# CS613 FINAL PROJECT

Connor Secen, Richard Strouss, Tristan Marshall

# PROBLEM STATEMENT

- COMPARE THE RESULTS OF MULTIPLE METHODS ON PREDICTING THE DISCRETIZED CRIME RATE OF VARIOUS LOCALITIES GIVEN A RELEVANT FEATURE SET

- TO INCREASE PRACTICAL APPLICABILITY, DIMENSIONALITY WILL BE REDUCED WHILE MAINTAINING ACCEPTABLE VALIDATION SCORES

# PRIOR WORK

- Ingilevich, V. & Ivanov, S. (2018). Crime Rate Prediction in the Urban Environment Using Social Factors. *Procedia Computer Science, 136,* 472-478. https://doi.org/10.1016/j.procs.2018.08.261
  - Compared results from linear regression, logistic regression, and gradient boosted decision trees
  - Found gradient boosting to be the most appropriate technique of the three
- Alves, L. G. A., Ribiero, H. V., & Rodrigues, F. A. (2018). Crime Prediction Through Urban Metrics and Statistical Learning. *Physica A, 55,* 435-443. https://doi.org/10.1016/j.physa.2018.03.084
  - Used a random forest model
  - Out performed previous linear models
- Kshatri, S. S., Singh, D., Narain, B., Bhatia, S., Quasim, M. T., & Sinha, G. R. (2021). An Empirical Analysis of Machine Learning Algorithms for Crime Prediction using Stacked Generalization: An Ensemble Approach. *IEEE Access.* Advanced online publication. https://doi.org/10.1109/ACCESS.2021.3075140
  - Compared J48 decision tree algorithm, random forests, sequential minimal optimization (SVM), bagging, and naïve bayes classifiers
  - Best results from stacking the models
  - Random forest was nearly as good as the stacked model on all validation metrics

# PRIOR WORK

- Zhu, J., S. Rosset, H. Zou, and T. Hastie. 2009. Multi-Class AdaBoost. Statistics and its Interface2 (3):349–360.
  - Multi-class Adaboost extension
  - Fixes problem with high multi-class error
- General lessons
  - Features are important
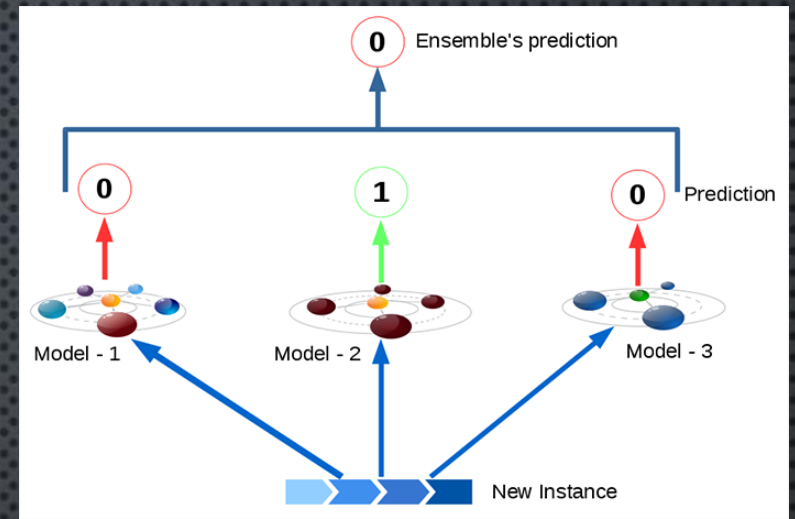  - Subject similarity does not guarantee learning performance similarity

# DATASET

- Communities and Crimes Dataset from the UC Irvine machine learning repository

  - http://archive.ics.uci.edu/ml/datasets/Communities+and+Crime

- Combines socio-economic data from 1990 US Census, law enforcement data from 1990 US LEMAS (*Law Enforcement Management and Administrative Statistics*) survey, and crime data from 1995 FBI UCR (*Uniform Crime Reporting*)

