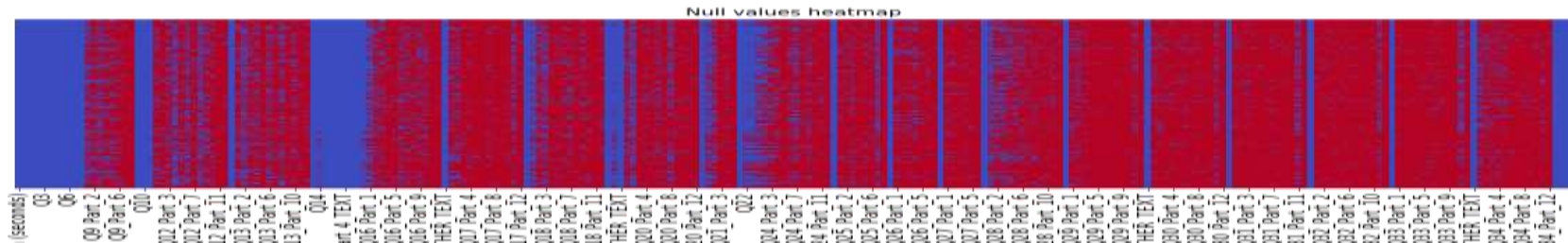# MIE 1624
# ASSIGNMENT 1
# SALARY CLASSIFICATION

Anshul Verma
1004730703
Date: 28$^{th}$. February. 2020
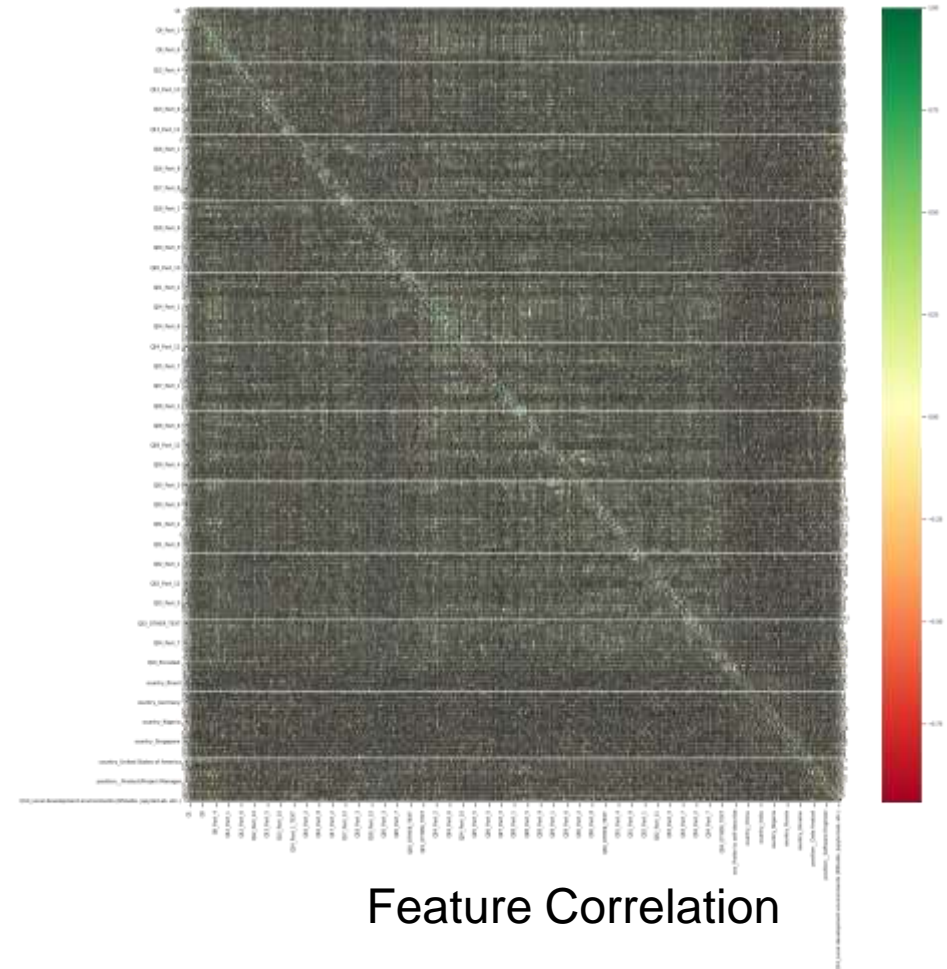
# Q1. Data Cleaning

- Most of the features in the data-set where binary either they had an entry or were null. So I encoded all these columns in 0-1 form.

- Then columns like sex, country, job-positions were treated as nominal categories and one-hot encoding was used for these types of columns.

