

Tree methods

Denis Tsvetkov and Lukas Ostrovskis

January 2023

1 A single decision tree model

5 feature model

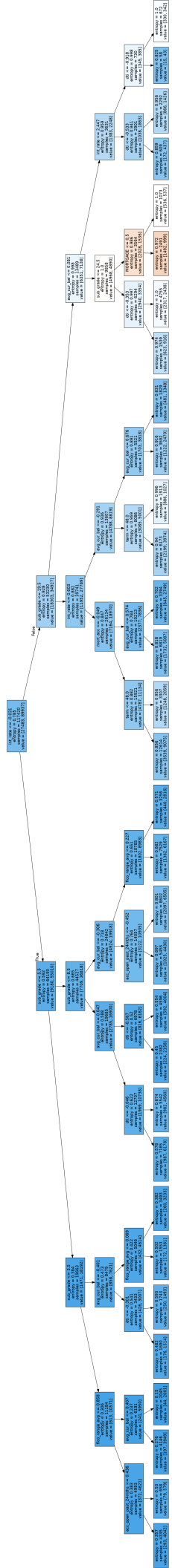
We trained the decision tree model with the 5 most important features identified by the regression model with the elastic net regularization as their importance looks the most plausible based on both model performance and feature definitions. The features are: *sub_grade*, *int_rate*, *WY*, *pub_rec_bankruptcies*, *mo_sin_old_il_acct*.

The feature that the model chose for splitting at the root node was *int_rate*. Below is a table that shows the information gain that could be obtained by each feature at the root node.

Feature	Threshold value	Expected entropy	Information gain
sub_grade	11.5	0.74483	0.04018
int_rate	-0.031	0.74456	0.04045
WY	N.A.	0.78499	0.0000173
pub_rec_bankruptcies	0.838	0.78428	0.00074
mo_sin_old_il_acct	-0.879	0.78439	0.000627

Table 1: Information gain from features to determine the root node

The plot of the final decision tree is provided on the next page. It is important to note that allowing the model to reach a depth of 5 results in a pretty big tree model, so it might be necessary to zoom in to inspect the individual nodes.



Furthermore, we decided to test 5 different prediction thresholds on our tree classifier - 0.5, 0.6, 0.7, 0.8, 0.9. The table with the confusion matrices for each of the thresholds can be found below.

0.5		0.6		0.7		0.8		0.9	
1.15	22.29	5.17	18.27	9.39	14.06	17.88	5.57	22.49	0.96
1.02	75.53	6.15	70.4	13.58	62.97	38.41	38.14	63.69	12.86

Table 2: Confusion matrices for the decision tree model with different prediction thresholds

This table is quite interesting to analyze in the context of our scenario.

On the one hand, it might be tempting to conclude that the model with a threshold of 0.5 is the best since it results in the highest accuracy (76.68%). However, it is important to acknowledge that it also produces a significant amount of false positives (predicting no default but the loan actually defaulting). This can be explained by the imbalance in training data: over 70% of the training samples are loans that do not default, hence the model is biased to guess that a loan is usually good.

On the other hand, we argue that it is obvious that a high percentage of false positives is expensive for loan providers but we can also support this claim by referring to Hull [1], who states at the end of Section 3.11 that revenue can be maximized by picking a threshold such that $1000TP - 4000FP$ is maximized. We tried to maximize the formula proposed by Hull by trying each of the thresholds and the resulting values are -13630, -2680, 6730, 15860 and 9020 respectively. Therefore, the performance metric proposed by Hull seems to indicate that it is most optimal to pick a threshold of 0.8 for the model.

While this looks counter-intuitive, as the accuracy is only 56%, it results in significantly more false negatives which cost less than false positives and thus the resulting model would most likely maximize the revenue for a loan provider but result in a stricter process to give out a loan.

Consequently, we decided to pick the 0.8 threshold. Table with the confusion matrix, true positive rate, and precision can be found below.

Confusion matrix		TPR	Precision
17.88	5.57	0.4982	0.8725
38.41	38.14		

Table 3: Confusion matrix, TPR, and precision for threshold = 0.8

To gain even more insight into the performance of the model we also generated the ROC curve for the model. The ROC curve plot is shown below.

