

References

<https://andhint.github.io/machine-learning/nlp/Feature-Extraction-From-Text/> (<https://andhint.github.io/machine-learning/nlp/Feature-Extraction-From-Text/>)

Tweet Sentiment Classification

Objective

Determine whether the sentiment on the Tweets related to a product or company is positive, negative or neutral.

Classification Techniques

A1. Supervised Learning

* Step 1: Using Manually Labelled Tweet Sentiments for Supervised Training*

```
In [1]: 1 from IPython.core.interactiveshell import InteractiveShell
        2 InteractiveShell.ast_node_interactivity = "all"
        3 from IPython.display import HTML, display
```

```
In [7]: 1 import pandas as pd
        2 import numpy as np
        3
        4 np.random.seed(7)
        5
        6 print("[INFO] Data Frame created with Manually-labelled data for individual stocks")
        7 aapl=pd.read_csv("/Users/hardeepsingh/Desktop/IS/University/sem-4/FE-520/Project/Manually_labelle
        8
        9 print("\n[INFO] Apple Data Frame - Manually Labelled Sentiments")
        10 aapl.info()
        11
        12 aapl=aapl.dropna()
```

```
[INFO] Data Frame created with Manually-labelled data for individual stocks
```

```
[INFO] Apple Data Frame - Manually Labelled Sentiments
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 10481 entries, 0 to 10480
```

```
Data columns (total 3 columns):
```

```
date          3977 non-null object
```

```
sentiment     3969 non-null float64
```

```
text          3977 non-null object
```

```
dtypes: float64(1), object(2)
```

```
memory usage: 245.7+ KB
```

* Step 2: Splitting the Data into Train (for Model training) and Test (for validation)*

```
In [9]: 1 from sklearn.model_selection import train_test_split
        2
        3 text=aapl.text
        4 target=aapl.sentiment
        5
        6 X_train,X_test,y_train,y_test=train_test_split(text,target, random_state=27,test_size=0.2)
```

* Step 3: Selecting the best Model*

WordCloud Of Tweets


```

In [13]: 1 # importing required libraries
2 from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
3 from sklearn.pipeline import Pipeline
4 from sklearn.linear_model import LogisticRegression
5 import warnings
6 warnings.filterwarnings('ignore')
7
8 # import GridSearch
9 from sklearn.model_selection import GridSearchCV
10
11 # To store the results
12 classifier_results={}
13
14 # using count vectorizer instead of tfidf
15 # tokenizing only alpha numeric
16 tokenPatt = '[A-Za-z0-9]+(=?\\s+)'
17
18 # pipeline which does two steps all together:
19 # (1) generate CountVector, and (2) train classifier
20 # each step is named, i.e. "vec", "clf"
21 pl_1 = Pipeline([
22     ('tfidf', CountVectorizer(token_pattern = tokenPatt)),
23     ('clf', LogisticRegression())])
24
25 pl_1.fit(X_train,y_train)
26
27 # accuracy
28 accuracy = pl_1.score(X_test,y_test)
29 print ("Untuned Accuracy of Logistic Regression using CountVectorizer: ", accuracy)
30
31 classifier_results["Untuned Accuracy of Logistic Regression using CountVectorizer"]=accuracy
32
33
34 # Parameters to be used for Tuning
35 parameters = {'tfidf__min_df':[2,3],
36               'tfidf__token_pattern':[ '[A-Za-z0-9]+(=?\\s+)' ],
37               'tfidf__stop_words':[None,"english"],
38               'clf__solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']
39 }
40
41 # the metric used to select the best parameters
42 metric = "f1_macro"
43
44 # GridSearch also uses cross validation
45 gs_clf = GridSearchCV(pl_1, param_grid=parameters, scoring=metric, cv=5)
46 gs_clf = gs_clf.fit(text, target)
47
48 # gs_clf.best_params_ returns a dictionary
49 # with parameter and its best value as an entry
50
51 for param_name in gs_clf.best_params_:
52     print(param_name,": ",gs_clf.best_params_[param_name])
53
54 print("Best f1 score:", gs_clf.best_score_)
55
56

