Statistical Methods - Math 222

- 1. In each case, state the specific statistical procedure that is appropriate for the given situation. Be specific: identify the response variable and the explanatory variable(s). If there are any categorical variables present, state how many levels each categorical variable has.
 - (a) You want to study whether men and women get different average amounts of sleep at night.

Solution: Two-sample t-test. The explanatory variable is gender, and the response variable is hours of sleep at night.

(b) You want to predict life satisfaction based on several factors, including income, regional cost of living, commuting time, and number of children.

Solution: This is multiple regression. The response variable is life satisfaction, and the explanatory variables are income, cost of living, commuting time, and number of children.

(c) You want to determine if there are significant differences between the cost of housing in three different cities. You also look at differences in costs of condominiums, versus townhomes, versus stand-alone houses.

Solution: This is 2-way ANOVA. The response variable is cost of housing. The explanatory variables are city (3 levels) and type of house (3 levels).

2. The scatterplot below shows the relationship between size (in square feet) and price (in thousands of dollars) of a random sample of 20 houses sold recently in Arroyo Grande, CA.



Below is a summary of the least squares regression model for this scatterplot.

Coefficients: Estimate Std. Error t value Pr(>|t|) (Intercept) 265.22212 42.64202 6.220 7.21e-06 *** myData\$Size 0.16859 0.03188 5.288 5.00e-05 *** ---Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1 Residual standard error: 51.31 on 18 degrees of freedom Multiple R-squared: 0.6084,Adjusted R-squared: 0.5866

(a) Is the trend statistically significant? How can you tell?

Solution: The trend is significant because the p-value for the slope is 5.00 times 10^{-5} .

(b) If $SE_{\hat{\mu}} = 55.18$, find a 95% confidence interval for the mean home price of a 1200 square foot house.

Solution: The confidence interval for $\hat{\mu}$ is $\hat{y} \pm t^* SE_{\hat{\mu}}$ where t^* has N-2 = 18 degrees of freedom. Use the t-distribution chart, $t^* = 2.101$. Also, $\hat{y} = 265.22 + 0.16859(1200) = 467.5$. So the confidence interval is: 467.5 ± 115.9 or equivalently: 351.6 to 583.4 thousand dollars

(c) Find a 95% prediction interval for the price of a 1200 square foot house (recall that $SE_{\hat{y}}^2 = SE_{\hat{\mu}}^2 + s^2$ where s is the residual standard error).

Solution: The standard error in \hat{y} is $\sqrt{55.18^2 + s^2}$ where s^2 is the variance of the residuals and is equal to the square of the residual standard error which is 51.31 in the chart above. So $SE_{\hat{y}} = \sqrt{55.18^2 + 51.31^2} = 75.35$. Then $\hat{y} \pm t^*SE_{\hat{y}}$ is 467.5 ± 158.3 which is from 309.2 to 625.8 thousand dollars.

3. This example is based on data from 78 seventh-grade students in a rural midwestern school. The researcher was interested in the relationship between the students' "self-concept" and their academic performance. The data included each student's grade point average (GPA), score on a standard IQ test, and gender, taken from school records. Gender is coded as 1 for female and 2 for male. The final variable is each student's score on the Piers-Harris Children's Self-Concept Scale, a psychological test administered by the researcher. Below is a summary of the multiple linear regression model for this data in R.

Call: lm(formula = gpa ~ iq + gender + concept, data = myData) Residuals: Min 1Q Median 3Q Max -3.5769 -0.7493 0.1984 0.9577 2.4089 Coefficients: Estimate Std. Error t value Pr(>|t|)

```
(Intercept) -2.83463
                        1.28584 -2.205 0.030641 *
iq
             0.08079
                        0.01336
                                   6.045 5.78e-08 ***
            -0.82214
                        0.31354
                                  -2.622 0.010630 *
gender
concept
             0.05048
                        0.01396
                                   3.616 0.000548 ***
Signif. codes:
                0 *** 0.001 ** 0.01 * 0.05 . 0.1
                                                    1
Residual standard error: 1.323 on 73 degrees of freedom
```

(1 observation deleted due to missingness) Multiple R-squared: 0.561,Adjusted R-squared: 0.543 F-statistic: 31.1 on 3 and 73 DF, p-value: 4.643e-13

(a) What is the formula for predicting GPA from IQ, Gender, and Self-Concept using this regression model?

Solution:

 ${\rm GPA} = -2.83463 + 0.08079\,{\rm IQ} - 0.82214\,{\rm Gender} + 0.05048\,{\rm Self-Concept}$

(b) What percent of the variability in GPA is explained by this model?

Solution: The $R^2 = 56.1\%$.

(c) Describe in words the affect of each of the explanatory variables on the response variable.

Solution: IQ has a significant positive effect on GPA. For each extra IQ point, GPA tends to go up 0.08. Gender has a negative effect on GPA which means that male students tend to have a GPA 0.82 points lower than female students, on average. Self-concept has a significant positive effect of 0.05 higher GPA for each extra point on the Piers-Harris scale.

