# Final Exam Review

## Format of the Exam

This final exam is cumulative. You will be tested on material from Unit 1-Unit 20 (up until section 2). Like the last two exams, you will be expected to read a dataset and be able to answer the following questions involving this dataset. There may be some conceptual questions as well. See the attached documents for specific question examples of this nature. The attached Jupyter notebook is not comprehensive, in that it does not cover ALL things talked about in this document and what might be on the test.

1. ## Cumulative Final Exam: Be sure to review questions/items in the Midterm 1 Review and Midterm 2 Reviews!

2. ## More About Logistic Regression
   2.1. What objective function are optimizing when you find the "optimal" values of the slopes and intercept in a logistic regression model.
   - **2.1.1.** **Conducting Inference on Logistic Regression Slopes**
     - 2.1.1.1. Using p-values and confidence intervals.
   - 2.1.2. **Log-likelihood function**
     - 2.1.2.1. What is the log likelihood function? Aka, write out the function that you are trying to optimize that gives you the "best" values of the slopes and intercept in a logistic regression model.
     - 2.1.2.2. How to extract the optimal log-likelihood function value of a model in Python.

3. ## Making Predictions with a Classifier: Making 0/1 Predictions with a Threshold and Assessing the Accuracy of the Predictions
   3.1. Definitions/Calculations
   - **3.1.1.** **Predictive Probability Threshold**
     - **3.1.1.1.** How to define.
     - **3.1.1.2.** How to classify a set of explanatory variables values given their corresponding predictive probabilities and a predictive probability threshold.
   - **3.1.2.** **True positives, true negatives, false positives, false negatives**
     - 3.1.2.1. How to define
     - 3.1.2.2. How to calculate by hand (given a set of true labels and predicted labels).
     - 3.1.2.3. How to calculate with Python function.
   - 3.1.3. **Confusion matrix**
     - 3.1.3.1. How to define
     - 3.1.3.2. How to calculate with Python function.
   - **3.1.4.** **Sensitivity, Specificity, True Positive Rate, False Positive Rate, and Accuracy**
     - 3.1.4.1. How to define
     - 3.1.4.2. How to calculate by hand (given a set of true labels and predicted labels).
     - 3.1.4.3. The relationship between each of these terms.
     - 3.1.4.4. How to calculate with Python function.
   - **3.1.5.** **ROC Curve and AUC (Area under Curve)**
     - 3.1.5.1. How to define each.
     - 3.1.5.2. The relationship between each of these terms.
     - 3.1.5.3. How to plot the ROC curve with Python (given a set of true labels and predicted labels).
     - 3.1.5.4. How to calculate the AUC with Python (given a set of true labels and predicted labels).

6.2.3.**Comparison**

      6.2.3.1.     For the same model, (which uses a sample size n>8), which statistic will penalize a larger number of slopes more: the AIC or the BIC of the model?

# 6.3. Backwards Elimination Algorithm

    6.3.1. How to use it to try to find the model with the best BIC or AIC score.