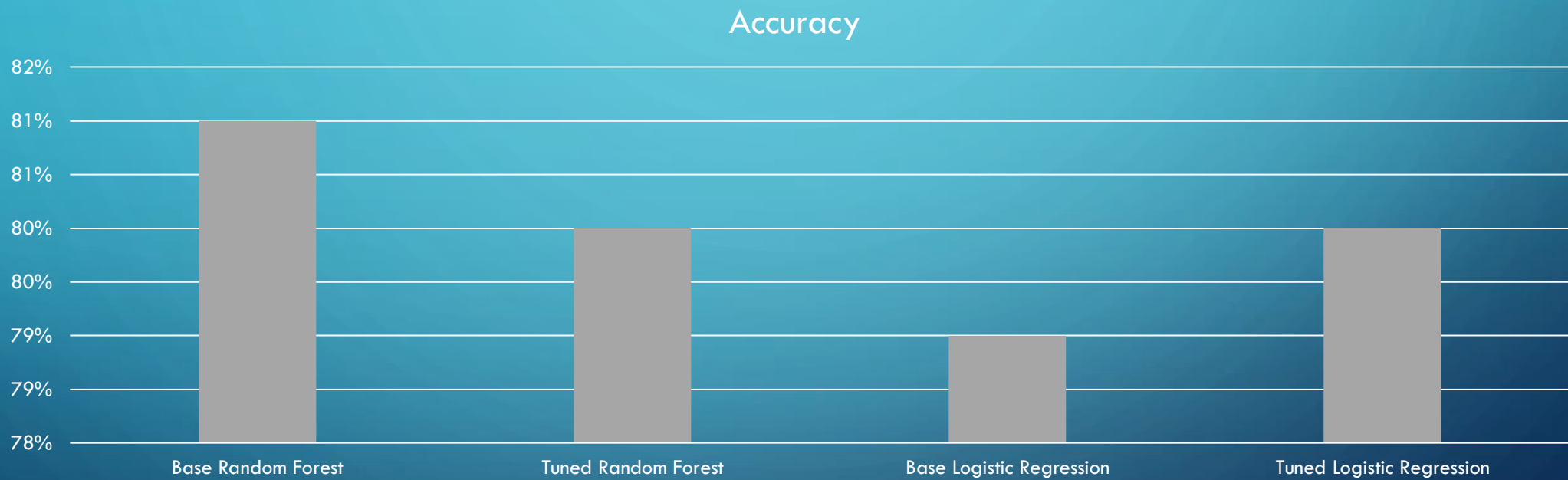# TWITTER HATE SPEECH ANALYSIS
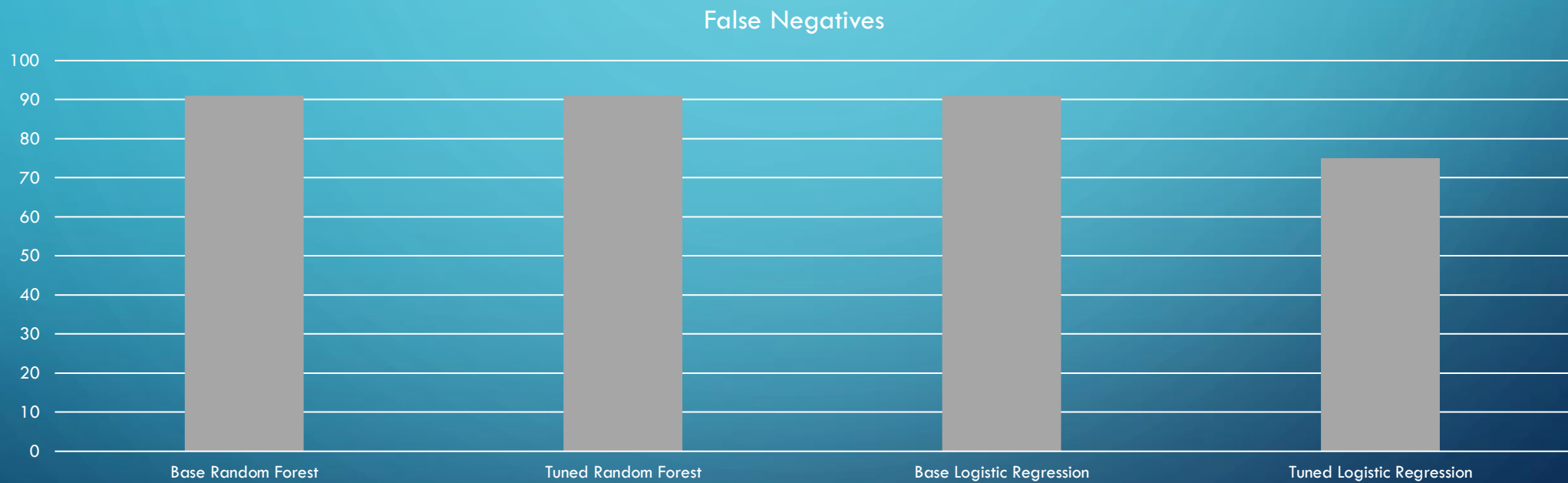
## MILESTONE THREE

COLIN GREEN & SEAN ZHANG

# BALANCED DATASET

- Ran Random Forest and Logistic Regression as a benchmark

- Compared with tuned versions of Random Forest and Logistic Regression

- Used the bagging method for those 4 models as well as untuned:
  - Extra Trees, KNN, SVC, Ridge Classifiers

