



The Seattle Public Library

KNN, RandomForest, and XGBoost Classifiers are doing the best job at the moment.

XGBoost Classifier

```
In [266]: 1 import xgboost as xgb

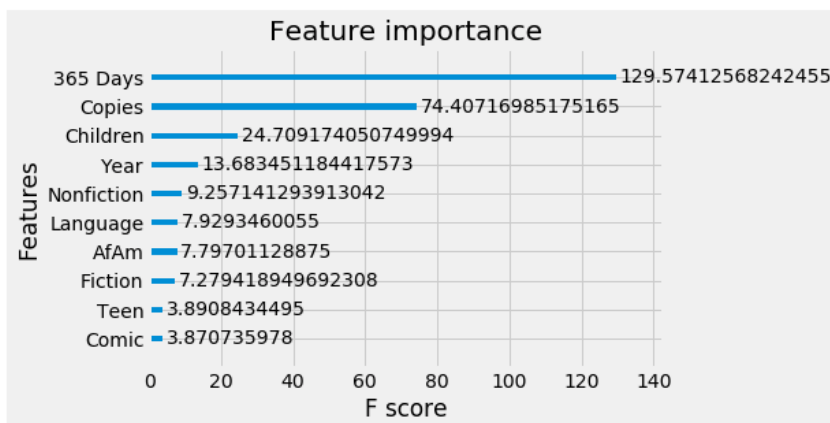
In [267]: 1 gbm = xgb.XGBClassifier(
2           n_estimators=30000,
3           max_depth=4,
4           objective='binary:logistic', #new objective
5           learning_rate=.05,
6           subsample=.8,
7           min_child_weight=3,
8           colsample_bytree=.8
9       )
10
11 eval_set= [(X_train,y_train),(X_val,y_val)]
12 fit_model = gbm.fit(
13     X_train, y_train,
14     eval_set=eval_set,
15     eval_metric='error', #new evaluation metric: classification error (could also use AUC, e.g.)
16     early_stopping_rounds=50,
17     verbose=False
18 )
19
20 # accuracy_score(y_test, gbm.predict(X_test, ntree_limit=gbm.best_ntree_limit))

In [268]: 1 y_val_predict = fit_model.predict(X_val)

In [269]: 1 print("Accuracy: ", accuracy_score(y_val, y_val_predict ))
2 print("Recall: ", recall_score(y_val, y_val_predict)) # Recall
3 print("Precision: ", precision_score(y_val, y_val_predict)) # Precision
4 print("f1: ", f1_score(y_val, y_val_predict))

Accuracy: 0.8200159277940006
Recall: 0.8055391868002357
Precision: 0.7970845481049562
f1: 0.8012895662368112
```

These are the feature importances sp far from XGBoost



My next step is to consider adding in more features