



Cumulative Practice Problem Set

DS 1000B Winter 2026

Notes

This practice problem set is intended to support your preparation for your midterm and exam. The questions reflect the range of difficulty you may encounter on the actual exam; some may be more challenging, while others may be less so.

- Some problems may be approached differently depending on the chosen methodology. A key example is the analysis of skewness: you may choose to either exclude outliers before describing the tails (Marieke's approach) or include the outliers in your description (Pavel's approach).

Neither method is objectively superior to the other; they are simply different conventions. Consequently, both approaches will be accepted for full credit.

Students whose answers align with Marieke's framework should not be concerned if their steps differ from the text below. Grading will account for both philosophical approaches.

- These problems are **not** representative of the overall length or coverage of the midterm or final exam.
- The midterm information sheet posted by February 14th, which will include exam logistics, the provided formula sheet, and general details about exam weighting.
- You may discuss these problems during office hours, on Piazza or with your classmates.

Chapter 1: Data Visualization

Circle one answer per question.

MC1. A survey involving 35 questions was given to a sample of 200 students attending a university with an enrollment of 10,000. How many individuals do the data describe?

- A.** 35 **B.** 200 **C.** 10,000 **D.** It is impossible to say.

MC2. A survey involving 35 questions was given to a sample of 200 students attending a university with an enrollment of 10,000. How many variables do the data contain?

- A.** 35 **B.** 200 **C.** 35×200 **D.** $35 \times 10,000$

MC3. A dataset of exam scores shows a histogram with a long tail to the right. Which statement is most likely true?

- A.** The mean is less than the median
B. The distribution is symmetric
C. The mean is greater than the median
D. The mode equals the mean
E. The distribution is uniform

MC4. A stemplot of test scores is provided below:

Stem	Leaf
3	5 7 8
4	2 6 7 8 9
5	8
6	1 5 5 8 9
7	0 2 2 4 6 7
8	1 3 4
9	0 2

Key: $5|8 = 58$. How many students scored in the 70s?

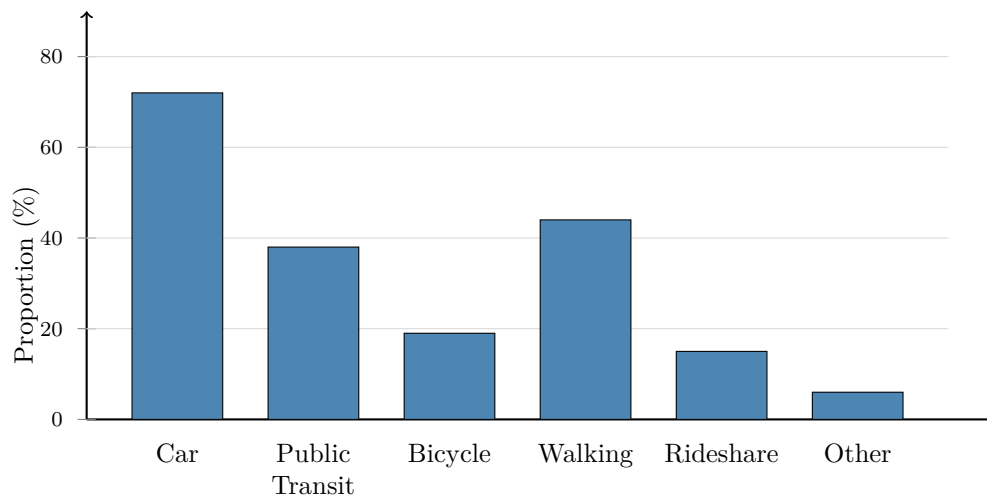
- A.** 2
B. 3
C. 4
D. 5
E. 6

MC5. Which statement is correct regarding the distribution of a variable?

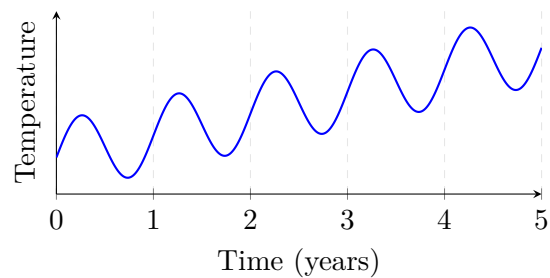
- A. The distribution of a quantitative variable can be displayed using a graphical display like a histogram.
- B. The distribution of a categorical variable can be displayed using a table of counts or percents.
- C. The distribution of a variable tells us what values the variable takes and how often it takes these values.
- D. All of these choices are correct.

