



relatively inexpensive to train which made it an ideal candidate for parameter tuning.

Finally, extreme gradient boosting required an acceptable amount of time in training since it supported GPU acceleration. Along with its comparatively strong performance, we decided to also use it for tuning alongside the Random Forest model.

2.2 Results Table

The table shows that the ensemble methods have an overall better performance than the other models in this dataset. We chose Extreme Gradient Boosting and Random Forest to refine their parameters given their lower test cross entropy and training time requirements.

Table 2.2 (test) Jaccard Index and (test) Cross Entropy

Algorithm	Parameters	(test) Jaccard Index	(test) Cross Entropy
2 Gradient Boosting	default	28.57%	1.0574
3 Extreme Gradient Boosting	default	27.71%	1.0449
1 Random Forest	default	27.26%	0.4045
0 Logistic Regression	solver=[lbfgs], multi_class=[multinomial]	23.63%	1.1063
5 K-Nearest Neighbors	n_neighbors = 10	20.31%	Unable to calculate
4 Naive Bayes	default	0.92%	13.4938

3. Hyperparameter Tuning

3.1 Chosen Parameters

Since cross validation required a tremendous amount of memory space, our resource constraints compelled us to split the process into 2 steps: first we find a reasonable number of trees (n\_estimators) given a fixed parameters according to the rules of thumb that we researched online; then we adjusted the max\_depth afterwards. Again, due to computational resource constraints, tuning over 3 parameters was found to be prohibitively expensive, so we decided to tune only these two parameters.

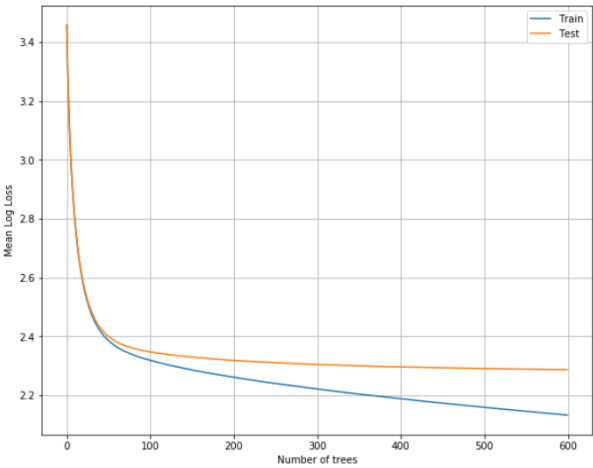


Figure 3.1 No significant improvement after 100 trees.

Although Random Forest had a slight advantage over other models on an efficiency basis, its’ memory requirements still exceeded our maximum memory space. Therefore, we again split the process into two steps and found our optimal parameters for XGBoosting and random forest to be 100 estimators and 50 estimators for XGBoosting and Random Forest respectively, both with max\_depth 3.

Table 3.1 Optimal Parameters for XGBoosting and Random Forest found by hyperparameter tuning

Parameters	XGBoosting	Random Forest
n_estimators	100	50
max_depth	3	3

3.2 Result Analysis

Fig. 3.2 below corroborates the observation in Fig. 1.1 that the geospatial coordinates (X, Y) of a crime incident plays an important role in classifying its crime type. Moreover, our engineered feature n\_days was also found to have a surprisingly significant role in predicting the crime class which suggests that there may be a significant temporal pattern in the crime incident types.

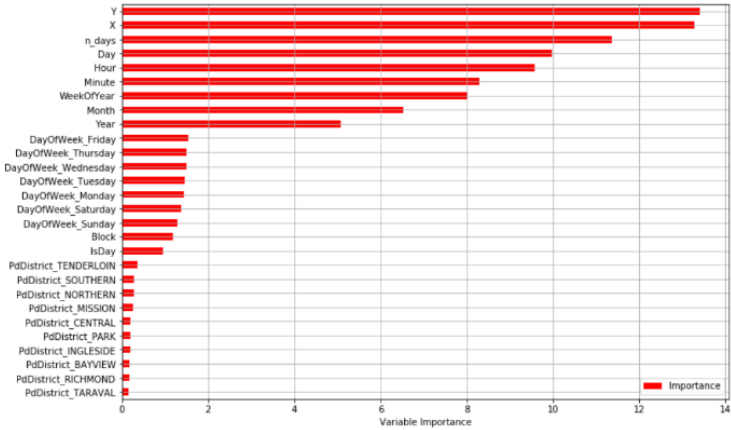


Figure 3.2 Feature importance weighting assigned by the Random Forest Model

3.3 Kaggle Result

After tuning the parameters, we used cross validation to compare the result of two models:

Random Forest: 2.91

Extreme Gradient Boosting: 5.33

Since we found that Random Forest had a lower mean log loss, we submitted the Random Forest model. The final score from Kaggle was determined to be: 2.58605.

submission.csv  
2 days ago by Kevin Tsai  
add submission details

For comparison purposes, we also decided to submit the result of XGB and found it yielded a better score of 2.35412:

Name: submission.csv  
Submitted: a minute ago  
Wait time: 0 seconds  
Execution time: 51 seconds  
Score: 2.35412

We suspect that the reason for this unexpected outcome may be due to the fact that we engineered a set of temporal features (n\_days, WeekofYear, etc.) that were better classifiers than many of the geographic features but that these engineered temporal features did not generalize well to the Kaggle evaluation data. In essence, our Random Forest model had “overfit” the provided test data and, in so doing, by excluding important geographic features (due to the Random Forests efforts to reduce overfitting of the train set) such as PdDistrict\_ABC etc., this had the effect of introducing omitted variable bias into the Random Forest model, thereby resulting in an underfitting of the true underlying data generating process. This is consistent with the observation of a temporal pattern in the crime classifications as suggested by the high feature importance weighting assigned to n\_days. In contrast, the XGBoosting model retained all features, including the Geographic features like PdDistrict\_ABC (to which XGBoosting assigned high feature importance) and was consequently able to achieve better performance in the final Kaggle submission.