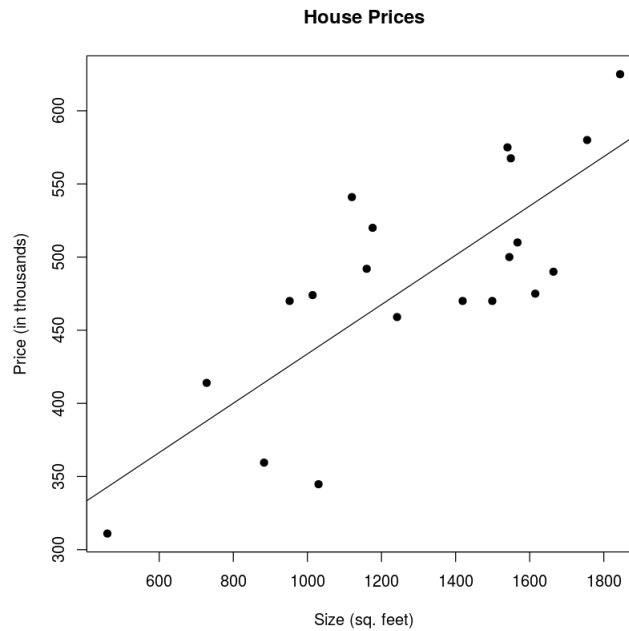


3. The scatterplot below shows the relationship between size (in square feet) and price (in thousands of dollars) of a random sample of 20 houses sold recently in Arroyo Grande, CA.



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).

- (d) Use the fact that $R^2 = 0.6084$ and $s_y = \$79,801.5$ and $n = 20$ houses to fill in the following ANOVA table for this example.

	DF	SS	MS	F
Size				
Residuals				
Total				

4. 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) on a ten-point scale, 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) Estimate the GPA of a male student who has an IQ of 110, and a self-concept score of 60.

5. 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. 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 are the null and alternative hypotheses 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) Why is it reasonable to assume that each group has the same population standard deviation in this example?
- (e) What number is the best estimate for the common standard deviation for each group?
- (f) Why don't we need to worry very much about whether the assumption of normality is met for this data?

- (g) 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 given that there are 10 possible pairs of groups to compare? *You don't need to compute the confidence interval.*

6. 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 multilinear regression models, removing variables always decreases the adjusted R^2 .
- (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.

7. Esophageal cancer can spread to the lymph nodes, and the larger the tumor is, the more likely it is to spread. Below is the R output for a logistic regression model based on a sample of 31 cancer patients. The explanatory variable is the size of the tumor in centimeters and the response variable is whether or not the cancer has spread to the lymph nodes.

```
## Call: glm(formula = spread ~ size, family = "binomial", data = cancer)
##
## Coefficients:
## (Intercept)      size
##      -2.086      0.5117
```

- (a) What is the formula for the log-odds in the model described above?
- (b) What are the odds that the cancer has spread to the lymph nodes if a patient has a 6 cm tumor?
- (c) What is the probability that the cancer has spread to the lymph nodes if a patient has a 6 cm tumor?
- (d) How many times higher are the odds of the cancer spreading for every extra centimeter of tumor?