MC6. The bar chart below displays the percentage of commuters by their primary modes of transportation. Is the total proportion of people who commute by bicycle or by walking greater than, less than, or equal to the proportion who commute by car?

- A. Greater than
- B. Equal to
- C. Less than
- D. Cannot be determined



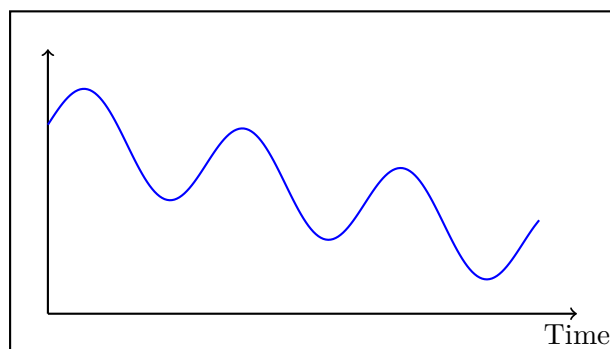
MC7. Where does the following time series plot exhibit deviation?



- A. 0
- B. Approximately at 1.4
- C. There are no deviations in this plot
- D. 5
- E. 2.7

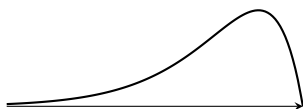
Short Answer

- (a) Comment on the following time series plot. What features does it exhibit?

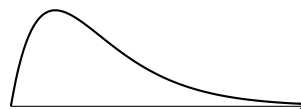


- (b) Two density plots are sketched below. One shows the distribution of age at death from natural causes (heart disease, cancer, and so forth). The other shows age at death from trauma (accident, murder, suicide). Which is which, and why?

(i)



(ii)



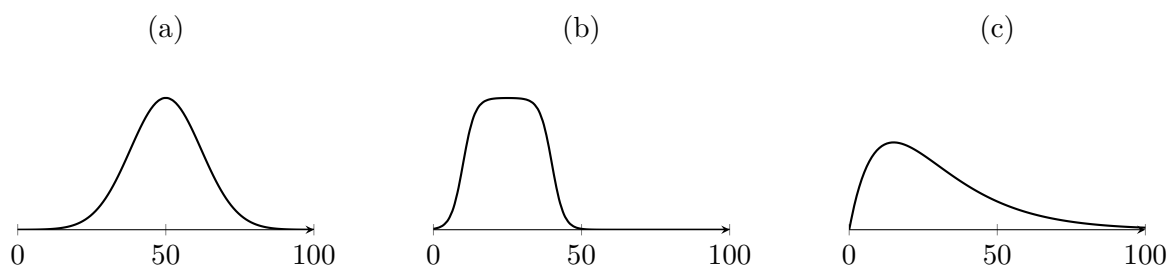
Chapter 2: Mean, Median, Quartiles, Boxplots, and Standard Deviation

- MC8.** A dataset has values: 10, 15, 20, 25, 100. If the observation 100 is removed, which measure of center changes the most?
- A. Median B. Mode C. Mean D. All change equally
- MC9.** Two classes have the same mean test score of 75. Class A has a standard deviation of 5, while Class B has a standard deviation of 15. Which statement is true?
- A. Class A scores are more spread out than Class B
B. Class B scores are more spread out than Class A
C. Both classes have identical score distributions
D. Class A has higher scores on average
E. Standard deviation does not measure spread
- MC10.** The five number summary of a dataset is given by: Min = 20, Q1 = 30, Median = 40, Q3 = 50, Max = 60. What is the interquartile range (IQR)?
- A. 10 B. 30 C. 20 D. 40
- MC11.** A dataset has mean 50 and standard deviation 10. If every value is multiplied by 2 and then 5 is added, what are the new mean and standard deviation?
- A. Mean = 50, SD = 10
B. Mean = 100, SD = 20
C. Mean = 105, SD = 20
D. Mean = 105, SD = 10
E. Mean = 55, SD = 25
- MC12.** For a right-skewed distribution, which inequality is most likely true?
- A. Mean < Median < Mode
B. Mode < Median < Mean
C. Median < Mode < Mean
D. Mean = Median = Mode
E. Cannot be determined without data

Short Answer

- (a) For each density plot below, indicate if the median is greater than ($>$), less than ($<$), or approximately equal to (\approx) the mean.