- Also used Ada Boost, Grad Boost, XG Boost and an Ensemble
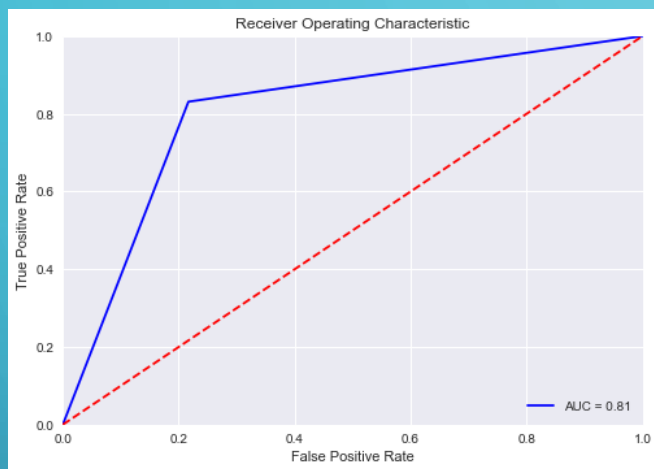
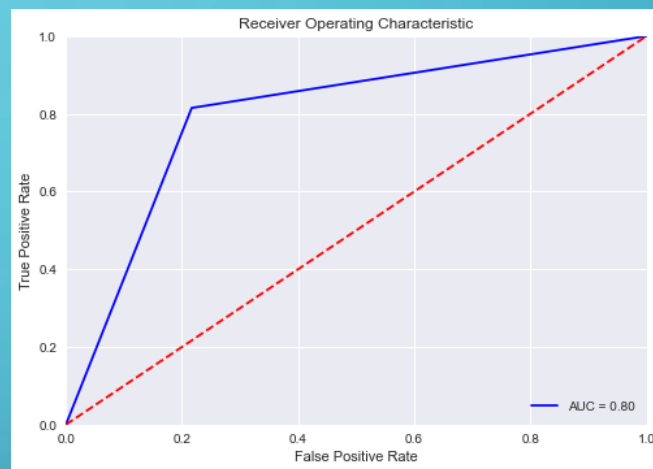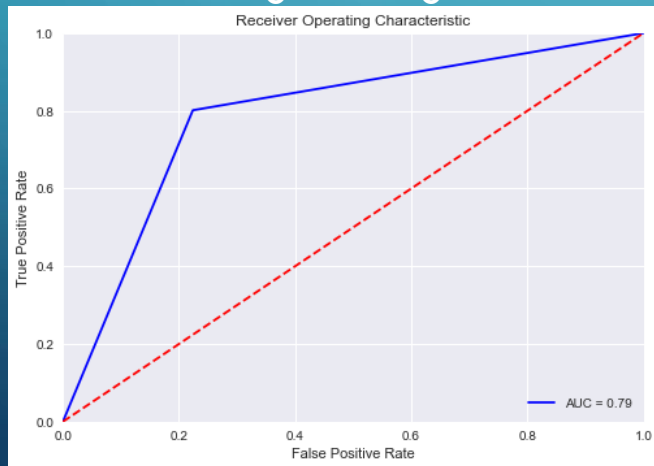# LOGISTIC REGRESSION VS RANDOM FOREST