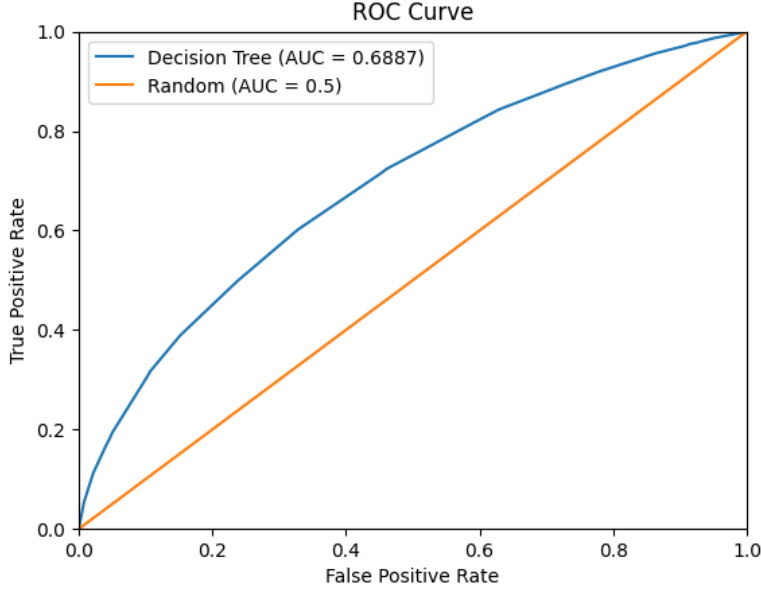


Figure 1: ROC for the Decision Tree (5 features)

It is visible that the model obtains an AOC of 0.6887.

10 feature model

To see if the number of features was an important factor in the analysis we redid it using 10 features. For this, we decided that it might be worthwhile to train a Random Forest model and look at the variable importance diagram that it generates. The features there are obviously effective for generating decision trees and might give us further improvement. According to the diagram, the 10 most important features are *int_rate*, *sub_grade*, *term*, *fico_range_avg*, *dti*, *acc_open_past_24mths*, *mort_acc*, *avg_cur_bal*, *Not Verified*, *MORTGAGE*. We now tried using these 10 features to train a decision tree model.

We again tried to maximize the formula proposed by Hull by trying each of the aforementioned thresholds and the resulting values are -13630, -2600, 8360, 16980, and 9380 respectively. We can again deduce that 0.8 is the optimal threshold. Table with the confusion matrix, true positive rate, and the precision for the optimal threshold can be found below.

Confusion matrix		TPR	Precision
16.76	6.68	0.5708	0.8673
32.85	43.7		

Table 4: Confusion matrix, TPR and precision for threshold = 0.8

While the optimal threshold does not change and remains 0.8, it is clear from the values that the new model is better at maximizing the formula and must therefore be better than the 5-feature decision tree. We can also support this claim by looking at the AUC of the model, which is 0.6957 and is greater than the 0.6887 obtained by the 5-feature model. The ROC curve plot is shown below.

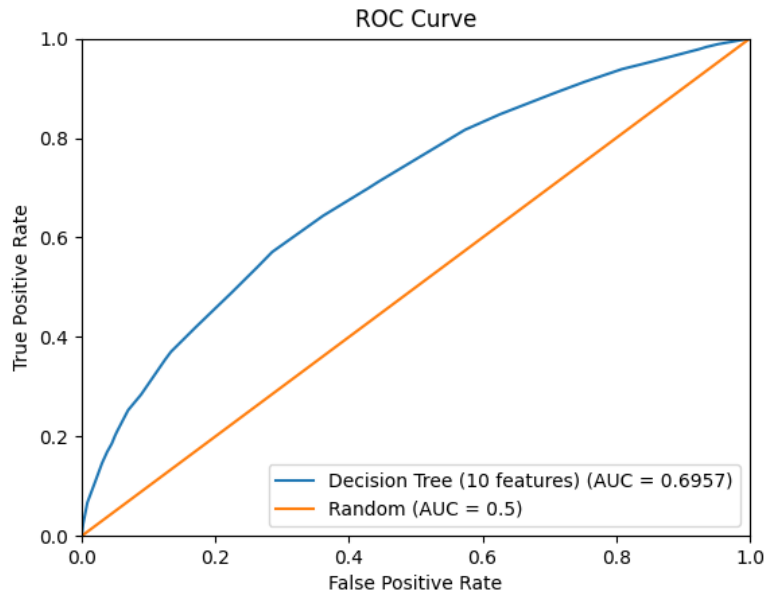


Figure 2: ROC for the Decision Tree (10 features)

Consequently, we conclude that the 10-feature decision tree model is the better one. It has a larger AUC and achieves a better performance as measured by the performance formula we're using.

