

Code

main.py

```
1
2 import warnings
3 warnings.filterwarnings("ignore")
4
5 import numpy as np
6 import seaborn as sns
7 import matplotlib.pyplot as plt
8 import pandas as pd
9 from tqdm import tqdm
10 import os
11 import math
12
13 from utils import *
14
15 np.random.seed(42)
16
17
18 # load data
19 train_df = pd.read_csv("./dataset/algerian_fires_train.csv")
20 test_df = pd.read_csv("./dataset/algerian_fires_test.csv")
21 test_y = test_df.loc[:, 'Classes'].values
22 train_idx = np.arange(train_df.shape[0])
23 test_idx = np.arange(train_df.shape[0], train_df.shape[0] +
24 test_df.shape[0])
25 df = pd.concat([train_df, test_df], axis=0).reset_index().drop(columns=
26 ['index'])
27 df.loc[test_idx, 'Classes'] = -1
28 df.tail(3)
29 df.info()
30
31
32 # ## **Data Process**
33
34 # feature normalization
35 Date_col = 'Date'
36 tar_col = 'Classes'
37 feat_cols = [col for col in train_df.columns if col != Date_col and col !=
38 tar_col]
39
40 df.loc[:, feat_cols] = (df.loc[:, feat_cols] - df.loc[:,
41 feat_cols].mean(axis=0)) / (df.loc[:, feat_cols].std(axis=0) + 1e-8)
42 df.head(3)
43
44 # data visualization
45 train_df = df.loc[train_idx, feat_cols+[tar_col]]
46
47 for col in feat_cols:
48     _ = sns.displot(train_df, x=col, hue=tar_col, kde=True)
```

```

47     plt.savefig("./figs/hist_{}.png".format(col), dpi=200)
48     plt.show()
49
50
51     # features correlation
52     corr = df.loc[:, feat_cols+[tar_col]].corr()
53     _ = plt.figure(figsize=(10, 8))
54     _ = plt.imshow(corr)
55     _ = plt.colorbar()
56     _ = plt.xticks(np.arange(len(feat_cols+[tar_col])), feat_cols+[tar_col],
57                     rotation=45)
58     _ = plt.yticks(np.arange(len(feat_cols+[tar_col])), feat_cols+[tar_col])
59     _ = plt.savefig("./figs/corr.png", dpi=200)
60     plt.show()
61
62     # time series feature consturction
63     df, time_cols = add_time_features(df, feat_cols)
64     df.head()
65
66
67     # ## **Logistic Regression**
68
69
70     from sklearn.linear_model import LogisticRegression
71
72     # train on the original dataset
73     log_clf = LogisticRegression()
74     params = {'penalty': ['l1', 'l2', 'elasticnet'], 'C': [0.001, 0.01, 0.1, 1,
75                 10], 'solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']}
76     log_clfs, log_res_df = finetune(df.loc[train_idx, feat_cols+[tar_col]],
77                                     feat_cols, tar_col, log_clf, params)
78     log_res_df.to_csv("./res/ori_log_res.csv", index=False)
79
80     # predict on the test dataset
81     best_log_clf = log_clfs.best_estimator_
82     best_log_clf.fit(df.loc[train_idx, feat_cols].values, df.loc[train_idx,
83                         tar_col].values)
84     log_f1, log_acc, log_cnf = get_clf_metrics(df.loc[test_idx, feat_cols+
85                                                 [tar_col]], feat_cols, test_y, best_log_clf)
86
87     print("Best LOGClassifier: ", log_clfs.best_estimator_)
88     print("F1-score of LogisticRegression is {:.3f}".format(log_f1))
89     print("Accuracy of LogisticRegression is {:.3f}".format(log_acc))
90     print("The confusion matrix of LogisticRegression:")
91     plot_cnf(log_cnf, name='ori_log')
92
93
94     # train on the augmented dataset
95     log_clf = LogisticRegression()
96     params = {'penalty': ['l1', 'l2', 'elasticnet'], 'C': [0.001, 0.01, 0.1, 1,
97                 10], 'solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']}
98     log_clfs, log_res_df = finetune(df.loc[train_idx, feat_cols+time_cols+
99                                     [tar_col]], feat_cols+time_cols, tar_col, log_clf, params)
100     log_res_df.to_csv("./res/aug_log_res.csv", index=False)
101
102     # predict on the test dataset
103     best_log_clf = log_clfs.best_estimator_

