

Homework 2

1. Download the classification output data set (attached in Blackboard to the assignment).
2. The data set has three key columns we will use: `class`: the actual class for the observation `scored.class`: the predicted class for the observation (based on a threshold of 0.5) `scored.probability`: the predicted probability of success for the observation Use the `table()` function to get the raw confusion matrix for this scored dataset. Make sure you understand the output. In particular, do the rows represent the actual or predicted class? The columns?
3. Write a function that takes the data set as a dataframe, with actual and predicted classifications identified, and returns the accuracy of the predictions.
4. Write a function that takes the data set as a dataframe, with actual and predicted classifications identified, and returns the classification error rate of the predictions. Verify that you get an accuracy and an error rate that sums to one.
5. Write a function that takes the data set as a dataframe, with actual and predicted classifications identified, and returns the precision of the predictions.
6. Write a function that takes the data set as a dataframe, with actual and predicted classifications identified, and returns the sensitivity of the predictions. Sensitivity is also known as recall.
7. Write a function that takes the data set as a dataframe, with actual and predicted classifications identified, and returns the specificity of the predictions.
8. Write a function that takes the data set as a dataframe, with actual and predicted classifications identified, and returns the F1 score of the predictions.
9. Before we move on, let's consider a question that was asked: What are the bounds on the F1 score? Show that the F1 score will always be between 0 and 1. (Hint: If $0 < \text{precision} < 1$ and $0 < \text{recall} < 1$ then $\text{precision} \times \text{recall} < \text{precision}$.)
10. Write a function that generates an ROC curve from a data set with a true classification column (class in our example) and a probability column (scored.probability in our example). Your function should return a list that includes the plot of the ROC curve and a vector that contains the calculated area under the curve (AUC). Note that I recommend using a sequence of thresholds ranging from 0 to 1 at 0.01 intervals.
11. Use your created R functions and the provided classification output data set to produce all of the classification metrics discussed above.
12. Investigate the `caret` package. In particular, consider the functions `confusionMatrix`, `sensitivity`, and `specificity`. Apply the functions to the data set. How do the results compare with your own functions?
13. Investigate the `pROC` package. Use it to generate an ROC curve for the data set. How do the results compare with your own functions?