



An Application of Natural Language Processing on forex gold spot News Analysis การประยุกต์ใช้การประมวลผลภาษาธรรมชาติเพื่อการวิเคราะห์ข่าวตลาด forex gold spot

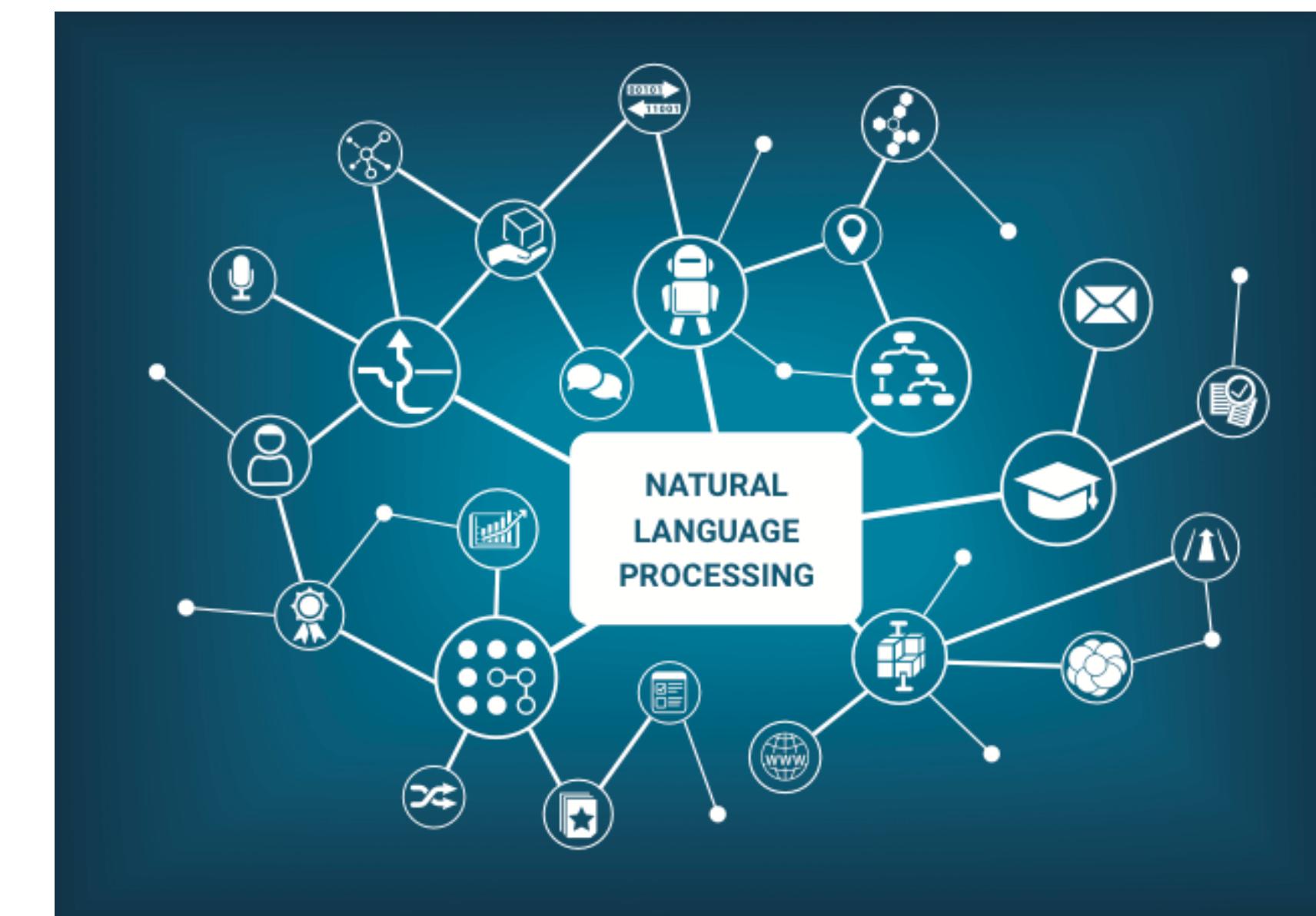
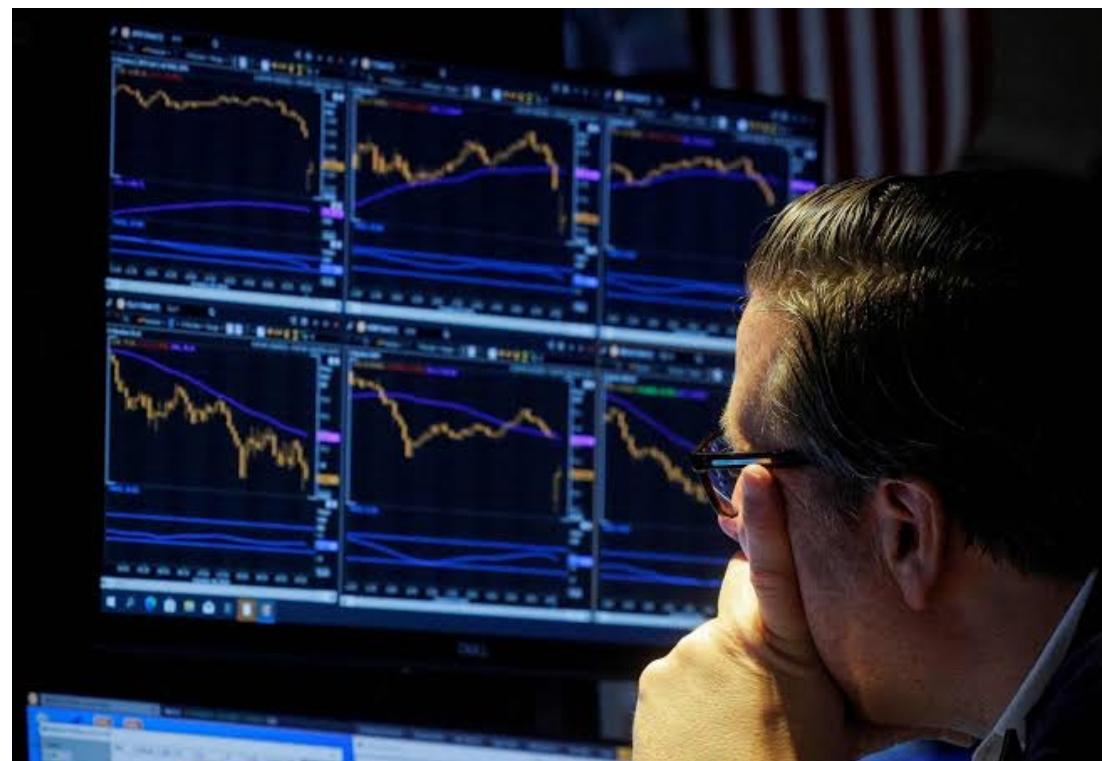
ผู้จัดทำโครงการ
นายชวัลชัย อภิชาติชูติวน์ รหัสนักศึกษา 6210110646

