

AP Statistics - Saturday Practice Exam – SOLUTIONS

- #1) E
- #2) E
- #3) C
- #4) C
- #5) C
- #6) B
- #7) B
- #8) A
- #9) D
- #10) D
- #11) C
- #12) B
- #13) C
- #14) A
- #15) B
- #16) B
- #17) A
- #18) A
- #19) B
- #20) B
- #21) A
- #22) D
- #23) C
- #24) C
- #25) D
- #26) E
- #27) B
- #28) E
- #29) C
- #30) D
- #31) D
- #32) B
- #33) E
- #34) C
- #35) E
- #36) C
- #37) E
- #38) B
- #39) A
- #40) A

Question 1

Intent of Question

The primary goals of this question were to assess a student's ability to (1) describe features of a distribution of sample data using information provided by a histogram; (2) identify potential outliers; (3) sketch a boxplot; and (4) comment on an advantage of displaying data as a histogram rather than as a boxplot.

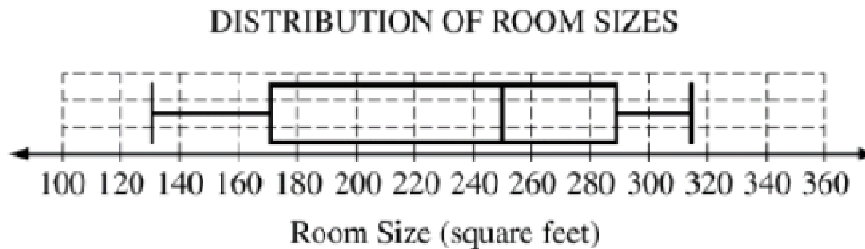
Solution

Part (a):

The distribution of the sample of room sizes is bimodal and roughly symmetric with most room sizes falling into two clusters: 100 to 200 square feet and 250 to 350 square feet. The center of the distribution is between 200 and 300 square feet. The range of the distribution is between 150 and 250 square feet. There are no apparent outliers.

Part (b):

The interquartile range is $IQR = 292 - 174 = 118$ square feet. There are no potential outliers because the minimum room size of 134 square feet does not fall below $Q_1 - 1.5(IQR) = -3$ square feet, and the maximum room size of 315 square feet does not exceed $Q_3 + 1.5(IQR) = 469$ square feet.



Part (c):

The histogram clearly shows the bimodal nature of the distribution of room sizes, but this is not apparent in the boxplot.

Question 1 (continued)

Scoring

This question is scored in three sections. Section 1 consists of part (a); Section 2 consists of the outlier determination in part (b); Section 3 consists of the boxplot sketch in part (b) and part (c). Each section is scored as essentially correct (E), partially correct (P), or incorrect (I).

Section 1 is scored as follows:

Essentially correct (E) if the description of the distribution of room sizes satisfies the following four components:

1. The shape is bimodal OR there are two peaks OR there are two clusters.
2. The center is between 200 and 300 square feet.
3. The spread is addressed by stating the range is a value between 150 and 250 square feet OR the interquartile range is a value between 50 and 150 square feet OR all room sizes are between 100 and 350 square feet.
4. The response includes context.

Partially correct (P) if the response satisfies two or three of the four components.

Incorrect (I) if the response does not satisfy the criteria for E or P.

Notes:

- Shape: Component 1 cannot be satisfied if a response describes the histogram as unimodal or describes the entire histogram as normal or approximately normal.
- Shape: A response that addresses symmetry, while appropriate, does not impact the scoring of section 1.
- Center: A response that states one cluster of the distribution is centered between 150 and 200 square feet and the other cluster is centered between 250 and 300 square feet satisfies both components 1 and 2.
- Center:
 - Responses that address center using interval language such as “the mean of the distribution is *between* 200 and 300” must, for any single measure of center, provide an interval with lower endpoint not below 200 square feet, and with upper endpoint not above 300 square feet to satisfy component 2.
 - Responses that address center using approximate language such as “the median of the distribution is *approximately* 225” must, for any single measure of center, specify a numeric value that is not less than 200 square feet, and that is not greater than 300 square feet to satisfy component 2.
 - Responses that use definitive language such as “the mean of the distribution *is* 231.4” must identify the corresponding numeric value correctly to satisfy component 2. Specifically, the median of the distribution can be correctly identified as any value between 250 and 253.5 square feet, inclusive; the mean of the distribution is 231.4 square feet; and the center (or average) of the distribution can be any value that is a correct median or mean.