```

```

98 best_log_clf.fit(df.loc[train_idx, feat_cols+time_cols].values,
99 df.loc[train_idx, tar_col].values)
100
101 log_f1, log_acc, log_cnf = get_clf_metrics(df.loc[test_idx,
102 feat_cols+time_cols+[tar_col]], feat_cols+time_cols, test_y, best_log_clf)
103
104 print("Best LOGClassifier: ", log_clfs.best_estimator_)
105 print("F1-score of LogisticRegression is {:.3f}".format(log_f1))
106 print("Accuracy of LogisticRegression is {:.3f}".format(log_acc))
107 print("The confusion matrix of LogisticRegression:")
108 plot_cnf(log_cnf, name='aug_log')
109
110
111 # ## **Support Vector Machine**
112
113 from sklearn.svm import SVC
114
115 # train on the original dataset
116 svm_clf = SVC()
117 params = {'kernel':['linear', 'poly', 'rbf', 'sigmoid'], 'C': [0.001, 0.01,
118 0.1, 1, 10], 'gamma': ['scale', 'auto']}
119 svm_clfs, svm_res_df = finetune(df.loc[train_idx, feat_cols+[tar_col]],
120 feat_cols, tar_col, svm_clf, params)
121 svm_res_df.to_csv("./res/ori_svm_res.csv", index=False)
122
123 # predict on the test dataset
124 best_svm_clf = svm_clfs.best_estimator_
125 best_svm_clf.fit(df.loc[train_idx, feat_cols].values, df.loc[train_idx,
126 tar_col].values)
127 svm_f1, svm_acc, svm_cnf = get_clf_metrics(df.loc[test_idx, feat_cols+
128 [tar_col]], feat_cols, test_y, best_svm_clf)
129
130
131 print("Best SVMClassifier: ", svm_clfs.best_estimator_)
132 print("F1-score of SupportVectorMachine is {:.3f}".format(svm_f1))
133 print("Accuracy of SupportVectorMachine is {:.3f}".format(svm_acc))
134 print("The confusion matrix of SupportVectorMachine:")
135 plot_cnf(svm_cnf, name='ori_svm')
136
137
138 # train on the augmented dataset
139 svm_clf = SVC()
140 params = {'kernel':['linear', 'poly', 'rbf', 'sigmoid'], 'C': [0.001, 0.01,
141 0.1, 1, 10], 'gamma': ['scale', 'auto']}
142 svm_clfs, svm_res_df = finetune(df.loc[train_idx, feat_cols+time_cols+
143 [tar_col]], feat_cols+time_cols, tar_col, svm_clf, params)
144 svm_res_df.to_csv("./res/aug_svm_res.csv", index=False)
145
146 # predict on the test dataset
147 best_svm_clf = svm_clfs.best_estimator_
148 best_svm_clf.fit(df.loc[train_idx, feat_cols+time_cols].values,
149 df.loc[train_idx, tar_col].values)
150 svm_f1, svm_acc, svm_cnf = get_clf_metrics(df.loc[test_idx,
151 feat_cols+time_cols+[tar_col]], feat_cols+time_cols, test_y, best_svm_clf)
152
153
154 print("Best SVMClassifier: ", svm_clfs.best_estimator_)
155 print("F1-score of SupportVectorMachine is {:.3f}".format(svm_f1))
156 print("Accuracy of SupportVectorMachine is {:.3f}".format(svm_acc))
157 print("The confusion matrix of SupportVectorMachine:")

