

Regression Analysis for Business (STAT 613)

SAMPLE EXAM

■ Instructions

Read these instructions carefully.

Questions 1-20 of this exam are used to place into STAT 621.

Questions 1-43 are used to waive STAT completely.

This is a **closed-book** exam. You are allowed to use a calculator and one page (8.5 by 11 inches or A4, both sides) of **handwritten** notes. No use of cellular telephones or other portable electronics is permitted.

You have **two hours** for the exam. There are **43** questions. The **computer output** associated with one or more questions should be considered an essential part of the questions. The multiple-choice questions are **equally weighted**; the number of correct answers determines your grade.

Throughout this exam, the word “significant” implies “statistically significant” and by “sample” we mean a simple random sample. Use 95% confidence intervals and a p -value threshold of 0.05 to determine statistical significance unless otherwise instructed. All logarithms are natural logs (that is, \ln or \log_e) unless otherwise noted.

All categorical predictor variables have been coded according to a dummy variable coding scheme with the reference category being the level of the variable that is last, using an alphanumeric sort.

Please note the following when filling in the answer form:

- Mark the answer form only using a **pencil**. Erase changes completely.
- **Fill in your name and student id number** on the answer form.
- **Mark the “bubbles”** under your name and student id number.
- Choose the **one best** answer by marking the item on the answer form.

When you have completed the exam, turn in the answer form and your exam with your name on it. Solutions will be posted in Canvas.

■ STOP

Do NOT turn the page until you are instructed to proceed.

1. The amount spent by a customer is normally distributed with mean $\mu = \$300$ and $\sigma = \$50$. The probability that a randomly selected customer spends less than \$200 is approximately
 - a. 0.167
 - b. 0.334
 - c. 0.050
 - d. 0.025
 - e. 0.010

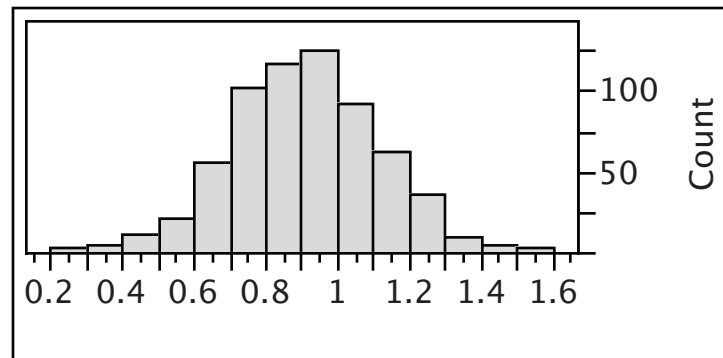
2. You toss a fair coin repeatedly. Of the following, which has the highest probability of getting exactly 50% heads.
 - a. Toss it once
 - b. Toss it 10 times
 - c. Toss it 100 times
 - d. Toss it 1000 times.
 - e. Toss it 1,000,000 times.

3. The IQs of a large population is Normally distributed with a mean of 100 and an SD of 15 points. Suppose you randomly choose 40 people from the population. What is the approximate chance that you get at least one person with an IQ of 130 or higher?
 - a. .025
 - b. .33
 - c. .36
 - d. .5
 - e. .64

4. The standard deviation in a population of incomes is $\sigma = \$20,000$. To obtain a 95% confidence interval for μ with total width (width = upper endpoint – lower endpoint) less than \$500 requires a sample size of about
 - a. 5,000
 - b. 100
 - c. 50
 - d. 2,000
 - e. More than 25,000

5. The standard error of the mean
- Estimates the SD of the population.
 - Increases with the size of a sample.
 - Measures the sample-to-sample standard deviation of sample means.
 - Determines the sample size needed in order to apply the central limit theorem.
 - Is the expected size of the deviation of \bar{x} -bar from μ .