- Other columns like education level, age, experience, Use of ML etc. were treated as ordinal categorical data and were encoded in a way that preserves the order.

- After encoding all the major columns there were some Nan values still present in the data-frame so the rows with Nan were removed from the dataset (I did this because the number of Nans were very small)

Heat map of data-frame without cleaning
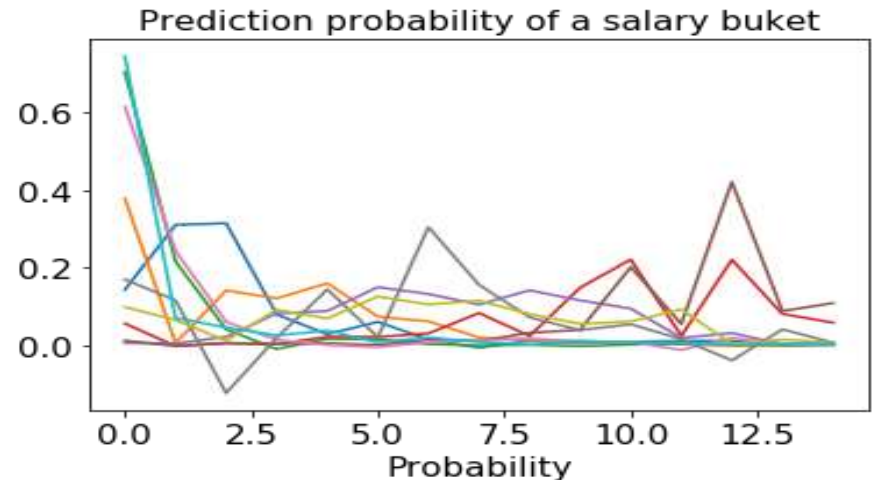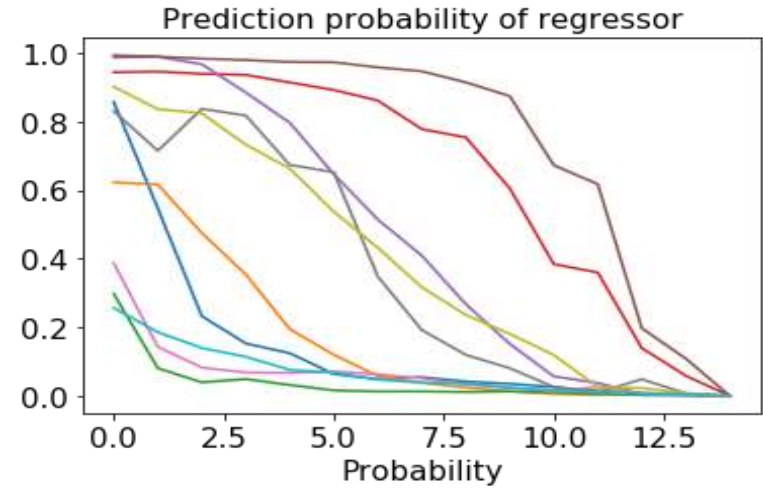


Null values heatmap

# Q2 & 3. Exploratory Analysis & Feature Selection

- For Exploratory data analysis the final salary bucket vs features which seem more influential is plotted and presented in the jupyter notebook.

- Correlation between features was used for feature selection

- Based on Correlation top 200 features were selected and used for the purpose of logistic regression.

- Finally the features and targets were seperated and features were further used to build an ordinal logistic regression model.
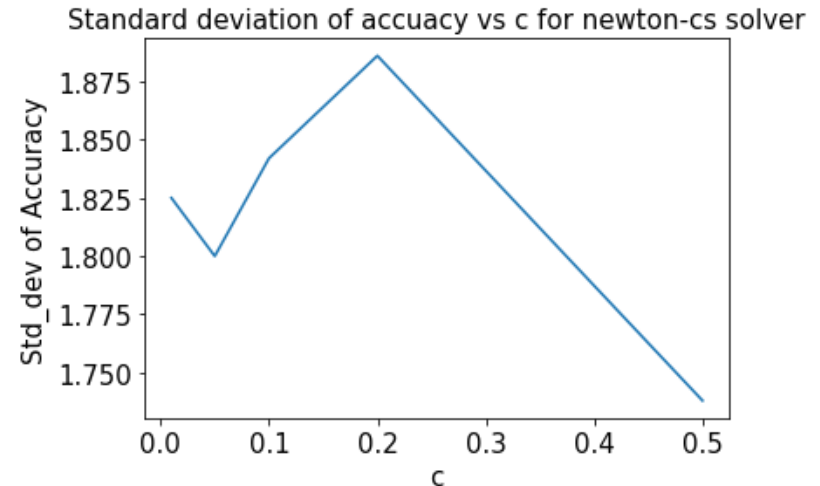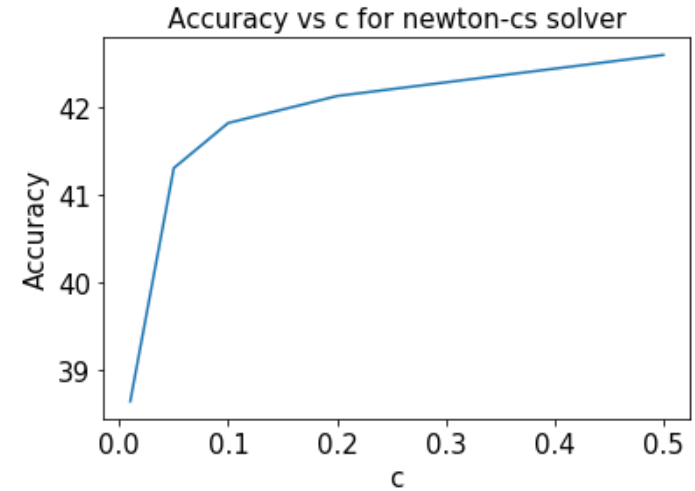