อาจารย์ที่ปรึกษาโครงการ
รศ.ดร.มนตรี กาญจนะเดชะ

OUTLINE

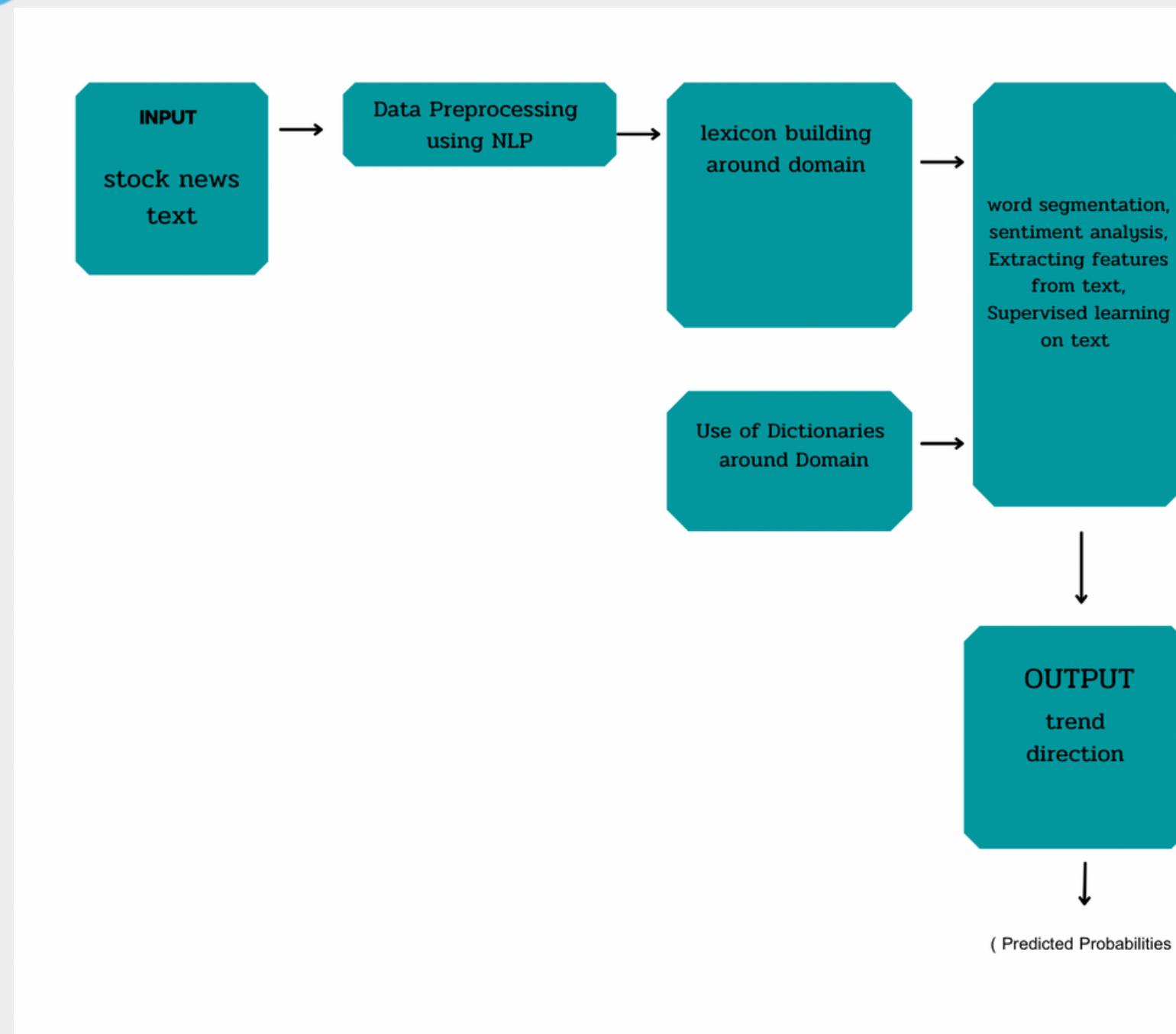
- บทคัดย่อ
- แผนการดำเนินงาน
- รายละเอียดการดำเนินงาน
- ความก้าวหน้าการดำเนินงาน
- สรุป

บทคัดย่อ



แผนการดำเนินงาน

รายละเอียดการดำเนินงาน



ขั้นตอน INPUT
ภาษา自然 text

ขั้นตอน Process
ใช้ algorithms

ขั้นตอน Output
คาดการณ์ล่วงหน้า
สรุปข้อความ

ข้อมูล Input

ข่าวสาร	ข่าวประเภท	อ้างอิงค์
1. Census Bureau	กระทรวงพาณิชย์สหรัฐ	[7]
2. Us Department of labo	กระทรวงแรงงาน	[8]
3. Energy information Administration	ข้อมูลด้านพลังงาน	[9]
4. กองทุน SPDR	การซื้อขายทองคำของ กองทุน SPDR	[10]
5. Federal Reserve	ธนาคารกลางสหรัฐ	[11]
6. Bloomberg	รายงานข่าวทั่วไปรอบโลก	[12]
7. Twitter	ข่าวทั่วไป	[13]

ขั้นตอน Process

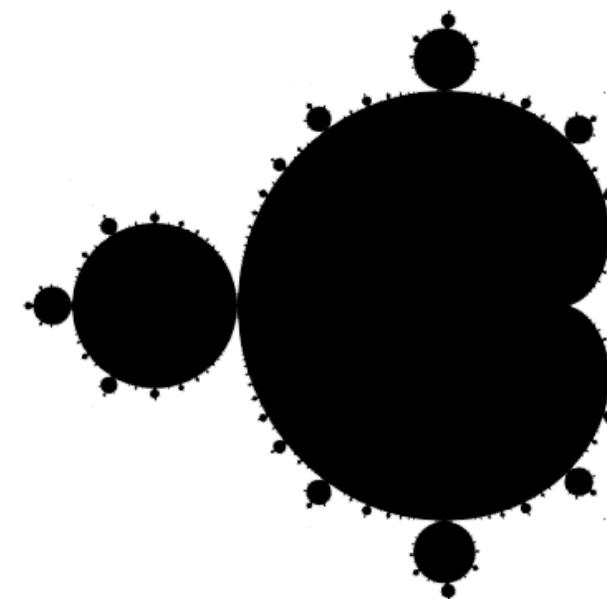
- sentiment analysis
 - tokenization
 - stopword
 - stemmer
- model naive baye

ข้อ^๑ บันตอน Output

- สามารถคาดการณ์ล่วงหน้าได้
- สามารถสรุปบทความได้

ความก้าวหน้าการดำเนินงาน

ความก้าวหน้า 1 ศึกษา Library Python NLP และทดสอบ



TextBlob

ความก้าวหน้าการดำเนินงาน

ความก้าวหน้า 2 ทดสอบ Library NLTK

```
stem.py > ...
1 from nltk.stem import PorterStemmer
2 from nltk.tokenize import sent_tokenize, word_tokenize
3
4 ps = PorterStemmer()
5
6 example_words = ["python","pythoner","pythoning","pythoned","pythonly,feeling"]
7
8 for w in example_words:
9     print(ps.stem(w))
10
11 new_text = "It is important to by very pythonly while you are pythoning with python. All pythonders have pythoned poorly at least onc
12 words = word_tokenize(new_text)
13
14 for w in words:
15     print(ps.stem(w))
```

PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE

```
MacBook-Pro-khxng-Chualchai:project whatbest$ cd /Users/whatbest/project ; /usr/bin/env /usr/local/bin/python3 /Users/whatbest/.vscode/extensions/ms-python.python-2021.10.1365161279/pythonFiles/lib/python/debugpy/launcher 53446 -- /Users/whatbest/project/stem.py
python
python
python
python
python,feel
it
is
import
to
by
veri
pythonli
while
you
are
pytho
with
pyt
all
pytho
have
pytho
portli
at
last
onc
MacBook-Pro-khxng-Chualchai:project whatbest$
```

```
ntknlp1.py > ...
1 from nltk.tokenize import sent_tokenize, word_tokenize
2
3 EXAMPLE_TEXT = "Hello Mr. Smith, how are you doing today? The weather is great, and Python is awesome. The sky is pinkish-blue. You
4
5 print(sent_tokenize(EXAMPLE_TEXT))
6 print(word_tokenize(EXAMPLE_TEXT))
```

PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE

```
MacBook-Pro-khxng-Chualchai:project whatbest$ cd /Users/whatbest/project ; /usr/bin/env /usr/local/bin/python3 /Users/whatbest/.vscode/extensions/ms-python.python-2021.10.1365161279/pythonFiles/lib/python/debugpy/launcher 52987 -- /Users/whatbest/project/ntknlp1.py
['Hello Mr. Smith, how are you doing today?', 'The weather is great, and Python is awesome.', 'The sky is pinkish-blue.', "You shouldn't eat cardboard."]
['Hello', 'Mr.', 'Smith', ',', 'how', 'are', 'you', 'doing', 'today', '?', 'The', 'weather', 'is', 'great', ',', 'and', 'Python', 'is', 'awesome', '.']
MacBook-Pro-khxng-Chualchai:project whatbest$
```

Stemmer words with NLTK

Tokenizing Words and Sentences with NLTK

```
stopwords.py > ...
1 from nltk.corpus import stopwords
2 from nltk.tokenize import word_tokenize
3
4 example_sent = "This is a sample sentence, showing off the stop words filtration."
5
6 stop_words = set(stopwords.words('english'))
7
8 word_tokens = word_tokenize(example_sent)
9
10 filtered_sentence = [w for w in word_tokens if not w in stop_words]
11
12 filtered_sentence = []
13
14 for w in word_tokens:
15     if w not in stop_words:
16         filtered_sentence.append(w)
17
18 print("stop_words = ", stop_words)
19 print("word_tokens = ", word_tokens)
20 print("filtered_sentence = ", filtered_sentence)
```

PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE

```
MacBook-Pro-khxng-Chualchai:project whatbest$ cd /Users/whatbest/project ; /usr/bin/env /usr/local/bin/python3 /Users/whatbest/.vscode/extensions/ms-python.python-2021.10.1365161279/pythonFiles/lib/python/debugpy/launcher 53097 -- /Users/whatbest/project/stopwords.py
stop_words = {'hadn', 'itself', 'just', 'shan', 'you', 'as', 'an', 'here', 'how', 'thatll', 'below', 'very', 'where', 'the', 'doesn', 't', 'theirs', 'won', 'wouldnt', 'is', 'any', 'nor', 'youve', 'he', 'them', 'yourselves', 'aren', 't', 'their', 'above', 'my', 'wasn', 't', 'while', 'y
oull', 'too', 'over', 'some', 'isn', 'aren', 'your', 'been', 'were', 'do', 'have', 'no', 'hsm', 'why', 'will', 'if', 'these', 'having', 'few', 'ma', 'or', 'verser', 'yourself', 'el', 'wouldve', 'mighthave', 'she', 'are', 'did', 'do', 'could', 'for', 'yo
u', 'of', 'our', 'only', 'who', 'it', 'when', 'we', 'had', 'll', 'that', 'from', 'they', 'about', 'could', 'al', 'weren', 'me', 'now', 'further', 'and', 'haven', 'hers', 'not', 'should', 'most', 'i', 're', 'wasn', 'once', 'during', 'myself', 'on
', 'misten', 'her', 'than', 'needin', 'in', 'again', 'between', 'haven', 'his', 'into', 'then', 'youd', 'but', 'themselves', 'don', 'does', 'at', 'own', 's', 'youre', 'because', 'himself', 'down', 'both', 'other', 'neednt', 'isnt', 'has', 'after', 'sheis', 'of', 'each', 'its', 'this', 'mighthave', 'to', 'its', 'through', 'doing', 'him', 'a', 'shan', 'what', 'whom', 'more', 'd', 'ain', 'mistr', 'same', 'which', 'didn', 't', 'out', 'hadn', 'before', 'don', 't', 'against', 'under', 'shouldn', 'wouldn', 'ours', 'until', 'there', 'yours', 'y', 'by', 'be', 'being', 'hasn', 't', 'such', 'doesnt', 've', 'those', 'with' ]
word_tokens = ['This', 'is', 'a', 'sample', 'sentence', ',', 'showing', 'off', 'the', 'stop', 'words', 'filtration', '.']
filtered_sentence = ['This', 'sample', 'sentence', ',', 'showing', 'stop', 'words', 'filtration', '.']
```

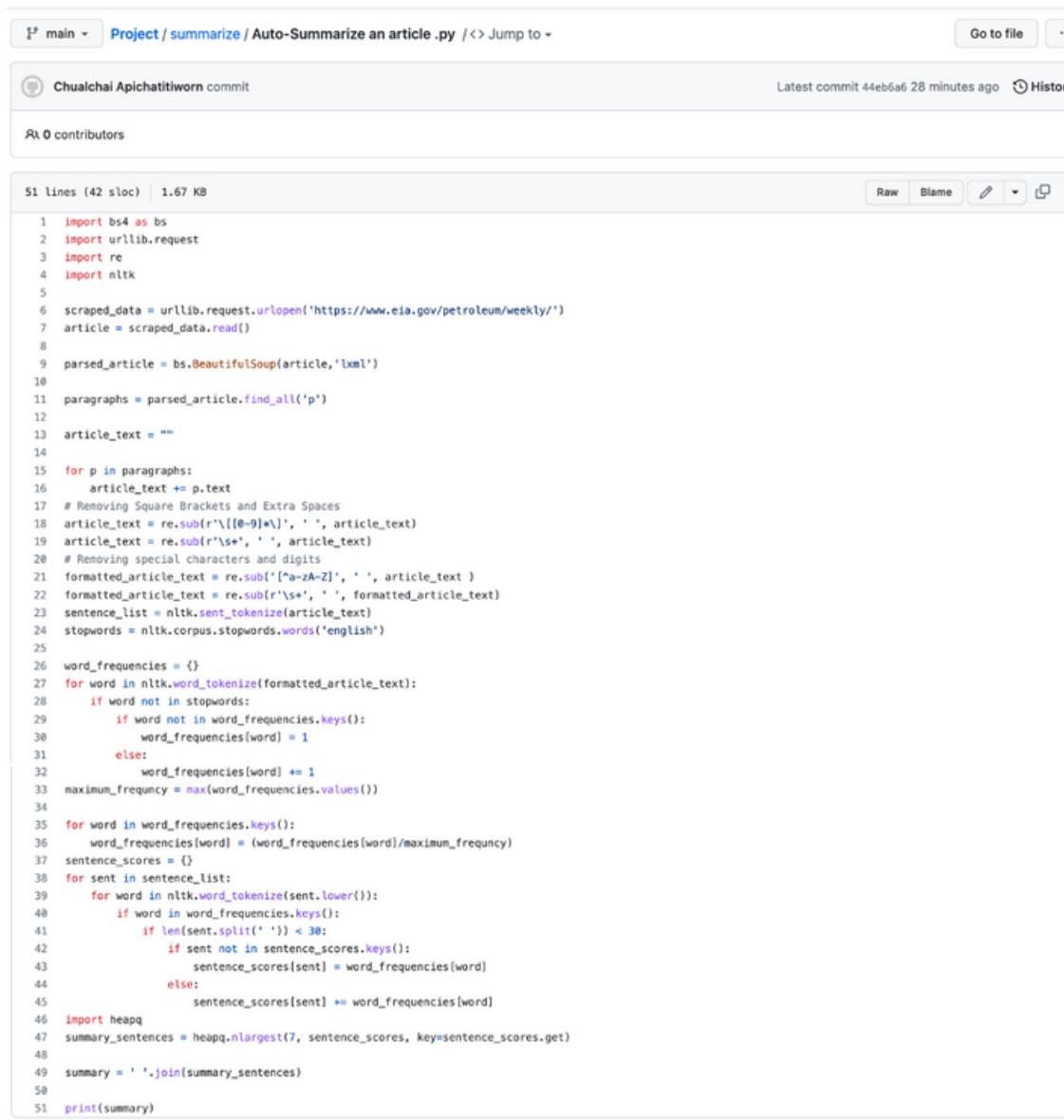
MacBook-Pro-khxng-Chualchai:project whatbest\$

