

```

1  # How to replicate the results
2
3  This project is divided into three parts: *data preprocessing*, *logistic regression*
4  (fit, predict, and accu_eval), and *testing*.
5  To understand this project properly, please refer to the complete report while
6  performing the experiment.
7
8  ---
9
10 ## Team Members
11
12 1. Yi Zhu, 260716006
13 2. Fei Peng, 260712440
14 3. Yukai Zhang, 260710915
15
16 ---
17
18 ## Python Version
19
20 1. **Writing**
21     Google Colab `Python 3.6`
22 2. **Testing**
23     Google Colab `Python 3.6`
24
25 ---
26
27 ## Imports
28
29 1. numpy
30 2. pandas
31 3. matplotlib.pyplot
32 4. time
33 5. seaborn
34 6. from google.colab import files
35 7. from google.colab import drive
36
37 ---
38
39 ## Data Preprocessing
40
41 Please run the following file:
42
43 **Data_Preprocessing.ipynb**
44 This file is for *data preprocessing*. Important characteristics of features and the
45 distribution of classes could be found in the result of this file.
46
47 Instructions on how to read the output plots:
48
49 > The first two plots are the distribution of the two classes of the two provided
50 datasets: *Hepatitis* and *Bankruptcy*.
51
52 > The following histogram plots are the distribution of all the features of the two
53 provided datasets.
54
55 ---
56
57 ## Logistic Regression
58
59 Please run the following file:
60
61 **Logistic_Regression.ipynb**
62 This file contains the required functions: *fit*, *perdict*, and *accu_eval*, which are
63 contained in a Python class called LogisticRegression, as well as the *k-fold
64 validation* class (KFoldValidation). The hyperparameters and models used in this file
65 are chosen based on the findings in the testing file.
66
67 Instructions on how to read the output plots:

```

```

60
61 > The first block of output plots are: `accuracy vs. iteration number plots` and
    `gradient vs. iteration number plots` of *Hepatitis* dataset.
62
63 > The second block of output plots are: `accuracy vs. iteration number plots` and
    `gradient vs. iteration number plots` of *Bankruptcy* dataset.
64
65 ---
66
67 ## Testing
68
69 The testing part is divided into two subparts:
70
71 1. To test the hyperparameters of gradient descent algorithm - `hyperparameter testing`
72 2. To test the effect of z-score normalization and adding/removing features -
    `normalization and feature testing`
73
74 **Hyperparameter_Testing.ipynb**
75 This file aims to find the best hyperparameters (of gradient descent algorithm) for the
    model.
76
77 Instructions on how to read the output plots:
78
79 The following testing are performed for both datasets.
80
81 1. Learning rate testing
82 2. Stopping criteria testing
83     * Maximum Iteration testing
84     * Epsilon (threshold for gradient) testing
85 3. Beta (Momentum Gradient Descent Constant) testing
86
87 > For every hyperparameter test, the first plot is the `mean validation accuracy vs.
    hyperparameter`, and the second plot is `processing time vs. hyperparameter`.
88 > You could uncomment some lines of code to explore the result more comprehensively,
    detailed instructions could be found in the code file.
89
90 **Normalization_Feature_Testing.ipynb**
91 This file aims to find the effect of normalization, increasing features on the model.
92
93 Instructions on how to read the output plots:
94
95 The following testing are performed for both datasets.
96
97 1. *z*-score `normalization` testing
98 2. Adding more features testing
99 3. Removing features testing
100
101 > In the output plot, the difference between with and without normalization could be
    seen. Moreover, `the mean validation accuracy vs. added feature order` is also plotted.
102 > You could uncomment some lines of code to explore the result more comprehensively,
    detailed instructions could be found in the code file.

```