not change as the algorithm is run. For hierarchical clustering the number of clusters is determined after the algorithm is completed.

(C) For  $K$ -means the algorithm needs to be re-run if the number of clusters is changed. This is not the case for hierarchical clustering.

(D) This is true for  $K$ -means clustering. Agglomerative hierarchical clustering starts with each data point being its own cluster.

17. DELETED

18. **Key: B**

$$TSS = (\text{Residual Sum of Squares}) / (1 - R^2) = 230 / (1 - 0.64) = 638.89$$

19. **Key: E**

The only two models that need to be considered are the full model with all four coefficients and the null model with only the intercept term. The test statistic is twice the difference of the log-likelihoods, which is 10.

The number of degrees of freedom is the difference in the number of coefficients in the two models, which is three.

At 5% significance with three degrees of freedom, the test statistic of 10 exceeds the 7.81 threshold, so the null hypothesis should be rejected.

20. **Key: C**

(A) is not appropriate because removing data will likely bias the model estimates.

(B) is not appropriate because altering data will likely bias the model estimates.

(C) is correct.

(D) is not appropriate because the canonical link function is the logarithm, which will not restrict values to the range zero to one.

21. **Key: C**

I is true because the mean  $E(y_t) = y_0 + t\mu_c$  depends on  $t$ .

II is false because the variance  $Var(y_t) = t\sigma_c^2 = 0$  does not depend in  $t$ .

III is true because the variance depends on  $t$ .

22. **Key: A**

The intercept term may be any value, hence (A) is false.

23. **Key: C**

The regression line is  $y = 0.40 + 0.05x$ . This can be obtained by either using the formula for the regression coefficients or by observing that the three points lie on a straight line and hence that line must be the solution. A 24-ounce cup costs 1.60. Two of them cost 3.20. A 48-ounce cup costs \$2.80. So the savings is 0.40.



24. **Key: C**

Even though the formula for  $R^2$  involves RSS and TSS, she just needs their ratio, which can be obtained from  $F$ .

$$F = \frac{(TSS - RSS) / 4}{RSS / (105 - 4 - 1)} = 20$$

$$\frac{TSS - RSS}{RSS} = 20(4) / 100 = 0.80$$

$$\frac{TSS}{RSS} = 1.80$$

$$\frac{RSS}{TSS} = \frac{1}{1.80}$$

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{1}{1.80} = \frac{0.80}{1.80} = 0.44$$

25. **Key: C**

I is true because the method optimizes with respect to the training set, but may perform poorly on the test set.

II is false because additional splits tends to increase variance due to adding to the complexity of the model.

III is true because in this case only the training error is measured.

26. **Key: A**

Each step must divide the one the existing regions into two parts, using either a vertical or a horizontal line. Only I can be created this way.

27. **Key: C**

Only variables with a  $p$ -value greater than 0.05 should be considered. Because only one of the variables (Number of weekend days) meets this criterion, it should be dropped.

28. DELETED

29. **Key: E**

All three statements are true. See Section 8.1 of *An Introduction to Statistical Learning*. The statement that trees are easier to explain than linear regression methods may not be obvious. For those familiar with regression but just learning about trees, the reverse may be the case. However, for those not familiar with regression, relating the dependent variable to the independent variables, especially if the dependent variable has been transformed, can be difficult to explain.

30. **Key: D**

I is false because the loadings are unique only up to a sign flip.

II is true. Principal components are designed to maximize variance. If there are no constraints on the magnitude of the loadings, the variance can be made arbitrarily large. The PCA algorithm's constraint is that the sum of the squares of the loadings equals 1.

III is true because four components can capture all the variation in four variables, provided there are at least four data points (note that the problem states that the data set is large).

31. **Key: D**

See Page 242 of *Regression Modeling with Actuarial and Financial Applications*.

I is true because a random walk is characterized by a linear trend and increasing variability.

II is true because differencing removes the linear trend and stabilizes the variance.

III is true as both the linear trend and the increasing variability contribute to a higher standard deviation.

32. **Key: E**

I and II are both true because the roles of rows and columns can be reversed in the clustering algorithm. (See Section 10.3 of *An Introduction to Statistical Learning*.)  
III is true. Clustering is unsupervised learning because there is no dependent (target) variable. It can be used in exploratory data analysis to learn about relationships between observations or features.