Stop words with NLTK

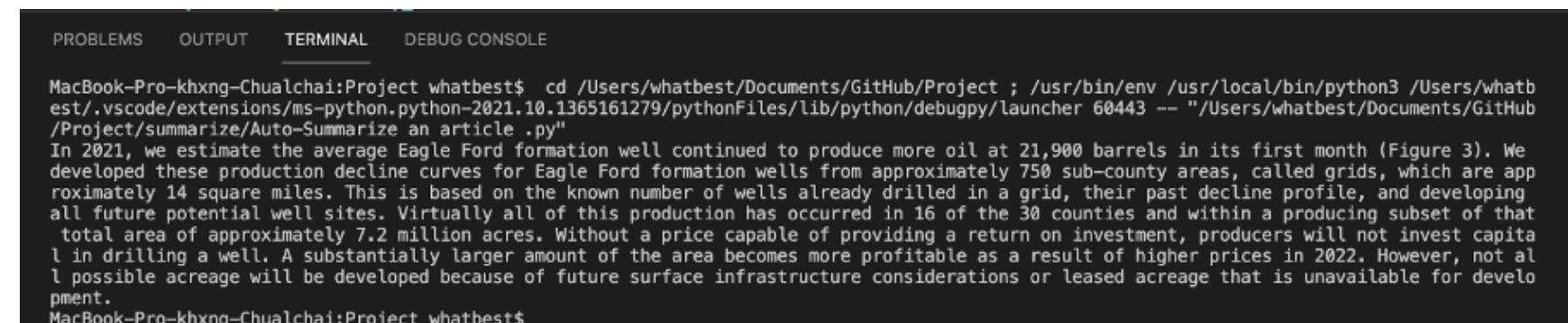


ความก้าวหน้าการดำเนินงาน

ความก้าวหน้า 3 Text Summarization with NLTK in Python



```
main Project / summarize / Auto-Summarize an article .py < > Jump to ▾ Go to file ...  
Chualchai Apichatitiworn commit Latest commit 44eb6a6 28 minutes ago History  
0 contributors  
51 lines (42 sloc) | 1.67 KB  
Raw Blame ⌂ ⌃ ⌄ ⌅ ⌆  
1 import bs4 as bs  
2 import urllib.request  
3 import re  
4 import nltk  
5  
6 scraped_data = urllib.request.urlopen('https://www.eia.gov/petroleum/weekly/').read()  
7 article = scraped_data.read()  
8  
9 parsed_article = bs.BeautifulSoup(article,'lxml')  
10 paragraphs = parsed_article.findAll('p')  
11  
12 article_text = ""  
13  
14 for p in paragraphs:  
15     article_text += p.text  
16  
17 # Removing Square Brackets and Extra Spaces  
18 article_text = re.sub(r'\[0-9]*\]', ' ', article_text)  
19 article_text = re.sub(r'\s+', ' ', article_text)  
20  
21 # Removing special characters and digits  
22 formatted_article_text = re.sub("[^a-zA-Z]", " ", article_text )  
23 formatted_article_text = re.sub(r'\s+', ' ', formatted_article_text)  
24 sentence_list = nltk.sent_tokenize(article_text)  
25 stopwords = nltk.corpus.stopwords.words('english')  
26  
27 word_frequencies = {}  
28 for word in nltk.word_tokenize(formatted_article_text):  
29     if word not in stopwords:  
30         if word not in word_frequencies.keys():  
31             word_frequencies[word] = 1  
32         else:  
33             word_frequencies[word] += 1  
34 maximum_frequency = max(word_frequencies.values())  
35  
36 for word in word_frequencies.keys():  
37     word_frequencies[word] = (word_frequencies[word]/maximum_frequency)  
38 sentence_scores = {}  
39 for sent in sentence_list:  
40     for word in nltk.word_tokenize(sent.lower()):  
41         if word in word_frequencies.keys():  
42             if len(sent.split(' ')) < 30:  
43                 if sent not in sentence_scores.keys():  
44                     sentence_scores[sent] = word_frequencies[word]  
45                 else:  
46                     sentence_scores[sent] += word_frequencies[word]  
47 import heapq  
48 summary_sentences = heapq.nlargest(7, sentence_scores, key=sentence_scores.get)  
49  
50 summary = ' '.join(summary_sentences)  
51 print(summary)
```



PROBLEMS OUTPUT TERMINAL DEBUG CONSOLE

```
MacBook-Pro-khxng-Chualchai:Project whatbest$ cd /Users/whatbest/Documents/GitHub/Project ; /usr/bin/env /usr/local/bin/python3 /Users/whatbest/.vscode/extensions/ms-python.python-2021.10.1365161279/pythonFiles/lib/python/debugpy/launcher 60443 -- "/Users/whatbest/Documents/GitHub/Project/summarize/Auto-Summarize an article .py"  
In 2021, we estimate the average Eagle Ford formation well continued to produce more oil at 21,900 barrels in its first month (Figure 3). We developed these production decline curves for Eagle Ford formation wells from approximately 750 sub-county areas, called grids, which are approximately 14 square miles. This is based on the known number of wells already drilled in a grid, their past decline profile, and developing all future potential well sites. Virtually all of this production has occurred in 16 of the 30 counties and within a producing subset of that total area of approximately 7.2 million acres. Without a price capable of providing a return on investment, producers will not invest capital in drilling a well. A substantially larger amount of the area becomes more profitable as a result of higher prices in 2022. However, not all possible acreage will be developed because of future surface infrastructure considerations or leased acreage that is unavailable for development.  
MacBook-Pro-khxng-Chualchai:Project whatbest$
```

Tokennization
stopword
stemmer

ความก้าวหน้าการดำเนินงาน

ความก้าวหน้า 4 สร้างเมืองข้อมูล

	Dates	URL	News	Price Direction Up	Price Direction Constant	Price Direction Down	Asset Comparision	Past Information	Future Information	PriceSentiment
0	28/1/16	http://www.marketwatch.com/story/april-gold-down-20-cents-to-settle-at-11161...	april gold down 20 cents to settle at \$1,116.1...	0.0	0.0	1.0	0.0	1.0	0.0	negative
1	13/9/17	http://www.marketwatch.com/story/gold-prices-s...	gold suffers third straight daily decline	0.0	0.0	1.0	0.0	1.0	0.0	negative
2	26/7/16	http://www.marketwatch.com/story/gold-futures-...	Gold futures edge up after two-session decline	1.0	0.0	0.0	0.0	1.0	0.0	positive
3	28/2/18	https://www.metalsdaily.com/link/277199/dent-re...	dent research : is gold's day in the sun comin...	0.0	0.0	0.0	0.0	0.0	1.0	none
4	6/9/17	http://www.marketwatch.com/story/gold-steadies...	Gold snaps three-day rally as Trump, lawmakers...	0.0	0.0	1.0	0.0	1.0	0.0	negative
5	16/8/16	http://www.marketwatch.com/story/dec-gold-clim...	Dec. gold climbs \$9.40, or 0.7%, to settle at ...	1.0	0.0	0.0	0.0	1.0	0.0	positive
6	24/9/13	https://economictimes.indiatimes.com/markets/c...	gold falls by rs 25 on sluggish demand, global...	0.0	0.0	1.0	0.0	1.0	0.0	negative

ความก้าวหน้าการดำเนินงาน

ความก้าวหน้า 5 ทดสอบนำเมืองข้อมูลมาใช้งาน Model Trainning and Test dataset

Step 1 - Loading the required libraries and modules.

Step 2 - Loading the data and performing basic data checks.

Step 3 - Pre-processing the raw text and getting it ready for machine learning.

Step 4 - Creating the Training and Test datasets.