```

```

Out[13]: Pipeline(memory=None,
    steps=[('tfidf', CountVectorizer(analyzer='word', binary=False, decode_error='strict',
    dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
    lowercase=True, max_df=1.0, max_features=None, min_df=1,
    ngram_range=(1, 1), preprocessor=None, stop_words=None,
    stri...penalty='l2', random_state=None, solver='warn',
    tol=0.0001, verbose=0, warm_start=False))])

```

```

Untuned Accuracy of Logistic Regression using CountVectorizer:  0.672544080604534
clf__solver :  newton-cg
tfidf__min_df :  3
tfidf__stop_words :  None
tfidf__token_pattern :  [A-Za-z0-9]+(=?\s+)
Best f1 score: 0.4268043195426816

```

```

In [14]: 1 # Using Parameters from above GridSearch Result
2 pl_2 = Pipeline([
3     ('tfidf', CountVectorizer(stop_words=None, token_pattern='[A-Za-z0-9]+(?=\s+)', min_df=2))
4     ('clf', LogisticRegression(solver='sag'))
5 ])
6
7 pl_2.fit(X_train, y_train)
8
9 # accuracy
10 accuracy = pl_2.score(X_test, y_test)
11 print("Tuned Accuracy of Logistic Regression using CountVectorizer: ", accuracy)
12
13 classifier_results["Tuned Accuracy of Logistic Regression using CountVectorizer"]=accuracy

```

```

Out[14]: Pipeline(memory=None,
    steps=[('tfidf', CountVectorizer(analyzer='word', binary=False, decode_error='strict',
    dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
    lowercase=True, max_df=1.0, max_features=None, min_df=2,
    ngram_range=(1, 1), preprocessor=None, stop_words=None,
    stri... penalty='l2', random_state=None, solver='sag',
    tol=0.0001, verbose=0, warm_start=False))])

```

Tuned Accuracy of Logistic Regression using CountVectorizer: 0.6662468513853904

**** ii. CountVectorizer + DecisionTreeClassifier ****

```

In [15]: 1 # importing required libraries
2 from sklearn.tree import DecisionTreeClassifier
3
4 p2_1 = Pipeline([
5     ('tfidf', CountVectorizer(token_pattern = tokenPatt)),
6     ('clf', DecisionTreeClassifier())
7 ])
8
9 p2_1.fit(X_train,y_train)
10
11 # accuracy
12 accuracy = p2_1.score(X_test,y_test)
13 print("Untuned Accuracy of Decision Tree using CountVectorizer: ", accuracy)
14
15 classifier_results["Untuned Accuracy of Decision Tree using CountVectorizer"]=accuracy
16
17 # Parameters to be used for Tuning
18 parameters = {'tfidf__min_df':[2,3],
19               'tfidf__token_pattern':['[A-Za-z0-9]+(=?\\s+)'],
20               'tfidf__stop_words':[None,"english"]}
21 }
22
23 # the metric used to select the best parameters
24 metric = "f1_macro"
25
26 # GridSearch also uses cross validation
27 gs_clf = GridSearchCV(p2_1, param_grid=parameters, scoring=metric, cv=5)
28 gs_clf = gs_clf.fit(text, target)
29
30 # gs_clf.best_params_ returns a dictionary
31 # with parameter and its best value as an entry
32
33 for param_name in gs_clf.best_params_:
34     print(param_name,": ",gs_clf.best_params_[param_name])
35
36 print("Best f1 score:", gs_clf.best_score_)

```

```

Out[15]: Pipeline(memory=None,
                  steps=[('tfidf', CountVectorizer(analyzer='word', binary=False, decode_error='strict',
                  dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
                  lowercase=True, max_df=1.0, max_features=None, min_df=1,
                  ngram_range=(1, 1), preprocessor=None, stop_words=None,
                  stri... min_weight_fraction_leaf=0.0, presort=False, random_state=None,
                  splitter='best'))])

Untuned Accuracy of Decision Tree using CountVectorizer:  0.6158690176322418
tfidf__min_df : 2
tfidf__stop_words : english
tfidf__token_pattern : [A-Za-z0-9]+(=?\s+)
Best f1 score: 0.40682321943433036

