
ECSE 551 Machine Learning for Engineers

Mini-project 1 Report

Fei Peng

260712440

fei.peng@mail.mcgill.ca

Yukai Zhang

260710915

yukai.zhang@mail.mcgill.ca

Yi Zhu

260716006

yi.zhu6@mail.mcgill.ca

Abstract

Logistic regression is one of the most popular linear classification techniques in machine learning and it has been widely used in binary classification and predictions. In this mini-project, we investigated and evaluated the performance of logistic regression by implementing it using two benchmark datasets: *Hepatitis* and *Bankruptcy*. Both two datasets were pre-processed to explore their characteristics, thereby achieving the best outcome in the following training. Fitting, predicting, and testing were then performed to train and compare different models. The accuracy evaluations were analysed, and detailed results will be presented in the following parts of this report. Among them, we found that the accuracy of the logistic regression classifier could be improved by normalizing the datasets, changing order of the model, or removing certain features of the numerical datasets.

1 Introduction

In recent years, machine learning has been adapted to many different fields, including medical research and financial analysis [1]. As one of the most feasible and efficient classification methods, logistic regression is widely used to predict binary outcomes. The main objective of this project is to implement a logistic regression classifier and investigate its accuracy using k -fold cross validation, as well as analysis the impact of data preprocessing on the validation accuracy. We took the probabilistic views to process the binary classification. In probabilistic approaches, we focused on discriminative learning which directly estimates the probability of y (i.e., output class/label) given x (i.e., input data). Using Bayes' Rule, the logistic function (also known as sigmoid function) could be obtained, which takes a linear function of x as independent variable $w^T x$, and produces the conditional probability $P(y|x)$.

The most critical task of implementing the logistic regression algorithm was to find the linear term w (i.e., weights) that produces the minimum error/maximum accuracy in prediction while ensuring the efficiency of this algorithm. This was achieved using gradient descent with momentum [2], and tweaking the hyperparameters (e.g., learning rate, momentum term, and stopping criteria) of the fit function. It was found that, with a carefully chosen set of hyperparameters, accuracy could be increased while remaining satisfying training efficiency.

During data pre-processing, normalization was performed and proved to increase the accuracy. While analysing the characteristics of the datasets, it was noticeable that some features of input x satisfies the null hypothesis [3], which means that these features are not helping discriminate between two classes and should be removed to improve performance. Moreover, through the distribution of the two classes of the two datasets, it was discovered that the entropy of classes could affect the accuracy of the model, resulting in the model accuracy for hepatitis data (with higher entropy) being higher

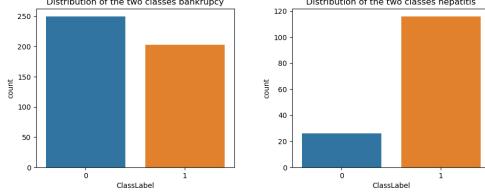


Figure 1: *Bankruptcy* data class distribution (left) and *Hepatitis* data class distribution (right)

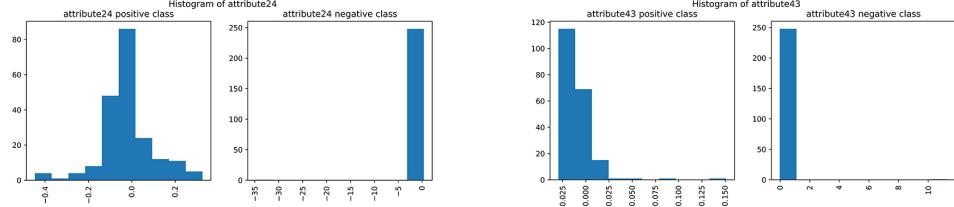


Figure 2: Attribute 24 (left) and attribute 43 (right) feature distribution for *Bankruptcy* data

than that of *bankruptcy* data. Last but not least, properly increasing the order of the model (order 3 for *bankruptcy* data) would also help increase accuracy. These aspects will be further discussed.

2 Datasets

Data preprocessing usually lays the groundwork for the later data analysis and helps it yield a higher accuracy efficiently. In this project, the following data preprocessing had been performed.

2.1 Entropy Analysis

Both datasets *Bankruptcy* and *Hepatitis* are examples of binary classification. The distributions of each class are shown in Figure 1. The uncertainty in prediction can be quantified by the entropy of each dataset with equation 1

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x) \quad (1)$$

2.2 Features Analysis

For the purpose of gaining a better understanding toward the dataset, the distributions of each feature in both classes are plotted. Some of the plots are shown in Figure 2

Note that both pairs of plots come from *Bankruptcy* dataset and in each pair, the distribution of positive class is shown on the left, otherwise right. When comparing positive and negative classes, it is noticeable that the data distribution differs greatly in attribute 24 but slightly in attribute 43. This indicates that features with similar distribution to attribute 43 will have little contribution in the prediction process and should be considered as null distribution. Kolmogorov-Smirnov method [4] was performed to efficiently find all null distributions. Attributes 43 and 60 in the *Bankruptcy* dataset were identified. The removal of these two features may improve the performance of the machine learning process. The Kolmogorov-Smirnov method, however, failed to yield a promising result in the *Hepatitis* dataset because of the existence of multiple binary features. The effect of removing features with null distribution will be further discussed in the result section.

2.3 Data Shuffling

It was found that the *Bankruptcy* dataset is ordered by class label. Proceeding without shuffling the data would result in some validation/training sets all filled by data from one class. Thus, rows will be shuffled every time after importation. Doing so will also help reduce the variance and generalize the algorithm, therefore the *Hepatitis* dataset should also be shuffled.

2.4 Data Normalization

Normalizing the data into z -score can centralize the value distribution of the data while keeping its precision [5]. When a calculation involves large float numbers, they can easily trigger overflow/underflow and lose precision without normalization. Calculating the gradient with high order data is an example of that. Luckily, using z -score should help prevent introducing such errors and thus increase the model accuracy.

