



Home Credit Default Risk

Can you predict how capable each applicant is of repaying a loan?

Home Loan Repayment

Group 4: Yuexuan Fu, Anika Kingsley, Henry Raymond, Xinye Xu, Yuxin Zhang

Presentation Contents

- Project Description
 - Data Description
 - Business insight
 - Project Scope
- Analysis
 - 4 Classification Models
- Summary of Findings
 - Conclusion
 - Recommendations



Project Description



Data

Binary target variable:

- 1 - client with payment difficulties: he/she had late payment more than X days on at least one of the first Y installments of the loan in our sample,
- 0 - all other cases

Features: client specific features (credit, employment, family size, etc.)

Data Size:

```
train.shape
```

```
(307511, 122)
```

Business Insight

Default Predictions

- False negative: costly
- False positive: not as costly

Loaning to those with poor credit history

- Analyze alternative data
- Minimize false negatives with recall score



Data Cleaning

- Dropped columns with N/A values
- Filled N/A with means/0
- Dropped features
 - Irrelevant to our business question
 - Mismatching data
- Turned categorical variables into binary
- Merged data (Add Previous Loans Count

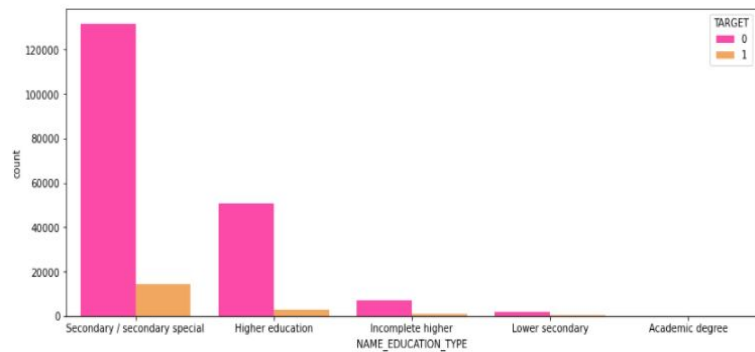
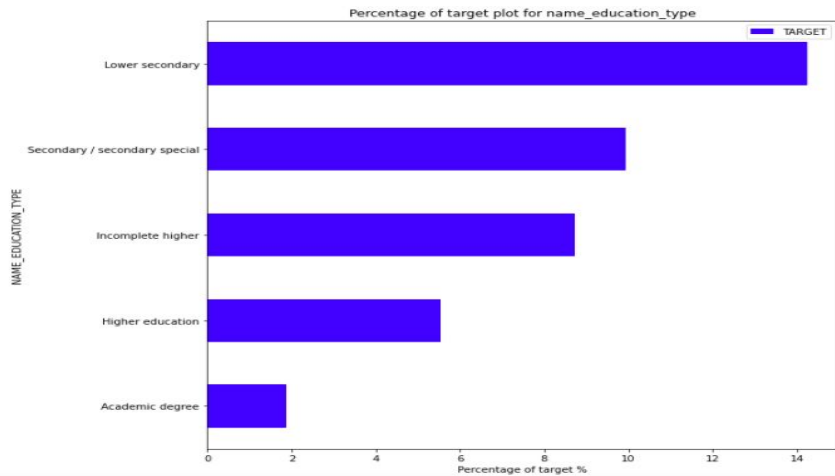
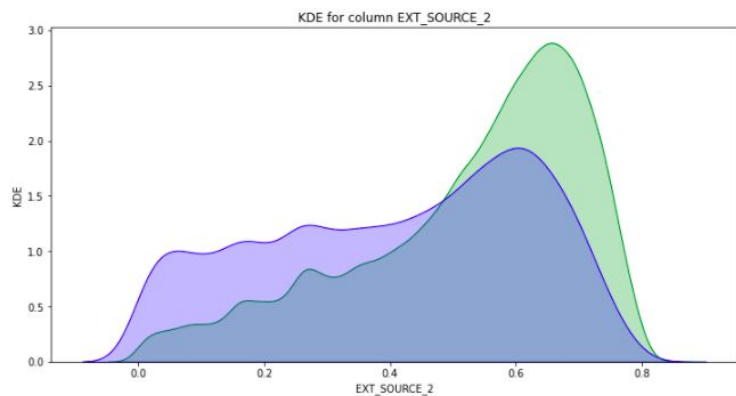
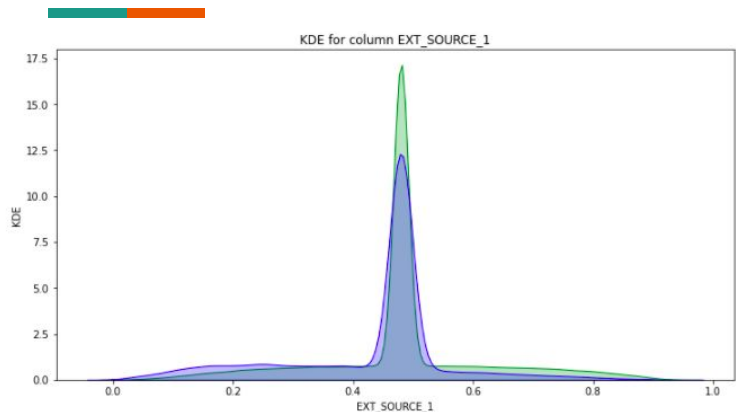
from Bureau.csv on SK-ID-CURR)

