

Australian Credit Approval Analysis

IBM Supervised Machine Learning: Classification

Course Assignment

Vladas Aleisius

2021-06-13

Main Objective

The purpose of this analysis is to examine how well different classifier methods are able to distinguish between approved and disapproved credit card applications. We are interested in performance metrics (accuracy, precision, sensitivity etc.) of different methods - we have selected logistic regression, decision trees and random forests. In other words, this analysis will focus on **predictability**. The analyzed dataset contains information about Australian credit card applications. All attribute names and values have been changed to meaningless symbols to protect confidentiality of the data. The source is this [link](#).

We will start with exploratory data analysis, to see how our variables are distributed and correlated, what transformations might be necessary. Since we will fit regularized logistic regression on our data, it is useful to have normally distributed, scaled features with no outliers.

After initial data pre-processing steps, it will be split into training and testing subsets. Yeo-Johnson transformation and standardization will be fit on the training subset, and applied to the testing subset as well. After that, cross-validation will be used to tune the hyperparameters for each model. The following measures will be assessed for each model:

- Classification report, containing precision, recall, accuracy and F1-score for both application classes
- Area Under (ROC) Curve (AUC) measure
- Confusion matrix
- Average prediction probabilities for each application group in the confusion matrix (true positive, false negative etc)
- True and false positive rates, and precision and recall values with multiple thresholds, used for model comparison

After that, we will discuss key insights and select the final model. Also, some suggestions about how to improve this study will be presented.

Data Description

We begin with importing necessary modules:

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
from sklearn.model_selection import StratifiedShuffleSplit, GridSearchCV
from sklearn.linear_model import LogisticRegressionCV
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix, classification_report, roc_auc_score, roc_curve, precision_recall_fscore_support
from sklearn.metrics import ConfusionMatrixDisplay, accuracy_score, precision_score, recall_score, f1_score
from sklearn.tree import DecisionTreeClassifier
from sklearn import tree
from sklearn.preprocessing import PowerTransformer
from matplotlib import pyplot as plt
from statsmodels.stats.outliers_influence import variance_inflation_factor
```

Let's read the credit dataset and examine the first 5 rows:

```
In [2]: credit = pd.read_table(r"C:\Users\Lenovo\Desktop\libm_machine_learning\3_supervised_machine_learning_classification\credit.head()")
credit.head()
```

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
0	1	22.08	1146	2	4	4	1.585	0	0	0	1	2	100	1213	0
1	0	22.67	7.00	2	8	4	0.165	0	0	0	0	2	160	1	0
2	0	29.58	1.75	1	4	4	1.250	0	0	0	1	2	280	1	0
3	0	21.67	11.50	1	5	3	0.000	1	1	11	1	2	0	1	1
4	1	20.17	8.17	2	6	4	1.960	1	1	14	0	2	60	159	1

In order to process dataset columns, we will assign them names as given in the description above (A1-A15):

```
In [3]: credit.columns = ['A1' + str(i) for i in range(1,16)]
credit.sample(5)
```

	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	A11	A12	A13	A14	A15
159	1	47.83	4165	2	14	5	0.085	0	0	0	1	2	520	1	0
143	1	64.08	20.000	2	14	8	17.500	1	1	9	1	2	0	1001	1
657	1	39.50	1.625	2	8	4	1.500	0	0	0	0	2	0	3171	0
598	1	24.08	0.875	2	7	4	0.085	0	1	4	0	2	254	1951	0
8	1	27.83	1.000	1	2	8	3.000	0	0	0	0	2	176	538	0

The dataset contains 6 numerical and 8 categorical attributes. According to the dataset description, the labels have been changed for the convenience of the statistical algorithms. For example, attribute 4 originally had 3 labels p.g.gg and these have been changed to labels 1,2,3.

- A1: 0,1 CATEGORICAL (formerly: a,b)
- A2: continuous.
- A3: continuous.
- A4: 1,2,3 CATEGORICAL (formerly: p,g,gg)
- A5: 1,2,3,4,5,6,7,8,9,10,11,12,13,14 CATEGORICAL (formerly: f,d,i,k,j,a,a,m,c,w,e,g,r,c,c)
- A6: 1,2,3,4,5,6,7,8,9 CATEGORICAL (formerly: f,d,d,j,b,b,v,a,o,h,z)
- A7: continuous.
- A8: 1,0 CATEGORICAL (formerly: t,f)
- A9: 1,0 CATEGORICAL (formerly: t,f)
- A10: continuous.
- A11: 1,0 CATEGORICAL (formerly: t,f)
- A12: 1,2,3 CATEGORICAL (formerly: s,g,p)
- A13: continuous.
- A14: continuous.
- A15: 1,2 class attribute (formerly: +,-) (negative class is presented as 0 instead of 2)

The summary below shows that this dataset contains only numeric (integer or float) variables, and no null values:

```
In [4]: credit.info()

Out[4]:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 690 entries, 0 to 689
Data columns (total 15 columns):
#   Column  Non-Null Count  Dtype
---  --
0   A1      690 non-null      int64
1   A2      690 non-null      float64
2   A3      690 non-null      float64
3   A4      690 non-null      int64
4   A5      690 non-null      int64
5   A6      690 non-null      float64
6   A7      690 non-null      float64
7   A8      690 non-null      int64
8   A9      690 non-null      int64
9   A10     690 non-null      float64
10  A11     690 non-null      int64
11  A12     690 non-null      int64
12  A13     690 non-null      int64
13  A14     690 non-null      int64
14  A15     690 non-null      int64 (12)
dtypes: float64(3), int64(12)
memory usage: 81.0 KB
```

However, as per the dataset description, 37 cases (5%) had one or more missing values. The missing values from particular attributes were:

- A1: 12
- A2: 12
- A4: 6
- A5: 6
- A6: 9
- A7: 9
- A14: 13

These were replaced by the mode of the attribute (categorical) and the mean of the attribute (continuous). The replacement was already done in the source data.

The dataset will be split into two - features (X) and the target variable (y):

```
In [5]: X = credit.drop('A15', axis='columns')
y = credit['A15']
```

The class distribution (variable A15) is given below, it is (almost) balanced:

```
+: 397 (44.5%)  CLASS 0
-: 383 (55.5%)  CLASS 1
```

As per the dataset description, the necessary features will be converted from integer type to categorical:

```
In [6]: col_list = ['A1', 'A4', 'A5', 'A6', 'A8', 'A9', 'A11', 'A12']
for col in col_list:
    if X[col].dtype == 'int64':
        X[col] = X[col].astype("category")
```

In [7]: X.dtypes

```
Out[7]:
A1      category
A2      float64
A3      float64
A4      category
A5      category
A6      category
A7      float64
A8      category
A9      category
A10     float64
A11     category
A12     category
A13     int64
A14     int64
dtype: object
```

Data Exploration

As per this [source](#), performing the following steps might improve the accuracy of logistic regression model:

- Removing potential outliers
- Making sure that the predictor variables are normally distributed. If not, you can use data transformations.
- Remove highly correlated predictors to minimize overfitting. The presence of highly correlated predictors might lead to an unstable model solution.

The following summary of the numeric variables shows that some of them might be skewed and have outliers - maximum values are large, compared to other parameters, and the standard deviations are often bigger than the means. For example, A14 has a few times bigger standard deviation than its mean:

