

▼ UCI ML Repo - credit card defaults

Training

Platform: Python 3, colab.research.google.com

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 import pandas as pd
4 from joblib import dump, load
5 from google.colab import drive
6 from sklearn.feature_selection import SelectKBest
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.metrics import accuracy_score, f1_score, recall_score, precision_score, log_loss
9 from sklearn.model_selection import train_test_split, GridSearchCV
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.pipeline import Pipeline
```

▼ Load data

```
1 drive.mount('/content/gdrive', force_remount=False)
```

☞ Drive already mounted at /content/gdrive; to attempt to forcibly remount, call drive.mount("/content/gdrive"

```
1 loc = "/content/gdrive/My Drive/Colab Notebooks/uci-credit-card-defaults/data/defaults_clean.csv"
2 data = pd.read_csv(loc, header=0)
3 data.shape
```

☞ (30000, 25)

```
1 data.head(5)
```

↗

	id	limit_bal	sex	education	marriage	age	pay_0	pay_2	pay_3	pay_4	...	bill_amt4	bill_amt5	bill_ar
0	1	20000	2	2	1	24	2	2	-1	-1	...	0	0	
1	2	120000	2	2	2	26	-1	2	0	0	...	3272	3455	3
2	3	90000	2	2	2	34	0	0	0	0	...	14331	14948	15
3	4	50000	2	2	1	37	0	0	0	0	...	28314	28959	29
4	5	50000	1	2	1	57	-1	0	-1	0	...	20940	19146	19

5 rows × 25 columns

```
1 data.tail(5)
```

↗

	id	limit_bal	sex	education	marriage	age	pay_0	pay_2	pay_3	pay_4	...	bill_amt4	bill_amt5	
29995	29996	220000	1	3	1	39	0	0	0	0	...	88004	31237	
29996	29997	150000	1	3	2	43	-1	-1	-1	-1	...	8979	5190	
29997	29998	30000	1	2	2	37	4	3	2	-1	...	20878	20582	
29998	29999	80000	1	3	1	41	1	-1	0	0	...	52774	11855	
29999	30000	50000	1	2	1	46	0	0	0	0	...	36535	32428	

5 rows × 25 columns

```
1 type_dict = {}
2 for i in ["limit_bal", "bill_amt1", "bill_amt2", "bill_amt3", "bill_amt4",
3          "bill_amt5", "bill_amt6", "pay_amt1", "pay_amt2", "pay_amt3",
4          "pay_amt4", "pay_amt5", "pay_amt6"]:
5     if i in data.columns:
6         type_dict[i] = "float64"
7 data = data.astype(type_dict)
8 data.dtypes
```

↗

id	int64
limit_bal	float64
sex	int64
education	int64
marriage	int64
age	int64
pay_0	int64
pay_2	int64
pay_3	int64
pay_4	int64
pay_5	int64
pay_6	int64
bill_amt1	float64
bill_amt2	float64
bill_amt3	float64
bill_amt4	float64
bill_amt5	float64
bill_amt6	float64
pay_amt1	float64
pay_amt2	float64
pay_amt3	float64
pay_amt4	float64
pay_amt5	float64
pay_amt6	float64
default	int64

```
1 data.describe().
```



	id	limit_bal	sex	education	marriage	age	pay_0	pay_2
count	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000	30000.000000
mean	15000.500000	167484.322667	1.603733	1.842267	1.557267	35.485500	-0.016700	-0.133767

▼ Training

```

1 def view_metrics(y_test, y_pred):
2     print("Accuracy: {}".format(accuracy_score(y_test, y_pred)))
3     print("Precision: {}".format(precision_score(y_test, y_pred)))
4     print("Recall: {}".format(recall_score(y_test, y_pred)))
5     print("F1: {}".format(f1_score(y_test, y_pred)))

