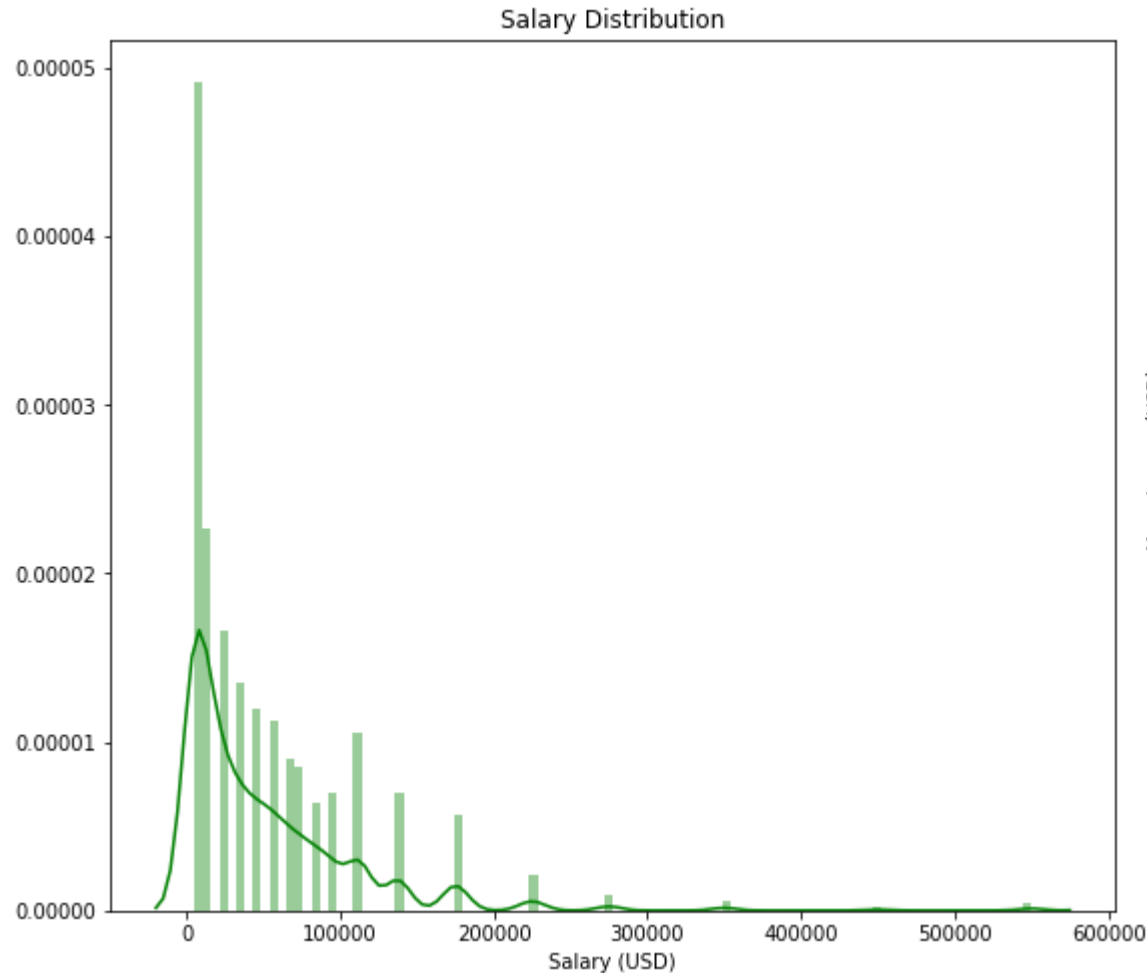# MIE 1624: Introduction to Data Science and Analytics