Question 1 (continued)

- Spread: A response recognizing all values in the sample fall between 100 and 350 square feet (or between 134 and 315 square feet) satisfies component 3 *only* for these exact endpoints and need not appeal to a specific measure of spread such as range or interquartile range (IQR).
- Spread:
 - Responses that appeal to a specific measure of spread using interval language, such as “the IQR is *between* 50 and 150,” must provide bounds appropriate to the corresponding measure of spread. For range, the lower endpoint must not be below 150 square feet and the upper endpoint cannot exceed 250 square feet; for IQR, the lower endpoint must not be below 50 square feet, with upper endpoint not to exceed 150 square feet; for standard deviation, the lower endpoint must not be below 25 square feet, with upper endpoint not to exceed 100 square feet.
 - Responses that appeal to a specific measure of spread using approximate language, such as “the range is *approximately* 250,” must specify a numeric value within the bounds appropriate to that measure of spread. For range, the value must be between 150 and 250 square feet (inclusive); for IQR, the value must be between 50 and 150 square feet (inclusive); for standard deviation, the value must be between 25 and 100 square feet (inclusive). Responses that appeal to a specific measure of spread using definitive language, such as “the range of the distribution *is* 181,” must identify the corresponding numeric value correctly to satisfy component 3. Specifically, the range of the distribution is 181 square feet; the IQR of the distribution is 118 square feet; and the standard deviation of the distribution is 68.12 square feet.

Section 2 is scored as follows:

Essentially correct (E) if the response satisfies the following three components:

1. Computation of both upper and lower outlier boundary fences that also shows the fences formulas either in words, symbols $Q_1 - 1.5(IQR)$ and $Q_3 + 1.5(IQR)$, or with values substituted from the table $174 - 1.5(118)$ and $292 + 1.5(118)$, or $(174 - 177)$ and $(292 + 177)$.
2. A correct decision regarding the presence of outliers.
3. Correct justification that compares the data with the fences.

Partially correct (P) if the response satisfies only two of the three components OR if the response omits exactly one of the fences but otherwise satisfies all three components.

Incorrect (I) if the response does not satisfy the requirements for E or P.

Notes:

- A response that identifies both fence formulas using symbols, but does not substitute values for all symbols, must also include the correct fence values of -3 and 469 to satisfy component 1.
- In place of an appeal to fences, a response may compute outlier bounds representing k standard deviations from the sample mean, where k is a number from 2 to 3 (inclusive), and must include formulas for both endpoints either in words, symbols $\bar{x} \pm k(\text{standard deviation})$, or with values substituted from the table. When $k = 2$ the outlier bounds are $(95.16, 367.64)$; when $k = 3$ the bounds are $(27.04, 435.76)$.

Question 1 (continued)

- A response that identifies the standard deviation bounds using symbols, but that does not substitute values for all symbols, does not satisfy component 1 unless the correct numeric bounds are provided.
- Component 3 is satisfied if the response states the outlier decision criterion: any data values falling outside of the interval from -3 to 469 are potential outliers.

Section 3 is scored as follows:

Essentially correct (E) if the response satisfies the following two components:

1. A correct sketch of the boxplot.
2. A response for part (c) that indicates the bimodal shape of the room size distribution is apparent in the histogram but not in the boxplot.

Partially correct (P) if the response satisfies only one of the two components.

Incorrect if the response does not meet the criteria for E or P.

Notes:

- The boxplot must be completely correct to satisfy component 1. Specifically:
 - The minimum is positioned between grid lines at 120 and 140 square feet.
 - Q_1 is positioned between grid lines at 160 and 180 square feet.
 - The median is positioned between grid lines at 240 and 260 square feet.
 - Q_3 is positioned between grid lines at 280 and 300 square feet.
 - The maximum is positioned between grid lines at 300 and 320 square feet.
- If a *mean* is included as a part of the boxplot, component 1 cannot be satisfied.
- A response based on skewness or symmetry does not satisfy component 2.
- A response stating the unimodal OR normal shape of the histogram of room sizes is apparent in the histogram but not in the boxplot will satisfy component 2 *only* if the shape description in section 1 component 1 was also unimodal OR normal, respectively.

4 Complete Response

Three sections essentially correct

3 Substantial Response

Two sections essentially correct and one section partially correct

2 Developing Response

Two sections essentially correct and no sections partially correct

OR

One section essentially correct and one or two sections partially correct

OR

Three sections partially correct

1 Minimal Response

One section essentially correct

OR

No sections essentially correct and two sections partially correct

Question 2

Intent of Question

The three primary goals of this question are to assess a student's ability to: (1) clearly explain the importance of a control group in the context of an experiment; (2) describe the randomization process required for three groups; and (3) reduce variability by grouping experimental units as homogeneously as possible.

