

Saranpat Prasertthum (655667271)

IE517 ML in Fin Lab

Module 7 Homework (Random Forest)

Out[1]:

[Click here to toggle on/off the raw code.](#)**Load Data**

Out[3]:

	LIMIT_BAL	SEX	EDUCATION	MARRIAGE	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5
ID										
1	20000	2	2	1	24	2	2	-1	-1	-2
2	120000	2	2	2	26	-1	2	0	0	0
3	90000	2	2	2	34	0	0	0	0	0
4	50000	2	2	1	37	0	0	0	0	0
5	50000	1	2	1	57	-1	0	-1	0	0

5 rows × 24 columns

Part 1: Random forest estimators

Out[5]:

```
GridSearchCV(cv=10, estimator=RandomForestClassifier(),
             param_grid={'n_estimators': [10, 50, 100]},
             return_train_score=True, scoring='accuracy')
```

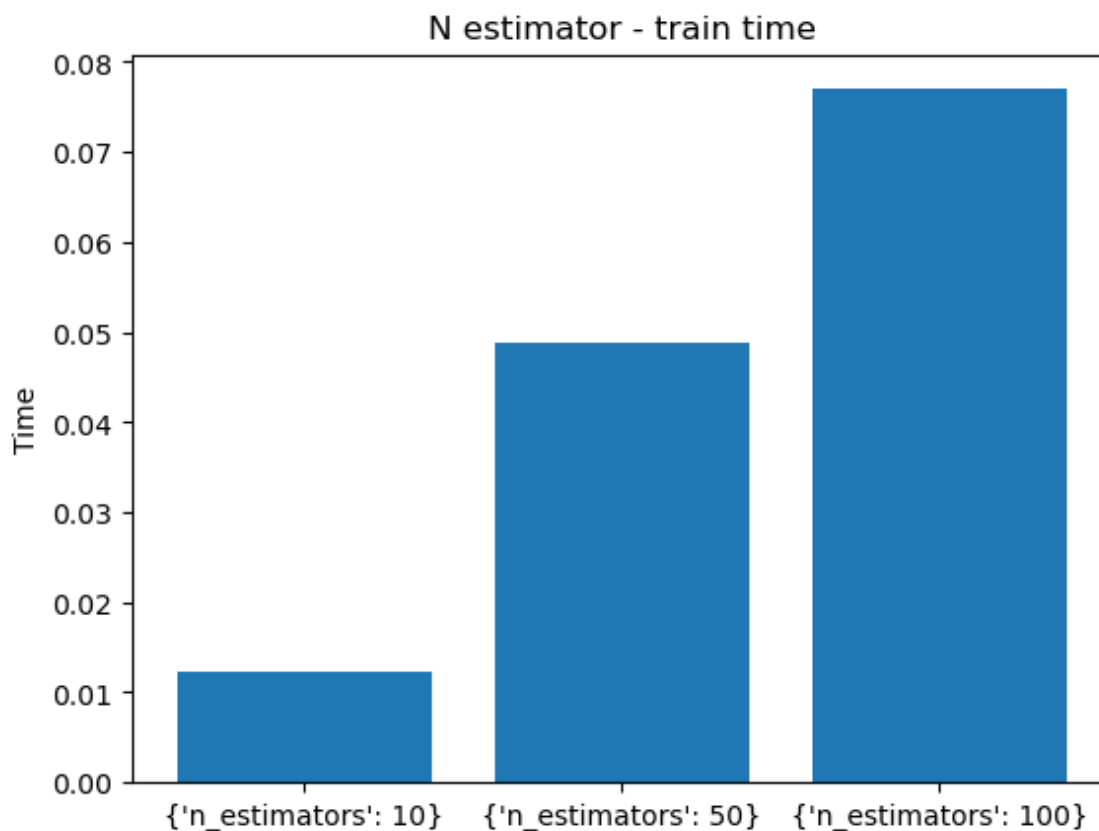
Out[6]:

	Params	Train Score	Test Score	Time
0	{'n_estimators': 10}	0.979978	0.807267	0.012284
1	{'n_estimators': 50}	0.998748	0.815400	0.048875
2	{'n_estimators': 100}	0.999315	0.817200	0.076984

Part 2: Random forest feature importance

```
1) PAY_0      0.08935668919578203
2) AGE       0.06696532281991693
3) BILL_AMT1 0.061449377578944954
4) LIMIT_BAL 0.05952946037937925
5) BILL_AMT2 0.05461007932687378
6) PAY_AMT1  0.05195440834030842
7) BILL_AMT3 0.05168286269081806
8) BILL_AMT6 0.05111619901788388
9) BILL_AMT4 0.050364675135451085
10) PAY_2     0.050191165474225745
11) BILL_AMT5 0.049101703951976765
12) PAY_AMT2  0.04708222260802307
13) PAY_AMT6  0.04617658105920907
14) PAY_AMT3  0.045396838633858405
15) PAY_AMT5  0.043521466552780864
16) PAY_AMT4  0.04303340041295496
17) PAY_3     0.02978275828837117
18) PAY_4     0.02146006628639994
19) PAY_5     0.021121021434264257
20) EDUCATION 0.0204601627837262
21) PAY_6     0.019689146426290622
22) MARRIAGE  0.014069990988382035
23) SEX       0.01188440061417859
```

Part 3: Conclusions



a) What is the relationship between `n_estimators`, in-sample CV accuracy and computation time?

ANS If N estimator increase it will also increase computational time.