Accuracy

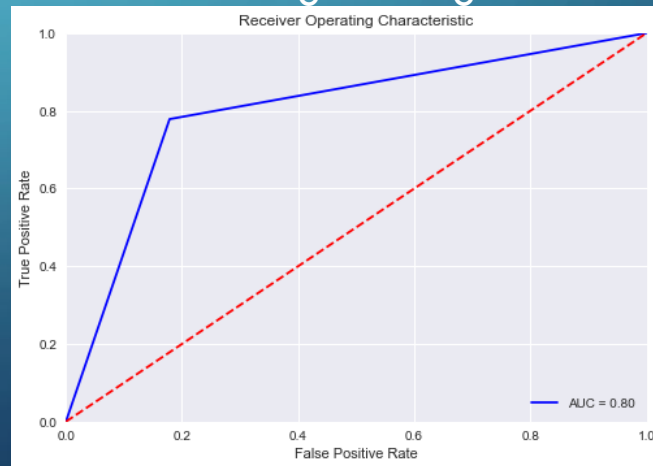# FALSE NEGATIVES

False Negatives
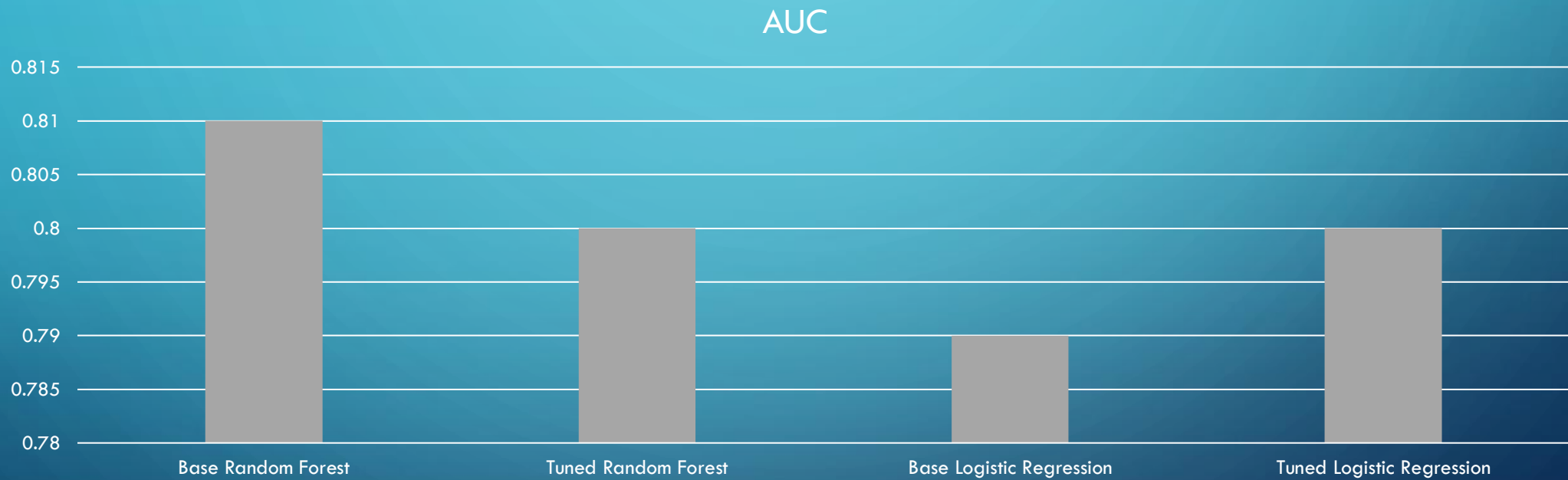
# ROC CURVES



Base Random Forest



Tuned Random Forest



Based Logistic Regression
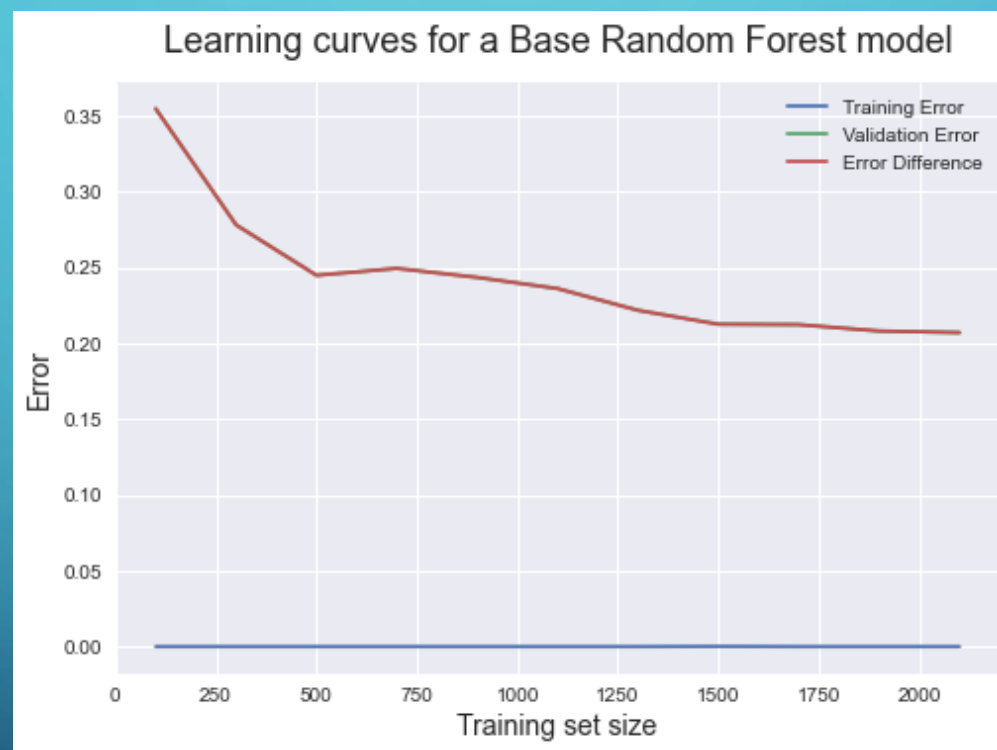