33. **Key: E**

LCA(I) = 8.146

LCA(II) = 8.028

LCA(III) = 7.771

34. **Key: B**

I is false.  $K$ -means clustering is subject to the random initial assignment of clusters.

II is true. Hierarchical clustering is deterministic, not requiring a random initial assignment.

III is false. The two methods differ in their approaches and hence may not yield the same clusters.

35. **Key: B**

The first PC explains about 62% of the variance. The second PC explains about 23% of the variance. Combined they explain about 85% of the variance, and hence two PCs are sufficient.

36. **Key: D**

I is false. All observations are assigned to a cluster.

II is true. By cutting the dendrogram at different heights, any number of clusters can be obtained.

III is true. Clustering methods have high variance, that is, having a different random sample from the population can lead to different clustering.

37. **Key: B**

I is true. Uniqueness up to a sign flip means all three signs must be flipped. This is true for X and Y.

II is true. The presence of absence of scaling can change the loadings.

III is false. Uniqueness up to a sign flip means all three signs must be flipped. For W and X only the second loading is flipped.

38. **Key: D**

I is true. See formula (7.8) in the Frees text.

II. is true. See formula (7.9) in the Frees text.

III is true. The only difference is the error terms, which are difficult to compare. See page 243 in the Frees text.

39. **Key: E**

(A) is false. Trees work better with qualitative data.

(B) is false. While trees accommodate nonlinear relations, as seen in (E) a linear model can work very well here.

(C) is false. The variance is constant, so that is not an issue here.

(D) is false. There is a clear relationship as noted in answer (E).

(E) is true. The points appear to lie on a quadratic curve so a model such as

$y = \beta_0 + \beta_1 x + \beta_2 x^2$  can work well here. Recall that linear models must be linear in the coefficients, not necessarily linear in the variables.

40. **Key: C**

I is true. At the lowest height, each observation is its own cluster. The number of clusters decreases as the height increases.

II is false. There is no need to plot the data to perform  $K$ -means clustering.

III is true.  $K$ -means does a fresh analysis for each value of  $K$  while for hierarchical clustering, reduction in the number of clusters is tied to clusters already made. This can miss cases where the clusters are not nested.

41. **Key: B**

I is true. Random forests differ from bagging by setting  $m < p$ .

II is true.  $p - m$  represents the splits not chosen.

III is false. Typical choices are the square root of  $p$  or  $p/3$ .

42. **Key: E**

The logit and probit models are similar (see page 307 of Frees, which also discusses items A-D).

43. **Key: D**

Item D is a statement about principal components analysis, not clustering.

44. **Key: C**

The model of Actuary 1 is the null model and hence values from it are not needed. The

$$\text{solution is } F = \frac{(TSS - RSS) / 1}{RSS / (n - 2)} = \frac{490,000 - 250,000}{250,000 / 98} = 94.08.$$

45. **Key: C**

Let  $T$  be the study time offered. The equation to solve is

$$6.0 = \exp[-0.1 + 0.5(2T) + 0.5 - 0.1 + 0.2] = \exp(0.5 + T)$$

$$1.79176 = 0.5 + T$$

$$T = 1.29176 \text{ or } 129 \text{ hours.}$$

46. **Key: D**

The smoothed forecast at 100 is  $0.2(100.2) + 0.8(95.1) = 96.1$ . This is also the forecast at 102.

The smoothed values are at 99:  $0.2(95.1) + 0.8(89.9) = 90.9$ . At 100:  $0.2(96.1) + 0.8(90.9) = 91.9$ . The trend is  $0.2(96.1 - 91.9) / .8 = 1.05$ . The intercept (value at 100) is  $2(96.1) - 91.9 = 100.3$ . The forecast at time 102 is  $100.3 + 2(1.05) = 102.4$ .

The difference is 6.3.

47. DELETED

48. **Key: B**

The first split is  $X_1 < t_1$ . This requires a horizontal line at  $t_1$  on the vertical axis. Graphs B, C, and E have such a line.

The second split is the case where the first split is true. That means all further action is below the line just described. All three graphs do that. The second split is  $x_2 < t_2$ . This requires a vertical line at  $t_2$  on the horizontal axis with the line only going up to  $t_1$ . Again, all three graphs have this.

