

```
In [2]: 1 import pandas as pd
        2 import string
        3 from collections import defaultdict
        4 import gzip
        5 from sklearn import linear_model
        6 import numpy as np
        7 import math
```

```
In [3]: 1 movies = pd.read_csv(r'C:\Users\Travi\Documents\GitHub\CSE158_HWs\AllMov
```

C:\Users\Travi\Anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3058: DtypeWarning: Columns (7) have mixed types.Specify dtype option on import or set low_memory=False.
interactivity=interactivity, compiler=compiler, result=result)

```
In [4]: 1 len(movies)
```

Out[4]: 329044

```
In [5]: 1 #dropping any empty overviews or genres
        2 movies = movies.dropna(subset=['overview', 'genres'])
        3 movies = movies[movies['overview'] != 'No overview found']
        4
        5 mult_genres_indices = movies[movies['genres'].str.find('|') != -1].index
        6 movies = movies.drop(mult_genres_indices)
        7
```

```
In [ ]: 1
```

```
In [6]: 1 movies['genres'].value_counts().index.values
```

Out[6]: array(['Drama', 'Documentary', 'Comedy', 'Music', 'Animation', 'Horror',
'Action', 'Western', 'Thriller', 'Family', 'Romance', 'Crime',
'Science Fiction', 'Adventure', 'Mystery', 'Fantasy', 'War',
'History', 'TV Movie', 'Foreign'], dtype=object)

```
In [7]: 1 #creating dictionary of the values in the genre column to have a corres
        2 genreID = dict(zip(list(movies['genres'].value_counts().index.values), range(1, len(movies['genres'].value_counts().index.values)+1)))
```

In [8]: 1 genreID

Out[8]: {'Drama': 0,
'Documentary': 1,
'Comedy': 2,
'Music': 3,
'Animation': 4,
'Horror': 5,
'Action': 6,
'Western': 7,
'Thriller': 8,
'Family': 9,
'Romance': 10,
'Crime': 11,
'Science Fiction': 12,
'Adventure': 13,
'Mystery': 14,
'Fantasy': 15,
'War': 16,
'History': 17,
'TV Movie': 18,
'Foreign': 19}

In [9]: 1 *#Replacing the genres so we can*
2 `movies['genres'] = movies['genres'].replace(genreID)`

In [10]: 1 *#Only taking columns we need*
2 `data = movies[['id', 'genres', 'overview']]`

```
In [11]: 1 data['genres'].value_counts()
```

```
Out[11]: 0      23869
         1      21417
         2      15734
         3       6334
         4       5463
         5       5329
         6       2560
         7       2369
         8       2174
         9       1515
        10       1081
        11        980
        12        923
        13        696
        14        467
        15        418
        16        280
        17        209
        18        138
        19        121
        Name: genres, dtype: int64
```

```
In [12]: 1 punct = string.punctuation
         2 bigram = defaultdict(int)
         3 unigram = defaultdict(int)
         4 totalbigrams = 0
         5 total_unigram = 0
         6
         7 def counts(d):
         8     t = str(d['overview'])
         9
        10
        11     t = t.lower() # Lowercase string
        12
        13     t = [c for c in t if not (c in punct)] # non-punct characters
        14     t = ''.join(t) # convert back to string
        15     words = t.strip().split() # tokenizes
        16
        17     for i in range(len(words)):
        18         unigram[str(words[i])] += 1 #unigram count
        19
        20         if (i+1) == len(words):
        21             break
        22         w1 = words[i]
        23         w2 = words[i + 1]
        24         bigram[w1 + ', ' + w2] += 1 #bigram count
        25
```

```
In [13]: 1 data.apply(counts, axis =1)
```

```
Out[13]: 4      None
          9      None
          16     None
          20     None
          26     None
          ...
          329033 None
          329034 None
          329035 None
          329037 None
          329043 None
          Length: 92077, dtype: object
```

```
In [14]: 1 unigram
```

```
Out[14]: defaultdict(int,
                        {'timo': 18,
                         'novotny': 7,
                         'labels': 31,
                         'his': 56209,
                         'new': 9549,
                         'project': 779,
                         'an': 27739,
                         'experimental': 422,
                         'music': 3481,
                         'documentary': 7318,
                         'film': 16983,
                         'in': 97200,
                         'a': 170329,
                         'remix': 37,
                         'of': 146108,
                         'the': 285280,
                         'celebrated': 262,
                         'megacities': 7,
                         '10071': 222})
```

