# Test 3

## STAT 021

## Swarthmore College

**Do not flip this page until instructed to do so.**

**Test organization:** There are 12 questions in total on this test and they are organized into three subsections: the first 4 questions are matching or True/False with explanation questions, the next 5 questions are free response short answer and should not require more than a sentence or two to answer. The last section contains 3 long answer free response questions that require more than a couple of sentences to answer fully. There are a total of 60 points possible on this test. The last section explains an extra credit opportunity. If you need additional scratch paper you may come to the front of the class and pick some up.

**Instructions:** Answer each question to the best of your ability and raise your hand if you are confused by any of the wording in the questions or suspect a typo. For the short and long answer questions show all your work and provide enough justification and/or explanation in order to get full credit or to be considered for partial credit. You do not need a calculator to evaluate any expressions. For any calculation problems, simply writing out the formula to find the answer will suffice.

**First and Last Name:** _William Vespde_

**Swarthmore Username:** _wvespol1_

**Take a deep breath.**

You have prepared for this test and with a clear and well-rested mind, you are ready to show me what you have learned this semester. The purpose of this test is to measure your understanding of the material we have covered this semester. This is nothing more than a metric for me to evaluate your preparedness to think statistically at this particular moment in time and in this particular setting. This is not a perfect measure of your knowledge and does not predict your future statistical skills.

# Section 1: Matching and True/False problems

**1.** (5 points)

Suppose we are modeling the weight of birds (in $kg$) as a linear function of a categorical predictor variable for bird type (with levels pigeon, sparrow, and finch) and a numeric predictor for bird age. Given a "full" model

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_1 x_3 + \beta_5 x_2 x_3 + \epsilon,$$

where $x_1 = \begin{cases} 1, & \text{if sparrow} \\ 0, & \text{otherwise} \end{cases}$, $x_2 = \begin{cases} 1, & \text{if finch} \\ 0, & \text{otherwise} \end{cases}$ and $x_3$ is the age of the bird (in months), match the questions below to their corresponding null hypotheses.

a) For newly hatched birds (of age zero months), is there a statistically discernible difference in the weights of these three different bird types?

b) Does the effect of age on a bird's weight depend on what type of bird it is?

c) Given we are only comparing birds of the same age, is there a statistically significant difference in the mean weight of sparrows and pigeons?

d) Given we are only comparing pigeons, is the effect of age on a bird's weight statistically significant?

e) Is there statistically discernible evidence of a linear relationship between bird age and type and bird weight?

1._____ C _____ $H_0 : \beta_1 = 0$

2._____ a. _____ $H_0 : \beta_1 = \beta_2 = 0$

3._____ d _____ $H_0 : \beta_3 = 0$

4._____ b _____ $H_0 : \beta_4 = \beta_5 = 0$

5._____ e _____ $H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4 = \beta_5 = 0$

**2.** (5 points)

Determine which of the following statements about MLR models are true and false. For each statement that is false, provide a brief explanation as to why it is false.

(a) If predictors are collinear, then removing one variable will have no influence on the point estimate of another variable's coefficient.

*False, because those variables have the same effect on the response, removing one may lessen the overall effect that both had on the predictor, or it may increase if there is omitted variable bias.*

2

(b) If a regression model's first variable has a coefficient of $\hat{\beta}_1 = 5.7$, then if we are able to influence the data so that an observation will have a value of $x_1$ be one unit larger than it was before, the value of $y_1$ for this observation would increase by 5.7 units.

*False, because $x_1$ may be included in other spots in the regression models, such as an interaction term, so the effect is ambiguous.*

(c) As the total sample size increases, the degrees of freedom for the residuals increases as well.

*True*

**3.** (5 points)

Determine which of the following statements about ANOVA models are true and false. For each statement that is false, provide a brief explanation as to why it is false.

If the null hypothesis that the means of four groups are all the same is rejected from an ANOVA model and overall F-test at a 5% significance level, then ...