Out[10]:

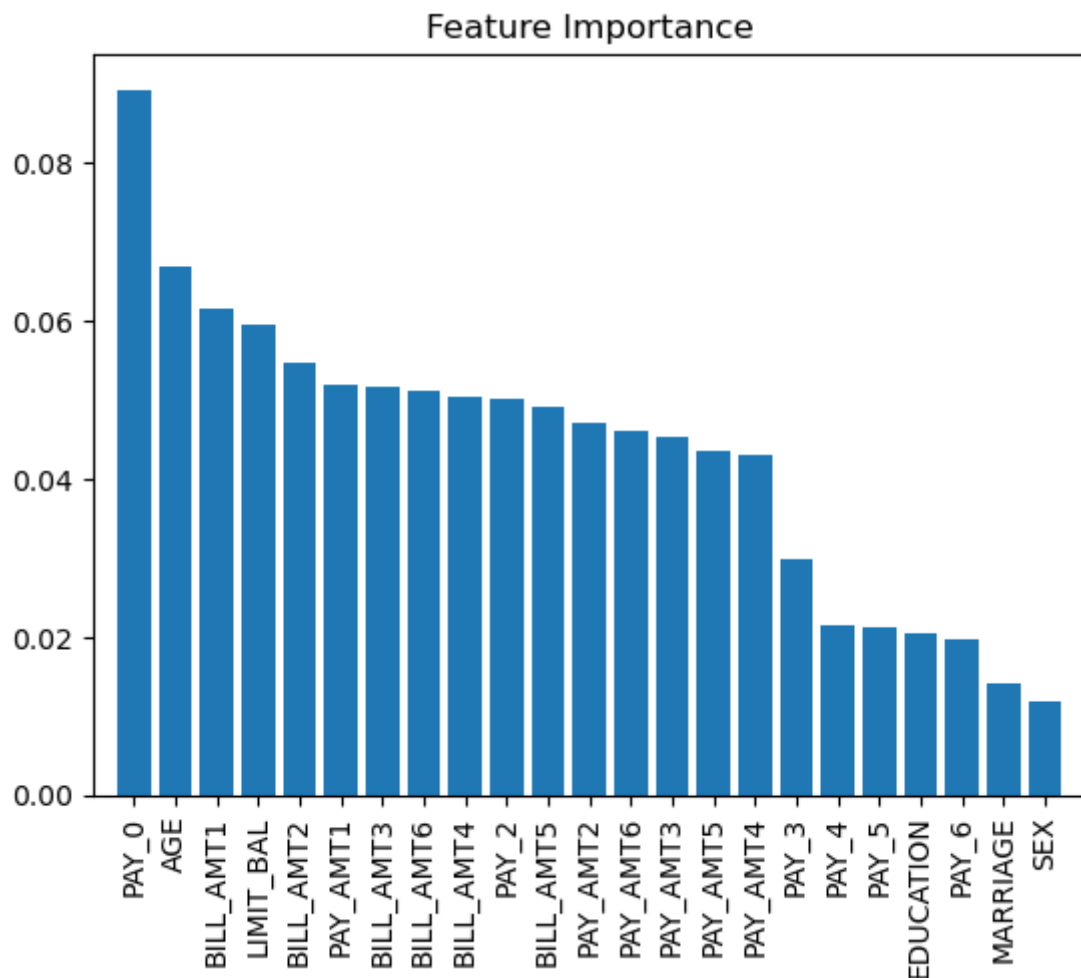
```
{'n_estimators': 100}
```

b) What is the optimal number of estimators for your forest?

ANS The best estimators is 100

Out[11]:

```
(-1.0, 23.0)
```



c) Which features contribute the most importance in your model according to scikit-learn function?

ANS The most importance feature in my model is "PAY_0"

d) What is feature importance and how is it calculated? (If you are not sure, refer to the Scikit-Learn.org documentation.)

ANS Feature importance in scikit-learn Random Forest models is calculated based on the decrease in node impurity caused by a feature, and the probability of reaching the node that uses the feature. This probability is determined by dividing the number of samples that reach the node by the total number of samples. The resulting values are then weighted and used to rank the importance of the features. Essentially, the higher the feature importance value, the more significant the feature is in predicting the target variable.

Impurity	Task	Formula	Description
Gini impurity	Classification	$\sum_{i=1}^C f_i(1 - f_i)$	f_i is the frequency of label i at a node and C is the number of unique labels.
Entropy	Classification	$\sum_{i=1}^C -f_i \log(f_i)$	f_i is the frequency of label i at a node and C is the number of unique labels.
Variance / Mean Square Error (MSE)	Regression	$\frac{1}{N} \sum_{i=1}^N (y_i - \mu)^2$	y_i is label for an instance, N is the number of instances and μ is the mean given by $\frac{1}{N} \sum_{i=1}^N y_i$
Variance / Mean Absolute Error (MAE) (Scikit-learn only)	Regression	$\frac{1}{N} \sum_{i=1}^N y_i - \mu $	y_i is label for an instance, N is the number of instances and μ is the mean given by $\frac{1}{N} \sum_{i=1}^N y_i$

Credit: <https://towardsdatascience.com/the-mathematics-of-decision-trees-random-forest-and-feature-importance-in-scikit-learn-and-spark-f2861df67e3#:~:text=Feature%20importance%20is%20calculated%20as,the%20more%20important%20the%20feature%20is,the%20more%20important%20the%20feature%20is>

Part 4: Appendix

My name is Saranpat Prasertthum

My NetID is: 655667271

I hereby certify that I have read the University policy on Academic Integrity and that I am not in violation.