Step 5 - Converting text to word frequency vectors with TfidfVectorizer.

Step 6 - Create and fit the classifier.

Step 7 - Computing the evaluation metrics.

Step 1 - Loading the Required Libraries and Modules

```
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
import re

stemmer = PorterStemmer()
words = stopwords.words("english")

df['processedtext'] = df['News'].apply(lambda x: " ".join([stemmer.stem(i) for i in re.sub("[^a-zA-Z]", " ", x).split() if i not in words])
```

ความก้าวหน้าการดำเนินงาน

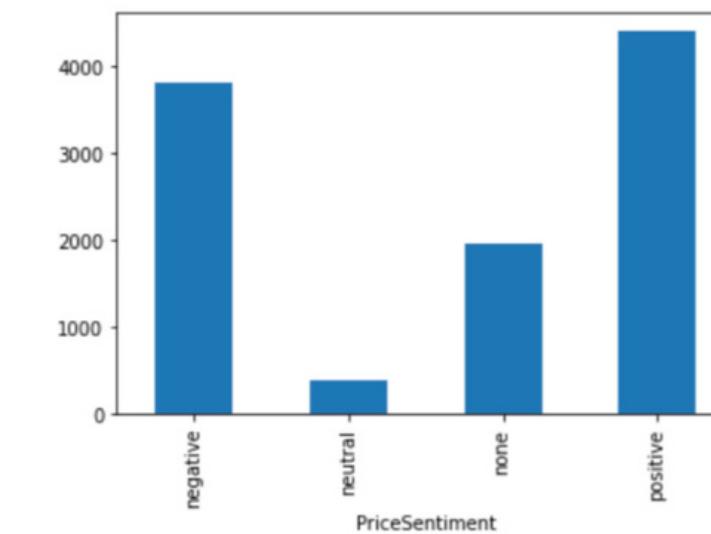
Step 2 - Loading the Data and Performing Basic Data Checks

```
df = pd.read_csv('gold-dataset-sinha-khandait.csv')
print(df.shape)
df.head()
```

(10570, 10)

	Dates	URL	News	Price Direction Up	Price Direction Constant	Price Direction Down	Asset Comparision	Past Information	Future Information	PriceSentiment
0	28/1/16	http://www.marketwatch.com/story/april-gold-down-20-cents-to-settle-at-\$11161...	april gold down 20 cents to settle at \$1,116.1...	0	0	1	0	1	0	negative
1	13/9/17	http://www.marketwatch.com/story/gold-prices-s...	gold suffers third straight daily decline	0	0	1	0	1	0	negative
2	26/7/16	http://www.marketwatch.com/story/gold-futures-...	Gold futures edge up after two-session decline	1	0	0	0	1	0	positive
3	28/2/18	https://www.metaldaily.com/link/277199/dent-re...	dent research : is gold's day in the sun comin...	0	0	0	0	0	1	none
4	6/9/17	http://www.marketwatch.com/story/gold-steadies...	Gold snaps three-day rally as Trump, lawmakers...	0	0	1	0	1	0	negative

```
df.groupby(df['PriceSentiment']).News.count().plot.bar(ylim=0)
plt.show()
print(4533/10570) #Baseline accuracy
```



0.4288552507095553

ความก้าวหน้าการดำเนินงาน

Step 3 – Pre-processing the Raw Text and Getting It Ready for Machine Learning

```
import nltk
nltk.download('stopwords')
```

```
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
import re

stemmer = PorterStemmer()
words = stopwords.words("english")

df['processedtext'] = df['News'].apply(lambda x: " ".join([stemmer.stem(i) for i in re.sub("[^a-zA-Z]", " ", x).split() if i not in words])
```

		URL	News	Price Up	Price Direction Constant	Price Direction Down	Asset Comprasion	Past Information	Future Information	PriceSentiment	processed
0	28/1/16	http://www.marketwatch.com/story/april-gold-do...	april gold down 20 cents to settle at \$1,116.1...	0	0	1	0	1	0	negative	april gold set
1	13/9/17	http://www.marketwatch.com/story/gold-prices...	gold suffers third straight daily decline	0	0	1	0	1	0	negative	gold si third stri d
2	26/7/16	http://www.marketwatch.com/story/gold-futures-...	Gold futures edge up after two-session decline	1	0	0	0	1	0	positive	gold futur two ses d
3	28/2/18	https://www.metaldaily.com/link/277199/dent-r...	dent research : is gold a day in the sun comin...	0	0	0	0	0	1	none	dent rese gold day come :
4	6/9/17	http://www.marketwatch.com/story/gold-steadies...	Gold snaps three-day rally as Trump, lawmakers...	0	0	1	0	1	0	negative	gold three day trump reac
5	16/8/16	http://www.marketwatch.com/story/dec-gold-clim...	Dec. gold futures end up \$9.40, or 0.7%, to settle at ...	1	0	0	0	1	0	positive	dec gold c set
6	24/9/13	https://economictimes.indiatimes.com/markets/c...	gold falls by rs 25 on sluggish demand, global...	0	0	1	0	1	0	negative	gold f slug demand gl
7	23/9/16	http://www.marketwatch.com/story/gold-futures-...	Gold futures fall for the session, but gain fo...	1	0	1	0	1	0	positive	gold futu session
8	21/10/12	https://www.thehindubusinessline.com/opinion/c...	Gold struggles; silver dips, base metals falter	0	1	0	1	1	0	neutral	gold str silver base n f
9	16/3/18	http://www.marketwatch.com/story/april-gold-ho...	april gold holds slight gain, up \$2.50, or 0.2...	1	0	0	0	1	0	positive	april gold slight ga

ความก้าวหน้าการดำเนินงาน

Step 4 - Creating the Training and Test Datasets

```
from sklearn.model_selection import train_test_split  
  
target = df['PriceSentiment']  
  
X_train, X_test, y_train, y_test = train_test_split(df['processedtext'], target, test_size=0.30, random_state=100)  
  
print(df.shape); print(X_train.shape); print(X_test.shape)  
  
(10570, 11)  
(7399,)  
(3171,)
```

ความก้าวหน้าการดำเนินงาน

Step 5 - Converting Text to Word Frequency Vectors with TfidfVectorizer.

```
from sklearn.feature_extraction.text import TfidfVectorizer

vectorizer_tfidf = TfidfVectorizer(stop_words='english', max_df=0.7)

train_tfidf = vectorizer_tfidf.fit_transform(X_train.values.astype('U'))

test_tfidf = vectorizer_tfidf.transform(X_test.values.astype('U'))

print(vectorizer_tfidf.get_feature_names()[:10])

['aayog', 'abat', 'abbrevi', 'abc', 'abn', 'acacia', 'acceler', 'access', 'account', 'accredit']
```

```
print(train_tfidf.shape); print(test_tfidf.shape)

(7399, 2698)
(3171, 2698)
```

ความก้าวหน้าการดำเนินงาน

Step 6 - Create and Fit the Classifier.