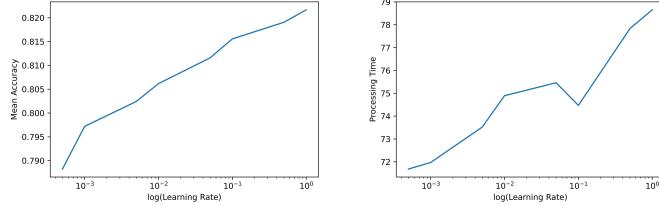


Figure 3: Mean accuracy (left) and processing time (right) vs. logarithm of learning rate

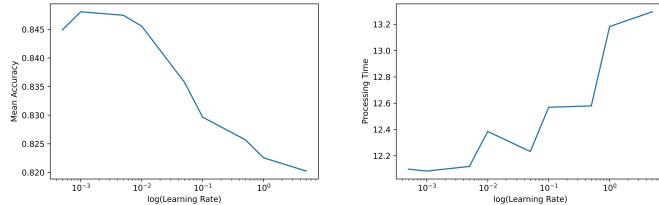


Figure 4: Mean accuracy (left) and processing time (right) vs. logarithm of learning rate

3 Results

Multiple experiments have been conducted to evaluate the performance of the logistic regression with respect to different hyperparameters and data preprocessing techniques. The detailed results are shown below.

3.1 Learning Rate

Learning rate is a parameter that determines the step size on the way approaching a (local) minimum. A large step size can accelerate the convergence but can also result in jumping over the desired solution (oscillate forever). A small step size can avoid skipping the solution but is very time consuming. To find a learning rate best fit for each model, values between 0.0005 and 5 are tested. The performance/processing time vs. logarithm of learning rate are shown in Figure 3 and Figure 4.

3.2 Stopping Criteria:

Theoretically, if a model suits the data distribution, the training accuracy will keep increasing with more iteration calculated. However, failing to set up a proper stopping criteria will not only waste the computational power but also lead to overfitting [6]. In this case, max iteration and epsilon(minimum gradient) are introduced as the upper limit of the iterations. In this experiment, the max iteration term is sampled from 500 to 25000 and the epsilon term is sampled from $1e^{-3}$ to 0.5. The plots of accuracy/processing time versus max iteration/epsilon are shown in Figure 5 and Figure 6.

Note that all the data displayed come from the *Bankruptcy* dataset. The *Hepatitis* dataset has a similar trend and will not be displayed. The plots confirm that in a reasonable iteration range, accuracy tends to rise as the maximum iteration increases and the minimum acceptable gradient decreases. For the trade-off between accuracy and efficiency, 25000 iterations with an epsilon of $1e^{-3}$ is selected for the *Bankruptcy* dataset and $(10000, 5e^{-3})$ is selected for the *Hepatitis* dataset.

3.3 Momentum Gradient Descent Constant - Beta

When a regular gradient descent method is performed, it is very unlikely that the algorithm can move directly from the starting point to the (local) minimum. Instead, the path usually oscillates back

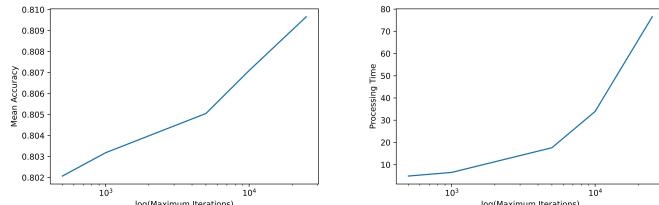


Figure 5: Mean accuracy (left) and Processing time (right) with respect to log of Max. Iteration

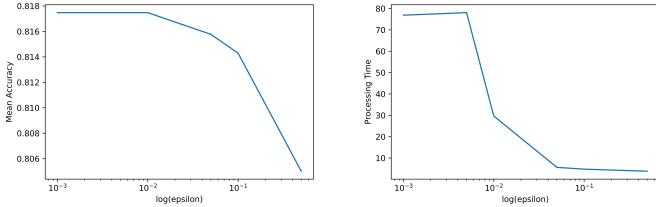


Figure 6: Mean accuracy (left) and processing time (right) with respect to log of epsilon

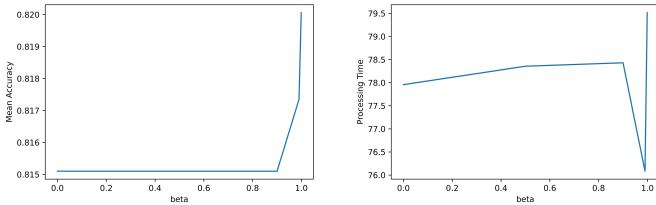


Figure 7: Mean accuracy (left) and processing time (right) with respect to beta

and forth around that optimum shortcut. As described in part 3.1, this oscillation will have a larger amplitude if the learning rate is high. A momentum gradient descent method can be used to flatten this oscillation and take a rather straightforward approach in the search of the minimum point. It does so by calculating the weighted sum of the current gradient and the previous step size and then using it as the new step size. β value is in charge of controlling the weight of the previous step size, if β is set to 0, it becomes equivalent to a regular gradient descent method [2]. In this experiment, value β is tested in the range of 0 to 0.999 and its results are shown in Figure 7.

The data displayed all come from *Bankruptcy* dataset. The trends of that in the *Hepatitis* dataset is also very similar to those above. According to the figures, a β of 0.99 yields the best overall performance regarding the accurateness and efficiency.

3.4 Normalization

During data preprocessing, it was discovered that most of the features in *Bankruptcy* data are floating points. Therefore, calculation may result in precision loss due to overflow/underflow. To improve model accuracy, data should be normalized before performing training or validation. The normalization technique used in this project is called *z-score* [5] which calculates the mean (μ_j) and standard deviation (σ_j) for each feature (x_j) of the training data, and perform $\frac{x_j - \mu_j}{\sigma_j}$ for both the training and validation data. Therefore the model accuracy could be improved by bringing more precision to the calculation of gradient. During the testing, it was found that the overflow warning indeed disappeared after normalization. It is shown in Figure 8 on the left that compared to that of without normalization, the validation accuracy of the normalized original dataset is higher. This advantage is even pronounced when more features (higher orders) are introduced, which will be discussed later.