```

```

In [52]: 1 # Using Parameters from above GridSearch Result
2 p2_2 = Pipeline([
3     ('tfidf', CountVectorizer(stop_words='english', token_pattern='[A-Za-z0-9]+(?=\\s+)', min_c
4     ('clf', DecisionTreeClassifier()))
5 ])
6
7 p2_2.fit(X_train,y_train)
8
9 # accuracy
10 accuracy = p2_2.score(X_test,y_test)
11 print("Tuned Accuracy of Decision Tree using CountVectorizer: ", accuracy)
12
13 classifier_results["Tuned Accuracy of Decision Tree using CountVectorizer"]=accuracy

```

```

Out[52]: Pipeline(memory=None,
      steps=[('tfidf', CountVectorizer(analyzer=u'word', binary=False, decode_error=u'strict',
      dtype=<type 'numpy.int64'>, encoding=u'utf-8', input=u'content',
      lowercase=True, max_df=1.0, max_features=None, min_df=3,
      ngram_range=(1, 1), preprocessor=None, stop_words='english',
      ...      min_weight_fraction_leaf=0.0, presort=False, random_state=None,
      splitter='best'))])

('Tuned Accuracy of Decision Tree using CountVectorizer: ', 0.5462012320328542)

** iii. CountVectorizer + MultinomialNB **

```

```

In [16]: 1 # importing required libraries
2 from sklearn.naive_bayes import MultinomialNB
3
4 p3_1 = Pipeline([
5     ('tfidf', CountVectorizer()),
6     ('clf', MultinomialNB())
7 ])
8
9 p3_1.fit(X_train,y_train)
10
11 # accuracy
12 accuracy = p3_1.score(X_test,y_test)
13 print("Untuned Accuracy of MultinomialNB using CountVectorizer: ", accuracy)
14
15 classifier_results["Untuned Accuracy of MultinomialNB using CountVectorizer"]=accuracy
16
17 # Parameters to be used for Tuning
18 parameters = {'tfidf__min_df':[2,3],
19               'tfidf__token_pattern':['[A-Za-z0-9]+(=?\\s+)'],
20               'tfidf__stop_words':[None,"english"],
21               'clf__alpha': [0.5,1.0,2.0],
22               }
23
24 # the metric used to select the best parameters
25 metric = "f1_macro"
26
27 # GridSearch also uses cross validation
28 gs_clf = GridSearchCV(p3_1, param_grid=parameters, scoring=metric, cv=5)
29 gs_clf = gs_clf.fit(text, target)
30
31 # gs_clf.best_params_ returns a dictionary
32 # with parameter and its best value as an entry
33
34 for param_name in gs_clf.best_params_:
35     print(param_name,": ",gs_clf.best_params_[param_name])
36
37 print("Best f1 score:", gs_clf.best_score_)

```

```

Out[16]: Pipeline(memory=None,
               steps=[('tfidf', CountVectorizer(analyzer='word', binary=False, decode_error='strict',
dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
lowercase=True, max_df=1.0, max_features=None, min_df=1,
ngram_range=(1, 1), preprocessor=None, stop_words=None,
strip_accents=None, token_pattern='(?u)\\b\\w\\w+\\b',
tokenizer=None, vocabulary=None)), ('clf', MultinomialNB(alpha=1.0, class_prior=None, fit_p
rior=True))])

Untuned Accuracy of MultinomialNB using CountVectorizer:  0.6649874055415617
clf__alpha :  1.0
tfidf__min_df :  3
tfidf__stop_words :  None
tfidf__token_pattern :  [A-Za-z0-9]+(=?\\s+)
Best f1 score: 0.4646847144630701

```

```

In [18]: 1 # Using Parameters from above GridSearch Result
2 p3_2 = Pipeline([
3     ('tfidf', CountVectorizer(stop_words=None, token_pattern='[A-Za-z0-9]+(?=\\s+)', min_df=2))
4     ('clf', MultinomialNB(alpha=0.5))
5 ])
6
7 p3_2.fit(X_train, y_train)
8
9 # accuracy
10 accuracy = p3_2.score(X_test, y_test)
11 print("Tuned Accuracy of MultinomialNB using CountVectorizer: ", accuracy)
12
13 classifier_results["Tuned Accuracy of MultinomialNB using CountVectorizer"]=accuracy

