## Question 6
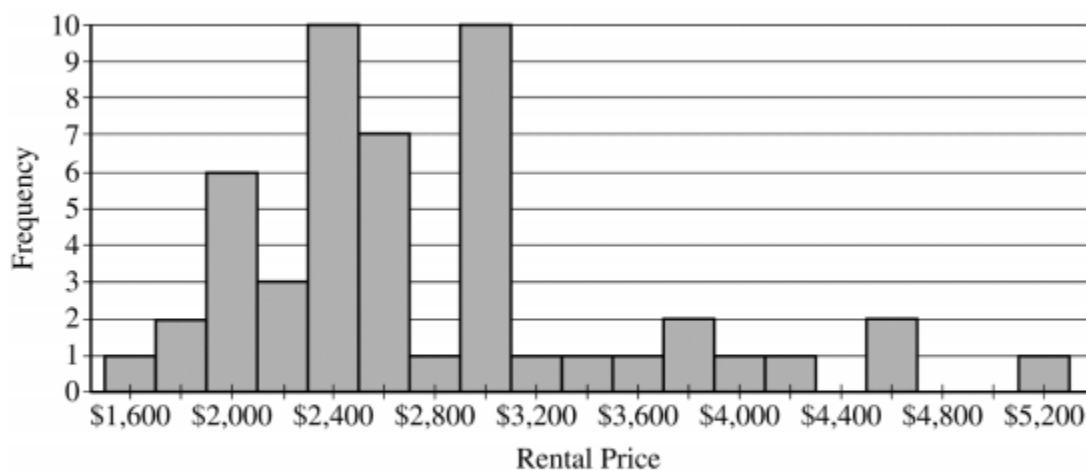### Spend about 25 minutes on this part of the exam.
### Percent of Section II score—25

**Directions:** Show all your work. Indicate clearly the methods you use, because you will be scored on the correctness of your methods as well as on the accuracy and completeness of your results and explanations.

6. Emma is moving to a large city and is investigating typical monthly rental prices of available one-bedroom apartments. She obtained a random sample of rental prices for 50 one-bedroom apartments taken from a Web site where people voluntarily list available apartments.

   (a) Describe the population for which it is appropriate for Emma to generalize the results from her sample.

The distribution of the 50 rental prices of the available apartments is shown in the following histogram.



   (b) Emma wants to estimate the typical rental price of a one-bedroom apartment in the city. Based on the distribution shown, what is a disadvantage of using the mean rather than the median as an estimate of the typical rental price?

   (c) Instead of using the sample median as the point estimate for the population median, Emma wants to use an interval estimate. However, computing an interval estimate requires knowing the sampling distribution of the sample median for samples of size 50. Emma has one point, her sample median, in that sampling distribution. Using information about rental prices that are available on the Web site, describe how someone could develop a theoretical sampling distribution of the sample median for samples of size 50.

Because Emma does not have the resources to develop the theoretical sampling distribution, she estimates the sampling distribution of the sample median using a process called bootstrapping. In the bootstrapping process, a computer program performs the following steps.

- Take a random sample, with replacement, of size 50 from the original sample.
- Calculate and record the median of the sample.
- Repeat the process to obtain a total of 15,000 medians.

Emma ran the bootstrap process, and the following frequency table is the bootstrap distribution showing her results of generating 15,000 medians.

| Bootstrap Distribution of Medians | | | | | |
|---|---|---|---|---|---|
| Median | Frequency | Median | Frequency | Median | Frequency |
| 2,345 | 1 | 2,585 | 1 | 2,825 | 247 |
| 2,390 | 13 | 2,587.5 | 171 | 2,837.5 | 7 |
| 2,395 | 18 | 2,600 | 22 | 2,847.5 | 1 |
| 2,400 | 56 | 2,612.5 | 1,190 | 2,872.5 | 317 |
| 2,445 | 4 | 2,625 | 174 | 2,885 | 10 |
| 2,447.5 | 56 | 2,672.5 | 5 | 2,950 | 700 |
| 2,450 | 55 | 2,675 | 1,924 | 2,962.5 | 93 |
| 2,475 | 3 | 2,687.5 | 1,341 | 2,972.5 | 6 |
| 2,495 | 66 | 2,700 | 2,825 | 2,975 | 65 |
| 2,497.5 | 136 | 2,735 | 35 | 2,985 | 12 |
| 2,500 | 1,899 | 2,747.5 | 619 | 2,987.5 | 1 |
| 2,522.5 | 2 | 2,750 | 2 | 2,995 | 6 |
| 2,525 | 945 | 2,795 | 278 | 3,000 | 2 |
| 2,550 | 1,673 | 2,812.5 | 16 | 3,062.5 | 3 |

The bootstrap distribution provides an approximation of the sampling distribution of the sample median. A confidence interval for the median can be constructed using a percentage of the values in the middle of the bootstrap distribution.