Plot	Mean vs Median
(a)	Mean ____ Median
(b)	Mean ____ Median
(c)	Mean ____ Median



- (b) For each density plot above, approximate the mode of the distribution.

a) Mode: _____

b) Mode: _____

c) Mode: _____

- (c) Compute the mean and median of the following data: 2, 4, 6, 8, 10.

- (d) Compute the mean and median of the following data: 1, 3, 3, 6, 7, 8, 9.

(e) Find the third quartile (Q_3) of the following data: 3, 5, 7, 9, 12, 15, 18.

Question 2

A dataset consists of the values 10, 15, 20, 25, and a fifth unknown value x . After including x , the median of the dataset is 18.

What are the possible values of x ? Explain your reasoning.

Chapter 4: Scatterplots and Correlation

MC13. A researcher finds that the correlation between hours of sleep and reaction time is $r = -0.72$. If a new dataset is collected where reaction times are measured in milliseconds instead of seconds, what happens to the correlation coefficient?

- A. It becomes positive
- B. It increases in magnitude
- C. It stays the same
- D. It becomes zero
- E. It decreases in magnitude

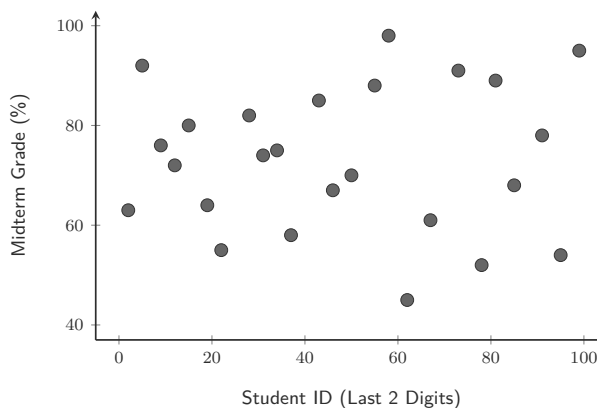
MC14. A study finds a correlation of $r = 0.89$ between ice cream sales and drowning deaths. Which conclusion is most appropriate?

- A. Ice cream consumption causes drowning
- B. Drowning causes people to buy ice cream
- C. A third variable (such as hot weather) may explain both
- D. The correlation proves a causal relationship
- E. The relationship must be nonlinear

Short Answer

(a) What is a scatterplot?

(b) Which features do we examine when commenting on a scatterplot? Comment on them for the scatterplot below.



Chapter 5: Linear Regression

MC15. In a regression predicting test scores from study hours, the equation is $\hat{y} = 45 + 5x$. A student who studies 8 hours has an actual score of 82. What is this student's residual?

- A. 37 B. 40 C. 74 D. -3 E. 3

MC16. Suppose the variables x and y have standard deviations $s_x = s_y = 1$. What is the slope of the regression line when regressing y on x ?

- A. The slope is 0
B. The slope is 1
C. The slope equals the correlation coefficient r
D. The slope is s_y/s_x
E. Cannot be determined without knowing the means

Question 3

A marketing analyst studies the relationship between advertising spending (in thousands of dollars) and monthly sales (in thousands of units). Data from 15 months yields the estimated regression:

$$\widehat{\text{Sales}} = 12.5 + 3.2 \times \text{Advertising}$$

The correlation coefficient is $r = 0.85$.

(a) Interpret the slope and intercept in context of the problem.

(b) Calculate R^2 and interpret it in context of the problem.

- (c) One month had advertising spending of \$8,000 and actual sales of 42,000 units. Calculate the residual. Is this an overestimate or underestimate?

Question 4

The table below shows the relationship between hours of study (x) and test scores (y) for four students:

Hours of Study (x)	Test Score (y)
1	50
2	70
3	65
4	85

- (a) Determine the least squares regression line for predicting test score from hours of study.

(b) Interpret the slope of the regression line in the context of the problem.

(c) A student studies for 2.5 hours. What is the predicted test score for this student?

