# STAT469 assignment #3

name : Bill Co

## I.   Overview

This Assignment is the extension of the previous one, in which I will introduce two more models, namely Multi-task Prediction(MTPS) and Random Forest. Furthermore, I have chosen

## II.   Data preparation and Model assessment

### a. Resampling

To ensure data stability, I have decided to use resampling to generate simulation data for model fittings. After some testing, I can make sure this method is identical to stratification, but some observations are repeated.

### b. Simulations

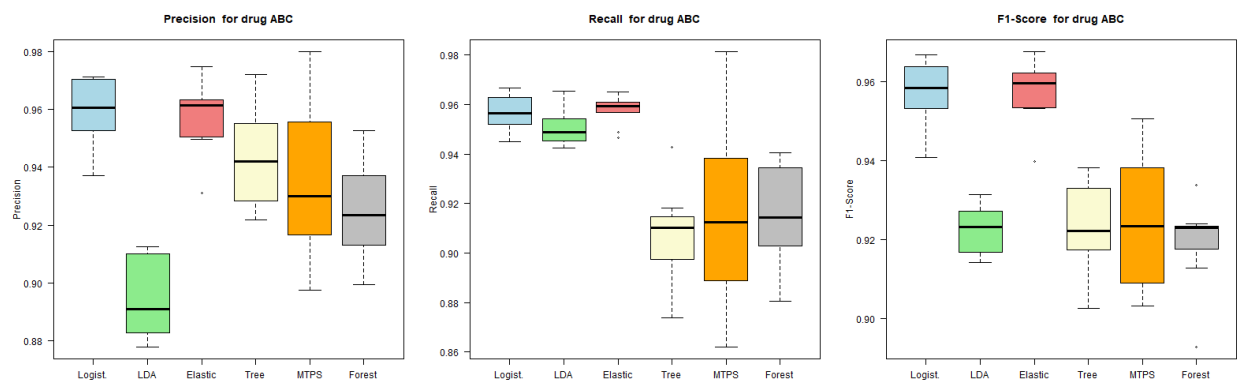There are going to be 50 simulations, each with 5 fold Cross Validation.

### c. Metrics

For productivity, some matrices are used to store metrics. Since the dataset shows class imbalance, metrics like Accuracy and Misclassification rate can be misleading. Thus, I excluded them from the assessment, keeping only Precision, Recall and F1 score.

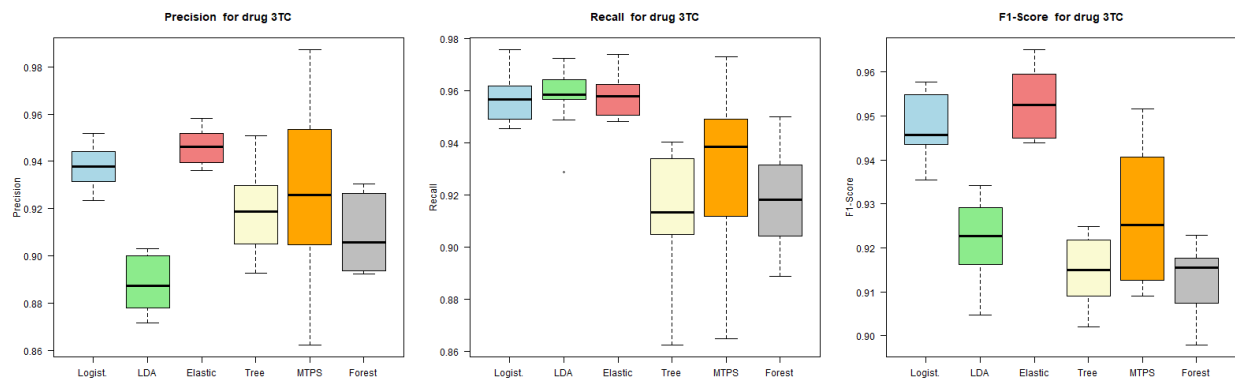# III. Result discussion (with visualization)
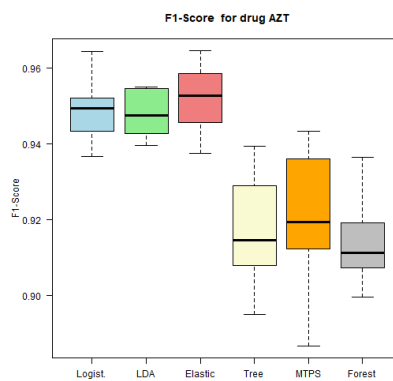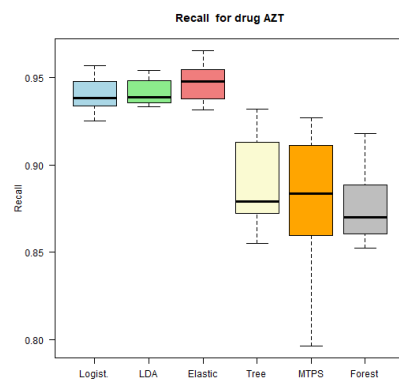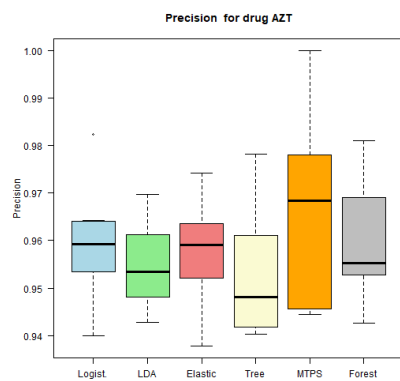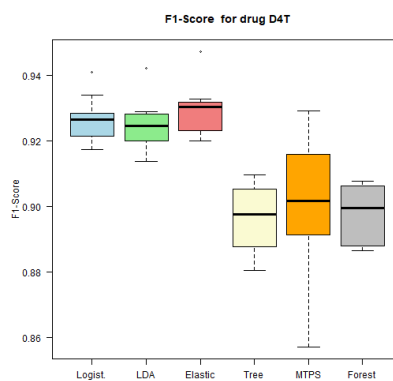
## 1. Individual Drug

