



UNIVERSITY OF TORONTO
MATHEMATICAL FINANCE

Kaggle Competition

HOME CREDIT DEFAULT RISK

Author

Keren Lai

Zhonghen Shen

Hongyi Wu

Tianyu(Shin) Ren

Contents

1	INTRODUCTION	3
1.1	Data	3
1.2	Exploratory Data Analysis	3
2	DATA TRANSFORMATION	5
2.1	Final Features Tables	5
2.2	Feature Engineering	8
2.2.1	Manual Generated Features	8
2.2.2	Features Selection	9
2.3	Documentation of QA	9
3	MODELS	9
3.1	Model Measures & Rationale	9
3.2	Model Development (Logistic Regression)	9
3.3	Model Development (LightGBM)	14
3.4	Model Development (LightGBM)	14
3.5	Model Tuning (LightGBM)	14
3.6	Model Prediction (LightGBM)	15
	Appendices	17

1. INTRODUCTION

Home Credit is a service dedicated to provided lines of loans to the unbanked population based on customer's previous records. This project is to help Home Credit to make predictions about clients capability of repayment using various statistical and machine learning methods.

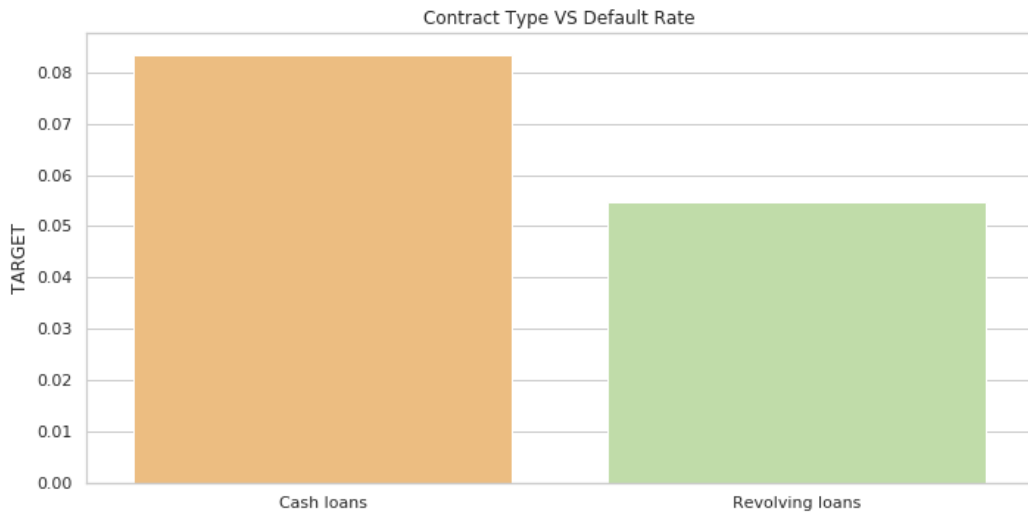
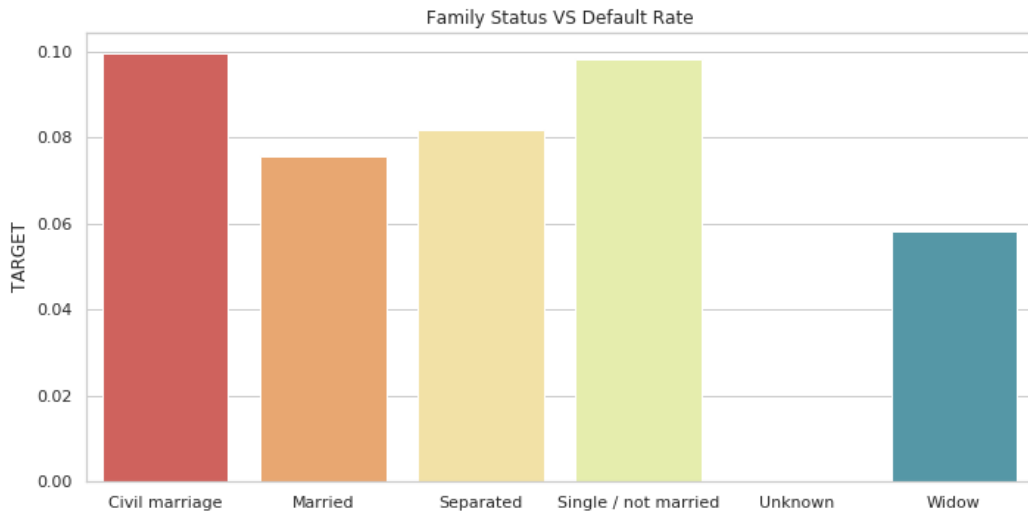
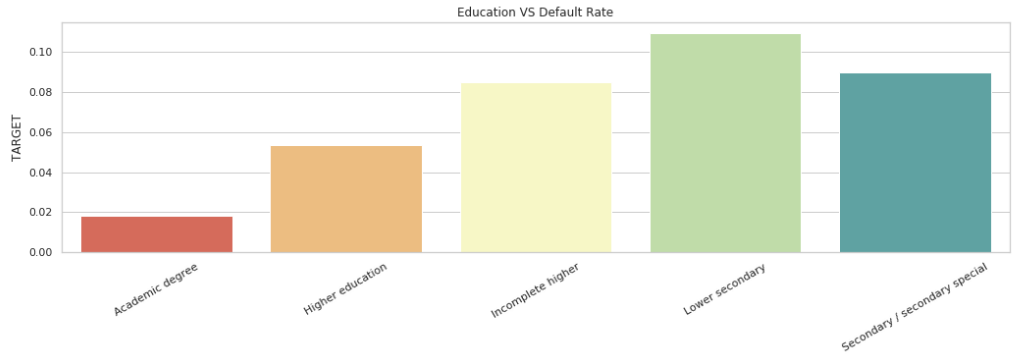
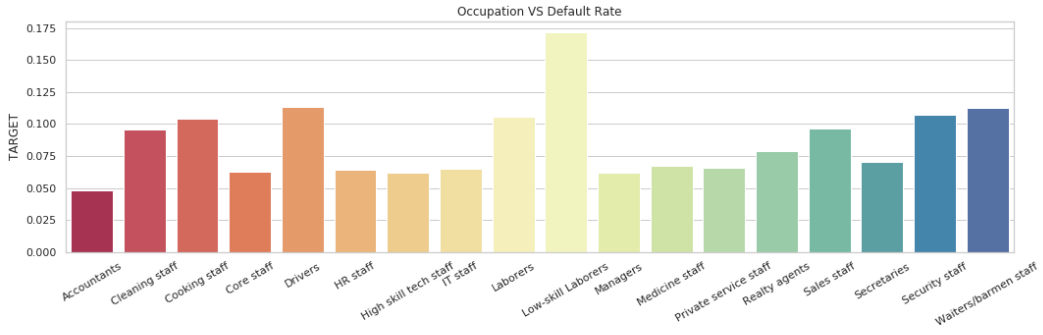
1.1. Data

The data we use is provided by Home Credit. There are 7 different sources of data:

- `application_train/application_test`: This is the main training and testing data we use to get information about each customers' loan application at Home Credit. We randomly choose 70% data from `application_train` to build our model, and use the rest 30% to validate the model; we produce Target for `application_test` indicating 0: the loan was likely to be repaid or 1: the load was likely to default;
- `bureau`: The data records clients's previous credits from other financial institutions;
- `bureau_balance`: This is monthly data records clients' previous credits from other financial institutions;
- `previous_application`: This is data records clients' previous applications for loans at Home Credit;
- `POS_CASH_balance`: This is monthly data about previous point of sale or cash loans clients have had at Home Credit;
- `credit_card_balance`: This is monthly data about previous credit cards clients have had at Home Credit;
- `installments_payment`: This is payment history of clients at Home Credit.