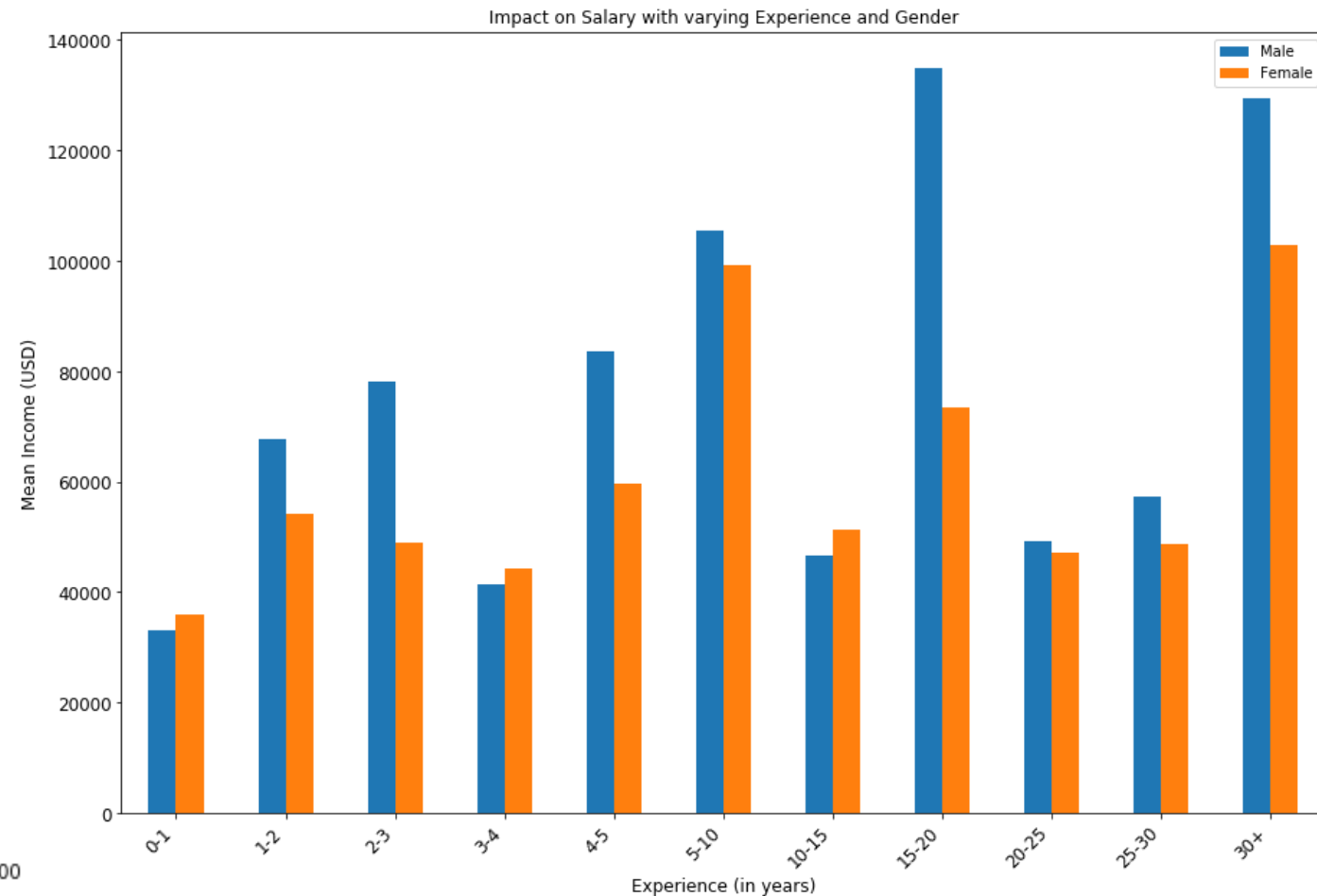# Assignment 1: Salary Classification

