

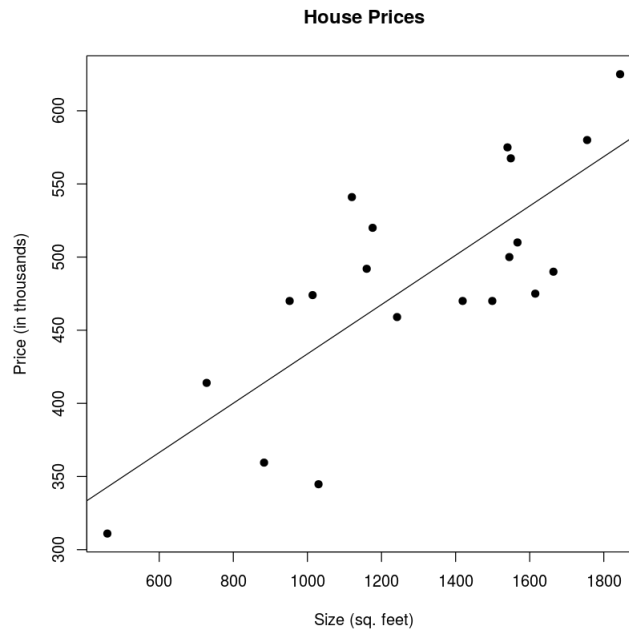
Statistical Methods - Math 222

Midterm 3 Review

This is not for a grade! This is just to test yourself and see what you know and what you need to study for the upcoming midterm.

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.
 - (b) You want to predict life satisfaction based on several factors, including income, regional 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.

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

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 ***
gender      -0.82214    0.31354   -2.622 0.010630 *
concept      0.05048    0.01396    3.616 0.000548 ***
---
Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1 1

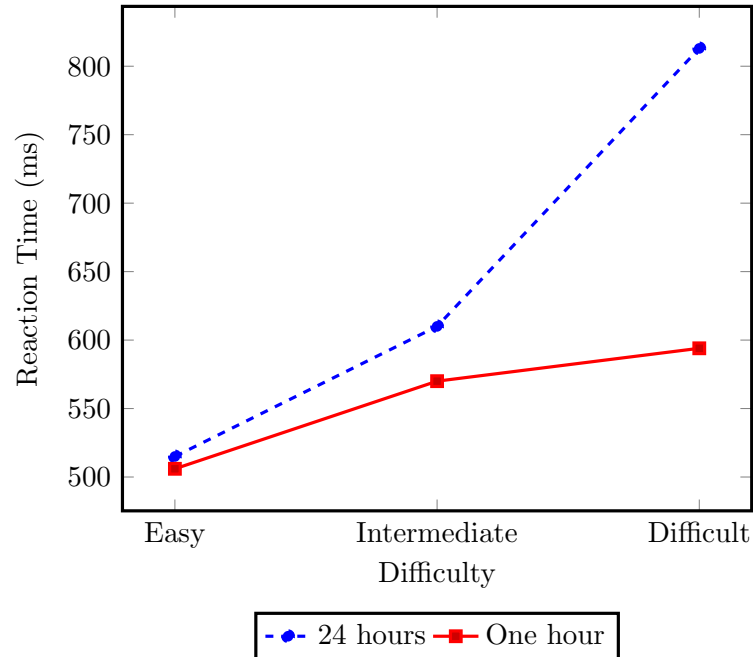
Residual standard error: 1.323 on 73 degrees of freedom
(1 observation deleted due to missingness)
Multiple R-squared:  0.561, Adjusted R-squared:  0.543
F-statistic: 31.1 on 3 and 73 DF, p-value: 4.643e-13
```

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

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

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

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.
- (b) Describe any interaction between the factors.
- (c) What should we do to determine if the interaction is statistically significant?

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.
 - (b) Linear regression assumes that the residuals are normally distributed.
 - (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).
 - (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.
 - (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.
 - (f) In two factor factorial design, factors A and B are said to have interaction if the effect on factor A is dependent on the level of factor B.
 - (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
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?

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

- (a) Complete the ANOVA table below for these results by filling in the five missing entries:

	Df	SS	MS	F
Culture		31.268		
Residuals			1.4044	n/a
Total	409	600.04	1.4671	n/a

- (b) What is are the null hypothesis and alternative hypothesis for this ANOVA test?
- (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?
- (d) Is it reasonable to used a pooled standard deviation for these data? Why or why not?
- (e) Why don't we need to worry very much about whether the assumption of normality is met for this data?
- (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.*