After:

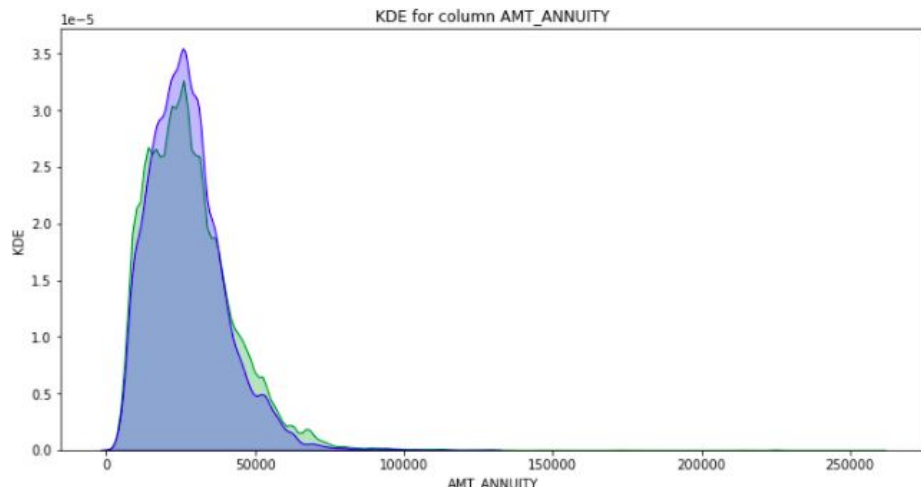
| | dtypes | column count |
|---|---------|--------------|
| 0 | object | 13 |
| 1 | float64 | 10 |
| 2 | int64 | 9 |

| | Number of Missing Values | Percentage of Entries Missing |
|--------------------------|--------------------------|-------------------------------|
| COMMONAREA_MEDI | 214865 | 69.87 |
| COMMONAREA_AVG | 214865 | 69.87 |
| COMMONAREA_MODE | 214865 | 69.87 |
| NONLIVINGAPARTMENTS_MODE | 213514 | 69.43 |
| NONLIVINGAPARTMENTS_AVG | 213514 | 69.43 |
| NONLIVINGAPARTMENTS_MEDI | 213514 | 69.43 |
| FONDKAPREMONT_MODE | 210295 | 68.39 |
| LIVINGAPARTMENTS_MODE | 210199 | 68.35 |
| LIVINGAPARTMENTS_AVG | 210199 | 68.35 |
| LIVINGAPARTMENTS_MEDI | 210199 | 68.35 |
| FLOORSMIN_AVG | 208642 | 67.85 |
| FLOORSMIN_MODE | 208642 | 67.85 |
| FLOORSMIN_MEDI | 208642 | 67.85 |

