

Name:

UTeid:

M339G Predictive Analytics
University of Texas at Austin
Mock In-Term Exam II
Instructor: Milica Ćudina

Notes: This is a closed book and closed notes exam. The maximal score on this exam is points.

All written work handed in by the student is considered to be
their own work, prepared without unauthorized assistance.

The University Code of Conduct

”The core values of The University of Texas at Austin are learning, discovery, freedom, leadership, individual opportunity, and responsibility. Each member of the university is expected to uphold these values through integrity, honesty, trust, fairness, and respect toward peers and community. As a student of The University of Texas at Austin, I shall abide by the core values of the University and uphold academic integrity.”

”I agree that I have complied with the UT Honor Code during my completion of this exam.”

Signature:

2.1. CONCEPTUAL QUESTIONS.

Problem 2.1. (10 points) Describe the differences and the similarities between LDA and the logistic regression for a two-class classification problem.

Problem 2.2. (10 points) Justify this statement from the textbook:

"Neither QDA nor naive Bayes is a special case of the other."

2.2. FREE RESPONSE PROBLEMS. Please, explain carefully all your statements and assumptions. Numerical results or single-word answers without an explanation (even if they're correct) are worth 0 points.

Problem 2.3. (10 points) *Source: Pitman's "Probability."*

Here is a summary of PSAT and SAT scores of a large group of students.

	mean	standard deviation
PSAT	1200	100
SAT	1300	90

Assume that the data are modeled as bivariate normal with the correlation coefficient equal to 0.6. Of the students who scored 1000 on the PSAT, about what percentage scored above average on the SAT?

Problem 2.4. (20 points) Consider the following observations of (X, Y) with X being the predictor and Y being the response:

$$(1, 4), \quad (2, 6), \quad (3, 5), \quad (6, 1).$$

After one iteration of recursive binary splitting, there are two groups of observations. Find the members of the two groups.

2.3. MULTIPLE CHOICE QUESTIONS.

Problem 2.5. (5 points) *Source: MAS-II, Fall 2019.*

You are given the following three statements about tree-based methods for regression and classification:

- I. The main difference between bagging and random forests is the number of predictors considered at each step in building individual trees.
 - II. Single decision tree models generally have a higher variance than random forest models.
 - III. Random forests provide an improvement over bagging because trees in a random forest are less correlated than those in bagged trees.
- (a) I only.
 (b) II only.
 (c) III only.
 (d) I, II and III.
 (e) The correct answer is not given above.

Problem 2.6. (5 points) Consider the following data set with the explanatory random variable X and the categorical response Y :

X	1	2	6	8	12	16	17	20	22
Y	N	N	L	N	L	N	L	L	L

Determine which of these splits is/are the best using classification error as the criterion.

- I. $R = \{X \leq 7\}$ and $R^c = \{X > 7\}$
- II. $R = \{X \leq 10\}$ and $R^c = \{X > 10\}$
- III. $R = \{X \leq 14\}$ and $R^c = \{X > 14\}$

- (a) I only.
 (b) II only.
 (c) I and II only.
 (d) I and III only.
 (e) II and III only.

Problem 2.7. (5 points) Consider the following statements involving classification trees.

- I. The use of Gini index or cross-entropy may result in split nodes with the same predicted class.
- II. Cross-validation can be used to prune the trees.
- III. In a single node with two possible classes, the Gini index always exceeds the cross-entropy.

Which of the statements above are true?

- (a) I only.
 (b) II only.
 (c) I and II only.
 (d) I and III only.
 (e) II and III only.