3.5 Removing Features

As discussed previously in datasets, null distributions in the *Bankruptcy* dataset were found using Kolmogorov-Smirnov method [4]. This means that the distribution of some attributes in two classes

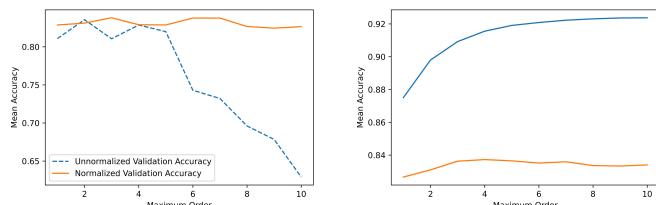


Figure 8: Comparison between with and without normalization (left) and the training and validation accuracy vs. order of *Bankruptcy* data (right)

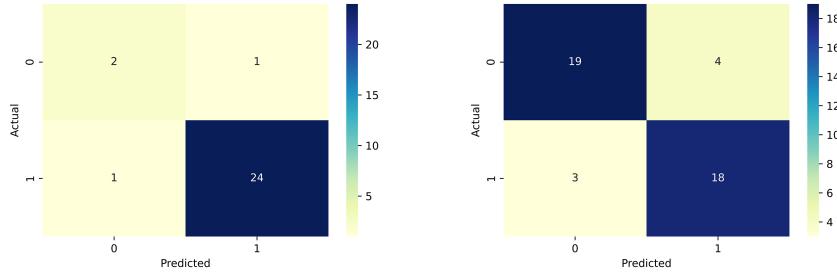


Figure 9: The Confusion Matrix. for *Hepatitis* data (left) and *Bankruptcy* data (right)

are similar and may not be useful for training and class prediction. Attribute 43 and 60 of *Bankruptcy* dataset were identified to be unhelpful, therefore were removed during training and validation. However, it came out that the mean validation accuracy and training efficiency did not vary significantly from before. One of the assumptions that could explain this result is that the linear term, w , of the log-odds ratio (i.e., the independent variable, $w^T x$, of the logistic function) as already taken care of the weight of each attribute, therefore removing any attribute would not have much impact on the accuracy.

3.6 Adding features

The approach of adding more features to the dataset was also adopted to improve model accuracy. The order of each attribute in the *Bankruptcy* dataset was increased and added to the original dataset. As shown in Figure 8 on the right, the blue curve indicates the training accuracy, while the orange curve represents the validation accuracy. It clearly shows that when the order is increased to around 3, the maximum validation accuracy is achieved. However, continuing increasing the order could cause overfitting [6], as the training accuracy remains increasing but the validation accuracy starts to decrease. Therefore, order 3 was chosen to be the most suitable model for training *Bankruptcy* dataset. This approach, however, was tested to be ineffective for *Hepatitis* dataset due to the fact that most of its features are binary.

4 Discussion and Conclusion

In conclusion, the logistic regression classifier implemented in this project performed as expected overall on the two datasets, *Bankruptcy* and *Hepatitis*. Here are the two representative confusion matrices, shown in Figure 9, calculated during testing to measure the accuracy, precision and sensitivity of the models. The selected model used for calculating the confusion matrix was the one that performed the best during testing (as explained in the previous section). The accuracy for the best *Hepatitis* model could achieve 92.86%, with 96% specificity. However the number of instances in this dataset is too small (only 142 total) to take any common measurements, therefore 5-fold validation was performed to obtain this confusion matrix. Moreover, this is one of the best case scenarios, thus in the future, more testing should be performed, ideally with more instances, to obtain a convincing average accuracy. For the best *Bankruptcy* model, the accuracy reached 84.09%, while the precision was 82.61% and the sensitivity was 86.36%. These measurements confirmed the previous hypothesis, that lower entropy can yield a higher accuracy.

Possible directions for future investigation could include choosing other types of model (e.g., with log, exponential, or interaction terms), exploring more efficient but complex gradient descent algorithms, and comparing logistic regression with genetic learning as *Hepatitis* is a relatively small dataset. These aspects should help not only improve model accuracy, but also accelerate the training process.

5 Statement of Contributions

The workload of this project is distributed equally between the team members (i.e., the three authors of this report). Yi Zhu implemented the fit function for training datasets, and explored various gradient descent algorithms to improve the computation speed. Fei Peng was responsible for developing the k -fold partition, predict and accuracy evaluation functions for predicting the classes/labels of the input data and evaluating the model accuracy, as well as performing thorough testing of the models. Yukai was in charge of data pre-processing and model selection, which helped to significantly improve the model accuracy.

References

- [1] T. Peña, S. Martínez, and B. Abudu, "Bankruptcy Prediction: A Comparison of Some Statistical and Machine Learning Techniques," in *Computational Methods in Economic Dynamics*, H. Dawid and W. Semmler Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 109-131.
- [2] S. Ruder, "An overview of gradient descent optimization algorithms," *arXiv preprint arXiv:1609.04747*, 2016.
- [3] R. W. Vitral, M. J. Campos, and M. R. Fraga, "The null hypothesis," *American journal of orthodontics and dentofacial orthopedics : official publication of the American Association of Orthodontists, its constituent societies, and the American Board of Orthodontics*, vol. 144, no. 4, pp. 498-9, 2013, doi: 10.1016/j.ajodo.2013.08.010.
- [4] T. Liu, "A Kolmogorov-Smirnov type test for two inter-dependent random variables," *arXiv preprint arXiv:1802.09899*, 2018.
- [5] A. E. Curtis, T. A. Smith, B. A. Ziganshin, and J. A. Elefteriades, "The Mystery of the Z-Score," (*in eng*), *Aorta (Stamford)*, vol. 4, no. 4, pp. 124-130, 2016, doi: 10.12945/j.aorta.2016.16.014.
- [6] AH. Allamy, "METHODS TO AVOID OVER-FITTING AND UNDER-FITTING IN SUPERVISED MACHINE LEARNING (COMPARATIVE STUDY)," 12/27 2014.