In [15]: 1 bigram

```
Out[15]: defaultdict(int,
                        {'timo, novotny': 1,
                         'novotny, labels': 1,
                         'labels, his': 1,
                         'his, new': 482,
                         'new, project': 13,
                         'project, an': 7,
                         'an, experimental': 125,
                         'experimental, music': 4,
                         'music, documentary': 7,
                         'documentary, film': 495,
                         'film, in': 279,
                         'in, a': 10892,
                         'a, remix': 7,
                         'remix, of': 9,
                         'of, the': 33303,
                         'the, celebrated': 33,
                         'celebrated, film': 2,
                         'film, megacities': 1,
                         'megacities, ...': 1})
```

In []: 1

```
In [16]: 1 #Taking only the top 1000 unigrams and bigrams
2 bigrams_1000 = sorted(bigram.items(), key=lambda x: x[1], reverse=True)[
3 unigram_1000 = sorted(unigram.items(), key=lambda x: x[1], reverse=True)
```

```
In [17]: 1 bigram_words = []
2 unigram_words = []
3 for i in bigrams_1000:
4     bigram_words.append(i[0])
5 for i in unigram_1000:
6     unigram_words.append(i[0])
```

In [18]: 1 len(bigram_words)

Out[18]: 1000

In []: 1

In []: 1

In []: 1

Bigram Count

```
In [19]: 1 wordId_bigram = dict(zip(bigram_words, range(len(bigrams_1000))))
2
3 def feature(datum):
4     datum = datum[1]
5     feat = [0]*len(bigrams_1000)
6     r = ''.join([c for c in str(datum['overview']).lower() if not c in ' , . : ; '])
7     text = r.split()
8     for w in range(len(text)):
9         if (w+1) == len(text):
10             continue
11         w1 = text[w]
12         w2 = text[w + 1]
13         word = str(w1) + ', ' + str(w2)
14         if word in bigram_words:
15             feat[wordId_bigram[word]] += 1
16     feat.append(1) #offset
17     return feat
18
```

```
In [20]: 1 train_movies = movies.sample(n = 50000)
2         vaildation_movies = movies.sample(n = 50000)
3         test_movies = movies.sample(n = 50000)
```

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [ ]: 1
```

```
In [21]: 1 X = [feature(d) for d in train_movies.iterrows()]
2         y = [d[1]['genres'] for d in train_movies.iterrows()]
```

```
In [ ]: 1
```

```
In [22]: 1 clf = linear_model.LogisticRegression(C=1, class_weight='balanced') # MS
2 clf.fit(X, y)
3 theta = clf.coef_
4 predictions = clf.predict(X)
```

C:\Users\Travi\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.
 FutureWarning)
 C:\Users\Travi\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:469: FutureWarning: Default multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
 "this warning.", FutureWarning)

```
In [ ]: 1
```

```
In [23]: 1 predictions
```

Out[23]: array([8, 1, 6, ..., 6, 4, 4])

```
In [24]: 1 correctPrediction = np.array(y) == np.array(predictions)
2 sum(correctPrediction) / len(correctPrediction)
```

Out[24]: 0.38218

```
In [ ]: 1
```

Unigram Count

```
In [25]: 1 wordId_unigram = dict(zip(unigram_words, range(len(unigram_1000))))
2
3 def feature(datum):
4     datum = datum[1]
5     feat = [0]*len(bigrams_1000)
6     r = ''.join([c for c in str(datum['overview']).lower() if not c in p
7     text = r.split()
8     for w in range(len(text)):
9         word = text[w]
10        if word in unigram_words:
11            feat[wordId_unigram[word]] += 1
12    feat.append(1) #offset
13    return feat
14
```

```
In [26]: 1 X = [feature(d) for d in train_movies.iterrows()]
2 y = [d[1]['genres'] for d in train_movies.iterrows()]
```

```
In [27]: 1 clf = linear_model.LogisticRegression(C=1, class_weight='balanced') # MS
2 clf.fit(X, y)
3 theta = clf.coef_
4 predictions = clf.predict(X)
```

C:\Users\Travi\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:432: FutureWarning: Default solver will be changed to 'lbfgs' in 0.22. Specify a solver to silence this warning.

FutureWarning)

C:\Users\Travi\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:469: FutureWarning: Default multi_class will be changed to 'auto' in 0.22. Specify the multi_class option to silence this warning.
"this warning.", FutureWarning)

```
In [28]: 1 correctPrediction = np.array(y) == np.array(predictions)
2 sum(correctPrediction) / len(correctPrediction)
```

Out[28]: 0.55972

Unigram Count Accuracy = .55

Bigram Count Accuracy = .38

```
In [ ]: 1
```

```
In [29]: 1 vaildation = pd.DataFrame(columns = ["TF-IDF", 'Unigram', 'Bigram', 'Acc
```