Tuned Logistic Regression

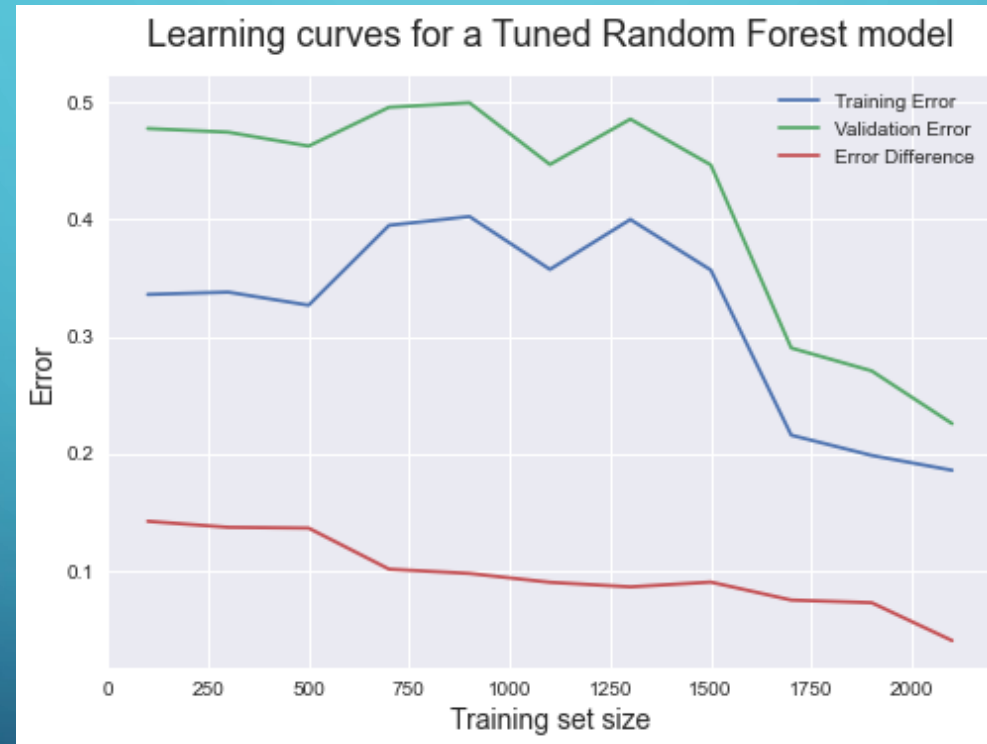# AREA UNDER THE CURVE

AUC



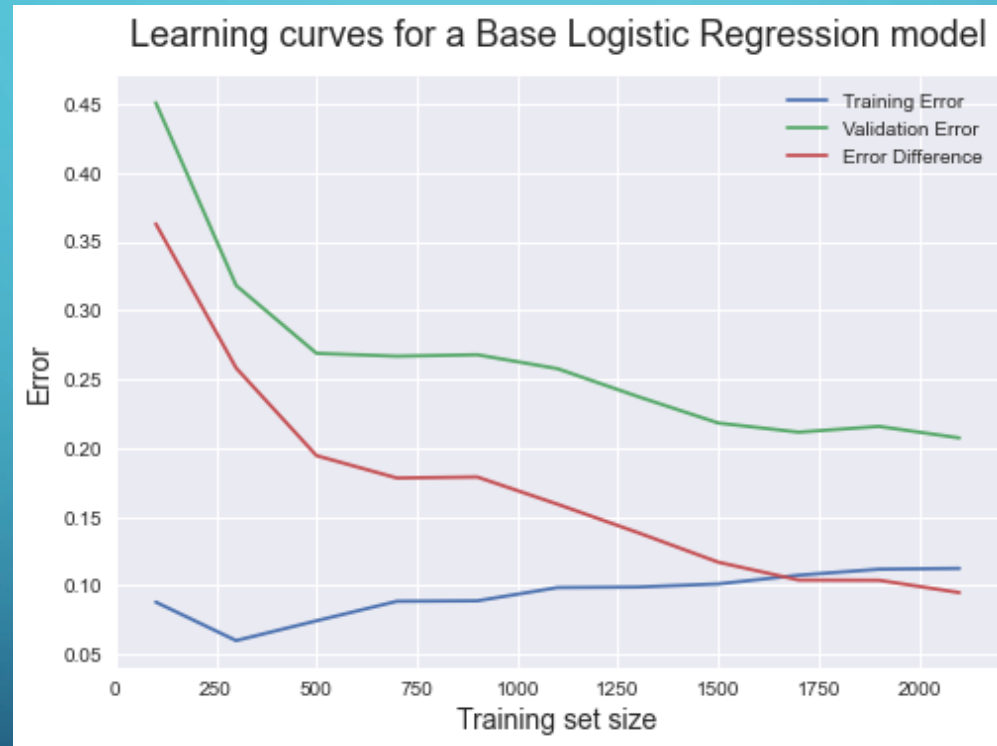| | Base Random Forest | Tuned Random Forest | Base Logistic Regression | Tuned Logistic Regression |