QUESTIONS 6-7



6. Consider the histogram shown immediately above this question. Based on this histogram, it is evident that
- The population mean $\mu \leq 0$.
 - The sample size is less than 200.
 - The 95% confidence interval for the mean includes 0.
 - The sample is too small to produce a normally distributed sampling distribution.
 - The sample mean is about the same as the sample median.
7. Refer again to the histogram shown before Question 6. The standard deviation of these data is approximately
- 0.01
 - 0.2
 - 2.0
 - 0.05
 - 1.0

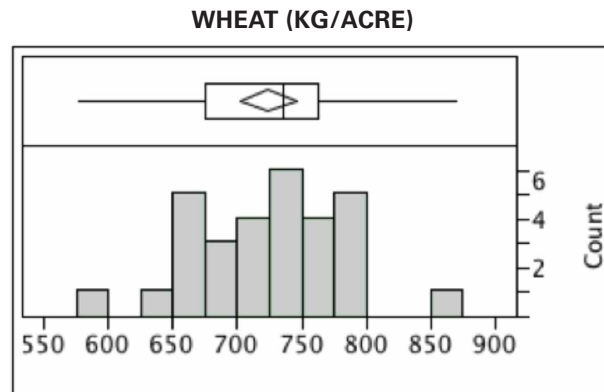
QUESTIONS 8–10

A retail web site gathered data about customers who bought either a camera or a phone (but not both). Within this population, 40% of the customers bought a camera. Of those customers who bought cameras, 25% reported incomes less than \$50,000. Among those who bought a mobile phone, 50% reported incomes less than \$50,000.

8. These results show that the proportion of this population of customers who have incomes of at least \$50,000 is
- a. 0.50
 - b. 0.40
 - c. 0.25
 - d. 0.60
 - e. 0.35
9. The web site profits \$10 for each camera and \$5 for each mobile phone sold. The expected value of the profit produced by the purchase choice of a randomly selected customer from this population is
- a. \$5
 - b. \$8
 - c. \$6
 - d. \$7
 - e. \$9
10. The probability that a randomly chosen customer from this population who has income above \$50,000 buys a phone is
- a. 0.25
 - b. 0.40
 - c. 0.50
 - d. 0.60
 - e. 0.35

QUESTIONS 11–16

An agricultural corporation grows wheat on 10,000 acres. To estimate the total harvest, the corporation measures the yield (in kilograms per acre, kg/acre) just prior to harvest on a sample of 30 one-acre, non-overlapping tracts randomly located over this land.



100.0%	maximum	870	Mean	724.133
75.0%	quartile	763	Std Dev	60.146
50.0%	median	741		
25.0%	quartile	675	N	30
0.0%	minimum	576		

- 11.** These data imply that the probability of observing a tract selected at random from this population with yield less than 741 kg/acres is
- About 25%.
 - About 50%.
 - About 50%, provided the population is normally distributed.
 - About 95%.
 - Less than 50%.
- 12.** Had the corporation gathered a larger sample of 100 tracts, then we should expect to find
- Smaller average yield.
 - Larger average yield.
 - Smaller standard deviation.
 - Smaller maximum yield.
 - Smaller standard error of the mean.

- 13.** Sales agreements were designed under the expectation that the average yield this corporation would produce is 750 kg/acre when fully harvested. Assuming the appropriate conditions are met, this analysis indicates that the average yield will be
- Statistically significantly less than expected.
 - Less than expected, but not significantly.
 - More than expected, but not significantly.
 - Statistically significantly more than expected.
 - Within 1% of the amount planners expected.
- 14.** Do these data satisfy the conditions of a one-sample confidence interval for the mean?
- Yes.
 - No, the sample size $n = 30$ is too small.
 - No, the data are not normally distributed.
 - No, the data were not sampled with replacement.
 - No, outlying yields have produced an unreliable estimated standard deviation.
- 15.** When the full 10,000 acres is harvested, this analysis (assuming the usual conditions are met and with 95% confidence) can be expected to provide
- Between 7,020,000 to 7,460,000 kg wheat.
 - The sample does not produce an estimate of the total yield.
 - Between 5,760,000 to 8,700,000 kg wheat
 - Between 6,750,000 to 7,630,000 kg wheat
 - Between 6,040,000 to 8,440,000 kg wheat
- 16.** Prior to this season, the firm from which the corporation purchased seeds promised that this variety would produce on average more than 700 kg wheat per acre. Assuming the necessary conditions are met, a statistical test of the null hypothesis $H_0: \mu \leq 700$ kg/acre has p -value
- About 1/6
 - About 1/3
 - Less than 0.05
 - About 1/2
 - Larger than 1/2

QUESTIONS 17–20

The Transportation Department (TDP) is concerned about over-weight trucks damaging highways. TDP maintains that no more than 10% of trucks on the roads are over-weight and has publically indicated that the percentage is less. As a precaution, TDP plans to weigh a sample of trucks operating on major routes. Unless the data reject its beliefs, TDP will continue normal operations. If the data reject its beliefs, TDP will institute reforms. Let p denote the population proportion of overweight trucks operating on roads.