### a. ABC



### b. 3TC



### c. AZT

Precision for drug AZT · Recall for drug AZT · F1-Score for drug AZT

## d. D4T


Precision for drug D4T · Recall for drug D4T · F1-Score for drug D4T

## e. DDI


Precision for drug DDI · Recall for drug DDI · F1-Score for drug DDI
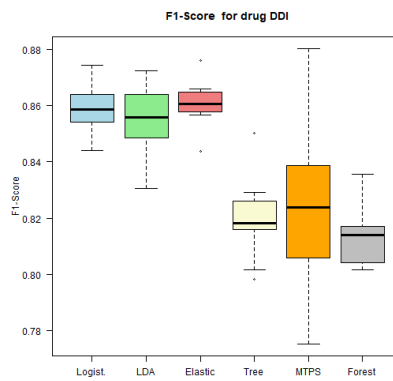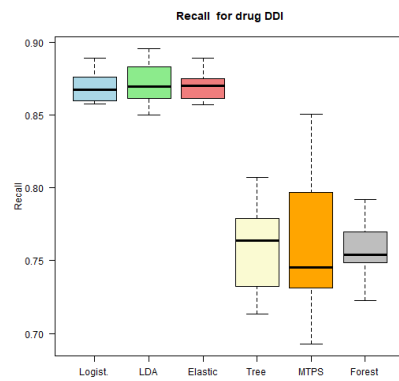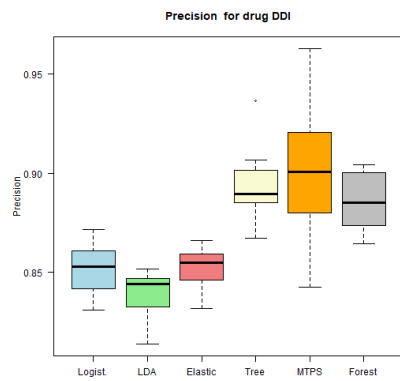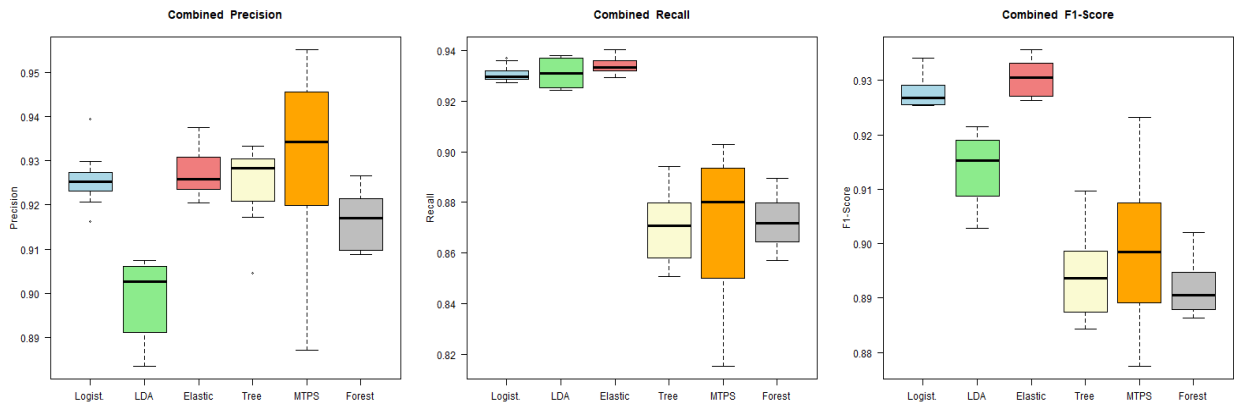
## 2. Combined Drug



The above visualizations are the assessments grouped by Drugs. As shown, Elastic Net and Logistic Regression are very stable across all measurements.

While MTPS shows a very high precision, the Recall and F1-Score are overkilled. LDA has a quite reasonable Recall but performs not very well in the Precision and F1-Score.

Tree and Forest are not good candidates for the best model in this scenario.

## 3. Hypothesis Testing

After conducting hypothesis, I concluded that:
- At the probability of 23%, there is no clear difference between the 2 models in Precision. However, Elastic Net performs better overall.

Thus, combining 2 criteria above, I decided **Elastic Nets should be used for predicting Drug resistance for HIV disease.**