# LEARNING CURVES – BASE RANDOM FOREST

# LEARNING CURVES – TUNED RANDOM FOREST

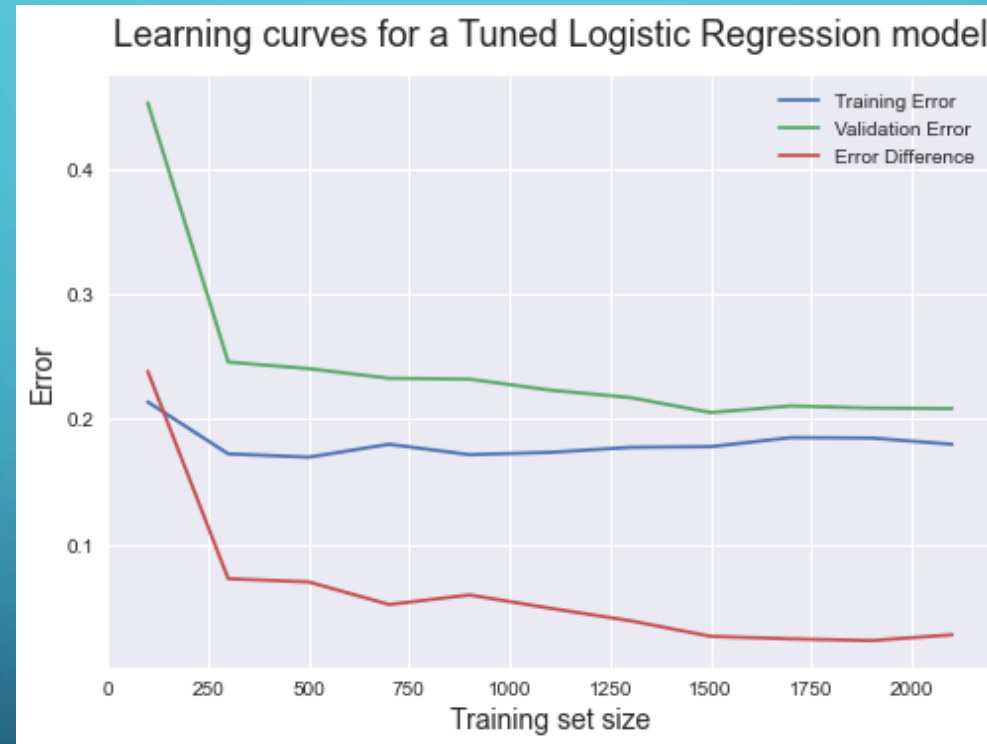