

## By hand Complex

```
In [49]: 1 import pandas as pd
2 pd.set_option("max_colwidth", None)
3
4 import pycaret
5 import numpy as np
6 import matplotlib.pyplot as plt
7 from pycaret.classification import *
8 from sklearn.model_selection import train_test_split
9 from sklearn.metrics import accuracy_score
10
11 from functions.homebrew import *
12 import numpy as np
13 import pandas as pd
14
15 from sklearn.linear_model import LogisticRegression
16 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis as LDA, QuadraticDiscriminantAnalysis as QDA
17 from sklearn.naive_bayes import GaussianNB
18 from sklearn.neighbors import KNeighborsClassifier
19 from sklearn.preprocessing import StandardScaler
20 from sklearn.model_selection import train_test_split, cross_val_score
21 from sklearn.metrics import accuracy_score
22 from tqdm import tqdm
23 from itertools import combinations
24 import pickle
25 import os
26
27 # If you're using statsmodels or ISLP for specific tasks, keep these imports
28 import statsmodels.api as sm
29 # Assuming ISLP and homebrew are custom modules specific to your project
30 from ISLP import load_data, confusion_table
31 from ISLP.models import ModelSpec as MS, summarize, contrast
32 import statsmodels.api as sm
33 from scipy import stats
```

## Helper Functions

```
In [50]: 1 def add_transformations(data, cont_cols):
2     for var in cont_cols:
3         data[f'log_{var}'] = np.log(data[var] + 1)
4         data[f'sq_{var}'] = data[var]**2
5         data[f'sqrt_{var}'] = np.sqrt(data[var])
6         data[f'inv_{var}'] = 1 / (data[var] + 1)
7         data[f'boxcox_{var}'], _ = stats.boxcox(data[var] + 1)
8         data[f'sigmoid_{var}'] = 1 / (1 + np.exp(-data[var]))
9         data[f'sin_{var}'] = np.sin(data[var])
10        data[f'cos_{var}'] = np.cos(data[var])
```

```
In [51]: 1 def convert_confusion_matrix(df, name):
2         """
3         Converts a confusion matrix dataframe into a format with columns for model name, TP, TN, FP, FN.
4
5         Args:
6         df (pd.DataFrame): Confusion matrix dataframe with multi-index (Truth, Predicted) and columns [0, 1].
7
8         Returns:
9         pd.DataFrame: Reformatted dataframe with model evaluation metrics.
10        """
11        # Extracting the values from the confusion matrix
12        tn, fp, fn, tp = df.iloc[0, 0], df.iloc[0, 1], df.iloc[1, 0], df.iloc[1, 1]
13        acc = (tp + tn) / (tp + tn + fp + fn)
14        prec = tp / (tp + fp)
15        recall = tp / (tp + fn)
16        f1 = 2 * ((prec * recall) / (prec + recall))
17        # Creating a new dataframe with the desired format
18        metrics_df = pd.DataFrame({
19            "name": name,
20            "tp": [tp],
21            "tn": [tn],
22            "fp": [fp],
23            "fn": [fn],
24            'acc': acc,
25            'prec': prec,
26            'recall': recall,
27            'f1': f1
28        })
29
30        return metrics_df
```

```
In [52]: 1 def format_results(df):
2         df = np.where(df == 1, 'Donor', 'No Donor')
3         return df
```

## LOAD DATA

```
In [53]: 1 df = pd.read_csv('./data/df.csv').drop('Unnamed: 0', axis=1)
```

```
In [54]: 1 train = df[df['type'] == 'train'].drop('type', axis=1)
2         dev = df[df['type'] == 'dev'].drop('type', axis=1)
3         test = df[df['type'] == 'test'].drop('type', axis=1)
```

## VIF

```
In [7]: 1 dummies = pd.get_dummies(df, drop_first=True)
2
3         cat_cols = [
4             'zipconvert2_Yes', 'zipconvert3_Yes', 'zipconvert4_Yes', 'boxcox_zipconvert5_Yes',
5             'homeowner_Yes', 'female_Yes', 'type_train', 'type_dev', 'type_test'
6         ]
7
8         cont_cols = [col for col in dummies.columns if col not in cat_cols + ['target']]
9         add_transformations(dummies, cont_cols)
10
11        kept, removed = remove_high_vif_features(X=dummies.drop('target_No Donor', axis=1), y=dummies['target_No Donor'], vif_threshold=5)
12        print('REMOVED:', removed)
```