```
In [8]: X.describe()
```

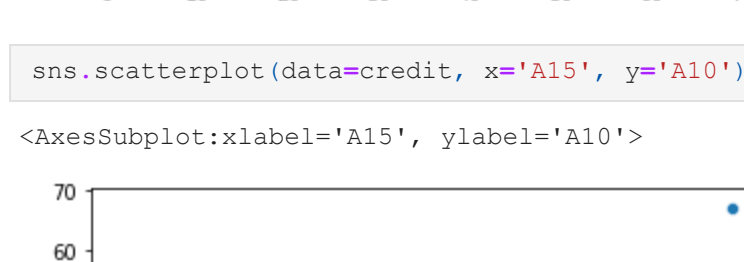
	A2	A3	A7	A10	A13	A14
count	690.000000	690.000000	690.000000	690.000000	690.000000	690.000000
mean	31.568203	4.758725	2.223406	2.400000	184.014493	1018.385507
std	13.853273	4.978163	3.346513	4.86294	172.159274	5210.102598
min	13.750000	0.000000	0.000000	0.000000	0.000000	1.000000
25%	22.670000	1.000000	0.165000	0.000000	80.000000	1.000000
50%	28.625000	2.750000	1.000000	0.000000	160.000000	6.000000
75%	37.707500	7.207500	2.625000	3.000000	272.000000	396.500000
max	80.250000	28.000000	28.500000	67.000000	2000.000000	100001.000000

The histograms below confirm that the numeric features are heavily skewed and have some outliers. The skew will be later dealt with Yeo-Johnson transformation and data standardization. However, the outliers will be kept in the dataset, since we are unaware of their nature.

Scatterplots demonstrate the differences in feature distribution for each application class. The biggest difference is for variable A14 - class 1 values vary much more, the outliers are very large. This might negatively impact the predictive power for class 1 cases.

```
In [32]: sns.distplot(x=credit['A2'], kde=False)
```

C:\Users\Lenovo\Anaconda3\lib\site-packages\seaborn\distributions.py:2551: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).



```
In [42]: sns.scatterplot(data=credit, x='A15', y='A2')
```

```
Out[42]: <AxesSubplot: xlabel='A15', ylabel='A2'>
```



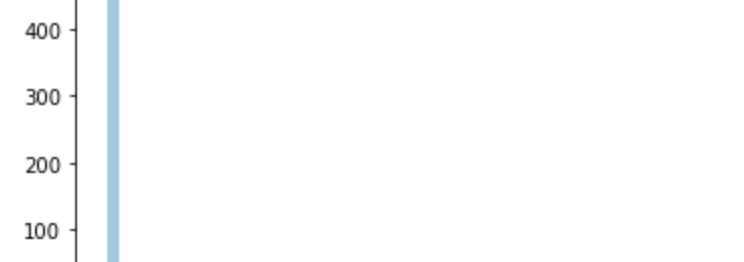
```
In [33]: sns.distplot(x=credit['A3'], kde=False)
```

```
Out[33]: <AxesSubplot: >
```



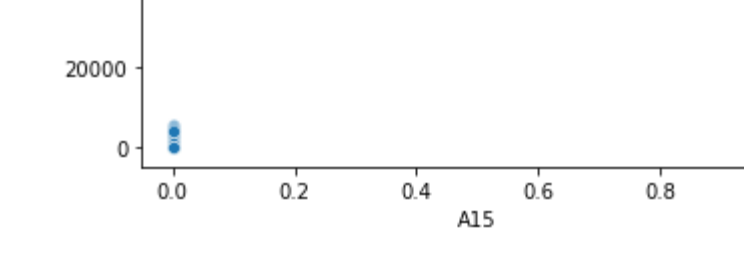
```
In [41]: sns.scatterplot(data=credit, x='A15', y='A3')
```

```
Out[41]: <AxesSubplot: xlabel='A15', ylabel='A3'>
```



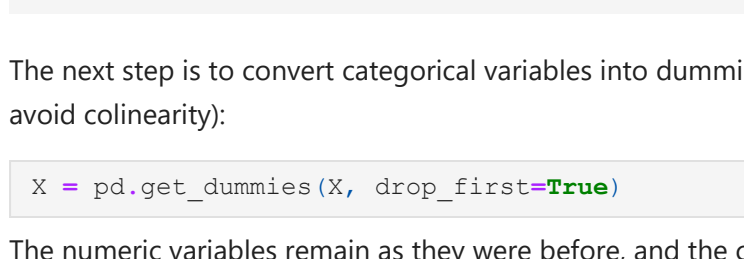
```
In [34]: sns.distplot(x=credit['A7'], kde=False)
```

```
Out[34]: <AxesSubplot: >
```



```
In [40]: sns.scatterplot(data=credit, x='A15', y='A7')
```

```
Out[40]: <AxesSubplot: xlabel='A15', ylabel='A7'>
```



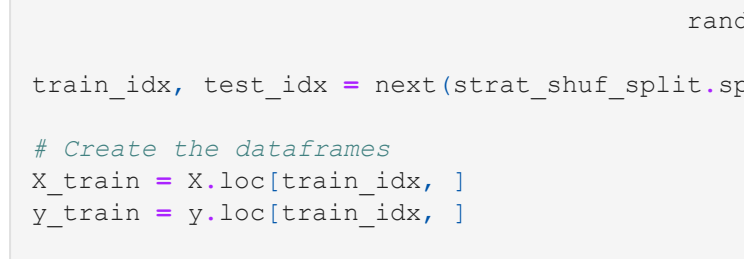
```
In [35]: sns.distplot(x=credit['A10'], kde=False)
```

```
Out[35]: <AxesSubplot: >
```



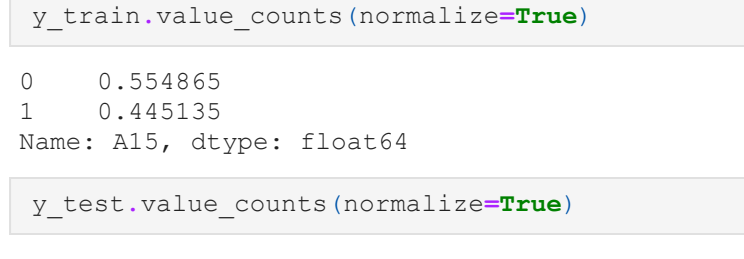
```
In [39]: sns.scatterplot(data=credit, x='A15', y='A10')
```

```
Out[39]: <AxesSubplot: xlabel='A15', ylabel='A10'>
```



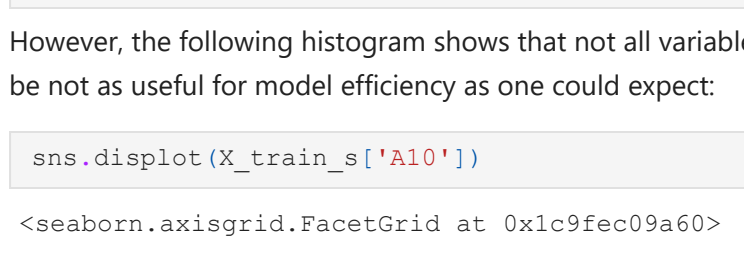
```
In [36]: sns.distplot(x=credit['A13'], kde=False)
```

```
Out[36]: <AxesSubplot: >
```



```
In [38]: sns.scatterplot(data=credit, x='A15', y='A13')
```

```
Out[38]: <AxesSubplot: xlabel='A15', ylabel='A13'>
```



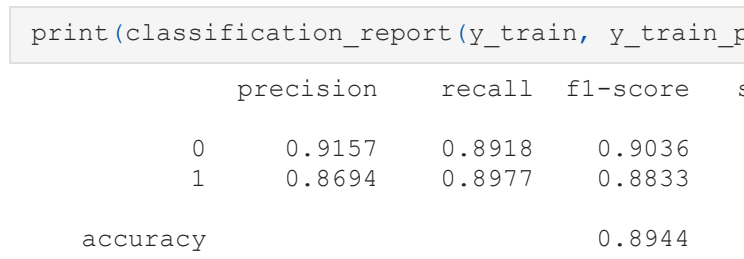
```
In [37]: sns.distplot(x=credit['A14'], kde=False)
```