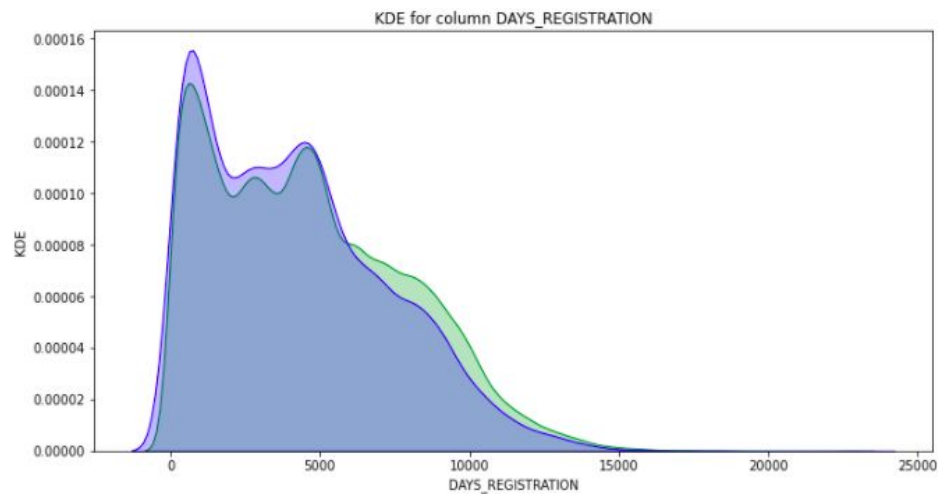
Exploratory Data Analysis - Important Features



Exploratory Data Analysis - Less Important Features

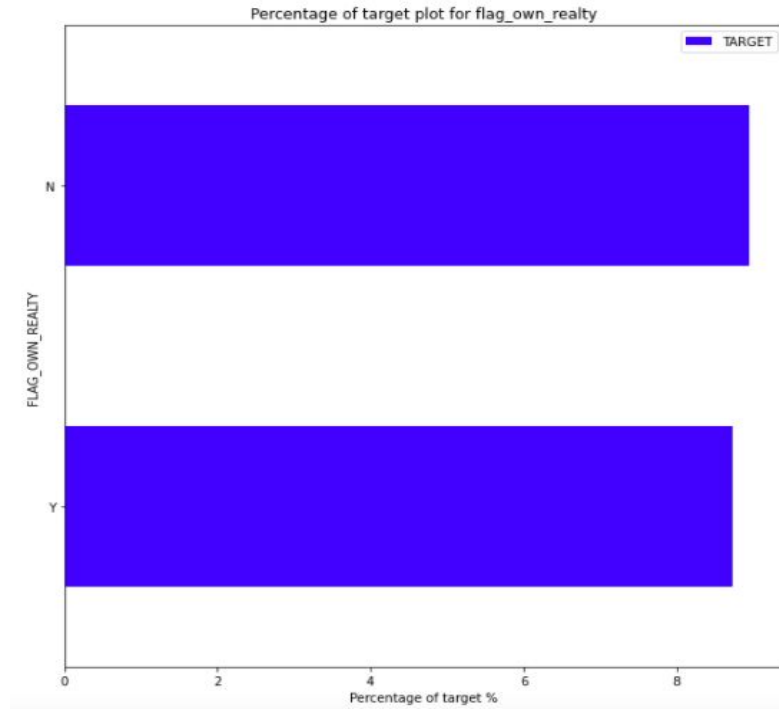


Loan annuity



How many days before the application did client change his registration

Exploratory Data Analysis - Features that seems with no impact



Flag if client owns a house or flat



Business Question:

What attributes of a client are associated with loan repayment?