```

```

146 plot_cnf(svm_cnf, name='aug_svm')
147
148
149 # ## **Neural Networks**
150
151
152 from sklearn.neural_network import MLPClassifier
153
154 # train on the original dataset
155 mlp_clf = MLPClassifier()
156 params = {'activation':['tanh', 'relu'], 'solver':['lbfgs', 'sgd', 'adam'],
157           'learning_rate':['constant', 'invscaling', 'adaptive'], 'alpha':[0.001,
158           0.01, 0.1, 0.5]}
159
160 mlp_clfs, mlp_res_df = finetune(df.loc[train_idx, feat_cols+[tar_col]],
161                                feat_cols, tar_col, mlp_clf, params)
162 mlp_res_df.to_csv("./res/ori_mlp_res.csv", index=False)
163
164 # predict on the test dataset
165 best_mlp_clf = mlp_clfs.best_estimator_
166 best_mlp_clf.fit(df.loc[train_idx, feat_cols].values, df.loc[train_idx,
167 tar_col].values)
168 mlp_f1, mlp_acc, mlp_cnf = get_clf_metrics(df.loc[test_idx, feat_cols+
169 [tar_col]], feat_cols, test_y, best_mlp_clf)
170
171
172 print("Best MLPClassifier: ", mlp_clfs.best_estimator_)
173 print("F1-score of NeuralNetwork is {:.3f}".format(mlp_f1))
174 print("Accuracy of NeuralNetwork is {:.3f}".format(mlp_acc))
175 print("The confusion matrix of NeuralNetwork:")
176 plot_cnf(mlp_cnf, name='ori_mlp')
177
178
179 # train on the augmented dataset
180 mlp_clf = MLPClassifier()
181 params = {'activation':['tanh', 'relu'], 'solver':['lbfgs', 'sgd', 'adam'],
182           'learning_rate':['constant', 'invscaling', 'adaptive'], 'alpha':[0.001,
183           0.01, 0.1, 0.5]}
184
185 mlp_clfs, mlp_res_df = finetune(df.loc[train_idx, feat_cols+time_cols+
186 [tar_col]], feat_cols+time_cols, tar_col, mlp_clf, params)
187 mlp_res_df.to_csv("./res/aug_mlp_res.csv", index=False)
188
189 # predict on the test dataset
190 best_mlp_clf = mlp_clfs.best_estimator_
191 best_mlp_clf.fit(df.loc[train_idx, feat_cols+time_cols].values,
192 df.loc[train_idx, tar_col].values)
193 mlp_f1, mlp_acc, mlp_cnf = get_clf_metrics(df.loc[test_idx,
194 feat_cols+time_cols+[tar_col]], feat_cols+time_cols, test_y, best_mlp_clf)
195
196
197 print("Best MLPClassifier: ", mlp_clfs.best_estimator_)
198 print("F1-score of NeuralNetwork is {:.3f}".format(mlp_f1))
199 print("Accuracy of NeuralNetwork is {:.3f}".format(mlp_acc))
200 print("The confusion matrix of NeuralNetwork:")
201 plot_cnf(mlp_cnf, name='aug_mlp')
202
203
204 # ## **Trivial System**
205
206
207 from model import TrivialClassifier

```

```

194
195 # train on the original dataset
196 trv_clf = TrivialClassifier()
197 mean_val, std_val = cross_val(df.loc[train_idx, feat_cols+[tar_col]],
198                               feat_cols, tar_col, trv_clf)
199 print("Mean accuracy on validation dataset: {:.3f}".format(mean_val))
200 print("Std accuracy on validation dataset: {:.3f}".format(std_val))
201
202 # predict on the test dataset
203
204 trv_acc = trv_clf.score(df.loc[test_idx, feat_cols].values, test_y)
205 y_pred = trv_clf.predict(df.loc[test_idx, feat_cols].values)
206 trv_f1 = f1_score(test_y, y_pred, average='macro')
207 trv_cnf = confusion_matrix(test_y, y_pred)
208 print("F1-score of TrivialClassifier is {:.3f}".format(trv_f1))
209 print("Accuracy of TrivialClassifier is {:.3f}".format(trv_acc))
210 print("The confusion matrix of TrivialClassifier:")
211 plot_cnf(trv_cnf, name='ori_trv')
212
213 # ## **Baseline: Nearest Mean Classifier**
214
215
216 from model import KernelNearestMeansClassifier
217
218 # train on the original dataset
219 knm_clf = KernelNearestMeansClassifier()
220 mean_val, std_val = cross_val(df.loc[train_idx, feat_cols+[tar_col]],
221                               feat_cols, tar_col, knm_clf)
222 print("Mean accuracy on validation dataset: {:.3f}".format(mean_val))
223 print("Std accuracy on validation dataset: {:.3f}".format(std_val))
224
225 # predict on the test dataset
226
227 knm_acc = knm_clf.score(df.loc[test_idx, feat_cols].values, test_y)
228 y_pred = knm_clf.predict(df.loc[test_idx, feat_cols].values)
229 knm_f1 = f1_score(test_y, y_pred, average='macro')
230 knm_cnf = confusion_matrix(test_y, y_pred)
231 print("F1-score of KernelNearestMeansClassifier is {:.3f}".format(knm_f1))
232 print("Accuracy of KernelNearestMeansClassifier is {:.3f}".format(knm_acc))
233 print("The confusion matrix of KernelNearestMeansClassifier:")
234 plot_cnf(knm_cnf, name='ori_knm')
235
236 # train on the augmented dataset
237 knm_clf = KernelNearestMeansClassifier()
238 mean_val, std_val = cross_val(df.loc[train_idx, feat_cols+time_cols+
239                               [tar_col]], feat_cols+time_cols, tar_col, knm_clf)
240 print("Mean accuracy on validation dataset: {:.3f}".format(mean_val))
241 print("Std accuracy on validation dataset: {:.3f}".format(std_val))
242
243 # predict on the test dataset
244
245 knm_acc = knm_clf.score(df.loc[test_idx, feat_cols+time_cols].values,
246                          test_y)
247 y_pred = knm_clf.predict(df.loc[test_idx, feat_cols+time_cols].values)
248 knm_f1 = f1_score(test_y, y_pred, average='macro')
249 knm_cnf = confusion_matrix(test_y, y_pred)