Solution

Part (a):

A control group gives the researchers a comparison group to be used to evaluate the effectiveness of the treatments. The control group allows the impact of the normal aging process on joint and hip health to be measured with appropriate response variables. The effects of glucosamine and chondroitin can be assessed by comparing the responses for these two treatment groups with those for the control group.

Part (b):

Each dog will be assigned a unique random number, 001–300, using a random number generator on a calculator, statistical software, or a random number table. The numbers will be sorted from smallest to largest. The dogs assigned the first 100 numbers in the ordered list will receive glucosamine. The dogs with the next 100 numbers in the ordered list will be assigned to the control group. Finally, the dogs with the numbers 201–300 will receive chondroitin.

Part (c):

The key question is which variable has the strongest association with joint and hip health. The goal of blocking is to create groups of homogeneous experimental units. It is reasonable to assume that most clinics will see all kinds and breeds of dogs so there is no reason to suspect that joint and hip health will be strongly associated with a clinic. On the other hand, different breeds of dogs tend to come in different sizes. The size of a dog is associated with joint and hip health, so it would be better to form homogeneous groups of dogs by blocking on breed.

Scoring

Parts (a), (b), and (c) are scored as essentially correct (E), partially correct (P), or incorrect (I).

Part (a) is scored as essentially correct (E) if an advantage of using a comparison group is described in the context of this study.

Part (a) is scored as partially correct (P) if an advantage of using a control group is described but not in the context of this study.

Part (a) is scored as incorrect (I) if the student says that control groups should always be used but gives no further explanation *OR* an incorrect explanation.

Note: Since “treatment” and “control” are standard terms in design, a comparison of specific aspects of the study is needed to establish context.

Question 2 (continued)

Part (b) is scored as essentially correct (E) if randomization is used correctly, and the method of randomization can be implemented after reading the student response (so that two knowledgeable statistics users would use the same method to assign dogs to treatment groups).

Part (b) is scored as partially correct (P) if randomization or chance is used, but the method could not be implemented after reading the student response.

Part (b) is scored as incorrect (I) if randomization or chance is not used in a planned way *OR* the solution does not yield a completely randomized design.

Part (c) is scored as essentially correct (E) if:

the student argues that the variable with the stronger relationship to joint and hip health should be used as the blocking variable;

OR

the student states that the variable with the larger anticipated variability in the response measure should be used as the blocking variable so that units within blocks are as homogeneous as possible. A rationale is required, but a variable does not have to be selected.

Part (c) is scored as partially correct (P) if:

the student indicates that the purpose of blocking is to create groups of homogeneous experimental units but makes an error in the application to this experiment;

OR

the student does not acknowledge that there may be more variability associated in the response variable with one of the variables (breed or clinic) than the other;

OR

the student does not recognize that both variables are associated with variation in the response variable.

Part (c) is scored as incorrect (I) if the student does not exhibit an understanding of the purpose of blocking.

4 Complete Response

All three parts essentially correct

3 Substantial Response

Two parts essentially correct and one part partially correct

2 Developing Response

Two parts essentially correct and no parts partially correct

OR

One part essentially correct and two parts partially correct

OR

Three parts partially correct

1 Minimal Response

One part essentially correct and either zero or one part partially correct

OR

No parts essentially correct and two parts partially correct

Question 3

Solution

Part (a):

Yes, the linear model is appropriate for these data. The scatterplot shows a strong, positive, linear association between the number of railcars and fuel consumption, and the residual plot shows a reasonably random scatter of points above and below zero.

Part (b):

According to the regression output, fuel consumption will increase by 2.15 units for each additional railcar. Since the fuel consumption cost is \$25 per unit, the average cost of fuel per mile will increase by approximately $\$25 \times 2.15 = \53.75 for each railcar that is added to the train.

Part (c):

The regression output indicates that $r^2 = 96.7\%$ or 0.967. Thus, 96.7% of the variation in the fuel consumption values is explained by using the linear regression model with number of railcars as the explanatory variable.

Part (d):

No, the data set does not contain any information about fuel consumption for any trains with more than 50 cars. Using the regression model to predict the fuel consumption for a train with 65 railcars, known as extrapolation, is not reasonable.

