

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import missingno as msno
```

▶ Load the dirty churn data [...]

▶ Load the clean churn rate data [...]

▼ Data Preprocessing

```
In [5]: 1 df = df.drop(columns="userid")
```

▼ One Hot Encoder (before tts)

```
In [6]: 1 df = pd.get_dummies(data=df, drop_first=True)
2 df
```

```
Out[6]:
```

	churn	age	trivia_played	trivia_shared_results	trivia_view_unlocked	trivia_view_results	cards_share	cards_viewed	cards_helpful	cards_not_helpful
0	1	22.0	0	0	0	0	0	48	0	
1	1	25.0	0	0	5	0	0	5	0	
2	1	32.0	0	0	0	0	0	49	0	
3	1	26.0	0	0	169	0	0	184	0	
4	1	28.0	0	0	11	0	0	65	0	
...
48074	0	42.0	0	0	1	0	0	6	0	
48075	0	38.0	0	0	0	0	0	207	0	
48076	0	20.0	5	0	51	5	0	132	0	
48077	0	27.0	0	0	77	0	0	143	0	
48078	0	33.0	3	2	18	3	0	20	0	

48079 rows × 42 columns

```
In [7]: 1 # Don't know how to work OneHotEncoder
2
3 # from sklearn.compose import ColumnTransformer
4 # from sklearn.preprocessing import OneHotEncoder
5
6 # encoder = OneHotEncoder(handle_unknown="ignore",)
7 # df2 = encoder.fit_transform(df)
```

```
In [8]: 1 X = df.drop(labels="churn", axis=1).to_numpy()
2 X
```

```
Out[8]: array([[22., 0., 0., ..., 1., 0., 0.],
 [25., 0., 0., ..., 0., 0., 0.],
 [32., 0., 0., ..., 0., 1., 0.],
 ...,
 [20., 5., 0., ..., 0., 0., 0.],
 [27., 0., 0., ..., 0., 1., 0.],
 [33., 3., 2., ..., 0., 0., 0.]])
```

```
In [9]: 1 y = df["churn"].values
2 y
```

```
Out[9]: array([1, 1, 1, ..., 0, 0, 0], dtype=int64)
```

▼ tts

```
In [10]: 1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
```

▼ Feature Scaling

```
In [11]: 1 from sklearn.preprocessing import StandardScaler, MinMaxScaler
2
3 sc = StandardScaler()
4 X_train = sc.fit_transform(X_train)
5 X_test = sc.transform(X_test)
```

▼ Model Training

```
In [12]: 1 from sklearn.linear_model import LogisticRegression
2
3 clf = LogisticRegression(random_state=0)
4 clf.fit(X_train, y_train)
```

Out[12]: LogisticRegression(random_state=0)

```
In [13]: 1 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
2
3 predictions = clf.predict(X_test)
4 print(classification_report(y_test, predictions))
5 print(accuracy_score(y_test, predictions))
6 print(confusion_matrix(y_test, predictions))
```

	precision	recall	f1-score	support
0	0.67	0.77	0.72	5620
1	0.60	0.48	0.53	3996
accuracy			0.65	9616
macro avg	0.64	0.62	0.62	9616
weighted avg	0.64	0.65	0.64	9616

0.6491264559068219
[[4343 1277]
[2097 1899]]

cv

```
In [14]: 1 from sklearn.model_selection import cross_val_score
2
3 accuracies = cross_val_score(estimator=clf, X=X_train, y=y_train, cv=10, verbose=1) # k = 10
4 print(accuracies)
5 print(f"Mean: {accuracies.mean() * 100} %")
6 print(f"Std: {accuracies.std() * 100} %")
```

[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[0.64803743 0.64231869 0.65115675 0.63520541 0.6474259 0.6450858
0.65444618 0.66016641 0.65106604 0.63806552]
Mean: 64.72974125619561 %
Std: 0.7123905317268312 %

[Parallel(n_jobs=1)]: Done 10 out of 10 | elapsed: 1.2s finished

GridSearchCV