```

```

Out[18]: Pipeline(memory=None,
      steps=[('tfidf', CountVectorizer(analyzer='word', binary=False, decode_error='strict',
      dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
      lowercase=True, max_df=1.0, max_features=None, min_df=2,
      ngram_range=(1, 1), preprocessor=None, stop_words=None,
      strip_accents=None, token_pattern='[A-Za-z0-9]+(?=\\s+)',
      tokenizer=None, vocabulary=None)), ('clf', MultinomialNB(alpha=0.5, class_prior=None, fit_p
      rior=True))])

Tuned Accuracy of MultinomialNB using CountVectorizer:  0.6435768261964736

```

iv. CountVectorizer + LinearSVC

In [19]:

```
1 # importing required libraries
2 from sklearn.svm import LinearSVC
3
4 p4_1 = Pipeline([
5     ('tfidf', CountVectorizer()),
6     ('clf', LinearSVC())
7 ])
8
9 p4_1.fit(X_train,y_train)
10
11 # accuracy
12 accuracy = p4_1.score(X_test,y_test)
13 print("Untuned Accuracy of LinearSVC using CountVectorizer: ", accuracy)
14
15 classifier_results["Untuned Accuracy of LinearSVC using CountVectorizer"]=accuracy
16
17 # Parameters to be used for Tuning
18 parameters = {'tfidf__min_df':[2,3],
19               'tfidf__token_pattern':['[A-Za-z0-9]+(=?\\s+)' ],
20               'tfidf__stop_words':[None,"english"],
21               'clf__loss':['hinge','squared_hinge']}
22 }
23
24 # the metric used to select the best parameters
25 metric = "f1_macro"
26
27 # GridSearch also uses cross validation
28 gs_clf = GridSearchCV(p4_1, param_grid=parameters, scoring=metric, cv=5)
29 gs_clf = gs_clf.fit(text, target)
30
31 # gs_clf.best_params_ returns a dictionary
32 # with parameter and its best value as an entry
33
34 for param_name in gs_clf.best_params_:
35     print(param_name,": ",gs_clf.best_params_[param_name])
36
37 print("Best f1 score:", gs_clf.best_score_)
```

Out[19]: Pipeline(memory=None,
 steps=[('tfidf', CountVectorizer(analyzer='word', binary=False, decode_error='strict',
 dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
 lowercase=True, max_df=1.0, max_features=None, min_df=1,
 ngram_range=(1, 1), preprocessor=None, stop_words=None,
 stri...ax_iter=1000,
 multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,
 verbose=0)))

Untuned Accuracy of LinearSVC using CountVectorizer: 0.6511335012594458

clf__loss : hinge

tfidf__min_df : 3

tfidf__stop_words : english

tfidf__token_pattern : [A-Za-z0-9]+(=?\\s+)

Best f1 score: 0.43101033304838277

```

In [20]: 1 # Using Parameters from above GridSearch Result
2 p4_2 = Pipeline([
3     ('tfidf', CountVectorizer(stop_words=None, token_pattern='[A-Za-z0-9]+(?=\\s+)', min_df=2))
4     ('clf', LinearSVC(loss='hinge'))
5 ])
6
7 p4_2.fit(X_train, y_train)
8
9 # accuracy
10 accuracy = p4_2.score(X_test, y_test)
11 print("Tuned Accuracy of LinearSVC using CountVectorizer: ", accuracy)
12
13 classifier_results["Tuned Accuracy of LinearSVC using CountVectorizer"] = accuracy

```

```

Out[20]: Pipeline(memory=None,
      steps=[('tfidf', CountVectorizer(analyzer='word', binary=False, decode_error='strict',
      dtype=<class 'numpy.int64'>, encoding='utf-8', input='content',
      lowercase=True, max_df=1.0, max_features=None, min_df=2,
      ngram_range=(1, 1), preprocessor=None, stop_words=None,
      stri...e', max_iter=1000, multi_class='ovr',
      penalty='l2', random_state=None, tol=0.0001, verbose=0))])