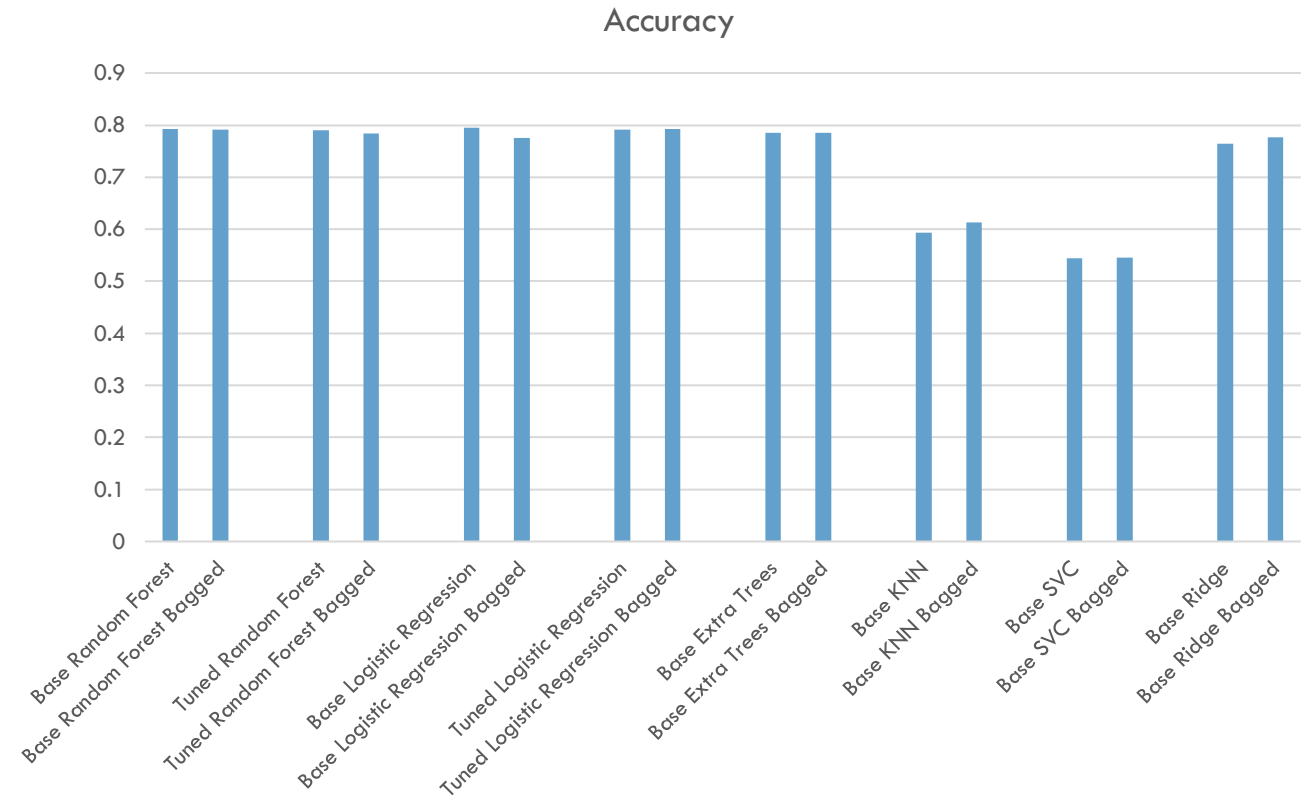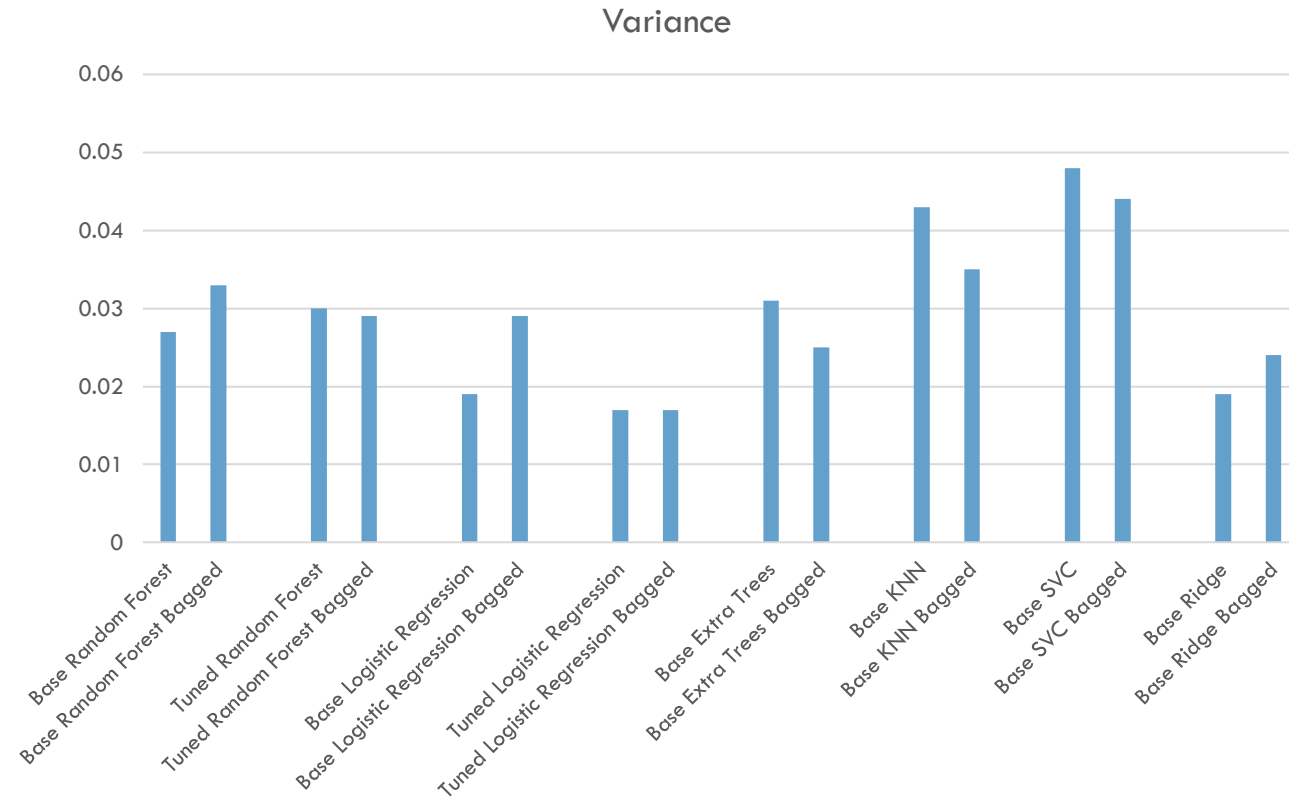# LEARNING CURVES – BASE LOGISTIC REGRESSION