```
In [15]: 1 from sklearn.model_selection import GridSearchCV
2
3 parameters = {"C": [0.001, 0.01, 0.1, 1, 10, 100, 1000], "penalty": ['l1', 'l2']}
4
5 grid_search = GridSearchCV(estimator=clf, param_grid=parameters,
6                             scoring="accuracy",
7                             cv=10, # cv = k = 10
8                             n_jobs=-1, # all cpu
9                             verbose=1 # print
10                            )
11 grid_search.fit(X_train, y_train)
12 print(grid_search.best_estimator_)
13 print(grid_search.best_params_)
14 print(grid_search.best_score_) # precision
15 print(grid_search.best_index_)
```

Fitting 10 folds for each of 14 candidates, totalling 140 fits

C:\Users\Perry\anaconda3\lib\site-packages\sklearn\model_selection_search.py:922: UserWarning: One or more of the test scores are non-finite: [nan 0.6434494 nan 0.64690737 nan 0.64706342
nan 0.64729741 nan 0.64721942 nan 0.64727143
nan 0.64721942]
warnings.warn(

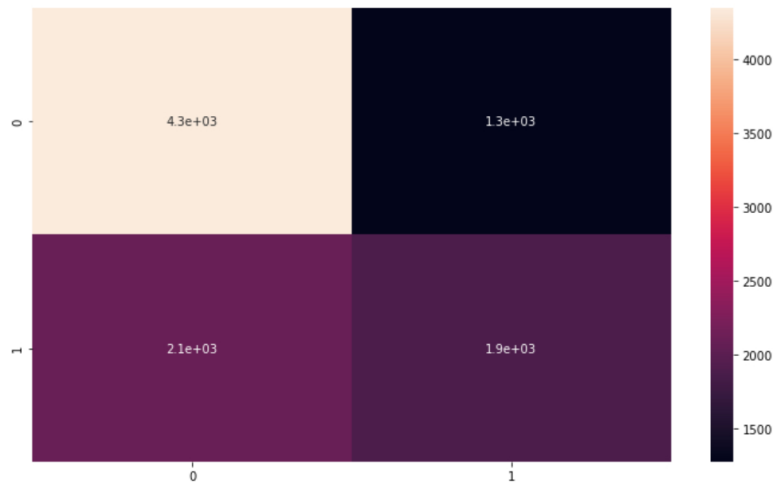
LogisticRegression(C=1, random_state=0)
{'C': 1, 'penalty': 'l2'}
0.647297412561956
7

```
In [16]: 1 grid_predictions = grid_search.predict(X_test)
2
3 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
4
5 print(classification_report(y_test, grid_predictions))
6 print(accuracy_score(y_test, grid_predictions))
7 print(confusion_matrix(y_test, grid_predictions))
```

	precision	recall	f1-score	support
0	0.67	0.77	0.72	5620
1	0.60	0.48	0.53	3996
accuracy			0.65	9616
macro avg	0.64	0.62	0.62	9616
weighted avg	0.64	0.65	0.64	9616

0.6491264559068219
[[4343 1277]
[2097 1899]]

```
In [17]: 1 plt.figure(figsize=(12, 7))
2 sns.heatmap(data=confusion_matrix(y_test, grid_predictions), annot=True);
```



Feature Selection 递归特征消除 (Recursive feature elimination) and Re-do Model Training

```
In [18]: 1 from sklearn.feature_selection import RFE # 递归特征消除 (Recursive feature elimination)
2
3 clf = LogisticRegression()
4 rfe = RFE(estimator=clf, n_features_to_select=10) # find 10 columns that are most valuable
5 rfe.fit(X_train, y_train)
```

Out[18]: RFE(estimator=LogisticRegression(), n_features_to_select=10)

```
In [19]: 1 rfe.support_
```

Out[19]: array([True, True, False, True, True, False, False, False, False,
 False, False, True, False, False, True, False, False, True,
 True, False, True, False, False, False, False, False, False,
 True, False, False, False, False, False, False, False, False,
 False, False, False, False, False])

```
In [20]: 1 rfe.ranking_
```

Out[20]: array([1, 1, 20, 1, 1, 9, 2, 13, 28, 12, 14, 1, 3, 22, 1, 10, 4,
 1, 1, 7, 6, 1, 11, 5, 21, 19, 29, 1, 8, 30, 25, 32, 23, 31,
 15, 18, 26, 17, 27, 24, 16])

```
In [21]: 1 X = df.drop(labels="churn", axis=1)
2 # using rfe.support_ as selected features
3 X = X[X.columns[rfe.support_]]
4 y = df["churn"]
5
6 from sklearn.model_selection import train_test_split
7
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
9
10 from sklearn.preprocessing import StandardScaler, MinMaxScaler
11
12 sc = StandardScaler()
13 X_train = sc.fit_transform(X_train)
14 X_test = sc.transform(X_test)
15
16 from sklearn.linear_model import LogisticRegression
17
18 clf = LogisticRegression(random_state=0)
19 clf.fit(X_train, y_train)
20
21 predictions = clf.predict(X_test)
22
23 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
24
25 predictions = clf.predict(X_test)
26 print(classification_report(y_test, predictions))
27 print(accuracy_score(y_test, predictions))
28 print(confusion_matrix(y_test, predictions))
```

	precision	recall	f1-score	support
0	0.67	0.77	0.72	5620
1	0.59	0.46	0.52	3996
accuracy			0.64	9616
macro avg	0.63	0.62	0.62	9616
weighted avg	0.64	0.64	0.63	9616

0.6424708818635607
[[4344 1276]
[2162 1834]]

```
In [22]: 1 clf.coef_
```

```
Out[22]: array([[ -0.15039154,  0.27296565, -0.71876625, -0.32931913,  0.15731766,  
                0.1145698 ,  0.10112571,  0.11232106, -0.13146983,  0.08687853]])
```