Question 5

Suppose that we know that the regression line of y onto x is given by

$$\hat{y} = 3 + 0.4x$$

If $s_y = 2$, $s_x = 4$, $\bar{x} = 0$, find an expression for the regression line of x onto y .

Question 6

Reference: DS1000A Midterm Questions 23–24

Context: Suppose it was found that the least-squares regression line predicting house price (y) from size in square feet (x) is given by

$$\hat{y} = 20000 + 150x$$

Q23: Which of the following gives the least-squares regression line predicting square footage (y) based on house price (x)?

It was claimed that the correct answer is:

The line cannot be determined from the information provided.

Explain why this is the correct answer and indicate what additional information would be needed to find the regression line of square footage on house price.

Chapter 6: Contingency Tables

Question 7

The popularity of computer, video, online, and virtual reality games has raised concerns about their ability to negatively affect youth. The data below are based on a recent survey of 14 to 18 year olds in Connecticut high schools.

	Grade Average		
	A's and B's	C's	D's and F's
Played games	736	450	193
Never played games	205	144	80

- (a) How many people does this table describe? How many of these have played video games?

- (b) Find the marginal distribution of grades. What percentage of students received a grade of C or lower?

- (c) Determine the conditional distribution of grades for those who play video games.

- (d) Find the relative risk of receiving a grade of C or lower for those who play video games compared to those who do not.

- (e) Find the relative odds of playing video games for those who received A's and B's compared to those who received D's and F's.

- (f) Interpret the relative odds calculated in part (e) in the context of this problem.

Chapters 12 and 13: Probability

MC17. A weighted four-sided die has the following probabilities: 1 shows with probability x/c , 2 with $2x/c$, 3 with $4y/c$, and 4 with $8z/c$, where $x, y, z > 0$. Which of the following must be true of c ?

A. $c = x + 2x + 4y + 8z$

B. $c = x + y + z$

C. $c = 2x + 4y + 8z$

D. $c = 3x + 4y + 8z$

Short Answer

- (a) What are the two types of random variables? Explain the difference(s) between them.

Question 8

The symmetric difference of two events A and B is defined as $A\Delta B = (A \cap B^c) \cup (B \cap A^c)$. Let S be the outcome space consisting of all people. Suppose that A and B are the events

A: people who like apples.

B: people who like bananas.

- (a) Draw a Venn diagram to illustrate the event $A\Delta B$.

(b) Interpret the symmetric difference $A\Delta B$ in the context of this problem.

(c) Explain why $P(A\Delta B) = P(A \cup B) - P(A \cap B)$.

(d) If $P(A) = 0.4$, $P(B) = 0.5$, and $P(A \cap B) = 0.2$, calculate $P(A\Delta B)$.

(e) Suppose A and B are independent and that $P(A) = 0.5$ and $P(B^c) = 0.1$. Compute $P(A\Delta B)$.

Question 9

A coin is tossed 3 times and the result of each toss is recorded as H (heads) or T (tails).

(a) What is an outcome of this random phenomenon?

(b) List the sample space S .

(c) Define the event A as “at most 1 tail observed.” List the outcomes in event A .

(d) What is the probability of event A ?

(e) Let B be the event “each face appears at least once.” List the outcomes in event B .

--

(f) Are events A and B independent? Justify your answer.

(g) Are the events A^c and B^c independent? Justify your answer.

Question 10

Two cards are drawn without replacement from a standard deck of 52 playing cards. Let A be the event that the first card drawn is red, and let B be the event that the second card drawn is red.

(a) Draw a tree diagram representing this scenario, labeling all branches with their probabilities.

(b) Are events A and B independent? Justify your answer.

(c) Calculate $P(A \cap B)$.



(d) Compute $P(B \mid A)$ and $P(A \mid B)$. Compare the two values.

(e) Compute $P(A)$ and $P(B)$. Compare the two values.

Question 11

In a match between Brazil and Germany, let X, Y denote the number of goals scored by Brazil and Germany, respectively. The distributions are:

x	0	1	2	3	4	5	6	7
$P(X = x)$	0	c	0.2	0.2	0.2	0.2	0	0

y	0	1	2	3	4	5	6	7
$P(Y = y)$	$2c$	c	c	$c/2$	0	0	0	d

Assume X and Y are independent.

- (a) Find values of c and d that make X and Y valid random variables.