1 # split train/test
2 y = data.loc[:, "default"]
3 X = data.drop(["id", "default"], axis=1)
4 X_train, X_test, y_train, y_test = train_test_split(
5     X, y, test_size=0.3, random_state=42)
6 loc_data = "/content/gdrive/My Drive/Colab Notebooks/uci-credit-card-defaults/data"
7 # save test file for evaluation
8 X_test.to_csv("{}defaults_clean_test_X.csv".format(loc_data), index=False)
9 y_test.to_csv("{}defaults_clean_test_y.csv".format(loc_data), index=False)
10 X_train.head(5)

```

	limit_bal	sex	education	marriage	age	pay_0	pay_2	pay_3	pay_4	pay_5	...	bill_amt3	bill_amt4
28465	240000.0	2	1	1	40	-2	-2	-2	-2	-2	...	0.0	0.0
27622	50000.0	2	1	2	23	-1	-1	-1	-1	-1	...	2299.0	4800.0
28376	50000.0	2	2	1	36	2	2	2	2	0	...	49125.0	47956.0
10917	200000.0	2	3	1	54	6	5	4	3	2	...	104686.0	102549.0
27234	240000.0	1	1	1	35	-1	-1	-1	0	-1	...	21790.0	17102.0

5 rows x 23 columns

```
1 transformed_features = None
```

▼ Tranformation pipeline classes

```
1 class StdCol():
2     """
3     Standardizes column 'col_name' in a pipeline
4     """
5     def __init__(self, col_name):
6         self.col_name = col_name
7
8     def fit(self, X, y=None):
9         self.ss = StandardScaler()
10        self.ss.fit(X.loc[:, self.col_name].values.reshape(-1, 1))
11        return self
12
13    def transform(self, X):
14        X = X.copy(deep=True)
15        X.loc[:, self.col_name] = self.ss.transform(X.loc[:, self.col_name].values.reshape(-1, 1))
16        global transformed_features
17        transformed_features = X
18        return X
```

```
1 class LogCol():
2     """
3     Log transforms column 'col_name' in a pipeline
4     """
5     def __init__(self, col_name):
6         self.col_name = col_name
7
8     def fit(self, X, y=None):
9         return self
10
11    def transform(self, X):
12        X = X.copy(deep=True)
13        X.loc[X[self.col_name]<1, self.col_name] = 1
14        X.loc[:, self.col_name] = np.log(X.loc[:, self.col_name])
15        global transformed_features
16        transformed_features = X
17        return X
```

```
1 class CategoricalColInt():
2     """
3     Tranforms column 'col_name' into n-1 categorical columns
4     """
5     def __init__(self, col_name):
```

```

6         self.col_name = col_name
7
8     def fit(self, X, y=None):
9         return self
10
11    def transform(self, X):
12        X = X.copy(deep=True)
13        dummies = pd.get_dummies(X.loc[:, self.col_name], prefix=self.col_name)
14        dummies_cols = list(dummies.columns)[1:] # drop last new category to avoid feature correlation
15        X = X.merge(dummies.loc[:, dummies_cols], left_index=True, right_index=True)
16        X = X.drop(columns=self.col_name, axis=1)
17        global transformed_features
18        transformed_features = X
19        return X

```

```

1 class AveColumns():
2     """
3     Calc average of columns
4     """
5     def __init__(self, feature_name, cols):
6         self.feature_name = feature_name
7         self.cols = cols
8         pass
9
10    def fit(self, X, y=None):
11        return self
12
13    def transform(self, X):
14        X = X.copy(deep=True)
15        X.loc[:, self.feature_name] = X.loc[:, self.cols].mean(axis=1)
16        global transformed_features
17        transformed_features = X
18        return X

```

```

1 class StDevColumns():
2     """
3     Calc standard deviation of columns cols
4     """
5     def __init__(self, feature_name, cols):
6         self.feature_name = feature_name
7         self.cols = cols
8         pass
9
10    def fit(self, X, y=None):
11        return self
12
13    def transform(self, X):
14        X = X.copy(deep=True)
15        X.loc[:, self.feature_name] = X.loc[:, self.cols].std(axis=1).

```

```

16         global transformed_features
17         transformed_features = X
18         return X

1 class SelectKBestFeatures():
2     """
3     Selects num of features to K_best
4     """
5     def __init__(self):
6         pass
7
8     def fit(self, X, y=None):
9         self.selector = SelectKBest(k=20).fit(X, y)
10        return self
11
12    def transform(self, X):
13        X = X.copy(deep=True)
14        return self.selector.transform(X)

```

▼ Train pipeline