```

Tuned Accuracy of LinearSVC using CountVectorizer: 0.6511335012594458

v. TfidfVectorizer + Logistic Regression

```

In [21]: 1 p5_1 = Pipeline([
2         ('tfidf', TfidfVectorizer()),
3         ('clf', LogisticRegression())
4     ])
5
6 p5_1.fit(X_train,y_train)
7
8
9 # accuracy
10 accuracy = p5_1.score(X_test,y_test)
11 print ("Untuned Accuracy of Logistic Regression using TfidfVectorizer: ", accuracy)
12
13 classifier_results["Untuned Accuracy of Logistic Regression using TfidfVectorizer"]=accuracy
14
15 # Parameters to be used for Tuning
16 parameters = {'tfidf__min_df':[2,3],
17               'tfidf__token_pattern':['[A-Za-z0-9]+(=?\\s+)'],
18               'tfidf__stop_words':[None,"english"],
19               'clf__solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga']}
20 }
21
22 # the metric used to select the best parameters
23 metric = "f1_macro"
24
25 # GridSearch also uses cross validation
26 gs_clf = GridSearchCV(p5_1, param_grid=parameters, scoring=metric, cv=5)
27 gs_clf = gs_clf.fit(text, target)
28
29 # gs_clf.best_params_ returns a dictionary
30 # with parameter and its best value as an entry
31
32 for param_name in gs_clf.best_params_:
33     print(param_name,": ",gs_clf.best_params_[param_name])
34
35 print("Best f1 score:", gs_clf.best_score_)

```

```

Out[21]: Pipeline(memory=None,
      steps=[('tfidf', TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
      dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
      lowercase=True, max_df=1.0, max_features=None, min_df=1,
      ngram_range=(1, 1), norm='l2', preprocessor=None, smooth_idf=True,...penalty='l2', random_s
      tate=None, solver='warn',
      tol=0.0001, verbose=0, warm_start=False))])

Untuned Accuracy of Logistic Regression using TfidfVectorizer:  0.672544080604534
clf__solver :  saga
tfidf__min_df :  3
tfidf__stop_words :  None
tfidf__token_pattern :  [A-Za-z0-9]+(=?\s+)
Best f1 score: 0.3880592998717089

```

```

In [22]: 1 # Using Parameters from above GridSearch Result
2 p5_2 = Pipeline([
3     ('tfidf', TfidfVectorizer(stop_words=None, token_pattern='[A-Za-z0-9]+(?=\\s+)', min_df=2))
4     ('clf', LogisticRegression(solver='newton-cg'))
5 ])
6
7 p5_2.fit(X_train, y_train)
8
9 # accuracy
10 accuracy = p5_2.score(X_test, y_test)
11 print("Tuned Accuracy of Logistic Regression using TfidfVectorizer: ", accuracy)
12
13 classifier_results["Tuned Accuracy of Logistic Regression using TfidfVectorizer"]=accuracy

```

```

Out[22]: Pipeline(memory=None,
      steps=[('tfidf', TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
      dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
      lowercase=True, max_df=1.0, max_features=None, min_df=2,
      ngram_range=(1, 1), norm='l2', preprocessor=None, smooth_idf=True, ...ty='l2', random_state=
      None, solver='newton-cg',
      tol=0.0001, verbose=0, warm_start=False))])