Scoring

Each part is scored as essentially correct (E), partially correct (P), or incorrect (I).

Part (a) is essentially correct (E) if the model is deemed appropriate AND the explanation clearly indicates:

- There is a linear pattern in the scatterplot; OR
- There is no pattern in the residual plot.

Part (a) is partially correct (P) if the:

- Model is deemed appropriate AND the student refers to the scatterplot or residual plot but fails to state the relevant characteristic of the plot; OR
- Student refers to the relevant characteristic of the scatterplot or residual plot without deeming model appropriate.

Part (a) is incorrect (I) if the student:

- States that the model is appropriate without an explanation; OR
- States that the model is inappropriate; OR
- Makes a decision based only on numeric values from the computer output.

Question 3 (continued)

Part (b) is essentially correct (E) if the point estimate for the slope (2.15 or 2.1495) and the fuel consumption cost per unit (\$25) are used to calculate the correct point estimate (\$53.75 or $\$53.7375 \approx \53.74).

Part (b) is partially correct (P) if only the point estimate for the slope (2.15 or 2.1495) is stated with a supporting calculation or interpretation.

Part (c) is essentially correct (E) if the student states:

- 96.7% of the variation in fuel consumption is explained by the linear regression model; OR
- 96.7% of the variation in fuel consumption is explained by the number of railcars.

Part (c) is partially correct (P) if the student makes one of the above statements using $R\text{-Sq}(\text{adj}) = 96.3\%$.

Part (d) is essentially correct (E) if the student states that this is unreasonable due to extrapolation.

Part (d) is partially correct (P) if the student states this is:

- Unreasonable but provides a weak explanation; OR
- Reasonable even though it is considered a slight extrapolation.

Note: Any answer appearing without supporting work is scored as incorrect (I).

Each essentially correct (E) response counts as 1 point, each partially correct (P) response counts as $\frac{1}{2}$ point.

- 4 Complete Response
- 3 Substantial Response
- 2 Developing Response
- 1 Minimal Response

Note: If a response is in between two scores (for example, $2\frac{1}{2}$ points), use a holistic approach to determine whether to score up or down depending on the strength of the response and communication.

Question 4

Intent of Question

The primary goals of this question were to assess a student's ability to (1) calculate a probability using basic probability rules or the geometric distribution; (2) recognize that a probability calculation for independent events does not depend on the previous outcomes of those events; and (3) assess whether a claim about the probability of a single event is reasonable based on a calculated probability of a series of those events.

Solution

Part (a):

If the failure rate for the super igniters is 15 percent, then the probability that each igniter fails is 0.15, and the probability that it does not fail is 0.85. Therefore the probability that the first 30 igniters tested do not fail is $(0.85)^{30} \approx 0.0076$. The solution can also be written as $(1 - 0.15)^{30} \approx 0.0076$.

Part (b):

Given that there are no failures in the first 30 trials, the probability that the first failure occurs on the 31st trial is 0.15, and the probability that it does not occur on the 31st but occurs on the 32nd trial is $(0.85)(0.15) = 0.1275$. Therefore the probability that the first failure occurs on the 31st or 32nd super igniter tested is $0.15 + 0.1275 = 0.2775$.

Note that this is equivalent to asking for the probability that the first failure occurs on the first or second trial, which is $0.15 + (0.85)(0.15) = 0.2775$.

Part (c):

The result of the probability calculation in part (a) provides a reason to believe that the failure rate of the super igniters is less than 15 percent. The calculated probability of 0.0076 shows that there is less than a 1 percent chance that 30 or more igniters in a row would not fail if the failure rate was 15 percent. This probability is smaller than conventional significance levels such as $\alpha = 0.05$ or $\alpha = 0.01$, and thus is small enough to make it reasonable to believe that the failure rate of the super igniters is less than 15 percent.

Question 4 (continued)

Scoring

Parts (a), (b), and (c) are scored as essentially correct (E), partially correct (P), or incorrect (I).

Part (a) is scored as follows:

Essentially correct (E) if the response gives the correct probability *AND* correct justification.

Partially correct (P) if the response correctly notes that the answer is the probability that there will be 30 successes in 30 attempts, but does not carry out a correct probability calculation;

OR

if the response defines the random variable X as the trial with the first failure, identifies X as having a geometric distribution with $p = 0.15$, and writes the desired probability as $P(X > 30)$, but does not carry out a correct probability calculation;

OR

if the response defines the random variable X as the number of failures in the first 30 attempts, identifies X as a binomial random variable with $p = 0.15$ and $n = 30$, and writes the desired probability as $P(X = 0)$, but does not carry out a correct probability calculation;

OR

if the response gives the correct probability but, in specifying a geometric or binomial distribution, has an incorrect or incomplete definition of parameters or value(s) of the random variable.

