


```

['technology',
 'technology',
 'sport',
 'business',
 'business',
 'technology',
 'technology',
 'technology',
 'business',
 'business',
 'technology',
 'internet']

In [166]: from nltk.stem.porter import PorterStemmer
nltk.download('averaged_perceptron_tagger')
nltk.download('wordnet')

[nltk_data] Downloading package averaged_perceptron_tagger to
[C:\Users\WZNDI\AppData\Roaming\nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\WZNDI\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!

Out[166]: True

In [167]: from sklearn.feature_extraction.text import CountVectorizer
import nltk
#-----
#----- Tokenizer and WordNet Lemmatizer with POS -----
#-----#
# define the function
def lemma_tokenizer_v_pos_tag(text):
    # define a nested function for converting POS tag for Lemmatizer
    def convert_tags(tag):
        if tag == 'vbd' or tag == 'vbg' or tag == 'vzb':
            return 'v'
        else:
            return 'n'

    #-----#
    # use the standard scikit-learn tokenizer first
    standard_tokenizer = CountVectorizer().build_tokenizer()
    tokens = standard_tokenizer(text)
    tokens_with_pos_tag = nltk.pos_tag(tokens)
    # print(tokens_with_pos_tag)
    # then use NLTK to perform lemmatisation on each token
    lemmatizer = nltk.WordNetLemmatizer()
    lemma_tokens = []
    for token in tokens_with_pos_tag:
        new_tag = convert_tags(token[1].lower())
        lemma_tokens.append(lemmatizer.lemmatize(token[0], new_tag))
    # print(lemma_tokens)
    return lemma_tokens

```

Text preprocessing and Term Weighting

```

In [168]: from sklearn.feature_extraction.text import TfidfVectorizer
# we can pass in some preprocessing parameters
vectorizer = TfidfVectorizer(stop_words='english', min_df = 10, tokenizer=lemma_tokenizer_v_pos_tag) #ทำการแปลงที่มันคือที่เราต้องการ
X = vectorizer.fit_transform(new_doc)
# display some sample weighted values
print(X)

C:\Users\WZNDI\anaconda3\lib\site-packages\sklearn\feature_extraction\text.py:383: UserWarning: Your stop_words may be inconsistent with your preprocessing. Tokenizing the stop words generated tokens '{'l', 'u'} not in stop_words.
warnings.warn('Your stop_words may be inconsistent with')

(0, 405) 0.2138012607563919
(0, 170) 0.106379046797080929
(0, 2941) 0.129707980384709363
(0, 810) 0.084264128572759518
(0, 2652) 0.08488553054357058
(0, 2908) 0.0927638146137926
(0, 257) 0.1123140131595962
(0, 830) 0.07259828010976334
(0, 229) 0.084264128572759518
(0, 2239) 0.09312961733756304
(0, 966) 0.07428022173629981
(0, 782) 0.10764288080181825
(0, 1662) 0.106950863012645595
(0, 2019) 0.1123140131595962
(0, 3217) 0.084264128572759518
(0, 2431) 0.13166542067711
(0, 2431) 0.1083571132439569
(0, 2470) 0.0786798318621317
(0, 1121) 0.11146499270808044
(0, 1387) 0.097563012324214
(0, 3445) 0.084264128572759518
(0, 1247) 0.1251770046159905
(0, 1714) 0.129707980384709363
(0, 2976) 0.12970798038470937
(0, 2770) 0.06967900046826814
:
(1407, 1982) 0.086490096099221325
(1407, 2123) 0.2837623205323384
(1407, 322) 0.129707980384709363
(1407, 2325) 0.1346779126801964
(1407, 2744) 0.11113978330351282
(1407, 2168) 0.08134498334585134
(1407, 1437) 0.07102657237652493
(1407, 1258) 0.11113978330351282
(1407, 2322) 0.13260252474886698
(1407, 2322) 0.068025142176869525
(1407, 3121) 0.07935521497193686
(1407, 2101) 0.09099613242567584
(1407, 1568) 0.10182245410904453
(1407, 3071) 0.10921239916836216
(1407, 3389) 0.129707980384709363
(1407, 3389) 0.060727642807973203
(1407, 963) 0.10158127850023103
(1407, 822) 0.07833567005804882
(1407, 406) 0.08171046527132773
(1407, 2731) 0.03326977703469115
(1407, 3371) 0.10921239916836216
(1407, 3468) 0.12812199968062039
(1407, 3468) 0.0421055639828279
(1407, 1243) 0.10711203044654932
(1407, 3175) 0.052398879049388045