Learning curves for a Base Logistic Regression model

# LEARNING CURVES – TUNED LOGISTIC REGRESSION
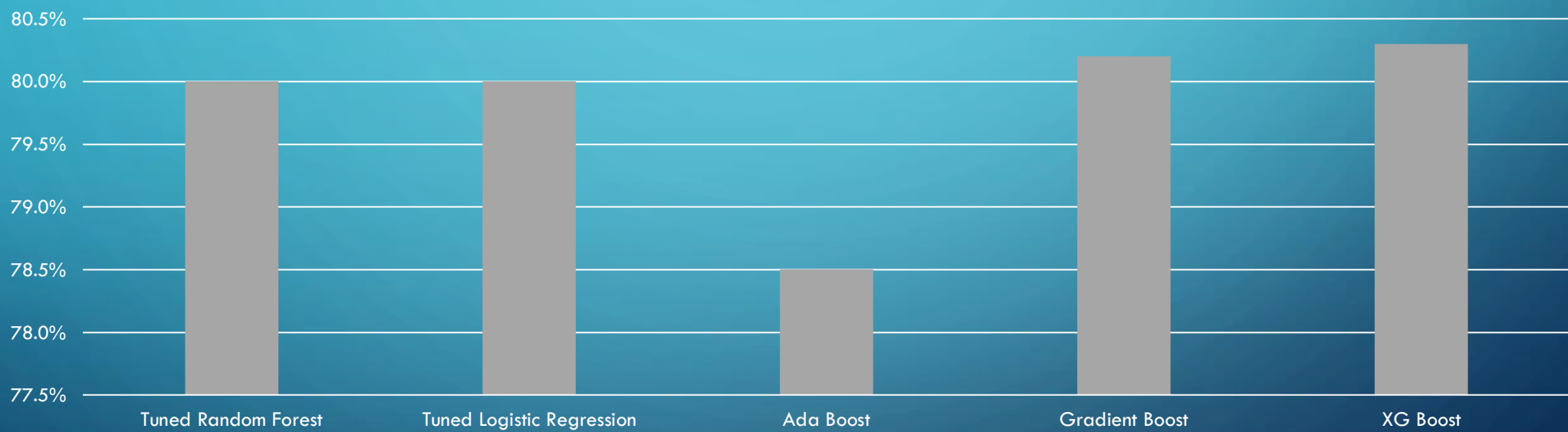
# BAGGING



Accuracy

# BAGGING TO REDUCE VARIANCE
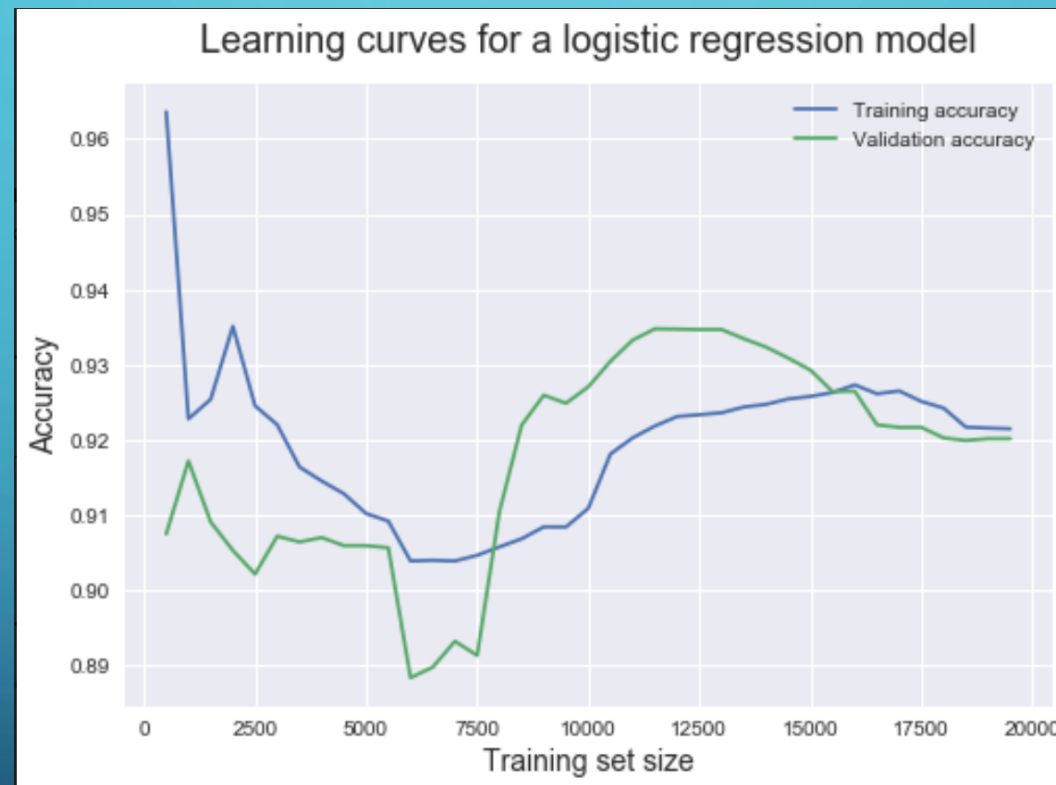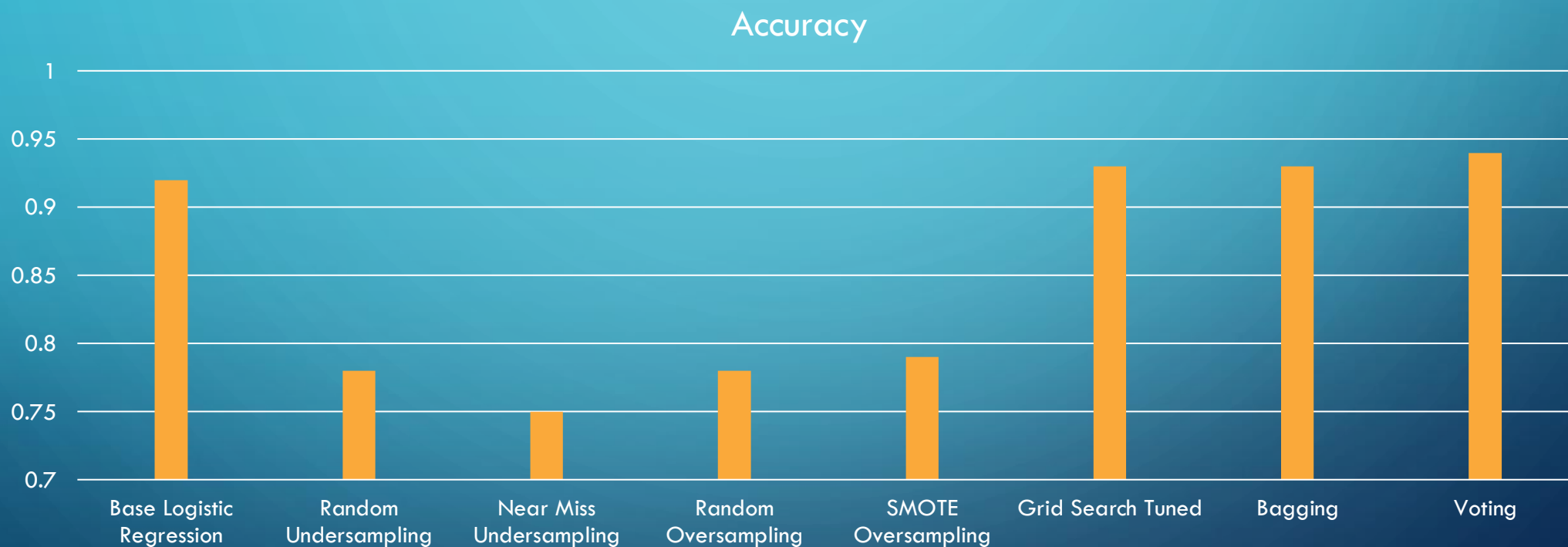


Variance

# BOOSTING



Accuracy

# FULL DATASET

- Used Logistic Regression as Benchmark

- Ran to see whether learning curves would converge

- Comparison of: Base, under/oversampling, hyperparameter tuning, bagging, ensemble (voting)