Appendix

(Please see following pages for code implementations attached)

```

1 # -*- coding: utf-8 -*-
2 """Data_Preprocessing
3
4 Automatically generated by Colaboratory.
5
6 Original file is located at
7     https://colab.research.google.com/drive/1cqi0tupPBH0ZZWkk2pkLA-k84lVnQgSB
8
9 <center><h1>Mini Project 1 - Logistic Regression</h1>
10 <h4>This file is for data preprocessing. Important characteristics of features
and the distribution of classes could be found in the result of this file.</h4>
11 </center>
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 seaborn as sns
26 import matplotlib.pyplot as plt
27 import random
28 from scipy import stats
29 from google.colab import files
30
31 path1 = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/hepatitis.csv"
32 path2 = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/bankrupcy.csv"
33
34 hepatitis_data = pd.read_csv(path1)
35 bankrupcy_data = pd.read_csv(path2)
36
37 class Data_Processing:
38     def __init__(self, data, name = 'New Data'):
39         self.data = data
40         self.name = name
41
42     def partition_by_class(self):
43         pos, neg = [], []
44         data = self.data
45         for row in range(1, data.shape[0]):
46             if data[row, -1] == 1:
47                 pos.append(list(data[row]))
48             else:
49                 neg.append(list(data[row]))
50         self.pos = pos
51         self.neg = neg
52
53     def show_y_dist(self, ydata):
54         plt.figure(figsize=(5,4))
55         plt.subplot(111), sns.countplot(x='ClassLabel', data=ydata)
56         plt.title('Distribution of the two classes {}'.format(self.name))
57         plt.savefig("Distribution of the two classes {}.png".format(self.name),
58 ), dpi = 1200)
58         files.download("Distribution of the two classes {}.png".format(self.
name))
59         plt.show()

```

```

60
61     def show_x_dist(self, xdata):
62         pos = self.pos
63         neg = self.neg
64         for i in range(0,len(pos[1])):
65             fig, (ax1, ax2) = plt.subplots(1, 2)
66             fig.set_size_inches(10,4)
67             plt.setp(ax1.xaxis.get_majorticklabels(), rotation=90)
68             plt.setp(ax2.xaxis.get_majorticklabels(), rotation=90)
69             fig.suptitle('Histogram of {}'.format(xdata.keys()[i]))
70             ax1.title.set_text('{} positive class'.format(xdata.keys()[i]))
71             ax1.hist([pos[j][i] for j in range(len(pos))])
72             ax2.title.set_text('{} negative class'.format(xdata.keys()[i]))
73             ax2.hist([neg[j][i] for j in range(len(neg))])
74
75     def find_null_data(self):
76         data, pos, neg = self.data, self.pos, self.neg
77         num_data = min(len(pos), len(neg))
78         num_feature = len(pos[0])
79         null_feature_count = np.zeros(num_feature)
80         pos = random.sample(pos, k = num_data)
81         neg = random.sample(neg, k = num_data)
82         for i in range(len(pos[1])):
83             posi_list = []
84             nega_list = []
85             for j in range(len(pos)):
86                 posi_list.append(pos[j][i])
87                 nega_list.append(neg[j][i])
88             a, b = stats.ks_2samp(posi_list, nega_list)
89             if(b > 0.35):
90                 null_feature_count[i] += 1
91             elif ((b > 0.10) and (a > 0.10)):
92                 null_feature_count[i] += 0.5
93         if sum(null_feature_count) > num_feature * 0.1:
94             print("cannot remove this many features")
95         else:
96             print("Here are the features we recomend you to delete")
97             for i in range(len(null_feature_count)):
98                 if null_feature_count[i] > 0:
99                     print("{}: {}".format(i, null_feature_count[i]))
100
101 """## Plot the classes and features distribution for Hepatitis data"""
102
103 hepatitis_data = pd.read_csv(path1)
104 data = hepatitis_data
105 data1 = Data_Processing(data.values, 'hepatitis')
106 data1.partition_by_class()
107 data1.show_y_dist(data)
108 data1.show_x_dist(data)
109 data1.find_null_data()
110
111 """## Plot the classes and features distribution for Bankruptcy data"""
112
113 bankrupcy_data = pd.read_csv(path2)
114 data = bankrupcy_data
115 data1 = Data_Processing(data.values, 'bankrupcy')
116 data1.partition_by_class()
117 data1.show_y_dist(data)
118 data1.show_x_dist(data)
119 data1.find_null_data()
120
121 """After tests, we find that in bankrupcy.csv, if features: 43, 60, (36) are
 deleted, it may deliver a better result.

```

```
122 <br>
123 43, 60 very high p-value in ks test
124 <br>
125 36 relative high but steady p-value
126 <hr>
127
128 Hepatitis cannot use this method since we have many features are only in 1 or
     0.
129 """
```

```

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
11   findings in the testing file.</h4></center>
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, predict and
34         accu_eval functions,
35             as well as many other useful functions.
36     """
37
38     def __init__(self, data, folds, lr=0.01, max_iter=10000, beta=0.99, epsilon
39 =5e-3):
40         self.data = data
41         self.folds = folds
42         self.lr = lr
43         self.max_iter = max_iter
44         self.beta = beta
45         self.epsilon = epsilon
46
47     def shuffle_data(self):
48         """
49             This function randomly shuffles the input dataset.
50         """
51         # Load data from data file.
52         self.data.insert(0, column='Bias', value=1)
53         self.data = self.data.sample(frac=1)
54
55     def partition(self, fold):
56         """
57             This function divides the dataset into training and validation set.
58
59             fold - the current fold
60         """
61         data = self.data
62         # to exclude last term in previous partition for training data

