# Definitions:

## Pearson's correlation coefficient:

The covariance of two variables divided by the product of their standard deviations.

## For a population:

$p_{x,y} = \frac{Cov(X, Y)}{\sigma_X \sigma_Y}$
where
$Cov(X, Y) = E[(X - E[X])(Y - E[Y])]$

## For a sample:

It's often referred to as the sample correlation coefficient, commonly abbreviated to just "r"

$$r = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n} (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^{n} (y_i - \bar{y})^2}}$$

(Above: the sample covariance divided by the product of the sample standard deviations)
which can be manipulated to get:

$$r = r_{xy} = \frac{1}{n-1} \sum_{i=1}^{n} \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

The correlation coefficient always takes a value between -1 and 1, with 1 or -1 indicating perfect correlation (all points would lie along a straight line in this case).

  a. A positive correlation indicates a positive association between the variables (increasing values in one variable correspond to increasing values in the other variable).

  b. A negative correlation indicates a negative association between the variables (increasing values is one variable correspond to decreasing values in the other variable).

  c. A correlation value close to 0 indicates no association between the variables.

The square of the correlation coefficient, $R^2$ , is a useful value in linear regression. This value represents the fraction of the variation in one variable that may be explained by the other variable. Thus, if a correlation of r $= 0.8$ is observed between two variables (say, height and weight, for example), then a linear regression model attempting to explain either variable in terms of the other variable will account for 64% ($r^2 = 0.8^2$ $= .64$) of the variability in the data.
The correlation coefficient also relates directly to the regression line Y $= a + bX$ for any two variables, where
b $= r\frac{s_x}{s_y}$
I found this info here:
http://www.stat.yale.edu/Courses/1997-98/101/correl.htm
and on the wikipedia page.

Future Notes

|            | Q1 | Q2 | Q3 | Q4 | Q5 |
| ---------- | -- | -- | -- | -- | -- |
| 50 Points  | 10 | 14 | 8  | 8  | 10 |

# Question 1

For multiple regression

$$y = X\beta + \epsilon, \ \epsilon \sim N(0, \sigma^2)$$

Derive or show that

a. $\hat{\beta} = (X'X)^{-1}X'Y$

b. $E[\hat{\beta}] = \beta$

c. $V[\hat{\beta}] = \sigma^2(X'X)^{-1}$

d. $E[\hat{Y}] = X\beta$

e. $V[\hat{Y}] = \sigma^2 H$, where H is the hat matrix and $H = X(X'X)^{-1}X'$

# Question 2 (problems 3.1 and 3.3 on page 121)

a. Fit a multiple linear regression model relating the number of games won to the team's passing yardage ($x_2$), the percentage of rushing plays ($x_7$), and the opponents' yards rushing ($x_8$).

b. Construct the analysis-of-variance table and test for significance of regression.

c. Calculate t statistics for testing the hypotheses $H_0$: $\beta_2 = 0$, $H_0$: $\beta_7 = 0$, $H_0$: $\beta_8 = 0$. What conclusions can you draw about the roles the variables $x_2$, $x_7$, and $x_8$ play in the model?

d. Calculate $R^2$ and $R^2_{adj}$ for this model.

e. Using the partial F test, determine the contribution of $x_7$ to the model. How is this partial F statistic related to the t test for $\beta_7$ calculated in part c above?

f. Find a 95% CI on $\beta_7$. (This is part a of problem 3.3, and the following one is part b of problem 3.3.)

g. Find a 95% CI on the mean number of games won by a team when $x_2 = 2300$, $x_7 = 56.0$, and $x_8 = 2100$.

Note: For c, d, f, and g, please show two versions of your results: (1) obtained using R code and (2) based on your manual calculation (please show detailed step for your manual calculation. You can use the partial output from the lm or ANOVA, e.g., the $SS_{reg}$, $SS_{res}$, the estimated value of $\beta$ and its variance or standard deviation).

# Question 3 (Exercise 3.4 on page 122

Reconsider the National Football League data from Problem 3.1. Fit a model to this data using only $x_7$ and $x_8$ as the regressors.

a. Test for significance of the regression.

   b. Calculate $R^2$ and $R^2{}_{adj}$. How do these quantities compare to the values computed for the model in problem 3.1, which included an additional regressor $(x^2)$?

   c. Calculate a 95% CI on $\beta_7$. Also, find a 95% CI on the mean number of games won by a team when $x_7 = 56.0$ and $x_8 = 2100$. Compare the lengths of these CIs to the lengths of the corresponding CIs from problem 3.3 (that is, the above part f and g in question 2)

   d. What conclusions can you draw from this problem about the consequences of omitting an important regressor from a model?

# Question 4 (exercise 4.2 on page 165

Consider the multiple regression model fit to the National Football League (NFL) team performance data in problem 3.1.

   a. Construct a normal probability plot of the residuals. Does there seem to be any problem with the normality assumption?

   b. Construct and interpret a plot of the residuals versus the predicted response.

   c. Construct plots of the residuals versus each of the regressor variables. Do these plots imply that the regressor is correctly specified?

   d. Construct the partial regression plots for this model. Compare the plots with the plots of residuals versus regressors from part c above. Discuss the type of information provided by these plots.

# Question 5

Show that the hat matrix $H = X(X'X)^{-1}X'$ and $I - H$ (where I is the identity matrix) are symmetric and idempotent. That is, please show:

   a. $H' = H$ and $HH = H$ ($H'$ means the transpose of H, HH means H * H)

   b. $(I - H)' = I - H$ and $(I - H)(I - H) = I - H$

Hint: $A = X'X$ is a symmetric matrix, and for a symmetric matrix, $(A')^{-1} = (A^{-1})'$. You can use this property directly in your proof of **(a)** and **(b)**. If you are interested in the proof of this property, you may check the following web page:

https://math.stackexchange.com/questions/325082/is-the-inverse-of-a-symmetric-matrix-also-symmetric