

```

1 # -*- coding: utf-8 -*-
2 """Additional Classifiers.ipynb
3
4 Automatically generated by Colaboratory.
5
6 Original file is located at
7     https://colab.research.google.com/drive/1jQsRTZihNlPC99pBzRYqvTbowAY-tqi1
8
9 <center><h1>Mini Project 2 - Bernoulli Naïve Bayes</h1>
10 <h4>This file records the accuracies of the combinations of 8 classifiers and 5
vectorizers.</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     normalizer_train = Normalizer().fit(X=vectors_train)
90     vectors_train = normalizer_train.transform(vectors_train)
91     vectors_test = normalizer_train.transform(vectors_test)
92
93     return vectors_train, vectors_test
94
95 """2. CountVectorizer with stop word
96 1) Use "CountVectorizer" with stop word to transform text data to vector.
97 2) Normalize your feature vectors
98 """
99
100 def count_vec_with_sw(X_train, X_test):
101     stop_words = text.ENGLISH_STOP_WORDS
102     vectorizer = CountVectorizer(stop_words=stop_words)
103     vectors_train_stop = vectorizer.fit_transform(X_train)
104     vectors_test_stop = vectorizer.transform(X_test)
105
106     normalizer_train = Normalizer().fit(X=vectors_train_stop)
107     vectors_train_stop= normalizer_train.transform(vectors_train_stop)
108     vectors_test_stop = normalizer_train.transform(vectors_test_stop)
109
110     return vectors_train_stop, vectors_test_stop
111
112 """3. TF-IDF
113 1) use "TfidfVectorizer" to weight features based on your train set.
114 2) Normalize your feature vectors
115 """
116
117 def tfidf_vectorizer(X_train, X_test):
118     tf_idf_vectorizer = TfidfVectorizer()
119     vectors_train_idf = tf_idf_vectorizer.fit_transform(X_train)
120     vectors_test_idf = tf_idf_vectorizer.transform(X_test)
121
122     normalizer_train = Normalizer().fit(X=vectors_train_idf)
123     vectors_train_idf= normalizer_train.transform(vectors_train_idf)
124     vectors_test_idf = normalizer_train.transform(vectors_test_idf)
125

```

```

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

```

```

188 ### 1. CountVectorizer
189 """
190
191 accuracies = []
192 clf = LogisticRegression(C=1.0, max_iter=1000, random_state=0)
193 kf = KFold(n_splits=5, shuffle=True)
194 for train_index, test_index in kf.split(X):
195     vectors_train, vectors_test = count_vectorizer(X[train_index], X[
196         test_index])
197     clf.fit(vectors_train, y[train_index])
198     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
199         vectors_test)))
200
201 print(np.mean(accuracies))
202
203 """## 2. CountVectorizer with stop word"""
204
205 accuracies = []
206 clf = LogisticRegression(C=1.0, max_iter=1000, random_state=0)
207 kf = KFold(n_splits=5, shuffle=True)
208 for train_index, test_index in kf.split(X):
209     vectors_train, vectors_test = count_vec_with_sw(X[train_index], X[
210         test_index])
211     clf.fit(vectors_train, y[train_index])
212     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
213         vectors_test)))
214
215 print(np.mean(accuracies))
216
217 """## 3. TF-IDF"""
218
219 accuracies = []
220 clf = LogisticRegression(C=40.0, max_iter=1000, random_state=0)
221 kf = KFold(n_splits=5, shuffle=True)
222 for train_index, test_index in kf.split(X):
223     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
224         test_index])
225     clf.fit(vectors_train, y[train_index])
226     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
227         vectors_test)))
228
229 print(np.mean(accuracies))
230
231 """accuracy vs. hyperparameter C"""
232
233 c_vecs = np.logspace(0, 2, 5)
234 acc = []
235 for c_vec in c_vecs:
236     accuracies = []
237     clf = LogisticRegression(C=c_vec, max_iter=1000, random_state=0)
238     kf = KFold(n_splits=5, shuffle=True)
239     for train_index, test_index in kf.split(X):
240         vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
241             test_index])
242         clf.fit(vectors_train, y[train_index])
243         accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
244             vectors_test)))
245     acc.append(np.mean(accuracies))
246
247 plt.xlabel('C')
248 plt.ylabel('Accuracy')
249 plt.plot(c_vecs, acc)