```

File - E:\Study\ECSE551\Mini_Project_1\code\logistic_regression.py

```

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
102        labels vector y
103        as well as other hyperparameters (such as learning rate) as input,
104        and execute the model training through modifying the model
105        parameters (i.e. W).
106
107        X - training data
108        y - class of training data
109        v_X - validation data
110        v_y - class of validation data
111        epsilon - the threshold value for gradient descent
112        normalize - whether to perform normalization
113
114        gradient_values, t_acc_val, v_acc_val = [], [], []
115
116        if normalize:
117            X, v_X = self.normalization(X, v_X)
118
119            # Retrive the learning rate, maximum iteration, momentum (beta)
120            lr, max_iter, beta, epsilon = self.lr, self.max_iter, self.beta, self.
epsilon

```

```

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
130      # maximum iteration, stop
131      while (np.linalg.norm(dw) > epsilon and iteration <= max_iter):
132          dw = self.gradient(X, y, w)
133          gradient_values.append(np.linalg.norm(dw))
134          # if beta = 0, it will be the same as general gradient descent
135          step = beta * step + (1 - beta) * dw # gradient descent with
136          # momentum
137          w = w - lr * step
138
139          # predict once every 10 interations
140          if iteration % 10 == 0:
141              t_y_pred = self.predict(X, w)
142              t_acc = self.accu_eval(t_y_pred, y)
143              v_y_pred = self.predict(v_X, w)
144              v_acc = self.accu_eval(v_y_pred, v_y)
145
146          # record the next best value
147          if v_acc >= v_acc_max:
148              v_acc_max = v_acc
149              best_w = w
150              self.marker = iteration # move the iteration marker
151
152          t_acc_val.append(t_acc)
153          v_acc_val.append(v_acc)
154
155          iteration = iteration + 1
156
157      return gradient_values, t_acc_val, v_acc_val, best_w
158
159  def predict(self, X, w):
160      """
161          This function takes a set of data as input and outputs predicted
162          labels for the input points.
163      """
164      result = self.log_func(np.dot(X, w))
165      # the prediction result converted to binary
166      predict_bin = []
167      for i in result:
168          if i>=0.5:
169              predict_bin.append(1)
170          else:
171              predict_bin.append(0)
172
173      return predict_bin
174
175  def accu_eval(self, y_pred, y):
176      """
177          This function evaluates the models' accuracy.
178      """
179      count = 0
180      for i in range(len(y_pred)):
181          if y_pred[i] == y[i]:
182              count = count + 1
183
184      # return the accuracy ratio: #corret prediction / #data points

```

```

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,
216                                     max_iter=self.max_iter, beta=self.beta, epsilon=self.epsilon)
217
218         log_reg.shuffle_data()
219
220         for fold in range(folds):
221             t_X, t_y, v_X, v_y = log_reg.partition(fold)
222             # t_X --> test value X, v_X --> validation value X
223             gradient_val, t_acc_val, v_acc_val, best_w = log_reg.fit(t_X, t_y
224             , v_X, v_y, normalize=normalize)
225
226             accuracies.append(np.max(v_acc_val))
227
228             # Uncomment this block to display the accuracy diagram
229             plt.figure()
230             plt.plot(t_acc_val, label='Training accuracy')
231             plt.plot(v_acc_val, label='Validation accuracy')
232             plt.axvline(log_reg.marker, color='r', label='Best Weights')
233             plt.xlabel('Iteration Number')
234             plt.ylabel('Accuracy')
235             plt.legend()
236             plt.show()
237             print("Learning Rate: " + str(log_reg.lr))
238             print("Average Accuracy: "+str(np.mean(accuracies)))
239
240             # Uncomment this block to display the gradient diagram
241             plt.figure()
242             plt.plot(gradient_val)

```

```
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             dataPowered = data.pow(i)
253             ret_val = ret_val.iloc[:, :-1]
254             ret_val = pd.concat([ret_val, dataPowered], 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,
268                                         max_iter=default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
269 # the input is the optimum model found during testing
270 mean_acc = hepatitis_learning.k_fold_validation()
271
272 """## Perform 10-fold validation for Bankruptcy dataset"""
273
274 path = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/bankrupcy.csv"
275 dataset_name = "Bankruptcy"
276 defult_lr = 0.1
277 default_max_iter = 25000
278 default_epsilon = 1e-3
279 defulat_beta = 0.99
280
281 # the input is the optimum hyperparameters found during testing
282 bankruptcy_learning = KFoldValidation(folds=10, path=path, lr=defult_lr,
283                                         max_iter=default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
284 # the input is the optimum model found during testing
285 mean_acc = bankruptcy_learning.k_fold_validation(normalize=True, inc_od=True)
```

```

1 # -*- coding: utf-8 -*-
2 """Hyperparameter_Testing
3
4 Automatically generated by Colaboratory.
5
6 Original file is located at
7     https://colab.research.google.com/drive/1Tkjs6AN6HaRH09FqtLdYzPa6BQHDaYwC
8
9 <center><h1>Mini Project 1 - Logistic Regression</h1>
10 <h4>This is a testing file aiming to find the best hyperparameters for 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
29 class LogisticRegression:
30     """
31         This is the logistic regression class, containing fit, predict and
accu_eval functions,
32             as well as many other useful functions.
33     """
34
35     def __init__(self, data, folds, lr=0.01, max_iter=10000, beta=0.99, epsilon
=5e-3):
36         self.data = data
37         self.folds = folds
38         self.lr = lr
39         self.max_iter = max_iter
40         self.beta = beta
41         self.epsilon = epsilon
42
43     def shuffle_data(self):
44         """
45             This function randomly shuffles the input dataset.
46         """
47         # Load data from data file.
48         self.data.insert(0, column='Bias', value=1)
49
50         self.data = self.data.sample(frac=1)
51
52     def set_learning_rate(self, lr):
53         self.lr = lr
54
55     def set_max_iter(self, max_iter):
56         self.max_iter = max_iter
57
58     def set_epsilon(self, epsilon):
59         self.epsilon = epsilon
60

