

```

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/1bRmLzkFmGP7xgU0UI6iv\_p6dW1hX6R6P
8
9  <center><h1>Mini Project 2 - Bernoulli Naïve Bayes</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_02/'
28
29 import numpy as np
30 import pandas as pd
31 import matplotlib.pyplot as plt
32 import seaborn as sns
33 import random
34 from scipy import stats
35 from google.colab import files
36 from time import time
37
38 from sklearn.model_selection import train_test_split
39 from sklearn.preprocessing import Normalizer
40 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
41 from sklearn.feature_extraction import text
42 from sklearn import metrics
43 from sklearn.model_selection import GridSearchCV, cross_val_score, KFold
44 from sklearn.pipeline import make_pipeline
45
46 !pip install nltk
47 import nltk
48 nltk.download('punkt')
49 nltk.download('wordnet')
50 nltk.download('averaged_perceptron_tagger')
51
52 from nltk.stem import PorterStemmer
53 from nltk import word_tokenize
54 from nltk import word_tokenize
55 from nltk.stem import WordNetLemmatizer
56 from nltk.corpus import wordnet
57
58 from sklearn.linear_model import LogisticRegression
59 from sklearn.naive_bayes import MultinomialNB
60 from sklearn import svm
61 from sklearn.ensemble import RandomForestClassifier
62 from sklearn.tree import DecisionTreeClassifier

```

```

63 from sklearn.ensemble import AdaBoostClassifier
64 from sklearn.neighbors import KNeighborsClassifier
65 from sklearn.neural_network import MLPClassifier
66
67 """# Data Preprocessing"""
68
69 reddit_dataset = pd.read_csv(path+"train.csv")
70 reddit_test = pd.read_csv(path+"test.csv")
71
72 X = reddit_dataset['body']
73 y = reddit_dataset['subreddit']
74
75 class Data_Processing:
76     def __init__(self, data, name='New Data'):
77         self.data = data
78         self.name = name
79
80     def show_y_dist(self, ydata):
81         plt.figure(figsize=(8,4))
82         plt.subplot(111), sns.countplot(x='subreddit', data=ydata)
83         plt.title('Distribution of Subreddit in {}'.format(self.name))
84         plt.savefig("Distribution of Subreddit in {}.png".format(self.name),
85 dpi = 1200)
86         files.download("Distribution of Subreddit in {}.png".format(self.name
87 ))
88         plt.show()
89
90 data_analysis = Data_Processing(reddit_dataset.values, 'train.csv')
91 data_analysis.show_y_dist(reddit_dataset)
92
93 # calculate the data entropy
94 from sklearn.preprocessing import LabelEncoder
95 le = LabelEncoder() # encoder for classes
96 le.fit(y)
97 y_label = le.transform(y)
98 n_k = len(le.classes_)
99 N = len(y)
100 theta_k = np.zeros(n_k) # probability of class k
101 # compute theta values
102 for k in range(n_k):
103     count_k = (y_label==k).sum()
104     theta_k[k] = count_k / N
105
106 from scipy.stats import entropy
107 print("Data entropy is", entropy(theta_k, base=2))
108
109 """# Define Vectorizer
110 (To vectorize the text-based data to numerical features)
111
112 1. CountVectorizer
113 1) Use "CountVectorizer" to transform text data to feature vectors.
114 2) Normalize your feature vectors
115 """
116
117 def count_vectorizer(X_train):
118     vectorizer = CountVectorizer()
119     vectors_train = vectorizer.fit_transform(X_train)
120     return vectors_train
121
122 """2. CountVectorizer with stop word
123 1) Use "CountVectorizer" with stop word to transform text data to vector.
124 2) Normalize your feature vectors
125 """

```

```

124
125 def count_vec_with_sw(X_train):
126     stop_words = text.ENGLISH_STOP_WORDS
127     vectorizer = CountVectorizer(stop_words=stop_words)
128     vectors_train_stop = vectorizer.fit_transform(X_train)
129     return vectors_train_stop
130
131 """3. TF-IDF
132 1) use "TfidfVectorizer" to weight features based on your train set.
133 2) Normalize your feature vectors
134 """
135
136 def tfidf_vectorizer(X_train):
137     tf_idf_vectorizer = TfidfVectorizer()
138     vectors_train_idf = tf_idf_vectorizer.fit_transform(X_train)
139     return vectors_train_idf
140
141 """4. CountVectorizer with stem tokenizer
142 1) Use "StemTokenizer" to transform text data to vector.
143 2) Normalize your feature vectors
144 """
145
146 class StemTokenizer:
147     def __init__(self):
148         self.wnl =PorterStemmer()
149     def __call__(self, doc):
150         return [self.wnl.stem(t) for t in word_tokenize(doc) if t.isalpha()]
151
152
153 def count_vec_stem(X_train):
154     vectorizer = CountVectorizer(tokenizer=StemTokenizer())
155     vectors_train_stem = vectorizer.fit_transform(X_train)
156     return vectors_train_stem
157
158 """5. CountVectorizer with lemma tokenizer
159 1) Use "LemmaTokenizer" to transform text data to vector.
160 2) Normalize your feature vectors
161 """
162
163 def get_wordnet_pos(word):
164     """Map POS tag to first character lemmatize() accepts"""
165     tag = nltk.pos_tag([word])[0][1][0].upper()
166     tag_dict = {"J": wordnet.ADJ,
167                 "N": wordnet.NOUN,
168                 "V": wordnet.VERB,
169                 "R": wordnet.ADV}
170     return tag_dict.get(tag, wordnet.NOUN)
171
172
173 class LemmaTokenizer:
174     def __init__(self):
175         self.wnl = WordNetLemmatizer()
176     def __call__(self, doc):
177         return [self.wnl.lemmatize(t,pos =get_wordnet_pos(t)) for t in
178                 word_tokenize(doc) if t.isalpha()]
179
180 def count_vec_lemma(X_train):
181     vectorizer = CountVectorizer(tokenizer=LemmaTokenizer())
182     vectors_train_lemma = vectorizer.fit_transform(X_train)
183     return vectors_train_lemma
184
185 """# Measure the time required for each vectorizer to perform vectorization

```

```
186
187 ### 1. CountVectorizer
188 """
189
190 tic = time()
191 X_vec = count_vectorizer(X)
192 print("\t\t- Count Vectorizer - \nfeature number: ", X_vec.shape[1], "\t\tTime
    spent: ", time()-tic, "s.")
193
194 ##### 2. CountVectorizer with stop word"""
195
196 tic = time()
197 X_vec = count_vec_with_sw(X)
198 print("\t\t- Count Vectorizer with stop word - \nfeature number: ", X_vec.
    shape[1], "\t\tTime spent: ", time()-tic, "s.")
199
200 ##### 3. TF-IDF"""
201
202 tic = time()
203 X_vec = tfidf_vectorizer(X)
204 print("\t\t- TF-IDF Vectorizer - \nfeature number: ", X_vec.shape[1], "\t\t
    Time spent: ", time()-tic, "s.")
205
206 ##### 4. CountVectorizer with stem tokenizer"""
207
208 tic = time()
209 X_vec = count_vec_stem(X)
210 print("\t\t- CountVectorizer with stem tokenizer - \nfeature number: ", X_vec.
    shape[1], "\t\tTime spent: ", time()-tic, "s.")
211
212 ##### 5. CountVectorizer with lemma tokenizer"""
213
214 tic = time()
215 X_vec = count_vec_lemma(X)
216 print("\t\t- CountVectorizer with lemma tokenizer - \nfeature number: ", X_vec
    .shape[1], "\t\tTime spent: ", time()-tic, "s.")
```