การประเมิน Naïve Bayes Model

```
from sklearn.naive_bayes import MultinomialNB
from sklearn import metrics

nb_classifier = MultinomialNB()

nb_classifier.fit(train_tfidf, y_train)

pred2 = nb_classifier.predict(test_tfidf)
print(pred2[:10])

['positive' 'negative' 'positive' 'positive' 'negative' 'none' 'none'
 'positive' 'negative' 'neutral']
```

ความก้าวหน้าการดำเนินงาน

Step 7 - Computing the Evaluation Metrics

```
# Calculate the accuracy score: score
accuracy_tfidf = metrics.accuracy_score(y_test, pred2)
print(accuracy_tfidf)

Conf_metrics_tfidf = metrics.confusion_matrix(y_test, pred2, labels=['positive', 'negative', 'none', 'neutral'])
print(Conf_metrics_tfidf)

0.7303689687795648
[[1142 151 29 0]
 [ 297 785 37 0]
 [ 172 64 377 0]
 [ 62 40 3 12]]
```

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix

classifier = RandomForestClassifier(n_estimators = 10, criterion = 'entropy', random_state = 100)
classifier.fit(train_tfIdf, y_train)
```

```
predRF = classifier.predict(test_tfIdf)
print(predRF[:10])

# Calculate the accuracy score
accuracy_RF = metrics.accuracy_score(y_test, predRF)
print(accuracy_RF)

Conf_metrics_RF = metrics.confusion_matrix(y_test, predRF, labels=['positive', 'negative', 'none', 'neutral'])
print(Conf_metrics_RF)

['positive' 'negative' 'positive' 'positive' 'negative' 'none' 'none'
 'negative' 'positive' 'neutral']
0.7723115736360769
[[1042 206 66 8]
 [ 173 891 50 5]
 [ 95 58 456 4]
 [ 23 29 5 60]]
```

สรุปผล

ระหว่าง RandomForest กับ naive baye
ผลที่ได้

Random Forest มีประสิทธิภาพมากกว่า
ที่เยอะกว่า 77%
และ naive baye 73 %

ความก้าวหน้าการดำเนินงาน

ความก้าวหน้า 6 ทดลองทำ Sentiment Analysis Naïve Bayes Classifier

Read into Python

Let's first read the required data from CSV file using Pandas library.

```
In [309...  
import pandas as pd  
import csv  
from sklearn.model_selection import train_test_split, cross_val_score  
from sklearn.utils import shuffle  
from sklearn.feature_extraction.text import CountVectorizer  
import numpy as np  
import pandas as pd  
# data processing, CSV file I/O (e.g. pd.read_csv)  
import matplotlib.pyplot as plt  
import seaborn as sns  
#For better Visualisation  
from bs4 import BeautifulSoup  
  
%matplotlib inline  
import warnings  
warnings.filterwarnings('ignore')
```

```
In [310...  
data = pd.read_csv('main.csv')  
data = data[['Dates','News','PriceSentiment']]  
print(data.shape)  
data.head(7)
```

```
(9999, 3)  
Out[310...  
          Dates           News  PriceSentiment  
0  28/1/16  april gold down 20 cents to settle at $1,116.1...  negative  
1  13/9/17      gold suffers third straight daily decline  negative  
2  26/7/16  Gold futures edge up after two-session decline  positive  
3  28/2/18  dent research : is gold's day in the sun comin...   none  
4  6/9/17  Gold snaps three-day rally as Trump, lawmakers...  negative  
5  16/8/16  Dec. gold climbs $9.40, or 0.7%, to settle at ...  positive  
6  24/9/13  gold falls by rs 25 on sluggish demand, global...  negative
```

Now, show the data how looks like...



ความก้าวหน้าการดำเนินงาน

```
In [314... import nltk
nltk.download('stopwords')

[nltk_data] Downloading package stopwords to
[nltk_data]   /Users/whatbest/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
True

Out[314...]

In [315... from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
import re

stemmer = PorterStemmer()
words = stopwords.words("english")

data['processedtext'] = data['News'].apply(lambda x: " ".join([stemmer.stem(i) for i in re.sub("[^a-zA-Z]", " ", x).lower()])

In [316... print(data.shape)
data.head(10)

(9999, 4)

Out[316...    Dates          News  PriceSentiment  processedtext
0  28/1/16  april gold down 20 cents to settle at $1,116.1...  negative      april gold cent settl oz
1  13/9/17        gold suffers third straight daily decline  negative      gold suffer third straight daili declin
2  26/7/16  Gold futures edge up after two-session decline  positive      gold futur edg two session declin
3  28/2/18  dent research : is gold's day in the sun comin...  none      dent research gold day sun come soon
4  6/9/17  Gold snaps three-day rally as Trump, lawmakers...  negative      gold snap three day ralli trump lawmak reach d...
5  16/8/16  Dec. gold climbs $9.40, or 0.7%, to settle at ...  positive      dec gold climb settl oz
6  24/9/13  gold falls by rs 25 on sluggish demand, global...  negative      gold fall rs sluggish demand global cue
7  23/9/16  Gold futures fall for the session, but gain fo...  positive      gold futur fall session gain week
8  21/10/12  Gold struggles; silver slides, base metals falter  neutral      gold struggl silver slide base metal falter
9  16/3/18  april gold holds slight gain, up $2.50, or 0.2...  positive      april gold hold slight gain oz
```

Pre-process Data

We need to remove package name as it's not relevant. Then convert text to lowercase for CSV data. So, this is data pre-process stage.

```
In [317... def preprocess_data(data):
# Remove package name as it's not relevant
#data = data.drop('News', axis=1)

# Convert text to lowercase
data['processedtext'] = data['processedtext'].str.strip().str.lower()
return data

In [318... data = preprocess_data(data)
```

Splitting Data

First, separate the columns into dependent and independent variables (or features and label). Then you split those variables into train and test set.

```
In [319... df = data
# Split into training and testing data
x = data['processedtext']
y = data['PriceSentiment']
x, x_test, y, y_test = train_test_split(x,y, stratify=y, test_size=0.25, random_state=42)

Vectorize text reviews to numbers.

In [320... # Vectorize text reviews to numbers
vec = CountVectorizer(stop_words='english')
x = vec.fit_transform(x).toarray()
x_test = vec.transform(x_test).toarray()
```

โดยใช้เทสในระบบ 25% และ เทวนะระบบ 75%

ความก้าวหน้าการดำเนินงาน

Model Generation

After splitting and vectorize text reviews to number, we will generate a random forest model on the training set and perform prediction on test set features.

```
In [321... from sklearn.naive_bayes import MultinomialNB  
  
model = MultinomialNB()  
model.fit(x, y)  
  
Out[321... MultinomialNB()  
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.  
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.
```

ผลที่ได้คือความแม่นยำจากการทดสอบและเท่านั้น 76%

```
['none'] [False]  
Gold heading for worst week since November on rate hike worries  
['negative'] [ True]  
64.3859649122807  
  
In [146... model.predict(vec.transform(['Changes in non-farm payrolls increase.']))  
  
Out[146... array(['none'], dtype='<U8')  
  