1.2. Exploratory Data Analysis

Before feature engineering, we calculated some statistics and make graphs to see trends, anomalies, patterns, or relationships with the data. To decide which features to use, we make figures showing the relationships between several variables and default rate.



2. DATA TRANSFORMATION

2.1. Final Features Tables

After features engineering, we use 233 features as final features:

FEATURES		
FLAG_OWN_CAR	NAME_CONTRACT_TYPE	CODE_GENDER
AMT_INCOME_TOTAL	FLAG_OWN_REALTY	CNT_CHILDREN
AMT_GOODS_PRICE	AMT_CREDIT	AMT_ANNUITY
NAME_EDUCATION_TYPE	NAME_TYPE_SUITE	NAME_INCOME_TYPE
REGION_POPULATION_RELATIVE	NAME_FAMILY_STATUS	NAME_HOUSING_TYPE
DAYS_REGISTRATION	DAYS_BIRTH	DAYS_EMPLOYED
FLAG_MOBIL	DAYS_ID_PUBLISH	OWN_CAR_AGE
FLAG_CONT_MOBILE	FLAG_EMP_PHONE	FLAG_WORK_PHONE
OCCUPATION_TYPE	FLAG_PHONE	FLAG_EMAIL
REGION_RATING_CLIENT_W_CITY	CNT_FAM_MEMBERS	REGION_RATING_CLIENT
REG_REGION_NOT_LIVE_REGION	WEEKDAY_APPR_PROCESS_START	HOUR_APPR_PROCESS_START
REG_CITY_NOT_LIVE_CITY	REG_REGION_NOT_WORK_REGION	LIVE_REGION_NOT_WORK_REGION
! ORGANIZATION_TYPE	REG_CITY_NOT_WORK_CITY	LIVE_CITY_NOT_WORK_CITY
EXT_SOURCE_3	EXT_SOURCE_1	EXT_SOURCE_2
OBS_60_CNT_SOCIAL_CIRCLE	OBS_30_CNT_SOCIAL_CIRCLE	DEF_30_CNT_SOCIAL_CIRCLE
FLAG_DOCUMENT_2	DEF_60_CNT_SOCIAL_CIRCLE	DAYS_LAST_PHONE_CHANGE
FLAG_DOCUMENT_5	FLAG_DOCUMENT_3	FLAG_DOCUMENT_4
FLAG_DOCUMENT_8	FLAG_DOCUMENT_6	FLAG_DOCUMENT_7
FLAG_DOCUMENT_11	FLAG_DOCUMENT_9	FLAG_DOCUMENT_10
FLAG_DOCUMENT_14	FLAG_DOCUMENT_12	FLAG_DOCUMENT_13
FLAG_DOCUMENT_17	FLAG_DOCUMENT_15	FLAG_DOCUMENT_16
FLAG_DOCUMENT_20	FLAG_DOCUMENT_18	FLAG_DOCUMENT_19
AMT_REQ_CREDIT_BUREAU_DAY	FLAG_DOCUMENT_21	AMT_REQ_CREDIT_BUREAU_HOUR
AMT_REQ_CREDIT_BUREAU_QRT	AMT_REQ_CREDIT_BUREAU_WEEK	AMT_REQ_CREDIT_BUREAU_MON
Still_Owed_min	AMT_REQ_CREDIT_BUREAU_YEAR	Still_Owed_max
Overdue_Days_max	Still_Owed_mean	Overdue_Days_std
	Overdue_Days_min	Overdue_Days_mean

StillOwed_std	NUM.INSTALMENT_NUMBER_std	Overdue_Days_sum
NUM.INSTALMENT_NUMBER_mean	NUM.INSTALMENT_VERSION_mean	NUM.INSTALMENT_VERSION_min
NUM.INSTALMENT_VERSION_max	NUM.INSTALMENT_VERSION_std	NUM.INSTALMENT_VERSION_sum
NUM.INSTALMENT_NUMBER_sum	NUM.INSTALMENT_NUMBER_min	NUM.INSTALMENT_NUMBER_max
StillOwed_sum	Excessive_Repay_std	Excessive_Repay_max
Excessive_Repay_min	Excessive_Repay_mean	AMT_DRAWINGS.OTHER.CURRENT_mean
AMT_DRAWINGS.POS.CURRENT_std	AMT_DRAWINGS.ATM.CURRENT_max	AMT_DRAWINGS.ATM.CURRENT_min
AMT_DRAWINGS.ATM.CURRENT_mean	AMT_DRAWINGS.OTHER.CURRENT_min	AMT_DRAWINGS.OTHER.CURRENT_max
AMT_DRAWINGS.OTHER.CURRENT_std	AMT_DRAWINGS.POS.CURRENT_mean	AMT_DRAWINGS.POS.CURRENT_min
AMT_DRAWINGS.POS.CURRENT_max	AMT_DRAWINGS.ATM.CURRENT_std	AMT.CREDIT.LIMIT.ACTUAL_min
AMT_DRAWINGS.ATM.CURRENT_sum	AMT.BALANCE_sum	AMT.BALANCE_std
AMT.BALANCE_max	AMT.CREDIT.LIMIT.ACTUAL_max	AMT.BALANCE_min
AMT.BALANCE_mean	MONTHS.BALANCE_sum	AMT.CREDIT.LIMIT.ACTUAL_std
AMT.CREDIT.LIMIT.ACTUAL_sum	MONTHS.BALANCE_std	MONTHS.BALANCE_max
MONTHS.BALANCE_min	AMT.CREDIT.LIMIT.ACTUAL_mean	AMT_DRAWINGS.CURRENT_mean
AMT.RECIVABLE_sum	AMT.RECIVABLE_std	AMT.RECIVABLE_max
AMT.RECIVABLE_min	AMT.RECIVABLE_mean	AMT_DRAWINGS.CURRENT_min
AMT_DRAWINGS.POS.CURRENT_sum	AMT_DRAWINGS.OTHER.CURRENT_sum	Excessive_Repay_sum
AMT_DRAWINGS.CURRENT_sum	AMT_DRAWINGS.CURRENT_std	AMT_DRAWINGS.CURRENT_max
MONTHS.BALANCE_mean	AMT.ANNUITY_mean	AMT.ANNUITY_min
AMT.ANNUITY_max	AMT.ANNUITY_std	AMT.ANNUITY_sum
AMT.APPLICATION_mean	AMT.APPLICATION_min	AMT.APPLICATION_max
AMT.APPLICATION_std	AMT.APPLICATION_sum	AMT.CREDIT_mean
AMT.CREDIT_min	AMT.CREDIT_max	AMT.CREDIT_std
AMT.CREDIT_sum	RATE.DOWN.PAYMENT_mean	RATE.DOWN.PAYMENT_min
RATE.DOWN.PAYMENT_max	RATE.DOWN.PAYMENT_std	RATE.DOWN.PAYMENT_sum
AMT.DOWN.PAYMENT_mean	AMT.DOWN.PAYMENT_min	AMT.DOWN.PAYMENT_max
AMT.DOWN.PAYMENT_std	AMT.DOWN.PAYMENT_sum	AMT.GOODS.PRICE_mean
AMT.GOODS.PRICE_min	AMT.GOODS.PRICE_max	AMT.GOODS.PRICE_std
AMT.GOODS.PRICE_sum	DAYS.DECISION_mean	DAYS.DECISION_min
DAYS.DECISION_max	DAYS.DECISION_std	DAYS.DECISION_sum
CNT.PAYMENT_mean	CNT.PAYMENT_min	CNT.PAYMENT_max
CNT.PAYMENT_std	CNT.PAYMENT_sum	b_feature_1
b_feature_2	b_feature_3	b_feature_4
b_feature_5	b_feature_6	b_feature_7
b_feature_auto_DAYS.CREDIT.ENDDATE_mean	b_feature_auto_DAYS.CREDIT.ENDDATE_min	b_feature_auto_DAYS.CREDIT.ENDDATE_max
b_feature_auto_DAYS.CREDIT.ENDDATE_std	b_feature_auto_DAYS.CREDIT.ENDDATE_sum	b_feature_auto_AMT.CREDIT.SUM_mean
b_feature_auto_AMT.CREDIT.SUM_min	b_feature_auto_AMT.CREDIT.SUM_max	b_feature_auto_AMT.CREDIT.SUM_std
b_feature_auto_AMT.CREDIT.SUM_sum	b_feature_auto_AMT.CREDIT.SUM.DEBT_mean	b_feature_auto_AMT.CREDIT.SUM.DEBT_min
b_feature_auto_AMT.CREDIT.SUM.DEBT_max	b_feature_auto_AMT.CREDIT.SUM.DEBT_std	b_feature_auto_AMT.CREDIT.SUM.DEBT_sum
b_feature_auto_AMT.CREDIT.MAX.OVERDUE_mean	b_feature_auto_AMT.CREDIT.MAX.OVERDUE_min	b_feature_auto_AMT.CREDIT.MAX.OVERDUE_max
b_feature_auto_AMT.CREDIT.MAX.OVERDUE_std	b_feature_auto_AMT.CREDIT.MAX.OVERDUE_sum	MONTHS.BALANCE_mean_pos