```
Out[37]: <AxesSubplot: >
```



```
In [27]: sns.scatterplot(data=credit, x='A15', y='A14')
```

```
Out[27]: <AxesSubplot: xlabel='A15', ylabel='A14'>
```



The correlation matrix below shows that the numeric variables are only weakly correlated, so it is unnecessary to remove any of them:

```
In [15]: X.corr()
```

```
Out[15]:
          A2      A3      A7      A10      A13      A14
A2    1.000000  0.201315  0.392788  0.185574  -0.077159  0.018539
A3    0.201315  1.000000  0.298902  0.271207  -0.222346  0.123121
A7    0.392788  0.298902  1.000000  0.322330  -0.076389  0.051345
A10   0.185574  0.271207  0.322330  1.000000  -0.119808  0.063692
A13   -0.077159  -0.222346  -0.076389  -0.119808  1.000000  0.065609
A14   0.018539  0.123121  0.051345  0.063692  0.065609  1.000000
```

The next step is to convert categorical variables into dummies - variables with 0/1 value for each category (except the first value in order to avoid collinearity):

```
In [16]: X = pd.get_dummies(X, drop_first=True)
```

The numeric variables remain as they were before, and the categorical variables are now converted. The total number of features increased from 14 to 34:

```
In [17]: X
```

```
Out[17]:
          A2      A3      A7      A10      A13      A14      A1_A2_0      A1_A2_1      A1_A2_2      A1_A2_3      A1_A2_4      A1_A2_5      A1_A2_6      A1_A2_7      A1_A2_8      A1_A2_9      A1_A2_10      A1_A2_11      A1_A2_12      A1_A2_13      A1_A2_14      A1_A2_15      A1_A2_16      A1_A2_17      A1_A2_18      A1_A2_19      A1_A2_20      A1_A2_21      A1_A2_22      A1_A2_23      A1_A2_24      A1_A2_25      A1_A2_26      A1_A2_27      A1_A2_28      A1_A2_29      A1_A2_30      A1_A2_31      A1_A2_32      A1_A2_33      A1_A2_34      A1_A2_35      A1_A2_36      A1_A2_37      A1_A2_38      A1_A2_39      A1_A2_40      A1_A2_41      A1_A2_42      A1_A2_43      A1_A2_44      A1_A2_45      A1_A2_46      A1_A2_47      A1_A2_48      A1_A2_49      A1_A2_50      A1_A2_51      A1_A2_52      A1_A2_53      A1_A2_54      A1_A2_55      A1_A2_56      A1_A2_57      A1_A2_58      A1_A2_59      A1_A2_60      A1_A2_61      A1_A2_62      A1_A2_63      A1_A2_64      A1_A2_65      A1_A2_66      A1_A2_67      A1_A2_68      A1_A2_69      A1_A2_70      A1_A2_71      A1_A2_72      A1_A2_73      A1_A2_74      A1_A2_75      A1_A2_76      A1_A2_77      A1_A2_78      A1_A2_79      A1_A2_80      A1_A2_81      A1_A2_82      A1_A2_83      A1_A2_84      A1_A2_85      A1_A2_86      A1_A2_87      A1_A2_88      A1_A2_89      A1_A2_90      A1_A2_91      A1_A2_92      A1_A2_93      A1_A2_94      A1_A2_95      A1_A2_96      A1_A2_97      A1_A2_98      A1_A2_99      A1_A2_100      A1_A2_101      A1_A2_102      A1_A2_103      A1_A2_104      A1_A2_105      A1_A2_106      A1_A2_107      A1_A2_108      A1_A2_109      A1_A2_110      A1_A2_111      A1_A2_112      A1_A2_113      A1_A2_114      A1_A2_115      A1_A2_116      A1_A2_117      A1_A2_118      A1_A2_119      A1_A2_120      A1_A2_121      A1_A2_122      A1_A2_123      A1_A2_124      A1_A2_125      A1_A2_126      A1_A2_127      A1_A2_128      A1_A2_129      A1_A2_130      A1_A2_131      A1_A2_132      A1_A2_133      A1_A2_134      A1_A2_135      A1_A2_136      A1_A2_137      A1_A2_138      A1_A2_139      A1_A2_140      A1_A2_141      A1_A2_142      A1_A2_143      A1_A2_144      A1_A2_145      A1_A2_146      A1_A2_147      A1_A2_148      A1_A2_149      A1_A2_150      A1_A2_151      A1_A2_152      A1_A2_153      A1_A2_154      A1_A2_155      A1_A2_156      A1_A2_157      A1_A2_158      A1_A2_159      A1_A2_160      A1_A2_161      A1_A2_162      A1_A2_163      A1_A2_164      A1_A2_165      A1_A2_166      A1_A2_167      A1_A2_168      A1_A2_169      A1_A2_170      A1_A2_171      A1_A2_172      A1_A2_173      A1_A2_174      A1_A2_175      A1_A2_176      A1_A2_177      A1_A2_178      A1_A2_179      A1_A2_180      A1_A2_181      A1_A2_182      A1_A2_183      A1_A2_184      A1_A2_185      A1_A2_186      A1_A2_187      A1_A2_188      A1_A2_189      A1_A2_190      A1_A2_191      A1_A2_192      A1_A2_193      A1_A2_194      A1_A2_195      A1_A2_196      A1_A2_197      A1_A2_198      A1_A2_199      A1_A2_200      A1_A2_201      A1_A2_202      A1_A2_203      A1_A2_204      A1_A2_205      A1_A2_206      A1_A2_207      A1_A2_208      A1_A2_209      A1_A2_210      A1_A2_211      A1_A2_212      A1_A2_213      A1_A2_214      A1_A2_215      A1_A2_216      A1_A2_217      A1_A2_218      A1_A2_219      A1_A2_220      A1_A2_221      A1_A2_222      A1_A2_223      A1_A2_224      A1_A2_225      A1_A2_226      A1_A2_227      A1_A2_228      A1_A2_229      A1_A2_230      A1_A2_231      A1_A2_232      A1_A2_233      A1_A2_234      A1_A2_235      A1_A2_236      A1_A2_237      A1_A2_238      A1_A2_239      A1_A2_240      A1_A2_241      A1_A2_242      A1_A2_243      A1_A2_244      A1_A2_245      A1_A2_246      A1_A2_247      A1_A2_248      A1_A2_249      A1_A2_250      A1_A2_251      A1_A2_252      A1_A2_253      A1_A2_254      A1_A2_255      A1_A2_256      A1_A2_257      A1_A2_258      A1_A2_259      A1_A2_260      A1_A2_261      A1_A2_262      A1_A2_263      A1_A2_264      A1_A2_265      A1_A2_266      A1_A2_267      A1_A2_268      A1_A2_269      A1_A2_270      A1_A2_271      A1_A2_272      A1_A2_273      A1_A2_274      A1_A2_275      A1_A2_276      A1_A2_277      A1_A2_278      A1_A2_279      A1_A2_280      A1_A2_281      A1_A2_282      A1_A2_283      A1_A2_284      A1_A2_285      A1_A2_286      A1_A2_287      A1_A2_288      A1_A2_289      A1_A2_290      A1_A2_291      A1_A2_292      A1_A2_293      A1_A2_294      A1_A2_295      A1_A2_296      A1_A2_297      A1_A2_298      A1_A2_299      A1_A2_300      A1_A2_301      A1_A2_302      A1_A2_303      A1_A2_304      A1_A2_305      A1_A2_306      A1_A2_307      A1_A2_308      A1_A2_309      A1_A2_310      A1_A2_311      A1_A2_312      A1_A2_313      A1_A2_314      A1_A2_315      A1_A2_316      A1_A2_317      A1_A2_318      A1_A2_319      A1_A2_320      A1_A2_321      A1_A2_322      A1_A2_323      A1_A2_324      A1_A2_325      A1_A2_326      A1_A2_327      A1_A2_328      A1_A2_329      A1_A2_330      A1_A2_331      A1_A2_332      A1_A2_333      A1_A2_334      A1_A2_335      A1_A2_336      A1_A2_337      A1_A2_338      A1_A2_339      A1_A2_340      A1_A2_341      A1_A2_342      A1_A2_343      A1_A2_344      A1_A2_345      A1_A2_346      A1_A2_347      A1_A2_348      A1_A2_349      A1_A2_350      A1_A2_351      A1_A2_352      A1_A2_353      A1_A2_354      A1_A2_355      A1_A2_356      A1_A2_357      A1_A2_358      A1_A2_359      A1_A2_360      A1_A2_361      A1_A2_362      A1_A2_363      A1_A2_364      A1_A2_365      A1_A2_366      A1_A2_367      A1_A2_368      A1_A2_369      A1_A2_370      A1_A2_371      A1_A2_372      A1_A2_373      A1_A2_374      A1_A2_375      A1_A2_376      A1_A2_377      A1_A2_378      A1_A2_379      A1_A2_380      A1_A2_381      A1_A2_382      A1_A2_383      A1_A2_384      A1_A2_385      A1_A2_386      A1_A2_387      A1_A2_388      A1_A2_389      A1_A2_390      A1_A2_391      A1_A2_392      A1_A2_393      A1_A2_394      A1_A2_395      A1_A2_396      A1_A2_397      A1_A2_398      A1_A2_399      A1_A2_400      A1_A2_401      A1_A2_402      A1_A2_403      A1_A2_404      A1_A2_405      A1_A2_406      A1_A2_407      A1_A2_408      A1_A2_409      A1_A2_410      A1_A2_411      A1_A2_412      A1_A2_413      A1_A2_414      A1_A2_415      A1_A2_416      A1_A2_417      A1_A2_418      A1_A2_419      A1_A2_420      A1_A2_421      A1_A2_422      A1_A2_4
```


Predictions for both subsets will be made using the first decision tree classifier.