2 Tree ensemble methods

To create a *RandomForestClassifier*, we looked at the documentation of the random forest classifier class¹ and decided to explicitly provide these parameters:

- `criterion` (the measure for splitting; we'll use "entropy" as we did for decision trees)
- `max_depth` (maximum depth of the tree; important for the bias-variance tradeoff)
- `oob_score` (whether to use out-of-bag samples; necessary for calculating the OOB error)
- `random_state` (controls the randomness; necessary for consistent results)
- `max_features` (number of features to consider for a split; we picked 2 values as defined in the assignment)
- `n_estimators` (number of trees in the forest; 100 by default, worthwhile to try more)
- `warm_start` (reuse the solution of the previous call to reuse and add more estimators; useful to see how OOB error changes based on the number of estimators)

First of all, we ran a grid search² for the *max_depth*, *max_features*, and *n_estimators* parameters to get an initial idea of good values for the hyperparameters. The search resulted in the following values:

- `max_depth = 20`

¹<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>

²https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html

- `max_features = "sqrt"`
- `n_estimators = 500`

Based on these findings, we decided to try the following parameter combinations for growing forests:

- `max_depth = 10` and `max_features = "sqrt"`
- `max_depth = 20` and `max_features = "sqrt"`
- `max_depth = 10` and `max_features = "log2"`
- `max_depth = 20` and `max_features = "log2"`

Furthermore, to investigate the effect of the bootstrap sample size we chose an interval of [100, 500] for the sample sizes with a step of 100. In order to evaluate the performance of the models, we created a graph utilizing OOB estimates in relation to the sample size. The resulting plot is provided below.

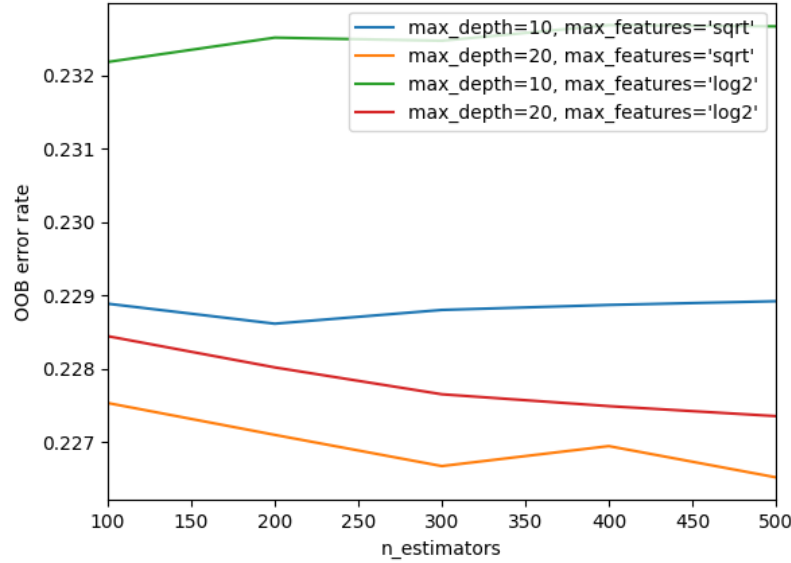


Figure 3: Performance of Random Forest models based on their OOB error estimates

Firstly, it is immediately noticeable that the models with a maximum depth of 10 are inferior to those with a depth of 20. This could be explained by the large amount of data and features we are working with: lower tree depth can result in slight underfitting and is not enough to capture the essence of all the information.

Secondly, we can see that the optimal bootstrap sample size for trees of depth 10 is around 200. On the other hand, for trees of depth 20, it keeps decreasing for the entire interval and most likely converges at an optimal value that is larger than 500.

As a result, the best Random Forest classifier seems to be the one with a maximum depth of 20, maximum features equal to the square root of all features, and a bootstrap sample size equal to 500.

Moreover, to evaluate the OOB error estimation, we compared the OOB error of the best model to the error estimate of our test set. Evaluating the model on the test set results in an error of 0.228, which is only marginally higher than the observed OOB error of 0.2267. We can thus conclude that the OOB error estimation was sufficiently accurate to be a reasonable metric in choosing the best model and did not result in us overfitting the dataset.

Lastly, after finding the best model, we looked at which variables had the largest influence on the model. The importance of a feature is calculated as the mean criterion decrease obtained by that feature. We used a convenient attribute of the *RandomForestClassifier*, *feature_importances_*³ which does exactly that. The variable importance diagram with the 10 most important variables can be found below.

