

```

1  # -*- coding: utf-8 -*-
2  """Normalization_Feature_Testing
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/103D1WTRS\_XIebFMiTu3iE6eSALqeQn\_S
8
9  <center><h1>Mini Project 1 - Logistic Regression</h1>
10 <h4>This is a testing file aiming to find the effect of normalization,
    increasing feacures on the model.</h4></center>
11
12 <h3>Team Members:</h3>
13 <center>
14 Yi Zhu, 260716006<br>
15 Fei Peng, 260712440<br>
16 Yukai Zhang, 260710915
17 </center>
18 """
19
20 from google.colab import drive
21 drive.mount('/content/drive')
22
23 import numpy as np
24 import pandas as pd
25 import matplotlib.pyplot as plt
26 import seaborn as sns
27 import time
28 from google.colab import files
29
30 class LogisticRegression:
31     '''
32         This is the logistic regression class, containing fit, perdict and
33         accu_eval functions,
34         as well as many other useful functions.
35     '''
36     def __init__(self, data, folds, lr=0.01, max_iter=10000, beta=0.99, epsilon
37 =5e-3):
38         self.data = data
39         self.folds = folds
40         self.lr = lr
41         self.max_iter = max_iter
42         self.beta = beta
43         self.epsilon = epsilon
44     def shuffle_data(self):
45         '''
46             This function randomly shuffles the input dataset.
47         '''
48         # Load data from data file.
49         self.data.insert(0, column='Bias', value=1)
50         self.data = self.data.sample(frac=1)
51     def partition(self, fold):
52         '''
53             This function divides the dataset into training and validation set.
54
55             fold - the current fold
56         '''
57         data = self.data
58         # to exclude last term in previous partition for training data
59         train_add = 1 if fold < self.folds else 0
60

```

```

61     # to exclude last term in previous partition for testing data
62     test_add = 1 if fold > 0 else 0
63
64     # number of data sets
65     n = len(self.data)
66
67     train_set_1 = data.iloc[0:int((fold)/self.folds*n), :]
68     train_set_2 = data.iloc[int((fold+1)/self.folds*n)+train_add:n, :]
69     train_set = pd.concat([train_set_1, train_set_2])
70
71     test_set = data.iloc[int((fold)/self.folds*n+test_add):int((fold+1)/
self.folds*n), :]
72
73     train_X = train_set.iloc[:, :-1].values
74     train_y = train_set.iloc[:, -1].values
75     train_y = np.reshape(train_y, (-1,1))
76
77     test_X = test_set.iloc[:, :-1].values
78     test_y = test_set.iloc[:, -1].values
79     test_y = np.reshape(test_y, (-1,1))
80
81     return train_X, train_y, test_X, test_y
82
83     def normalization(self, X, v_X):
84         '''
85             This function performs the z-score normalization
86
87             X - training data
88             v_X - validation data
89         '''
90         mean = np.mean(X[:,1:], axis = 0)
91         sigma = np.std(X[:,1:], axis = 0)
92         mean = np.reshape(mean, (1,-1))
93         sigma = np.reshape(sigma, (1,-1))
94         X[:,1:] = (X[:,1:] - mean) / sigma
95         v_X[:,1:] = (v_X[:,1:] - mean) / sigma
96         return X, v_X
97
98     def fit(self, X, y, v_X, v_y, normalize=False):
99         '''
100             This function takes the training data X and its corresponding
labels vector y
101             as well as other hyperparameters (such as learning rate) as input,
102             and execute the model training through modifying the model
parameters (i.e. W).
103
104             X - training data
105             y - class of training data
106             v_X - validation data
107             v_y - class of validation data
108             epsilon - the threshold value for gradient descent
109             normalize - whether to perform normalization
110         '''
111         gradient_values, t_acc_val, v_acc_val = [], [], []
112
113         if normalize:
114             X, v_X = self.normalization(X, v_X)
115
116         # Retrive the learning rate, maximum iteration, momentum (beta)
117         lr, max_iter, beta, epsilon = self.lr, self.max_iter, self.beta, self.
epsilon
118
119         # initial weight vector