```
In [42]: y_train_pred = dt.predict(X_train_s)
y_test_pred = dt.predict(X_test_s)
```

All metrics equal 1 for the training subset, and are significantly lower for the testing subset, which indicates overfitting.

```
In [43]: print(classification_report(y_train, y_train_pred, digits=4))
```

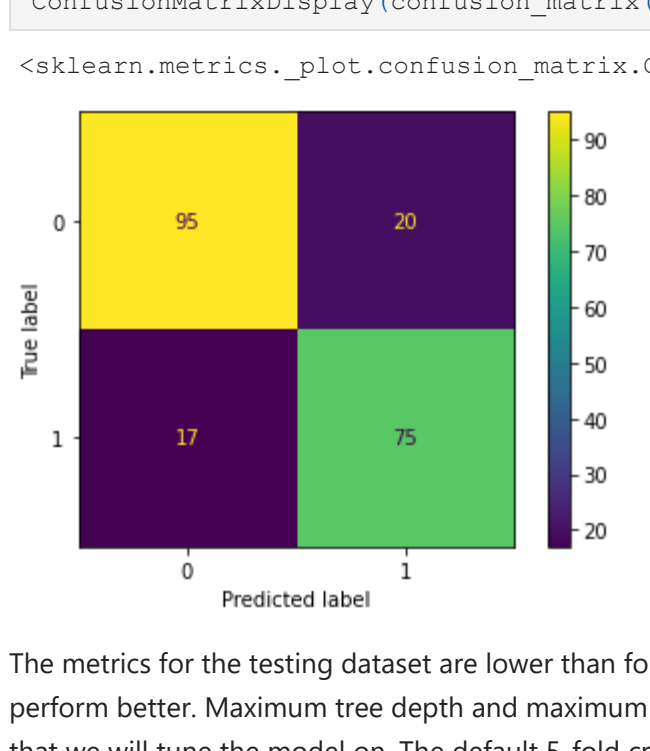
	precision	recall	f1-score	support
0	1.0000	1.0000	1.0000	268
1	1.0000	1.0000	1.0000	215
accuracy			1.0000	483
macro avg	1.0000	1.0000	1.0000	483
weighted avg	1.0000	1.0000	1.0000	483

```
In [44]: print(classification_report(y_test, y_test_pred, digits=4))
```

	precision	recall	f1-score	support
0	0.8482	0.8261	0.8370	115
1	0.7895	0.8352	0.8021	92
accuracy			0.8213	207
macro avg	0.8188	0.8207	0.8196	207
weighted avg	0.8221	0.8213	0.8215	207

```
In [45]: ConfusionMatrixDisplay(confusion_matrix(y_test, y_test_pred)).plot()
```

```
Out[45]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0xc9815a74c0>
```



The metrics for the testing data are lower than for the logistic regression model, so we will apply grid search to find a tree that will perform better. Maximum tree depth and maximum number of features considered when splitting the node will be the hyperparameters that we will tune the model on. The default 5-fold cross validation will be used for the selection:

```
In [46]: param_grid = {'max_depth':range(1, dt.tree_max_depth+1, 2),
                    'max_features': range(1, len(dt.feature_importances_)+1)}
best_tree = GridSearchCV(DecisionTreeClassifier(random_state=101),
                        param_grid=param_grid,
                        scoring='accuracy',
                        n_jobs=-1)
best_tree = best_tree.fit(X_train_s, y_train)
```

The number of nodes and tree depth for the best tree are given below. As we can see, the tree is much simpler now:

```
In [47]: best_tree.best_estimator_.tree_.node_count, best_tree.best_estimator_.tree_.max_depth
```

```
Out[47]: (3, 1)
```

Here are the most important features of the pruned tree. The same as for the logistic regression model, A8_1 is the most important (and the only) feature:

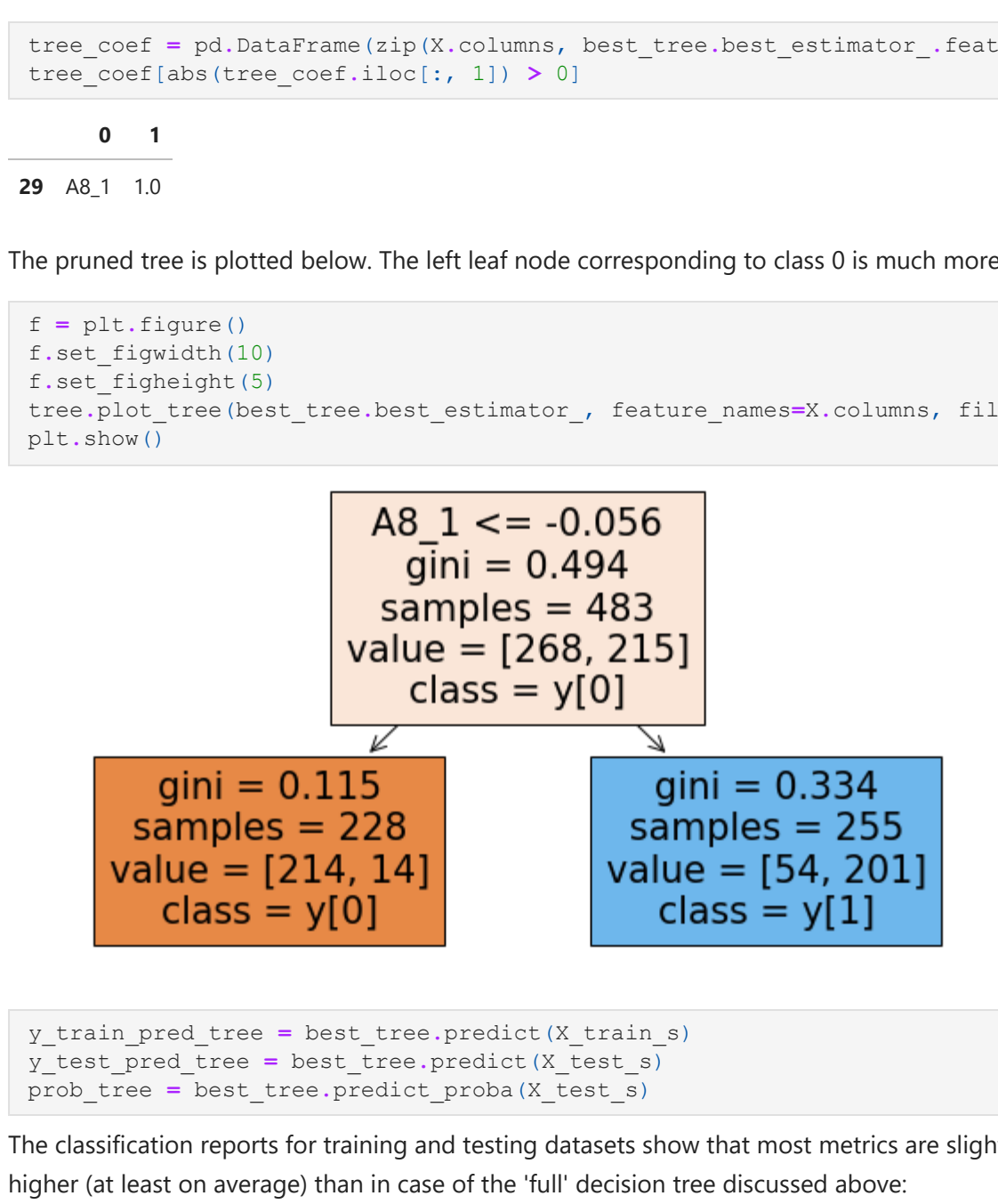
```
In [48]: tree_coef = pd.DataFrame(zip(X.columns, best_tree.best_estimator_.feature_importances_)).sort_values(by=1)
tree_coef[abs(tree_coef.iloc[:, 1]) > 0]
```

