```python
 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/1cqiOtupPBH0ZZWkk2pkLA-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
   ></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 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
   ), dpi = 1200)
58         files.download("Distribution of the two classes {}.png".format(self.
   name))
59         plt.show()
```

```python
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 """
```