The third split is when the second split is true. That means all further action is to the left of the line just described. That rules out graph C. The only difference between graphs B and E is which part relates to node C and which to node D. The third split indicates that node C is the case when  $x_1 < t_3$ . Only graph B has this region marked as C.

49. **Key: C**

I is true. A confidence interval reflects the error in estimating the expected value. With an infinite amount of data, the error goes to zero.

II is false. A prediction interval also reflects the uncertainty in the predicted observation. That uncertainty is independent of the sample size and thus the interval cannot go to zero.

III is true. The additional uncertainty in making predictions about a future value as compared to estimating its expected value leads to a wider interval.

50. **Key: B**

I is false. Highly flexible models are harder to interpret. For example, a ninth degree polynomial is harder to interpret than a straight line.

II is true. Inference is easier when using simple and relatively inflexible methods. Lasso is simpler and less flexible than bagging.

III is false. Flexible methods tend to overfit the training set and be less accurate when applied to unseen data.

51. **Key: C**

For duck X: Age = 7 => go left. Gender = Male => go right. Prediction = 0.90 kg.

For duck Y: Age = 5 => go left. Gender = Female => go left. Prediction = 0.80 kg.

For duck Z: Age = 8 => go right. Gender = Male => go right. Wing Span = 5.7 => go right. Prediction = 1.25.

$Y < X < Z$ .

52. **Key: D**

I is true. See Page 348 of Frees.

II is true. See Page 347 of Frees.

III is true. See Page 347 of Frees.

53. **Key: A**

(A) is false.  $\beta_0$  is the expected value of  $Y$  when  $X = 0$ .

The other four statements are true (see Page 63 of James, et al.)

54. **Key: C**

The variables TOPSCHOOL and LARGECITY both lack significance and are candidates for removal. However, only one variable should be removed at a time, and it should be the one with the highest  $p$ -value, which is TOPSCHOOL. After removing TOPSCHOOL and rerunning the model it is possible that LARGECITY will become significant.

55. **Key: B**

The three one-step predicted values are  $12+2.25 = 14.25$ ;  $15+2.25 = 17.25$ ; and  $21+2.25 = 23.25$ . The errors are  $15 - 14.25 = 0.75$ ;  $21 - 17.25 = 3.75$ ; and  $22 - 23.25 = -1.25$ .

$$F = (0.75 + 3.75 - 1.25)/3 = 1.083$$

$$G = (0.75^2 + 3.75^2 + 1.25^2)/3 = 5.396$$

The absolute difference is 4.313.

56. **Key: B**

A is false because different models are likely to produce different results.

B is true because using the fitted model implies that this model continues to apply.

C is false as there is no easy way to compare reliability of the two approaches.

D is false in that a narrower interval provides more useful information about the true value.

E is false in that the interval contains the most likely values with the point prediction being the single most likely point.

57. **Key: A**

T1 has observations with  $Z \leq 3$  and  $Y = A$  or  $B$ , which are observations 1, 5, and 9. The values are 4.75, 4.53, and 3.89, which average to 4.39.

T2 has observations with  $Z \leq 3$  and  $Y = C$  or  $D$ , which are observations 3, 4, and 6. The values are 4.67, 4.56, and 3.91, which average to 4.38.

T3 has observations with  $Z > 3$  and  $X = F$ , which are observations 2 and 7. The values are 4.67 and 3.90, which average to 4.29.

T4 has observations with  $Z > 3$  and  $X = M$ , which is observation 8. The value is 3.90.



58. **Key: C**

The estimate of  $\beta_1$  is

$$b_1 = r_1 = \frac{(31-40)(35-40) + (35-40)(37-40) + (37-40)(41-40) + (41-40)(45-40) + (45-40)(51-40)}{(31-40)^2 + (35-40)^2 + (37-40)^2 + (41-40)^2 + (45-40)^2 + (51-40)^2} = \frac{117}{262} = 0.4466.$$

The estimator of  $\beta_0$  is  $b_0 = \bar{y}(1 - r_1) = 22.136$ .