Project Scope

Determine significant predictors of whether a client will repay a loan by using the following models:

- Logistic regression
- Linear SVC
- Decision Tree
- Random Forest

Analysis



```
1 safe= train[train.TARGET == 0]
2 bad= train[train.TARGET == 1]
3 percentage = len(bad)/float(len(safe))
4 percentage
5 risky_loans = bad
6 safe_loans = safe.sample(frac = percentage, random_state = 1)
7 print( "Number of safe loans : " , len(safe_loans))
8 print( "Number of risky loans : " , len(risky_loans))
9 # Append the risky_loans with the downsampled version of safe_loans
10 train = risky_loans.append(safe_loans)
```

Number of safe loans : 3689

Number of risky loans : 3689



Logistic Regression

```
1 col_cat = ['NAME_CONTRACT_TYPE', 'CODE_GENDER', 'FLAG_OWN_CAR', 'FLAG_OWN_REALTY',  
2           'NAME_TYPE_SUITE', 'NAME_INCOME_TYPE', 'NAME_EDUCATION_TYPE', 'NAME_FAMILY_STATUS',  
3           'NAME_HOUSING_TYPE', 'OCCUPATION_TYPE', 'ORGANIZATION_TYPE']  
4 col_num = ['SK_ID_CURR', 'CNT_CHILDREN', 'AMT_INCOME_TOTAL', 'AMT_CREDIT', 'AMT_ANNUITY',  
5           'AMT_GOODS_PRICE', 'REGION_POPULATION_RELATIVE', 'DAYS_BIRTH', 'DAYS_EMPLOYED',  
6           'DAYS_REGISTRATION', 'DAYS_ID_PUBLISH', 'CNT_FAM_MEMBERS', 'REGION_RATING_CLIENT',  
7           'REGION_RATING_CLIENT_W_CITY', 'LIVE_REGION_NOT_WORK_REGION', 'EXT_SOURCE_1',  
8           'EXT_SOURCE_2', 'EXT_SOURCE_3', 'AMT_REQ_CREDIT_BUREAU_QRT', 'previous_loan_counts']  
9  
10 X_train_cat = X_train[col_cat]  
11 X_train_num = X_train[col_num]  
12 X_test_cat = X_test[col_cat]  
13 X_test_num = X_test[col_num]
```



Logistic Regression

```
9 pipe_cat = make_pipeline(SimpleImputer(strategy='constant'), OneHotEncoder(handle_unknown='ignore'))
10 pipe_num = make_pipeline(SimpleImputer(), MinMaxScaler())
11 preprocessor = make_column_transformer((pipe_cat, col_cat), (pipe_num, col_num))

1 make_pipeline(preprocessor, LogisticRegression(solver='lbfgs', multi_class='auto', random_state=42, max_iter=1000))

1 pipe.fit(X_train, y_train)
2 accuracy = pipe.score(X_test, y_test)
3 print('Accuracy score of the {} with MinMaxScaler is {:.3f}'.format(pipe.__class__.__name__, accuracy))
```

Accuracy score of the Pipeline with MinMaxScaler is 0.671

Cross Validation

```
1 from sklearn.model_selection import cross_validate
2
3 pipe = make_pipeline(preprocessor, LogisticRegression(solver='lbfgs', multi_class='auto',
4                                                       random_state=42, max_iter=1000))
5 scores = cross_validate(pipe, X, y, cv=3, return_train_score=True)
```

```
1 df_scores = pd.DataFrame(scores)
2 df_scores
```

| | fit_time | score_time | test_score | train_score |
|--|----------|------------|------------|-------------|
|--|----------|------------|------------|-------------|

| | | | | |
|---|------|------|------|------|
| 0 | 8.32 | 0.44 | 0.91 | 0.91 |
|---|------|------|------|------|

| | | | | |
|---|------|------|------|------|
| 1 | 8.18 | 0.41 | 0.91 | 0.91 |
|---|------|------|------|------|

| | | | | |
|---|------|------|------|------|
| 2 | 9.23 | 0.41 | 0.91 | 0.91 |
|---|------|------|------|------|