```

```

120     w = np.zeros((len(X[0]), 1))
121     # record the best weight vector
122     best_w = w
123     # iteration number, validation accuracy, last validation accuracy
124     # the step to take in gradient descent, maximum validation accuracy
125     iteration, v_acc, step, v_acc_max = 0, 0, 0, 0
126
127     dw = np.inf
128     # if the gradient delta w is smaller than threshold or achieved
    maximum iteration, stop
129     while (np.linalg.norm(dw) > epsilon and iteration <= max_iter):
130         dw = self.gradient(X, y, w)
131         gradient_values.append(np.linalg.norm(dw))
132         # if beta = 0, it will be the same as general gradient descent
133         step = beta * step + (1 - beta) * dw # gradient descent with
    momentum
134         w = w - lr * step
135
136         # predict once every 10 iterations
137         if iteration % 10 == 0:
138             t_y_pred = self.predict(X, w)
139             t_acc = self.accu_eval(t_y_pred, y)
140             v_y_pred = self.predict(v_X, w)
141             v_acc = self.accu_eval(v_y_pred, v_y)
142
143         # record the next best value
144         if v_acc >= v_acc_max:
145             v_acc_max = v_acc
146             best_w = w
147             self.marker = iteration # move the iteration marker
148
149         t_acc_val.append(t_acc)
150         v_acc_val.append(v_acc)
151
152         iteration = iteration + 1
153
154         # Uncomment to display the confusion matrix
155         # self.confusion_matrix(v_y_pred, v_y)
156         return gradient_values, t_acc_val, v_acc_val, best_w
157
158     def predict(self, X, w):
159         '''
160         This function takes a set of data as input and outputs predicted
    labels for the input points.
161         '''
162         result = self.log_func(np.dot(X, w))
163         # the prediction result converted to binary
164         predict_bin = []
165         for i in result:
166             if i >= 0.5:
167                 predict_bin.append(1)
168             else:
169                 predict_bin.append(0)
170         return predict_bin
171
172     def accu_eval(self, y_pred, y):
173         '''
174         This function evaluates the models' accuracy.
175         '''
176         count = 0
177         for i in range(len(y_pred)):
178             if y_pred[i] == y[i]:
179                 count = count + 1

```

```

180         # return the accuracy ratio: #corret prediction / #data points
181         return count / len(y)
182
183     def log_func(self, alpha):
184         return 1 / (1 + np.exp(-alpha))
185
186     def gradient(self, X, y, w):
187         N = len(X[0])
188         y_hat = self.log_func(np.dot(X, w))
189         delta = np.dot(X.T, y_hat - y) / N
190         return delta
191
192     def confusion_matrix(self, y_pred, y):
193         y = y.reshape(-1)
194         data = {'Actual_y': y,
195                'Predicted_y': y_pred}
196         df = pd.DataFrame(data, columns=['Actual_y', 'Predicted_y'])
197         confusion_matrix = pd.crosstab(df['Actual_y'], df['Predicted_y'],
rownames=['Actual'], colnames=['Predicted'])
198         svm = sns.heatmap(confusion_matrix, annot=True, cmap="YlGnBu")
199
200         # Uncomment this part to download the confusion matrix.
201         # temp_time = time.time()
202         # figure = svm.get_figure()
203         # figure.savefig("Confusion_Matrix-{}.png".format(temp_time), dpi =
1200)
204         # files.download("Confusion_Matrix-{}.png".format(temp_time))
205
206 class KFoldValidation:
207     def __init__(self, folds, path, lr, max_iter, epsilon, beta):
208         self.folds = folds
209         self.path = path
210         self.lr = lr
211         self.max_iter = max_iter
212         self.epsilon = epsilon
213         self.beta = beta
214
215     def k_fold_validation(self, normalize=False, inc_od=False, order=3,
rm_features=False):
216         '''
217             This function performs the k-fold validation
218
219             normalize - whether to perform normalization
220             inc_od - whether to increase the feature order
221             order - the order of the added feature
222         '''
223         folds = self.folds
224         data = pd.read_csv(self.path)
225
226         accuracies = []
227         accuracies_train = []
228         tic = time.time()
229
230         if rm_features:
231             data.drop(columns=['attribute43', 'attribute60'])
232
233         if inc_od:
234             print("order rise to {}".format(order))
235             data = self.rise_order(data, order)
236
237         log_reg = LogisticRegression(data=data, folds=self.folds, lr=self.lr,
max_iter=self.max_iter, beta=self.beta, epsilon=self.epsilon)
238