MONTHS_BALANCE_min_pos	MONTHS_BALANCE_max_pos	MONTHS_BALANCE_std_pos
MONTHS_BALANCE_sum_pos	CNT_INSTALMENT_mean_pos	CNT_INSTALMENT_min_pos
CNT_INSTALMENT_max_pos	CNT_INSTALMENT_std_pos	CNT_INSTALMENT_sum_pos
CNT_INSTALMENT_FUTURE_mean_pos	CNT_INSTALMENT_FUTURE_min_pos	CNT_INSTALMENT_FUTURE_max_pos
! CNT_INSTALMENT_FUTURE_std_pos	CNT_INSTALMENT_FUTURE_sum_pos	SK_DPD_mean_pos
SK_DPD_min_pos	SK_DPD_max_pos	SK_DPD_std_pos
SK_DPD_sum_pos	SK_DPD_DEF_mean_pos	SK_DPD_DEF_min_pos
SK_DPD_DEF_max_pos	SK_DPD_DEF_std_pos	SK_DPD_DEF_sum_pos
amt_diff	num_previous_loan	insurance_ratio

2.2. Feature Engineering

2.2.1. Manual Generated Features

Based on intuitive understanding, we have created 13 new features as following. Please see Appendix for the code details.

Features Name	Data Source	Methodology to Generate
amt_diff	previous_application	AMT_CREDIT - AMT_APPLICATION
num_previous_loan	previous_application	groupby SK_ID_CURR and count SK_ID_PREV
insurance_ratio	previous_application	percentage that previous loan at Home Credit was insured
Overdue_Days	installments_payments	DAYS_INSTALMENT - DAYS_ENTRY_PAYMENT
Still_Owed	installments_payments	AMT_INSTALMENT - AMT_PAYMENT
Excessive_Repay	credit_card_balance	AMT_PAYMENT_CURRENT - AMT_INST_MIN_REGULARITY
b_feature.1	bureau_data	how many loans will expire in the future
b_feature.2	bureau_data	number of loan this ID has (currently and in the past)
b_feature.3	bureau_data	the ratio of debt to the sum of debt and credit
b_feature.4	bureau_data	whether this ID has prolonged his loan before
b_feature.5	bureau_data	the furthest due date of this ID's loan in the future
b_feature.6	bureau_data	the ratio of overdue to debt
b_feature.7	bureau_data & application_data	the ratio of income to debt

- amt_diff is the difference between ask and given of previous loan at Home Credit.
- num_previous_loan is the total number of previous loan at Home Credit.
- insurance_ratio is the percentage that previous loan at Home Credit was insured.
- Overdue_Days is the delay of repaying installment.
- Still_Owed is the amount still owed after repayment
- Excessive_Repay is the excessive repayment

2.2.2. Features Selection

First, features with too many missing values should be dropped. For the logistic regression, we need to avoid multicollinearity problems and features with high correlation are dropped. Next, we drop columns of building/construction information which add many noises to our model. One method of this is computing feature importance. More details can be found in section 3.2. Model Development (Logistic Regression).

2.3. Documentation of QA

Final Features contains both automated generated features and manual generated features. For the automated generated features, we select useful numerical columns based on our domain knowledge on all datasets except for the 'application_train.csv'. Manual generated features and their methodologies have been illustrated above. Please see Appendix for more details.

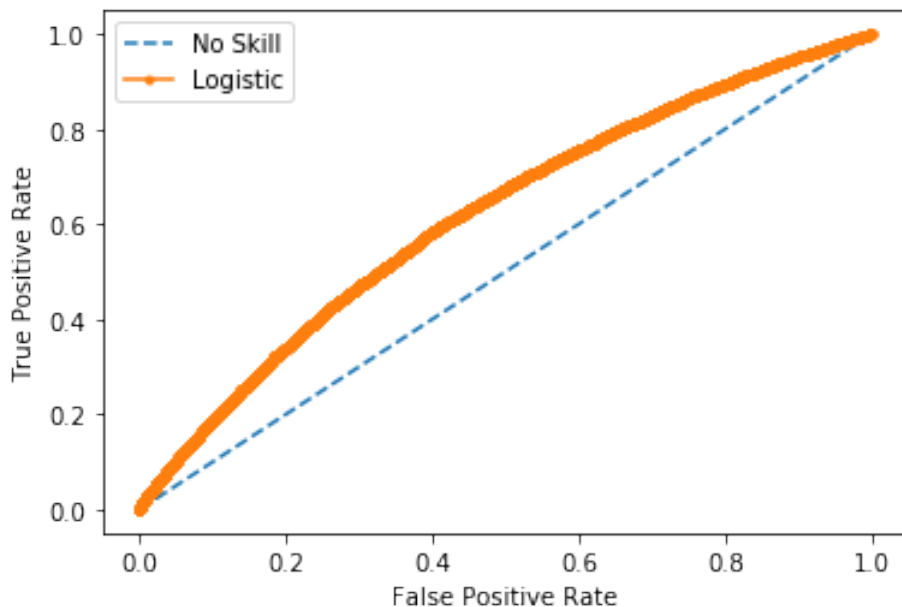
3. MODELS

3.1. Model Measures & Rationale

We use AUC as our model measure because for a binary classification problem, AUC is pretty good measure since it is easy to be interpreted and demonstrated. Also, submissions of this Kaggle competition are evaluated on AUC.