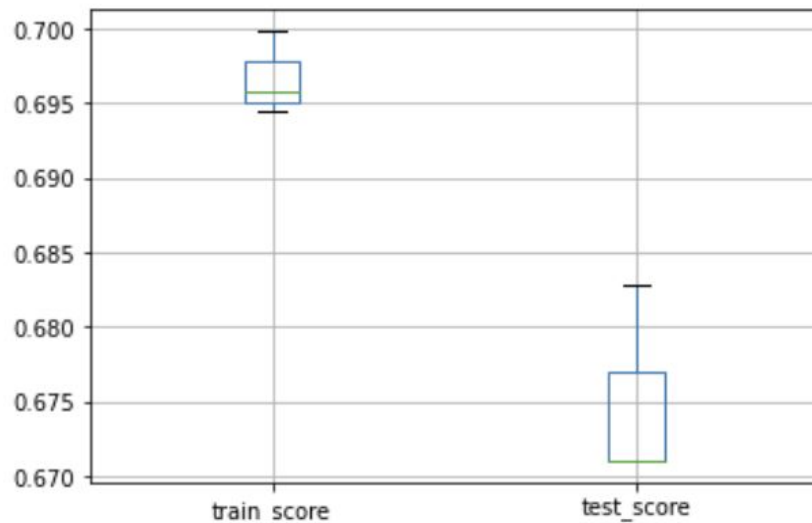

```
: 1 print("Mean times and scores:\n", df_scores.mean())
```

Mean times and scores:

```
fit_time    0.39
score_time   0.03
test_score   0.67
train_score  0.70
dtype: float64
```

```
: 1 df_scores[['train_score', 'test_score']].boxplot()
```

: <AxesSubplot:>



In testing set, the interquartile lies between 0.672 to 0.677, the mean is about 0.672.

```

1 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, balanced_accuracy_score, roc_auc_score
2
3 model_log_predicted = model_log.predict(X_test)
4
5 print('Accuracy of model_rf: {:.2f}'.format(accuracy_score(y_test, model_log_predicted)))
6 print('Precision of model_rf: {:.2f}'.format(precision_score(y_test, model_log_predicted)))
7 print('Recall of model_rf: {:.2f}'.format(recall_score(y_test, model_log_predicted)))
8 print('F1 of model_rf: {:.2f}'.format(f1_score(y_test, model_log_predicted)))
9 print('Balanced accuracy score of model_rf: {:.2f}'.format(balanced_accuracy_score(y_test, model_log_predicted)))
10 print('Roc auc score of model_rf: {:.2f}'.format(roc_auc_score(y_test, model_log_predicted)))