```

```

248 print("F1-score of KernelNearestMeansClassifier is {:.3f}".format(knm_f1))
249 print("Accuracy of KernelNearestMeansClassifier is {:.3f}".format(knm_acc))
250 print("The confusion matrix of KernelNearestMeansClassifier:")
251 plot_cnf(knm_cnf, name='aug_knm')
252
253
254 # ## **Feature Importance**
255
256
257 importance = best_log_clf.coef_[0]
258 plt.figure(figsize=(12, 6))
259 _ = plt.bar([x for x in range(len(importance))], importance, color='red')
260 _ = plt.xticks(ticks=list(range(len(importance))),
261               labels=feat_cols+time_cols, rotation=80)
261 _ = plt.xlabel("features")
262 _ = plt.ylabel("importance")
263 _ = plt.tight_layout()
264 _ = plt.savefig("./figs/imp_log.png", dpi=300)
265 _ = plt.show()
266

```

model.py

```

1  import numpy as np
2
3
4  class TrivialClassifier(object):
5      def __init__(self, random_state=42):
6          np.random.seed(random_state)
7
8      def fit(self, X, y):
9          classes = np.unique(y)
10         self.classes = classes
11         prob = np.zeros(len(classes))
12         for i, c in enumerate(classes):
13             prob[i] = np.sum(np.where(y==c, 1, 0))
14         self.prob = prob / np.sum(prob)
15
16     def predict(self, X):
17         y_pred = np.random.choice(self.classes, size=X.shape[0],
18                                   p=self.prob)
19         return y_pred
20
21     def score(self, X, y):
22         y_pred = self.predict(X)
23         return np.mean(y_pred.ravel() == y.ravel())
24
25 class KernelNearestMeansClassifier(object):
26     def __init__(self, gamma=0.01, kernel_type='linear'):
27         self.gamma = gamma
28         self.kernel_type = kernel_type
29
30     def fit(self, X, y):
31         self.X_train = X
32         self.y_train = y

```

```

32
33
34     def predict(self, X):
35         X1 = self.X_train[self.y_train==0]
36         X2 = self.X_train[self.y_train==1]
37         if self.kernel_type == 'rbf':
38             dist1 = np.zeros((X1.shape[0], X1.shape[0]))
39             for i in range(X1.shape[0]):
40                 dist1[i, :] = np.sum((X1[i] - X1)**2, axis=1)
41             self.K1 = np.mean(np.exp(-self.gamma * dist1))
42
43             dist2 = np.zeros((X2.shape[0], X2.shape[0]))
44             for i in range(X2.shape[0]):
45                 dist2[i, :] = np.sum((X2[i] - X2)**2, axis=1)
46             self.K2 = np.mean(np.exp(-self.gamma * dist2))
47
48             res = np.zeros((X.shape[0], 2))
49             for i in range(X.shape[0]):
50                 res[i, 0] = -self.K1 + 2 * np.mean(np.exp(-self.gamma *
np.sum((X[i] - X1)**2, axis=1)))
51                 res[i, 1] = -self.K2 + 2 * np.mean(np.exp(-self.gamma *
np.sum((X[i] - X2)**2, axis=1)))
52             elif self.kernel_type == "linear":
53                 self.K1 = np.mean(X1 @ X1.T)
54                 self.K2 = np.mean(X2 @ X2.T)
55
56                 res = np.zeros((X.shape[0], 2))
57                 for i in range(X.shape[0]):
58                     res[i, 0] = -self.K1 + 2 * np.mean(X[i] @ X1.T)
59                     res[i, 1] = -self.K2 + 2 * np.mean(X[i] @ X2.T)
60             else:
61                 raise NotImplementedError
62             y_pred = np.argmax(res, axis=1)
63             return y_pred
64
65     def score(self, X, y):
66         y_pred = self.predict(X)
67         acc = np.mean(y_pred==y)
68         return acc