- 17.** The appropriate null hypothesis for TDP to test is
- $H_0: p \leq 0.10$
 - $H_0: \text{sample proportion} \leq 0.10$
 - $H_0: p = 0.10$
 - $H_0: p = 0.50$
 - $H_0: p \geq 0.10$
- 18.** TDP plans to estimate p by the sample proportion of overweight trucks from a random sample of size $n = 100$. If in the population, $p = 0.10$ then the probability that more than 20% of the trucks in a sample are overweight
- Is about 0.13.
 - Is less than 0.05.
 - Is about 0.37.
 - Cannot be determined without further information.
 - Is more than 0.84.
- 19.** A critic of TDP gathered a sample of $n = 400$ trucks and used these to test $H_0: p \leq 0.09$ with $n = 400$. In her random sample, the sample proportion of overweight trucks was 0.10 (10%). She should conclude that
- She needs a larger sample to test H_0 .
 - She should reject H_0 .
 - She cannot reject H_0 .
 - She should change H_0 to match the data, revising it to $H_0: p \leq 0.10$.
 - The data do not meet the conditions for using a one-sample t -test.

- 20.** An analyst dislikes hypothesis testing and produced a 95% confidence interval from an independent sample of data, with the interval found to be 0.093 to 0.157. Based on this confidence interval (assuming the usual conditions), the analyst should
- Reject $H_0: p = 0.10$.
 - Reject $H_0: p \neq 0.10$.
 - Use a 90% confidence interval instead.
 - Gather more data until the confidence interval omits 0.10.
 - Not reject $H_0: p = 0.10$.

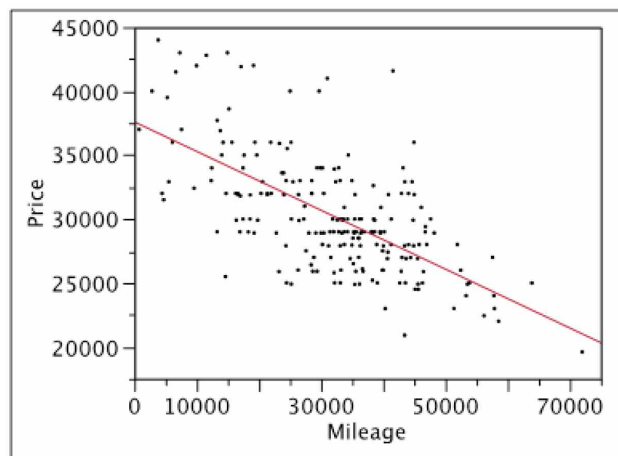
Only answer the rest of the questions if you wish to waive STAT completely.

- 21.** Which of the following summaries from a fitted regression best measures the degree of collinearity associated with an estimated coefficient?
- R^2 .
 - RMSE.
 - The VIF for the coefficient.
 - The standard error of the slope.
 - The 95% confidence band for the true regression line.
- 22.** If the variances of the error terms in a simple regression increase with increasing values of the explanatory variable X and one incorrectly assumes the SRM, then the
- Estimated slope will be too steep.
 - Estimated intercept will be too close to zero.
 - Explanatory variable should be re-expressed on a log scale.
 - Residuals from the fit will appear to be autocorrelated.
 - Prediction intervals for small values of X will be too wide.
- 23.** When building a regression model with a categorical explanatory variable, a common diagnostic plot shows side-by-side comparison boxplots of the model residuals grouped by the levels of the categorical variable. This plot is most useful to
- Check the assumption of independence.
 - Check assumption of equal variance.
 - Check the assumption of normality.
 - Identify the presence of leveraged outliers.
 - Determine the statistical significance of the categorical variable.

24. Regression models assume that some of the variability in the response is due to random sources not identified in the model's equation. The estimated standard deviation of this unexplained variation is known as
- RMSE.
 - Standard error of Y .
 - Standard error of X .
 - R^2 .
 - F - ratio.
25. If the variable X_2 is added to a simple regression that includes X_1 , then which of the following must happen in the multiple regression if X_1 and X_2 are uncorrelated?
- RMSE must get smaller.
 - The partial regression coefficient for X_2 will be smaller than that for X_1 .
 - The overall ANOVA F -statistic will have a larger p -value.
 - The partial regression coefficient for X_1 will be the same as the marginal coefficient.
 - The overall ANOVA F -statistic will have a smaller p -value.