3.2. Model Development (Logistic Regression)

Logistic Regression is used as a benchmark since it is easy to train and interpret. However, we need to fill in NAs, which increases bias drastically to our model. We use mode to fill in categorical variables and median to fill in numerical variables. Before this, we transform binary variables via label encoding while multi-class variables via one-hot encoding. Also, features with more than 0.75 missing values are dropped. The following ROC curve is trained by logistic regression using this raw dataset.



LR Without Feature Engineering, AUC=0.620

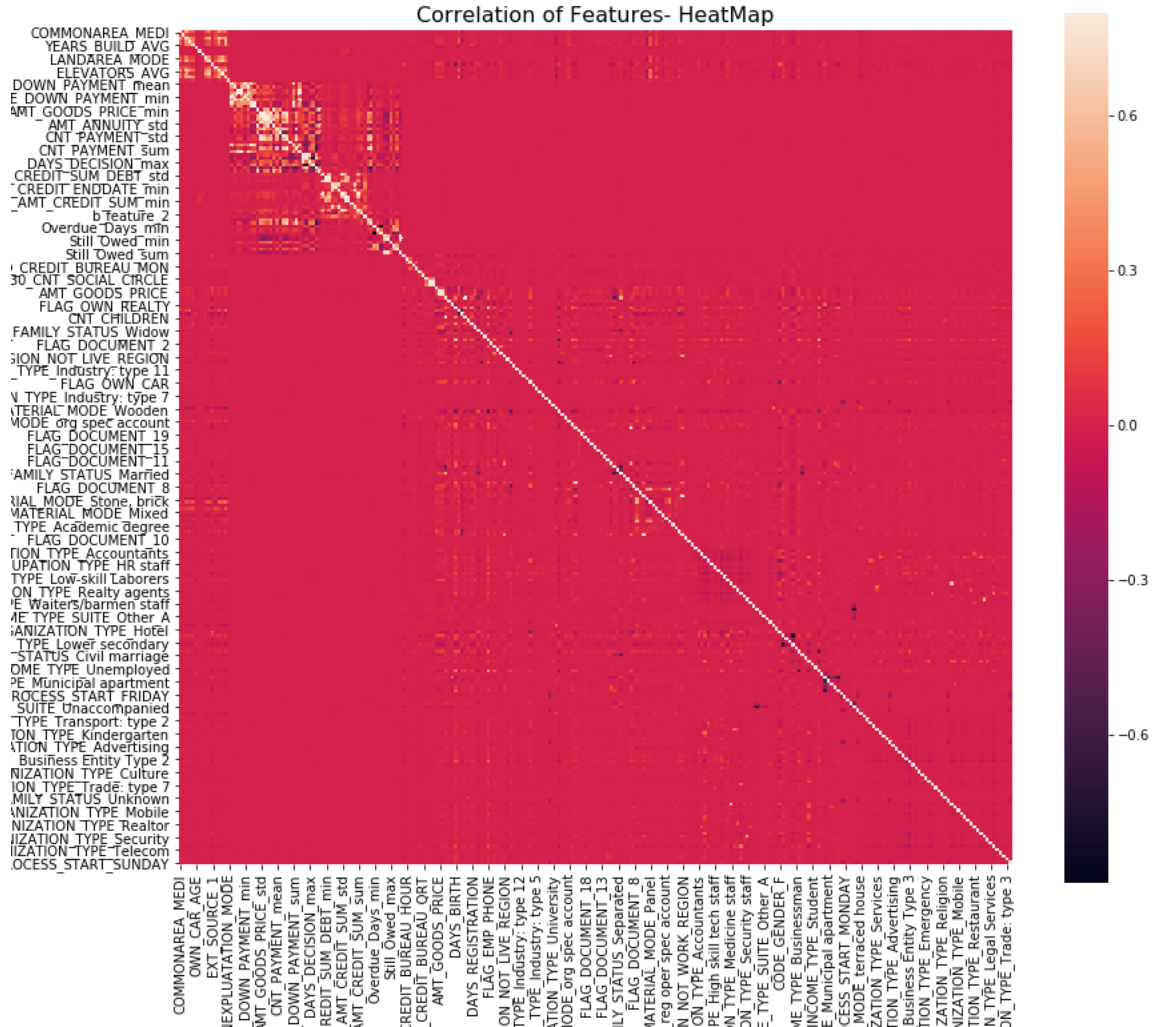
As we can see, the first AUC is not quite satisfying. Our next step is to perform additional data engineering and normalization to increase AUC.

To improve the performance of LR, further data engineering is conducted. It is possible that some features in provided datasets have extreme values, which might distort the information of the data. Therefore, we conduct extreme value detecting and transfer those values into missing values. Following table shows some examples of this processing:

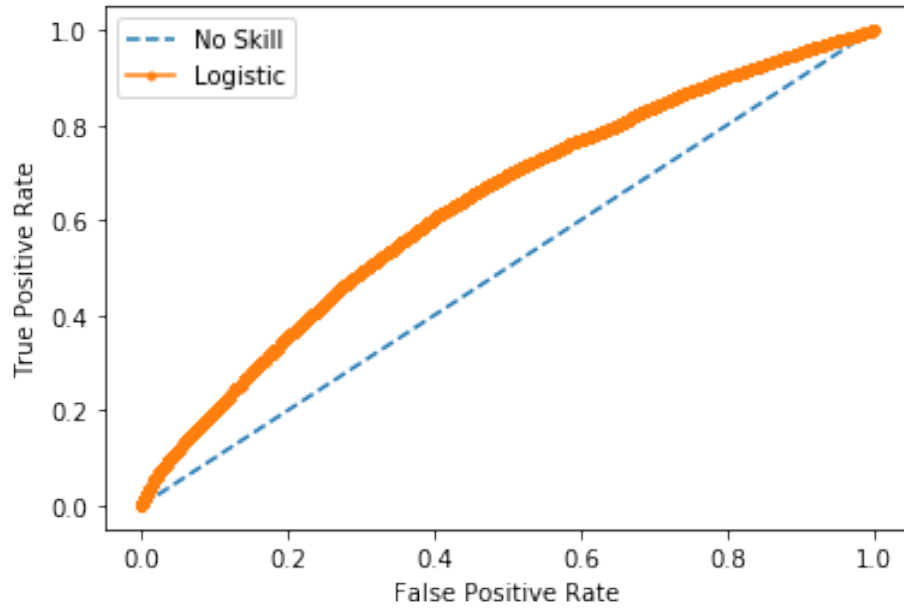
feature name	value range	processing method
OBS_30_CNT_SOCIAL_CIRCLE	(0, 354)	Turning values more than 100 into missing value
DAYS_CREDIT_ENDDATE	(-42020, 31199)	Turning values less than -2000 into missing value
AMT_CREDIT_SUM	(0, 585,000,000)	Turning values more than 10000000 into missing value

Besides, collinearity is considered. We first calculate the correlation matrix of our feature dataframe and then select the features with high correlation(≥ 0.9). Only one feature in the group with high correlation is kept.

