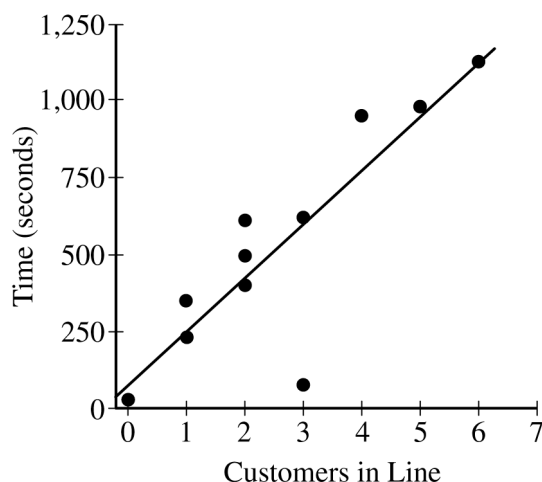**STATISTICS**
**SECTION II**
**Part A**
**Questions 1-5**
**Spend about 1 hour and 5 minutes on this part of the exam.**
**Percent of Section II score—75**

**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.

1. The manager of a grocery store selected a random sample of 11 customers to investigate the relationship between the number of customers in a checkout line and the time to finish checkout. As soon as the selected customer entered the end of a checkout line, data were collected on the number of customers in line who were in front of the selected customer and the time, in seconds, until the selected customer was finished with the checkout. The data are shown in the following scatterplot along with the corresponding least-squares regression line and computer output.



| Predictor | Coef | SE Coef | T | P |
|---|---|---|---|---|
| Constant | 72.95 | 110.36 | 0.66 | 0.525 |
| Customers in line | 174.40 | 35.06 | 4.97 | 0.001 |
| S = 200.01 | | R-Sq = 73.33% | R-Sq (adj) = 70.37% | |

   (a) Identify and interpret in context the estimate of the intercept for the least-squares regression line.

   (b) Identify and interpret in context the coefficient of determination, $r^2$.

   (c) One of the data points was determined to be an outlier. Circle the point on the scatterplot and explain why the point is considered an outlier.

**GO ON TO THE NEXT PAGE.**

2. An environmental science teacher at a high school with a large population of students wanted to estimate the proportion of students at the school who regularly recycle plastic bottles. The teacher selected a random sample of students at the school to survey. Each selected student went into the teacher's office, one at a time, and was asked to respond yes or no to the following question.

> Do you regularly recycle plastic bottles?

Based on the responses, a 95 percent confidence interval for the proportion of all students at the school who would respond yes to the question was calculated as $(0.584, 0.816)$.

(a) How many students were in the sample selected by the environmental science teacher?

(b) Given the method used by the environmental science teacher to collect the responses, explain how bias might have been introduced and describe how the bias might affect the point estimate of the proportion of all students at the school who would respond yes to the question.

(c) The statistics teacher at the high school was concerned about the potential bias in the survey. To obtain a potentially less biased estimate of the proportion, the statistics teacher used an alternate method for collecting student responses. A random sample of 300 students was selected, and each student was given the following instructions on how to respond to the question.

  - In private, flip a fair coin.
  - If heads, you must respond no, regardless of whether you regularly recycle.
  - If tails, please truthfully respond yes or no.

  (i) What is the expected number of students from the sample of 300 who would be required to respond no because the coin flip resulted in heads?

  (ii) The results of the sample showed that 213 of the 300 selected students responded no. Based on the results of the sample, give a point estimate for the <u>proportion</u> of all students at the high school who would respond <u>yes</u> to the question.

**GO ON TO THE NEXT PAGE.**

## Question 1

**Intent of Question**

The primary goals of this question were to assess a student's ability to (1) identify various values in regression computer output; (2) interpret the intercept of a regression line in context; (3) interpret the coefficient of determination $(r^2)$ in context; and (4) identify an outlier from a scatterplot.

**Solution**

**Part (a):**

The estimate of the intercept is 72.95. It is estimated that the average time to finish checkout if there are no other customers in line is 72.95 seconds.