```
In [30]: 1 dicta = dict(zip(["TF-IDF", 'Unigram', 'Bigram', 'Accuracy', 'n_estimatc
```

```
In [31]: 1 dicta
```

Out[31]: {'TF-IDF': [], 'Unigram': [], 'Bigram': [], 'Accuracy': [], 'n_estimators': []}

Unigram TF-IDF


```
In [33]: 1 import nltk
2 from nltk.corpus import stopwords
3 nltk.download("stopwords")
4 from nltk.tokenize import word_tokenize
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\Travi\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

```
In [34]: 1 stop_words = set(stopwords.words('english'))
```

```
In [35]: 1 wordId = dict(zip(unigram_words, range(len(unigram_words))))
2
3
4
5 def feature(datum):
6     datum = datum[1]
7     feat = [0]*len(unigram_words)
8     r = ''.join([c for c in datum['overview'].lower() if not c in punct])
9     text = r.split()
10    unigram_wordId = dict(zip(unigram_words, [0]*1000))
11    text = [w for w in text if not (w in stop_words)]
12    for w in text:
13        if w in unigram_words:
14            unigram_wordId[w] += 1
15
16    for i in unigram_wordId:
17        if len(text) == 0 or unigram[i] == 0:
18            feat[wordId[i]] = 0
19            continue
20
21        tf = unigram_wordId[i]
22        idf = math.log(len(train_movies)/unigram[i], 10)
23        feat[wordId[i]] = tf/idf
24    feat.append(1) #offset
25    return feat
26
27
```

```
In [36]: 1 X = [feature(d) for d in train_movies.iterrows()]
2 y = [d[1]['genres'] for d in train_movies.iterrows()]
```

```
In [ ]: 1
```

```
1 Testing new classifiers
```

```
In [37]: 1 # clf = linear_model.LogisticRegression(C=1, class_weight='balanced') #
2 # clf.fit(X, y)
3 # theta = clf.coef_
4 # predictions = clf.predict(X)
```

```
In [38]: 1 # correctPrediction = np.array(y) == np.array(predictions)
2 # sum(correctPrediction) / len(correctPrediction)
```

```
In [39]: 1 from sklearn.ensemble import RandomForestClassifier
```

```
In [40]: 1 # clf = RandomForestClassifier()
2 # clf.fit(X, y)
3 # predictions = clf.predict(X)
```

```
In [41]: 1 # correctPrediction = np.array(y) == np.array(predictions)
2 # sum(correctPrediction) / len(correctPrediction)
```

```
In [42]: 1 X_vaild= [feature(d) for d in vaildation_movies.iterrows()]
2 y_vaild = [d[1]['genres'] for d in vaildation_movies.iterrows()]
```

```
1 RandomForest best perfoming classifier
```

```
In [43]: 1 for i in range(10, 100, 20):
2     clf = RandomForestClassifier(n_estimators = i)
3     clf.fit(X, y)
4     predictions = clf.predict(X_vaild)
5     correctPrediction = np.array(y_vaild) == np.array(predictions)
6     dicta['TF-IDF'].append('Yes')
7     dicta['Unigram'].append('Yes')
8     dicta['Bigram'].append('No')
9     dicta['Accuracy'].append(sum(correctPrediction) / len(correctPrediction))
10    dicta['n_estimators'].append(i)
11
12
13
```

```
In [44]: 1 dicta  
2
```

```
Out[44]: {'TF-IDF': ['Yes', 'Yes', 'Yes', 'Yes', 'Yes'],  
          'Unigram': ['Yes', 'Yes', 'Yes', 'Yes', 'Yes'],  
          'Bigram': ['No', 'No', 'No', 'No', 'No'],  
          'Accuracy': [0.75206, 0.77082, 0.77172, 0.77382, 0.77462],  
          'n_estimators': [10, 30, 50, 70, 90]}
```

```
In [ ]: 1
```

```
In [ ]: 1
```

test set for bigram tf-idf

```
In [45]: 1 X = [feature(d) for d in test_movies.iterrows()]  
2 y = [d[1]['genres'] for d in test_movies.iterrows()]
```

```
In [46]: 1 predictions = clf.predict(X)
```

```
In [47]: 1 correctPrediction = np.array(y) == np.array(predictions)  
2 sum(correctPrediction) / len(correctPrediction)
```

```
Out[47]: 0.77366
```

```
In [ ]: 1
```

```
In [ ]: 1
```

bigram tf-idf

```

In [48]: 1 wordId = dict(zip(bigram_words, range(len(bigram_words))))
2
3
4
5 def feature(datum):
6     datum = datum[1]
7     feat = [0]*len(bigram_words)
8     r = ''.join([c for c in datum['overview'].lower() if not c in punct])
9     text = r.split()
10    bigram_wordId = dict(zip(bigram_words, [0]*1000))
11    text = [w for w in text if not (w in stop_words)]
12
13    for w in range(len(text)):
14        if (w+1) == len(text):
15            continue
16        w1 = text[w]
17        w2 = text[w + 1]
18        word = str(w1) + ', ' + str(w2)
19        if word in bigram_words:
20            bigram_wordId[word] += 1
21
22    for i in bigram_words:
23        if len(text) == 0 or bigram[i] == 0:
24            feat[wordId[i]] = 0
25            continue
26
27        tf = bigram_wordId[i]
28        idf = math.log(len(train_movies)/bigram[i], 10)
29        feat[wordId[i]] = tf/idf
30    feat.append(1) #offset
31    return feat
32

```

```

In [ ]: 1

```

```

In [49]: 1 X = [feature(d) for d in train_movies.iterrows()]
2         y = [d[1]['genres'] for d in train_movies.iterrows()]

```

```

1 Testing new classifiers

```

```

In [50]: 1 # clf = linear_model.LogisticRegression(C=1, class_weight='balanced') #
2         # clf.fit(X, y)
3         # theta = clf.coef_
4         # predictions = clf.predict(X)

```

```

In [51]: 1 # correctPrediction = np.array(y) == np.array(predictions)
2         # sum(correctPrediction) / len(correctPrediction)

```

```
In [52]: 1 X_vaild= [feature(d) for d in vaildation_movies.iterrows()]
2 y_vaild = [d[1]['genres'] for d in vaildation_movies.iterrows()]
```

```
In [53]: 1 # from sklearn.ensemble import RandomForestClassifier
```

```
In [54]: 1 # clf = RandomForestClassifier()
2 # clf.fit(X, y)
3 # predictions = clf.predict(X)
```

```
In [55]: 1 # correctPrediction = np.array(y) == np.array(predictions)
2 # sum(correctPrediction) / len(correctPrediction)
```

```
In [ ]: 1
```

```
1 Again RandomFroest best perfoming classifer
```

```
In [56]: 1 for i in range(10, 100, 20):
2     clf = RandomForestClassifier(n_estimators = i)
3     clf.fit(X, y)
4     predictions = clf.predict(X_vaild)
5     correctPrediction = np.array(y_vaild) == np.array(predictions)
6     dicta['TF-IDF'].append('Yes')
7     dicta['Unigram'].append('No')
8     dicta['Bigram'].append('Yes')
9     dicta['Accuracy'].append(sum(correctPrediction) / len(correctPredict
10    dicta['n_estimators'].append(i)
```

In [57]: `dicta`

```
Out[57]: {'TF-IDF': ['Yes',
                    'Yes',
                    'Yes',
                    'Yes',
                    'Yes',
                    'Yes',
                    'Yes',
                    'Yes',
                    'Yes',
                    'Yes'],
          'Unigram': ['Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No', 'No', 'No', 'No', 'No'],
          'Bigram': ['No', 'No', 'No', 'No', 'No', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes'],
          'Accuracy': [0.75206,
                       0.77082,
                       0.77172,
                       0.77382,
                       0.77462,
                       0.28576,
                       0.28622,
                       0.28622,
                       0.28614,
                       0.28626],
          'n_estimators': [10, 30, 50, 70, 90, 10, 30, 50, 70, 90]}
```

In [58]: `pd.DataFrame(dicta)`

Out[58]:

	TF-IDF	Unigram	Bigram	Accuracy	n_estimators
0	Yes	Yes	No	0.75206	10
1	Yes	Yes	No	0.77082	30
2	Yes	Yes	No	0.77172	50
3	Yes	Yes	No	0.77382	70
4	Yes	Yes	No	0.77462	90
5	Yes	No	Yes	0.28576	10
6	Yes	No	Yes	0.28622	30
7	Yes	No	Yes	0.28622	50
8	Yes	No	Yes	0.28614	70
9	Yes	No	Yes	0.28626	90

test set for bigram tf-idf

In []: 1

In [61]: 1 X = [feature(d) for d in test_movies.iterrows()]
2 y = [d[1]['genres'] for d in test_movies.iterrows()]

In [62]: 1 predictions = clf.predict(X)

In [63]: 1 correctPrediction = np.array(y) == np.array(predictions)
2 sum(correctPrediction) / len(correctPrediction)

Out[63]: 0.28752

In []: 1

Unigram TF-IDF is by far the better predictor for this dataset

In []: 1