(a) We can then conclude that all the means are different from one another.

*False, we just know that one mean is different*

(b) The standardized variability among the group averages is higher than the estimate of the variability of the data within each group.

*True,*

(c) A post-hoc pairwise analysis will identify if there is at least one pair of means that are significantly different.

*True*

**4.** (5 points)

Determine if the following statements about statistical modeling are true or false, and explain your reasoning. If false, state how it could be corrected.

(a) Decreasing the significance level ($\alpha$) will increase the probability of making a Type 1 Error. *False,*

(b) With large sample sizes, even small differences between the null value and the observed point estimate will be identified as statistically significant.

*False,*

(c) Correlation is a measure of the association between any two variables.

*True*

3

# Section 2: Short answer questions

**5.** (4 points)

Briefly describe a benefit of analyzing the studentized residuals of a regression model rather than just analyzing the observed residuals.

The studentized residuals allows us to easier examine the influence of data points because it deletes the observation from the regression altogether, which affects the model. The observed residuals do not have this benefit.

**6.** (3 points)

If you could only use one measure (among the studentized residuals, leverage values, and Cook's distance values) to identify potentially influential data points, which would you choose and why?

I would use leverage values because it measures the influence of x-values rather than y-values. Therefore, if we see that an observation has an extreme leverage value, we can observe the corresponding y-value and make a decision to delete the observation. The other two measures only give information about the response variables from the residuals.

4

For questions 7-9 consider the following random single-serving samples of $n = 76$ breakfast cereals. We are going to model the average calories per serving (in g) (`calories`) as a linear function of the cereal manufacturer (a categorical variable with levels: G=General Mills, K=Kelloggs, N=Nabisco, P=Post, Q=Quaker Oats, R=Ralston Purina). Below is the R summary output for this one-way ANOVA model.

```
##
## Call:
## lm(formula = calories ~ Manufacturer, data = cereal_dat)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -58.696   -8.696   -0.126    5.909   51.304
##
## Coefficients:
##                Estimate Std. Error t value Pr(>|t|)
## (Intercept)     111.364      3.959  28.126  < 2e-16 ***
## ManufacturerK    -2.668      5.538  -0.482  0.63149
## ManufacturerN   -24.697      8.553  -2.887  0.00516 **
## ManufacturerP    -2.475      7.348  -0.337  0.73729
## ManufacturerQ   -16.364      7.667  -2.134  0.03633 *
## ManufacturerR     3.636      7.667   0.474  0.63678
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 18.57 on 70 degrees of freedom
## Multiple R-squared:  0.1618, Adjusted R-squared:  0.102
## F-statistic: 2.703 on 5 and 70 DF,  p-value: 0.02724
```

**7.** (3 points)

(a) What are the error degrees of freedom based on this model?

(b) What is the reference level?

a. 70 df

b. General Mills

**8.** (6 points)

Suppose the average amount of calories for all these samples is 106.97 over all 76 data points. What is the estimated group effect for Quaker Oats cereal brand?

The estimated group effect Quaker oats would be $106.97 - 16.364$, which is roughly 91.6, calories. Therefore, this tells us that on average Quaker Oats had 16.364 fewer calories than the total average.

**9.** (4 points)

Consider two additional numeric predictors: sugars (in g) and protein (in g). If we were to fit a regression model including each of the three predictor variables (including manufacturer) and an interaction between the two numeric variables, explain the meaning of the coefficient for the interaction term within the context of this data. (You should be able to answer this in no more than two sentences.)

The coefficient of the interaction term captures the effect of sugars or protein while the value of the other numeric predictor is allowed to change. We would want to include this term if we suspect a meaningful relationship between the variables.

## Section 3: Long answer questions

**10.** (9 points)

Suppose you have access to a data set on a random sample of undergraduate-only institutions in the US. The variables included in this data set are a numeric variable for the average cost of tuition each semester, a binary categorical variable distinguishing private institutions from public ones, a numeric variable for the percentage of full-time instructional staff employed at the institution, and a categorical variable indicating whether the school is a liberal arts college, a community college, a technical/vocational school, or if they are institutionally affiliated with certain groups (e.g. historically Black, women's only, tribal, etc).