- Dataset contains 1994 instances and 128 Attributes

  - Mix of categorical and continuous features

- 102 attribute With removal of missing data including

- Converted the target value to categorical by binning the data into 11 different groups separated by .1
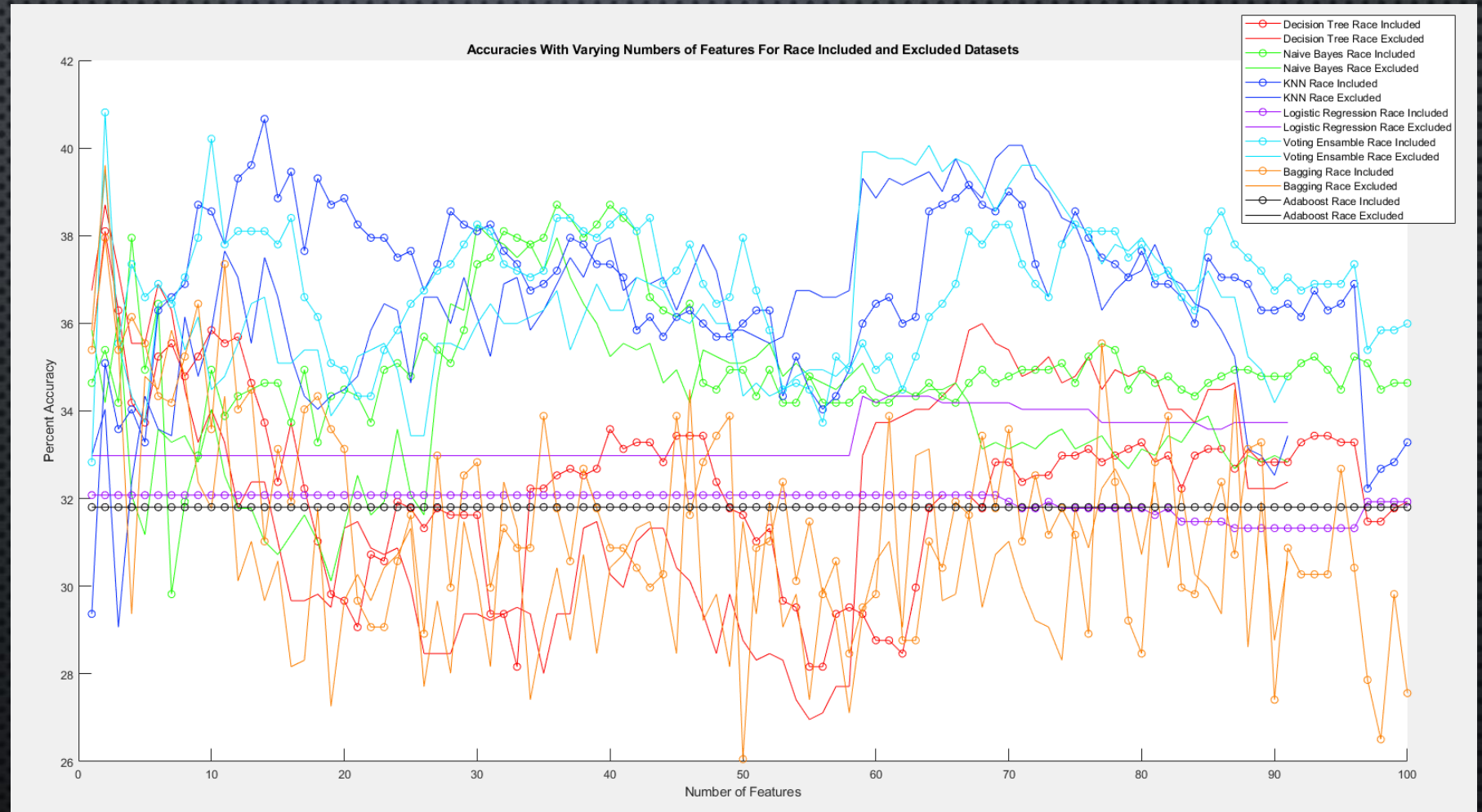
# APPROACH



- FOUND THE CORRELATION BETWEEN EACH FEATURE AND THE TARGET VALUE
  - SPLIT THE DATASET INTO ONE, ONE CONTAINING AND ONE NOT CONTAINING RACE DATA
  - SORTED THE FEATURE FROM MOST CORRELATED TO LEAST AND PROGRESSIVELY ADDED THE FEATURES TO THE DATASET CURRENTLY BEING SORTED
- COMPARE SIMPLE CLASSIFICATION METHODS
  - LOGISTIC REGRESSION VS DECISION TREE
- COMPARE MORE COMPLEX ENSEMBLE METHODS
  - VOTING
    - LOGISTIC REGRESSION
    - ID3 DECISION TREE
    - KNN
    - NAÏVE BAYES
  - BAGGING
  - BOOSTING
    - BETA VALUE FOR WEIGHT UPDATE CALCULATE AS FOLLOWS: BETA = LOG(1-ERROR/ERROR)+LOG(K-1) WHERE K IS THE CLASS
- DETERMINE WHICH METHOD WORKS BEST FOR THE GIVEN DATA BY COMPARING THE ACCURACY