```

utils.py

```

1  import numpy as np
2  import pandas as pd
3  from sklearn.model_selection import GridSearchCV, StratifiedKFold
4  from sklearn.metrics import f1_score, confusion_matrix
5  import matplotlib.pyplot as plt
6
7  def add_time_features(df, cols, duration=5):
8      '''
9      Add time series features to the dataframe
10     '''
11     aug_cols=[]
12     for col in cols:
13         max_col = '{}_{}max'.format(col, duration)

```

```

14     min_col = '{}_{}_min'.format(col, duration)
15     mean_col = '{}_{}_mean'.format(col, duration)
16     aug_cols.append(max_col)
17     aug_cols.append(min_col)
18     aug_cols.append(mean_col)
19     df[max_col] = 0
20     df[min_col] = 0
21     df[mean_col] = 0
22     for i in range(duration, df.shape[0]):
23         df.loc[i, max_col] = df.loc[i-duration:i, col].max()
24         df.loc[i, min_col] = df.loc[i-duration:i, col].min()
25         df.loc[i, mean_col] = df.loc[i-duration:i, col].mean()
26     for col in aug_cols:
27         df.loc[df[col]==0, col] = df[col].mean()
28     return df, aug_cols
29
30 def cross_val(train_df, cols, tar_col, classifier, seed=42):
31     '''
32     Cross validation on the training data
33     returns:
34         accuracy mean and std
35     '''
36     X, y = train_df[cols].values, train_df[tar_col].values
37     skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=seed)
38     acc = []
39     for tr_idx, val_idx in skf.split(X, y):
40         classifier.fit(X[tr_idx], y[tr_idx])
41         acc.append(classifier.score(X[val_idx], y[val_idx]))
42     return np.mean(acc), np.std(acc)
43
44 def finetune(train_df, cols, tar_col, classifier, params, seed=42):
45     '''
46     Model selection on the training data
47     returns:
48         the best model and validation results
49     '''
50     X, y = train_df[cols].values, train_df[tar_col].values
51     clfs = GridSearchCV(classifier, params)
52     clfs.fit(X, y)
53     # return results dataframe
54     results = pd.DataFrame()
55     idx = np.where(np.isnan(clfs.cv_results_['std_test_score'])==False)[0]
56     results['params'] = np.array(list(map(str, clfs.cv_results_['params'])))
57     [idx]
58     results['mean_val_score'] = clfs.cv_results_['mean_test_score'][idx]
59     results['std_val_score'] = clfs.cv_results_['std_test_score'][idx]
60     return clfs, results
61
62 def get_clf_metrics(test_df, cols, test_y, classifier):
63     '''
64     Performance of classifier on the test dataset
65     returns:
66         f1-score, accuracy, confusion_matrix
67     '''
68     X, y = test_df[cols].values, test_y
69     accuracy = classifier.score(X, y)
70     y_pred = classifier.predict(X)
71     f1 = f1_score(y, y_pred, average='macro')

```



```
71     conf_matrix = confusion_matrix(y, y_pred)
72     return f1, accuracy, conf_matrix
73
74 def plot_cnf(cnf, name):
75     plt.imshow(cnf)
76     for i in range(2):
77         for j in range(2):
78             plt.text(i, j, str(cnf[i, j]), color='red')
79     plt.xticks(np.arange(2), ['Positive', 'Negative'])
80     plt.yticks(np.arange(2), ['True', 'False'])
81     plt.savefig('./figs/cnf_{}.png'.format(name), dpi=200)
82     plt.show()
83
```