(d) Use the frequency table to find the following.

(i) Value of the 5th percentile:

(ii) Value of the 95th percentile:

(e) Find the percentage of bootstrap medians in the table that are equal to or between the values found in part (d).

(f) Use your values from parts (d) and (e) to construct and interpret a confidence interval for the median rental price.

# Question 6

**Spend about 25 minutes on this part of the exam.**
**Percent of Section II score—25**

**Directions:** Show all your work. Indicate clearly the methods you use, because you will be scored on the correctness of your methods as well as on the accuracy and completeness of your results and explanations.

6. Systolic blood pressure is the amount of pressure that blood exerts on blood vessels while the heart is beating. The mean systolic blood pressure for people in the United States is reported to be 122 millimeters of mercury (mmHg) with a standard deviation of 15 mmHg.

   The wellness department of a large corporation is investigating whether the mean systolic blood pressure of its employees is greater than the reported national mean. A random sample of 100 employees will be selected, the systolic blood pressure of each employee in the sample will be measured, and the sample mean will be calculated.

   Let $\mu$ represent the mean systolic blood pressure of all employees at the corporation. Consider the following hypotheses.

   $$H_0 : \mu = 122$$
   $$H_a : \mu > 122$$

   (a) Describe a Type II error in the context of the hypothesis test.

   (b) Assume that $\sigma$, the standard deviation of the systolic blood pressure of all employees at the corporation, is 15 mmHg. If $\mu = 122$, the sampling distribution of $\bar{x}$ for samples of size 100 is approximately normal with a mean of 122 mmHg and a standard deviation of 1.5 mmHg. What values of the sample mean $\bar{x}$ would represent sufficient evidence to reject the null hypothesis at the significance level of $\alpha = 0.05$ ?

The actual mean systolic blood pressure of all employees at the corporation is 125 mmHg, not the hypothesized value of 122 mmHg, and the standard deviation is 15 mmHg.

(c) Using the actual mean of 125 mmHg and the results from part (b), determine the probability that the null hypothesis will be rejected.

(d) What statistical term is used for the probability found in part (c) ?

(e) Suppose the size of the sample of employees to be selected is greater than 100. Would the probability of rejecting the null hypothesis be greater than, less than, or equal to the probability calculated in part (c) ? Explain your reasoning.

# Question 6
**Spend about 25 minutes on this part of the exam.**
**Percent of Section II score—25**

**Directions:** Show all your work. Indicate clearly the methods you use, because you will be scored on the correctness of your methods as well as on the accuracy and completeness of your results and explanations.

6. Consider an experiment in which two men and two women will be randomly assigned to either a treatment group or a control group in such a way that each group has two people. The people are identified as Man 1, Man 2, Woman 1, and Woman 2. The six possible arrangements are shown below.

| Arrangement A | |
|---|---|
| Treatment | Control |
| Man 1 | Woman 1 |
| Man 2 | Woman 2 |

| Arrangement B | |
|---|---|
| Treatment | Control |
| Man 1 | Man 2 |
| Woman 1 | Woman 2 |

| Arrangement C | |
|---|---|
| Treatment | Control |
| Man 1 | Man 2 |
| Woman 2 | Woman 1 |

| Arrangement D | |
|---|---|
| Treatment | Control |
| Woman 1 | Man 1 |
| Woman 2 | Man 2 |

| Arrangement E | |
|---|---|
| Treatment | Control |
| Man 2 | Man 1 |
| Woman 2 | Woman 1 |

| Arrangement F | |
|---|---|
| Treatment | Control |
| Man 2 | Man 1 |
| Woman 1 | Woman 2 |

Two possible methods of assignment are being considered: the sequential coin flip method, as described in part (a), and the chip method, as described in part (b). For each method, the order of the assignment will be Man 1, Man 2, Woman 1, Woman 2.

(a) For the sequential coin flip method, a fair coin is flipped until one group has two people. An outcome of tails assigns the person to the treatment group, and an outcome of heads assigns the person to the control group. As soon as one group has two people, the remaining people are automatically assigned to the other group.