```

Tuned Accuracy of Logistic Regression using TfidfVectorizer: 0.6700251889168766

vi. TfidfVectorizer + DecisionTreeClassifier

```
In [23]: 1 p6_1 = Pipeline([
2         ('tfidf', TfidfVectorizer()),
3         ('clf', DecisionTreeClassifier())
4     ])
5
6 p6_1.fit(X_train,y_train)
7
8 # accuracy
9 accuracy = p6_1.score(X_test,y_test)
10 print("Untuned Accuracy of Decision Tree using TfidfVectorizer: ", accuracy)
11
12 classifier_results["Untuned Accuracy of Decision Tree using TfidfVectorizer"]=accuracy
13
14 # Parameters to be used for Tuning
15 parameters = {'tfidf__min_df':[2,3],
16               'tfidf__token_pattern':['[A-Za-z0-9]+(=?\\s+)'],
17               'tfidf__stop_words':[None,"english"]}
18 }
19
20 # the metric used to select the best parameters
21 metric = "f1_macro"
22
23 # GridSearch also uses cross validation
24 gs_clf = GridSearchCV(p6_1, param_grid=parameters, scoring=metric, cv=5)
25 gs_clf = gs_clf.fit(text, target)
26
27 # gs_clf.best_params_ returns a dictionary
28 # with parameter and its best value as an entry
29
30 for param_name in gs_clf.best_params_:
31     print(param_name,": ",gs_clf.best_params_[param_name])
32
33 print("Best f1 score:", gs_clf.best_score_)
```

```
Out[23]: Pipeline(memory=None,
      steps=[('tfidf', TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
      dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
      lowercase=True, max_df=1.0, max_features=None, min_df=1,
      ngram_range=(1, 1), norm='l2', preprocessor=None, smooth_idf=True,...      min_weight_fract
      ion_leaf=0.0, presort=False, random_state=None,
      splitter='best'))])

Untuned Accuracy of Decision Tree using TfidfVectorizer:  0.6246851385390428
tfidf__min_df :  2
tfidf__stop_words :  english
tfidf__token_pattern :  [A-Za-z0-9]+(=?\\s+)
Best f1 score: 0.40299918571643156
```

```
In [24]: 1 # Using Parameters from above GridSearch Result
2 p6_2 = Pipeline([
3         ('tfidf', TfidfVectorizer(stop_words='english',token_pattern='[A-Za-z0-9]+(=?\\s+)',min_c
4         ('clf', DecisionTreeClassifier())
5     ])
6
7 p6_2.fit(X_train,y_train)
8
9 # accuracy
10 accuracy = p6_2.score(X_test,y_test)
11 print("Tuned Accuracy of Decision Tree using TfidfVectorizer: ", accuracy)
12
13 classifier_results["Tuned Accuracy of Decision Tree using TfidfVectorizer"]=accuracy
```

```
Out[24]: Pipeline(memory=None,
      steps=[('tfidf', TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
      dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
      lowercase=True, max_df=1.0, max_features=None, min_df=3,
      ngram_range=(1, 1), norm='l2', preprocessor=None, smooth_idf=True,...      min_weight_fract
      ion_leaf=0.0, presort=False, random_state=None,
      splitter='best'))])

Tuned Accuracy of Decision Tree using TfidfVectorizer:  0.628463476070529
```

```
In [25]: 1 p7_1 = Pipeline([
2         ('tfidf', TfidfVectorizer()),
3         ('clf', MultinomialNB())
4     ])
5
6 p7_1.fit(X_train,y_train)
7
8 # accuracy without parameter tuning
9 accuracy = p7_1.score(X_test,y_test)
10 print("Untuned Accuracy of MultinomialNB using TfidfVectorizer: ", accuracy)
11
12 classifier_results["Untuned Accuracy of MultinomialNB using TfidfVectorizer"]=accuracy
13
14 # Parameters to be used for Tuning
15 parameters = {'tfidf__min_df':[2,3],
16               'tfidf__token_pattern':['[A-Za-z0-9]+(?:\\s+)''],
17               'tfidf__stop_words':[None,"english"],
18               'clf__alpha': [0.5,1.0,2.0],
19             }
20
21 # the metric used to select the best parameters
22 metric = "f1_macro"
23
24 # GridSearch also uses cross validation
25 gs_clf = GridSearchCV(p7_1, param_grid=parameters, scoring=metric, cv=5)
26 gs_clf = gs_clf.fit(text, target)
27
28 # gs_clf.best_params_ returns a dictionary
29 # with parameter and its best value as an entry
30
31 for param_name in gs_clf.best_params_:
32     print(param_name,": ",gs_clf.best_params_[param_name])
33
34 print("Best f1 score:", gs_clf.best_score_)
```

```
Out[25]: Pipeline(memory=None,
      steps=[('tfidf', TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
      dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
      lowercase=True, max_df=1.0, max_features=None, min_df=1,
      ngram_range=(1, 1), norm='l2', preprocessor=None, smooth_idf=True,...rue,
      vocabulary=None)), ('clf', MultinomialNB(alpha=1.0, class_prior=None, fit_prior=True))])

Untuned Accuracy of MultinomialNB using TfidfVectorizer:  0.6599496221662469
clf__alpha :  0.5
tfidf__min_df :  3
tfidf__stop_words :  english
tfidf__token_pattern :  [A-Za-z0-9]+(?:\\s+)
Best f1 score: 0.40078895416122884
```

```

In [26]: 1 # Using Parameters from above GridSearch Result
2 p7_2 = Pipeline([
3     ('tfidf', TfidfVectorizer(stop_words='english', token_pattern='[A-Za-z0-9]+(?!\\s+)', min_
4     ('clf', MultinomialNB(alpha=0.5))
5 ])
6
7 p7_2.fit(X_train,y_train)
8
9 # accuracy
10 accuracy = p7_2.score(X_test,y_test)
11 print("Tuned Accuracy of MultinomialNB using TfidfVectorizer: ", accuracy)
12
13 classifier_results["Tuned Accuracy of MultinomialNB using TfidfVectorizer"]=accuracy