4. A study looked at how lack of sleep affects reaction times. Volunteers were randomly assigned to either complete a task one hour after waking up or after 24 hours without sleep. Reaction times were measured (in milliseconds) in a discrimination task. Three levels of task difficulty were used. The results are shown in the interaction plot below.



Use this plot to answer the following questions.

(a) Describe clearly the main effects of each factor in this experiment.

Solution: The main effect of difficulty is the more difficult tasks have a longer reaction time, and the main effect of sleep is that lack of sleep increases reaction time.

(b) Describe any interaction between the factors.

Solution: The effect of lack of sleep appear to be larger for more difficult tasks than it is for easy or intermediate tasks.

(c) What should we do to determine if the interaction is statistically significant?

Solution: We should carry out a 2-way ANOVA test, and look at the F-value for the interaction term.

- 5. Determine whether each statement below is True or False.
 - (a) In one way ANOVA the response variable is categorical and the explanatory variable is quantitative.

Solution: False. The response variable is quantitative and the explanatory variable is categorical.

(b) Linear regression assumes that the residuals are normally distributed.

Solution: True.

(c) One of the assumptions made in the application of the one-way ANOVA F test is homogeneity of variance (i.e., the variances for all populations are assumed to be the same).

Solution: True.

(d) If the data in each group is strongly right skewed, it is okay to do an ANOVA F-test as long as the sample sizes are large.

Solution: True.

(e) When testing differences between population means using the One-Way Analysis of Variance (ANOVA) statistical method, the region of rejection is always in the left tail of the F distribution.

Solution: False. The rejection region is the right tail of the F-distribution.

(f) In two factor factorial design, factors A and B are said to have interaction if the effect of factor A is dependent on the level of factor B.

Solution: True.

(g) If the null hypothesis is rejected when conducting a one-way ANOVA F-test, then there are statistically significant differences between all pairs of means.

Solution: False. Not all pairs have to have significant differences.

6. Suppose you are performing one-way ANOVA to test for a difference in means for 4 groups. Each group contains 10 individuals that are randomly selected from a large population. Before conducting the test, you conduct a quick power computation for a specific alternative hypothesis where $\mu_1 = 10$, $\mu_2 = 11$, $\mu_3 = 12$ and $\mu_4 = 13$. You need to estimate σ for the computation, and so you choose $\sigma = 3$, which seems reasonable. Would the power be larger, smaller, or about the same if the true σ was actually larger than 3?

Solution: The power would be larger if the variance were smaller.

7. Do people from different cultures experience emotions differently? One study designed to examine this question collected data from 410 college students from five different cultures. 9 The participants were asked to record, on a 1 (never) to 7 (always) scale, how much of the time they typically felt eight specific emotions. These were averaged to produce the global emotion score for each participant. Here is a summary of this measure:

Culture	n	\bar{x}	SD
European American	46	4.39	1.06
Asian American	33	4.35	1.18
Japanese	91	4.72	1.13
Indian	160	4.34	1.26
Hispanic American	80	5.04	1.16

	Df	SS	MS	F
Culture		31.268		
Residuals			1.4044	n/a
Total	409	600.04	14671	n/a

31.268

568.772

600.04

7.817

1.4044

1.4671

5.566

n/a

n/a

1	~)	Com	lata	+la a	ANO	578	table	halarr	for	there	magnita	h	£11:n m	:	+la a	6	minaina	antriage
L	aj	Ծմոր	nete	0116	ANU	vл	table	DEIOW	101	unese	resuits	Dy	mmig	111	une	nve	missing	enumes.

(b) What is are the null hypothesis and alternative hypothesis for this ANOVA test?

4

405

409

Culture

Total

Residuals

Solution:

Solution:

 H_0 : The mean is the same for all cultures.

 H_A : There are differences in the means.

(c) It turns out that the *p*-value for the F-statistic above is 2.27×10^{-4} . What does that mean in this situation?

Solution: We should reject the null hypothesis and conclude that there are statistically significant differences in the means.

(d) Is it reasonable to used a pooled standard deviation for these data? Why or why not?

Solution: Yes, because the sample standard deviations for each group are all very similar.

(e) Why don't we need to worry very much about whether the assumption of normality is met for this data?

Solution: The sample sizes are so large (the smallest is group has 30 people), so normality will not be a big concern.

(f) Recall that the confidence interval for the difference between the means of two groups is $\bar{x}_A - \bar{x}_B \pm t^{**} s_p \sqrt{\frac{1}{n_A} + \frac{1}{n_B}}$, where t^{**} is the adjusted critical value with the Bonferroni correction. According to the Bonferroni method, what adjusted confidence level should we use to be 95% certain that we capture the true difference in population mean for each pair of groups simultaneously? You don't need to compute the confidence interval.

Solution: There are ${}_{5}C_{2} = \frac{(5)(4)}{2} = 10$ different pairwise comparisons to consider. Therefore we need to use a 1 - 0.05/10 = 0.995 = 99.5% confidence level for each confidence interval.