QUESTIONS 26–31

The data shown in the following scatterplot and simple regression report the price and the number of miles driven for 218 used cars (all in the BMW 325 series) offered for sale in the San Francisco area.



$$\text{PRICE (\$)} = 37650 - 0.24 \text{ MILEAGE}$$

RSquare	0.42
Root Mean Square Error	3500
Mean of Response	30269
Observations	218

- 26.** According to the fitted equation, a car like these with 40,000 miles would be expected to cost on average
- \$37,650.
 - \$30,450.
 - \$28,050.
 - \$36,690.
 - \$25,050.
- 27.** A used car with 40,000 miles is offered for sale at a price that is \$3,500 above the prediction from this model. Given that the assumptions of the SRM hold, then this car's price is
- Larger than about 50% of cars with 40,000 miles.
 - Larger than about 95% of cars with 40,000 miles.
 - Larger than about 5% of cars with 40,000 miles.
 - Larger than about 84% of cars with 40,000 miles.
 - Larger than about 67% of cars with 40,000 miles.
- 28.** The equation of the fitted model implies that on average for cars such as these, an additional 1,000 miles of driving
- Has no effect on the expected price of the used car.
 - Increases the expected price by about \$240.
 - Decreases the expected price by about \$240.
 - Increases the expected price by about 24 percent.
 - Decreases the expected price by about 240 percent.
- 29.** An economist claims that the elasticity of price with respect to mileage driven is constant. How would you adapt this model to estimate the elasticity?
- The model does not need to be changed.
 - Take the log of price but not mileage.
 - Take the log of mileage but not price.
 - Add a quadratic term to the model to create a power function.
 - Regress the log of price against the log of mileage.

- 30.** To obtain a more precise estimate of the slope in this model, we should (assuming the assumptions of the SRM hold)
- a.** Remove the outlying expensive cars with prices above \$40,000.
 - b.** Remove all of the cars with mileage above 50,000.
 - c.** Add prices for 10 more cars which have the average mileage $\bar{x} = 32,000$ miles.
 - d.** Add prices for 10 low-mileage cars.
 - e.** Add prices for 5 low-mileage cars and 5 cars with 50,000 to 60,000 miles.
- 31.** If an additional variable, the age of the car, were added to the regression, then which of the following would you expect to happen.
- a.** Severe autocorrelation would be introduced.
 - b.** Collinearity would be eliminated.
 - c.** R-squared would decrease.
 - d.** The estimated regression coefficient for mileage would be closer to zero.
 - e.** The standard deviation of the residuals would increase.

QUESTIONS 32–43

A banking analyst was interested in predicting the yield on a ten-year bond for a new issue from a company in one of two South East Asian countries (Malaysia or Singapore). The analyst collected data on 85 previous bond issues from different companies in the region. Output on the next page shows a model that includes a country effect, the annual revenue of the bond-issuing company in millions of \$US in the previous year and a measure of the financial leverage of the company (calculated as the ratio of total liabilities to net worth). Note that the variable “Financial Leverage” has nothing to do whatsoever with the concept of statistical leverage.

SUMMARY OF FIT

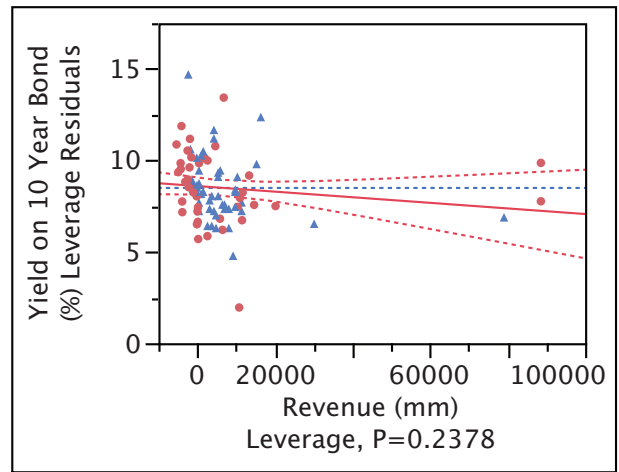
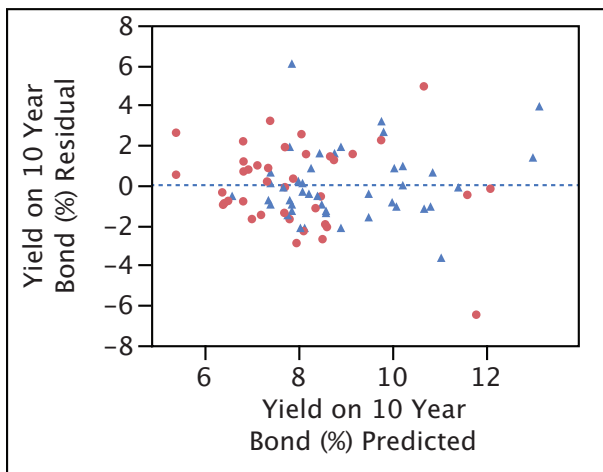
RSquare	0.407634
Root Mean Square Error	1.931565
Mean of Response	8.490706
Observations (or Sum Wgts)	85