State a research question that can be answered with the overall F-test for each of the following models, based on this data. (You do not need to use every variable, but you can.) Also provide a mathematical representation of the model and state the null hypothesis based on the notation you define for each model.

(a) a simple linear regression model;

(b) an ANOVA model;

(c) a multiple linear regression model (not **SLR** or **ANOVA**).

a. Does the percentage of full-time instructional staff employed at an institution affect the average cost of tuition for each semester?

$$Y = \beta_0 + \beta_1 X + \epsilon \text{, where } X = \% \text{ of full time staff}$$

$$H_0 : \beta_1 = 0$$

b. Is there a statistically significant difference in the average tuition each semester between different types of institutions (i.e. liberal arts, community, etc)?

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon \text{, where reference level} = \text{liberal arts, } X_1 \text{ is indicator var for}$$
community college, $X_2$ for technical, $X_3$ for affiliated w/ group

$$H_0 : \mu_0 = \mu_1 = \mu_2 = \mu_3 \text{ (group averages)}$$

c. Do the percentage of full-time staff or type of college (public/private) have a significant effect on average tuition price?

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon \text{, where } X_1 = \% \text{ of ft staff, } X_2 = \text{indicator var for type}$$
$$\begin{cases} 0, \text{ if private} \\ 1, \text{ if public} \end{cases}$$

$$H_0 : \beta_1 = \beta_2 = 0$$

**11.** (8 points)

Consider the ANOVA model for the cereal data you used in questions 7-8. Reference the R output on pg 5 and the plots on pg 10 to answer the following questions about this model.

(a) Check the conditions necessary for conducting a test to determine if the average calories (per serving) is significantly different for these six different cereal manufacturers. (You do not need to check the zero mean or linearity conditions but you do need to describe what it means for the group effects to be constant in this context.)

(b) Write out in words and in symbols the hypotheses that would be tested in part (a). (Clearly define your notation.)

(c) What can you conclude about the test in part (b)? Write a paragraph discussing your conclusions and reference any relevant statistics and/or plots as part of your discussion.

a. The residual plot tells us that the constant variance assumption is met because there are points evenly scattered around the zero line and display no clear pattern. There are more outliers for Kelloggs however. Overall, since the variance is constant, our test will be reliable because some groups aren't more variable than others.

b. $H_0 : \mu_0 = \mu_1 = \mu_2 = \mu_3 = \mu_4 = \mu_5$

$H_a : \mu_i \neq \mu_1, \mu_2, \mu_3, \mu_4, \mu_5$

The null hypothesis states that the group averages for calories is the same for all groups (i.e. brand of cereal). The alternative hypothesis states that at least one group mean is not equal to the rest of the group means.

c. The ANOVA test returns a p-value of 0.02724 and f-statistic of 2.703, which is significant at the 0.05 level. We can reject the null hypothesis and conclude that at the average of amount of calories in a brand of cereal is significantly different than the other brands tested. However, the effect was only significant at the 0.05 level for Nabisco. Nabisco, on average, has 24.697 fewer calories than the average for the reference level brand, General Mills

8

**12.** (8 points)

Suppose two people are studying the historic data set about the amount of arsenic (Arsenic) in local wells. This data contains $n = 70$ observations from a random selection of well water samples from across the state. In addition to the levels of arsenic, the data also records the year the data was collected (Year) and the distance from the well to the nearest mining site (Miles).

**Person A** fits the following MLR model to the data:

$$Arsenic = \beta_0 + \beta_1 Year + \beta_2 Miles + \epsilon$$

and computes an adjusted $R^2$ value of 0.26.

