

```

1  # -*- coding: utf-8 -*-
2  """Data_Preprocessing.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1ALc0\_L0DWWHcrm6Qcxv-KH4TzpSpAlak
8
9  <center><h1>Mini Project 3 - Convolutional Neural Network</h1>
10 <h3>Data Preprocessing</h3>
11 <h4>This file performs some of the operations on Data Preprocessing and
    Analysis.</h4></center>
12
13 <h3>Team Members:</h3>
14 <center>
15 Yi Zhu, 260716006<br>
16 Fei Peng, 260712440<br>
17 Yukai Zhang, 260710915
18 </center>
19
20 # Importations
21 """
22
23 from google.colab import drive
24 drive.mount('/content/drive')
25
26 # make path = './' in-case you are running this locally
27 path = '/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_03/'
28
29 import numpy as np
30 import pandas as pd
31 import matplotlib.pyplot as plt
32 import seaborn as sns
33 from scipy import stats
34 from google.colab import files
35 from sklearn.preprocessing import LabelEncoder
36 from scipy.stats import entropy
37
38 """# Data Preprocessing"""
39
40 dataset = pd.read_csv(path+"TrainLabels.csv")
41 # reddit_test = pd.read_csv(path+"test.csv")
42
43 y = dataset['class']
44
45 class Data_Processing:
46     def __init__(self, data, name='New Data'):
47         self.data = data
48         self.name = name
49
50     def show_y_dist(self, ydata):
51         plt.figure(figsize=(8,4))
52         plt.subplot(111), sns.countplot(x='class', data=ydata)
53         plt.title('Distribution of Class')
54         plt.savefig("Distribution of Class.png", dpi=1200)
55         files.download("Distribution of Class.png")
56         plt.show()
57
58 data_analysis = Data_Processing(dataset.values, 'TrainLabels.csv')
59 data_analysis.show_y_dist(dataset)
60
61 # calculate the data entropy
62 le = LabelEncoder() # encoder for classes

```

```
63 le.fit(y)
64 y_label = le.transform(y)
65 n_k = len(le.classes_)
66 N = len(y)
67 theta_k = np.zeros(n_k) # probability of class k
68 # compute theta values
69 for k in range(n_k):
70     count_k = (y_label==k).sum()
71     theta_k[k] = count_k / N
72
73 print("Data entropy is", entropy(theta_k, base=2))
```