

**THIS IS A PRACTICE MIDTERM EXAM -- YOUR EXAM WILL BE DIFFERENT.**

Online STAT 705  
Sample Midterm Exam

Name \_\_\_\_\_

**Question 1**

An experiment was conducted to explore the relationship between the temperature (in degrees) at which a chemical reaction occurs and the yield (in ppm) of the reaction.

Use the provided SAS output to answer the following questions. The SAS output contains information about two fitted models:

- a linear model  $(Y_i = \beta_0 + \beta_1 x_{1i} + \varepsilon_i)$ , and
  - a quadratic model  $(Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{1i}^2 + \varepsilon_i)$ .
- a. (5 points) Based on the results in the SAS output, does it appear that the linear regression model is appropriate for these data? Briefly explain.
- b. (5 points) Use the estimated linear regression equation to estimate the yield when the temperature is 110 degrees. (round your answer to 3 decimal places)
- c. (5 points) Use the estimated quadratic regression equation to estimate the yield when the temperature is 110 degrees. (round your answer to 3 decimal places)
- d. (5 points) Which of these two estimates is more reliable? Briefly explain.

**Question 2**

In a project to study age and growth characteristics of selected mussel species, researchers measured the age (in years) and weight (in pounds) of numerous mussels at three distinct locations (A, B and C).

- a. (10 points) Write the complete model specification for an interaction model. Include all subscripts and define all the variables.

**Use the SAS output to answer the following questions.**

- b. (5 points) Would an additive model be appropriate for these data? Explain.
  
  
  
  
  
  
  
  
  
  
- c. (5 points) Write the estimated equation for location B. (simplify, and round to 3 decimal places)

**Question 2, continued**

- d. (5 points) In the parameter estimates table, look at the line for 'location A'. The test statistic is 0.46, with p-value 0.6464. Does this mean "location A" can be removed from the model? Briefly explain.

- e. (5 points) Consider this table in the SAS output:

Source	DF	Type III SS	Mean Square	F Value	Pr > F
location	2	1.1045181	0.5522591	0.64	0.5315
age	1	379.9989681	379.9989681	438.34	<.0001
age*location	2	41.8750322	20.9375161	24.15	<.0001

In the first line, the test statistic is 0.64 and the p-value is 0.5315. Write the null and alternative hypotheses being testing, using the notation you defined in part (a).

- f. (5 points) According to the assumptions for a linear model, the errors are independent, normally distributed with mean 0 and common variance  $\sigma^2$ . What is the estimate for  $\sigma^2$ ?

**Question 3.**

Consider an experiment in which the researcher wants to determine a relationship between the seal strength of a bread wrapper stock (Y) and three predictor variables: sealing temperature ( $x_1$ ), cooling bar temperature ( $x_2$ ), and percent polyethylene in the stock ( $x_3$ ). The researcher felt that it might be appropriate to include interaction terms and terms that are quadratic in the predictors, so he generates the model

$$\text{Model 1: } Y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + \beta_4 x_{1i}^2 + \beta_5 x_{2i}^2 + \beta_6 x_{3i}^2 + \beta_7 x_{1i} x_{2i} + \beta_8 x_{1i} x_{3i} + \beta_9 x_{2i} x_{3i} + \varepsilon_i$$

After fitting the model and reviewing the SAS output, the researcher believes that he can remove all three of the interaction terms ( $x_1 x_2$ ,  $x_1 x_3$ , and  $x_2 x_3$ ). He removes these terms and fits the new model:

$$\text{Model 2: } Y_i = \tau_0 + \tau_1 x_{1i} + \tau_2 x_{2i} + \tau_3 x_{3i} + \tau_4 x_{1i}^2 + \tau_5 x_{2i}^2 + \tau_6 x_{3i}^2 + \varepsilon_i.$$

Now the researcher must choose which model he should use to complete his analysis.

Use the provided SAS code and output to conduct a nested model F test for these hypotheses:

$$H_0 : \beta_7 = \beta_8 = \beta_9 = 0$$

$$H_a : \text{at least one of } \beta_7, \beta_8, \text{ and } \beta_9 \text{ is not } 0$$

Please show your work in a logical fashion on the next page, and answer the specific questions below.

- (3 points) What is the value of the test statistic?
- (3 points) What is the critical value for this test? (use  $\alpha = 0.05$ )
- (3 points) Which hypothesis ( $H_0$  or  $H_a$ ) do you decide is true?
- (3 points) Based on the available information, which model (Model 1 or Model 2) should the researcher use?
- (3 points) What other information would you like to see about these two models before you choose a model?

**THIS IS A PRACTICE MIDTERM EXAM -- YOUR EXAM WILL BE DIFFERENT.**

Use this page to show your work for Question 3.

**Question 4.**

Obesity is a common, serious and costly disease. In 2010, Mayor Bloomberg of New York City supported a law restricting the sale of sugary drinks (Coke, Pepsi, etc.) because he believes that consuming sugary drinks contributes to the rate of obesity, so reducing the consumption of sugary drinks will cause a reduction in the rate of obesity. To support his position, Mayor Bloomberg surveyed residents in various neighborhoods across the city. The analysis consists of describing the relationship between two variables:

X = percent of adults in the neighborhood who drink at least one sugary drink per day, and

Y = percent of adults in the neighborhood who are obese.

Use the SAS output to answer the following questions.

- (5 points) How many neighborhoods are in the sample?
- (5 points) Since there is a relatively strong correlation between X and Y (the sample correlation is 0.76), Mayor Bloomberg argued that a reduction in X will cause a reduction in Y. Is this a valid argument? Briefly explain.
- (5 points) Suppose another neighborhood in New York City has 5 percent of adults that drink at least one sugary drink per day. Why would it be inappropriate to use the results of this analysis to estimate the percent obese for this neighborhood?

**Question 5.**

An assistant in the district sales office of a national cosmetics company is analyzing data on sales and expenditures in several of the district's territories. The data contain four variables:

Y = sales (in thousands of cases)

X1 = expenditures for point-of-sale displays in department stores

X2 = expenditures for local media advertising

X3 = expenditures for national media advertising

$X_1$ ,  $X_2$  and  $X_3$  are in thousands of dollars.

A multiple linear regression model was fit to the data, and the results are shown in the SAS output.

(Note: None of the variables have been transformed.)

- a. (5 points) Interpret the slope on local media expenditures.
- b. (5 points) Does it appear that there are outliers in the data? Explain.
- c. (5 points) Does it appear that multicollinearity is present? Explain.