```

Accuracy of model_rf: 0.69

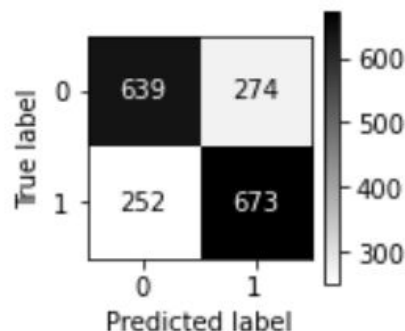
Precision of model_rf: 0.69

Recall of model_rf: 0.68

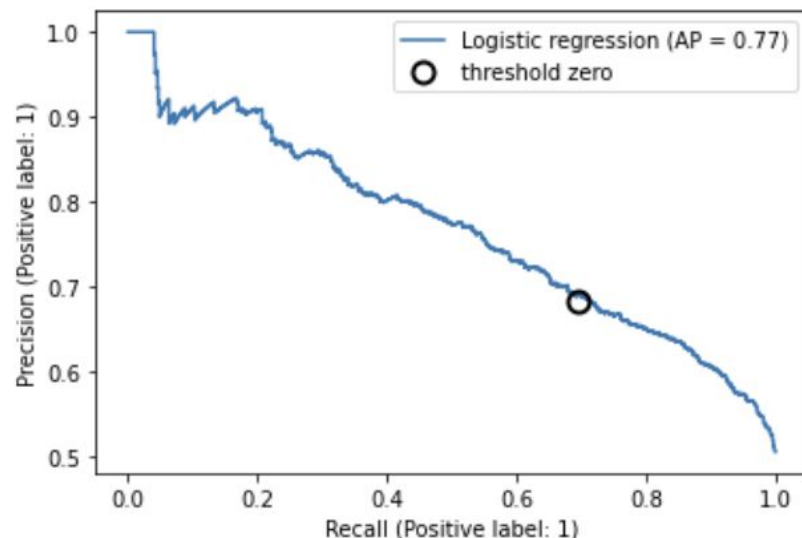
F1 of model_rf: 0.69

Balanced accuracy score of model_rf: 0.69

Roc auc score of model_rf: 0.69



<matplotlib.legend.Legend at 0x7fc23e30ddc0>





Linear SVC

```
1 from sklearn.svm import LinearSVC
2 clf = make_pipeline(preprocessor, LinearSVC(C=1.0, max_iter=500000))
3 clf.fit(X_train, y_train)
4 accuracy = clf.score(X_test, y_test)
5 print('Accuracy score of the {} is {:.3f}'.format(clf.__class__.__name__, accuracy))
```

Accuracy score of the Pipeline is 0.675



Decision Tree

```
from sklearn.tree import DecisionTreeClassifier, export_graphviz
pipe1 = make_pipeline(preprocessor, DecisionTreeClassifier(max_depth=2, criterion='entropy'))
pipe2 = make_pipeline(preprocessor, DecisionTreeClassifier(max_depth=6, criterion='entropy'))
pipe3 = make_pipeline(preprocessor, DecisionTreeClassifier(max_depth=10, criterion='entropy'))
pipe_cat = make_pipeline(SimpleImputer(strategy='constant'), OneHotEncoder(handle_unknown='ignore'))
pipe_num = make_pipeline(SimpleImputer(), MinMaxScaler())
preprocessor = make_column_transformer((pipe_cat, col_cat), (pipe_num, col_num))
```



Decision Tree

```
1 from sklearn.tree import DecisionTreeClassifier, export_graphviz
2 pipe1 = make_pipeline(preprocessor, DecisionTreeClassifier(max_depth=2, criterion='entropy'))
3 pipe2 = make_pipeline(preprocessor, DecisionTreeClassifier(max_depth=6, criterion='entropy'))
4 pipe3 = make_pipeline(preprocessor, DecisionTreeClassifier(max_depth=10, criterion='entropy'))
5
6
7 pipe1.fit(X_train, y_train)
8 accuracy = pipe1.score(X_test, y_test)
9 print('Accuracy score of the decision tree of depth 2 {} is {:.5f}'.format(pipe1.__class__.__name__, accuracy))
10
11 pipe2.fit(X_train, y_train)
12 accuracy = pipe2.score(X_test, y_test)
13 print('Accuracy score of the decision tree of depth 6 {} is {:.5f}'.format(pipe2.__class__.__name__, accuracy))
14
15 pipe3.fit(X_train, y_train)
16 accuracy = pipe3.score(X_test, y_test)
17 print('Accuracy score of the decision tree of depth 10 {} is {:.5f}'.format(pipe3.__class__.__name__, accuracy))
```

```
Accuracy score of the decision tree of depth 2 Pipeline is 0.62840
Accuracy score of the decision tree of depth 6 Pipeline is 0.65778
Accuracy score of the decision tree of depth 10 Pipeline is 0.62514
```



Random Forest

```
pipe_cat = make_pipeline(SimpleImputer(strategy='constant'), OneHotEncoder(handle_unknown='ignore'))
pipe_num = make_pipeline(SimpleImputer())
preprocessor = make_column_transformer((pipe_cat, col_cat), (pipe_num, col_num))
rf = RandomForestClassifier()
pipe_rf = make_pipeline(preprocessor, rf)
pipe_rf.get_params()
```

```
'randomforestclassifier__max_depth': None,
'randomforestclassifier__max_features': 'auto',
'randomforestclassifier__max_leaf_nodes': None,
'randomforestclassifier__max_samples': None,
'randomforestclassifier__min_impurity_decrease': 0.0,
'randomforestclassifier__min_impurity_split': None,
'randomforestclassifier__min_samples_leaf': 1,
'randomforestclassifier__min_samples_split': 2,
'randomforestclassifier__min_weight_fraction_leaf': 0.0,
'randomforestclassifier__n_estimators': 100,
'randomforestclassifier__n_jobs': None,
```




```
1 from sklearn.ensemble import RandomForestClassifier
2 np.set_printoptions(precision=6, suppress=True)
3 param_grid = {'randomforestclassifier__max_features': [2,4,6,8,10],
4               'randomforestclassifier__max_depth': [6,8,10,12,14]}
5 param_grid
```

```
{'randomforestclassifier__max_features': [2, 4, 6, 8, 10],
 'randomforestclassifier__max_depth': [6, 8, 10, 12, 14]}
```