```

```

243
244 """## 4. CountVectorizer with stem tokenizer"""
245
246 accuracies = []
247 clf = LogisticRegression(C=1.0, max_iter=1000, random_state=0)
248 kf = KFold(n_splits=5, shuffle=True)
249 for train_index, test_index in kf.split(X):
250     vectors_train, vectors_test = count_vec_stem(X[train_index], X[test_index])
251     clf.fit(vectors_train, y[train_index])
252     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
253         vectors_test)))
254 print(np.mean(accuracies))
255
256 """## 5. CountVectorizer with lemma tokenizer"""
257
258 accuracies = []
259 clf = LogisticRegression(C=1.0, max_iter=1000, random_state=0)
260 kf = KFold(n_splits=5, shuffle=True)
261 for train_index, test_index in kf.split(X):
262     vectors_train, vectors_test = count_vec_lemma(X[train_index], X[test_index])
263     clf.fit(vectors_train, y[train_index])
264     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
265         vectors_test)))
266 print(np.mean(accuracies))
267
268 """## 2. Multinomial Naïve Bayes
269
270 ### 1. CountVectorizer
271 """
272
273 accuracies = []
274 clf = MultinomialNB()
275 kf = KFold(n_splits=5, shuffle=True)
276 for train_index, test_index in kf.split(X):
277     vectors_train, vectors_test = count_vectorizer(X[train_index], X[
278         test_index])
279     clf.fit(vectors_train, y[train_index])
280     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
281         vectors_test)))
282 print(np.mean(accuracies))
283
284 """## 2. CountVectorizer with stop word"""
285
286 accuracies = []
287 clf = MultinomialNB()
288 kf = KFold(n_splits=5, shuffle=True)
289 for train_index, test_index in kf.split(X):
290     vectors_train, vectors_test = count_vec_with_sw(X[train_index], X[
291         test_index])
292     clf.fit(vectors_train, y[train_index])
293     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
294         vectors_test)))
295
296 """## 3. TF-IDF"""
297 accuracies = []

```

```

298 clf = MultinomialNB()
299 kf = KFold(n_splits=5, shuffle=True)
300 for train_index, test_index in kf.split(X):
301     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
302         test_index])
303     clf.fit(vectors_train, y[train_index])
304     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
305         vectors_test)))
306
307 """## 4. CountVectorizer with stem tokenizer"""
308
309 accuracies = []
310 clf = MultinomialNB()
311 kf = KFold(n_splits=5, shuffle=True)
312 for train_index, test_index in kf.split(X):
313     vectors_train, vectors_test = count_vec_stem(X[train_index], X[test_index
314         ])
314     clf.fit(vectors_train, y[train_index])
315     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
316         vectors_test)))
317
318 print(np.mean(accuracies))
319 """## 5. CountVectorizer with lemma tokenizer"""
320
321 accuracies = []
322 clf = MultinomialNB()
323 kf = KFold(n_splits=5, shuffle=True)
324 for train_index, test_index in kf.split(X):
325     vectors_train, vectors_test = count_vec_lemma(X[train_index], X[test_index
326         ])
326     clf.fit(vectors_train, y[train_index])
327     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
328         vectors_test)))
329
330 print(np.mean(accuracies))
331 """## 3. Support Vector Machine
332
333 ## 1. CountVectorizer
334 """
335
336 accuracies = []
337 clf = svm.SVC(kernel='linear', gamma='auto', C=1)
338 kf = KFold(n_splits=5, shuffle=True)
339 for train_index, test_index in kf.split(X):
340     vectors_train, vectors_test = count_vectorizer(X[train_index], X[
341         test_index])
341     clf.fit(vectors_train, y[train_index])
342     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
343         vectors_test)))
344
345 print(np.mean(accuracies))
346 """## 2. CountVectorizer with stop word"""
347
348 accuracies = []
349 clf = svm.SVC(kernel='linear', gamma='auto', C=1)
350 kf = KFold(n_splits=5, shuffle=True)
351 for train_index, test_index in kf.split(X):
352     vectors_train, vectors_test = count_vec_with_sw(X[train_index], X[

```

```

352 test_index])
353     clf.fit(vectors_train, y[train_index])
354     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
355         vectors_test)))
356 print(np.mean(accuracies))
357
358 """## 3. TF-IDF
359
360 Linear
361 """
362
363 accuracies = []
364 clf = svm.SVC(kernel='linear', gamma='auto', C=1)
365 kf = KFold(n_splits=5, shuffle=True)
366 for train_index, test_index in kf.split(X):
367     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
368         test_index])
369     clf.fit(vectors_train, y[train_index])
370     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
371         vectors_test)))
372
373 print(np.mean(accuracies))
374
375 """Non-Linear"""
376
377 accuracies = []
378 clf = svm.SVC(gamma='auto', C=1)
379 kf = KFold(n_splits=5, shuffle=True)
380 for train_index, test_index in kf.split(X):
381     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
382         test_index])
383     clf.fit(vectors_train, y[train_index])
384     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
385         vectors_test)))
386
387 print(np.mean(accuracies))
388
389 """## 4. CountVectorizer with stem tokenizer"""
390
391 accuracies = []
392 clf = svm.SVC(kernel='linear', gamma='auto', C=1)
393 kf = KFold(n_splits=5, shuffle=True)
394 for train_index, test_index in kf.split(X):
395     vectors_train, vectors_test = count_vec_stem(X[train_index], X[test_index]
396     ])
397     clf.fit(vectors_train, y[train_index])
398     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
399         vectors_test)))
400
401 print(np.mean(accuracies))
402
403 """## 5. CountVectorizer with lemma tokenizer"""
404
405 accuracies = []
406 clf = svm.SVC(kernel='linear', gamma='auto', C=1)
407 kf = KFold(n_splits=5, shuffle=True)
408 for train_index, test_index in kf.split(X):
409     vectors_train, vectors_test = count_vec_lemma(X[train_index], X[test_index]
410     ])
411     clf.fit(vectors_train, y[train_index])
412     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
413         vectors_test)))

```

```

406
407 print(np.mean(accuracies))
408
409 """## 4. Random Forest
410
411 ### 1. CountVectorizer
412 """
413
414 accuracies = []
415 clf = RandomForestClassifier(max_depth=2, random_state=0)
416 kf = KFold(n_splits=5, shuffle=True)
417 for train_index, test_index in kf.split(X):
418     vectors_train, vectors_test = count_vectorizer(X[train_index], X[
419         test_index])
420     clf.fit(vectors_train, y[train_index])
421     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
422         vectors_test)))
423
424 print(np.mean(accuracies))
425
426 """## 2. CountVectorizer with stop word"""
427
428 accuracies = []
429 clf = RandomForestClassifier(max_depth=2, random_state=0)
430 kf = KFold(n_splits=5, shuffle=True)
431 for train_index, test_index in kf.split(X):
432     vectors_train, vectors_test = count_vec_with_sw(X[train_index], X[
433         test_index])
434     clf.fit(vectors_train, y[train_index])
435     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
436         vectors_test)))
437
438 print(np.mean(accuracies))
439
440 """## 3. TF-IDF"""
441
442 accuracies = []
443 clf = RandomForestClassifier(max_depth=2, random_state=0)
444 kf = KFold(n_splits=5, shuffle=True)
445 for train_index, test_index in kf.split(X):
446     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
447         test_index])
448     clf.fit(vectors_train, y[train_index])
449     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
450         vectors_test)))
451
452 print(np.mean(accuracies))
453
454 """## 4. CountVectorizer with stem tokenizer"""
455
456 accuracies = []
457 clf = RandomForestClassifier(max_depth=2, random_state=0)
458 kf = KFold(n_splits=5, shuffle=True)
459 for train_index, test_index in kf.split(X):
460     vectors_train, vectors_test = count_vec_stem(X[train_index], X[test_index
461         ])
462     clf.fit(vectors_train, y[train_index])
463     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
464         vectors_test)))
465
466 print(np.mean(accuracies))
467
468 """## 5. CountVectorizer with lemma tokenizer"""

```

```

461
462 accuracies = []
463 clf = RandomForestClassifier(max_depth=2, random_state=0)
464 kf = KFold(n_splits=5, shuffle=True)
465 for train_index, test_index in kf.split(X):
466     vectors_train, vectors_test = count_vec_lemma(X[train_index], X[test_index])
467     clf.fit(vectors_train, y[train_index])
468     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
469         vectors_test)))
470 print(np.mean(accuracies))
471
472 """## 5. Decision Tree
473
474 ### 1. CountVectorizer
475 """
476
477 accuracies = []
478 clf = DecisionTreeClassifier(random_state=0)
479 kf = KFold(n_splits=5, shuffle=True)
480 for train_index, test_index in kf.split(X):
481     vectors_train, vectors_test = count_vectorizer(X[train_index], X[
482         test_index])
483     clf.fit(vectors_train, y[train_index])
484     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
485         vectors_test)))
486 print(np.mean(accuracies))
487 """
488 """## 2. CountVectorizer with stop word"""
489
490 accuracies = []
491 clf = DecisionTreeClassifier(random_state=0)
492 kf = KFold(n_splits=5, shuffle=True)
493 for train_index, test_index in kf.split(X):
494     vectors_train, vectors_test = count_vec_with_sw(X[train_index], X[
495         test_index])
496     clf.fit(vectors_train, y[train_index])
497     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
498         vectors_test)))
499 """
500 """## 3. TF-IDF"""
501
502 accuracies = []
503 clf = DecisionTreeClassifier(random_state=0)
504 kf = KFold(n_splits=5, shuffle=True)
505 for train_index, test_index in kf.split(X):
506     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
507         test_index])
508     clf.fit(vectors_train, y[train_index])
509     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
510         vectors_test)))
511 """
512 """## 4. CountVectorizer with stem tokenizer"""
513
514 accuracies = []
515 clf = DecisionTreeClassifier(random_state=0)
516 kf = KFold(n_splits=5, shuffle=True)

```

```

516 for train_index, test_index in kf.split(X):
517     vectors_train, vectors_test = count_vec_stem(X[train_index], X[test_index]
518     ])
518     clf.fit(vectors_train, y[train_index])
519     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
519         vectors_test)))
520
521 print(np.mean(accuracies))
522
523 """## 5. CountVectorizer with lemma tokenizer"""
524
525 accuracies = []
526 clf = DecisionTreeClassifier(random_state=0)
527 kf = KFold(n_splits=5, shuffle=True)
528 for train_index, test_index in kf.split(X):
529     vectors_train, vectors_test = count_vec_lemma(X[train_index], X[test_index]
529     ])
530     clf.fit(vectors_train, y[train_index])
531     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
531         vectors_test)))
532
533 print(np.mean(accuracies))
534
535 """## 6. Ada Boost
536
537 ### 1. CountVectorizer
538 """
539
540 accuracies = []
541 clf = AdaBoostClassifier(n_estimators=100, learning_rate=0.5, random_state=0)
542 kf = KFold(n_splits=5, shuffle=True)
543 for train_index, test_index in kf.split(X):
544     vectors_train, vectors_test = count_vectorizer(X[train_index], X[
544         test_index])
545     clf.fit(vectors_train, y[train_index])
546     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
546         vectors_test)))
547
548 print(np.mean(accuracies))
549
550 """## 2. CountVectorizer with stop word"""
551
552 accuracies = []
553 clf = AdaBoostClassifier(n_estimators=100, learning_rate=0.5, random_state=0)
554 kf = KFold(n_splits=5, shuffle=True)
555 for train_index, test_index in kf.split(X):
556     vectors_train, vectors_test = count_vec_with_sw(X[train_index], X[
556         test_index])
557     clf.fit(vectors_train, y[train_index])
558     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
558         vectors_test)))
559
560 print(np.mean(accuracies))
561
562 """## 3. TF-IDF"""
563
564 # Commented out IPython magic to ensure Python compatibility.
565 # find the best hyperparameters
566 clf = AdaBoostClassifier(random_state=0)
567 model = make_pipeline(clf)
568
569 param_grid = {'adaboostclassifier__n_estimators': [50, 100],
569     'adaboostclassifier__learning_rate': [0.1, 0.5, 1]}

```

File - E:\Study\ECSE551\Mini\_Project\_2\additional\_classifiers.py

```

571 grid = GridSearchCV(model, param_grid)
572
573 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.8,
574 test_size=0.2, random_state=0)
575 # %time grid.fit(vectors_train, y_train)
576 print(grid.best_params_)
577
578 accuracies = []
579 clf = AdaBoostClassifier(n_estimators=100, learning_rate=0.5, random_state=0)
580 kf = KFold(n_splits=5, shuffle=True)
581 for train_index, test_index in kf.split(X):
582     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
583         test_index])
584     clf.fit(vectors_train, y[train_index])
585     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
586         vectors_test)))
587
588 print(np.mean(accuracies))
589
590 """## 4. CountVectorizer with stem tokenizer"""
591
592 accuracies = []
593 clf = AdaBoostClassifier(n_estimators=100, learning_rate=0.5, random_state=0)
594 kf = KFold(n_splits=5, shuffle=True)
595 for train_index, test_index in kf.split(X):
596     vectors_train, vectors_test = count_vec_stem(X[train_index], X[test_index]
597         ])
598     clf.fit(vectors_train, y[train_index])
599     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
600         vectors_test)))
601
602 print(np.mean(accuracies))
603
604 """## 5. CountVectorizer with lemma tokenizer"""
605
606 accuracies = []
607 clf = AdaBoostClassifier(n_estimators=100, learning_rate=0.5, random_state=0)
608 kf = KFold(n_splits=5, shuffle=True)
609 for train_index, test_index in kf.split(X):
610     vectors_train, vectors_test = count_vec_lemma(X[train_index], X[test_index]
611         )
612     clf.fit(vectors_train, y[train_index])
613     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
614         vectors_test)))
615
616 print(np.mean(accuracies))
617
618 """## 7. k-Neighbors
619
620 ### 1. CountVectorizer
621 """
622
623 accuracies = []
624 neigh = KNeighborsClassifier(n_neighbors=3)
625 kf = KFold(n_splits=5, shuffle=True)
626 for train_index, test_index in kf.split(X):
627     vectors_train, vectors_test = count_vectorizer(X[train_index], X[
628         test_index])
629     neigh.fit(vectors_train, y[train_index])
630     accuracies.append(metrics.accuracy_score(y[test_index], neigh.predict(
631         vectors_test)))
632

```

```

625 print(np.mean(accuracies))
626
627 """### 2. CountVectorizer with stop word"""
628
629 accuracies = []
630 neigh = KNeighborsClassifier(n_neighbors=3)
631 kf = KFold(n_splits=5, shuffle=True)
632 for train_index, test_index in kf.split(X):
633     vectors_train, vectors_test = count_vec_with_sw(X[train_index], X[
634         test_index])
635     neigh.fit(vectors_train, y[train_index])
636     accuracies.append(metrics.accuracy_score(y[test_index], neigh.predict(
637         vectors_test)))
638
639 print(np.mean(accuracies))
640
641 """### 3. TF-IDF"""
642
643 accuracies = []
644 neigh = KNeighborsClassifier(n_neighbors=3)
645 kf = KFold(n_splits=5, shuffle=True)
646 for train_index, test_index in kf.split(X):
647     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
648         test_index])
649     neigh.fit(vectors_train, y[train_index])
650     accuracies.append(metrics.accuracy_score(y[test_index], neigh.predict(
651         vectors_test)))
652
653 """### 4. CountVectorizer with stem tokenizer"""
654
655 accuracies = []
656 neigh = KNeighborsClassifier(n_neighbors=3)
657 kf = KFold(n_splits=5, shuffle=True)
658 for train_index, test_index in kf.split(X):
659     vectors_train, vectors_test = count_vec_stem(X[train_index], X[test_index]
660     ])
661     neigh.fit(vectors_train, y[train_index])
662     accuracies.append(metrics.accuracy_score(y[test_index], neigh.predict(
663         vectors_test)))
664
665 """### 5. CountVectorizer with lemma tokenizer"""
666
667 accuracies = []
668 neigh = KNeighborsClassifier(n_neighbors=3)
669 kf = KFold(n_splits=5, shuffle=True)
670 for train_index, test_index in kf.split(X):
671     vectors_train, vectors_test = count_vec_lemma(X[train_index], X[test_index]
672     ])
673     neigh.fit(vectors_train, y[train_index])
674     accuracies.append(metrics.accuracy_score(y[test_index], neigh.predict(
675         vectors_test)))
676
677 """## 8. Neural Network
678 """
679

```

```

680 accuracies = []
681 kf = KFold(n_splits=5, shuffle=True)
682 for train_index, test_index in kf.split(X):
683     vectors_train, vectors_test = count_vectorizer(X[train_index], X[
684         test_index])
684     clf = MLPClassifier(random_state=0, max_iter=300).fit(vectors_train, y[
685         train_index])
685     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
686         vectors_test)))
686
687 print(np.mean(accuracies))
688
689 """## 2. CountVectorizer with stop word"""
690
691 accuracies = []
692 kf = KFold(n_splits=5, shuffle=True)
693 for train_index, test_index in kf.split(X):
694     vectors_train, vectors_test = count_vec_with_sw(X[train_index], X[
695         test_index])
695     clf = MLPClassifier(random_state=0, max_iter=300).fit(vectors_train, y[
696         train_index])
696     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
697         vectors_test)))
698
699 print(np.mean(accuracies))
700
701 """## 3. TF-IDF"""
702
703 accuracies = []
704 kf = KFold(n_splits=5, shuffle=True)
705 for train_index, test_index in kf.split(X):
706     vectors_train, vectors_test = tfidf_vectorizer(X[train_index], X[
707         test_index])
708     clf = MLPClassifier(random_state=0, max_iter=300).fit(vectors_train, y[
709         train_index])
710     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
711         vectors_test)))
712
713 print(np.mean(accuracies))
714
715 """## 4. CountVectorizer with stem tokenizer"""
716
717 accuracies = []
718 kf = KFold(n_splits=5, shuffle=True)
719 for train_index, test_index in kf.split(X):
720     vectors_train, vectors_test = count_vec_stem(X[train_index], X[test_index]
721         )
722     clf = MLPClassifier(random_state=0, max_iter=300).fit(vectors_train, y[
723         train_index])
724     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(
725         vectors_test)))
726
727 print(np.mean(accuracies))
728
729 """## 5. CountVectorizer with lemma tokenizer"""
730
731 accuracies = []
732 kf = KFold(n_splits=5, shuffle=True)
733 for train_index, test_index in kf.split(X):
734     vectors_train, vectors_test = count_vec_lemma(X[train_index], X[test_index]
735         )
736     clf = MLPClassifier(random_state=0, max_iter=300).fit(vectors_train, y[
737         train_index])

```

```
729     accuracies.append(metrics.accuracy_score(y[test_index], clf.predict(  
    vectors_test)))  
730  
731 print(np.mean(accuracies))
```