```python
1  # -*- coding: utf-8 -*-
2  """Logistic_Regression
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1qioJbplkgpPKdiDEP2SYKmu6t7V-Maoo
8
9  <center><h1>Mini Project 1 - Logistic Regression</h1>
10 <h4>The hyperparameters and models used in this file are chosen based on the
   findings in the testing file.</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
28 path1 = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/hepatitis.csv"
29 path2 = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/bankrupcy.csv"
30
31 class LogisticRegression:
32     '''
33         This is the logistic regression class, containing fit, perdict and
   accu_eval functions,
34         as well as many other useful functions.
35     '''
36
37     def __init__(self, data, folds, lr=0.01, max_iter=10000, beta=0.99, epsilon
   =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
45     def shuffle_data(self):
46         '''
47             This function randomly shuffles the input dataset.
48         '''
49         # Load data from data file.
50         self.data.insert(0, column='Bias', value=1)
51         self.data = self.data.sample(frac=1)
52
53     def partition(self, fold):
54         '''
55             This function divides the dataset into training and validation set.
56
57             fold - the current fold
58         '''
59         data = self.data
60         # to exclude last term in previous partition for training data
```

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

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

```python
180            return count / len(y)
181
182     def log_func(self, alpha):
183            return 1 / (1 + np.exp(-alpha))
184
185     def gradient(self, X, y, w):
186            N = len(X[0])
187            y_hat = self.log_func(np.dot(X, w))
188            delta = np.dot(X.T, y_hat - y) / N
189            return delta
190
191 class KFoldValidation:
192     def __init__(self, folds, path, lr, max_iter, epsilon, beta):
193            self.folds = folds
194            self.data = pd.read_csv(path)
195            self.lr = lr
196            self.max_iter = max_iter
197            self.epsilon = epsilon
198            self.beta = beta
199
200     def k_fold_validation(self, normalize=False, inc_od=False, order=3):
201            '''
202                This function performs the k-fold validation
203
204                normalize - whether to perform normalization
205                inc_od - whether to increase the feature order
206                order - the order of the added feature
207            '''
208            folds = self.folds
209            data = self.data
210            accuracies = []
211
212            if inc_od:
213                data = self.rise_order(data, order)
214
215            log_reg = LogisticRegression(data=data, folds=self.folds, lr=self.lr,
    max_iter=self.max_iter, beta=self.beta, epsilon=self.epsilon)
216
217            log_reg.shuffle_data()
218
219            for fold in range(folds):
220                t_X, t_y, v_X, v_y = log_reg.partition(fold)
221                # t_X --> test value X, v_X --> validation value X
222                gradient_val, t_acc_val, v_acc_val, best_w = log_reg.fit(t_X, t_y
    , v_X, v_y, normalize=normalize)
223
224                accuracies.append(np.max(v_acc_val))
225
226                # Uncommant this block to display the accuracy diagram
227                plt.figure()
228                plt.plot(t_acc_val, label = 'Training accuracy')
229                plt.plot(v_acc_val, label='Validation accuracy')
230                plt.axvline(log_reg.marker, color='r', label='Best Weights')
231                plt.xlabel('Iteration Number')
232                plt.ylabel('Accuracy')
233                plt.legend()
234                plt.show()
235                print("Learning Rate: " + str(log_reg.lr))
236                print("Average Accuracy: "+str(np.mean(accuracies)))
237
238                # Uncommant this block to display the gradiant diagram
239                plt.figure()
240                plt.plot(gradient_val)
```

```python
241                 plt.xlabel('Iteration Number')
242                 plt.ylabel('Gradiant')
243                 plt.show()
244                 print("------------------------------------------------------")
245
246             mean_acc = np.mean(accuracies)
247             return mean_acc
248
249         def rise_order(self, data, order=3):
250             ret_val = data
251             for i in range(2, order + 1):
252                 data_powered = data.pow(i)
253                 ret_val = ret_val.iloc[:, :-1]
254                 ret_val = pd.concat([ret_val, data_powered],axis=1)
255             return ret_val
256
257 """### Perform 10-fold validation for Hepatitis dataset"""
258
259 path = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/hepatitis.csv"
260 dataset_name = "Hepatitis"
261 defult_lr = 0.01
262 default_max_iter = 10000
263 default_epsilon = 5e-3
264 defulat_beta = 0.99
265
266 # the input is the optimum hyperparameters found during testing
267 hepatitis_learning = KFoldValidation(folds=10, path=path, lr=defult_lr,
    max_iter=default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
268 # the input is the optimum model found during testing
269 mean_acc = hepatitis_learning.k_fold_validation()
270
271 """### Perform 10-fold validation for Bankruptcy dataset"""
272
273 path = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/bankrupcy.csv"
274 dataset_name = "Bankruptcy"
275 defult_lr = 0.1
276 default_max_iter = 25000
277 default_epsilon = 1e-3
278 defulat_beta = 0.99
279
280 # the input is the optimum hyperparameters found during testing
281 bankruptcy_learning = KFoldValidation(folds=10, path=path, lr=defult_lr,
    max_iter=default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
282 # the input is the optimum model found during testing
283 mean_acc = bankruptcy_learning.k_fold_validation(normalize=True, inc_od=True)
```