

```
In [ ]: 1 from google.colab import drive
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```

```
In [205]:
```

```
In [ ]:
```

```
In [8]: 1 import nltk
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...  
[nltk_data]   Unzipping corpora/stopwords.zip.
```

```
Out[8]: True
```

```
In [9]:
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...  
[nltk_data]   Unzipping tokenizers/punkt.zip.
```

```
Out[9]: True
```

```
In [120]: 1 from google.colab import files
```

**Browse...** No files selected.

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

```
Saving Canada_COVID19_data_M1.csv to Canada_COVID19_data_M1 (4).csv  
Saving Canada_COVID19_data_M2.csv to Canada_COVID19_data_M2 (4).csv  
Saving DecreaseInCOVIDAprToSepM1.csv to DecreaseInCOVIDAprToSepM1 (4).  
csv  
Saving DecreaseInCOVIDAprToSepM2.csv to DecreaseInCOVIDAprToSepM2 (1).  
csv  
Saving IncreaseInTweetAprToSepM1.csv to IncreaseInTweetAprToSepM1 (4).  
csv  
Saving IncreaseInTweetAprToSepM2.csv to IncreaseInTweetAprToSepM2 (1).  
csv  
Saving NoChangeTotalM1.csv to NoChangeTotalM1.csv  
Saving NoChangeTotalM2.csv to NoChangeTotalM2 (1).csv
```

```
In [179]: 1 import pandas as pd
          2 import io
          3
          4 modelnumber='M2'
          5 df = pd.read_csv(io.StringIO(uploaded['Canada_COVID19_data_']+modelnumi
```

```
Out[179]: 0 NaN
          1 NaN
          2 NaN
          3 NaN
          4 NaN
          ...
          409 RT CPHOCanada This illustration adapted from ...
          410 In the past week there were over 300000 new C...
          411 Among people exposed to COVID19 in their home...
          412 As we continue to address the impacts of COVI...
          413 Want trusted up to date information on COVID1...
          Name: Tweet, Length: 414, dtype: object
```

```
In [180]: 1
          2 from gensim.models import Word2Vec
          3
          4 from nltk.cluster import KMeansClusterer
          5 import nltk
          6 from nltk.corpus import stopwords
          7 import numpy as np
          8 import pandas as pd
          9 import io
         10
         11 from sklearn import cluster
         12 from sklearn import metrics
         13
         14 df=df.dropna();
         15
         16 df['Tweet'] = df['Tweet'].str.lower()
         17 corpus = df['Tweet'].to_numpy()
         18 tweets1 = corpus[0:207] # tweets from Nov 2019 to Apr 2020
         19 tweets2 = corpus[207:360] # tweets from Apr 2020 to Sep 2020
         20 tweets3 = corpus[360:396] # tweets from Sep 2020 to Oct 2020
         21 tweetcategories1 = df['category_value'].to_numpy()[0:207]
         22 tweetcategories2 = df['category_value'].to_numpy()[207:360]
         23 tweetcategories3 = df['category_value'].to_numpy()[360:396]
         24
         25 tweets1 = [nltk.word_tokenize(tweet) for tweet in tweets1]
         26 tweets2 = [nltk.word_tokenize(tweet) for tweet in tweets2]
         27 tweets3 = [nltk.word_tokenize(tweet) for tweet in tweets3]
         28
         29 for i in range(len(tweets1)):
         30     tweets1[i] = [word for word in tweets1[i] if word not in stopwords]
         31
         32 for i in range(len(tweets2)):
         33     tweets2[i] = [word for word in tweets2[i] if word not in stopwords]
         34
         35 for i in range(len(tweets3)):
         36     tweets3[i] = [word for word in tweets3[i] if word not in stopwords]
         37
```