REMOVED: ['num\_child', 'income', 'cos\_target\_No Donor', 'sin\_target\_No Donor', 'sigmoid\_target\_No Donor', 'boxcox\_target\_No Donor', 'inv\_target\_No Donor', 'sqrt\_target\_No Donor', 'sq\_target\_No Donor', 'cos\_zipconvert5\_Yes', 'sin\_zipconvert5\_Yes', 'sigmoid\_zipconvert5\_Yes', 'boxcox\_zipconvert5\_Yes', 'inv\_zipconvert5\_Yes', 'sqrt\_zipconvert5\_Yes', 'sq\_zipconvert5\_Yes', 'log\_zipconvert5\_Yes', 'log\_income', 'sq\_income', 'log\_num\_child', 'sq\_num\_child', 'sigmoid\_avg\_fam\_inc', 'inv\_num\_child', 'sqrt\_months\_since\_donate', 'sqrt\_med\_fam\_inc', 'sq\_wealth', 'boxcox\_avg\_fam\_inc', 'boxcox\_num\_prom', 'log\_last\_gift', 'months\_since\_donate', 'log\_avg\_gift', 'inv\_home\_value', 'inv\_avg\_fam\_inc', 'boxcox\_time\_lag', 'log\_wealth', 'inv\_med\_fam\_inc', 'log\_largest\_gift', 'boxcox\_home\_value', 'sqrt\_num\_prom', 'log\_lifetime\_gifts', 'sqrt\_income', 'boxcox\_pct\_lt15k', 'sqrt\_wealth', 'sq\_months\_since\_donate', 'sqrt\_avg\_fam\_inc', 'sqrt\_avg\_gift', 'boxcox\_med\_fam\_inc', 'sqrt\_time\_lag', 'boxcox\_largest\_gift', 'sqrt\_pct\_lt15k', 'inv\_income', 'sqrt\_home\_value', 'sigmoid\_num\_child', 'log\_num\_prom', 'sqrt\_last\_gift', 'log\_months\_since\_donate', 'boxcox\_avg\_gift', 'sqrt\_lifetime\_gifts', 'avg\_fam\_inc', 'wealth', 'largest\_gift', 'med\_fam\_inc', 'boxcox\_last\_gift', 'inv\_pct\_lt15k', 'zipconvert5\_Yes', 'num\_prom', 'pct\_lt15k', 'home\_value', 'log\_avg\_fam\_inc', 'inv\_wealth', 'log\_time\_lag', 'boxcox\_lifetime\_gifts', 'sqrt\_num\_child', 'avg\_gift', 'inv\_time\_lag', 'log\_med\_fam\_inc', 'boxcox\_income', 'inv\_largest\_gift', 'inv\_last\_gift', 'sq\_avg\_fam\_inc', 'inv\_num\_prom', 'last\_gift']

```
In [55]: 1 final_vars = list(kept.corr().drop('target')[np.abs(kept.corr()['target']).drop('target') > .05].index)
```

```
In [56]: 1 regress = kept[final_vars]
```

```
In [57]: 1 regress['target'] = (df['target'] == 'Donor').astype(int)
```

```
In [72]: 1 kept = kept.drop('log_target_No Donor', axis =1)
```

```
In [59]: 1 train = kept[regress['type_train'] ==1]
2 dev = kept[(regress['type_test'] == 0) & (regress['type_train'] == 0)]
3 test = kept[regress['type_test'] ==1]
```

```
In [60]: 1 for data in [train, dev, test]:
2     data.drop('type_train', inplace = True, axis = 1)
3     data.drop('type_test', inplace = True, axis = 1)
4 test = test.drop('target', axis =1)
```

## Logistic Regression

```
In [61]: 1 results_df = pd.DataFrame()
```

```
In [62]: 1 # Selecting features and target variable for training data
2 X_train = train.drop(['target'], axis =1 )
3 y_train = train['target']
4 X_test = dev.drop(['target'], axis = 1)
5 y_test = dev['target']
6
7 # Fitting logistic regression model
8 glm = sm.GLM(y_train, X_train, family=sm.families.Binomial())
9 glm = glm.fit()
10
11 # Summarizing results
12 # print(results.summary())
```

```
In [63]: 1 log_preds = (glm.predict(X_test) >= 0.5).astype(int)
2 log_acc = accuracy_score(log_preds, y_test)
3 print(log_acc)
4
5 d = confusion_table(log_preds,y_test)
6 results_df = pd.concat([results_df,convert_confusion_matrix(d, 'Logistic Regression')])
7
8 log_test_preds = (glm.predict(test) >= 0.5).astype(int)
9 log_test_preds = format_results(log_test_preds)
10
11 save_df = pd.DataFrame(log_test_preds, columns=['values'])
12 save_df.to_csv('./preds/log.csv', index=False)
```