```

```

239     log_reg.shuffle_data()
240
241     for fold in range(folds):
242         t_X, t_y, v_X, v_y = log_reg.partition(fold)
243         # t_X --> test value X, v_X --> validation value X
244         gradient_val, t_acc_val, v_acc_val, best_w = log_reg.fit(t_X, t_y
, v_X, v_y, normalize=normalize)
245
246         accuracies.append(np.max(v_acc_val))
247         accuracies_train.append(np.max(t_acc_val))
248
249         # Uncomment this block to display the accuracy diagram
250         plt.figure()
251         plt.plot(t_acc_val, label = 'Training accuracy')
252         plt.plot(v_acc_val, label='Validation accuracy')
253         plt.axvline(log_reg.marker, color='r', label='Best Weights')
254         plt.xlabel('Iteration Number')
255         plt.ylabel('Accuracy')
256         plt.legend()
257         plt.show()
258         print("Learning Rate: " + str(log_reg.lr))
259         print("Average Accuracy: "+str(np.mean(accuracies)))
260
261         # Uncomment this block to display the gradient diagram
262         plt.figure()
263         plt.plot(gradient_val)
264         plt.xlabel('Iteration Number')
265         plt.ylabel('Gradient')
266         plt.show()
267         print("-----")
268
269         mean_acc = np.mean(accuracies)
270         mean_acc_train = np.mean(accuracies_train)
271         toc = time.time()
272         return mean_acc, mean_acc_train, toc-tic
273
274     def rise_order(self, data, order=3):
275         ret_val = data
276         for i in range(2, order + 1):
277             data_powered = data.pow(i)
278             ret_val = ret_val.iloc[:, :-1]
279             ret_val = pd.concat([ret_val, data_powered],axis=1)
280         return ret_val
281
282     """### Default Values"""
283
284     # path = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/hepatitis.
csv"
285     # dataset_name = "Hepatitis"
286     # default_lr = 0.01
287     # default_max_iter = 10000
288     # default_epsilon = 5e-3
289     # default_beta = 0.99
290
291     path = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/bankruptcy.csv"
292     dataset_name = "Bankruptcy"
293     default_lr = 0.1
294     default_max_iter = 25000
295     default_epsilon = 1e-3
296     default_beta = 0.99
297
298     unnorm_ordering = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
default_max_iter, epsilon=default_epsilon, beta=default_beta)

```

```

299 normed_ordering = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
    default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
300
301 order_min = 1
302 order_max = 10
303
304 orders = np.arange(order_min, order_max+1)
305
306 mean_acc = []
307 mean_acc_train = []
308 proc_time = []
309 mean_acc_norm = []
310 mean_acc_train_norm = []
311 proc_time_norm = []
312
313 for order in orders:
314     mean_acc_temp, mean_acc_train_temp, time_temp = unnorm_ordering.
    k_fold_validation(normalize=False, inc_od=True, order=order)
315     mean_acc.append(mean_acc_temp)
316     mean_acc_train.append(mean_acc_train_temp)
317     proc_time.append(time_temp)
318     mean_acc_temp, mean_acc_train_temp, time_temp = normed_ordering.
    k_fold_validation(normalize=True, inc_od=True, order=order)
319     mean_acc_norm.append(mean_acc_temp)
320     mean_acc_train_norm.append(mean_acc_train_temp)
321     proc_time_norm.append(time_temp)
322
323 plt.plot(orders, mean_acc_train, '--', label='Unnormalized Training Accuracy')
324 plt.plot(orders, mean_acc, '--', label='Unnormalized Validation Accuracy')
325 plt.plot(orders, mean_acc_train_norm, label='Normalized Training Accuracy')
326 plt.plot(orders, mean_acc_norm, label='Normalized Validation Accuracy')
327 plt.legend()
328 plt.xlabel("Maximum Order")
329 plt.ylabel("Mean Accuracy")
330 # Uncomment this to save the accuracy vs Order figure
331 # plt.savefig("Validation_Accuracy_vs_Maximum_Order.png", dpi = 1200)
332 # files.download("Validation_Accuracy_vs_Maximum_Order.png")
333 plt.show()
334
335 plt.plot(orders, proc_time, '--', label='Unnormalized Processing Time')
336 plt.plot(orders, proc_time_norm, label='Normalized Processing Time')
337 plt.legend()
338 plt.xlabel("Maximum Order")
339 plt.ylabel("Processing Time")
340 # Uncomment this to save the Processing time vs Order figure
341 # plt.savefig("Processing_Time_vs_Maximum_Order.png", dpi = 1200)
342 # files.download("Processing_Time_vs_Maximum_Order.png")
343 plt.show()
344
345 path = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/bankruptcy.csv"
346 dataset_name = "Bankruptcy"
347 default_lr = 0.1
348 default_max_iter = 25000
349 default_epsilon = 1e-3
350 defulat_beta = 0.99
351
352 full_feature_test = KFoldValidation(folds=10, path=path, lr=default_lr,
    max_iter=default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
353 rm_feature_test = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
    default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
354
355 mean_acc = []
356 proc_time = []

```

```
357
358 mean_acc_temp, mean_acc_train_temp, time_temp = full_feature_test.
    k_fold_validation(normalize=True, inc_od=True)
359 mean_acc.append(mean_acc_temp)
360 proc_time.append(time_temp)
361 mean_acc_temp, mean_acc_train_temp, time_temp = rm_feature_test.
    k_fold_validation(normalize=True, inc_od=True, rm_features=True)
362 mean_acc.append(mean_acc_temp)
363 proc_time.append(time_temp)
364
365 print("Without removing attribute 43 and 60\tmean accuracy is {} \tprocessing
    time is {}\nAfter removing those two attribute\tmean accuracy is {} \t
    processing time is {}".format(mean_acc[0],proc_time[0],mean_acc[1],proc_time[1
    ]))
```