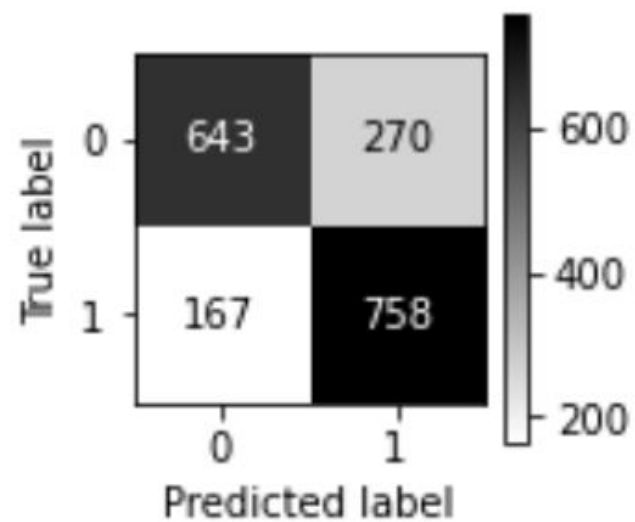
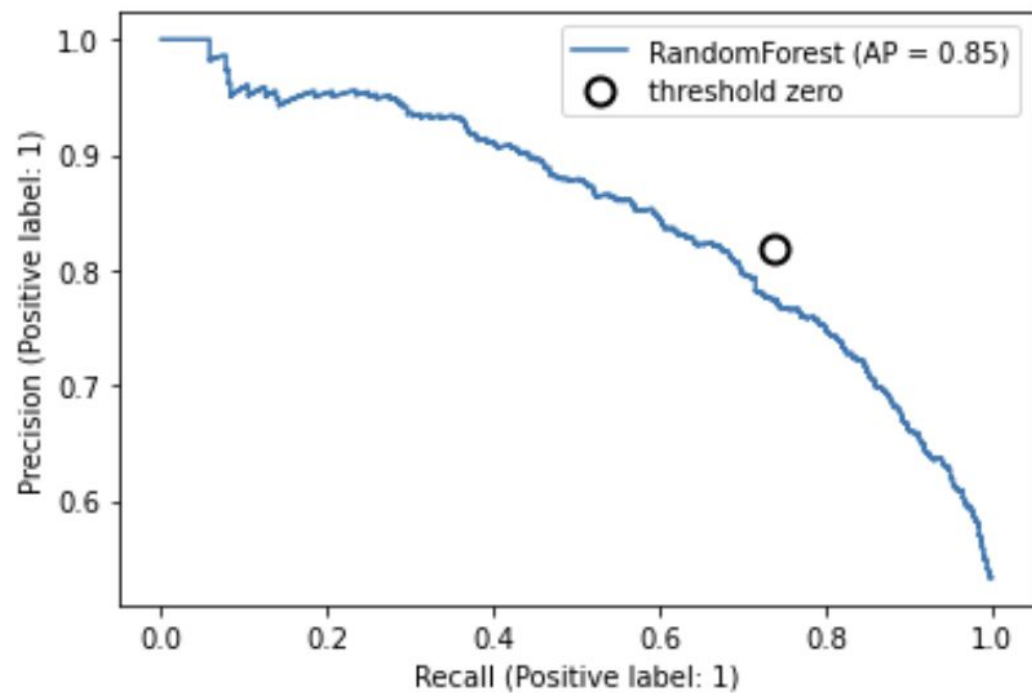
```
1 from sklearn.model_selection import GridSearchCV
2
3 grid_rf = GridSearchCV(pipe_rf, param_grid=param_grid, cv=3, return_train_score=True)
4 grid_rf.fit(X_train, y_train)
5
6 print(f"best parameters: {grid_rf.best_params_}")
7
8 print(f"test-set score: {grid_rf.score(X_test, y_test):.3f}")
```

```
best parameters: {'randomforestclassifier__max_depth': 6, 'randomforestclassifier__max_features': 10}
test-set score: 0.690
```



```
1 from sklearn.metrics import accuracy_score,precision_score,recall_score,f1_score,balanced_accuracy_score,roc_auc_score
2
3 model_rf_predicted = model_rf.predict(X_test)
4
5 print('Accuracy of model_rf: {:.2f}'.format(accuracy_score(y_test, model_rf_predicted)))
6 print('Precision of model_rf: {:.2f}'.format(precision_score(y_test, model_rf_predicted)))
7 print('Recall of model_rf: {:.2f}'.format(recall_score(y_test, model_rf_predicted)))
8 print('F1 of model_rf: {:.2f}'.format(f1_score(y_test, model_rf_predicted)))
9 print('Balanced accuracy score of model_rf: {:.2f}'.format(balanced_accuracy_score(y_test, model_rf_predicted)))
10 print('Roc auc score of model_rf: {:.2f}'.format(roc_auc_score(y_test, model_rf_predicted)))
```