Syed Hamdan Mustafa
1006193209
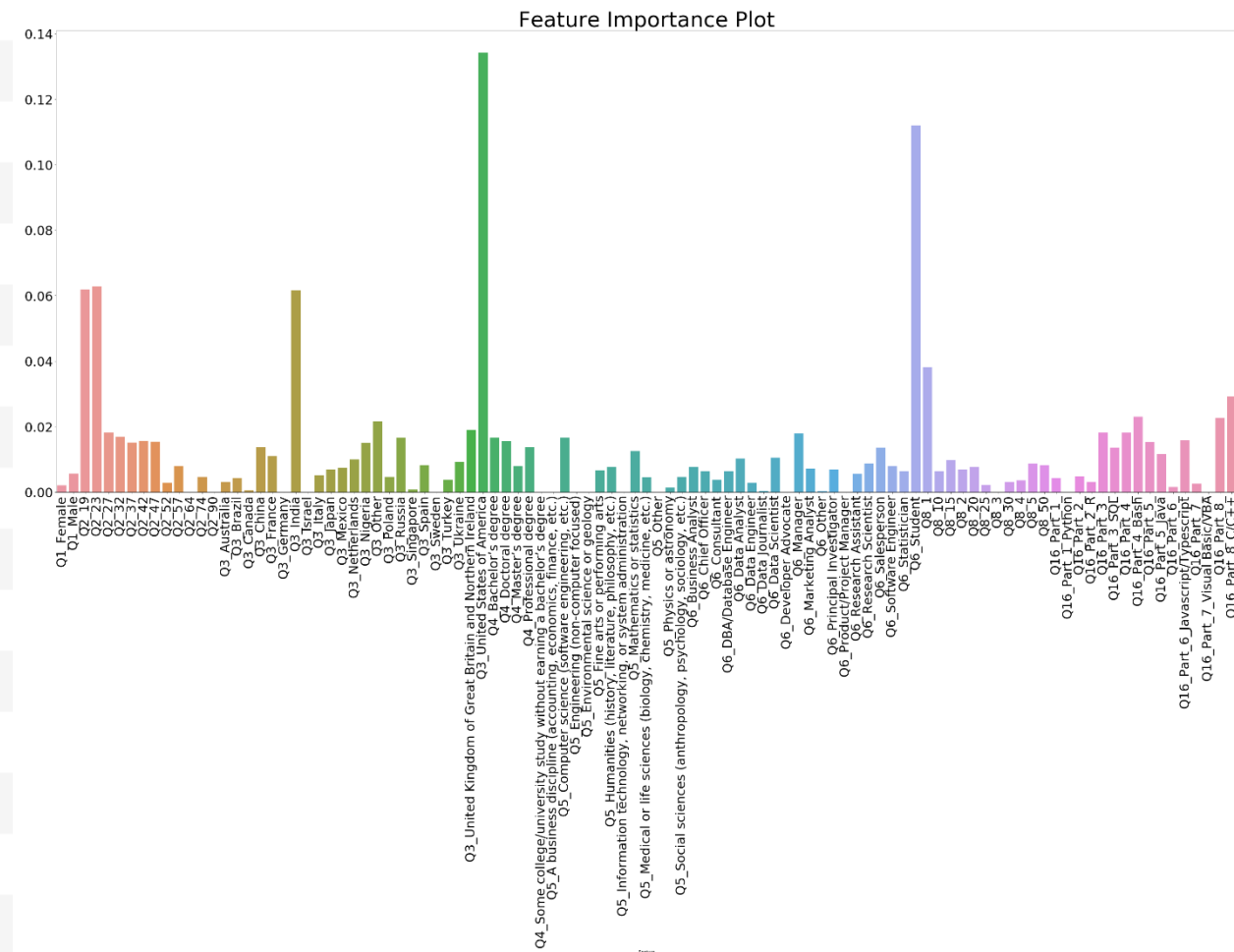
# Exploratory Analysis