Feature Correlation

# Q4. Model Implementation

- 14 different binary linear regressor were used to classify all the classes.

- The classification for every class was of the form 0(<=salary_bucket) and 1(>salary_bucket).

- Using probabilities for this mixed classification probabilities of each class was identified.

- Then Based on the each classes probability predictions were made.

- The models were then used for 10-fold cross validation to make sure that the model doesn't overfit. (Making sure standard deviations of accuracy is low )



Prediction probability of regressor



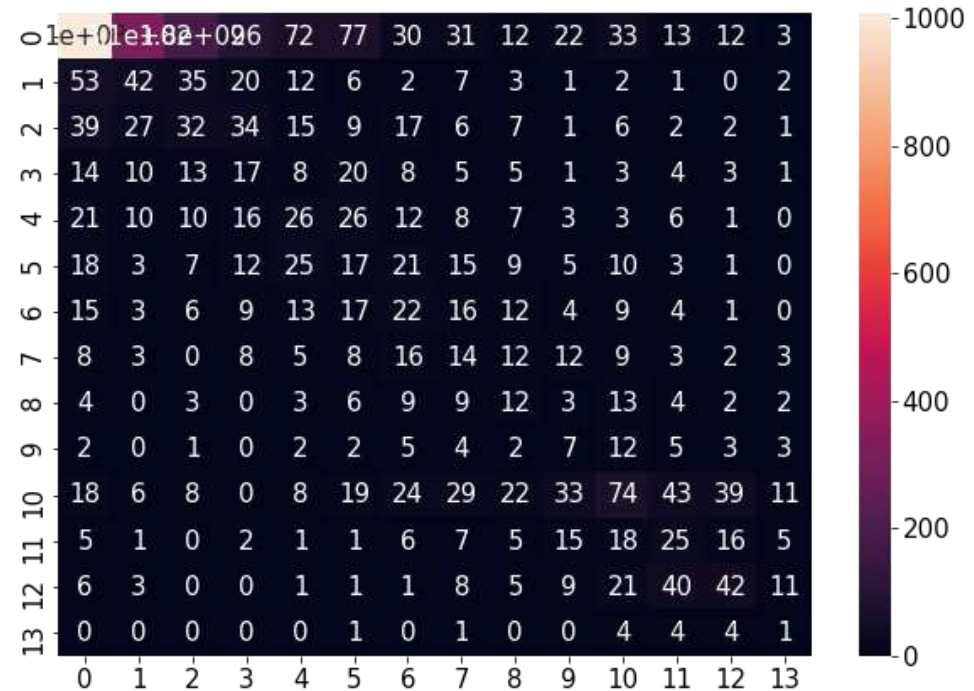Prediction probability of a salary buket

# Q5. Model Tuning

- Grid search on some possible hyper-parameter list was performed.

- There can be more hyper-parameter like regularize type, number of iterations. But for the purpose of assignment only L2 regularizer was used and max_iter was set to 100.

- The final selected model had a 10-fold average accuracy of 42.6% with standard deviation 1.73%.

- This model was finally used to make predictions on the test set.



Accuracy vs c for newton-cs solver



Standard deviation of accuacy vs c for newton-cs solver

# Q4. Testing

- The final selected model was then used to test the points in the test-set.

- The model was found to be 35.9% accurate which is pretty close to training error(42.6%). This shows that generalization was decent as test data was not seen while training

- The accuracy is not very high because of a lot of classes and simplicity of the ordinal logistic regression model.

- Using some other classification approach can increase the



Confusion Matrix of Test Prediction