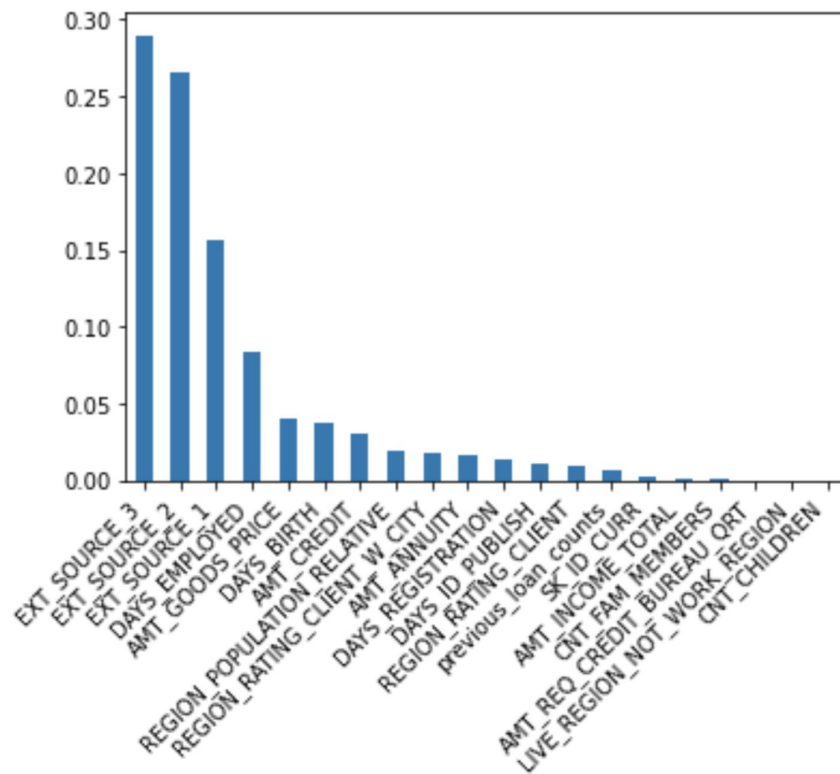
```
Accuracy of model_rf: 0.76
Precision of model_rf: 0.74
Recall of model_rf: 0.82
F1 of model_rf: 0.78
Balanced accuracy score of model_rf: 0.76
Roc auc score of model_rf: 0.76
```

Random Forest Feature Importance



Calculated with
feature_importances_



Summary of Findings



Conclusion / Recommendations

What're the results of the models?

Recall: 68% for logistic regression model, 82% for random forest

Attributes associated with loan repayment: price of goods for which the loan is given, days employed, credit amount, age

Loaning to applicants who banks deem at risk of defaulting is socially productive, yet extreme care must be taken in identifying those who should receive loans.

—

Thank you!