

Final Project Report

By Robert Syvarth & William Headrick | RS42999 & WH5965

1. How did your team go about tackling this problem?

We started out by reading through a lot of the advice on the Kaggle forums. After reading through and getting some ideas we determined that it would be most beneficial if we wrote our solution without using any existing code as a template. This way we would have to understand everything that we were doing – not just blindly copying it from elsewhere on the Internet.

We went with an iterative approach starting out with a very simple solution that just involved a single classifier to prove that our programming was correct and over time adding on additional processing steps in order to improve upon this baseline.

2. Which methods/algorithms did you try?

We ended up trying a lot of different classifiers from Scikit early on. We tried the Bernoulli Naive Bayes classifier, LogisticRegression, AdaBoost, GradientBoost, SVC, and RandomForest. For each of these classifiers we did around a hundred test trains with varying parameters in order to evaluate the efficacy of that particular classifier on our data. We also determined that we needed a method for combining our predicted results so we researched various options there and experimented with simple linear combinations which were simple and performed pretty well. Finally we investigated various forms of feature combination and selection. Many of the existing solutions used feature combination in order to expand the dimensionality of the problem followed by greedy feature selection. In addition to these options we also looked into manual, heuristic-based, combination and Recursive Feature Elimination.

3. What is your final methodology? Walk through it in detail, starting from data pre-processing. Explain all the machine learning algorithms(s) you used as well as the parameters you chose. Also discuss any external tools or libraries that you used.

Our final solution uses a combination of three different classifiers as well as feature combination and Recursive Feature Elimination in order to predict employee access requirements.

First we use Numpy's loadtxt function in order to read CSV data from disk. This gives us easy access to the data without having to worry about details of CSV parsing. The first step we take once we have all of the data loaded is to expand the dimensionality of the data. In order to achieve this we consider all pairs of features in the original data as well as the original features themselves. This allows us to coax a little more subtlety out of some of our classifiers by selecting specific sets of combined features.

The next step is feature extraction. Once we have expanded our set of features we need to remove those which are redundant or otherwise unhelpful to our classifiers. In order to do this we use Scikit's Recursive Feature Elimination class. This allows us to pass a classifier and a dataset it returns a list of the most relevant N features in that classifier, in our case we ended up selecting 30 features in total. We came to this result through simply running our cross validation with different numbers of features until we reached a maximum. We used the LogisticRegression classifier for this step since, as we will discuss later, it accounts for the greatest weight in our final model. Another thing worth noting is that we scale all of our features from 0-1 prior to passing them to the feature elimination step. We include this step because otherwise features with larger absolute values are preferred. The alternative would be to attempt a one-hot encoding at this step but the number of features produced by a one-hot encoding prevent any reasonable feature selection from occurring after it has occurred.

After we have selected our features we simply remove those not in the top 30 from our original datasets. The next step is to perform a one-hot encoding on this newly reduced dataset so that it is ready for our classifiers. Once the one-hot encoding is complete we are ready to start training our classifiers. In the end we decided to use three models for the machine learning step of our project: LogisticRegression, RandomForest, and SVC. These were chosen from the list of classifiers which we evaluated due to their performance on our cross validation tests.

While choosing the best parameters for these models it became pretty clear that there were only a few which actually significantly impacted our results. The most important for both the LogisticRegression model and SVC is C. The C value essentially defines how much regularization is applied to the model in order to prevent overfitting. As you can see in **Figure 1** there were some pretty clear trends in our cross validation results based on the value of C both these classifiers. We ended up choosing values of 4.4 and 1.4 for these two models respectively based on these tests.

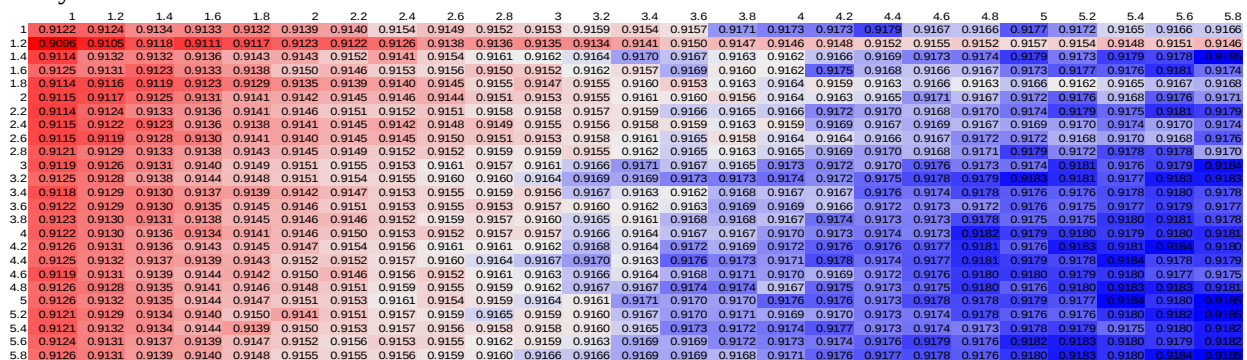


Figure 1. Heatmap of cross validation scores by Logistic Regression and SVC C values

After we had selected and tuned each of our classifiers we had to determine how to combine their predictions into a single set of values. Our solution was a simple weighted linear combination of our result values. We just needed to determine the weights. In order to do this we took a similar approach to parameter selection and simply tested out many different combinations in order to find a near optimal value, as you can see in **Figure 2**. In this case though we found that despite our cross validation we ended up with an overfit result, so we performed some manual adjustments and ended up using weights of 0.83 for LR, 0.12 for RF, and 0.05 of SVC.

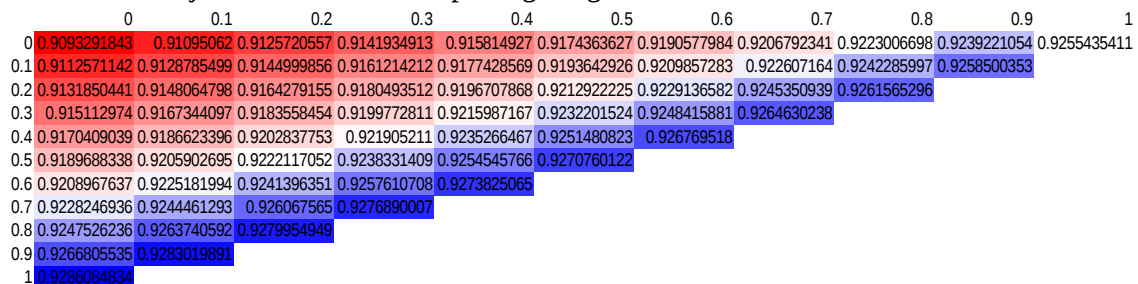


Figure 2. Heatmap of cross validations scores based on weights of individual models

The final step of our classifier is to then combine the predictions from each of our models using the weighted linear combination explained above. We then simple our final averaged results to a csv so that we can submit for grading on Kaggle. Using this approach we achieved a private score of 0.89177.

More detailed versions of the data provided here can be found in the Data Analysis.ods file