```

38 corpus = tweets1+tweets2+tweets3
39
40 # Create word embedding model for the corpus of all tweets
41 model = Word2Vec(corpus, min_count=1)
42
43 def tweet_vectorizer(tweet, model):
44     tweet_vec = []
45     numw = 0
46     for w in tweet:
47         try:
48             if numw == 0:
49                 tweet_vec = model[w]
50             else:
51                 tweet_vec = np.add(tweet_vec, model[w])
52             numw+=1
53         except:
54             pass
55
56     return np.asarray(tweet_vec) / numw
57
58 X1=[] #Vectorized tweets from Nov 2019 to Apr 2020
59 X2=[] #Vectorized tweets from Apr 2020 to Sep 2020
60 X3=[] #Vectorized tweets from Sep 2020 to Oct 2020
61
62 # Vectorize the tweets by using the above Word2Vec model
63 for tweet in tweets1:
64     X1.append(tweet_vectorizer(tweet, model))
65
66 for tweet in tweets2:
67     X2.append(tweet_vectorizer(tweet, model))
68
69 for tweet in tweets3:
70     X3.append(tweet_vectorizer(tweet, model))
71
72 # Corresponding tweet labels
73 Y1 = [0 if y == 0 else 1 if y > 0 else 2 for y in tweetcategories1]
74 Y2 = [0 if y == 0 else 1 if y > 0 else 2 for y in tweetcategories2]
75 Y3 = [0 if y == 0 else 1 if y > 0 else 2 for y in tweetcategories3]
76
77

```

/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:49: DeprecationWarning: Call to deprecated `\_\_getitem\_\_` (Method will be removed in 4.0.0, use self.wv.\_\_getitem\_\_() instead).

/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:51: DeprecationWarning: Call to deprecated `\_\_getitem\_\_` (Method will be removed in 4.0.0, use self.wv.\_\_getitem\_\_() instead).

In [181]:

```

1 # Tweets from dataset2 for increase in COVID
2 df_new = pd.read_csv('IncreaseInTweetAprToSep'+modelnumber+'.csv')
3 df_new["Tweet"] = df_new["Tweet"].str.lower()
4 tweetAprtoSept_Increase = df_new["Tweet"].to_numpy()
5 tweetAprtoSept_Increase= [nltk.word_tokenize(tweet) for tweet in tweetAprtoSept_Increase]
6 X_tweetAprtoSept_Increase=[]
7
8 for i in range(len(tweetAprtoSept_Increase)):

```

```

9         tweetAprtoSept_Increase[i] = [word for word in tweetAprtoSept_Inc
10
11 for tweet in tweetAprtoSept_Increase:
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:49: Depre
cationWarning: Call to deprecated `__getitem__` (Method will be remove
d in 4.0.0, use self.wv.__getitem__() instead).
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:51: Depre
cationWarning: Call to deprecated `__getitem__` (Method will be remove
d in 4.0.0, use self.wv.__getitem__() instead).

```

In [182]:

```

1 # Tweets from dataset2 for decrease in COVID
2 df_new2 = pd.read_csv('DecreaseInCOVIDAprToSep'+modelnumber+'.csv')
3 df_new2["Tweet"] = df_new2["Tweet"].str.lower()
4 df = df_new2["Tweet"].dropna()
5 tweetAprtoSept_Decrease = df.to_numpy()
6
7 tweetAprtoSept_Decrease= [nltk.word_tokenize(tweet) for tweet in twee
8 X_tweetAprtoSept_Decrease=[]
9
10 for i in range(len(tweetAprtoSept_Decrease)):
11     tweetAprtoSept_Decrease[i] = [word for word in tweetAprtoSept_Dec
12
13 for tweet in tweetAprtoSept_Decrease:
14     X_tweetAprtoSept_Decrease.append(tweet_vectorizer(tweet, model))

```

/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:49: Depre  
cationWarning: Call to deprecated `\_\_getitem\_\_` (Method will be remove  
d in 4.0.0, use self.wv.\_\_getitem\_\_() instead).  
/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:51: Depre  
cationWarning: Call to deprecated `\_\_getitem\_\_` (Method will be remove  
d in 4.0.0, use self.wv.\_\_getitem\_\_() instead).

In [183]:

```

1 # Tweets from dataset2 for no change in COVID
2 df_new2 = pd.read_csv('NoChangeTotal'+modelnumber+'.csv')
3 df_new2["Tweet"] = df_new2["Tweet"].str.lower()
4 df = df_new2["Tweet"].dropna()
5 print(df)
6 nochange = df.to_numpy()
7
8 nochange
9 nochange = [nltk.word_tokenize(tweet) for tweet in nochange if tweet
10 len(nochange)
11 X_nochange=[]
12
13 for i in range(len(nochange)):
14     nochange[i] = [word for word in nochange[i] if word not in stopwo
15
16 for tweet in nochange:

```

```
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:49: DeprecationWarning: Call to deprecated `__getitem__` (Method will be removed in 4.0.0, use self.wv.__getitem__() instead).
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:51: DeprecationWarning: Call to deprecated `__getitem__` (Method will be removed in 4.0.0, use self.wv.__getitem__() instead).
```

```
1 # Clustering
2 from nltk.cluster import KMeansClusterer
3 import nltk
4 NUM_CLUSTERS=3
5 clustermodel = KMeansClusterer(NUM_CLUSTERS, distance=nltk.cluster.util.distance_matrix)
6 assigned_clusters = clustermodel.cluster(X_train, assign_clusters=True)
7 print (assigned_clusters)
```

```

1 cluster_numbers = []
2 count0 = 0
3 count1 = 0
4 count2 = 0
5
6 for i in X_dev1:
7     cluster_numbers.append(clustermodel.classify_vectorspace(i))
8
9 for i in cluster_numbers:
10     if(i== 0):
11         count0 += 1
12     if(i == 1):
13         count1 += 1
14     if(i == 2):
15         count2 += 1

```

```
1 cluster_numbers = []
2 count0 = 0
```

```

3 count1 = 0
4 count2 = 0
5
6 for i in X_dev2:
7     cluster_numbers.append(clustermodel.classify_vectorspace(i))
8
9 for i in cluster_numbers:
10     if(i== 0):
11         count0 += 1
12     if(i == 1):
13         count1 += 1
14     if(i == 2):
15         count2 += 1

```

No change: 0: 54, 1:23, 2:7

```

In [ ]: 1 cluster_numbers = []
2 count0 = 0
3 count1 = 0
4 count2 = 0
5
6 for i in X_dev3:
7     cluster_numbers.append(clustermodel.classify_vectorspace(i))
8
9 for i in cluster_numbers:
10     if(i== 0):
11         count0 += 1
12     if(i == 1):
13         count1 += 1
14     if(i == 2):
15         count2 += 1

```

No change: 0: 55, 1:33, 2:14

```

In [199]: 1 # Classification
2
3 # KNN
4 from sklearn.neighbors import KNeighborsClassifier
5 neigh = KNeighborsClassifier(n_neighbors=4, metric='cosine')
6 X_train = np.array(X1+X2)
7 Y_train = np.array(Y1+Y2)
8 neigh.fit(X_train, Y_train)
9

```

```

Out[199]: KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='cosine',
                               metric_params=None, n_jobs=None, n_neighbors=4, p
                               =2,
                               weights='uniform')

```

```

In [200]: 1 #Precision, recall, and accuracy on the development dataset, before h
2
3 from sklearn.metrics import recall_score
4 from sklearn.metrics import precision_score
5 from sklearn.metrics import accuracy_score
6 import numpy as np
7
8

```