Incorrect (I) if the response does not meet the criteria for E or P.

Note: Justification can be given using the multiplication rule; *OR* by defining X to be the trial with the first failure, recognizing that X has a geometric distribution, and using that information to find $P(X > 30)$; *OR* by defining X to be the number of failures in the first 30 attempts, and then finding $P(X = 0)$ using either probability rules or the binomial distribution with $n = 30$ and $p = 0.15$.

Part (b) is scored as follows:

Essentially correct (E) if the response gives the correct probability *AND* correct justification.

Partially correct (P) if the response makes a reasonable attempt to calculate a geometric, binomial, or conditional probability, but does not successfully carry out the calculation;

OR

if the response gives the correct probability but, in specifying a geometric or binomial distribution, has an incorrect or incomplete definition of parameters or value(s) of the random variable.

Incorrect (I) if the response finds an incorrect probability resulting from an unreasonable attempt to calculate a geometric, binomial, or conditional probability or otherwise does not meet the criteria for E or P.

Note: Similar to part (a) justification can be given using probability rules; *OR* by stating that X is geometric where X is the trial with the first failure, then finding $P(X = 1 \text{ or } X = 2)$; *OR* by stating that X is the number of failures in two trials and finding $1 - P(X = 0)$ or $P(X = 1 \text{ or } X = 2)$ using the binomial distribution.

Question 4 (continued)

Part (c) is scored as follows:

Essentially correct (E) if the response states that it is reasonable to believe that the failure rate is less than 15 percent *AND* bases this decision on the fact that the probability of 30 consecutive successful launches with a failure rate of 15 percent (that is, answer from part (a)) is small *AND* does so in the context of the situation.

Partially correct (P) if the response otherwise satisfies the criteria for an (E) but does so without any context;

OR

if the response states a significance level and makes a decision in a context that is appropriate to the given probability in part (a) and the stated significance level but does not explicitly compare the probability and the significance level (no linkage).

Incorrect (I) if the response does not explicitly make a decision about whether it is reasonable to conclude that the failure rate is less than 15 percent (For example: "As seen in Part (a), if the failure rate is 15 percent then the probability of 30 successful launches in a row is very small.");

OR

if the response otherwise does not meet the criteria for E or P.

Notes:

- Justification based on the probability can come by stating a significance level and noting that the probability is smaller than the significance level *OR* by simply stating that the probability of 0.0076 is small *OR* by referring to the expected number of failures (4.5) as being very unlikely because zero failures is more than two standard deviations below 4.5.
- If the response bases the decision on the expected number of failures (4.5) for $n = 30$ and $p = 0.15$ without referencing why zero failures would be considered to be too far below 4.5 to give reason to doubt the stated 15 percent failure rate, the response is scored P.
- If the calculation in part (a) is incorrect, the answer in part (c) needs to be consistent with the answer in part (a), unless the value is recalculated in part (c).

4 Complete Response

Three parts essentially correct

3 Substantial Response

Two parts essentially correct and one part partially correct

2 Developing Response

Two parts essentially correct and no parts partially correct

OR

One part essentially correct and one or two parts partially correct

OR

Three parts partially correct

1 Minimal Response

One part essentially correct

OR

No parts essentially correct and two parts partially correct

Question 5

Intent of Question

The primary goals of this question were to assess a student's ability to (1) determine whether a cause-and-effect conclusion can be made based on how a study was conducted and (2) set up, perform, and interpret the results of a hypothesis test, in the context of the problem.

Solution

Part (a):

Yes, it would be reasonable to conclude that the new procedure causes a reduction in recovery time, for patients similar to those in the study. The patients in the study were randomly assigned to the two procedures, which reduces the chance that confounding variables will affect the results. Therefore the statistically significant reduction in mean recovery time can be attributed to the new procedure being superior to the standard procedure.

Part (b):

Step 1: State a correct pair of hypotheses.

Let μ_S represent the mean recovery time among all patients similar to those in the study if they were to receive the standard treatment.

Let μ_N represent the mean recovery time among all patients similar to those in the study if they were to receive the new treatment.

The hypotheses to be tested are $H_0 : \mu_S = \mu_N$ versus $H_a : \mu_S > \mu_N$.