Average hourly earnings, m/m, remain unchanged. รายได้เฉลี่ยต่อชั่วโมง m/m ยังคงไม่เปลี่ยนแปลง The change in non-farm payrolls increased from the previous time. การเปลี่ยนแปลงในการจ้างงานนอกภาคเกษตรเพิ่มขึ้นจากครั้งก่อน lower unemployment rate อัตราการว่างงานลดลง
```

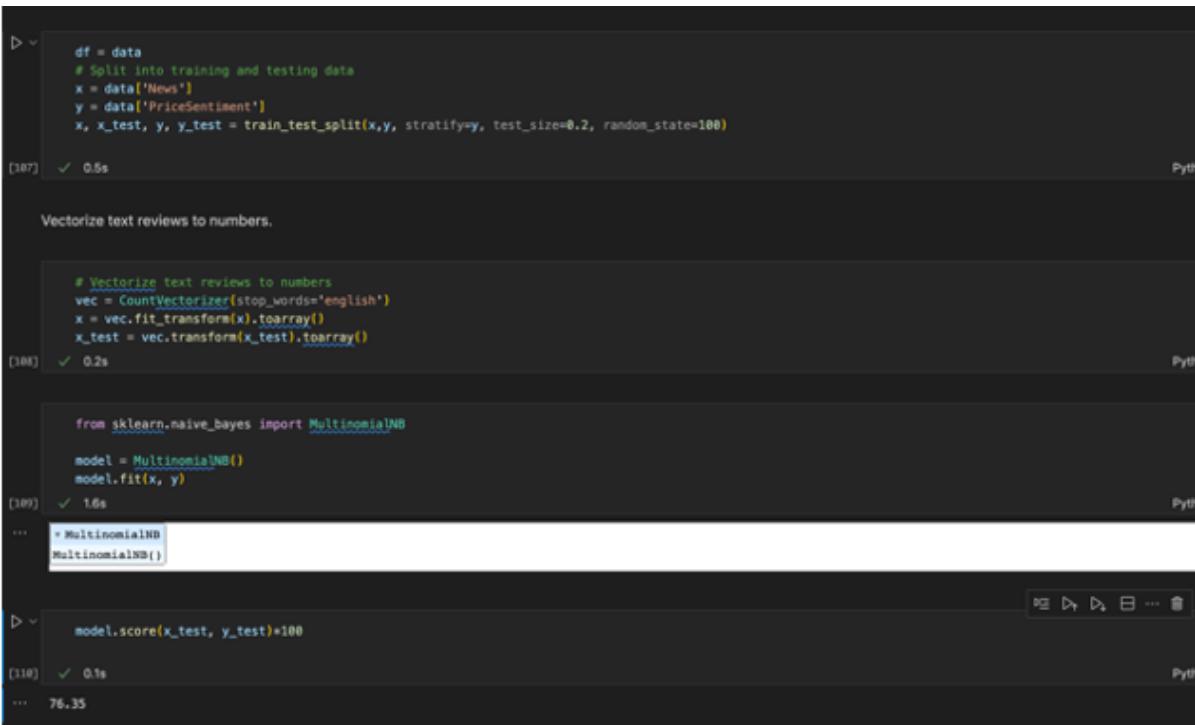
Evaluating Model

After model generation, check the accuracy using actual and predicted values.

```
In [322... model.score(x_test, y_test)*100  
  
Out[322... 75.92  
  
Then check prediction...  
  
In [323... from itertools import count  
import pandas as pd  
df = pd.read_csv('gold-dataset-sinha-khandait.csv', sep=',', header=None)  
start = 10000  
end = 10570  
df = df[start - 1:end - 1]  
correct = 0  
for i in range(len(df)):  
    print(df.values[i][2])  
    print(model.predict(vec.transform([df.values[i][2]])), df.values[i][9] == model.predict(vec.transform([df.values[i][2]])))  
  
    if df.values[i][9] == model.predict(vec.transform([df.values[i][2]])):  
        correct += 1  
  
print(correct / len(df) * 100 )  
  
Kerala govt studying tax reduction on gold  
['none'] [ True]  
gold gets a festive lift, tops rs 31,000-mark  
['positive'] [ True]  
gold ends at new record high of $837.50 an ounce on nygemex  
['positive'] [ True]  
Gold futures prices trade near the session's lows on Comex  
['negative'] [ True]  
gold settles at a more than 1-week high as dollar softens ahead of fed decision  
['positive'] [ True]  
Gold prices mark first climb in 7 sessions  
['positive'] [ True]  
Gold investment demand falls 34%, but WGC says it was still 'robust'  
['negative'] [False]  
gold futures end higher, but post a weekly decline  
['positive'] [ True]
```

ความก้าวหน้าการดำเนินงาน

ความก้าวหน้า 7 การเตรียมเทียบ โดยใช้ Model Naïve baye



```
df = data
# Split into training and testing data
x = data['News']
y = data['PriceSentiment']
x, x_test, y, y_test = train_test_split(x,y, stratify=y, test_size=0.2, random_state=100)

[107] ✓ 0.5s

Vectorize text reviews to numbers.

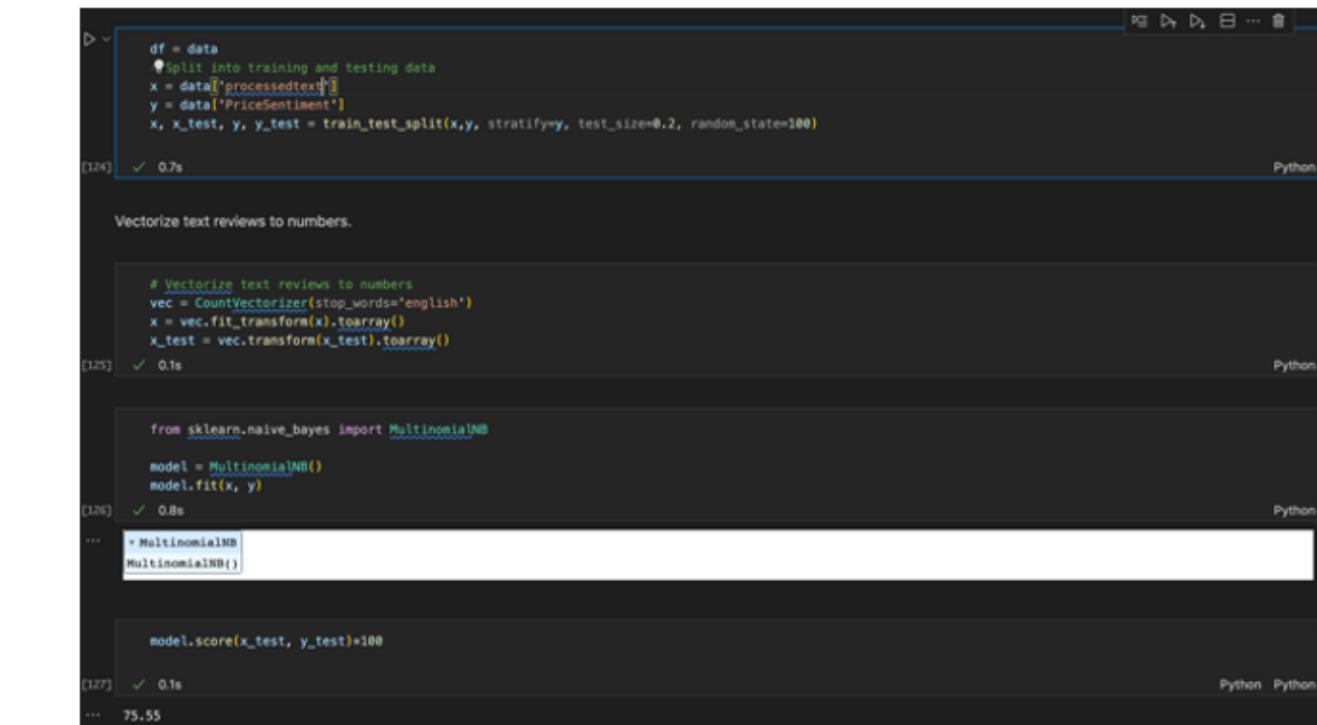
# Vectorize text reviews to numbers
vec = CountVectorizer(stop_words='english')
x = vec.fit_transform(x).toarray()
x_test = vec.transform(x_test).toarray()

[108] ✓ 0.2s

from sklearn.naive_bayes import MultinomialNB
model = MultinomialNB()
model.fit(x, y)

[109] ✓ 1.6s
... + MultinomialNB
MultinomialNB()

D: model.score(x_test, y_test)*100
[110] ✓ 0.1s
... 76.35
```



```
df = data
# Split into training and testing data
x = data['processedtext']
y = data['PriceSentiment']
x, x_test, y, y_test = train_test_split(x,y, stratify=y, test_size=0.2, random_state=100)

[124] ✓ 0.7s

Vectorize text reviews to numbers.