The residuals are then

$$e_2 = 35 - (22.136 + 0.4466 \times 31) = -0.9806$$

$$e_3 = 37 - (22.136 + 0.4466 \times 35) = -0.7670$$

$$e_4 = 41 - (22.136 + 0.4466 \times 37) = 2.3398$$

$$e_5 = 45 - (22.136 + 0.4466 \times 41) = 4.5534$$

$$e_6 = 51 - (22.136 + 0.4466 \times 45) = 8.7670.$$

The average residual is  $\bar{e} = 2.78252$  and then the mean square error is

$$s^2 = \frac{(-0.9806 - 2.78252)^2 + (-0.7670 - 2.78252)^2 + (2.3398 - 2.78252)^2 + (4.5534 - 2.78252)^2 + (8.7670 - 2.78252)^2}{6 - 3} = 21.969.$$

59. **Key: C**

The means for cluster 1 are  $(1 + 0 + 6)/3 = 2.3333$  for  $X_1$  and  $(3 + 4 + 2)/3 = 3$  for  $X_2$  and the variation is

$$(1 - 2.3333)^2 + (3 - 3)^2 + (0 - 2.3333)^2 + (4 - 3)^2 + (6 - 2.3333)^2 + (2 - 3)^2 = 22.6667.$$

The means for cluster 2 are  $(5 + 1)/2 = 3$  for  $X_1$  and  $(2 + 6)/2 = 4$  for  $X_2$  and the variation is  $(5 - 3)^2 + (2 - 4)^2 + (1 - 3)^2 + (6 - 4)^2 = 16$ .

The total within-cluster variation is (per equation (10.12) in the first edition of *ISLR*)  $2(22.6667 + 16) = 77.33$ .

60. **Key: C**

I is false. Setting  $K = n$  will almost certainly overfit as there are unlikely to be that many true clusters.

II is false. Within-cluster variation is minimized when  $K = n$ , which as noted above is unlikely to be optimal.

III is true. There is no exact method for determining the optimal value of  $K$ .

61. **Key: D**

I is false, AIC and BIC make an indirect estimate by adjusting the training error.

II is true, for a fixed value of the number of predictors, the two provide the same ranking.

III is true as for  $n > 7$ , BIC provides a greater penalty for additional variables, and hence will select a number less than equal to that selected by AIC.

62. **Key: A**

Using the rule of thumb  $t$ -value of 2 produces an interval of  $-1.03 \pm 2(0.06)$  for an interval of  $(-1.15, -0.91)$ . Note that if the more precise value of 1.96 is used, the same answer (to two decimal places) is obtained.

63. **Key: D**

For the first split,  $\text{ShelveLoc} = \text{Good}$ , the answer is no, and the branch to the right is taken. For the second split,  $\text{Price} = 120 \geq 110$ , the answer is yes, and the branch to the left is taken. The node has a predicted value of 9.2.

64. **Key: E**

From the definition of a stationary AR(1) process on Page 254 of Freed, the parameter  $\beta_0$  can be any fixed constant, ruling out answers (A) and (B). To be stationary it is necessary that  $-1 < \beta_1 < 1$ , which makes answer (E) correct.

65. **DELETED**

66. **Key:D**

To be in Region 1,  $X_1$  must be  $\leq 5$  and  $X_2$  must be  $\leq 40$ . This eliminates answer C

To be in Region 2,  $X_1$  must be  $\leq 5$  and  $X_2$  must be  $> 40$ . This also eliminates answer C.

To be in Region 3,  $X_1$  must be  $> 5$  and be  $\leq 10$ . This eliminates answers A, B, and E.

That leaves only answer D as possible. To confirm:

To be in Region 4,  $X_1$  must be  $> 10$  and  $X_2$  must be  $\leq 60$ .

To be in Region 5,  $X_1$  must be  $> 10$  and  $X_2$  must be  $> 60$ .

Both of these are satisfied by answer D.

67. **Key: D**

In the full model, there are 16 outcomes with a probability of 0.8 and 4 outcomes with a probability of 0.2, so the loglikelihood is  $16 \ln 0.8 + 4 \ln 0.2 = -10.01$ . In the reduced model, the probability of each of the 20 outcomes is 0.5, so the loglikelihood is  $20 \ln 0.5 = -13.86$ .

The likelihood ratio test statistic (equal to twice the difference in loglikelihoods) has a chi-square distribution with degrees of freedom equal to the difference in the number of explanatory variables, in this case 1. The test statistic is  $2(-10.01 - (-13.86)) = 7.71$ . The critical values for one degree of freedom are 6.635 and 7.879 for  $p = 0.010$  and 0.005 respectively, so the null hypothesis can be rejected at 0.010 but not 0.005.