(i) Complete the table below by calculating the probability of each arrangement occurring if the sequential coin flip method is used.

| Arrangement | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| Probability | | | | | | |

(ii) For the sequential coin flip method, what is the probability that Man 1 and Man 2 are assigned to the same group?

The six arrangements are repeated below.

| Arrangement A | |
| --- | --- |
| Treatment | Control |
| Man 1 | Woman 1 |
| Man 2 | Woman 2 |

| Arrangement B | |
| --- | --- |
| Treatment | Control |
| Man 1 | Man 2 |
| Woman 1 | Woman 2 |

| Arrangement C | |
| --- | --- |
| Treatment | Control |
| Man 1 | Man 2 |
| Woman 2 | Woman 1 |

| Arrangement D | |
| --- | --- |
| Treatment | Control |
| Woman 1 | Man 1 |
| Woman 2 | Man 2 |

| Arrangement E | |
| --- | --- |
| Treatment | Control |
| Man 2 | Man 1 |
| Woman 2 | Woman 1 |

| Arrangement F | |
| --- | --- |
| Treatment | Control |
| Man 2 | Man 1 |
| Woman 1 | Woman 2 |

(b) For the chip method, two chips are marked "treatment" and two chips are marked "control." Each person selects one chip at random without replacement.

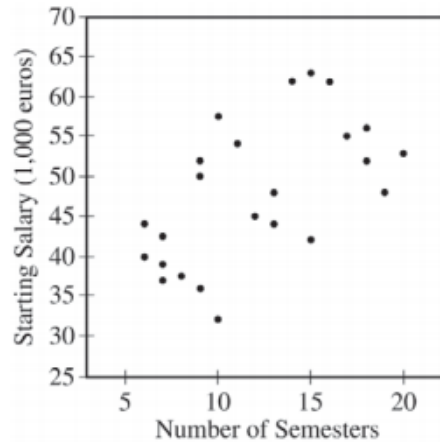(i) Complete the table below by calculating the probability of each arrangement occurring if the chip method is used.

| Arrangement | A | B | C | D | E | F |
| --- | --- | --- | --- | --- | --- | --- |
| Probability | | | | | | |

(ii) For the chip method, what is the probability that Man 1 and Man 2 are assigned to the same group?

(c) Sixteen participants consisting of 10 students and 6 teachers at an elementary school will be used for an experiment to determine lunch preference for the school population of students and teachers. As the participants enter the school cafeteria for lunch, they will be randomly assigned to receive one of two lunches so that 8 will receive a salad, and 8 will receive a grilled cheese sandwich. The students will enter the cafeteria first, and the teachers will enter next. Which method, the sequential coin flip method or the chip method, should be used to assign the treatments? Justify your choice.

6. A newspaper in Germany reported that the more semesters needed to complete an academic program at the university, the greater the starting salary in the first year of a job. The report was based on a study that used a random sample of 24 people who had recently completed an academic program. Information was collected on the number of semesters each person in the sample needed to complete the program and the starting salary, in thousands of euros, for the first year of a job. The data are shown in the scatterplot below.



(a) Does the scatterplot support the newspaper report about number of semesters and starting salary? Justify your answer.

The table below shows computer output from a linear regression analysis on the data.

| Predictor | Coef | SE Coef | T | P |
|---|---|---|---|---|
| Constant | 34.018 | 4.455 | 7.64 | 0.000 |
| Semesters | 1.1594 | 0.3482 | 3.33 | 0.003 |

S = 7.37702    R-Sq = 33.5%    R-Sq(adj) = 30.5%

The table below shows computer output from a linear regression analysis on the data.

| Predictor | Coef | SE Coef | T | P |
|---|---|---|---|---|
| Constant | 34.018 | 4.455 | 7.64 | 0.000 |
| Semesters | 1.1594 | 0.3482 | 3.33 | 0.003 |

S = 7.37702    R-Sq = 33.5%    R-Sq(adj) = 30.5%

(b) Identify the slope of the least-squares regression line, and interpret the slope in context.

An independent researcher received the data from the newspaper and conducted a new analysis by separating the data into three groups based on the major of each person. A revised scatterplot identifying the major of each person is shown below.
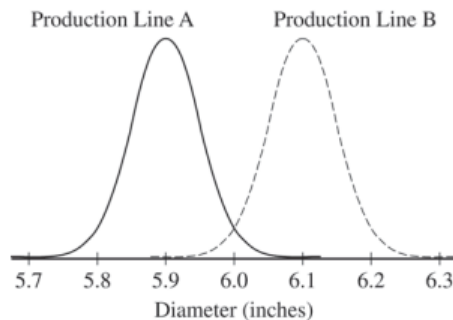


(c) Based on the people in the sample, describe the association between starting salary and number of semesters for the <u>business</u> majors.