```

1 pipe_list = []
2 ave_cols = {"avg_bill": ["bill_amt1", "bill_amt2", "bill_amt3", "bill_amt4", "bill_amt5", "bill_amt6"],
3             "avg_pay": ["pay_amt1", "pay_amt2", "pay_amt3", "pay_amt4", "pay_amt5", "pay_amt6"]}
4 st_dev_cols = {"st_dev_bill": ["bill_amt1", "bill_amt2", "bill_amt3", "bill_amt4", "bill_amt5", "bill_amt6"],
5               "st_dev_pay": ["pay_amt1", "pay_amt2", "pay_amt3", "pay_amt4", "pay_amt5", "pay_amt6"]}
6 categorical_cols = ["sex", "education", "marriage"]
7 log_cols = ["limit_bal", "avg_bill", "avg_pay",
8             "bill_amt1", "bill_amt2", "bill_amt3", "bill_amt4", "bill_amt5", "bill_amt6",
9             "pay_amt1", "pay_amt2", "pay_amt3", "pay_amt4", "pay_amt5", "pay_amt6"]
10 standardize_cols = ["age"]
11 for i in ave_cols:
12     pipe_list.append(("enc_ave_"+i, AveColumns(i, ave_cols[i])))
13 for i in st_dev_cols:
14     pipe_list.append(("enc_stdev_"+i, StDevColumns(i, st_dev_cols[i])))
15 for i in categorical_cols:
16     if i in X_train.columns:
17         pipe_list.append(("enc_cat_"+i, CategoricalColInt(col_name=i)))
18 for i in log_cols:
19     if i in X_train.columns:
20         pipe_list.append(("enc_log_"+i, LogCol(col_name=i)))
21 for i in standardize_cols:
22     if i in X_train.columns:
23         pipe_list.append(("enc_std_"+i, StdCol(col_name=i)))
24 pipe_list.append(("k_best_selector", SelectKBestFeatures()))
25 pipe_list.append(("model", LogisticRegression()))
26 pipeline = Pipeline(pipe_list)

```

```
27 model = pipeline.fit(X_train, y_train)
28 transformed_features.columns
```

```
↳ /usr/local/lib/python3.6/dist-packages/sklearn/utils/validation.py:595: DataConversionWarning: Data with inq
   warnings.warn(msg, DataConversionWarning)
   /usr/local/lib/python3.6/dist-packages/sklearn/utils/validation.py:595: DataConversionWarning: Data with inq
   warnings.warn(msg, DataConversionWarning)
   /usr/local/lib/python3.6/dist-packages/sklearn/linear_model/logistic.py:433: FutureWarning: Default solver w
   FutureWarning)
   Index(['limit_bal', 'age', 'pay_0', 'pay_2', 'pay_3', 'pay_4', 'pay_5',
         'pay_6', 'bill_amt1', 'bill_amt2', 'bill_amt3', 'bill_amt4',
         'bill_amt5', 'bill_amt6', 'pay_amt1', 'pay_amt2', 'pay_amt3',
         'pay_amt4', 'pay_amt5', 'pay_amt6', 'avg_bill', 'avg_pay',
         'st_dev_bill', 'st_dev_pay', 'sex_2', 'education_2', 'education_3',
         'education_4', 'marriage_2', 'marriage_3'],
        dtype='object')
```

```
1 y_pred = model.predict(X_train)
2 view_metrics(y_train, y_pred).
```

```
↳ Accuracy: 0.8031428571428572
   Precision: 0.6782894736842106
   Recall: 0.22048759623609923
   F1: 0.3327953518398967
   /usr/local/lib/python3.6/dist-packages/sklearn/utils/validation.py:595: DataConversionWarning: Data with inq
   warnings.warn(msg, DataConversionWarning)
```

```
1 # parameters = {
2 #     "model__penalty": ["l2"],
3 #     "model__solver": ["lbfgs", "liblinear"],
4 #     "model__max_iter": [50, 100],
5 #     "model__C": [0.7, 0.3, 0.1]}
6 # grid = GridSearchCV(pipeline, parameters, cv=4, scoring="f1")
7 # grid.fit(X_train, y_train)
8 # print(grid.best_params_)
9 # y_pred = grid.predict(X_train)
10 # view_metrics(y_train, y_pred)
```



```
1 # print("y_train")
2 # print(np.array(y_train)[:200])
3 # print("y_pred")
4 # print(y_pred[:200])
```

```
1 # loc_model = "/content/gdrive/My Drive/Colab Notebooks/uci-credit-card-defaults"
2 # dump(grid.best_estimator_, "{}model.joblib".format(loc_model))
```

```
1
```