0.5116666666666667

## LDA

```
In [64]: 1 lda = LDA(store_covariance=True)
2 lda.fit(X_train, y_train)
3
4 lda_preds = lda.predict(X_test)
5
6 lda_acc = accuracy_score(lda_preds,y_test)
7 print(lda_acc)
8
9 d = confusion_table(lda_preds,y_test)
10 results_df = pd.concat([results_df,convert_confusion_matrix(d, 'LDA')])
11
12
13 lda_test_preds = (lda.predict(test) >= 0.5).astype(int)
14 lda_test_preds = format_results(lda_test_preds)
15
16 save_df = pd.DataFrame(lda_test_preds, columns=['values'])
17 save_df.to_csv('./preds/lda.csv', index=False)
```

0.5066666666666667

## QDA.

```
In [65]: 1 qda = QDA(store_covariance=True)
2 qda.fit(X_train, y_train)
3
4 qda_preds = qda.predict(X_test)
5
6 qda_acc = accuracy_score(qda_preds,y_test)
7
8 print(qda_acc)
9
10 d = confusion_table(qda_preds,y_test)
11 results_df = pd.concat([results_df,convert_confusion_matrix(d, 'QDA')])
12
13 qda_test_preds = (qda.predict(test) >= 0.5).astype(int)
14 qda_test_preds = format_results(qda_test_preds)
15
16 save_df = pd.DataFrame(qda_test_preds, columns=['values'])
17 save_df.to_csv('./preds/qda.csv', index=False)

0.49166666666666664
```

KNN

```
In [66]: 1 knn1 = KNeighborsClassifier(n_neighbors=1)
2 knn1.fit(X_train, y_train)
3 knn1_pred = knn1.predict(X_test)
4 knn1_acc = accuracy_score(knn1_pred,y_test)
5
6 print(knn1_acc)
7
8 d = confusion_table(knn1_pred, y_test)
9 results_df = pd.concat([results_df,convert_confusion_matrix(d, 'KNN')])
10
11 knn1_test_preds = (knn1.predict(test) >= 0.5).astype(int)
12 knn1_test_preds = format_results(knn1_test_preds)
13
14 save_df = pd.DataFrame(knn1_test_preds, columns=['values'])
15 save_df.to_csv('./preds/knn1.csv', index=False)

0.48166666666666667
```

NB

```
In [67]: 1 nb = GaussianNB()
2 nb.fit(X_train, y_train)
3 nb_preds = nb.predict(X_test)
4 nb_acc = accuracy_score(nb_preds,y_test)
5
6 print(nb_acc)
7 save_df = pd.DataFrame(nb_preds, columns=['values'])
8 save_df.to_csv('./preds/nb.csv', index=False)
9
10 d = confusion_table(nb_preds, y_test)
11 results_df = pd.concat([results_df,convert_confusion_matrix(d, 'Naïve Bayes')])
12
13 nb_test_preds = (nb.predict(test) >= 0.5).astype(int)
14 nb_test_preds = format_results(nb_test_preds)
15
16 save_df = pd.DataFrame(nb_test_preds, columns=['values'])
17 save_df.to_csv('./preds/nb.csv', index=False)

0.495
```

```
In [68]: 1 results_df
```

Out[68]:

	name	tp	tn	fp	fn	acc	prec	recall	f1
0	Logistic Regression	159	148	153	140	0.511667	0.509615	0.531773	0.520458
0	LDA	156	148	156	140	0.506667	0.500000	0.527027	0.513158
0	QDA	26	269	286	19	0.491667	0.083333	0.577778	0.145658
0	KNN	140	149	172	139	0.481667	0.448718	0.501792	0.473773
0	Naïve Bayes	18	279	294	9	0.495000	0.057692	0.666667	0.106195

```
In [69]: 1 # test_acc = {
2 #       'log': 0.5333333,
3 #       'lda': 0.5583333,
4 #       'qda': 0.525,
5 #       'knn': 0.475,
6 #       'nb': 0.5166667,
7 # }
```

```
In [70]: 1 # results_df['test_acc'] = test_acc.values()
```

```
In [71]: 1 results_df
```

Out[71]:

	name	tp	tn	fp	fn	acc	prec	recall	f1	test_acc
0	Logistic Regression	159	148	153	140	0.511667	0.509615	0.531773	0.520458	0.533333
0	LDA	156	148	156	140	0.506667	0.500000	0.527027	0.513158	0.558333
0	QDA	26	269	286	19	0.491667	0.083333	0.577778	0.145658	0.525000
0	KNN	140	149	172	139	0.481667	0.448718	0.501792	0.473773	0.475000
0	Naïve Bayes	18	279	294	9	0.495000	0.057692	0.666667	0.106195	0.516667