

## Report

4a. Accuracy that I got from decision tree itself is: (Accuracy: 0.9174511223750905)

Accuracy that I got from Random Forest tree itself is: (Accuracy: 0.9377262853005068)

Total runtime of the Random Forest from scratch is 1828.573861837387 seconds

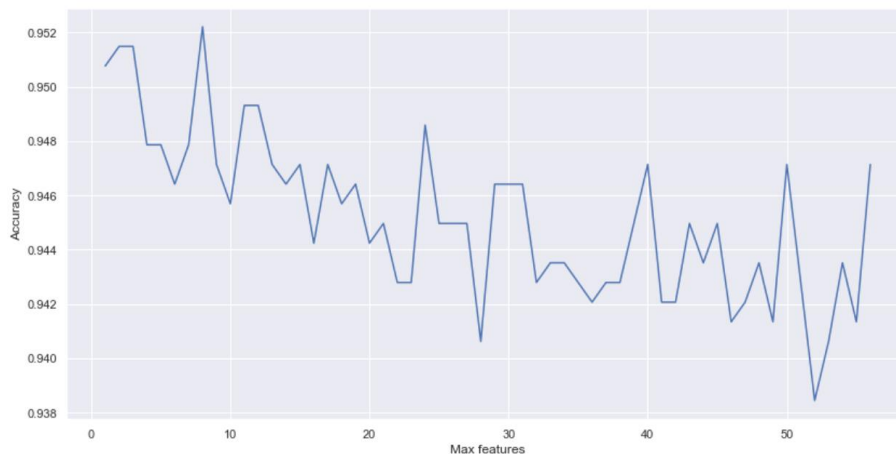
Accuracy that I got from Random Forest using sklearn itself is:

(Accuracy:0.9500362056480811)

Total runtime of the Random Forest sklearn is at most 4 seconds

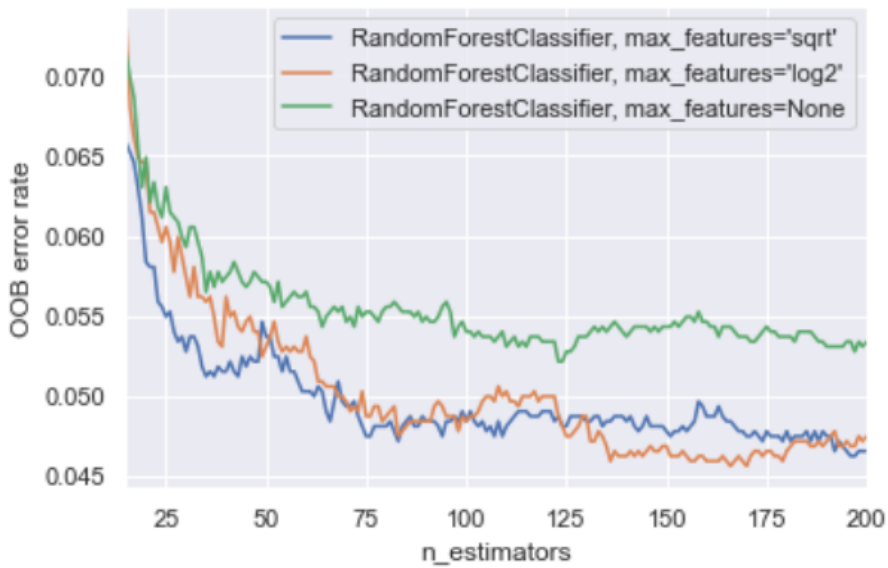
From the above results My random forest performed better while I implemented using sklearn as my accuracy is more and runtime is less for 100 trees using information gain.

4b.

