

# STATISTICS 452/652: Statistical Learning and Prediction

October 29, 2020

## PROJECT 1

Due Friday, Nov 13, 11pm, Submission details to be announced

### POLICY

1. This project is to be completed *independently*, with no outside help except for a named teammate, if you have been assigned one. You may use whatever class materials you wish in completing this assignment. **BUT DO NOT DISCUSS QUESTIONS OR RESULTS WITH ANYONE ELSE, WITHIN OR OUTSIDE OF THE CLASS.** Failure to follow this directive will result in a failing grade.
2. Late projects will be accepted at a penalty of 2 points/hour (it's a 100 point project), *strictly enforced*.

### ASSIGNMENT

The data for this project will be posted to Canvas at 5PM on November 6. This file will be a comma-separated file (csv) that contains  $n$  observations of  $p$  explanatory variables (labeled  $X_1$ – $X_p$ ) and one response ( $Y$ ). That's all I'm telling you about the data.

Your job is to develop a method for predicting  $Y$  based on  $X$ . You may use any of the techniques covered in class.

### DELIVERABLES

You will produce two required items and one optional item.

1. I will post a test set of explanatory variables without the response variable attached. You will return a list of predicted values, *in the same order*. The list should be *one column of numbers with no row numbers and no column header*. Use `write.table(predictions, file.name, sep = ",", row.names = F, col.names = F)` to submit your code, where `predictions` is the name of the vector containing your predicted values and `file.name` is the location on your computer where you want to

store the results. These will later be uploaded into a Shiny app that we will introduce later.

**Look at your file before you submit it** to make sure that the format is correct (and also to make sure that the predicted values are sensible!).

2. You will supply a written report of the steps you took to create your model and predictions. **This gets posted to Crowdmark.** Details are given below.

## REPORT

Your report, which is submitted to Crowdmark, should answer these questions, as numbered below:

1. *What models or machines did you attempt to fit?* For each one, paste the R code from your program for the initial successful model fit. I want to see what you tried. For example, `"fit1 = lm(y~., data=train)"` is what I would list if I used multiple linear regression on all of the variables and my training data were called "train". **The answer should be a list (e.g., with bullets) of nothing but the code for each of these model fits.** Don't list code that did not run. If tuning was involved in the initial fitting process, you can paste the function with variable names for the tuning parameters (e.g., your function might have `"mtry=mm"` if you looped over a variable called "mm").
2. *What process(es) did you use to evaluate and compare models and to select your final model?* I am thinking of Lecture 3, specifically: **Give 1-2 sentences explaining the method, the quantity, how results were turned into decisions.** For example, "I used 50,000 bootstrap resamples, fit all models to each resample, and used largest training error from last resample as my best model." (This example answer is complete, but represents something rather stupid to do...)
3. *Did you tune any methods?* If so, (a) what process(es) did you use to evaluate and compare models and to select your final model (i.e., **I want to see an answer like to the previous question, but relating to how TUNING was done**), and (b) **for each method list all parameter values that were considered** (e.g., "For "Blasting" I use a grid of values with  $A=(1, 2, 3, \dots, 60)$  and  $B=(0.00317, \sqrt{3.14159})$ ). For "Blooming" I used combinations of  $(z, \gamma)=(0.1, 3), (0.5, 6),$  and  $(1.1, 12)$  ).
4. *What was your chosen prediction machine?* **Paste the code that produced your predicted values, including all values of tuning parameters if any, random number seeds, and explaining any variable names that are not obvious.** I should be able to run your code and produce the same results (or extremely similar if randomization is used). If I try and it doesn't work, there will be a major deduction.
5. *(optional)* **List the variables that you believe are important.** A positive bonus will be given for each correct results. A deduction will be made for each variable listed that was not important.

The main thing here is that I should be able to see what your thought process was and whether you considered (or failed to consider) important ideas.

## GRADES

Your grade will be based partly on how well your model performs, and partly on the steps you took to get there. I will compute a form of  $R^2$  between your predicted values and the test set responses. *I will scale these against the best model produced by a member of the class, so this is a competition!* If your  $R^2$  is only 80% as large as the best, your mark for this part will be 80%.

Your report as described above. This portion will count for 60% of your grade. The remaining 40% will come from your model's performance.

If you supply a list of variables that are in the model, I will give a small bonus for each variable that you correctly identify. *I will subtract the same amount for each variable that you incorrectly list!*

## FINAL COMMENTS

GOOD LUCK, HAVE FUN, and remember: in real life an employer will take action based on the results you provide them. These may be million-dollar decisions which rely extensively on YOUR expertise. This is practice...