```

```

Out[26]: Pipeline(memory=None,
      steps=[('tfidf', TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
      dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
      lowercase=True, max_df=1.0, max_features=None, min_df=3,
      ngram_range=(1, 1), norm='l2', preprocessor=None, smooth_idf=True,...rue,
      vocabulary=None)), ('clf', MultinomialNB(alpha=0.5, class_prior=None, fit_prior=True))])

```

Tuned Accuracy of MultinomialNB using TfidfVectorizer: 0.6662468513853904

viii. TfidfVectorizer + LinearSVC

In [27]:

```
1 p8_1 = Pipeline([
2     ('tfidf', TfidfVectorizer()),
3     ('clf', LinearSVC())
4 ])
5
6 p8_1.fit(X_train,y_train)
7
8 # accuracy
9 accuracy = p8_1.score(X_test,y_test)
10 print("Untuned Accuracy of LinearSVC using TfidfVectorizer: ", accuracy)
11
12 classifier_results["Untuned Accuracy of LinearSVC using TfidfVectorizer"]=accuracy
13
14 # Parameters to be used for Tuning
15 parameters = {'tfidf__min_df':[2,3],
16               'tfidf__token_pattern':['[A-Za-z0-9]+(=?\\s+)'],
17               'tfidf__stop_words':[None,"english"],
18               'clf__loss':['hinge','squared_hinge']}
19 }
20
21 # the metric used to select the best parameters
22 metric = "f1_macro"
23
24 # GridSearch also uses cross validation
25 gs_clf = GridSearchCV(p8_1, param_grid=parameters, scoring=metric, cv=5)
26 gs_clf = gs_clf.fit(text, target)
27
28 # gs_clf.best_params_ returns a dictionary
29 # with parameter and its best value as an entry
30
31 for param_name in gs_clf.best_params_:
32     print(param_name,": ",gs_clf.best_params_[param_name])
33
34 print("Best f1 score:", gs_clf.best_score_)
```

Out[27]: Pipeline(memory=None,
 steps=[('tfidf', TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
 dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
 lowercase=True, max_df=1.0, max_features=None, min_df=1,
 ngram_range=(1, 1), norm='l2', preprocessor=None, smooth_idf=True,...ax_iter=1000,
 multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,
 verbose=0))])

Untuned Accuracy of LinearSVC using TfidfVectorizer: 0.6687657430730478
clf__loss : squared_hinge
tfidf__min_df : 3
tfidf__stop_words : english
tfidf__token_pattern : [A-Za-z0-9]+(=?\s+)
Best f1 score: 0.4317829573509598


```
In [28]: 1 # Using Parameters from above GridSearch Result
2 p8_2 = Pipeline([
3     ('tfidf', TfidfVectorizer(stop_words=None, token_pattern='[A-Za-z0-9]+(?=\s+)', min_df=2))
4     ('clf', LinearSVC(loss='hinge'))
5 ])
6
7 p8_2.fit(X_train, y_train)
8
9 # accuracy
10 accuracy = p8_2.score(X_test, y_test)
11 print("Tuned Accuracy of LinearSVC using TfidfVectorizer: ", accuracy)
12
13 classifier_results["Tuned Accuracy of LinearSVC using TfidfVectorizer"]=accuracy
```

```
Out[28]: Pipeline(memory=None,
      steps=[('tfidf', TfidfVectorizer(analyzer='word', binary=False, decode_error='strict',
      dtype=<class 'numpy.float64'>, encoding='utf-8', input='content',
      lowercase=True, max_df=1.0, max_features=None, min_df=2,
      ngram_range=(1, 1), norm='l2', preprocessor=None, smooth_idf=True,...e', max_iter=1000, mul
      ti_class='ovr',
      penalty='l2', random_state=None, tol=0.0001, verbose=0))])

Tuned Accuracy of LinearSVC using TfidfVectorizer:  0.6738035264483627
```