(d) Based on the people in the sample, compare the median starting salaries for the three majors.

(e) Based on the analysis conducted by the independent researcher, how could the newspaper report be modified to give a better description of the relationship between the number of semesters and the starting salary for the people in the sample?

6. Corn tortillas are made at a large facility that produces 100,000 tortillas per day on each of its two production lines. The distribution of the diameters of the tortillas produced on production line A is approximately normal with mean 5.9 inches, and the distribution of the diameters of the tortillas produced on production line B is approximately normal with mean 6.1 inches. The figure below shows the distributions of diameters for the two production lines.

Production Line A    Production Line B

5.7    5.8    5.9    6.0    6.1    6.2    6.3
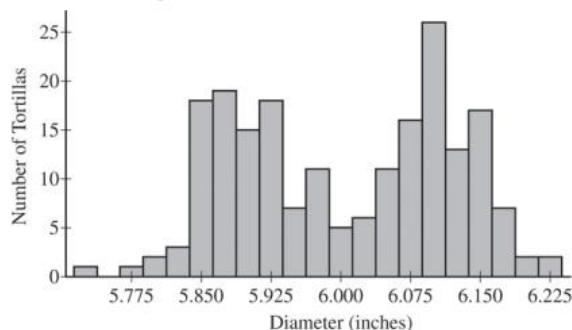Diameter (inches)

The tortillas produced at the factory are advertised as having a diameter of 6 inches. For the purpose of quality control, a sample of 200 tortillas is selected and the diameters are measured. From the sample of 200 tortillas, the manager of the facility wants to estimate the mean diameter, in inches, of the 200,000 tortillas produced on a given day. Two sampling methods have been proposed.

Method 1: Take a random sample of 200 tortillas from the 200,000 tortillas produced on a given day. Measure the diameter of each selected tortilla.

Method 2: Randomly select one of the two production lines on a given day. Take a random sample of 200 tortillas from the 100,000 tortillas produced by the selected production line. Measure the diameter of each selected tortilla.

(a) Will a sample obtained using Method 2 be representative of the population of all tortillas made that day, with respect to the diameters of the tortillas? Explain why or why not.

(b) The figure below is a histogram of 200 diameters obtained by using one of the two sampling methods described. Considering the shape of the histogram, explain which method, Method 1 or Method 2, was most likely used to obtain a such a sample.

5.775    5.850    5.925    6.000    6.075    6.150    6.225
Diameter (inches)

(c) Which of the two sampling methods, Method 1 or Method 2, will result in less variability in the diameters of the 200 tortillas in the sample on a given day? Explain.

Each day, the distribution of the 200,000 tortillas made that day has mean diameter 6 inches with standard deviation 0.11 inch.
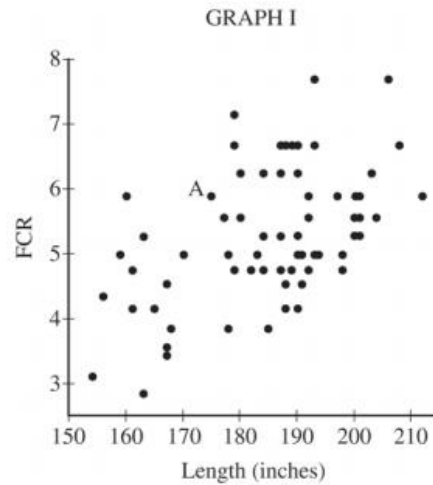
(d) For samples of size 200 taken from one day's production, describe the sampling distribution of the sample mean diameter for samples that are obtained using Method 1.

(e) Suppose that one of the two sampling methods will be selected and used every day for one year (365 days). The sample mean of the 200 diameters will be recorded each day. Which of the two methods will result in less variability in the distribution of the 365 sample means? Explain.

(f) A government inspector will visit the facility on June 22 to observe the sampling and to determine if the factory is in compliance with the advertised mean diameter of 6 inches. The manager knows that, with both sampling methods, the sample mean is an unbiased estimator of the population mean. However, the manager is unsure which method is more likely to produce a sample mean that is close to 6 inches on the day of sampling. Based on your previous answers, which of the two sampling methods, Method 1 or Method 2, is more likely to produce a sample mean close to 6 inches? Explain.