Salary was heavily skewed to the right. Hence it was difficult to predict the higher classes. Salaries were divided into 4 brackets.

Salary varied for the different features selected. They were also separately displayed for males and females as seen.
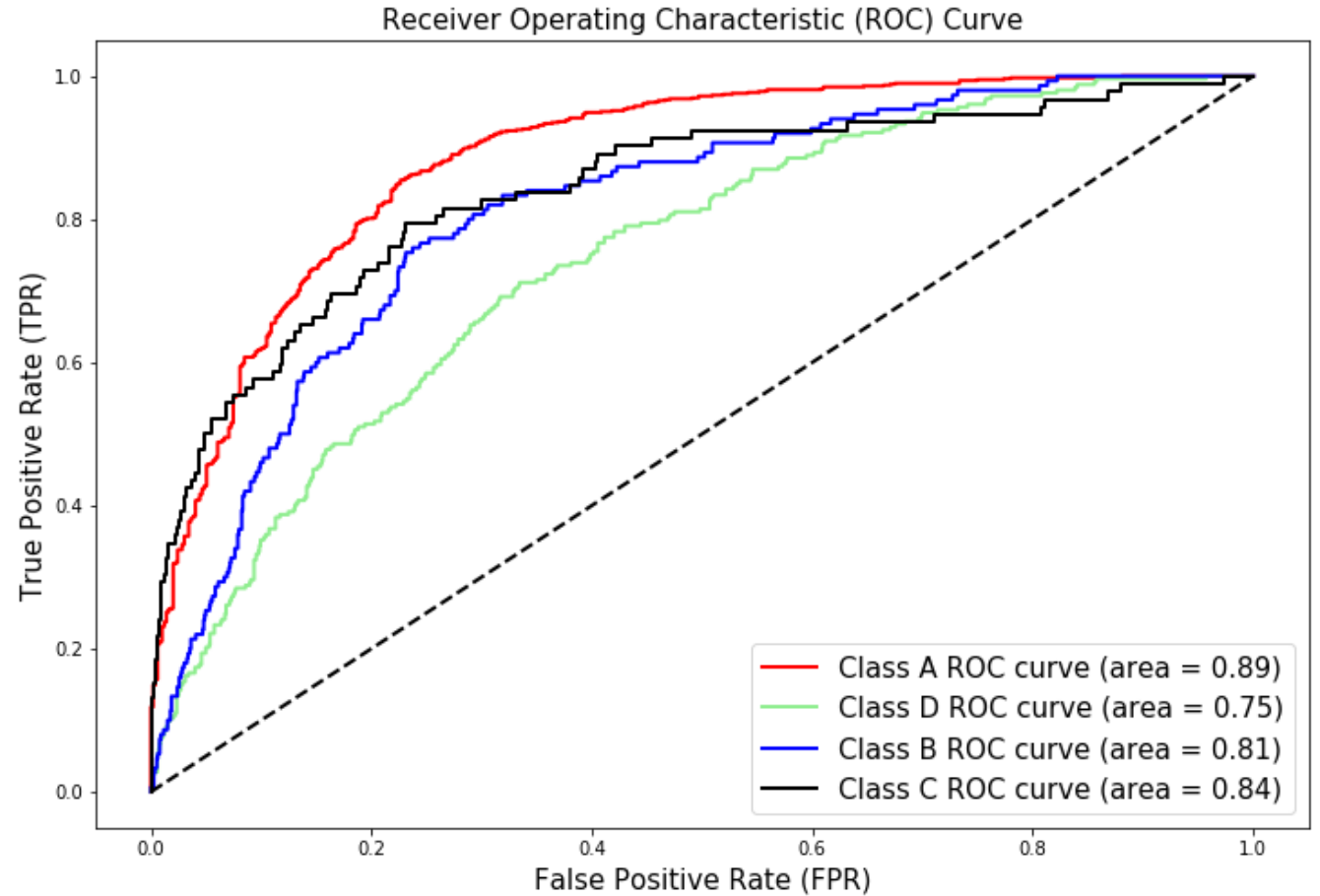
# Feature Selection and Importance
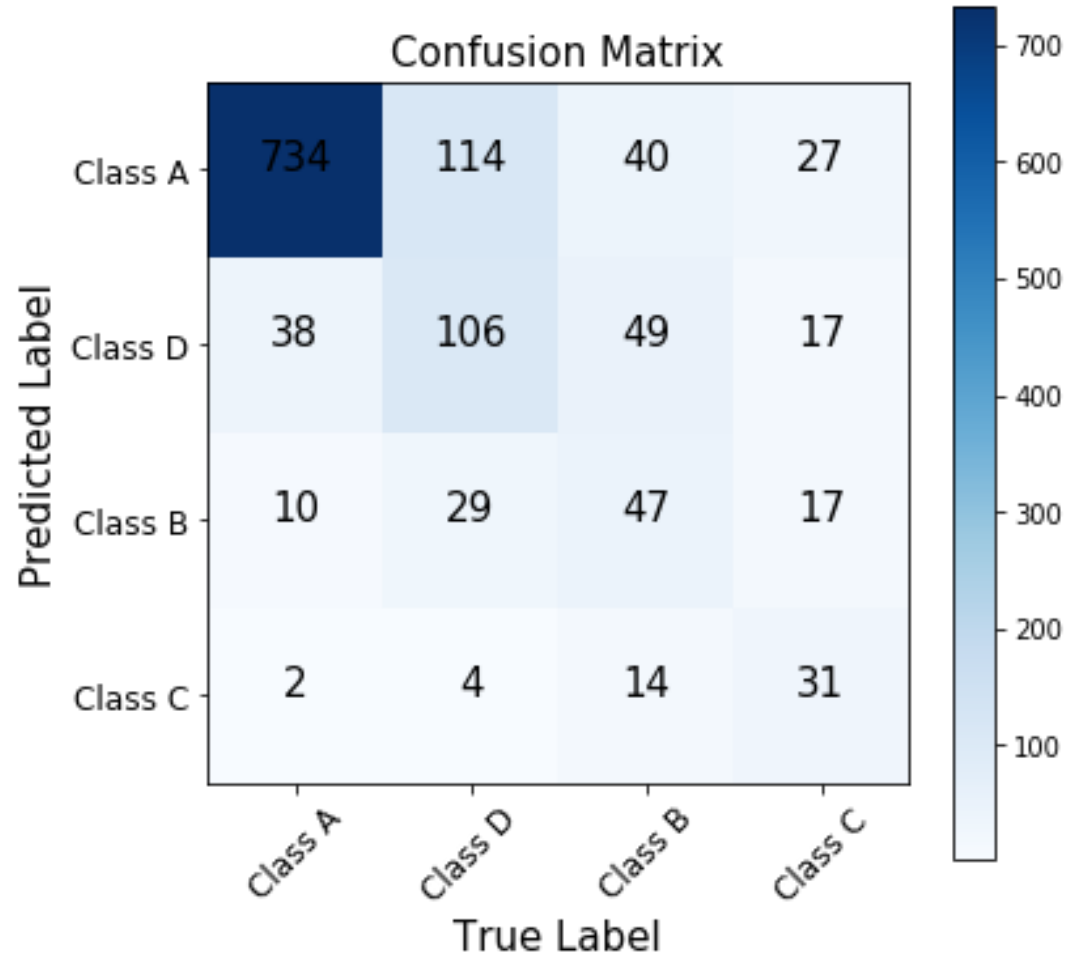