```

```

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

```

```

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

```

```

179     return predict_bin
180
181     def accu_eval(self, y_pred, y):
182         """
183             This function evaluates the models' accuracy.
184         """
185         count = 0
186         for i in range(len(y_pred)):
187             if y_pred[i] == y[i]:
188                 count = count + 1
189         # return the accuracy ratio: #correct prediction / #data points
190         return count / len(y)
191
192     def log_func(self, alpha):
193         return 1 / (1 + np.exp(-alpha))
194
195     def gradient(self, X, y, w):
196         N = len(X[0])
197         y_hat = self.log_func(np.dot(X, w))
198         delta = np.dot(X.T, y_hat - y) / N
199         return delta
200
201 class KFoldValidation:
202     def __init__(self, folds, path, lr, max_iter, epsilon, beta):
203         self.folds = folds
204         self.data = pd.read_csv(path)
205         self.lr = lr
206         self.max_iter = max_iter
207         self.epsilon = epsilon
208         self.beta = beta
209         self.log_reg = LogisticRegression(data=self.data, folds=folds, lr=lr,
210                                         max_iter=max_iter, beta=beta, epsilon=epsilon)
211         self.log_reg.shuffle_data()
212
213     def set_learning_rate(self, lr):
214         self.log_reg.set_learning_rate(lr)
215
216     def set_max_iter(self, max_iter):
217         self.log_reg.set_max_iter(max_iter)
218
219     def set_epsilon(self, epsilon):
220         self.log_reg.set_epsilon(epsilon)
221
222     def set_beta(self, beta):
223         self.log_reg.set_beta(beta)
224
225     def k_fold_validation(self):
226         """
227             This function performs the k-fold validation
228
229             normalize - whether to perform normalization
230             inc_od - whether to increase the feature order
231             order - the order of the added feature
232         """
233         folds = self.folds
234         data = self.data
235         log_reg = self.log_reg
236         accuracies = []
237         tic = time.time()
238
239         for fold in range(folds):
240             t_X, t_y, v_X, v_y = log_reg.partition(fold)
# t_X --> test value X, v_X --> validation value X

```

```

241         gradient_val, t_acc_val, v_acc_val, best_w = log_reg.fit(t_X, t_y
242             , v_X, v_y)
243             accuracies.append(np.mean(v_acc_val))
244
245             # # Uncomment this block to display the accuracy diagram
246             # plt.figure()
247             # plt.plot(t_acc_val, label = 'Training accuracy')
248             # plt.plot(v_acc_val, label='Validation accuracy')
249             # plt.axvline(log_reg.marker, color='r', label='Best Weights')
250             # plt.xlabel('Iteration Number')
251             # plt.ylabel('Accuracy')
252             # plt.legend()
253             # plt.show()
254             # print("Learning Rate: " + str(log_reg.ln))
255             # print("Average Accuracy: "+str(np.mean(accuracies)))
256
257             # # Uncomment this block to display the gradiant diagram
258             # plt.figure()
259             # plt.plot(gradient_val)
260             # plt.xlabel('Iteration Number')
261             # plt.ylabel('Gradiant')
262             # plt.show()
263             # print("-----")
264
265     mean_acc = np.max(accuracies)
266     toc = time.time()
267     return mean_acc, toc-tic
268
269 def rise_order(self, data, order=3):
270     ret_val = data
271     for i in range(2, order + 1):
272         dataPowered = data.pow(i)
273         ret_val = ret_val.iloc[:, :-1]
274         ret_val = pd.concat([ret_val, dataPowered], axis=1)
275     return ret_val
276
277 def plot_fig(x, y, xlabel, ylabel, log_x=False, plt_name="Untitled_Figure.png"
278             , download=False):
279     plt.plot(x, y)
280     if log_x:
281         plt.xscale("log")
282     plt.xlabel(xlabel)
283     plt.ylabel(ylabel)
284     if download:
285         plt.savefig(plt_name, dpi = 1200)
286         files.download(plt_name)
287     plt.show()
288 """
289 """## Default Values for Hepatitis Analysis"""
290 path = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/hepatitis.csv"
291 dataset_name = "Hepatitis"
292 defult_lr = 0.01
293 default_max_iter = 10000
294 default_epsilon = 5e-3
295 defulat_beta = 0.99
296 """
297 """## Learning rate testing"""
298
299 lr_testing = KFoldValidation(folds=10, path=path, lr=defult_lr, max_iter=
300     default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
300