```
Out[48]:
```

```
0 1
29 A8_1 1.0
```

The pruned tree is plotted below. The left leaf node corresponding to class 0 is much more pure:

```
In [49]: f = plt.figure()
f.set_figwidth(10)
f.set_figheight(5)
tree.plot_tree(best_tree.best_estimator_, feature_names=X.columns, filled=True, class_names=['0', '1'])
plt.show()
```



```
In [50]: y_train_pred_tree = best_tree.predict(X_train_s)
y_test_pred_tree = best_tree.predict(X_test_s)
prob_tree = best_tree.predict_proba(X_test_s)
```

The classification reports for training and testing datasets show that most metrics are slightly worse for the testing dataset, but they are higher (at least on average) than in case of the 'full' decision tree discussed above:

```
In [51]: print(classification_report(y_train, y_train_pred_tree, digits=4))
```

	precision	recall	f1-score	support
0	0.9386	0.7985	0.8629	268
1	0.7882	0.9349	0.8553	215
accuracy			0.8592	483
macro avg	0.8634	0.8667	0.8591	483
weighted avg	0.8717	0.8592	0.8595	483

```
In [52]: print(classification_report(y_test, y_test_pred_tree, digits=4))
```

	precision	recall	f1-score	support
0	0.9109	0.8000	0.8519	115
1	0.7830	0.9022	0.8384	92
accuracy			0.8454	207
macro avg	0.8470	0.8511	0.8451	207
weighted avg	0.8541	0.8454	0.8459	207

```
In [53]: auc_tree = accuracy_score(y_test, y_test_pred_tree)
prec_tree = precision_score(y_test, y_test_pred_tree)
rec_tree = recall_score(y_test, y_test_pred_tree)
f1_tree = f1_score(y_test, y_test_pred_tree)
```

The decision tree model has a higher precision for class 0 (and lower for class 1), compared to logistic regression. Aggregated metrics are slightly better for logistic regression. Area Under Curve is significantly lower for the decision tree model than for logistic regression:

```
In [53]: auc_tree = roc_auc_score(y_test, prob_tree[:, 1])
```

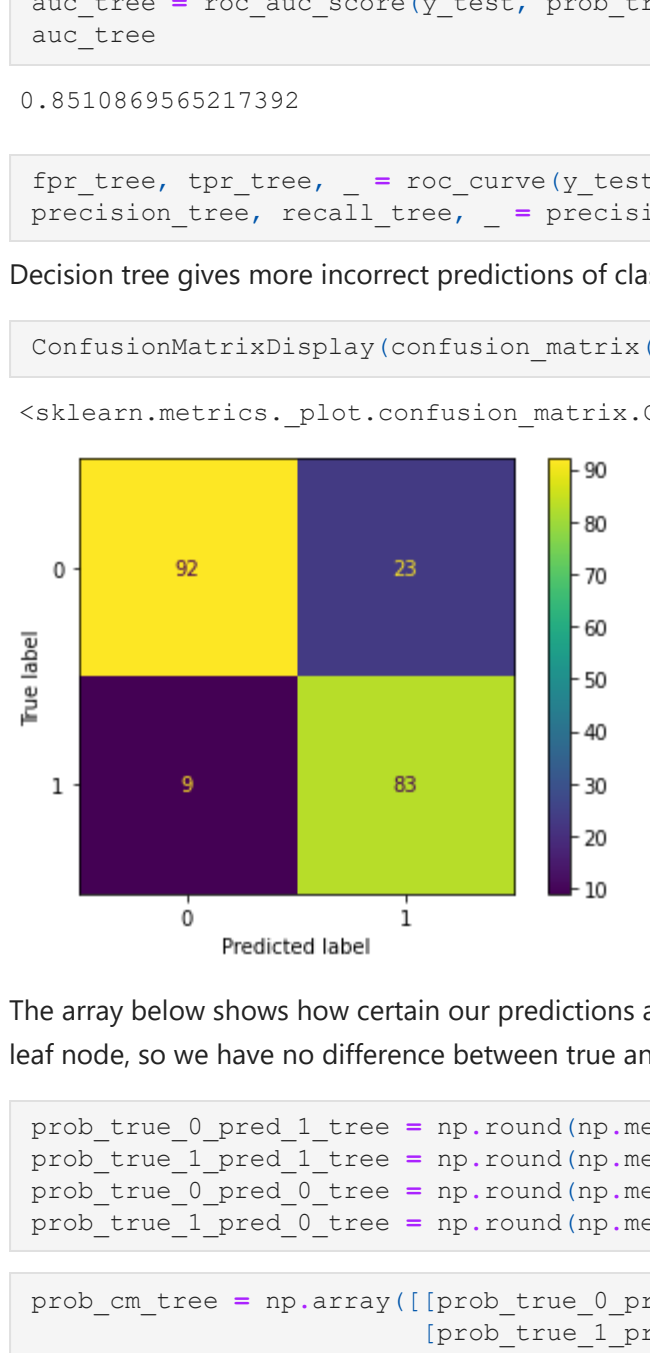
```
Out[53]: 0.8510869565217392
```

```
In [76]: fpr_tree, tpr_tree, _ = roc_curve(y_test, prob_tree[:, 1])
precision_tree, recall_tree, _ = precision_recall_curve(y_test, prob_tree[:, 1])
```

Decision tree gives more incorrect predictions of class 1 and fewer incorrect predictions of class 0 than logistic regression:

```
In [55]: ConfusionMatrixDisplay(confusion_matrix(y_test, y_test_pred_tree)).plot()
```

