

```

1  # -*- coding: utf-8 -*-
2  """Export CSV.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1uNTFaEcHjLGUkGUojaKE8gCcied9FaLQ
8
9  <center><h1>Mini Project 2 - Bernoulli Naïve Bayes</h1>
10 <h4>This file is for file exporting.</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 # make path = './' in-case you are running this locally
24 path = '/content/drive/My Drive/ECSE_551_F_2020/Mini_Project_02/'
25
26 import numpy as np
27 import pandas as pd
28 import matplotlib.pyplot as plt
29
30 from sklearn.model_selection import train_test_split
31 from sklearn.preprocessing import Normalizer
32 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
33 from sklearn.feature_extraction import text
34 from sklearn import metrics
35 from sklearn.model_selection import GridSearchCV, cross_val_score, KFold
36 from sklearn.pipeline import make_pipeline
37
38 !pip install nltk
39 import nltk
40 nltk.download('punkt')
41 nltk.download('wordnet')
42 nltk.download('averaged_perceptron_tagger')
43
44 from nltk.stem import PorterStemmer
45 from nltk import word_tokenize
46 from nltk import word_tokenize
47 from nltk.stem import WordNetLemmatizer
48 from nltk.corpus import wordnet
49
50 """Additional classifiers:
51 1. Logistic Regression
52 2. Multinomial Naïve Bayes
53 3. Support Vector Machine
54 4. Random Forest
55 5. Decision Tree
56 6. Ada Boost
57 7. k-Neighbors
58 8. Neural Network
59 """
60
61 from sklearn.linear_model import LogisticRegression
62 from sklearn.naive_bayes import MultinomialNB
63 from sklearn import svm

```

```

64 from sklearn.ensemble import RandomForestClassifier
65 from sklearn.tree import DecisionTreeClassifier
66 from sklearn.ensemble import AdaBoostClassifier
67 from sklearn.neighbors import KNeighborsClassifier
68 from sklearn.neural_network import MLPClassifier
69
70 reddit_dataset = pd.read_csv(path+"train.csv")
71 reddit_test = pd.read_csv(path+"test.csv")
72
73 X = reddit_dataset['body']
74 y = reddit_dataset['subreddit']
75
76 """# Define Vectorizer
77 ### (To vectorize the text-based data to numerical features)
78
79 1. CountVectorizer
80 1) Use "CountVectorizer" to transform text data to feature vectors.
81 2) Normalize your feature vectors
82 """
83
84 def count_vectorizer(X_train, X_test):
85     vectorizer = CountVectorizer()
86     vectors_train = vectorizer.fit_transform(X_train)
87     vectors_test = vectorizer.transform(X_test)
88
89     # z-score normalization
90     normalizer_train = Normalizer().fit(X=vectors_train)
91     vectors_train = normalizer_train.transform(vectors_train)
92     vectors_test = normalizer_train.transform(vectors_test)
93
94     return vectors_train, vectors_test
95
96 """2. CountVectorizer with stop word
97 1) Use "CountVectorizer" with stop word to transform text data to vector.
98 2) Normalize your feature vectors
99 """
100
101 def count_vec_with_sw(X_train, X_test):
102     stop_words = text.ENGLISH_STOP_WORDS
103     vectorizer = CountVectorizer(stop_words=stop_words)
104     vectors_train_stop = vectorizer.fit_transform(X_train)
105     vectors_test_stop = vectorizer.transform(X_test)
106
107     # z-score normalization
108     normalizer_train = Normalizer().fit(X=vectors_train_stop)
109     vectors_train_stop = normalizer_train.transform(vectors_train_stop)
110     vectors_test_stop = normalizer_train.transform(vectors_test_stop)
111
112     return vectors_train_stop, vectors_test_stop
113
114 """3. TF-IDF
115 1) use "TfidfVectorizer" to weight features based on your train set.
116 2) Normalize your feature vectors
117 """
118
119 def tfidf_vectorizer(X_train, X_test):
120     tf_idf_vectorizer = TfidfVectorizer()
121     vectors_train_idf = tf_idf_vectorizer.fit_transform(X_train)
122     vectors_test_idf = tf_idf_vectorizer.transform(X_test)
123
124     # z-score normalization
125     normalizer_train = Normalizer().fit(X=vectors_train_idf)
126     vectors_train_idf = normalizer_train.transform(vectors_train_idf)

