

```
In [3]: 1 !wget --header="Host: storage.googleapis.com" --header="User-Agent: Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7; rv:68.0) Gecko/20100101 Firefox/68.0" https://storage.googleapis.com/kaggle-competition-s-data/kaggle-v2/9120/860599/bundle/archive.zip?GoogleAccessId=web-data@kaggle-161607.iam.gserviceaccount.com&Expires=1581611378&Signature=XRrXJdQkSswHY5iMTnjVSXiaJaaVY7zrXqa0liII2EaggoXI3pDpfrY5DyIJhiYaobl1s1ZE6PVCr0H3%2BzSfb6ryuWmE02n9fqsAGN7wRxGn6u7xW8Y7YC258RVjvTVKLvp2Ls50Df4XexGqMTqteAhZ9nZIOPieUTCiqN3sqAPr65JRgghDoYBx4eAnvXqMi7vVUx%2B5wMk86qQV1Qv64H4%2BVbPn3%2BZ2BDJ0xcrc%2FK2W7lxs6VuSppzjbMH9MSJ9pqZMF2mI%2FfuBqDwf1ZwcaawK2GG7Gt20XZrCyuaVctcRT2GGShvIWpbjtj4q4G9Txo6Ry2o2N4YTosVuuP1pdR0A%3D%3D&response-content-disposition=attachment%3B+filename%3Dhome-credit-default-risk.zip (https://storage.googleapis.com/kaggle-competitions-data/kaggle-v2/9120/860599/bundle/archive.zip?GoogleAccessId=web-data@kaggle-161607.iam.gserviceaccount.com&Expires=1581611378&Signature=XRrXJdQkSswHY5iMTnjVSXiaJaaVY7zrXqa0liII2EaggoXI3pDpfrY5DyIJhiYaobl1s1ZE6PVCr0H3%2BzSfb6ryuWmE02n9fqsAGN7wRxGn6u7xW8Y7YC258RVjvTVKLvp2Ls50Df4XexGqMTqteAhZ9nZIOPieUTCiqN3sqAPr65JRgghDoYBx4eAnvXqMi7vVUx%2B5wMk86qQV1Qv64H4%2BVbPn3%2BZ2BDJ0xcrc%2FK2W7lxs6VuSppzjbMH9MSJ9pqZMF2mI%2FfuBqDwf1ZwcaawK2GG7Gt20XZrCyuaVctcRT2GGShvIWpbjtj4q4G9Txo6Ry2o2N4YTosVuuP1pdR0A%3D%3D&response-content-disposition=attachment%3B+filename%3Dhome-credit-default-risk.zip)
Resolving storage.googleapis.com (storage.googleapis.com)... 172.217.203.128, 2607:f8b0:400c:c04::80
Connecting to storage.googleapis.com (storage.googleapis.com)|172.217.203.128|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 721616255 (688M) [application/zip]
Saving to: 'home-credit-default-risk.zip'

home-credit-default 100%[=====] 688.19M  143MB/s   in 5.0s

2020-02-10 16:30:03 (138 MB/s) - 'home-credit-default-risk.zip' saved [721616255/721616255]
```

```
In [4]: 1 !unzip 'home-credit-default-risk.zip'

Archive: home-credit-default-risk.zip
  inflating: HomeCredit_columns_description.csv
  inflating: POS_CASH_balance.csv
  inflating: application_test.csv
  inflating: application_train.csv
  inflating: bureau.csv
  inflating: bureau_balance.csv
  inflating: credit_card_balance.csv
  inflating: installments_payments.csv
  inflating: previous_application.csv
  inflating: sample_submission.csv
```

```
In [5]: 1 !pip install lightgbm

Requirement already satisfied: lightgbm in /usr/local/lib/python3.6/dist-packages (2.2.3)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.6/dist-packages (from lightgbm) (0.22.1)
Requirement already satisfied: scipy in /usr/local/lib/python3.6/dist-packages (from lightgbm) (1.4.1)
Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from lightgbm) (1.17.5)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.6/dist-packages (from scikit-learn->lightgbm) (0.14.1)
```

```
In [0]: 1 import pandas as pd
        2
```

```

In [0]: 1 import gc
2 import pickle
3 import numpy as np
4 from lightgbm import LGBMClassifier
5 from sklearn.metrics import roc_auc_score, roc_curve, auc
6 from sklearn.model_selection import KFold, StratifiedKFold
7 import seaborn as sns
8 import warnings
9 warnings.simplefilter(action='ignore', category=FutureWarning)
10
11 def final_fun_1(t_df):
12     # code taken from https://www.kaggle.com/gemartin/load-data-reduce-m
13     def reduce_mem_usage(df):
14         """ iterate through all the columns of a dataframe and modify the
15             to reduce memory usage.
16         """
17         start_mem = df.memory_usage().sum() / 1024**2
18         #print('Memory usage of dataframe is {:.2f} MB'.format(start_mem))
19
20         for col in df.columns:
21             col_type = df[col].dtype
22
23             if col_type != object:
24                 c_min = df[col].min()
25                 c_max = df[col].max()
26                 if str(col_type)[:3] == 'int':
27                     if c_min > np.iinfo(np.int8).min and c_max < np.iinfo(
28                         df[col] = df[col].astype(np.int8)
29                     elif c_min > np.iinfo(np.int16).min and c_max < np.iin
30                         df[col] = df[col].astype(np.int16)
31                     elif c_min > np.iinfo(np.int32).min and c_max < np.iin
32                         df[col] = df[col].astype(np.int32)
33                     elif c_min > np.iinfo(np.int64).min and c_max < np.iin
34                         df[col] = df[col].astype(np.int64)
35                 else:
36                     if c_min > np.finfo(np.float16).min and c_max < np.fin
37                         df[col] = df[col].astype(np.float16)
38                     elif c_min > np.finfo(np.float32).min and c_max < np.f
39                         df[col] = df[col].astype(np.float32)
40                     else:
41                         df[col] = df[col].astype(np.float64)
42
43         end_mem = df.memory_usage().sum() / 1024**2
44         #print('Memory usage after optimization is: {:.2f} MB'.format(end_
45         #print('Decreased by {:.1f}%'.format(100 * (start_mem - end_mem) /
46
47         return df
48
49
50
51 def one_hot(df):
52     original_col = list(df.columns)
53     c = [c for c in df.columns if df[c].dtype == 'object']
54     df = pd.get_dummies(df, columns= c, dummy_na= True)
55     new_col = [c for c in df.columns if c not in original_col]
56     return df, new_col
57
58 df = reduce_mem_usage(pd.read_csv('application_train.csv'))
59 test_df = reduce_mem_usage(t_df)
60 df = df.append(test_df).reset_index()
61
62 df = df[df['CODE_GENDER'] != 'XNA']
63 df = df[df['AMT_GOODS_PRICE'].notnull()]
64 df = df[df['NAME_INCOME_TYPE'] != 'Maternity leave']
65 df = df[df['DAYS_LAST_PHONE_CHANGE'].notnull()]
66 df['DAYS_EMPLOYED'].replace(365243, np.nan, inplace= True)
67
68 df['DAYS BIRTH / DAYS EMPLOYED'] = df['DAYS BIRTH'] / df['DAYS EMPLO