```

9 X_dev1 = np.array(X_nochange)
10 X_dev2 = np.array(X_tweetAprtoSept_Increase)
11 X_dev3 = np.array(X_tweetAprtoSept_Decrease)
12
13 p1 = neigh.predict(X_dev1)
14 p2 = neigh.predict(X_dev2)
15 p3 = neigh.predict(X_dev3)
16
17 Y_dev_true = [0]*len(X_dev1) + [1]*len(X_dev2) + [2]*len(X_dev3)
18 Y_dev_predict = list(p1) + list(p2) + list(p3)
19 precision = precision_score(Y_dev_true, Y_dev_predict, average='micro')
20 recall_score = recall_score(Y_dev_true, Y_dev_predict, average='micro')
21 accuracy_score = accuracy_score(Y_dev_predict, Y_dev_true)
22 print(precision, recall_score, accuracy_score)
23
24 print("No change: 0: {}, 1:{}, 2:{}".format(np.count_nonzero(p1 == 0),
25 print("Increase: 0: {}, 1:{}, 2:{}".format(np.count_nonzero(p2 == 0),
26 print("Decrease: 0: {}, 1:{}, 2:{}".format(np.count_nonzero(p3 == 0),

```

0.754601226993865 0.754601226993865 0.754601226993865

No change: 0: 127, 1:8, 2:7

Increase: 0: 5, 1:62, 2:15

Decrease: 0: 5, 1:40, 2:57

In [192]:

```

1 # Hyperparameter tuning for KNN with grid search
2
3 from sklearn.model_selection import GridSearchCV
4
5 hyperparams = {'n_neighbors': [3,4,5,6,9],
6                 'weights': ['distance'],
7                 'metric': ['euclidian', 'cosine', 'manhattan', 'jaccard', ''],
8 gs = GridSearchCV(KNeighborsClassifier(), hyperparams, verbose=1, cv=5)
9 X_train = np.array(X1+X2)
10 Y_train = np.array(Y1+Y2)
11 goodmodel = gs.fit(X_train, Y_train)

```

Fitting 3 folds for each of 25 candidates, totalling 75 fits

[Parallel(n\_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.

[Parallel(n\_jobs=-1)]: Done 75 out of 75 | elapsed: 0.4s finished

Out[192]: 0.6388888888888889

In [193]:

```

1 #Precision, recall, and accuracy on the development dataset, with the
2
3 from sklearn.metrics import recall_score
4 from sklearn.metrics import precision_score
5 from sklearn.metrics import accuracy_score
6 import numpy as np
7
8 X_dev1 = np.array(X_nochange)
9 X_dev2 = np.array(X_tweetAprtoSept_Increase)
10 X_dev3 = np.array(X_tweetAprtoSept_Decrease)
11
12 p1 = goodmodel.predict(X_dev1)
13 p2 = goodmodel.predict(X_dev2)

```

```

14 p3 = goodmodel.predict(X_dev3)
15
16 Y_dev_true = [0]*len(X_dev1) + [1]*len(X_dev2) + [2]*len(X_dev3)
17 Y_dev_predict = list(p1) + list(p2) + list(p3)
18 precision = precision_score(Y_dev_true, Y_dev_predict, average='micro')
19 recall_score = recall_score(Y_dev_true, Y_dev_predict, average='micro')
20 accuracy_score = accuracy_score(Y_dev_predict, Y_dev_true)
21 print(precision, recall_score, accuracy_score)
22
23 print("No change: 0: {}, 1:{}, 2:{}".format(np.count_nonzero(p1 == 0),
24 print("Increase: 0: {}, 1:{}, 2:{}".format(np.count_nonzero(p2 == 0),
25 print("Decrease: 0: {}, 1:{}, 2:{}".format(np.count_nonzero(p3 == 0),

0.9846625766871165 0.9846625766871165 0.9846625766871165
No change: 0: 142, 1:0, 2:0
Increase: 0: 1, 1:80, 2:1
Decrease: 0: 2, 1:1, 2:99

```

In [144]:

Out[144]: {'metric': 'manhattan', 'n\_neighbors': 9, 'weights': 'distance'}

In [202]:

```

1 # Precision, Recall, and Accuracy for Test Dataset
2
3 from sklearn.metrics import recall_score
4 from sklearn.metrics import precision_score
5 from sklearn.metrics import accuracy_score
6 import numpy as np
7
8 X_test = np.array(X3)
9
10 Y_test_true = Y3
11 Y_test_predict = list(neigh.predict(X_test))
12 precision = precision_score(Y_test_true, Y_test_predict, average='micro')
13 recall_score = recall_score(Y_test_true, Y_test_predict, average='micro')
14 accuracy_score = accuracy_score(Y_test_predict, Y_test_true)
15 print(precision, recall_score, accuracy_score)

