

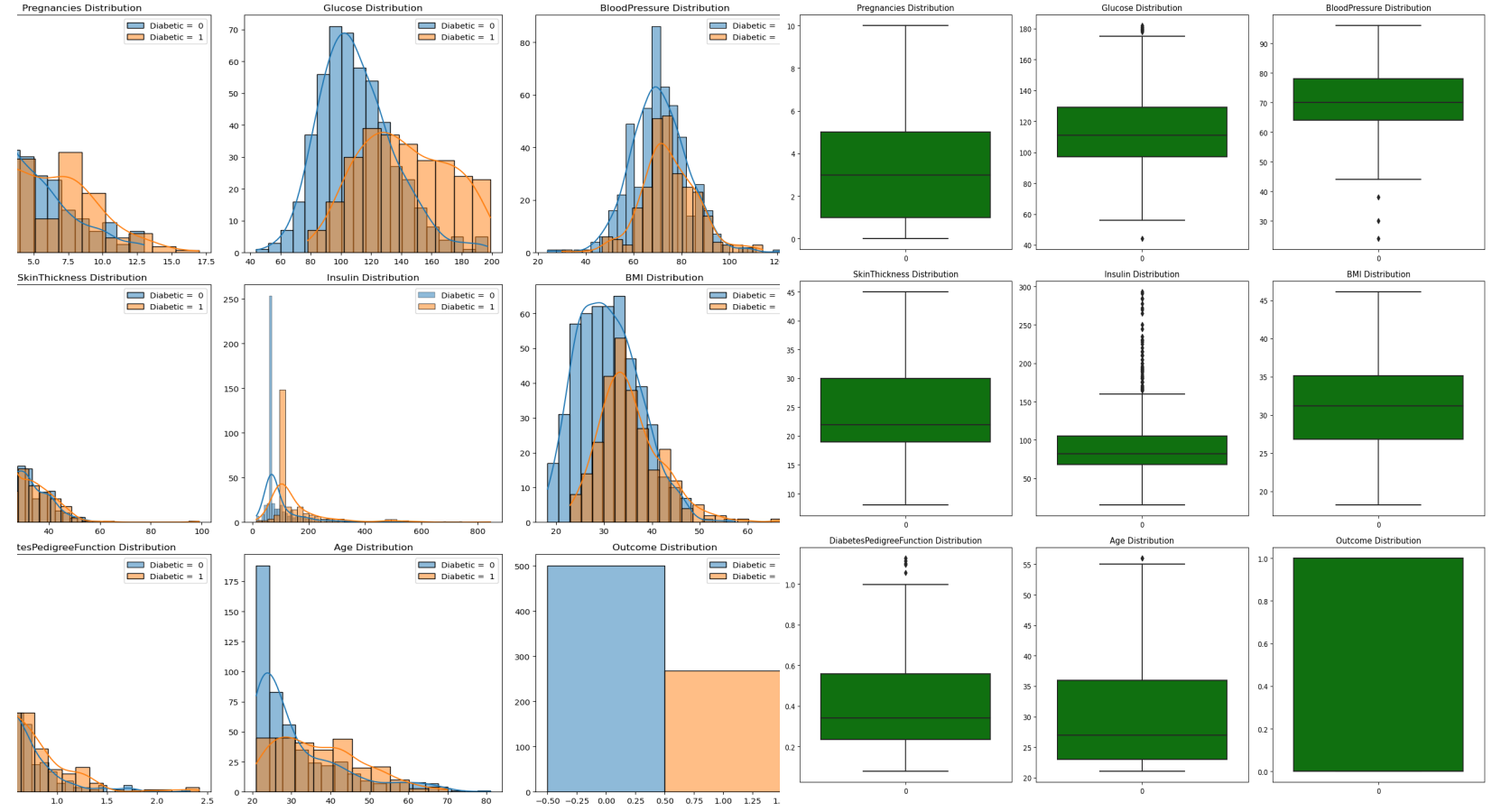
By Skyler Wilson

Supervised Learning Project

Exploratory Data Analysis

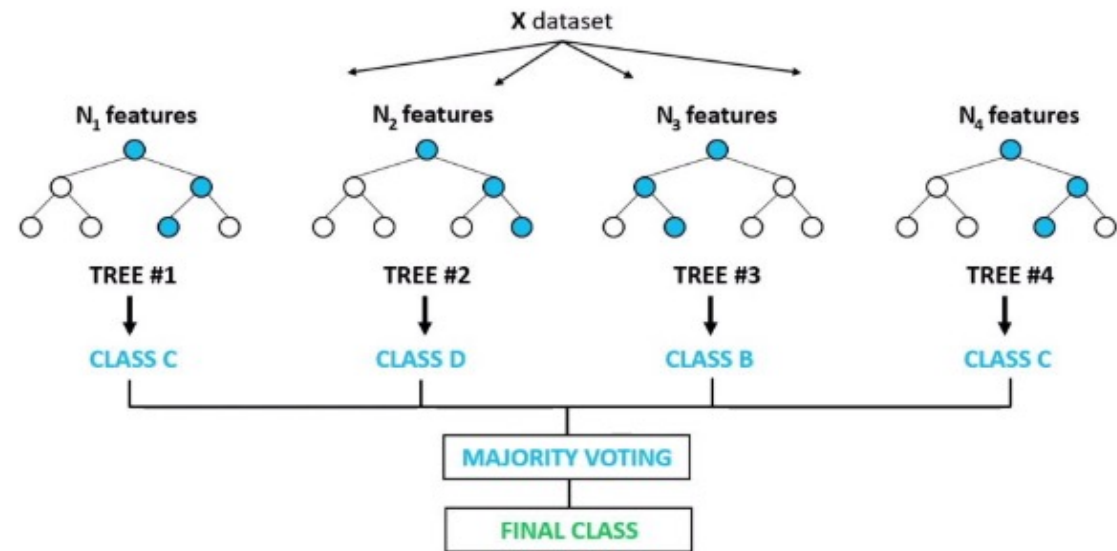
Correlation Heatmap

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction	Age	Outcome
Pregnancies	1.00	0.13	0.21	0.02	-0.00	0.02	-0.03	0.54	0.22
Glucose	0.13	1.00	0.22	0.18	0.43	0.24	0.14	0.27	0.50
BloodPressure	0.21	0.22	1.00	0.14	0.02	0.28	0.00	0.33	0.17
SkinThickness	0.02	0.18	0.14	1.00	0.26	0.55	0.16	0.03	0.22
Insulin	-0.00	0.43	0.02	0.26	1.00	0.22	0.17	0.05	0.26
BMI	0.02	0.24	0.28	0.55	0.22	1.00	0.15	0.03	0.32
DiabetesPedigreeFunction	-0.03	0.14	0.00	0.16	0.17	0.15	1.00	0.03	0.17
Age	0.54	0.27	0.33	0.03	0.05	0.03	0.03	1.00	0.24
Outcome	0.22	0.50	0.17	0.22	0.26	0.32	0.17	0.24	1.00