```

สร้างโมเดล โดยใช้วิธีการ Holdout Splitting และ cross validation 5 fold

นำข้อมูลมาแบ่งเป็นชุดทดสอบและชุดทดสอบ ทดสอบผลใน Train 70.30 ค่า Test

```

In [169]: from sklearn.model_selection import train_test_split
from sklearn import metrics
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=2) #ใช้ส่วนที่ด้านบนเป็นแบบ train, test แล้วที่นี่จะ

```

K nearest neighbor

```

In [170]: from sklearn.neighbors import KNeighborsClassifier
for i in range(100):
    classifier = KNeighborsClassifier(n_neighbors=61)
    classifier.fit(X_train, y_train)
    classifier_pred = classifier.predict(X_test)
    print("Accuracy of KNN = %.2f percent" % (metrics.accuracy_score(y_test, classifier_pred)*100))

Accuracy of KNN = 98.35 percent

```

```

In [171]: from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import cross_val_score
kneigh = KNeighborsClassifier(n_neighbors=20)
cross_knn = cross_val_score(kneigh, X, y, cv=5)
print("KNN: Mean cross-validation accuracy = %.2f percent" % (cross_knn.mean()*100))

KNN: Mean cross-validation accuracy = 97.02 percent

```

Naive Bayes

```

In [172]: from sklearn.naive_bayes import BernoulliNB
naive = BernoulliNB(alpha=0.6)
naive.fit(X_train, y_train)
naive_pred = naive.predict(X_test)
print("Accuracy of Naive Bayes = %.2f percent" % (metrics.accuracy_score(y_test, naive_pred)*100))

Accuracy of Naive Bayes = 97.64 percent

```

```
In [173]: from sklearn.naive_bayes import BernoulliNB
from sklearn.model_selection import cross_val_score
model_nb = BernoulliNB(alpha=0.3)
cross_naive = cross_val_score(model_nb, X, Y, cv=5)
print("Naive Bayes: Mean cross-validation accuracy = %.2f percent" %(cross_naive.mean()*100))

Naive Bayes: Mean cross-validation accuracy = 97.73 percent
```

SVM

```
In [174]: from sklearn import svm
svm = svm.SVC(C=1, kernel='linear')
svm.fit(X_train, y_train)
svm_pred = svm.predict(X_test)
print("Accuracy of SVM = %.2f percent" %(metrics.accuracy_score(y_test,svm_pred)*100))

Accuracy of SVM = 99.05 percent
```

```
In [175]: from sklearn.svm import SVC
from sklearn.model_selection import cross_val_score
model = SVC(gamma='scale')
acc_scores = cross_val_score(model,X,Y, cv=5, scoring="accuracy")
print("SVM: Mean cross-validation accuracy = %.2f percent" %(acc_scores.mean()*100))

SVM: Mean cross-validation accuracy = 98.30 percent
```

Decisiontree

```
In [176]: from sklearn.tree import DecisionTreeClassifier
decision_tree = DecisionTreeClassifier()
decision_tree.fit(X_train, y_train)
decision_tree_pred = decision_tree.predict(X_test)
print("Accuracy of Decision tree = %.2f percent" %(metrics.accuracy_score(y_test,decision_tree_pred)*100))

Accuracy of Decision tree = 86.76 percent
```

```
In [177]: from sklearn.tree import DecisionTreeClassifier
decision_tree = DecisionTreeClassifier(random_state=1)
cross_decision_tree = cross_val_score(decision_tree, X, Y, cv=5)
print("Decision tree: Mean cross-validation accuracy = %.2f percent" %(cross_decision_tree .mean()*100))

Decision tree: Mean cross-validation accuracy = 86.79 percent
```

สรุปผล

เมื่อใช้ cross validation 5 fold เมื่อเรา Crossvalidation เป็นการตรวจสอบอย่างใดอย่างหนึ่งว่าโมเดลที่เราได้จากการtrain นั้นดีหรือไม่ดี ค่าความแม่นยำของ Crossvalidation ได้บันทึก SVM มีค่าAccuracy บนสูงที่สุด 98.30% ทำให้เรามั่นใจว่าโมเดลที่เราได้มาจะสามารถใช้ในการจำแนกหมวดหมู่ระหว่างคนหนุ่มสาวช่วงวัยรุ่นได้ดีที่สุด เมื่อบรรลุเป้าหมายที่ตั้งไว้