0.5 0.5 0.5

```

In [203]:

```

1 # Confusion Matrix
2
3 from sklearn.metrics import multilabel_confusion_matrix
4 from sklearn.metrics import plot_confusion_matrix
5 import matplotlib.pyplot as plt
6
7 multilabel_confusion_matrix(Y_test_true, Y_test_predict, labels=[0, 1])
8 np.set_printoptions(precision=2)
9
10 # Plot non-normalized confusion matrix
11 titles_options = [("Confusion matrix, without normalization", None),
12                  ("Normalized confusion matrix", 'true')]
13 for title, normalize in titles_options:
14     disp = plot_confusion_matrix(neigh, X_test, Y_test_predict,
15                                 display_labels=np.array(["NO_CHANGE",
16                                 cmap=plt.cm.Blues,

```



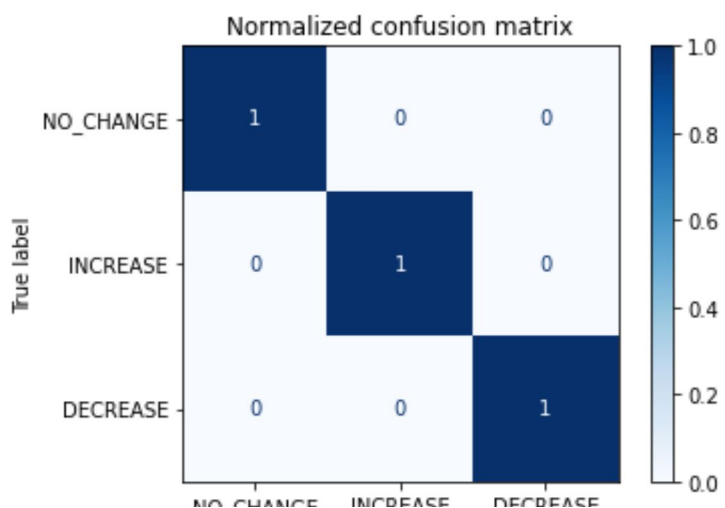
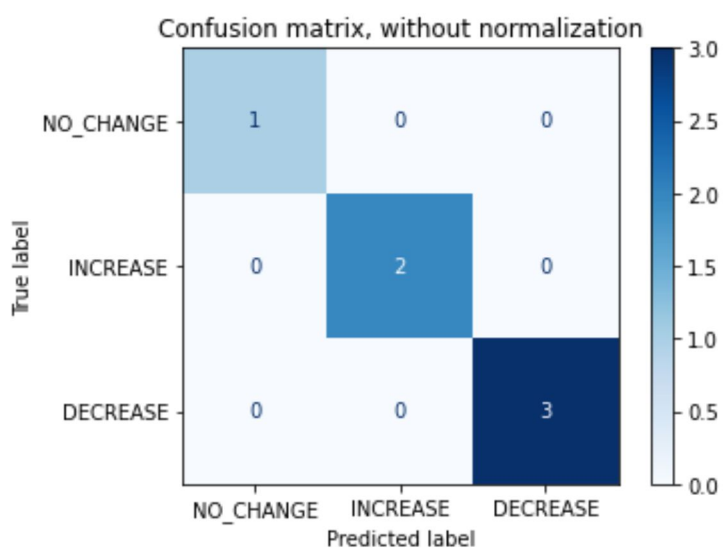
```
17                                     normalize=normalize)
18     disp.ax_.set_title(title)
19
20
21     print(title)
22     print(disp.confusion_matrix)
23
24
25 plt.show()
26
27
28
29
30
31
32
```

Confusion matrix, without normalization

```
[[1 0 0]
 [0 2 0]
 [0 0 3]]
```

Normalized confusion matrix

```
[[1. 0. 0.]
 [0. 1. 0.]
 [0. 0. 1.]]
```



In [105]:

Out[105]: [1, 1, 2, 2, 1, 1, 2, 1, 1, 2, 1, 2, 1, 2, 2, 2, 1]

In [ ]: