

```

1  # -*- coding: utf-8 -*-
2  """Stacking Classifier.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/10KCCLcD9Qf2Z59egC0iTBdwc2xkamhsM
8
9  <center><h1>Mini Project 2 - Bernoulli Naïve Bayes</h1>
10 <h4>The hyperparameters and models used in this file are chosen based on the
    findings in the testing 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 # 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 """1. CountVectorizer
77 1) Use "CountVectorizer" to transform text data to feature vectors.
78 2) Normalize your feature vectors
79 """
80
81 def count_vectorizer(X_train, X_test):
82     vectorizer = CountVectorizer()
83     vectors_train = vectorizer.fit_transform(X_train)
84     vectors_test = vectorizer.transform(X_test)
85
86     normalizer_train = Normalizer().fit(X=vectors_train)
87     vectors_train = normalizer_train.transform(vectors_train)
88     vectors_test = normalizer_train.transform(vectors_test)
89
90     return vectors_train, vectors_test
91
92 """2. CountVectorizer with stop word
93 1) Use "CountVectorizer" with stop word to transform text data to vector.
94 2) Normalize your feature vectors
95 """
96
97 def count_vec_with_sw(X_train, X_test):
98     stop_words = text.ENGLISH_STOP_WORDS
99     vectorizer = CountVectorizer(stop_words=stop_words)
100     vectors_train_stop = vectorizer.fit_transform(X_train)
101     vectors_test_stop = vectorizer.transform(X_test)
102
103     normalizer_train = Normalizer().fit(X=vectors_train_stop)
104     vectors_train_stop = normalizer_train.transform(vectors_train_stop)
105     vectors_test_stop = normalizer_train.transform(vectors_test_stop)
106
107     return vectors_train_stop, vectors_test_stop
108
109 """3. TF-IDF
110 1) use "TfidfVectorizer" to weight features based on your train set.
111 2) Normalize your feature vectors
112 """
113
114 def tfidf_vectorizer(X_train, X_test, binary=False):
115     stop_words = text.ENGLISH_STOP_WORDS
116     tf_idf_vectorizer = TfidfVectorizer(binary=binary, stop_words=stop_words)
117     vectors_train_idf = tf_idf_vectorizer.fit_transform(X_train)
118     vectors_test_idf = tf_idf_vectorizer.transform(X_test)
119
120     normalizer_train = Normalizer().fit(X=vectors_train_idf)
121     vectors_train_idf = normalizer_train.transform(vectors_train_idf)
122     vectors_test_idf = normalizer_train.transform(vectors_test_idf)
123
124     return vectors_train_idf, vectors_test_idf
125

```

```

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

```

```

188 accuracies = []
189 kf = KFold(n_splits=5, shuffle=True)
190 for train_index, test_index in kf.split(X):
191     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
test_index], binary=True)
192     estimators = [
193         ('mlp', MLPClassifier(max_iter=1000, learning_rate="adaptive",
learning_rate_init=0.0001)),
194         ('svc', svm.SVC(kernel='linear', gamma='auto', C=1, probability=True
)),
195         ('lr', LogisticRegression(C=40.0, max_iter=1000))
196     ]
197     clf = StackingClassifier(
198         estimators=estimators, final_estimator=LogisticRegression()
199     )
200     clf.fit(vectors_train, y[train_index])
201     y_test = clf.predict(vectors_test)
202     accuracies.append(metrics.accuracy_score(y[test_index], y_test))
203     print(accuracies[-1])
204 print("Average accuracy of Stacking Classification = {}".format(np.mean(
accuracies)));
205 toc = time()
206 print("Time spent for Stacking Classification (with 5 fold validation) = {}".
format(toc - tic))
207
208 """## 10. Voting Classifier"""
209
210 from sklearn.ensemble import VotingClassifier
211 tic = time()
212 accuracies = []
213 kf = KFold(n_splits=5, shuffle=True)
214 for train_index, test_index in kf.split(X):
215     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
test_index], binary=True)
216     clf1 = MLPClassifier(max_iter=1000, learning_rate="adaptive",
learning_rate_init=0.0001)
217     clf2 = LogisticRegression(C=40.0, max_iter=1000)
218     clf3 = svm.SVC(kernel='linear', gamma='auto', C=1, probability=True)
219     eclf = VotingClassifier(estimators=[('mlp', clf1), ('lr', clf2), ('svc',
clf3)],
220                             voting='soft', weights=[1,1,1])
221     eclf = eclf.fit(vectors_train, y[train_index])
222     y_test = eclf.predict(vectors_test)
223     accuracies.append(metrics.accuracy_score(y[test_index], y_test))
224     print(accuracies[-1])
225 print("Average accuracy of Volting Classification = {}".format(np.mean(
accuracies)));
226 toc = time()
227 print("Time spent for Volting Classification (with 5 fold validation) = {}".
format(toc - tic))

```