Step 2: Identify a correct test procedure (by name or by formula) and check appropriate conditions.

The appropriate procedure is a two-sample t -test for a difference between means.

Because this is an experiment, the first condition is that subjects were randomly assigned to one treatment group or the other. In this case the condition is satisfied because we were told that the subjects were randomly assigned to either the standard or new procedure.

The second condition is that the recovery times of the two populations are normally distributed or the sample sizes are sufficiently large to presume that the distribution of the difference in the sample means is approximately normal. In this case the condition is met because the sample sizes of 110 and 100 are both sufficiently large.

Step 3: Correct mechanics, including the value of the test statistic, degrees of freedom, and p -value (or rejection region).

$$\text{The test statistic is } t = \frac{\bar{x}_S - \bar{x}_N}{\sqrt{\frac{s_S^2}{n_S} + \frac{s_N^2}{n_N}}} = \frac{217 - 186}{\sqrt{\frac{34^2}{110} + \frac{29^2}{100}}} \approx 7.13.$$

The p -value is the area greater than 7.13 for a t -distribution with $df = 207.18$, which is essentially 0 (8.36×10^{-12}).

Step 4: State a correct conclusion in the context of the problem, using the result of the statistical test.

Because the p -value is very small, we have sufficient evidence to conclude that for patients similar to the ones in the study, those receiving the new procedure would have less recovery time, on average, than those receiving the standard procedure.

Question 5 (continued)

Scoring

This question is scored in three sections. Section 1 consists of part (a); section 2 consists of step 1, step 2, and the test statistic in step 3 in part (b); and section 3 consists of the p -value in step 3 and step 4 in part (b). Sections 1, 2, and 3 are each scored essentially correct (E), partially correct (P), or incorrect (I).

Section 1 is scored as follows:

Essentially correct (E) if the response satisfies the following three components:

1. Correctly states that it is reasonable to make a causal conclusion.
2. Justifies the causal conclusion based on random assignment of patients to procedures (or procedures to patients);

OR

justifies the causal conclusion by stating that a randomized experiment was conducted.

3. Includes the context of the situation.

Partially correct (P) if the response satisfies component 1 *AND* provides WEAK justification of the causal conclusion by stating that there was random assignment or a randomized experiment was conducted, but with no context;

OR

by stating that an experiment was conducted or there was assignment (without the word “randomized”) *AND* the response includes context of the situation;

OR

by stating that the study design reduces the chance of confounding variables or balances the effects of uncontrolled variables across both groups in context without explicitly referring to the random assignment.

Incorrect (I) if the response does not meet the criteria for E or P.

Notes:

- If the response states that it is *not* reasonable to make a causal conclusion because the result could have been due to random chance *AND* explains that there is evidence for a causal conclusion based on random assignment of patients to procedures or by stating that a randomized experiment was conducted, then the response is scored E.
- If the response discusses aspects of an experiment other than random assignment (such as, control, replication, or large samples), then those aspects are considered extraneous and the response can be scored E unless those aspects are incorrect for this study (such as, blocking is a requirement, or the study used blocking, or the study used a placebo) in which case the score should be lowered one level (that is, from E to P, or from P to I).
- If the response correctly states in context that it is reasonable to make a causal conclusion but includes incorrect or contradictory justification (such as, random selection of patients), then the response is scored I.

Question 5 (continued)

Section 2 is scored as follows:

Essentially correct (E) if the response satisfies the following four components:

1. Parameters are defined correctly.
2. Hypotheses imply equality in the null and correct direction in the alternative.
3. Correct test is identified by name or formula.
4. Correct test statistic for a difference in means is calculated.

Partially correct (P) if the response satisfies only two or three of the four components.

Incorrect (I) if the response satisfies at most one of the four components.

Notes:

- If standard symbols are used for the parameters with appropriate group labels (such as, μ_S, μ_N), component 1 is satisfied.
- If the correct test is identified, but the response states an incorrect formula or uses incorrect notation in the formula, component 3 is not satisfied.
- A pooled two-sample t -test is acceptable for component 3, but the student must also state and comment on the plausibility of the equal population variances assumption.
- If the response identifies a z -test for equal means as the correct test identification, component 3 is not satisfied but component 4 could be satisfied.