6. Jamal is researching the characteristics of a car that might be useful in predicting the fuel consumption rate (FCR); that is, the number of gallons of gasoline that the car requires to travel 100 miles under conditions of typical city driving. The length of a car is one explanatory variable that can be used to predict FCR. Graph I is a scatterplot showing the lengths of 66 cars plotted with the corresponding FCR. One point on the graph is labeled A.

GRAPH I



Length (inches)

Jamal examined the scatterplot and determined that a linear model would be a reasonable way to express the relationship between FCR and length. A computer output from a linear regression is shown below.

Linear Fit
FCR = −1.595789 + 0.0372614 * Length

Summary of Fit

| | |
|---|---|
| RSquare | 0.250401 |
| Root Mean Square Error | 0.902382 |
| Observations | 66 |

(a) The point on the graph labeled A represents one car of length 175 inches and an FCR of 5.88. Calculate and interpret the residual for the car relative to the least squares regression line.

Jamal knows that it is possible to predict a response variable using more than one explanatory variable. He wants to see if he can improve the original model of predicting FCR from length by including a second explanatory variable in addition to length. He is considering including engine size, in liters, or wheel base (the length between axles), in inches. Graph II is a scatterplot showing the engine size of the 66 cars plotted with the corresponding residuals from the regression of FCR on length. Graph III is a scatterplot showing the wheel base of the 66 cars plotted with the corresponding residuals from the regression of FCR on length.
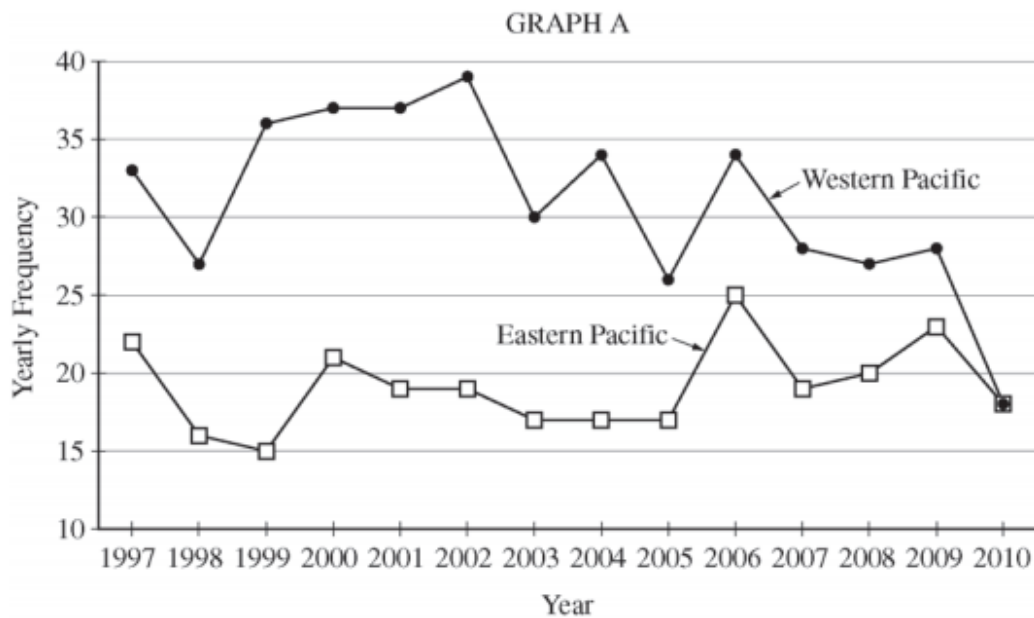
GRAPH II                                        GRAPH III



(b) In graph II, the point labeled A corresponds to the same car whose point was labeled A in graph I. The measurements for the car represented by point A are given below.

| FCR | Length (inches) | Engine Size (liters) | Wheel Base (inches) |
|-----|-----------------|----------------------|---------------------|
| 5.88 | 175 | 3.6 | 93 |

   (i) Circle the point on graph III that corresponds to the car represented by point A on graphs I and II.

   (ii) There is a point on graph III labeled B. It is very close to the horizontal line at 0. What does that indicate about the FCR of the car represented by point B?

(c) Write a few sentences to compare the association between the variables in graph II with the association between the variables in graph III.

(d) Jamal wants to predict FCR using length and one of the other variables, engine size or wheel base. Based on your response to part (c), which variable, engine size or wheel base, should Jamal use in addition to length if he wants to improve the prediction? Explain why you chose that variable.

