

ENGS/QBS 108 Fall 2017 Assignment 3

Due October 17, 2017

Instructors: George Cybenko and Saeed Hassanpour

Prepared by: Benjamin Priest

Problem: Porto Seguro Dataset [70 points]. In this problem, you will explore and attempt to solve a realistic machine learning problem hosted on the [Kaggle](#) competition platform. In order to download the data and access other competition resources, you will need to create a kaggle account. The [Porto Seguro dataset](#) includes a number of anonymized features of vehicle drivers collected over the course of a year, as well as a target flag indicating whether the driver filed an insurance claim during the year. Your task is to train models to predict driver safety using this data.

1. [10 points] Data analysis is an important first step. The first two fields, “id” and “target” are an index and the label, respectively. The other 57 fields are features with partially anonymized names. The elements of the “target” column are in $\{0, 1\}$, where 1 indicates the driver filed an insurance claim during the year. Proceed to explore this problem’s dataset by addressing the following:
 - (a) What are the dimensions of the training and testing datasets?
 - (b) Features ending in “_cat” are categorical, meaning that their entries are categories rather than numerals. How many categories does each such feature include?
 - (c) Features ending in “_bin” are binary, taking entries in $\{0, 1\}$. Are any of these features “one-hot” encodings of a single categorical feature? If so, which ones, and how many categories does it include?
 - (d) Explore the [discussion](#) page of the competition. What else can you learn about the features or dataset from what others have done? Include code if relevant.
2. [20 points] In this part, you will apply an SVM to solve a classification problem using the Porto Seguro dataset. You will design an SVM that reads a feature vector and predicts whether that driver will output a
 - (a) Separate the training data into a [k-fold cross validation](#) testing framework for $k=4$.
 - (b) Train SVM models using your validation framework. Try a few different kernels. What kernel yields the best performance? Note: Training will most likely take a *long* time.
 - (c) Apply your best model to the testing set. Report your classification accuracy.
3. [20 points] In this part, you will attempt proscribed a proscribed solution to the competition. You will build a (simple) deep neural net whose desired output is a probability that a driver will file a claim given the input feature vector. You will report results in terms of a [ROC curve](#), where you vary the decision threshold between 0 and 1.
 - (a) Construct a deep neural network with three dense layers of 100, 50, and 20 neurons, respectively, as well as a single sigmoidally activated output layer.
 - (b) Train your model using [cross-entropy loss](#). Report your ROC curve on the training data.
 - (c) Apply your model to the testing set. Report your ROC curve.

4. [20 points] In this part, you will attempt your own solution to the competition. You may use any model with which you are familiar. You are free to browse [existing kernels](#) for ideas, but your work should ultimately be your own. Report your results in terms of ROC curves.
- (a) Construct a model of your choice.
 - (b) Train your model and report your ROC curve on the training data.
 - (c) Apply your model to the testing set. Report your ROC curve.
 - (d) How does your solution compare to the solution in part 3? Explain the comparison.