```

```

In [0]: 1 import gc
2 import pickle
3 import numpy as np
4 from lightgbm import LGBMClassifier
5 from sklearn.metrics import roc_auc_score, roc_curve, auc
6 from sklearn.model_selection import KFold, StratifiedKFold
7 import seaborn as sns
8 import warnings
9 import pandas as pd
10 warnings.simplefilter(action='ignore', category=FutureWarning)
11
12 def final_fun_2(t_df,y):
13     # code taken from https://www.kaggle.com/gemartin/load-data-reduce-m
14     def reduce_mem_usage(df):
15         """ iterate through all the columns of a dataframe and modify the
16             to reduce memory usage.
17         """
18         start_mem = df.memory_usage().sum() / 1024**2
19         #print('Memory usage of dataframe is {:.2f} MB'.format(start_mem))
20
21         for col in df.columns:
22             col_type = df[col].dtype
23
24             if col_type != object:
25                 c_min = df[col].min()
26                 c_max = df[col].max()
27                 if str(col_type)[:3] == 'int':
28                     if c_min > np.iinfo(np.int8).min and c_max < np.iinfo(
29                         df[col] = df[col].astype(np.int8)
30                     elif c_min > np.iinfo(np.int16).min and c_max < np.iin
31                         df[col] = df[col].astype(np.int16)
32                     elif c_min > np.iinfo(np.int32).min and c_max < np.iin
33                         df[col] = df[col].astype(np.int32)
34                     elif c_min > np.iinfo(np.int64).min and c_max < np.iin
35                         df[col] = df[col].astype(np.int64)
36                 else:
37                     if c_min > np.finfo(np.float16).min and c_max < np.fin
38                         df[col] = df[col].astype(np.float16)
39                     elif c_min > np.finfo(np.float32).min and c_max < np.f
40                         df[col] = df[col].astype(np.float32)
41                     else:
42                         df[col] = df[col].astype(np.float64)
43
44         end_mem = df.memory_usage().sum() / 1024**2
45         #print('Memory usage after optimization is: {:.2f} MB'.format(end_
46         #print('Decreased by {:.1f}%'.format(100 * (start_mem - end_mem) /
47
48         return df
49
50
51
52 def one_hot(df):
53     original_col = list(df.columns)
54     c = [c for c in df.columns if df[c].dtype == 'object']
55     df = pd.get_dummies(df, columns= c, dummy_na= True)
56     new_col = [c for c in df.columns if c not in original_col]
57     return df, new_col
58
59 #df = reduce_mem_usage(pd.read_csv('application_train.csv'))
60 df = reduce_mem_usage(t_df)
61 #df = df.reset_index()
62
63 df = df[df['CODE_GENDER'] != 'XNA']
64 df = df[df['AMT_GOODS_PRICE'].notnull()]
65 df = df[df['NAME_INCOME_TYPE'] != 'Maternity leave']
66 df = df[df['DAYS_LAST_PHONE_CHANGE'].notnull()]
67 df['DAYS_EMPLOYED'].replace(365243, np.nan, inplace= True)
68

```

```
In [32]: 1 test_df=pd.read_csv('application_test.csv')
          2 test_df=test_df[:1]
          3 test=final_fun_1(test_df)
          4 test[['SK_ID_CURR', 'TARGET']]
```

```
Out[32]:
```

	SK_ID_CURR	TARGET
	307511	100001 0.035038

```
In [33]: 1 train_df=pd.read_csv('application_train.csv')
          2 train_df=train_df.reset_index()
          3 y=train_df['TARGET']
          4 print(final_fun_2(train_df.v))

0.7836583223570655
```

```
In [0]: 1
```