6. Tropical storms in the Pacific Ocean with sustained winds that exceed 74 miles per hour are called typhoons. Graph A below displays the number of recorded typhoons in two regions of the Pacific Ocean—the Eastern Pacific and the Western Pacific—for the years from 1997 to 2010.

GRAPH A



(a) Compare the distributions of yearly frequencies of typhoons for the two regions of the Pacific Ocean for the years from 1997 to 2010.

(b) For each region, describe how the yearly frequencies changed over the time period from 1997 to 2010.

A moving average for data collected at regular time increments is the average of data values for two or more consecutive increments. The 4-year moving averages for the typhoon data are provided in the table below. For example, the Eastern Pacific 4-year moving average for 2000 is the average of 22, 16, 15, and 21, which is equal to 18.50.

| Year | Number of Typhoons in the Eastern Pacific | Eastern Pacific 4-year moving average | Number of Typhoons in the Western Pacific | Western Pacific 4-year moving average |
|------|------|------|------|------|
| 1997 | 22 | | 33 | |
| 1998 | 16 | | 27 | |
| 1999 | 15 | | 36 | |
| 2000 | 21 | 18.50 | 37 | 33.25 |
| 2001 | 19 | 17.75 | 37 | 34.25 |
| 2002 | 19 | 18.50 | 39 | 37.25 |
| 2003 | 17 | 19.00 | 30 | 35.75 |
| 2004 | 17 | 18.00 | 34 | 35.00 |
| 2005 | 17 | 17.50 | 26 | 32.25 |
| 2006 | 25 | 19.00 | 34 | 31.00 |
| 2007 | 19 | 19.50 | 28 | 30.50 |
| 2008 | 20 | 20.25 | 27 | 28.75 |
| 2009 | 23 | 21.75 | 28 | 29.25 |
| 2010 | 18 | 20.00 | 18 | |

(c) Show how to calculate the 4-year moving average for the year 2010 in the Western Pacific. Write your value in the appropriate place in the table.

(d) Graph B below shows both yearly frequencies (connected by dashed lines) and the respective 4-year moving averages (connected by solid lines). Use your answer in part (c) to complete the graph.



GRAPH B

(e)  Consider graph B.

   i)  What information is more apparent from the plots of the 4-year moving averages than from the plots of the yearly frequencies of typhoons?

   ii)  What information is less apparent from the plots of the 4-year moving averages than from the plots of the yearly frequencies of typhoons?

6. Two students at a large high school, Peter and Rania, wanted to estimate $\mu$, the mean number of soft drinks that a student at their school consumes in a week. A complete roster of the names and genders for the 2,000 students at their school was available. Peter selected a simple random sample of 100 students. Rania, knowing that 60 percent of the students at the school are female, selected a simple random sample of 60 females and an independent simple random sample of 40 males. Both asked all of the students in their samples how many soft drinks they typically consume in a week.

(a) Describe a method Peter could have used to select a simple random sample of 100 students from the school.

Peter and Rania conducted their studies as described. Peter used the sample mean $\overline{X}$ as a point estimator for $\mu$.

Rania used $\overline{X}_{overall} = (0.6)\overline{X}_{female} + (0.4)\overline{X}_{male}$ as a point estimator for $\mu$, where $\overline{X}_{female}$ is the mean of the sample of 60 females and $\overline{X}_{male}$ is the mean of the sample of 40 males.

Summary statistics for Peter's data are shown in the table below.

| Variable | N | Mean | Standard Deviation |
|---|---|---|---|
| Number of soft drinks | 100 | 5.32 | 4.13 |

(b) Based on the summary statistics, calculate the estimated standard deviation of the sampling distribution (sometimes called the standard error) of Peter's point estimator $\overline{X}$.

Summary statistics for Rania's data are shown in the table below.

| Variable | Gender | N | Mean | Standard Deviation |
|---|---|---|---|---|
| Number of soft drinks | Female | 60 | 2.90 | 1.80 |
| | Male | 40 | 7.45 | 2.22 |

(c) Based on the summary statistics, calculate the estimated standard deviation of the sampling distribution of Rania's point estimator $\overline{X}_{overall} = (0.6)\overline{X}_{female} + (0.4)\overline{X}_{male}$.

A dotplot of Peter's sample data is given below.



Number of Soft Drinks

Comparative dotplots of Rania's sample data are given below.



Number of Soft Drinks

(d) Using the dotplots above, explain why Rania's point estimator has a smaller estimated standard deviation than the estimated standard deviation of Peter's point estimator.