**Person B** considers the following correlations:

$$Cor(Arsenic, Year) = \rho_1; \quad Cor(Arsenic, Miles) = \rho_2$$

and estimates each with their sample correlations $r_1 = 0.77$ and $r_2 = -0.34$. Are the two people's conclusions contradictory? Explain your answer.
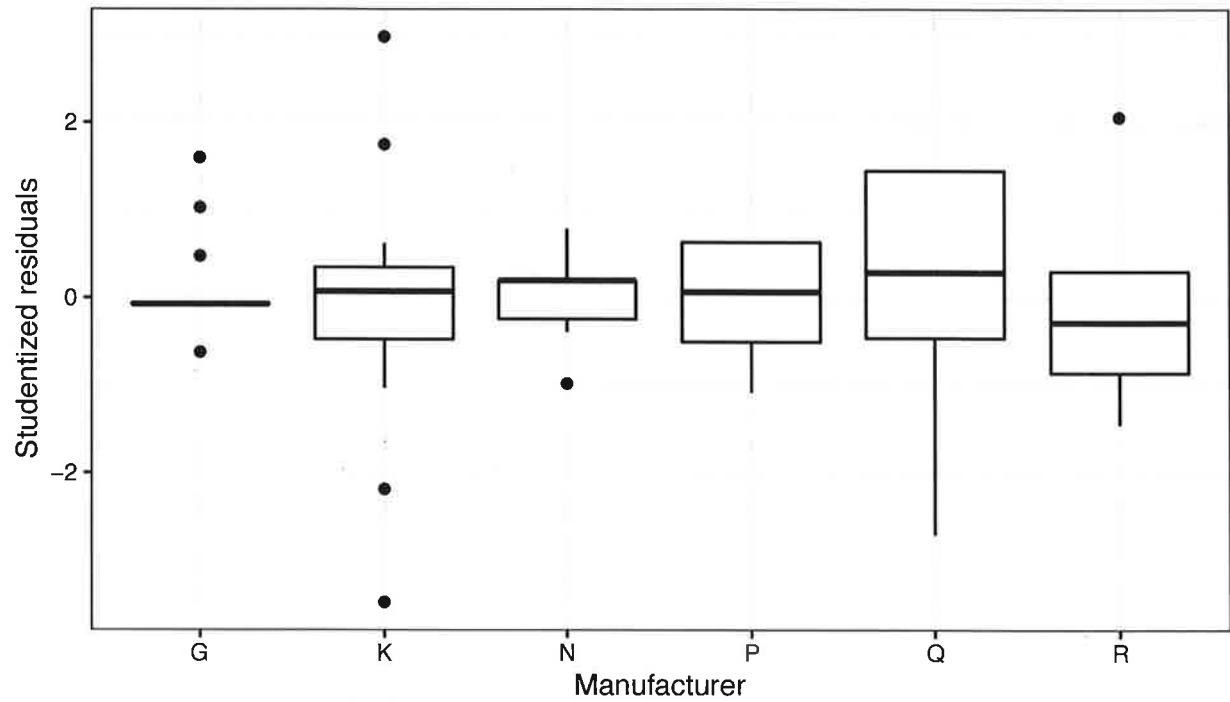
No, not necessarily. The $R^2$ is the amount of variability in the response that is explained to both predictors. Therefore, Year and Miles can both explain the variability in Arsenic, regardless of the direction of their impact. $r_1$ is negative because miles and Arsenic are inversely correlated while Year and Arsenic are positively correlated. In the MLR, the $\beta_1$ coefficient is likely positive while $\beta_2$ is likely negative.

# Section 4: Extra credit opportunity

If the response rate to my end of the semester evaluation form (on Moodle under Week 13 and 14) is at least 85% of our class size (over both sections), two percentage points will be added to everyone's Test 3 grade (up to 100 total possible points). **Hint:** You may not know how to or want to contact everyone in my class but you do know your group mates pretty well.

9

# Cereal ANOVA Model

## Residual plot for ANOVA model



## Normal quantile plot for ANOVA model