Confidence Interval approach:

- If a single two-sample t -interval for the difference in means is used, components 3 and 4 can be satisfied. Component 3 is satisfied if the t -interval is correctly identified by name or formula. Component 4 is satisfied if the correct interval is calculated. If an alpha level is stated, then an appropriate adjustment to the confidence level must be made because the appropriate test is one-sided.
- If two one-sample t -intervals are used, while not a recommended approach, component 3 is not satisfied but component 4 could be satisfied. Component 4 is satisfied if both intervals are calculated correctly.

Question 5 (continued)

Section 3 is scored as follows:

Essentially correct (E) if the response satisfies the following three components:

1. Makes reference to an approximately correct p -value that is consistent with the test statistic and alternative hypothesis for a difference in means.
2. Correctly justifies the conclusion based on the size of the p -value or the test statistic.
3. Correctly states the conclusion in context.

Partially correct (P) if the response satisfies only two of the three components.

Incorrect (I) if the response does not meet the criteria for E or P or includes a justification not based on the inferential results.

Notes:

Component 1:

- Is satisfied if the response makes reference to a large test statistic without referring to a p -value.

Component 2:

- No alpha level is needed to provide justification of the conclusion based on the size of the p -value.
- Is satisfied if the response states the p -value without reference to size, but it is contiguous to the conclusion and clearly indicates a continuous train of thought.
- A correct interpretation of the p -value with a complete explanation that obtaining a test statistic at least this extreme is unlikely due to chance alone is considered justification based on the size of the p -value.
- If an incorrect interpretation of the p -value is given, the score is lowered one level (that is, from E to P, or from P to I).
- A decision about the null hypothesis (reject H_0 or fail to reject H_0) is not required, but if an incorrect decision is stated based on the given p -value then component 2 is not satisfied.
- If a rejection region approach is used, a reasonable critical value replaces the p -value.

Component 3:

- A correct conclusion must be related to the alternative hypothesis in order to satisfy component 3.
- The following responses do not satisfy component 3:
 - States or implies that the null hypothesis is *accepted*
 - States or implies that the alternative hypothesis has been *proven*
 - States the conclusion in past tense (unless the response did not satisfy a component of section 2 for the use of past tense)

Question 5 (continued)

Confidence Interval approach:

- If a single two-sample t -interval for the difference in means is used:
 - Component 1 is satisfied if the response indicates that zero is either included or not included in the calculated interval.
 - Component 2 is satisfied if the response indicates that the bounds are either both above or both below zero (consistent with alternative hypothesis) and uses that as justification for the conclusion.
 - Component 3 is satisfied if the conclusion is stated in context.
- If two one-sample t -intervals are used (which is not recommended) the response is scored at most P if all three components are satisfied, otherwise scored I:
 - Component 1 is satisfied if the response states that the intervals do not overlap.
 - Component 2 is satisfied if the conclusion indicates that the confidence interval for the new procedure lies below the confidence interval for the standard procedure.
 - Component 3 is satisfied if the conclusion is stated in context.

Note: If the three sections of the response are scored as E, to earn a score of 4 as a complete response, both conditions in step 2 must be correctly stated and justified. Additional condition(s) inappropriate for a two-sample t -test must not be stated. Otherwise, the response earns a score of 3 a substantial response.

4 Complete Response

Three sections essentially correct with conditions for inference

3 Substantial Response

Three sections essentially correct without conditions for inference

OR

Two sections essentially correct and one section partially correct

2 Developing Response

Two sections essentially correct and no sections partially correct

OR

One section essentially correct and one or two sections partially correct

OR

Three sections partially correct

1 Minimal Response

One section essentially correct

OR

No sections essentially correct and one or two sections partially correct

Question 6

Intent of Question

The primary goals of this question were to assess a student's ability to (1) recognize the population to which results from a random sample may be generalized; (2) describe a disadvantage of using a sample mean rather than a sample median to indicate typical values when the sample distribution is skewed; (3) describe how the theoretical sampling distribution of the sample median could be constructed; (4) construct an approximate confidence interval for a population median using results from a bootstrap procedure; and (5) interpret a confidence interval.

Solution

Part (a):

Because random sampling was used, the results of the sample may be generalized to the population of rental prices for one-bedroom apartments in the city that are listed on this particular website at the time the sample was taken.

Part (b):

Because the distribution of the 50 rental prices in the sample is skewed to the right, the sample median provides a better indicator of typical rental prices than the sample mean. Some very large rental prices results in a sample mean that is substantially larger than the more typical rental prices. As a result the sample mean would overestimate the typical rental price, whereas the sample median would be a more accurate representation of typical rental prices.