**Part (b):**

The coefficient of determination is $r^2 = 73.33\%$. This value indicates that 73.33% of the variability in the times it takes customers to finish checkout, including time waiting in line, can be explained by knowing how many customers are in line in front of the selected customer.

**Part (c):**

The outlier is the point with $x = 3$ and $y$ close to 0. This point is considered an outlier because the combination of $x$ and $y$ values differs from the pattern of the rest of the data. Specifically, the value of $y$ (time to finish checkout) is much lower than would be expected when there are $x = 3$ customers in line in front of the selected customer, given the remaining data.

**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 satisfies the following three components:
1. Correctly identifies 72.95 as the intercept.
2. Communicates the concept of a $y$-intercept in a context that includes both time and zero customers.
3. Indicates that the value of the intercept is a prediction by using language such as "predicted," "estimated," or "average" value of $y$.

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

Incorrect (I) if the response includes at most one of the three components.

*Notes:*
- Regression equations (such as $\hat{y} = 72.95 + 174.40x$) cannot be used to satisfy identification of the intercept in component 1, unless the intercept is explicitly labeled as such.
- A regression equation cannot be used to satisfy component 3.
- Incorrect regression equations are treated as extraneous and do not affect the scoring of any component.
- A response that interprets 72.95 as a slope does not satisfy components 1 or 2.

**Part (b)** is scored as follows:

Essentially correct (E) if the response satisfies the following three components:
1. Correctly identifies 73.33% as the coefficient of determination.
2. Provides a correct (possibly generic) interpretation of $r^2$.
3. Interpretation includes context.

Partially correct (P) if the response satisfies only two of the three components;
    *OR*
if the response satisfies the three components, but reverses the roles of number of customers in line and time to finish checkout in the interpretation.

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

*Notes:*
- In component 2 the correct interpretation of the coefficient of determination can take any of several equivalent forms, such as:
  - The percent variability in $y$ that is attributed to the linear relationship between $y$ and $x$ or between $x$ and $y$.
  - The proportion of the total variability in the dependent variable $y$ that is explained by the independent variable $x$.
  - The proportion of variation in $y$ that is accounted for by the linear model.
  - The proportionate reduction of total variation of the $y$ values that is associated with the use of the independent variable $x$.
  - The proportionate reduction in the sum of the squares of vertical deviations obtained by using the least-squares line instead of the naïve prediction of $\overline{y}$.
- In component 2 common *incorrect* interpretations of the coefficient of determination include:
  - The percent variability in the *predicted* $y$ values that is explained by the linear relationship between $y$ and $x$.
  - The percent variability in the *data* that is explained by the linear relationship between $y$ and $x$.
  - The percent variability that is explained by the linear relationship between $y$ and $x$.
  - The percent variability in $y$ that is *on average* explained by the linear relationship between $y$ and $x$.
- For component 3 context must include mention of time or customers.

**Part (c)** is scored as follows:

Essentially correct (E) if the response satisfies the following two components:
1. Correctly identifies the outlier.
2. Describes an unusual feature of the identified scatter plot point, relative to the remaining data points, that is sufficient to identify it as the outlier. Examples include:
   - The combination of *x* and *y* values is unusual compared to the other points.
   - The value of *y* is much lower than would be expected (or predicted), given the remaining data.
   - The residual for the point is unusually large relative to the other residuals.

Partially correct (P) if the response satisfies component 1 but does not satisfy component 2.

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

*Notes:*
- In the absence of any point being circled on the graph, component 1 can still be satisfied by explicitly referring to the coordinates of the outlier. Valid coordinates for outlier identification must specify an *x* value of 3 and a *y* value that is strictly between 0 and 250.
- A response that does not make a comparison to the remaining data points, such as stating the outlier has a large residual or is nowhere near the regression line, does not satisfy component 2.
- A response that makes a comparison to the remaining data points based upon an unusual feature that is *insufficient* for outlier identification, such as stating the point is the only point with that particular *y* value, does not satisfy component 2.
- In the absence of explicit numerical calculation, a response that appeals to the influence that the outlier has on the regression coefficient estimates or on the sample correlation coefficient does not satisfy component 2.