Lasso (l1) Regression was used to remove all the features with zero coefficients. Note that the features were one-hot encoded.

Feature Importance Plot using Mutual Information (MI).
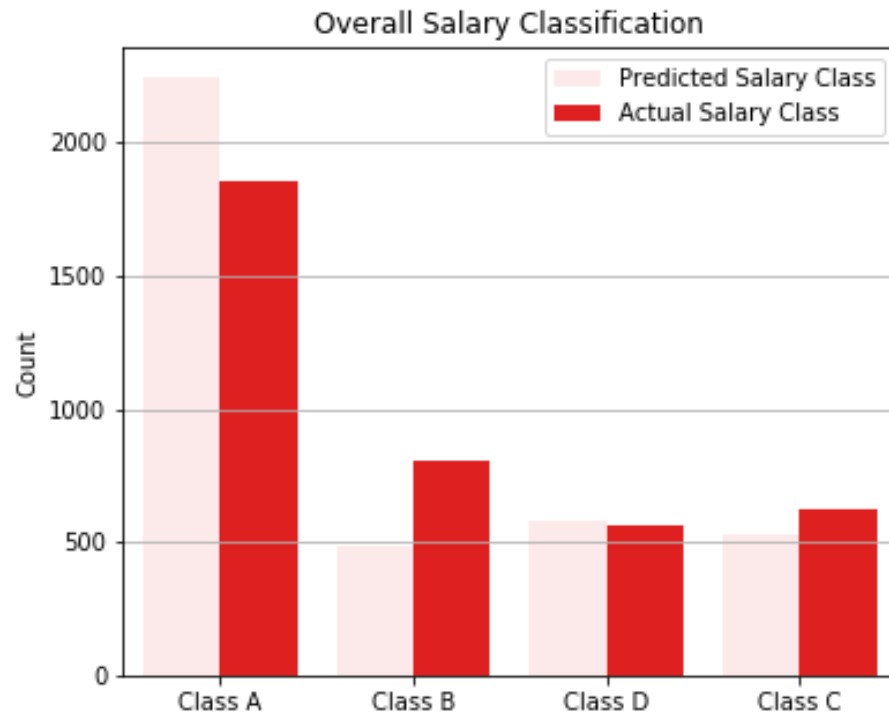Student and USA had the biggest impact on Salary.

# Model Implementation



Multinomial Logisitic regression was used. The accuracy from the 10-folds was **68%**

OneVsRest Classifier was used to classify the different salary brackets. Class D showed the worst performance whereas Class A showed the best.

# Model Testing



Overall Salary Classification

- Predicted Salary Class
- Actual Salary Class

Class A- 0-30k
Class B- 30k-60k
Class C- 60k-95k
Class D- 95k+

This model got an accuracy of **69.12%** on the training set
This model got an accuracy of **66.95%** on the testing set



Error for the model

- Test Data
- Training Data

C=1 onwards gives low error.
In our case, the difference between the training and test error is not great, hence the variance is low. However, the bias is high since the training error itself is quite high. This is understood and expected given the bias in the data for the lower salary classes.

# Summary

- 8 questions from the original survey were considered for the model.

- The features were cleaned and one-hot encoded.

- The salary target variable was differentiated into 4 classes.
  - This was done because the original 18 salary brackets were showing poor accuracy.
  - Salary distribution is heavily skewed to the left and has limited data for higher salary classes.
  - Grouping the salary brackets made the model accurate; however, the model was not able to predict exact classes (from the original count of 18).

- Multinomial logistic regression is used and the accuracy achieved is around 68% via 10-fold.

- OneVsRest Classifier is used to classify the different brackets and obtain the ROC curve.

- Model is not overfitting nor underfitting looking at the training and test error.

- Model has exhibited low variance and high bias (towards the lower salary bracket).