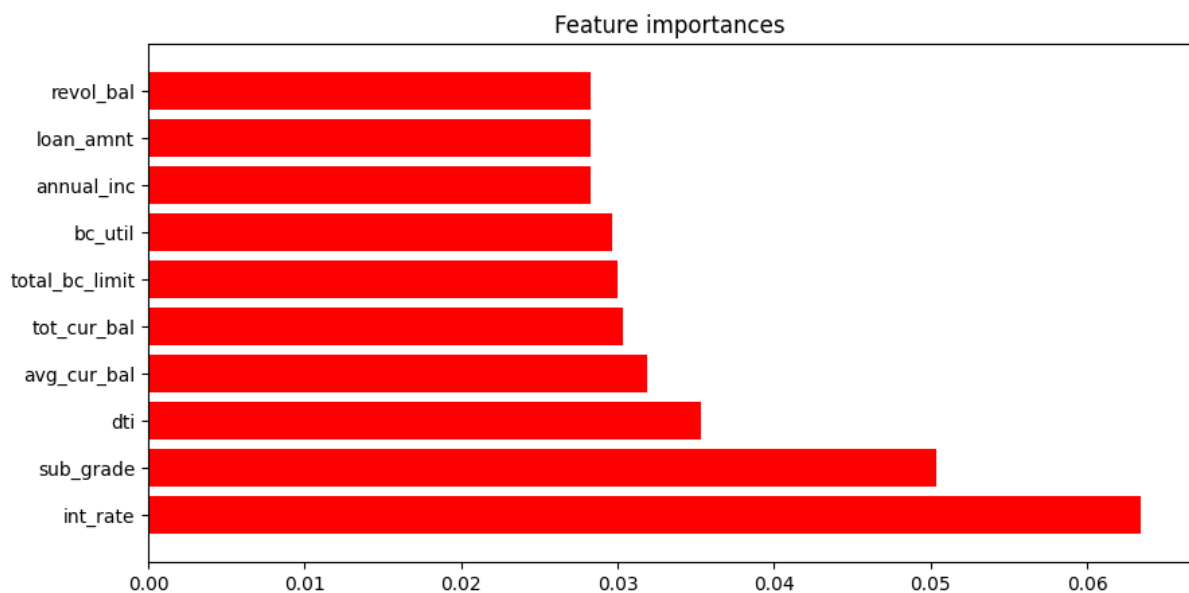


Figure 4: Variable importance diagram

As we can see, *int_rate* and *sub_grade* are the most important variables according to the model. This confirms that these are crucial features in determining the outcome of a loan. Furthermore, since these variables seem the most influential on the Random Forest model, we thought that it can be a good idea to try using them on a single Decision Tree model (1) as well, and as we already discussed, it did result in better performance.

3 Best final tree model

After comparing single Decision Tree and Random Forest models, we came to the conclusion that Random Forest models achieve better performance. For the Random Forest model, we again focused on maximizing the formula proposed by Hull [1] by trying each of the thresholds we used for single Decision Tree models and the resulting values are -10710, 460, 12720, 19220, 10020 respectively. We can again deduce that 0.8 is the optimal threshold but the obtained values are significantly larger

³https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html#sklearn.ensemble.RandomForestClassifier.feature_importances_

than the one found by a single Decision Tree using any of the thresholds. Table with the confusion matrix, true positive rate, and the precision for the optimal threshold can be found below.

Confusion matrix		TPR	Precision
18.23	5.22	0.5238	0.8849
36.45	40.1		

Table 5: Confusion matrix, TPR and precision for a Random Forest model with threshold = 0.8

Moreover, the AUC of the model is 0.7184, which is greater than the 0.6957 obtained by the 10-feature single Decision Tree model. The ROC curve plot is shown below.

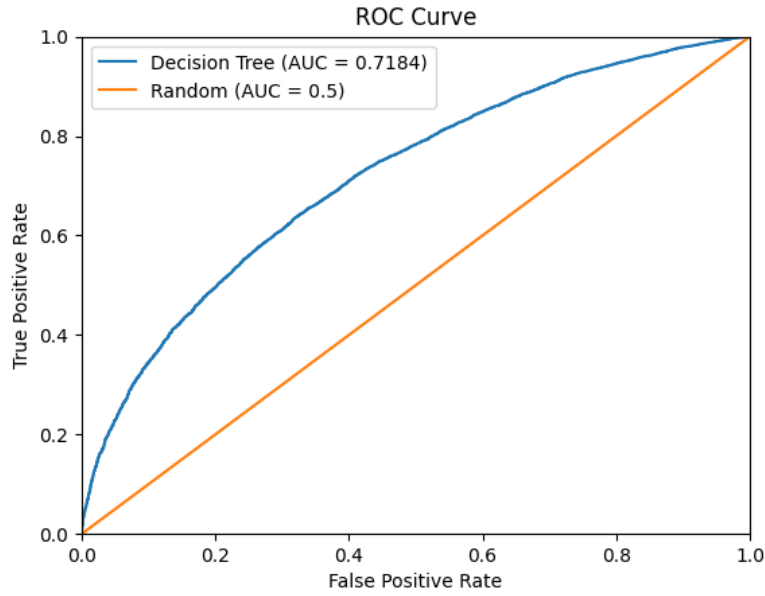


Figure 5: ROC for the Random Forest

References

- [1] Hull, J. C. (2019). Machine Learning in Business: An Introduction to the World of Data Science.