ANALYSIS OF VARIANCE

Source	DF	Sum of Squares	Mean Square	F Ratio
Model	3	207.96174	69.3206	18.5799
Error	81	302.20642	3.7309	Prob > F
C. Total	84	510.16816		<.0001*

INDICATOR FUNCTION PARAMETERIZATION

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	6.4595152	0.434899	14.85	<.0001*
Revenue (mm)	-1.531e-5	1.287e-5	-1.19	0.2378
Financial Leverage	0.7118301	0.11191	6.36	<.0001*
country[Malaysia]	1.0083861	0.423381	2.38	0.0196*



- 32.** The estimated yield on a bond based on the fitted model for a Singaporean company with revenues of \$25 (mm) and a Financial Leverage of 2 is approximately
- 15.4
 - 40.6
 - 7.9
 - 7.5
 - 8.5
- 33.** If we remove the variable Financial Leverage from this multiple regression, then (given the assumptions of the Multiple Regression Model [MRM] hold)
- The value of R^2 would increase by a significant amount.
 - The value of R^2 would increase.
 - The value of R^2 would decrease.
 - The value of R^2 would decrease by a significant amount.
 - The change in the value of R^2 cannot be determined.
- 34.** A bond from a company with \$20 (mm) revenue and Financial Leverage of 3 is issued. In which country would it be expected to have the lower yield?
- Singapore.
 - Malaysia.
 - The model predicts the same yield in both countries.
 - There is no information available in the model to answer this question.
 - This question cannot be answered without an interaction in the model.
- 35.** If a company in Singapore were to increase its Financial Leverage by 1, keeping revenues constant, this model implies that the yield would be expected to
- Decrease on average by between (0.488, 0.936) with 95% confidence
 - Decrease on average by between (-3.15, 4.58) with 95% confidence.
 - Increase on average by between (0.488, 0.936) with 95% confidence.
 - Increase on average by between (-3.15, 4.58) with 95% confidence
 - Decrease on average by between (0.598, 0.825) with 95% confidence

- 36.** The fitted regression model includes a positive estimated coefficient for Financial Leverage. The best interpretation of this coefficient is that
- a.** If this variable were removed from the model, then R-squared would increase but not by a significant amount.
 - b.** Within a specific country and for those with identical Financial Leverages but different Revenues, companies are expected to have identical yields.
 - c.** Within a specific country and for those with identical Revenues, companies with greater Financial Leverage are expected to have higher yields, but the difference is not significant.
 - d.** Within a specific country and for those with identical Revenues, companies with greater Financial Leverage are expected to have higher yields and the difference is significant.
 - e.** Companies with more leverage have higher yields.
- 37.** If the analysts had measured Revenues in Japanese Yen rather than dollars and rerun the above model then which of the following regression summaries would change in value?
- a.** R^2 .
 - b.** Standard error of the estimated Revenue slope.
 - c.** t-statistic for the Revenue slope.
 - d.** RMSE.
 - e.** Overall Anova F-ratio.
- 38.** With reference to the leverage plot for Revenue, which of the following is a reasonable conclusion?
- a.** There is severe collinearity in this dataset.
 - b.** Autocorrelation is likely.
 - c.** There are no leveraged outliers in the dataset.
 - d.** Leveraged observations produce a significant effect for Revenue in the model.
 - e.** Leveraged observations reduce the standard error of Revenue in the model.

Output below on this page comes from a model that drops three outliers in Revenue and adds an interaction between Financial Leverage and the variable Country.