```
Out[55]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0xc981211310>
```

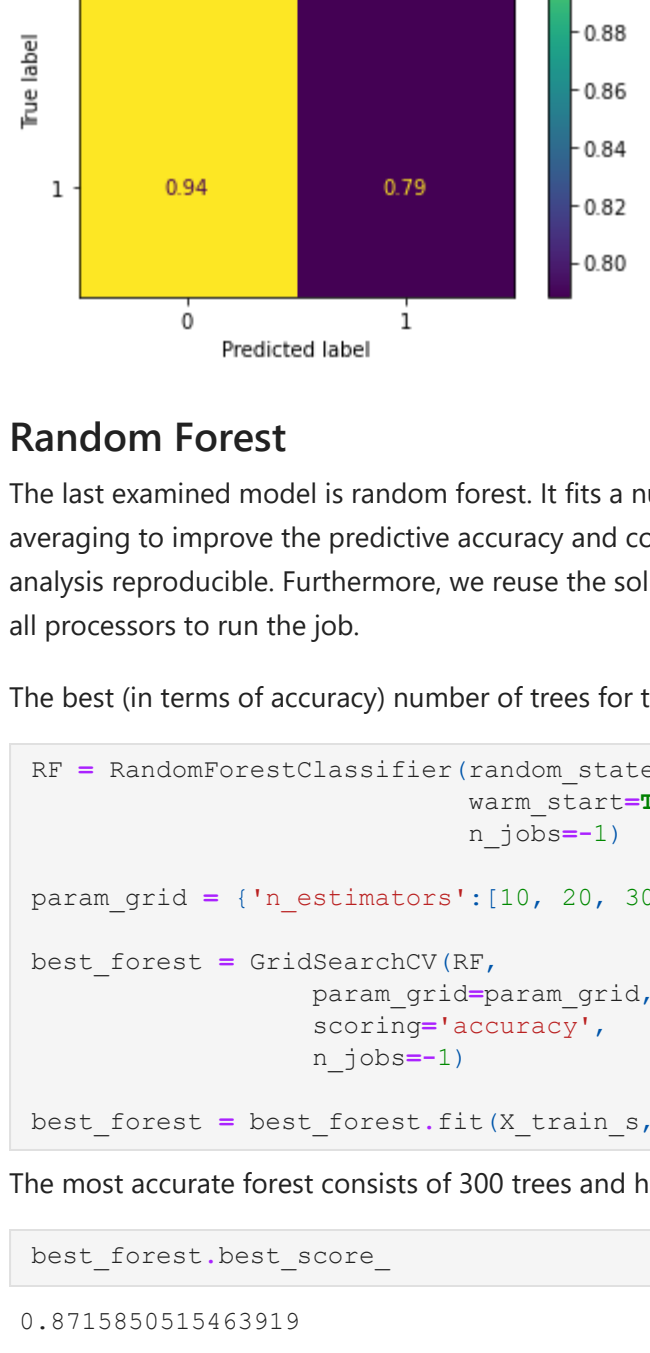


The array below shows how certain our predictions are for each confusion matrix entry. We assign the same probability to each entry in the leaf node, so we have no difference between true and false predictions:

```
In [56]: prob_true_0_pred_1_tree = np.round(np.mean(np.max(prob_tree[(y_test==0) & (y_test_pred_tree==1), :, axis=1]),
prob_true_1_pred_1_tree = np.round(np.mean(np.max(prob_tree[(y_test==1) & (y_test_pred_tree==1), :, axis=1]),
prob_true_0_pred_0_tree = np.round(np.mean(np.max(prob_tree[(y_test==0) & (y_test_pred_tree==0), :, axis=1]),
prob_true_1_pred_0_tree = np.round(np.mean(np.max(prob_tree[(y_test==1) & (y_test_pred_tree==0), :, axis=1]),
```

```
In [57]: prob_cm_tree = np.array([[prob_true_0_pred_0_tree, prob_true_0_pred_1_tree],
                             [prob_true_1_pred_0_tree, prob_true_1_pred_1_tree]])
ConfusionMatrixDisplay(prob_cm_tree).plot()
```

```
Out[57]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0xc981211310>
```



Random Forest

The last examined model is random forest. It fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. As for previous models, we select a random state to keep the analysis reproducible. Furthermore, we reuse the solution of the previous call to fit and add more estimators to the ensemble. Also, we use all processors to run the job.

The best (in terms of accuracy) number of trees for the created random forest classifier will be searched using GridSearchCV:

```
In [58]: RF = RandomForestClassifier(random_state=101,
                                warm_start=True,
                                n_jobs=-1)
```

```
param_grid = {'n_estimators':[10, 20, 30, 40, 50, 100, 200, 300, 400, 500]}
best_forest = GridSearchCV(RF,
                          param_grid=param_grid,
                          scoring='accuracy',
                          n_jobs=-1)
best_forest = best_forest.fit(X_train_s, y_train)
```

The most accurate forest consists of 300 trees and has an accuracy of about 0.87:

```
In [59]: best_forest.best_score_
```

```
Out[59]: 0.8715850515463919
```

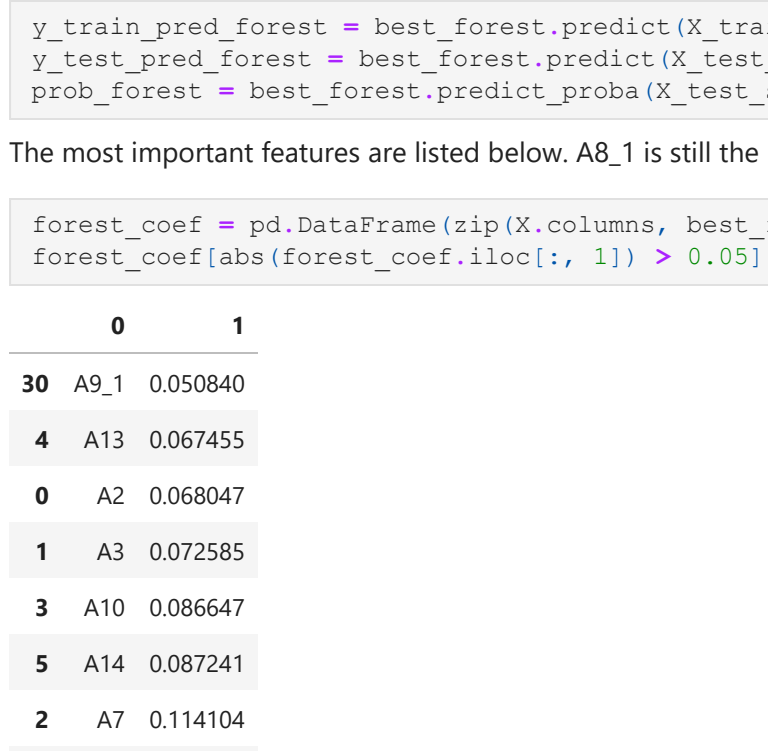
```
In [60]: best_forest.best_params_
```

```
Out[60]: {'n_estimators': 300}
```

The accuracies of other forests are not much different, above 0.85:

```
In [61]: g = sns.lineplot(x=range(len(param_grid['n_estimators'])), y=best_forest.cv_results_['mean_test_score'])
g.set_xticks(range(len(param_grid['n_estimators'])))
g.set_xticklabels(param_grid['n_estimators'])
g.set_title("Random Forest Accuracies")
g.set_xlabel("Tree Number", ylabel="Accuracy")
```

```
Out[61]: [Text(0.5, 0, 'Tree Number'), Text(0, 0.5, 'Accuracy')]
```



```
In [62]: y_train_pred_forest = best_forest.predict(X_train_s)
y_test_pred_forest = best_forest.predict(X_test_s)
prob_forest = best_forest.predict_proba(X_test_s)
```

The most important features are listed below. A8_1 is still the most important, followed by the numeric features:

```
In [63]: forest_coef = pd.DataFrame(zip(X.columns, best_forest.best_estimator_.feature_importances_)).sort_values(by=1)
forest_coef[abs(forest_coef.iloc[:, 1]) > 0.05]
```

```
Out[63]:
```

```
0 1
30 A9_1 0.050840
4 A13 0.067455
0 A2 0.068047
1 A3 0.072585
5 A10 0.086647
5 A14 0.087241
2 A7 0.116103
29 A8_1 0.291123
```

Although all metric values are 1 for the training subset (which might indicate overfitting), the testing metrics are high as well:

```
In [64]: print(classification_report(y_train, y_train_pred_forest, digits=4))
```

	precision	recall	f1-score	support
0	1.0000	1.0000	1.0000	268
1	1.0000	1.0000	1.0000	215
accuracy			1.0000	483
macro avg	1.0000	1.0000	1.0000	483
weighted avg	1.0000	1.0000	1.0000	483

```
In [65]: print(classification_report(y_test, y_test_pred_forest, digits=4))
```

	precision	recall	f1-score	support
0	0.8696	0.8696	0.8696	115
1	0.8370	0.8370	0.8370	92
accuracy			0.8551	207
macro avg	0.8533	0.8533	0.8533	207
weighted avg	0.8531	0.8531	0.8551	207

```
In [81]: auc_forest = accuracy_score(y_test, y_test_pred_forest)
prec_forest = precision_score(y_test, y_test_pred_forest)
rec_forest = recall_score(y_test, y_test_pred_forest)
f1_forest = f1_score(y_test, y_test_pred_forest)
```

In this case, classification report metrics are only slightly lower than for logistic regression, and predictions for class 0 are again a bit more accurate than for class 1. However, AUC score is a bit higher than for logistic regression:

```
In [66]: auc_forest = roc_auc_score(y_test, prob_forest[:, 1])
```

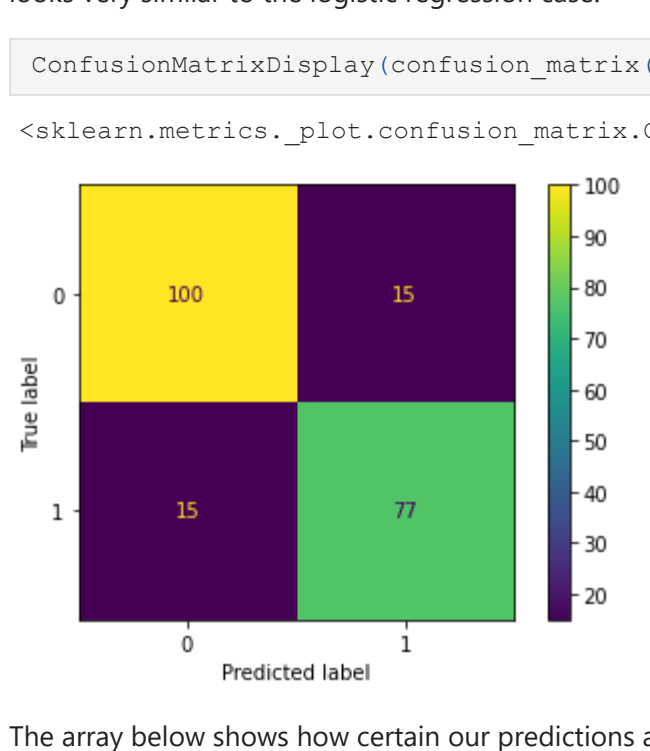
```
Out[66]: 0.9247637051039697
```

```
In [67]: fpr_forest, tpr_forest, _ = roc_curve(y_test, prob_forest[:, 1])
precision_forest, recall_forest, _ = precision_recall_curve(y_test, prob_forest[:, 1])
```

Identical precision and recall values for both classes are explained by the same number of false positives and false negatives. But the matrix looks very similar to the logistic regression case:

```
In [68]: ConfusionMatrixDisplay(confusion_matrix(y_test, y_test_pred_forest)).plot()
```

```
Out[68]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0xc984ac74c0>
```

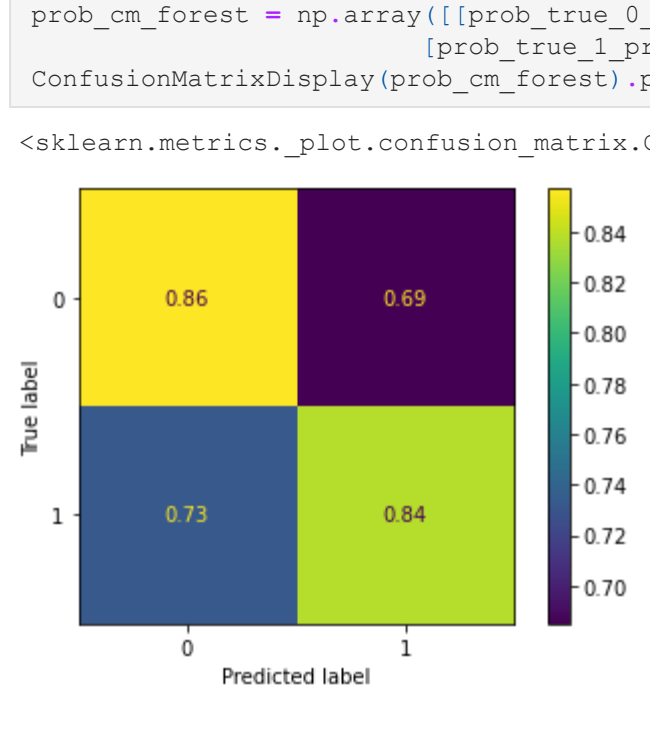


The array below shows how certain our predictions are for each confusion matrix entry. Similarly to the logistic regression model (and contrary to the decision tree), incorrect predictions are less certain on average:

```
In [69]: prob_true_0_pred_1_forest = np.round(np.mean(np.max(prob_forest[(y_test==0) & (y_test_pred_forest==1), :, axis=1]),
prob_true_1_pred_1_forest = np.round(np.mean(np.max(prob_forest[(y_test==1) & (y_test_pred_forest==1), :, axis=1]),
prob_true_0_pred_0_forest = np.round(np.mean(np.max(prob_forest[(y_test==0) & (y_test_pred_forest==0), :, axis=1]),
prob_true_1_pred_0_forest = np.round(np.mean(np.max(prob_forest[(y_test==1) & (y_test_pred_forest==0), :, axis=1]),
```

```
In [70]: prob_cm_forest = np.array([[prob_true_0_pred_0_forest, prob_true_0_pred_1_forest],
                                [prob_true_1_pred_0_forest, prob_true_1_pred_1_forest]])
ConfusionMatrixDisplay(prob_cm_forest).plot()
```