Part (c):

To determine the sampling distribution of median rental prices for random samples of 50 one-bedroom apartments from this population, Emma would need to obtain every possible sample of 50 one-bedroom apartments from this website and compute the median of each sample. The collection of all possible sample medians is the theoretical sampling distribution for sample median.

Part (d):

- (i) $(0.05)(15,000) = 750$ and $(0.95)(15,000) = 14,250$. The 5th percentile is a value, say $x_{0.05}$, such that at least 750 values in the table are less than or equal to $x_{0.05}$ and at least 14,250 are greater than or equal to $x_{0.05}$. Cumulate frequencies starting with the smallest sample median listed in the table and going toward the largest (going down columns) until you first reach 750 values, to obtain $x_{0.05} = \$2,500$.
- (ii) Similarly, $x_{0.95} = \$2,950$.

Part (e):

The percentage of bootstrap medians between (and including) the values found in part (d) for the 5th and 95th percentiles is

$$\frac{14,404}{15,000} \times 100\% \approx 96.03\%$$

Part (f):

From the results in part (d) and part (e), an approximate 96 percent confidence interval for the median rental price of all one-bedroom apartments listed on this website for this city is $(\$2,500, \$2,950)$. We are approximately 96 percent confident that the median rental price of all one-bedroom apartments listed on this website for this city is between \$2,500 and \$2,950.

Question 6 (continued)

Scoring

This question is scored in four sections. Section 1 consists of parts (a) and (b), section 2 consists of part (c), section 3 consists of parts (d) and (e), and section 4 consists of part (f). Sections 1, 2, 3, and 4 are each scored as essentially correct (E), partially correct (P), or incorrect (I).

Section 1 is scored as follows:

Essentially correct (E) if the response satisfies the following three components:

1. The correct population (listings of one-bedroom apartments on the website) is identified in part (a).
2. In part (b), identifying that using the sample mean instead of the sample median overestimates the typical rental price.
3. The disadvantage of using the sample mean that is reported in part (b) is correctly linked to some feature of the distribution (e.g. skewness) that is evident in the histogram.

Partially correct (P) if the response satisfies only two of the three components.

Incorrect (I) if the response does not meet the criteria for E or P.

Note: Responses that refer to the mean being larger than the median in a skewed right distribution alone is not sufficient to satisfy component 2.

Section 2 is scored as follows:

Essentially correct (E) if the response satisfies the following two components:

1. Indicates that Emma would need to obtain every possible sample of 50 one-bedroom apartments.
2. Indicates that Emma would need to compute the median rental price for each sample.

Partially correct (P) if the response satisfies only one of the two components.

Incorrect (I) if the response does not satisfy the criteria for E or P.

Section 3 is scored as follows:

Essentially correct (E) if the response satisfies the following two components:

1. Correct values for the 5th percentile and the 95th percentile are reported in part (d).
2. The correct percentage of bootstrap samples that produced sample medians at or between the two values, if they are plausible, reported in part (d) is reported in part (e).

Partially correct (P) if the response satisfies only one of the two components.

Incorrect (I) if the response does not satisfy the criteria for E or P.

Note: Plausible values for part (d) will be considered values between 2,345 and 3,062.5.

Question 6 (continued)

Section 4 is scored as follows:

Essentially correct if the response in part (f) satisfies the following three components:

1. Uses \$2500 and \$2950 or the values of the percentiles reported in part (d) as the endpoints of the confidence interval.
2. Indicates an approximate 90 or 96 percent level of confidence or a level of confidence consistent with part (e).
3. Makes a correct statement in context indicating that the confidence interval is for the median.

Partially correct if the response satisfies only two of the three components.

Incorrect if the response does not satisfy the criteria for E or P.

Note: Since rental prices from the population are discrete values, the true confidence level of the interval from part (d) is unknown. A correctly calculated part (e) is a way to estimate the confidence level; from Emma's sample the confidence level is estimated to be approximately 96 percent. The process described in part (d) for calculating the interval will result in a confidence level of at least 90 percent. For these reasons, confidence levels of either 90 or 96 percent satisfy component 2.

Each essentially correct (E) section counts as 1 point, and a partially correct (P) section counts as $\frac{1}{2}$ point.

- | | |
|----------|-----------------------------|
| 4 | Complete Response |
| 3 | Substantial Response |
| 2 | Developing Response |
| 1 | Minimal Response |

If a response is between two scores (for example, $2\frac{1}{2}$ points), use a holistic approach to decide whether to score up or down depending on the strength of the response and communication.