```

```

127     vectors_test_idf = normalizer_train.transform(vectors_test_idf)
128
129     return vectors_train_idf, vectors_test_idf
130
131 """4. CountVectorizer with stem tokenizer
132 1) Use "StemTokenizer" to transform text data to vector.
133 2) Normalize your feature vectors
134 """
135
136 class StemTokenizer:
137     def __init__(self):
138         self.wnl = PorterStemmer()
139     def __call__(self, doc):
140         return [self.wnl.stem(t) for t in word_tokenize(doc) if t.isalpha()]
141
142
143 def count_vec_stem(X_train, X_test):
144     vectorizer = CountVectorizer(tokenizer=StemTokenizer())
145     vectors_train_stem = vectorizer.fit_transform(X_train)
146     vectors_test_stem = vectorizer.transform(X_test)
147
148     # z-score normalization
149     normalizer_train = Normalizer().fit(X=vectors_train_stem)
150     vectors_train_stem = normalizer_train.transform(vectors_train_stem)
151     vectors_test_stem = normalizer_train.transform(vectors_test_stem)
152
153     return vectors_train_stem, vectors_test_stem
154
155 """5. CountVectorizer with lemma tokenizer
156 1) Use "LemmaTokenizer" to transform text data to vector.
157 2) Normalize your feature vectors
158 """
159
160 def get_wordnet_pos(word):
161     """Map POS tag to first character lemmatize() accepts"""
162     tag = nltk.pos_tag([word])[0][1][0].upper()
163     tag_dict = {"J": wordnet.ADJ,
164                 "N": wordnet.NOUN,
165                 "V": wordnet.VERB,
166                 "R": wordnet.ADV}
167     return tag_dict.get(tag, wordnet.NOUN)
168
169
170 class LemmaTokenizer:
171     def __init__(self):
172         self.wnl = WordNetLemmatizer()
173     def __call__(self, doc):
174         return [self.wnl.lemmatize(t, pos=get_wordnet_pos(t)) for t in
175 word_tokenize(doc) if t.isalpha()]
176
177
178 def count_vec_lemma(X_train, X_test):
179     vectorizer = CountVectorizer(tokenizer=LemmaTokenizer())
180     vectors_train_lemma = vectorizer.fit_transform(X_train)
181     vectors_test_lemma = vectorizer.transform(X_test)
182
183     # z-score normalization
184     normalizer_train = Normalizer().fit(X=vectors_train_lemma)
185     vectors_train_lemma = normalizer_train.transform(vectors_train_lemma)
186     vectors_test_lemma = normalizer_train.transform(vectors_test_lemma)
187
188     return vectors_train_lemma, vectors_test_lemma

```

```
189 """# Export csv"""
190
191 # test set id
192 X_id = reddit_test['id']
193 # test set features
194 X_test = reddit_test['body']
195
196 # vectorize the training and testing data
197 vectors_train, vectors_test = tfidf_vectorizer(X, X_test)
198 # perform MLP classification
199 clf = MLPClassifier(random_state=0, max_iter=1000, learning_rate="adaptive",
    learning_rate_init=0.0001).fit(vectors_train, y)
200
201 # predict the result
202 y_test = clf.predict(vectors_test)
203
204 # put the result into a pandas dataframe
205 result = {'id': X_id, 'subreddit': y_test}
206 df = pd.DataFrame(data=result)
207
208 # export to csv
209 df.to_csv('result.csv', index=False)
210
211 # download csv
212 from google.colab import files
213 files.download('result.csv')
```