

17.Random Forest:

Now we will utilize random forest classification instead of logistic regression. To do this we will import xgboost utilizing the built in method in xgboost to do our random forest modeling. We will conduct the same sequence that we did for Logistic Regression besides finding the best z-score as we already previously did that and found it to be 1.7 and 1.4 for english and universal respectively. Like previously we get the z-score, remove the outliers, use random forest on the cleaned up data set, then get our confusion matrices. It should look like this for the imports:

```
1 import pickle
2 import numpy as np
3 import pandas as pd
4 from matplotlib import pyplot as plt
5 from sklearn.metrics import ConfusionMatrixDisplay, accuracy_score
6 from sklearn import model_selection, preprocessing
7 from xgboost import XGBClassifier
```

We use the same methods as before for remove outliers and getting the z-score so they will not be gone over again refer to the previous logistic regression overview to refresh your knowledge on how to implemented those functions. The main difference is obviously how random forest is implemented and how it is used for the confusion matrices which we will now cover.

The random forest algorithm works by first selecting random samples from a given dataset then the algorithm constructs a decision tree for each training datapoint. A type of collective voting is done by averaging these decision trees out (mathematically). Finally the final prediction result is the most voted on decision tree. We implement this in python by way of XGBClassifier. We then create a dictionary which holds the attributes max_depth, n_estimators, and learning rate which corresponds to the max depth of the tree we go down which we assign to "2" the estimators which is the max amount of decision trees we make until final decision averages take place, then the learning rate which corresponds to the speed at which the model learns. This is from a scale of 0 to 1 which the closer to 0 the number is the more careful your model is when in the training phase. We then use GridSearchCV which helps us find the best parameters for our model. We pass into GridSearchCV our created dictionary then use pickle to open the model name and then we print the result like [this](#):

```
47 def random_forest(input_df, print_result=False):
48     x = input_df.drop(["spammer", "label"], axis=1).values
49     if "id" in input_df.columns:
50         x = input_df.drop(["id", "spammer", "label"], axis=1).values
51     y = input_df["label"].values
52     x_train, x_test, y_train, y_test = model_selection.train_test_split(x, y, test_size=0.30, random_state=100)
53
54     x_train_scaled = preprocessing.scale(x_train)
55
56     model = XGBClassifier(use_label_encoder=False)
57     optimization_dict = {
58         'max_depth': [2],
59         'n_estimators': [200],
60         'learning_rate': [0.1],
61     }
62     model = model_selection.GridSearchCV(model, optimization_dict,
63                                         scoring='accuracy', verbose=1)
64     model.fit(x_train_scaled, y_train, eval_metric="error")
65
66     pickle.dump(model, open(f'rand_forest_eng.dat', 'wb'))
67
68     if print_result:
69         print(model.best_score_)
70         print(model.best_params_)
```

For confusion matrix nearly the same is done akin to what we did for the logistic regression except now we do it for random forest. We still do the 4 types of normalization for the confusion matrices and all. Like last time we need to call our methods at the bottom to see our work. It should look like [this](#) (GitHub code base for random forest):

```
73 def confusion_matrix(lang, input_df, show=False):
74     x = input_df.drop(["spammer", "label"], axis=1).values
75     if "id" in input_df.columns:
76         x = input_df.drop(["id", "spammer", "label"], axis=1).values
77     y = input_df["label"].values
78     x_train, x_test, y_train, y_test = model_selection.train_test_split(x, y, test_size=0.30, random_state=100)
79
80     x_train_scaled = preprocessing.scale(x_train)
81     x_test_scaled = preprocessing.scale(x_test)
82
83     model = XGBClassifier(use_label_encoder=False)
84     optimization_dict = {
85         'max_depth': (2),
86         'n_estimators': (200),
87         'learning_rate': (0.1),
88     }
89     model = model_selection.GridSearchCV(model, optimization_dict,
90                                         scoring='accuracy', verbose=1)
91     model.fit(x_train_scaled, y_train, eval_metric="error")
92
93     title_options = [(f"Confusion Matrix, without normalization\n({lang})", None),
94                     (f"Normalization: true\n({lang})", "true"),
95                     (f"Normalization: pred\n({lang})", "pred"),
96                     (f"Normalization: all\n({lang})", "all")]
97
98     for title, normalize in title_options:
99         disp = ConfusionMatrixDisplay.from_estimator(model, x_test_scaled, y_test, display_labels=["bot", "human"],
100                                                    cmap=plt.cm.Blues, normalize=normalize)
101         disp.ax_.set_title(title)
102         # plt.savefig(fr'C:\Users\Rachel\Documents\Twitter Data\normalize_{lang}_test')
103
104     if show:
105         plt.show()
106
107
108 dff = orgs_to_bots(input_path_en)
109 confusion_matrix("eng", remove_outliers(1.7, dff), True)
```