**Part B**

**Question 6**

**Spend about 25 minutes on this part of the exam.**

**Percent of Section II score—25**

**Directions:** Show all your work. Indicate clearly the methods you use, because you will be scored on the correctness of your methods as well as on the accuracy and completeness of your results and explanations.

6. Every year, each student in a nationally representative sample is given tests in various subjects. Recently, a random sample of 9,600 twelfth-grade students from the United States were administered a multiple-choice United States history exam. One of the multiple-choice questions is below. (The correct answer is C.)

> In 1935 and 1936 the Supreme Court declared that important parts of the New Deal were unconstitutional. President Roosevelt responded by threatening to
>
> (A) impeach several Supreme Court justices
> (B) eliminate the Supreme Court
> (C) appoint additional Supreme Court justices who shared his views
> (D) override the Supreme Court's decisions by gaining three-fourths majorities in both houses of Congress

Of the 9,600 students, 28 percent answered the multiple-choice question correctly.

(a) Let $p$ be the proportion of all United States twelfth-grade students who would answer the question correctly. Construct and interpret a 99 percent confidence interval for $p$.

Assume that students who actually know the correct answer have a 100 percent chance of answering the question correctly, and students who do not know the correct answer to the question guess completely at random from among the four options.

Let $k$ represent the proportion of all United States twelfth-grade students who actually know the correct answer to the question.

(b) A tree diagram of the possible outcomes for a randomly selected twelfth-grade student is provided below. Write the correct probability in each of the five empty boxes. Some of the probabilities may be expressions in terms of $k$.

TREE DIAGRAM OF OUTCOMES FOR A
RANDOMLY SELECTED TWELFTH-GRADE STUDENT

| | Outcome | Probability |
|---|---|---|
| Conditional probability = 1   Answers correctly | Knows answer and answers correctly | Probability = $k$ |
| Probability = $k$   Knows correct answer | | |
| Conditional probability = 0   Answers incorrectly | Knows answer and answers incorrectly | Probability = 0 |
| Conditional probability = ☐   Answers correctly | Guesses at random and answers correctly | Probability = ☐ |
| Probability = ☐   Guesses at random | | |
| Conditional probability = ☐   Answers incorrectly | Guesses at random and answers incorrectly | Probability = ☐ |

(c) Based on the completed tree diagram, express the probability, in terms of $k$, that a randomly selected twelfth-grade student would correctly answer the history question.

(d) Using your interval from part (a) and your answer to part (c), calculate and interpret a 99 percent confidence interval for $k$, the proportion of all United States twelfth-grade students who actually know the answer to the history question. You may assume that the conditions for inference for the confidence interval have been checked and verified.

6. Hurricane damage amounts, in millions of dollars per acre, were estimated from insurance records for major hurricanes for the past three decades. A stratified random sample of five locations (based on categories of distance from the coast) was selected from each of three coastal regions in the southeastern United States. The three regions were Gulf Coast (Alabama, Louisiana, Mississippi), Florida, and Lower Atlantic (Georgia, South Carolina, North Carolina). Damage amounts in millions of dollars per acre, adjusted for inflation, are shown in the table below.

HURRICANE DAMAGE AMOUNTS IN MILLIONS OF
DOLLARS PER ACRE

| | Distance from Coast | | | | |
|---|---|---|---|---|---|
| | < 1 mile | 1 to 2 miles | 2 to 5 miles | 5 to 10 miles | 10 to 20 miles |
| Gulf Coast | 24.7 | 21.0 | 12.0 | 7.3 | 1.7 |
| Florida | 35.1 | 31.7 | 20.7 | 6.4 | 3.0 |
| Lower Atlantic | 21.8 | 15.7 | 12.6 | 1.2 | 0.3 |

(a) Sketch a graphical display that compares the hurricane damage amounts per acre for the three different coastal regions (Gulf Coast, Florida, and Lower Atlantic) and that also shows how the damage amounts vary with distance from the coast.

(b) Describe differences and similarities in the hurricane damage amounts among the three regions.

Because the distributions of hurricane damage amounts are often skewed, statisticians frequently use rank values to analyze such data.

(c) In the table below, the hurricane damage amounts have been replaced by the ranks 1, 2, or 3. For each of the distance categories, the highest damage amount is assigned a rank of 1 and the lowest damage amount is assigned a rank of 3. Determine the missing ranks for the 10-to-20-miles distance category and calculate the average rank for each of the three regions. Place the values in the table below.