```

```

301 Learning_rates = [0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1]
302
303 mean_acc = []
304 proc_time = []
305
306 for lr in Learning_rates:
307     lr_testing.set_learning_rate(lr)
308     mean_acc_temp, time_temp = lr_testing.k_fold_validation()
309     mean_acc.append(mean_acc_temp)
310     proc_time.append(time_temp)
311
312 plot_fig(Learning_rates, mean_acc, "log(Learning Rate)", "Mean Accuracy",
313           log_x=True, plt_name="Validation_Accuracy_vs_Learning_Rate_for_{}.png".format(
314           dataset_name))
313 plot_fig(Learning_rates, proc_time, "log(Learning Rate)", "Processing Time",
314           log_x=True, plt_name="Validation_Accuracy_vs_Learning_Rate_for_{}.png".format(
314           dataset_name))
314
315 """### Maximum Iteration Test"""
316
317 max_iter_testing = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
318                                     default_max_iter, epsilon=1e-3, beta=defulat_beta)
318
319 max_iters = [500, 1000, 5000, 10000, 25000]
320
321 mean_acc = []
322 proc_time = []
323
324 for max_iter in max_iters:
325     max_iter_testing.set_max_iter(max_iter)
326     mean_acc_temp, time_temp = max_iter_testing.k_fold_validation()
327     mean_acc.append(mean_acc_temp)
328     proc_time.append(time_temp)
329
330 plot_fig(max_iters, mean_acc, "log(Maximum Iterations)", "Mean Accuracy",
331           log_x=True, plt_name="Validation_Accuracy_vs_Maximum_Iterations_for_{}.png".
331           format(dataset_name))
331 plot_fig(max_iters, proc_time, "log(Maximum Iterations)", "Processing Time",
332           log_x=True, plt_name="Processing_Time_vs_Maximum_Iterations_for_{}.png".format(
332           dataset_name))
332
333 """### Epsilon Test"""
334
335 epsilon_testing = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
336                                     default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
336
337 epsilons = [1e-3, 5e-3, 1e-2, 5e-2, 1e-1, 0.5]
338
339 mean_acc = []
340 proc_time = []
341
342 for epsilon in epsilons:
343     epsilon_testing.set_epsilon(epsilon)
344     mean_acc_temp, time_temp = epsilon_testing.k_fold_validation()
345     mean_acc.append(mean_acc_temp)
346     proc_time.append(time_temp)
347
348 plot_fig(epsilons, mean_acc, "log(Epsilons)", "Mean Accuracy", log_x=True,
349           plt_name="Validation_Accuracy_vs_Epsilons_for_{}.png".format(dataset_name))
349 plot_fig(epsilons, proc_time, "log(Epsilons)", "Processing Time", log_x=True,
350           plt_name="Processing_Time_vs_Epsilons_for_{}.png".format(dataset_name))
350
351 """### Momentum Gradient Descent Constant - Beta Testing"""

```

```

352
353 beta_testing = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
354     default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
355 betas = [0, 0.5, 0.9, 0.99, 0.999]
356
357 mean_acc = []
358 proc_time = []
359
360 for beta in betas:
361     beta_testing.set_beta(beta)
362     mean_acc_temp, time_temp = beta_testing.k_fold_validation()
363     mean_acc.append(mean_acc_temp)
364     proc_time.append(time_temp)
365
366 plot_fig(betas, mean_acc, "betas", "Mean Accuracy", log_x=False, plt_name="
367 Validation_Accuracy_vs_betas_for_{}.png".format(dataset_name))
368 plot_fig(betas, proc_time, "betas", "Processing Time", log_x=False, plt_name="
369 Processing_Time_vs_betas_for_{}.png".format(dataset_name))
370 """
371 ### Default Values for Bankruptcy Analysis
372 """
373
374 path = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/bankrupcy.csv"
375 dataset_name = "Bankruptcy"
376 default_lr = 0.1
377 default_max_iter = 25000
378 default_epsilon = 1e-3
379 defulat_beta = 0.99
380
381 """
382 ### Learning rate testing
383
384 lr_testing = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
385     default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
386
387 mean_acc = []
388 proc_time = []
389
390 for lr in Learning_rates:
391     lr_testing.set_learning_rate(lr)
392     mean_acc_temp, time_temp = lr_testing.k_fold_validation()
393     mean_acc.append(mean_acc_temp)
394     proc_time.append(time_temp)
395
396 plot_fig(Learning_rates, mean_acc, "log(Learning Rate)", "Mean Accuracy",
397     log_x=True, plt_name="Validation_Accuracy_vs_Learning_Rate_for_{}.png".format(
398         dataset_name))
399 plot_fig(Learning_rates, proc_time, "log(Learning Rate)", "Processing Time",
400     log_x=True, plt_name="Validation_Accuracy_vs_Learning_Rate_for_{}.png".format(
401         dataset_name))
402
403 """
404 ### Maximum Iteration Test
405
406 max_iter_testing = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
407     default_max_iter, epsilon=1e-3, beta=defulat_beta)
408
409 max_iters = [500, 1000, 5000, 10000, 25000]
410
411 mean_acc = []

```

```

406 proc_time = []
407
408 for max_iter in max_iters:
409     max_iter_testing.set_max_iter(max_iter)
410     mean_acc_temp, time_temp = max_iter_testing.k_fold_validation()
411     mean_acc.append(mean_acc_temp)
412     proc_time.append(time_temp)
413
414 plot_fig(max_iters, mean_acc, "log(Maximum Iterations)", "Mean Accuracy",
415           log_x=True, plt_name="Validation_Accuracy_vs_Maximum_Iterations_for_{}.png".
416           format(dataset_name))
417 plot_fig(max_iters, proc_time, "log(Maximum Iterations)", "Processing Time",
418           log_x=True, plt_name="Processing_Time_vs_Maximum_Iterations_for_{}.png".format(
419             dataset_name))
420
421 """
422     ### Epsilon Test
423
424 epsilon_testing = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
425                                     default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
426
427 epsilons = [1e-3, 5e-3, 1e-2, 5e-2, 1e-1, 0.5]
428
429 mean_acc = []
430 proc_time = []
431
432 for epsilon in epsilons:
433     epsilon_testing.set_epsilon(epsilon)
434     mean_acc_temp, time_temp = epsilon_testing.k_fold_validation()
435     mean_acc.append(mean_acc_temp)
436     proc_time.append(time_temp)
437
438 plot_fig(epsilons, mean_acc, "log(Epsilons)", "Mean Accuracy", log_x=True,
439           plt_name="Validation_Accuracy_vs_Epsilons_for_{}.png".format(dataset_name))
440 plot_fig(epsilons, proc_time, "log(Epsilons)", "Processing Time", log_x=True,
441           plt_name="Processing_Time_vs_Epsilons_for_{}.png".format(dataset_name))
442
443 """
444     ### Momentum Gradient Descent Constant - Beta Testing
445
446 beta_testing = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
447                                 default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
448
449 betas = [0, 0.5, 0.9, 0.99, 0.999]
450
451 mean_acc = []
452 proc_time = []
453
454 for beta in betas:
455     beta_testing.set_beta(beta)
456     mean_acc_temp, time_temp = beta_testing.k_fold_validation()
457     mean_acc.append(mean_acc_temp)
458     proc_time.append(time_temp)
459
460 plot_fig(betas, mean_acc, "betas", "Mean Accuracy", log_x=False, plt_name="
461 Validation_Accuracy_vs_betas_for_{}.png".format(dataset_name))
462 plot_fig(betas, proc_time, "betas", "Processing Time", log_x=False, plt_name="
463 Processing_Time_vs_betas_for_{}.png".format(dataset_name))