- (b) What is the probability that Germany scores 7 goals and Brazil scores 1 goal?

(c) Determine the probability the game ends in a tie.

(d) Determine the probability that Germany wins given that Brazil scored at least 4 goals.

Question 12

During press conferences, President Donald Trump is sometimes asked questions about geopolitics. When this occurs, there is a 50% chance that he goes on a tangent. He derides Joe Biden or boasts about having the “hottest economy” only if he is on a tangent.

- If he goes on a tangent, there is an 80% chance that he mentions Joe Biden.
- If he mentions Joe Biden, there is a 40% chance that he also mentions the “hottest economy.”
- The overall probability that he mentions the “hottest economy” is 25%.

Let T, B, H be the events defined as follows

T : Trump goes on a tangent

B : Trump derides Joe Biden

H : Trump boasts about having the “hottest economy”

(a) Draw a tree diagram representing this scenario.

(b) Find the probability $P(B \cup T)$.



(c) Find the probability $P(T \cap B \cap H)$.

(d) Find the probability that Trump boasts about having the “hottest economy” given that he did not deride Joe Biden.

(e) Find the probability $P(H \cup T)$.

- (f) Determine $P(T \cap B^c \cap H^c)$, the probability that Trump goes on a tangent but mentions neither Joe Biden nor the “hottest economy.”

Chapter 3: Normal Distribution

Question 13

Argentina is testing javelin throwers for the Olympics. Distances are normally distributed with mean 70 meters and standard deviation 5 meters.

- (a) What is the probability an athlete breaks the Olympic record of 92.97 meters?

- (b) Find x such that 99% of throws do not land farther than x .

- (c) Athletes are shortlisted if their throw exceeds 80 meters. Given that an athlete was shortlisted, what is the probability their throw was at most 85 meters?

Chapter 15: Sampling Distributions

Question 14

A farm ships crates of potatoes. Potato weight has mean $\mu = 300\text{g}$ and standard deviation $\sigma = 50\text{g}$. Each crate contains 100 potatoes. An inspector randomly samples crates and calculates the average weight per potato \bar{X} .

- (a) What are the mean and standard deviation of the average potato weight per crate \bar{X} ?

- (b) Approximate the probability that the average weight per potato in a crate is between 295g and 305g.

- (c) If the inspector wants to be 95% confident that the average weight per potato in a sampled crate is within 3g of the true mean, how many potatoes should they sample?

- (d) Find the probability that the total weight of potatoes in a crate exceeds 31 kg.

Chapter 16: Confidence Intervals

Question 15

A battery manufacturer wants to estimate the average lifetime of their AA batteries. The population standard deviation is known to be $\sigma = 12$ hours. A random sample of 49 batteries has a sample mean lifetime of 98 hours.

- (a) Calculate the standard error of the sample mean.

- (b) Construct a 99% confidence interval for the mean battery lifetime.

- (c) How many batteries must be sampled to achieve a 95% confidence interval with margin of error no more than 2 hours?

Question 16

A random sample of 64 light bulbs has a mean lifetime of $\bar{x} = 1200$ hours. Assume the population standard deviation is $\sigma = 150$ hours.

- (a) Construct a 90% confidence interval for the mean lifetime of all bulbs.

- (b) Interpret your interval from part (a).

- (c) Recall that the total width of a confidence interval is the distance between the upper and lower bounds. If the interval has total width 40 hours and the sample size is still 64, what is the confidence level?

Chapters 7 and 8: Sampling and Experiments

Question 17

A researcher is interested in learning the proportion of university students that use recreational drugs. They decide to survey students who visit the campus library at opening. A simple random sample of 50 students is taken, and each student is asked whether they use drugs recreationally.

- (a) Define the population, sample, parameter, and statistic in this study.

- (b) Identify two potential sources of bias in this sampling process.

- (c) Suppose the population standard deviation is $\sigma = 0.5$. Calculate the standard error for a sample of $n = 50$ students.

- (d) Define stratified random sampling.

Question 18

A high school wants to investigate whether using flashcards improves students' math exam scores compared to traditional study methods.

- (a) Identify the subjects, factor, and treatments in this experiment.
- (b) Describe a block design that could reduce variance in this study.
- (c) Identify two potential lurking variables.
- (d) Define double blinding and explain its importance. Is it feasible in this experiment?