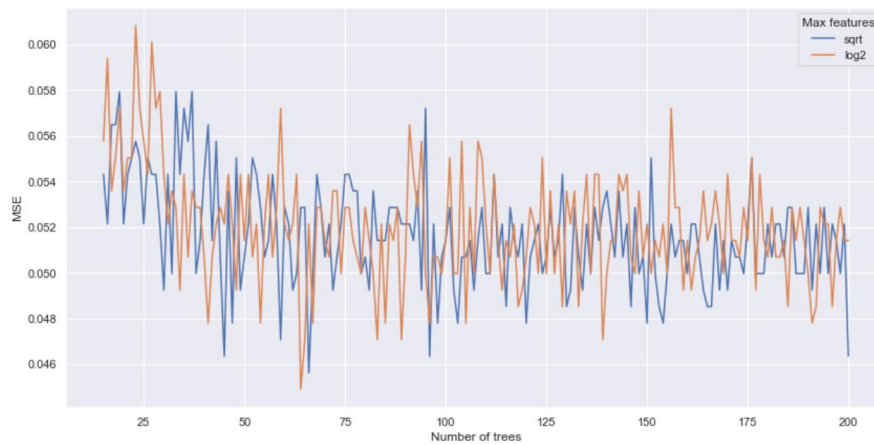


I drew the graph between maximum features and accuracy by implementing for loop. One take-away is Random Forest works better just like what Breiman suggested. When  $m \ll \text{number of predictor}(M) \text{ variables}$ . Breiman suggests three possible values for  $m$ :  $\frac{1}{2}\sqrt{M}$ ,  $\sqrt{M}$ , and  $2\sqrt{M}$ . Highest is observed at  $\sqrt{M}$  which is 'sqrt' aka "auto" for max\_features in sklearn.

4c. I have performed analysis for OOB and error and drew a graph for 3 m values which are “auto”, “log2” and “None” I.e., 3 inbuilt max\_features of sklearn. The plot is something like this:



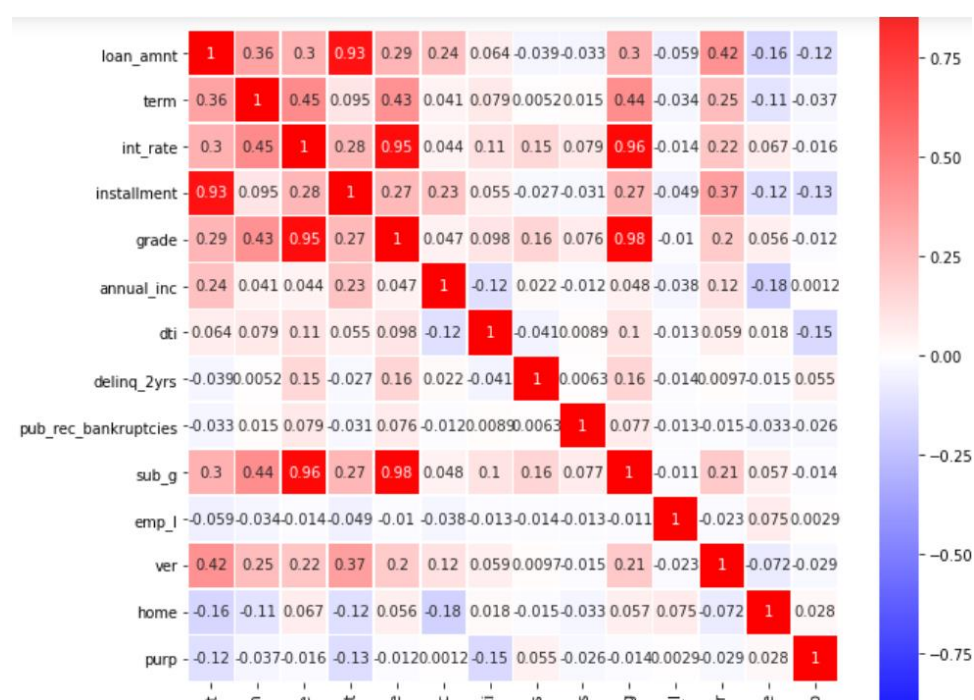
Plot between MSE and N\_estimators for selected m features:



5a. My first step in Data-cleaning is removing the columns that 24999 null values. I have removed the columns "tax\_liens","chargeoff\_within\_12\_mths","collections\_12\_mths\_ex\_med" as it has only one unique value which is 0 and it wouldn't be useful, Also dropping the columns which aren't useful for analysis 'url', 'desc', 'policy\_code','earliest\_cr\_line', 'emp\_title','id','next\_pymnt\_d', 'title', 'total\_rec\_int', 'total\_rec\_late\_fee', 'total\_rec\_prncp', 'zip\_code','issue\_d'.

-> 'loan\_amnt','term','int\_rate','installment','grade','sub\_grade','emp\_length','home\_ownership', 'annual\_inc','verification\_status','purpose','dti','delinq\_2yrs','pub\_rec\_bankruptcies','loan\_status', These are the columns that I felt plays the crucial role for predicting whether the loan is Default or not.

I plotted heatmap for these features by changing the object type columns to category type columns



I found grade and sub-grade column have large correlation with int\_rate so I have dropped them.

Null values that I found in my new data frame are:

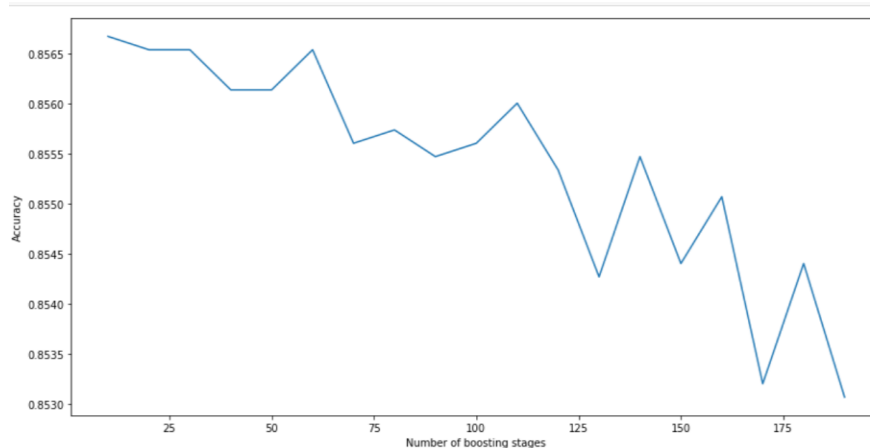
	Count	Percent
emp_length	41	0.164007
pub_rec_bankruptcies	417	1.668067

In emp\_length ,I filled the null values with 0 assuming that the borrower hasn't worked many years for his data to be recorded. For Pub\_rec\_bankruptcies I filled the Nan values with 0.0 Which are 94.3% values.

I have 'Charged Off', 'Fully Paid', 'Current values' in my loan\_status , I gave default to charged-off ' and rest all I assumed it to be not-Default.

By hyperparameter tuning I got best results when learning\_rate = 0.05, n\_estimators = 10, subsample = 0.5, max\_depth=5 which 0.862964

I found accuracy decreased as the no. Of trees increasing. Here's the graph for it.



Classification result for the best model:

	precision	recall	f1-score	support
0	0.86	1.00	0.92	6425
1	0.00	0.00	0.00	1075
accuracy			0.86	7500
macro avg	0.43	0.50	0.46	7500
weighted avg	0.73	0.86	0.79	7500

Classification accuracy using simple Decision Tree

	precision	recall	f1-score	support
0	0.87	0.86	0.86	6425
1	0.19	0.20	0.20	1075
accuracy			0.76	7500
macro avg	0.53	0.53	0.53	7500
weighted avg	0.77	0.76	0.77	7500

Gradient Boosting has an accuracy of 86 percent while Simple decision tree gave 76 percent accuracy this show how powerful Gboosting is