```

```

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
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
61         train_add = 1 if fold < self.folds else 0

```

```

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)/
72      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
98  def fit(self, X, y, v_X, v_y, normalize=False):
99      """
100         This function takes the training data X and its corresponding
101        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
104        parameters (i.e. W).
105
106        X - training data
107        y - class of training data
108        v_X - validation data
109        v_y - class of validation data
110        epsilon - the threshold value for gradient descent
111        normalize - whether to perform normalization
112
113        ...
114        gradient_values, t_acc_val, v_acc_val = [], [], []
115
116        if normalize:
117            X, v_X = self.normalization(X, v_X)
118
119        # Retrive the learning rate, maximum iteration, momentum (beta)
120        lr, max_iter, beta, epsilon = self.lr, self.max_iter, self.beta, self.
121        epsilon
122
123        # 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
129         # 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
135             momentum
136             w = w - lr * step
137
138             # predict once every 10 interations
139             if iteration % 10 == 0:
140                 t_y_pred = self.predict(X, w)
141                 t_acc = self.accu_eval(t_y_pred, y)
142                 v_y_pred = self.predict(v_X, w)
143                 v_acc = self.accu_eval(v_y_pred, v_y)
144
145             # record the next best value
146             if v_acc >= v_acc_max:
147                 v_acc_max = v_acc
148                 best_w = w
149                 self.marker = iteration # move the iteration marker
150
151             t_acc_val.append(t_acc)
152             v_acc_val.append(v_acc)
153
154             iteration = iteration + 1
155
156             # Uncomment to display the confusion matrix
157             # self.confusion_matrix(v_y_pred, v_y)
158             return gradient_values, t_acc_val, v_acc_val, best_w
159
160     def predict(self, X, w):
161         """
162             This function takes a set of data as input and outputs predicted
163             labels for the input points.
164         """
165         result = self.log_func(np.dot(X, w))
166         # the prediction result converted to binary
167         predict_bin = []
168         for i in result:
169             if i>=0.5:
170                 predict_bin.append(1)
171             else:
172                 predict_bin.append(0)
173         return predict_bin
174
175     def accu_eval(self, y_pred, y):
176         """
177             This function evaluates the models' accuracy.
178         """
179         count = 0
180         for i in range(len(y_pred)):
181             if y_pred[i] == y[i]:
182                 count = count + 1

```



```

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
245 , v_X, v_y, normalize=normalize)
246
247         accuracies.append(np.max(v_acc_val))
248         accuracies_train.append(np.max(t_acc_val))
249
250         # Uncomment this block to display the accuracy diagram
251         plt.figure()
252         plt.plot(t_acc_val, label = 'Training accuracy')
253         plt.plot(v_acc_val, label='Validation accuracy')
254         plt.axvline(log_reg.marker, color='r', label='Best Weights')
255         plt.xlabel('Iteration Number')
256         plt.ylabel('Accuracy')
257         plt.legend()
258         plt.show()
259         print("Learning Rate: " + str(log_reg.lr))
260         print("Average Accuracy: "+str(np.mean(accuracies)))
261
262         # Uncomment this block to display the gradient diagram
263         plt.figure()
264         plt.plot(gradient_val)
265         plt.xlabel('Iteration Number')
266         plt.ylabel('Gradient')
267         plt.show()
268         print("-----")
269
270     mean_acc = np.mean(accuracies)
271     mean_acc_train = np.mean(accuracies_train)
272     toc = time.time()
273     return mean_acc, mean_acc_train, toc-tic
274
275 def rise_order(self, data, order=3):
276     ret_val = data
277     for i in range(2, order + 1):
278         dataPowered = data.pow(i)
279         ret_val = ret_val.iloc[:, :-1]
280         ret_val = pd.concat([ret_val, dataPowered], axis=1)
281     return ret_val
282 """
283     Default Values
284
285     path = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/hepatitis.
286     csv"
287     dataset_name = "Hepatitis"
288     default_lr = 0.01
289     default_max_iter = 10000
290     default_epsilon = 5e-3
291     defulat_beta = 0.99
292
293     path = "/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_01/bankruptcy.csv"
294     dataset_name = "Bankruptcy"
295     default_lr = 0.1
296     default_max_iter = 25000
297     default_epsilon = 1e-3
298     defulat_beta = 0.99
299
300     unnorm_ordering = KFoldValidation(folds=10, path=path, lr=default_lr, max_iter=
301     default_max_iter, epsilon=default_epsilon, beta=defulat_beta)

```

```

299 normed_ordering = KFoldValidation(folds=10, path=path, lr=defult_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 defult_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=defult_lr,
    max_iter=default_max_iter, epsilon=default_epsilon, beta=defulat_beta)
353 rm_feature_test = KFoldValidation(folds=10, path=path, lr=defult_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.
359 k_fold_validation(normalize=True, inc_od=True)
360 mean_acc.append(mean_acc_temp)
361 proc_time.append(time_temp)
362 mean_acc_temp, mean_acc_train_temp, time_temp = rm_feature_test.
363 k_fold_validation(normalize=True, inc_od=True, rm_features=True)
364 mean_acc.append(mean_acc_temp)
365 proc_time.append(time_temp)
366
367 print("Without removing attribute 43 and 60\tmean accuracy is {} \tprocessing
368 time is {}\\nAfter removing those two attribute\tmean accuracy is {} \t
369 processing time is {}".format(mean_acc[0],proc_time[0],mean_acc[1],proc_time[1]
370 )))
```