The Heat Map of correlations are displayed below:

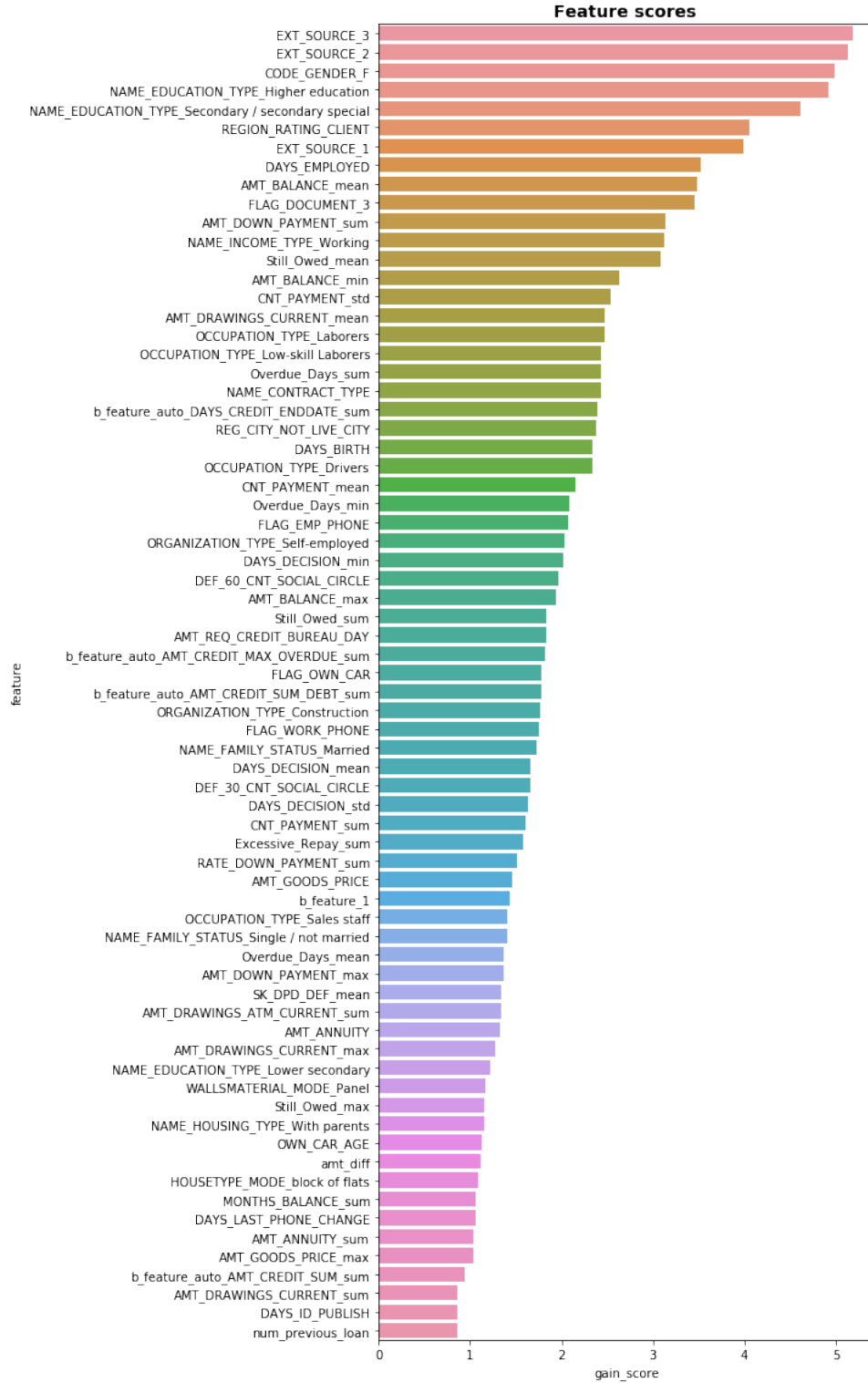


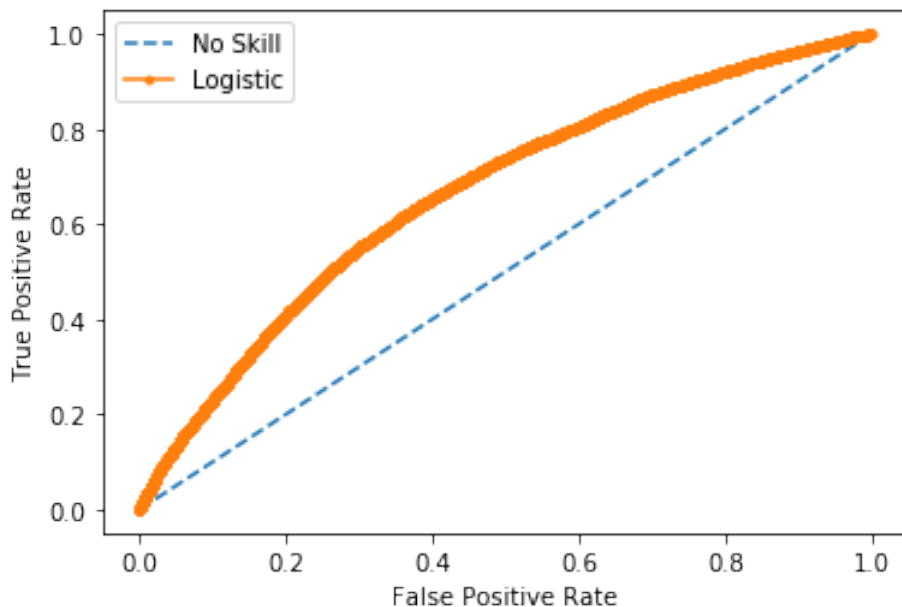
From the ROC diagram below, it is clear that the extreme value and correlation processing has improved the LR performance.



LR after further data processing, AUC=0.632

Finally, feature importance is utilized to select features. This is because some features in the former dataframe may have no contribution to the model performance. LightGBM is used to calculate the feature importance:





LR Without Feature Engineering, AUC=0.667

From the ROC diagram above we can know that after feature selection using feature importance, the AUC of logistic regression is further increased to 0.667.

3.3. Model Development (LightGBM)

In terms of the model selection, we choose LightGBM as our classifier. Currently the two most prevailing classification algorithms in Kaggle and other Data Science competitions are XGBoost and LightGBM. The reason that we choose LightGBM is it's 2 - 10 times faster than XGBoost.

3.4. Model Development (LightGBM)

3.5. Model Tuning (LightGBM)

Different hyper-parameter sets have quite large impact on the final result. We first initialized a parameters set, then applied to determine the best set of parameters.

(i) LightGBM Parameters Introduction

In LightGBM model, certain parameters plays very important roles and they are the targets of

hyper-parameter tuning. The following is a table demonstrating three important hyper-parameters and their functions.

hyper-parameters	function
num_leave	control the complexity of the tree model
min_data_in_leaf	prevent over-fitting in a leaf-wise tree
max_depth	similar to num_leave

Meanwhile, there are some universal hyper-parameters which exist in most machine learning models that worth of being tuned, such as learning-rate, regularization term, etc.