- Decision boundary visualization

# LEARNING CURVES – TUNED LOGISTIC REGRESSION (FULL DATA)

# ACCURACY METRICS (FULL DATA)

Accuracy



Imbalanced sampling done with **imblearn** package

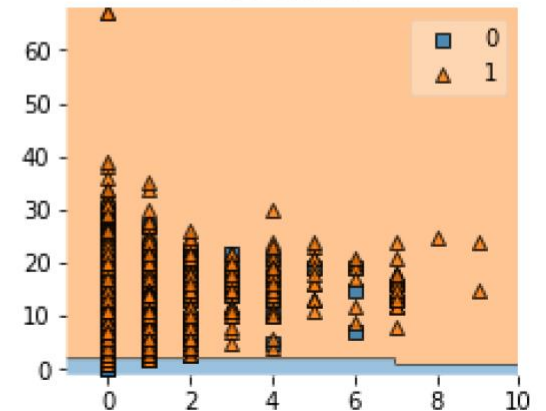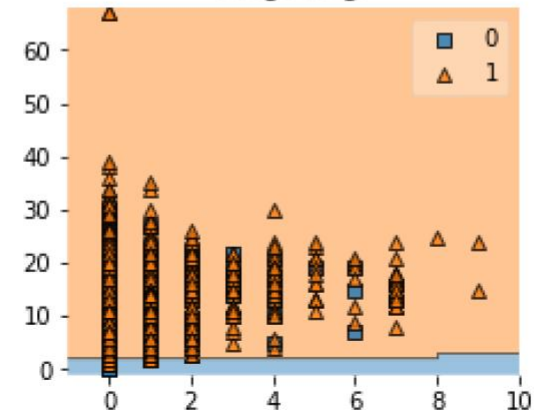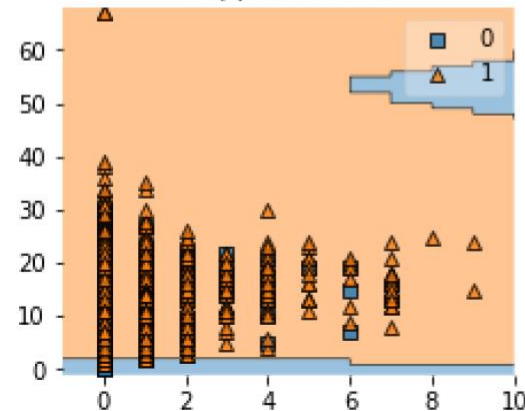Voting ensemble: Random Forest, Extra Trees, KNN, Support Vector Machine, Logistic Regression
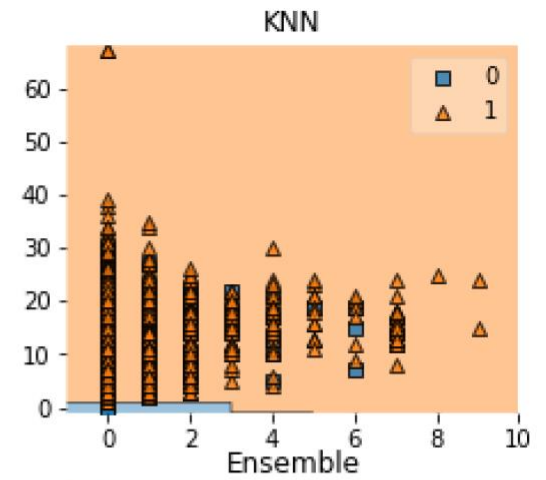
# DECISION BOUNDARIES

X: hashtag count
Y: number of tokens

0: Hate speech
1: Non-hate speech
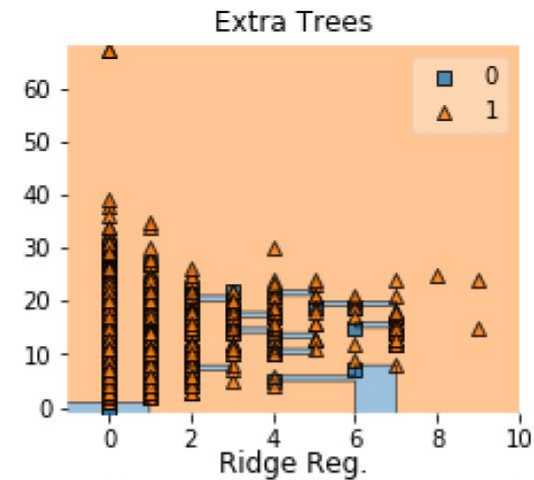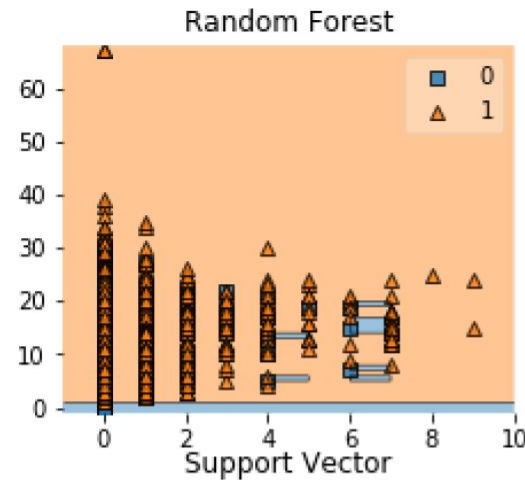
- Chose numerical variables based on importance and the fact that word features are too sparse

- Observations with lower number of tokens tends to be classified as hate speech (this was something noticed in previous EDA, feature importance, and past research)

# CONCLUSIONS

- Algorithm converged in full dataset compared to balanced subset

- Tuning slightly increased accuracy

- Bagging slightly increased accuracy and decreased variance in most cases

- Imbalanced sampling significantly decreased accuracy with no improvement in sensitivity or precision

- Ensemble learning slightly improved accuracy