```
Out[70]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0xc984ac74c0>
```



Final Model Recommendation

After fitting three classification models, we will compare them and select the best one.

Classification report metrics (precision, recall, accuracy, F1-score) are very similar for logistic regression and random forest models (about 0.85). Confusion matrices also have almost the same counts.

The results for decision tree are a bit different. Its great advantage is interpretability - it is easy to visualize and understand, especially in our case when only one feature is selected. But contrary to other models, precision and recall values differ for the positive category by more than 0.1 (0.78 vs 0.9). Although average metric values are similar to other models, its AUC value is significantly lower (0.85 vs 0.92). These differences will be illustrated graphically below.

Also, we compared prediction probabilities for all three models. Logistic regression and random forest assign lower prediction probabilities for incorrect predictions (about 0.85 for correct vs 0.7 for incorrect). For the decision tree model however, the assigned probabilities depend only on the predicted label, regardless of whether the prediction is correct. So that is another drawback of the decision tree model.

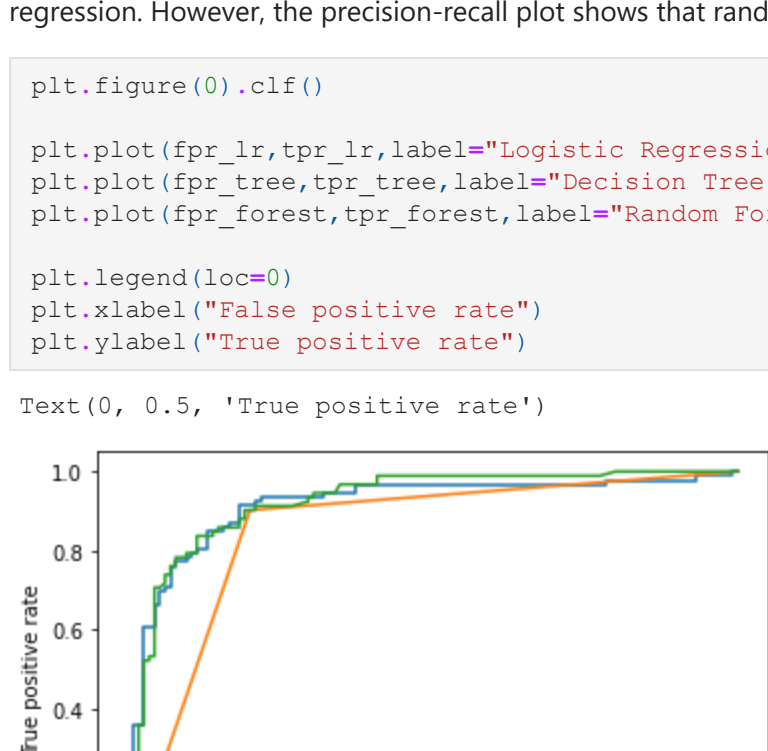
Below, ROC and precision-recall curves are plotted for all three models. Again, decision tree is doing the worst - its area under the curve is the lowest, since it has the only 2 thresholds (the same number as leaf nodes) for which to calculate FPR and TPR (the same for the precision-recall plot). Logistic regression and random forest are similar. AUC for the random forest model is slightly higher than for logistic regression. However, the precision-recall plot shows that random forest has worse precision rates for low recall values.

```
In [75]: plt.figure(0).clf()
```

```
plt.plot(fpr_lr, tpr_lr, label='Logistic Regression', auc="str(np.round(auc_lr, 3))")
plt.plot(fpr_tree, tpr_tree, label='Decision Tree', auc="str(np.round(auc_tree, 3))")
plt.plot(fpr_forest, tpr_forest, label='Random Forest', auc="str(np.round(auc_forest, 3))")

plt.legend(loc=0)
plt.xlabel("False positive rate")
plt.ylabel("True positive rate")
```

```
Out[75]: Text(0.5, 0, 'True positive rate')
```



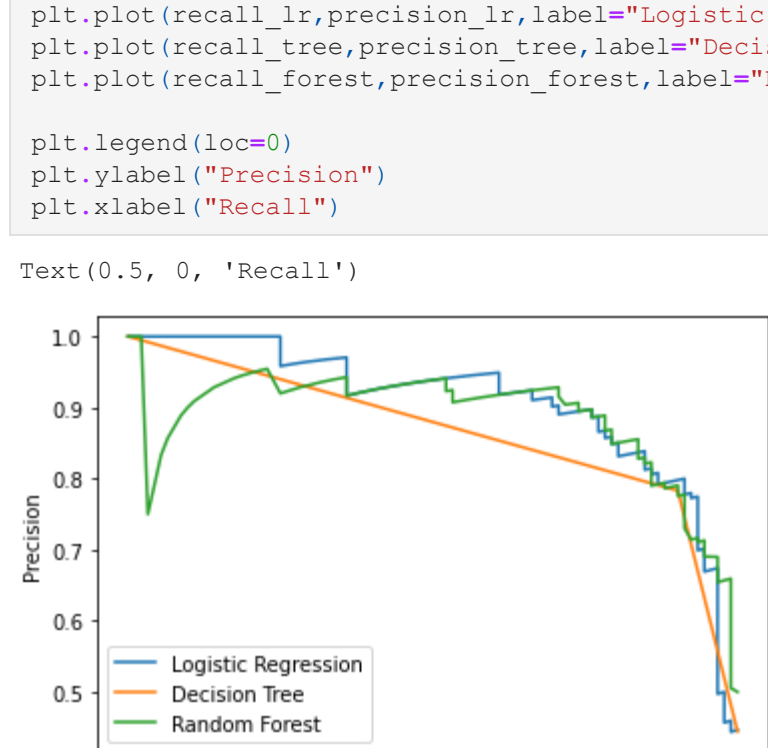
```
Out[75]:
```

```
In [73]: plt.figure(0).clf()
```

```
plt.plot(recall_lr, precision_lr, label='Logistic Regression')
plt.plot(recall_tree, precision_tree, label='Decision Tree')
plt.plot(recall_forest, precision_forest, label='Random Forest')

plt.legend(loc=0)
plt.xlabel("Precision")
plt.ylabel("Recall")
```

```
Out[73]: Text(0.5, 0, 'Recall')
```



In the previous part, we examined whether the logistic regression model adheres the necessary assumptions. We had doubts about it, so selecting the random forest model might be sensible, since it does not rely on any formal distributional assumptions.

So the final recommendation is to use the random forest model.

Key Findings And Insights

As we discussed above, random forest model is considered to be the best for this analysis. However, all three models give fairly good results in terms of prediction accuracy. Below is the summary of classification metrics for all the models (precision, recall and F1-score are given for class 1):

```
In [92]: results_df=pd.DataFrame(data={'Logistic Regression': [auc_lr, prec_lr, rec_lr, f1_lr, auc_lr],
                                'Decision Tree': [auc_tree, prec_tree, rec_tree, f1_tree, auc_tree],
                                'Random Forest': [auc_forest, prec_forest, rec_forest, f1_forest, auc_forest]},
                                results_df.set_index('Metrics'))
```

```
Out[92]:
```

	Logistic Regression	Decision Tree	Random Forest
Metrics			
Accuracy	0.859903	0.845411	0.855072
Precision	0.838710	0.783019	0.836957
Recall	0.847826	0.902174	0.836957
F1 Score	0.843243	0.838384	0.836957
AUC	0.917013	0.851087	0.924764

Except precision for the decision tree model, all metrics exceed 0.8, and accuracy for logistic regression and random forest is above 0.85. Inconsistency of metric values for the decision tree model was discussed above and is most likely caused by excessive tree pruning.

The metrics tend to have lower values for class 1 cases across all our models. This could be explained by large values of outliers for class 1 that were observed in the data exploration part.

Apart from predictability, let's analyze which variables are found to be the most important in our models.

The most important variable in all three models is A8_1. For logistic regression and random forest, it is several times more important than other features. Also, it is the only variable in the decision tree model, and is enough for a quite good prediction accuracy.

The other similarity between the logistic regression and random forest models is a relatively high importance of numeric variables - A13, A14, A7 and A10 are common for both models. Other numeric variables are more important for the random forest model, while the logistic regression gives more importance to some dummy variables (A4_2, A5_9, A6_8, A5_14). On the other hand, some of these variables have high VIF values (see the analysis above), and thus might be useful to remove from the dataset.

So the analysis of variable importance justifies the selection of random forest for this study.

Suggestions

The initial analysis shows that the numeric features had some outliers that we did not remove, because we were unaware of their nature. Examining them is much harder due to the missing feature names. A deeper analysis of potential outliers could be useful to decide whether to preserve them in the dataset. Model predictability might improve upon their removal, especially of logistic regression, since it relies on linearity assumption.

Although entry numbers for both classes are similar (55% and 45%), we could apply oversampling and/or undersampling techniques to make the dataset perfectly balanced. These approaches would have their own disadvantages, such as information loss or a risk of overfitting.

Feature multicollinearity analysis after fitting logistic regression showed that some (mostly categorical) variables have high variance inflation factors. Since information provided by them is redundant, re-fitting without them might be another way to improve our models.

Also, some other classification models could be experimented with. We used random forest as an averaging method, based on many independent trees. Boosting methods (such as AdaBoost) rely on building estimators sequentially, increasing weights on samples that were previously predicted incorrectly.

```
In [ ]:
```