# RESULTS



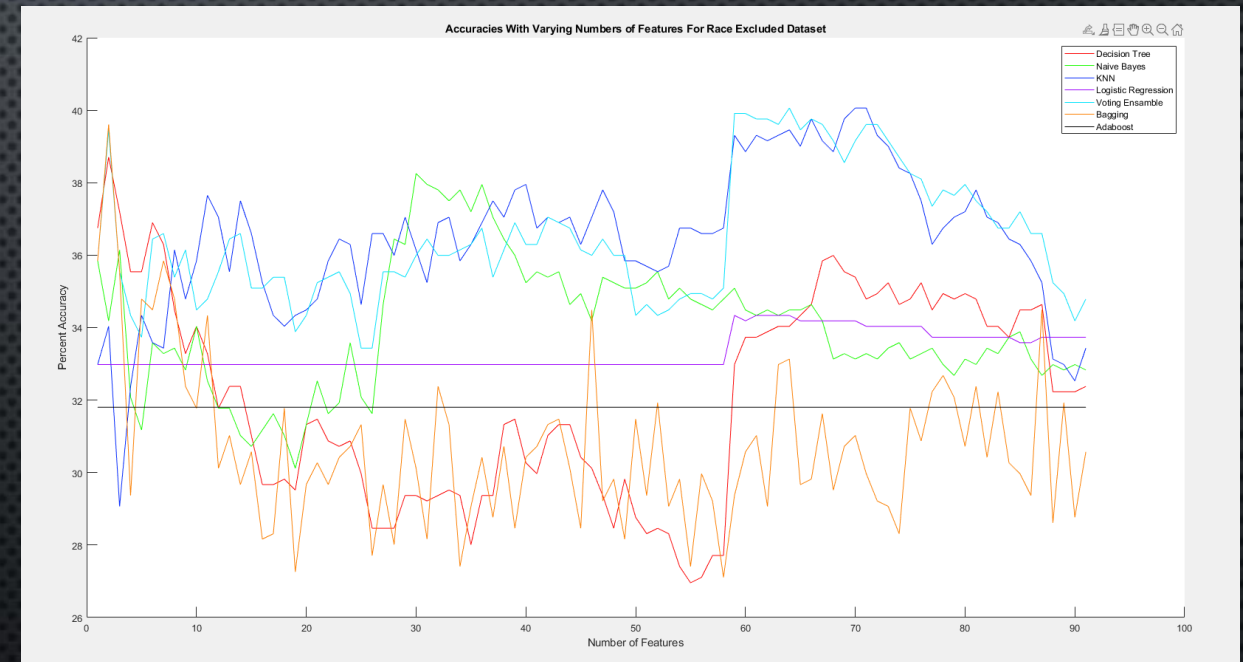Accuracies With Varying Numbers of Features For Race Included and Excluded Datasets

# OBSERVATIONS AND DISCUSSION

- MAXIMUM ACCURACY 40.06%

- BEST CLASSIFIERS VOTING ENSEMBLE AND KNN

- RACE DATA DID NOT ADD SIGNIFICANTLY TO ACCURACY AND IN MANY CASES DECREASED ACCURACY

- LOGISTIC REGRESSION AND ADABOOST HAD CONSISTENT ACCURACIES, BUT ADHERED TO SPECIFIC CLASSES



Accuracies With Varying Numbers of Features For Race Excluded Dataset

# OBSERVATIONS AND DISCUSSION

- Low accuracy potentially a result of the data itself

- Dataset includes:

  - census data which has a 10% margin of error

  - Categories that should not hold non-zero data but do

  - Potentially subjective categories

# FUTURE WORK/EXTENSIONS

- ATTEMPT TO USE MORE COMPLEX CLASSIFIERS
  - ARTIFICIAL NEURAL NETWORKS
- CHANGE THE LEARNING TASK TO BE A REGRESSION STYLE TASK
  - THE TARGET IN THE DATASET IS CONTINUOUS LENDING ITSELF TO REGRESSION
- COLLECT MORE DATA
  - MORE EXAMPLES TO REDUCE OVERFITTING
  - DIFFERENT FEATURES TO AVOID BIAS
  - REDUCE THE NUMBER OF FEATURES BY COLLECTING FEATURES THAT ARE BETTER PREDICTORS

# THANK YOU!