(ii) **Parameter Set**

We tuned our hyper-parameters from the following set:

(iii) **RandomizedSearchCV**

We use RandomizedSearchCV to sample a given number of candidates from a parameter set. Compared with grid search, randomized search is slightly worse in terms of performance, but due to the large amount of our dataset, we sacrificed a little accuracy for lowering the running time. The best parameter set we got is the following (for other parameters, we use the default values):

hyper-parameters	value
num_leave	30
reg_alpha	0.01
learning_rate	0.05

3.6. Model Prediction (LightGBM)

We split the whole dataset into training dataset and test dataset by the rule of "80-20". After training our model, we evaluated it on the test dataset. The confusion matrix and AUC score we got is the following:

	negative	positive
true	42241	3269
false	1606	14387
AUC	0.78059	

Appendices

```
1 !pip install kaggle
2 !mkdir .kaggle
3 import json
4 token = {"username":"xyz","key":"xxxxx"}
5 with open('/content/.kaggle/kaggle.json', 'w') as file:
6     json.dump(token, file)
7 !cp /content/.kaggle/kaggle.json ~/.kaggle/kaggle.json
8 !kaggle config set -n path -v{/content}
9 !cp /content/.kaggle/kaggle.json ~/.kaggle/kaggle.json
10 !kaggle config set -n path -v{/content}
11 !chmod 600 /root/.kaggle/kaggle.json
12 !mkdir Kaggle_data
13 !kaggle competitions download -c home-credit-default-risk -p /content/Kaggle_data
14 !cd Kaggle_data; unzip \*.zip
15 !cd Kaggle_data; rm *.zip
16 !rm -r sample_data/
17
18 # code of lightGBM
19 # data preprocessing
20 import pandas as pd
21 import lightgbm as lgb
22 from sklearn.model_selection import GridSearchCV, train_test_split
23
24 train_data = pd.read_csv('Kaggle_data/application_train.csv')
25 train_data.index = train_data['SK_ID_CURR']
26 del train_data['SK_ID_CURR']
27
28 # Label encoding of binary categorical variables
29 from sklearn.preprocessing import LabelEncoder
30 le = LabelEncoder()
```

```

31 le_count = 0
32
33 for col in train_data.columns:
34     if train_data[col].dtypes == 'object':
35         le.fit(train_data[col].astype(str))
36         train_data[col] = le.transform(train_data[col].astype(str))
37         le_count += 1
38
39 print('%d columns were label encoded.' % le_count)
40
41 # one-hot encoding of muti-categorical variables
42 # train_data = pd.get_dummies(train_data)
43 print('Training Features shape: ', train_data.shape)
44
45 # integrate multiple dataframes by left join
46 final_app_df = train_data.join(instal_pay_df, sort=False)
47 final_app_df = final_app_df.join(credit_card_df, sort=False)
48 final_app_df = final_app_df.join(pre_app_df, sort=False)
49 final_app_df = final_app_df.join(bureau_df, sort=False)
50 final_app_df = final_app_df.join(pos_cash_df, sort=False)
51 final_app_df = final_app_df.join(pre_manual_features_df, sort=False)
52
53 # specify categorical features
54 f = open('categorical_features.txt', 'r')
55 x = f.readlines()
56 f.close()
57 categorical_features = [f[:-1] for f in x]
58 categorical_features = categorical_features[:-1]
59
60 X_train, X_test, y_train, y_test = train_test_split(train_data_small.loc[:,
61 train_data_small.columns != 'TARGET'], train_data_small['TARGET'], test_size=0.2)
62
63 clf = lgb.LGBMClassifier(n_estimators=10000, objective = 'binary', is_unbalance=True,
64 metric='auc', learning_rate = 0.05,

```

```

65 reg_alpha = 0.1, reg_lambda = 0.1,
66 subsample = 0.8, n_jobs = -1, random_state = 50)
67
68 clf.fit(X_train, y_train, eval_metric = 'auc',
69         eval_set = [(X_test, y_test), (X_train, y_train)],
70         eval_names = ['test', 'train'], categorical_feature = categorical_features,
71         early_stopping_rounds = 100, verbose = 200)
72
73 # hyper-parameters tuning
74 from sklearn.metrics import confusion_matrix
75 y_pred = clf.predict(X_test)
76 print(confusion_matrix(y_test, y_pred))
77
78 from sklearn.model_selection import RandomizedSearchCV
79 def LgmParamTuning(X_train, y_train):
80     # application
81     # the last column of train_data is Y
82     param_dict = {"num_leaves": [25, 30, 35],
83                  "reg_alpha": [0, 0.05, 0.1],
84                  }
85     # "min_data_in_leaf": [300, 500, 800]}
86     gsearch = GridSearchCV(estimator = lgb.LGBMClassifier(num_iterations=10000,
87     objective = 'binary', learning_rate=0.05,
88     is_unbalance=True, metric='auc',
89     subsample = 0.8, n_jobs = -1, random_state=50, feature_fraction=0.8),
90     param_grid = param_dict, scoring='roc_auc')
91
92     gsearch.fit(X_train,y_train)
93     return gsearch.best_params_, gsearch.best_score_
94
95 # bureau_data analysis and generate new features
96 def bureau_data_feature_1(df):
97     # how many loans will expire in the future
98     return sum(df['DAYS_CREDIT_ENDDATE'] > 0) / len(df)

```

```

99
100 def bureau_data_feature_2(df):
101     # how many loans this person has now and past
102     return len(df)
103
104 def bureau_data_feature_3(df):
105     # the ratio of debt and the sum of debt and credit
106     return max(df['AMT_CREDIT_SUM_DEBT'] / (df['AMT_CREDIT_SUM_DEBT'] + df['AMT_CREDIT_SUM']))
107
108 def bureau_data_feature_4(df):
109     # whether this person has prolonged his loan
110     return 1 if np.sum(df['CNT_CREDIT_PROLONG']) > 0 else 0
111
112 def bureau_data_feature_5(df):
113     # the farrest due date of a person's loan in the future
114     pos_due_date_ind = df['DAYS_CREDIT_ENDDATE'] > 0
115     tmp = df[pos_due_date_ind]['DAYS_CREDIT_ENDDATE']
116     return 0 if tmp.empty else max(tmp)
117
118 def bureau_data_feature_6(df):
119     # the ratio of overdue and debt
120     return max(df['AMT_CREDIT_SUM_OVERDUE'] / df['AMT_CREDIT_SUM_DEBT'])
121
122 def bureau_data_feature_7_pre(df):
123     # ratio of income to debt
124     # use bureau data and application data
125     # df is a left join table by bureau_df and application_df
126     existing_loan = df['DAYS_CREDIT_ENDDATE'] > 0
127     tmp = df[existing_loan]['AMT_CREDIT_SUM_DEBT']
128
129     debt_sum = np.nan if tmp.empty else sum(tmp)
130     return debt_sum
131
132 def gen_new_feature_from_bureau_data(bureau_data, application_data):