# Vectorize text reviews to numbers
vec = CountVectorizer(stop_words='english')
x = vec.fit_transform(x).toarray()
x_test = vec.transform(x_test).toarray()

[125] ✓ 0.1s

from sklearn.naive_bayes import MultinomialNB
model = MultinomialNB()
model.fit(x, y)

[126] ✓ 0.8s
... + MultinomialNB
MultinomialNB()

D: model.score(x_test, y_test)*100
[127] ✓ 0.1s
... 75.55
```

เป็นการเตรียมเทียบระหว่างใช้ NLP และ ไม่ใช้ NLP ของในระบบ

ความแม่นยำโดยไม่ใช้NLP 76.35%

ความแม่นยำโดยผ่านกระบวนการ NLP 75.55%

ความก้าวหน้าการดำเนินงาน

ความก้าวหน้า 7 การปรีปเทียบ โดยใช้ Model Naïve baye

```
Then check prediction...

from itertools import count
import pandas as pd
df = pd.read_csv('gold-dataset-sinha-khandait.csv', sep=',', header=None)
start = 10000
end = 10570
df = df[start - 1:end - 1]
correct = 0
for i in range(len(df)):
    if df.values[i][2]:
        (model.predict(vec.transform([df.values[i][2]])), df.values[i][9] == model.predict(vec.transform([df.values[i][2]])))
        if df.values[i][9] == model.predict(vec.transform([df.values[i][2]])):
            correct += 1
print(correct / len(df) * 100 )
[166] ✓ 0.5s
... 75.26315789473685
```

```
✓ from itertools import count
import pandas as pd
df = pd.read_csv('gold-dataset-sinha-khandait.csv', sep=',', header=None)
start = 10000
end = 10570
df = df[start - 1:end - 1]
correct = 0
for i in range(len(df)):
    (df.values[i][2])
    (model.predict(vec.transform([df.values[i][2]])), df.values[i][9] == model.predict(vec.transform([df.values[i][2]])))
    if df.values[i][9] == model.predict(vec.transform([df.values[i][2]])):
        correct += 1
print(correct / len(df) * 100 )
[181] ✓ 0.5s
... 64.3859649122807
```

การทดสอบ โดยการนำป่าวสาร 570 record มาทดสอบจริง

ความแม่นยำโดยผ่านกระบวนการ NLP 64%

ความแม่นยำโดยไม่ใช้NLP 75.26%

ความก้าวหน้าการดำเนินงาน

ความก้าวหน้า 8 ทดลองทำ pywebio

```
64 | ##pywebio##
65 | from pywebio.input import input, FLOAT, TEXT
66 | from pywebio.output import put_text
67 | def main():
68 |     ## Model Generation
69 |     from sklearn.naive_bayes import MultinomialNB
70 |
71 |     model = MultinomialNB()
72 |     model.fit(x, y)
73 |     from itertools import count
74 |     import pandas as pd
75 |     df = pd.read_csv('gold-dataset-sinha-khandait.csv', sep=',', header=None)
76 |     start = 10000
77 |     end = 10570
78 |     df = df[start - 1:end - 1]
79 |     correct = 0
80 |     str = input('This is label', type=TEXT, placeholder='This is placeholder',
81 |                help_text='This is help text', required=True)
82 |     put_text(model.predict(vec.transform([str])))
83 |     for i in range(len(df)):
84 |         print(df.values[i][2])
85 |         put_text(model.predict(vec.transform([df.values[i][2]])), df.values[i][9] == model.predict(vec.transform([df.values[i][2]])))
86 |
87 |         if df.values[i][9] == model.predict(vec.transform([df.values[i][2]])):
88 |             correct += 1
89 |
90 |     print(correct / len(df) * 100)
91 |
92 |
93 |     if __name__ == '__main__':
94 |         main()
95 |
96 |
```

This is Text

Final Services PMI

This is help text

Submit Reset

PyWebIO Application

Gmail YouTube บล็อกนี้ python - จัดการข้อมูล Nested Lists

['negative']
['none'] [True]
['positive'] [True]
['positive'] [True]
['negative'] [True]
['positive'] [True]
['positive'] [True]
['none'] [True]

สรุป

สรุปผลการดำเนินงาน

- สามารถสรุปบทความได้ ใช้ Library NLTK NLP
- สามารถวิเคราะห์ บทความได้ระดับเบื้องต้น
- จัดทำเหมืองข้อมูล Dataset
- เปรียบเทียบอัลกอริทึมระหว่าง Random Forest Classifier กับ Naïve Bayes Classifier
- ทำ Sentiment Analysis Naïve Bayes Classifier
- ทำการเปรียบเทียบระหว่างข่าวสารดั้งเดิม และ ข่าวสารที่ผ่านระบบ NLP
- มีการนำ pywebio เข้ามาในการทำ Web Application

ปัญหาและอุปสรรค

- Library ที่จะใช้ค่อนข้างเยอะ จึงทำให้เริ่ม Project ช้า
- ไม่สามารถนำ ตัวเลขที่ข่าวประกาศไปวิเคราะห์ หรือ เก็บข้อมูลได้เนื่องจากตอนนี้การ trenนิ่ง ต้องใช้ตัวอักษรภาษาอังกฤษเท่านั้น
- ผลจากการทำ Sentiment Analysis ผลที่ได้รับ ไม่ตรงกับที่เหลื่ในระบบ
- เนื่องจากข่าวสารที่ใช้ในการทดสอบ มีตัวเลขและบางประโยคสำคัญทำให้ NLP ตัดบางส่วนที่สำคัญไป

งานที่จะดำเนินการต่อไป

- จัดทำ Web Applicationให้สมบูรณ์
- ศึกษา Framework Django

เอกสารอ้างอิง

- 
- [1] S. Shalev-Shwartz, S. Ben-David, **Understanding Machine Learning: From Theory to Algorithms** (2014), Cambridge University Press , เข้าถึงล่าสุด 15 มกราคม 2565
 - [2] T. Mikolov, I. Sutskever, K. Chen, G. S Corrado, and J. Dean. **Distributed Representations of Words and Phrases and their Compositionality** (2013), Advances in Neural Information Processing Systems 26 เข้าถึงล่าสุด 15 มกราคม 2565
 - [3] J. Pennington, R. Socher, and C. D. Manning, **GloVe: Global Vectors for Word Representation** (2014), In EMNLP. เข้าถึงล่าสุด 16 มกราคม 2565
 - [4] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov. **Enriching word vectors with subword information** (2016), arXiv preprint เข้าถึงล่าสุด 17 มกราคม 2565
 - [5] NLP Implementations : URL :
<https://towardsdatascience.com/gentle-start-to-natural-language-processing-using-python-6e46c07addf3> เข้าถึงล่าสุด 18 มกราคม 2565

เอกสารอ้างอิง

- [6] The theory you need to know before you start an NLP : URL :
<https://towardsdatascience.com/the-theory-you-need-to-know-before-you-start-an-nlp-project-1890f5bbb793> เข้าถึงล่าสุด 12 มีนาคม 2565
- [7] Us Department of labor : URL : <https://www.dol.gov/> เข้าถึงล่าสุด 12 มีนาคม 2565
- [8] Energy information Administration : URL : <https://www.eia.gov/> เข้าถึงล่าสุด 12 มีนาคม 2565
- [9] กองทุน SPDR : URL : <https://traderider.com/forex/spdr-%EO%B8%81%EO%B8%AD%EO%B8%87%EO%B8%97%EO%B8%B8%EO%B8%99%EO%B8%97%EO%B8%AD%EO%B8%87%EO%B8%84%EO%B8%B3%EO%B9%81%EO%B8%97%EO%B9%88%EO%B8%87> เข้าถึงล่าสุด 12 มีนาคม 2565
- [10] Federal Reserve : URL : <https://www.federalreserve.gov/> เข้าถึงล่าสุด 12 มีนาคม 2565
- [11] Bloomberg : URL : <https://www.bloomberg.com/asia> เข้าถึงล่าสุด 12 มีนาคม 2565
- [12] Twitter : URL : <https://twitter.com/> เข้าถึงล่าสุด 12 มีนาคม 2562

คำถ้ามและข้อเสนอแนะ