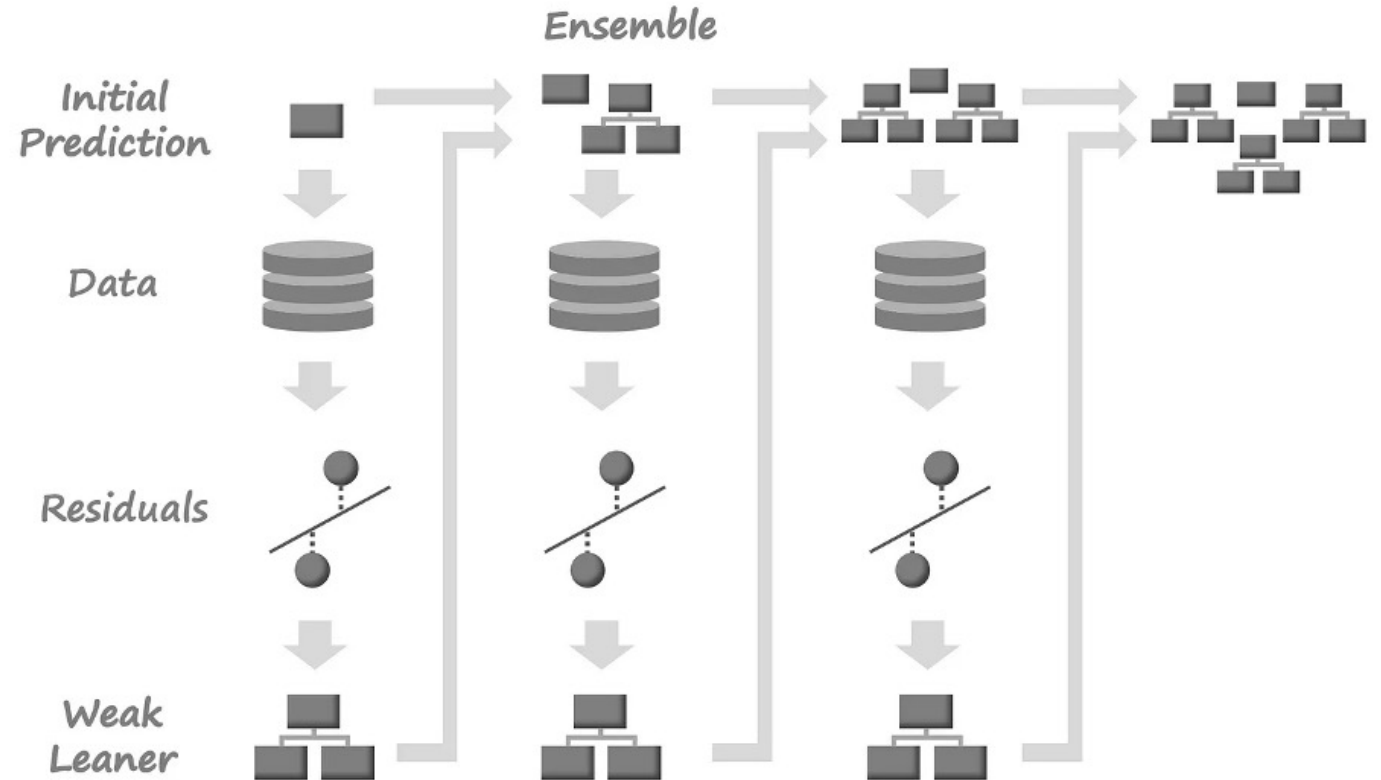
Random Forest Classifier

Random Forest Classifier



A random forest classifier is an ensemble method that aggregates several decision trees together to yield a more accurate machine learning model that is less biased and smaller variation

XGBoost Classifier



XGBoost classification is an ensemble method that aggregates multiple 'weak models' to form a more robust and accurate final model. XGBoost uses extreme gradient boosting to iteratively improve the models fit

Key Takeaways



The model accuracy was slightly better for the RandomForest than XGBoost at predicting whether an individual has diabetes based on the feature variables provided



The three most important features for prediction according to the RandomForestClassifier were Insulin, Glucose, and SkinThickness whereas in the XGBoostClassifier the three most important features for prediction were Insulin, Age, and BMI



While the RandomForestClassifier has a higher accuracy score and precision, XGBoost has a better recall score, and the F1 scores are very close in value. The reason the recall value difference is important is because it means that there are fewer false negatives, this is important because a false negative classification for an individual with diabetes is life threatening, since they will presumably not receive adequate treatment.



According to the heatmap of correlation the most highly correlated feature with the outcome was the level of Glucose, this is interesting because it was not the most important feature in either of the models