```

```

133 grouped_data = bureau_data.groupby('SK_ID_CURR')
134 result = pd.DataFrame(index=bureau_data['SK_ID_CURR'].unique())
135 # handcrafted features
136 result['b_feature_1'] = grouped_data.apply(bureau_data_feature_1)
137 result['b_feature_2'] = grouped_data.apply(bureau_data_feature_2)
138 result['b_feature_3'] = grouped_data.apply(bureau_data_feature_3)
139 result['b_feature_4'] = grouped_data.apply(bureau_data_feature_4)
140 result['b_feature_5'] = grouped_data.apply(bureau_data_feature_5)
141 result['b_feature_6'] = grouped_data.apply(bureau_data_feature_6)
142 debt_sum = grouped_data.apply(bureau_data_feature_7_pre)
143 result['b_feature_7'] = application_df.loc[bureau_data['SK_ID_CURR'].unique()]
144 ['AMT_INCOME_TOTAL'] / debt_sum
145 # auto_generated features
146 auto_feature_list = ['DAYS_CREDIT_ENDDATE', 'AMT_CREDIT_SUM', 'AMT_CREDIT_SUM_DEBT',
147 'AMT_CREDIT_MAX_OVERDUE']
148 stats_dict = {'mean': np.mean, 'min': min, 'max': max, 'std': np.std, 'sum': sum}
149 for i in tqdm(range(0, len(auto_feature_list))):
150     for k in stats_dict.keys():
151         name = 'b_feature_auto_' + auto_feature_list[i] + '_' + k
152         result[name] = grouped_data[auto_feature_list[i]].apply(stats_dict[k])
153     result
154 return result

1 # Manual Feature: the difference between ask and given of previous loan at Home Credit
2 pre_app_df["amt_diff"] = pre_app_df["AMT_CREDIT"] - pre_app_df["AMT_APPLICATION"]
3 amt_diff = pre_app_df.groupby("SK_ID_CURR").sum()["amt_diff"].to_frame()
4 amt_diff.columns = ["amt_diff"]
5 app_train_df = app_train_df.join(amt_diff, on="SK_ID_CURR")
6 app_train_df["amt_diff"] = app_train_df["amt_diff"].fillna(value=0)
7
8 # Construct new features for installments_payments
9 feature_select = ['SK_ID_PREV', 'SK_ID_CURR', 'NUM_INSTALMENT_VERSION', 'NUM_INSTALMENT_NUMBER']
10 instal_pay_df = installments_payments[feature_select]

```

```

11 instal_pay_df['Overdue_Days'] = installments_payments['DAYS_INSTALMENT']
12 -installments_payments['DAYS_ENTRY_PAYMENT']
13 instal_pay_df['Still_Owed'] = installments_payments['AMT_INSTALMENT']
14 -installments_payments['AMT_PAYMENT']
15
16 # Construct new features for credit_card_balance
17 feature_select = ['SK_ID_PREV', 'SK_ID_CURR', 'MONTHS_BALANCE', 'AMT_BALANCE',
18                  'AMT_CREDIT_LIMIT_ACTUAL', 'AMT_DRAWINGS_ATM_CURRENT', 'AMT_DRAWINGS_CURRENT',
19                  'AMT_DRAWINGS_OTHER_CURRENT', 'AMT_DRAWINGS_POS_CURRENT', 'SK_DPD_DEF', 'AMT_RECIVABLE',
20                  'NAME_CONTRACT_STATUS']
21 cre_card_df = credit_card_balance[feature_select]
22 cre_card_df['Excessive_Repay'] = credit_card_balance['AMT_PAYMENT_CURRENT']
23 -credit_card_balance['AMT_INST_MIN_REGULARITY']
24
25 # Manual Feature: the total number of previous loan at Home Credit
26 num_previous_loan = pre_app_df.groupby("SK_ID_CURR").count()["SK_ID_PREV"].to_frame()
27 num_previous_loan.columns = ["num_previous_loan"]
28 app_train_df = app_train_df.join(num_previous_loan, on="SK_ID_CURR")
29 app_train_df["num_previous_loan"] = app_train_df["num_previous_loan"].fillna(value=0)
30
31 # Manual Feature: percentage that previous loan at Home Credit was insured
32 def insurance_ratio(df):
33     return sum(df['NFLAG_INSURED_ON_APPROVAL'] == 1) / len(df['NFLAG_INSURED_ON_APPROVAL'])
34 insurance_ratio = pre_app_df.groupby('SK_ID_CURR').apply(insurance_ratio).to_frame()
35 insurance_ratio.columns = ["insurance_ratio"]
36 app_train_df = app_train_df.join(insurance_ratio, on="SK_ID_CURR")

```

Listing 1: Manual Generated Features

```

1
2 def del_missing_much(df):
3     total = df.isnull().sum().sort_values(ascending = False)
4     percent = (df.isnull().sum()/df.shape[0]*100).sort_values(ascending = False)

```

```

5  missing_instal_pay_df = pd.concat([total, percent], axis=1, keys=['Total', 'Percent'])
6  missing_instal_pay_df = missing_instal_pay_df.loc[missing_instal_pay_df['Percent']<75]
7  column = missing_instal_pay_df.index
8  df_new = df[column]
9  df_new = df_new[list(df_new.columns)[:-1]]
10 return df[column]
11
12 bureau_df = del_missing_much(bureau_df)
13 pre_app_df = del_missing_much(pre_app_df)
14 cre_card_df = del_missing_much(cre_card_df)
15 instal_pay_df = del_missing_much(instal_pay_df)
16 app_train_df = del_missing_much(app_train_df)
17 extra_df = del_missing_much(extra_df)
18
19 app_train_df['AMT_INCOME_TOTAL'].replace({117000000: np.nan}, inplace = True)
20 app_train_df['DAYS_EMPLOYED'].replace({365243: np.nan}, inplace=True)
21 app_train_df['REGION_RATING_CLIENT_W_CITY'].replace({-1: np.nan}, inplace = True)
22 app_train_df['OBS_30_CNT_SOCIAL_CIRCLE'].replace({354: np.nan}, inplace = True)
23 app_train_df['OBS_60_CNT_SOCIAL_CIRCLE'].replace({344: np.nan}, inplace = True)
24 app_train_df['AMT_REQ_CREDIT_BUREAU_QRT'].replace({261: np.nan}, inplace = True)
25 idx = bureau['DAYS_CREDIT_ENDDATE'].loc[bureau['DAYS_CREDIT_ENDDATE']<-20000].index
26 bureau['DAYS_CREDIT_ENDDATE'].loc[idx] = np.nan
27 idx = bureau['DAYS_ENDDATE_FACT'].loc[bureau['DAYS_ENDDATE_FACT']<-4000].index
28 bureau['DAYS_ENDDATE_FACT'].loc[idx] = np.nan
29 idx = bureau['AMT_CREDIT_SUM'].loc[bureau['AMT_CREDIT_SUM']>10000000].index
30 bureau['AMT_CREDIT_SUM'].loc[idx] = np.nan
31 idx = bureau['AMT_CREDIT_SUM_DEBT'].loc[bureau['AMT_CREDIT_SUM_DEBT']<0].index
32 bureau['AMT_CREDIT_SUM'].loc[idx] = 0
33 idx = bureau['AMT_CREDIT_SUM_DEBT'].loc[bureau['AMT_CREDIT_SUM_DEBT']>50000000].index
34 bureau['AMT_CREDIT_SUM'].loc[idx] = np.nan
35 idx = bureau['AMT_ANNUITY'].loc[bureau['AMT_ANNUITY']>10000000].index
36 bureau['AMT_ANNUITY'].loc[idx] = np.nan
37 idx = pre_app_df['AMT_DOWN_PAYMENT'].loc[bureau['AMT_DOWN_PAYMENT']<0].index
38 pre_app_df['AMT_DOWN_PAYMENT'].loc[idx] = 0