A2. Supervised Learning - Model Comparison Analysis

```
In [30]: 1 for key in classifier_results:
2     print('{}={}'.format(key, classifier_results[key]))
```

```
Untuned Accuracy of Logistic Regression using CountVectorizer=0.672544080604534
Tuned Accuracy of Logistic Regression using CountVectorizer=0.6662468513853904
Untuned Accuracy of Decision Tree using CountVectorizer=0.6158690176322418
Untuned Accuracy of MultinomialNB using CountVectorizer=0.6649874055415617
Tuned Accuracy of MultinomialNB using CountVectorizer=0.6435768261964736
Untuned Accuracy of LinearSVC using CountVectorizer=0.6511335012594458
Tuned Accuracy of LinearSVC using CountVectorizer=0.6511335012594458
Untuned Accuracy of Logistic Regression using TfidfVectorizer=0.672544080604534
Tuned Accuracy of Logistic Regression using TfidfVectorizer=0.6700251889168766
Untuned Accuracy of Decision Tree using TfidfVectorizer=0.6246851385390428
Tuned Accuracy of Decision Tree using TfidfVectorizer=0.628463476070529
Untuned Accuracy of MultinomialNB using TfidfVectorizer=0.6599496221662469
Tuned Accuracy of MultinomialNB using TfidfVectorizer=0.6662468513853904
Untuned Accuracy of LinearSVC using TfidfVectorizer=0.6687657430730478
Tuned Accuracy of LinearSVC using TfidfVectorizer=0.6738035264483627
```

B. Unsupervised Learning

**** Unsupervised Sentiment Classification using NLTK Vader****

*** Step B1: Analyze Sentiments***

```

In [33]: 1 import pandas as pd
2 from nltk.sentiment.vader import SentimentIntensityAnalyzer
3
4 def sentiment_analysis_vader_validation(df, filepath):
5     sid = SentimentIntensityAnalyzer()
6     # print df.head()
7     d = []
8     sentiment_map = {'pos': 4, 'neg': 0, 'neu': 2}
9     for index, tweet in df.iterrows():
10
11         if len(str(tweet['text']).split()) > 4:
12             tweet_txt = tweet['text']
13             tweet_date = tweet['date']
14             tweet_manual_label = tweet['sentiment']
15
16             ss = sid.polarity_scores(tweet_txt)
17
18             '''MAX LOGIC'''
19             score_sentiment = max(ss['neg'], ss['neu'], ss['pos'])
20
21             '''
22             # COMPLEX LOGIC
23             if ss['neg']>0 and ss['pos']>0 and ss['neu']>0:
24                 score_sentiment = max(ss['neg'], ss['neu'], ss['pos'])
25             elif ss['neg']==0 and ss['pos']>0 and ss['neu']>0:
26                 score_sentiment = ss['pos']
27             elif ss['pos'] == 0 and ss['neg'] > 0 and ss['neu'] > 0:
28                 score_sentiment = ss['neg']
29             elif ss['pos'] == 0 and ss['neg'] == 0 and ss['neu'] > 0:
30                 score_sentiment = ss['neu']
31             '''
32             sentiment = [k for k, v in ss.items() if v == score_sentiment][0]
33             sentiment_mapping = sentiment_map[sentiment]
34             if tweet_manual_label == sentiment_mapping:
35                 validation_result='Match'
36             else:
37                 validation_result='Mismatch'
38
39             d.append({'date': tweet_date, 'text': tweet_txt, 'polarity_score_neg':ss['neg'], 'pol
40
41         df_processed = pd.DataFrame(d)
42         #df_processed.to_csv(filepath, index=False)
43         print(df_processed.groupby(['validation_result'])['validation_result'].count())
44
45
46 # Using merged_df created in Step A1
47 # merged_df has all the labelled tweets for MSFT and AAPL
48 output_file = 'vader_predictions.csv'
49 sentiment_analysis_vader_validation(aapl, output_file)

```

```

validation_result
Match      2269
Mismatch   1389
Name: validation_result, dtype: int64

```

OUTCOME OF THE SELECTION PROCESS

For our project, it is important to achieve the highest accuracy in classifying the tweet sentiments. With the results of the above analysis, we observed that LinearSVC alongwith TfidfVectorizer gave the best accuracy of 60.99%.

Therefore, we have selected the untuned LinearSVC alongwith TfidfVectorizer for the tweet sentiment classification.