### ASSIGNED RANKS WITHIN DISTANCE CATEGORIES

| | Distance from Coast | | | | | Average Rank |
|---|---|---|---|---|---|---|
| | < 1 mile | 1 to 2 miles | 2 to 5 miles | 5 to 10 miles | 10 to 20 miles | |
| Gulf Coast | 2 | 2 | 3 | 1 | | |
| Florida | 1 | 1 | 1 | 2 | | |
| Lower Atlantic | 3 | 3 | 2 | 3 | | |

(d) Consider testing the following hypotheses.

$H_0$: There is no difference in the distributions of hurricane damage amounts among the three regions.

$H_a$: There is a difference in the distributions of hurricane damage amounts among the three regions.

If there is no difference in the distribution of hurricane damage amounts among the three regions (Gulf Coast, Florida, and Lower Atlantic), the expected value of the average rank for each of the three regions is 2. Therefore, the following test statistic can be used to evaluate the hypotheses above:

$$Q = 5\left[\left(\bar{R}_G - 2\right)^2 + \left(\bar{R}_F - 2\right)^2 + \left(\bar{R}_A - 2\right)^2\right]$$

where $\bar{R}_G$ is the average rank over the five distance categories for the Gulf Coast (and $\bar{R}_F$ and $\bar{R}_A$ are similarly defined for the Florida and Lower Atlantic coastal regions).
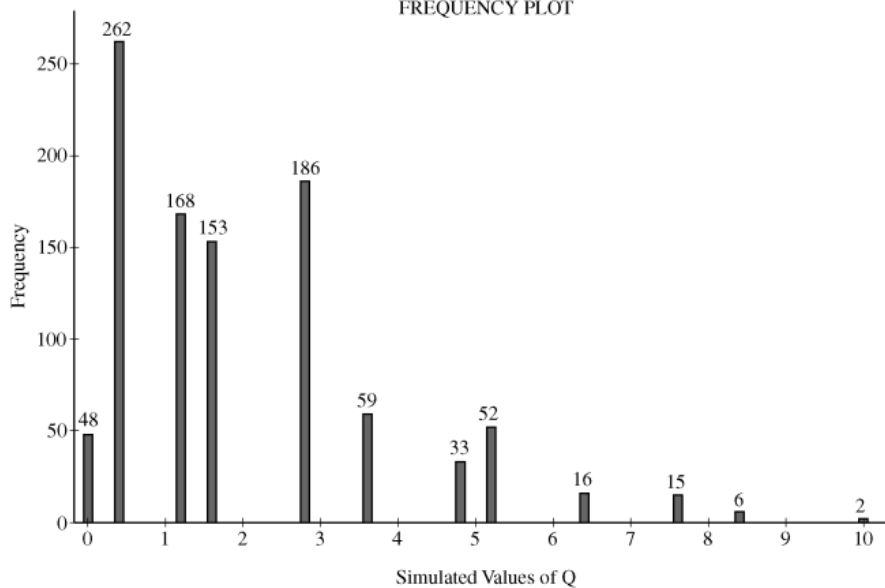
Calculate the value of the test statistic $Q$ using the average ranks you obtained in part (c).

(e) One thousand simulated values of this test statistic, $Q$, were calculated, assuming no difference in the distributions of hurricane damage amounts among the three coastal regions. The results are shown in the table below. These data are also shown in the frequency plot where the heights of the lines represent the frequency of occurrence of simulated values of $Q$.

Frequency Table for Simulated Values of Q

| Q | Frequency | Cumulative Frequency | Percent | Cumulative Percent |
|------|-----------|----------------------|---------|--------------------|
| 0.0  | 48        | 48                   | 4.80    | 4.80               |
| 0.4  | 262       | 310                  | 26.20   | 31.00              |
| 1.2  | 168       | 478                  | 16.80   | 47.80              |
| 1.6  | 153       | 631                  | 15.30   | 63.10              |
| 2.8  | 186       | 817                  | 18.60   | 81.70              |
| 3.6  | 59        | 876                  | 5.90    | 87.60              |
| 4.8  | 33        | 909                  | 3.30    | 90.90              |
| 5.2  | 52        | 961                  | 5.20    | 96.10              |
| 6.4  | 16        | 977                  | 1.60    | 97.70              |
| 7.6  | 15        | 992                  | 1.50    | 99.20              |
| 8.4  | 6         | 998                  | 0.60    | 99.80              |
| 10.0 | 2         | 1000                 | 0.20    | 100.00             |

FREQUENCY PLOT



Use these simulated values and the test statistic you calculated in part (d) to determine if the observed data provide evidence of a significant difference in the distributions of hurricane damage amounts among the three coastal regions. Explain.