SUMMARY OF FIT

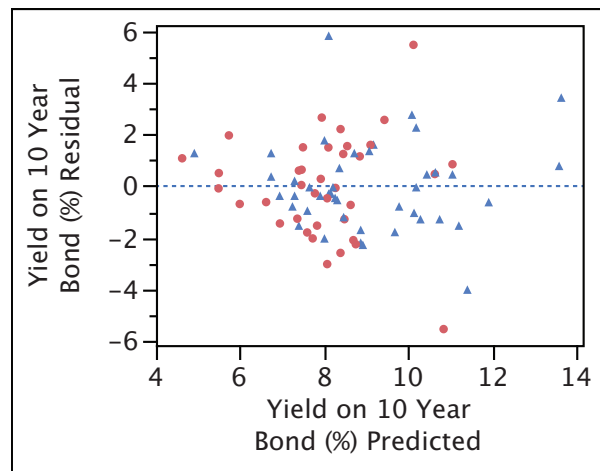
RSquare	0.465887
Root Mean Square Error	1.856047
Mean of Response	8.558659
Observations (or Sum Wgts)	82

ANALYSIS OF VARIANCE

Source	DF	Sum of Squares	Mean Square	F Ratio	Prob > F
Model	4	231.37523	57.8438	16.7911	
Error	77	265.25812	3.4449		
C. Total	81	496.63335			<.0001*

INDICATOR FUNCTION PARAMETERIZATION

Term	Estimate	Std Error	t Ratio	Prob> t
Intercept	7.3640367	0.545934	13.49	<.0001*
Revenue (mm)	-0.000115	3.748e-5	-3.08	0.0029*
Financial Leverage	0.469127	0.161105	2.91	0.0047*
Country[Malaysia]	0.3603138	0.640399	0.56	0.5753
Financial Leverage*country[Malaysia]	0.2875725	0.210909	1.36	0.1767



39. Which of the following features is the new model (that excludes the three leveraged points and adds an interaction) able to address that the original model (shown on page 171) could not?
- That the slope for Revenue depends on the county.
 - That yields in Singapore are lower than in Malaysia.
 - A differential impact of Financial Leverage on yields across countries.
 - Autocorrelation of the residuals due to the time series nature of the data.
 - Collinearity between Revenue and Financial Leverage.

- 40.** Comparing the new model (that excludes the three leveraged points and adds an interaction) to the original model, a fair interpretation of the results is that
- a.** The interaction term adds significant explanatory power to the model.
 - b.** The interaction term removes collinearity from the model.
 - c.** The Country term should be removed from the model since not significant.
 - d.** The removal of the three outliers reveals the importance of Revenue.
 - e.** The removal of the three outliers improperly inflates R^2 .
- 41.** Assuming the MRM holds, what does the p-value for Revenue in the interaction model tell you?
- a.** The Revenue variable should be removed from the model.
 - b.** The probability that the true Revenue partial slope is one is less than 0.0001.
 - c.** If the true Revenue partial slope were equal to zero, then it is extremely unlikely that we would have observed an estimate so far from zero.
 - d.** That the confidence interval for the partial slope of Revenue contains zero.
 - e.** The addition of Revenue to a model containing the others adds little to R^2 .
- 42.** Based on the new fitted model above with the interaction term, the estimated yield on a bond for a Singaporean company with revenues of \$25 (mm) and a Financial Leverage of 2 is approximately
- a.** 8.3
 - b.** 8.9
 - c.** 8.6
 - d.** 9.2
 - e.** 0.93
- 43.** If a company with the characteristics presented in the prior question (#42) issued a bond and its yield was less than 0.5 below the predicted yield, then based on the fitted interaction model would you be surprised?
- a.** No, as this event is almost certain to happen.
 - b.** No, as this event lies within the statistical uncertainty of the model.
 - c.** No, as the high value of R^2 makes this likely with 95% confidence.
 - d.** Yes, because the large sample size makes this unlikely.
 - e.** Yes, as this yield is more than twice the RMSE from the fitted value.

Sample Exam Answers STAT 613

- | | |
|--------|--------|
| 1. d. | 23. b. |
| 2. b. | 24. a. |
| 3. e. | 25. d. |
| 4. e. | 26. c. |
| 5. c. | 27. d. |
| 6. e. | 28. c. |
| 7. b. | 29. e. |
| 8. d. | 30. e. |
| 9. d. | 31. d. |
| 10. c. | 32. c. |
| 11. b. | 33. d. |
| 12. e. | 34. a. |
| 13. a. | 35. c. |
| 14. a. | 36. d. |
| 15. a. | 37. b. |
| 16. c. | 38. e. |
| 17. a. | 39. c. |
| 18. b. | 40. d. |
| 19. c. | 41. c. |
| 20. e. | 42. a. |
| 21. c. | 43. b. |
| 22. e. | |