```

```

39 idx = cre_card_df['AMT_DRAWINGS_ATM_CURRENT'].loc[bureau['AMT_DRAWINGS_ATM_CURRENT']<0].index
40 cre_card_df['AMT_DRAWINGS_ATM_CURRENT'].loc[idx] = np.nan
41 idx = cre_card_df['AMT_DRAWINGS_CURRENT'].loc[bureau['AMT_DRAWINGS_CURRENT']<0].index
42 cre_card_df['AMT_DRAWINGS_CURRENT'].loc[idx] = np.nan
43
44 # Label encoding of binary categorical variables
45 from sklearn.preprocessing import LabelEncoder
46 le = LabelEncoder()
47
48 def Transfer_category(df):
49
50     for col in df:
51         if df[col].dtype == 'object':
52             # If 2 or fewer unique categories
53             if len(list(df[col].unique())) <= 2:
54                 # Train on the training data
55                 le.fit(df[col])
56                 # Transform both training and testing data
57                 df[col] = le.transform(df[col])
58     # one-hot encoding of muti-categorical variables
59     df = pd.get_dummies(df)
60     try:
61         Obj_columns = df.select_dtypes(['object']).columns
62     except Exception:
63         pass
64     df[Obj_columns] = df[Obj_columns].fillna(df[Obj_columns].mode())
65     return df
66
67 app_train_df = Transfer_category(app_train_df)
68 instal_pay_df = Transfer_category(instal_pay_df)
69 cre_card_df = Transfer_category(cre_card_df)
70
71 # Threshold for removing correlated variables
72 threshold = 0.9

```



```

73
74 # Absolute value correlation matrix
75 corr_matrix = final_app_df.iloc[0:10000,:].corr().abs()
76 corr_matrix.head()
77
78 upper = corr_matrix.where(np.triu(np.ones(corr_matrix.shape), k=1).astype(np.bool))
79 upper.shape
80 # Select columns with correlations above threshold
81 to_drop = [column for column in upper.columns if any(upper[column] > threshold)]
82 final_app_df = final_app_df.drop(columns = to_drop)

```

Listing 2: Data Cleaning

```

1 from sklearn.feature_selection import RFE
2 rfe_selector = RFE(estimator=LogisticRegression(), n_features_to_select=150, step=10, verbose=5)
3 rfe_selector.fit(app_train_df, target)
4 rfe_support = rfe_selector.get_support()
5 rfe_feature = app_train_df.loc[:,rfe_support].columns.tolist()
6 print(str(len(rfe_feature)), 'selected features')

```

Listing 3: Feature Selection

```

1 final_app_df_original = final_app_df.fillna(final_app_df.mean())
2 final_app_df_original = final_app_df_original.replace([np.inf, -np.inf], np.nan)
3 final_app_df_original.dropna(axis =1, inplace =True )
4 # feature importance
5 def get_feature_importances(data, shuffle, seed=None):
6     # Gather real features
7     train_features = [f for f in data if f not in ['TARGET', 'SK_ID_CURR']]
8     # Go over fold and keep track of CV score (train and valid) and feature importances
9
10    # Shuffle target if required

```

```

11 y = data['TARGET'].copy()
12 if shuffle:
13     # Here you could as well use a binomial distribution
14     y = data['TARGET'].copy().sample(frac=1.0)
15
16 # Fit LightGBM in RF mode, yes it's quicker than sklearn RandomForest
17 dtrain = lgb.Dataset(data[train_features], y, free_raw_data=False, silent=True)
18 lgb_params = {
19     'objective': 'binary',
20     'boosting_type': 'rf',
21     'subsample': 0.623,
22     'colsample_bytree': 0.7,
23     'num_leaves': 127,
24     'max_depth': 8,
25     'seed': seed,
26     'bagging_freq': 1,
27     'n_jobs': 4
28 }
29
30 # Fit the model
31 clf = lgb.train(params=lgb_params, train_set=dtrain, num_boost_round=200)
32
33 # Get feature importances
34 imp_df = pd.DataFrame()
35 imp_df["feature"] = list(train_features)
36 imp_df["importance_gain"] = clf.feature_importance(importance_type='gain')
37 imp_df["importance_split"] = clf.feature_importance(importance_type='split')
38 imp_df['trn_score'] = roc_auc_score(y, clf.predict(data[train_features]))
39
40 return imp_df
41
42 df=final_app_df_original.drop('TARGET', 1)
43 f , ax = plt.subplots(figsize = (14,12))
44 plt.title('Correlation of Features- HeatMap',y=1,size=16)

```

```

45 sns.heatmap(df.corr(),square = True, vmax=0.8)
46
47
48 def logistic_regression(df, y=None):
49     if y is None:
50         y = df.pop("TARGET")
51
52     X_train, X_test, y_train, y_test = train_test_split(df, y, test_size=0.3, random_state=0)
53     ns_probs = [0 for _ in range(len(y_test))]
54     lr = LogisticRegression(random_state=0).fit(X_train, y_train)
55     lr_probs = lr.predict_proba(X_test)
56     lr_probs = lr_probs[:, 1]
57     ns_auc = roc_auc_score(y_test, ns_probs)
58     lr_auc = roc_auc_score(y_test, lr_probs)
59     ns_auc = roc_auc_score(y_test, ns_probs)
60     lr_auc = roc_auc_score(y_test, lr_probs)
61     print('Logistic: ROC AUC=%.3f' % (lr_auc))
62     ns_fpr, ns_tpr, _ = roc_curve(y_test, ns_probs)
63     lr_fpr, lr_tpr, _ = roc_curve(y_test, lr_probs)
64     plt.plot(ns_fpr, ns_tpr, linestyle='--', label='No Skill')
65     plt.plot(lr_fpr, lr_tpr, marker='.', label='Logistic')
66     plt.xlabel('False Positive Rate')
67     plt.ylabel('True Positive Rate')
68     plt.legend()
69     plt.show()

```

Listing 4: Logistic Regression

```

1 def auto_gen_new_feature(df):
2     grouped_data = df.groupby('SK_ID_CURR')
3     result = pd.DataFrame(index=app_train_df['SK_ID_CURR'].unique())
4     # auto_feature_list = ['AMT_ANNUITY', 'AMT_APPLICATION', 'AMT_CREDIT',
5     'RATE_DOWN_PAYMENT', 'AMT_DOWN_PAYMENT', 'AMT_GOODS_PRICE', 'DAYS_DECISION', 'CNT_PAYMENT']

```

```

6 auto_feature_list = ['MONTHS_BALANCE', 'CNT_INSTALMENT',
7 'CNT_INSTALMENT_FUTURE', 'SK_DPD', 'SK_DPD_DEF']
8 stats_dict = {'mean': np.mean, 'min': min, 'max': max, 'std': np.std, 'sum': sum}
9 for i in range(len(auto_feature_list)):
10     for k in tqdm(stats_dict.keys()):
11         name = auto_feature_list[i] + '_' + k
12         result[name] = grouped_data[auto_feature_list[i]].apply(stats_dict[k])
13 return result

```

Listing 5: Auto Generated Features