

Case Study - Sentiment Analysis

June 3, 2020

You are currently looking at **version 1.0** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks in the Coursera platform, visit the [Jupyter Notebook FAQ](#) course resource.

Note: Some of the cells in this notebook are computationally expensive. To reduce runtime, this notebook is using a subset of the data.

1 Case Study: Sentiment Analysis

1.0.1 Data Prep

```
In [1]: import pandas as pd
import numpy as np
```

```
# Read in the data
df = pd.read_csv('Amazon_Unlocked_Mobile.csv')

# Sample the data to speed up computation
# Comment out this line to match with lecture
df = df.sample(frac=0.1, random_state=10)

df.head()
```

```
Out[1]:
```

		Product Name	Brand Name	Price	\
394349	Sony XPERIA Z2 D6503 FACTORY UNLOCKED Internat...		NaN	244.95	
34377	Apple iPhone 5c 8GB (Pink) - Verizon Wireless		Apple	194.99	
248521	Motorola Droid RAZR MAXX XT912 M Verizon Smart...		Motorola	174.99	
167661	CNPGD [U.S. Office Extended Warranty] Smartwat...		CNPGD	49.99	
73287	Apple iPhone 7 Unlocked Phone 256 GB - US Vers...		Apple	922.00	

	Rating	Reviews	\
394349	5	Very good one! Better than Samsung S and iphon...	
34377	1	The phone needed a SIM card, would have been n...	
248521	5	I was 3 months away from my upgrade and my Str...	
167661	1	an experience i want to forget	

```
73287          5          GREAT PHONE WORK ACCORDING MY EXPECTATIONS.
```

```

Review Votes
394349      0.0
34377       1.0
248521      3.0
167661      0.0
73287       1.0

```

```

In [2]: # Drop missing values
df.dropna(inplace=True)

# Remove any 'neutral' ratings equal to 3
df = df[df['Rating'] != 3]

# Encode 4s and 5s as 1 (rated positively)
# Encode 1s and 2s as 0 (rated poorly)
df['Positively Rated'] = np.where(df['Rating'] > 3, 1, 0)
df.head(10)

```

```

Out[2]:
Product Name  Brand Name  Price \
34377      Apple iPhone 5c 8GB (Pink) - Verizon Wireless      Apple  194.99
248521  Motorola Droid RAZR MAXX XT912 M Verizon Smart...  Motorola  174.99
167661  CNPGD [U.S. Office Extended Warranty] Smartwat...  CNPGD    49.99
73287      Apple iPhone 7 Unlocked Phone 256 GB - US Vers...  Apple   922.00
277158  Nokia N8 Unlocked GSM Touch Screen Phone Featu...  Nokia    95.00
100311  Blackberry Torch 2 9810 Unlocked Phone with 1...  BlackBerry  77.49
251669  Motorola Moto E (1st Generation) - Black - 4 G...  Motorola    89.99
279878  OtterBox 77-29864 Defender Series Hybrid Case ...  OtterBox     9.99
406017  Verizon HTC Rezound 4G Android Smarphone - 8MP...  HTC     74.99
302567  RCA M1 Unlocked Cell Phone, Dual Sim, 5Mp Came...  RCA   159.99

```

```

Rating  Reviews \
34377    1  The phone needed a SIM card, would have been n...
248521    5  I was 3months away from my upgrade and my Str...
167661    1              an experience i want to forget
73287     5          GREAT PHONE WORK ACCORDING MY EXPECTATIONS.
277158    5  I fell in love with this phone because it did ...
100311    5  I am pleased with this Blackberry phone! The p...
251669    5  Great product, best value for money smartphone...
279878    5          I've bought 3 no problems. Fast delivery.
406017    4              Great phone for the price...
302567    5  My mom is not good with new technoloy but this...

```

```

Review Votes  Positively Rated
34377         1.0                0
248521        3.0                1
167661        0.0                0

```

73287	1.0	1
277158	0.0	1
100311	0.0	1
251669	0.0	1
279878	0.0	1
406017	0.0	1
302567	4.0	1

```
In [3]: # Most ratings are positive
df['Positively Rated'].mean()
```

```
Out[3]: 0.74717766860786672
```

```
In [5]: from sklearn.model_selection import train_test_split
```

```
    # Split data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(df['Reviews'],
                                                    df['Positively Rated'],
                                                    random_state=0)
```

```
In [6]: print('X_train first entry:\n\n', X_train.iloc[0])
        print('\n\nX_train shape: ', X_train.shape)
```

X_train first entry:

Everything about it is awesome!

X_train shape: (23052,)

2 CountVectorizer

```
In [7]: from sklearn.feature_extraction.text import CountVectorizer
```

```
    # Fit the CountVectorizer to the training data
vect = CountVectorizer().fit(X_train)
```

```
In [8]: vect.get_feature_names()[::2000]
```

```
Out[8]: ['00',
        'arroja',
        'comapañas',
        'dvds',
        'golden',
        'lands',
        'oil',
        'razonable',
        'smallsliver',
        'tweak']
```

```

In [9]: len(vect.get_feature_names())
Out[9]: 19601

In [10]: # transform the documents in the training data to a document-term matrix
X_train_vectorized = vect.transform(X_train)

X_train_vectorized

Out[10]: <23052x19601 sparse matrix of type '<class 'numpy.int64'>'
        with 613289 stored elements in Compressed Sparse Row format>

In [11]: from sklearn.linear_model import LogisticRegression

# Train the model
model = LogisticRegression()
model.fit(X_train_vectorized, y_train)

Out[11]: LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
        intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
        penalty='l2', random_state=None, solver='liblinear', tol=0.0001,
        verbose=0, warm_start=False)

In [12]: from sklearn.metrics import roc_auc_score

# Predict the transformed test documents
predictions = model.predict(vect.transform(X_test))

print('AUC: ', roc_auc_score(y_test, predictions))

AUC: 0.897433277667

In [13]: # get the feature names as numpy array
feature_names = np.array(vect.get_feature_names())

# Sort the coefficients from the model
sorted_coef_index = model.coef_[0].argsort()

# Find the 10 smallest and 10 largest coefficients
# The 10 largest coefficients are being indexed using [-11:-1]
# so the list returned is in order of largest to smallest
print('Smallest Coefs:\n{}\n'.format(feature_names[sorted_coef_index[:10]]))
print('Largest Coefs: \n{}\n'.format(feature_names[sorted_coef_index[-11:-1]]))

Smallest Coefs:
['worst' 'terrible' 'slow' 'junk' 'poor' 'sucks' 'horrible' 'useless'
 'waste' 'disappointed']

Largest Coefs:
['excellent' 'excelente' 'excellent' 'perfectly' 'love' 'perfect' 'exactly'
 'great' 'best' 'awesome']

```

3 Tfidf

```
In [14]: from sklearn.feature_extraction.text import TfidfVectorizer
```

```
# Fit the TfidfVectorizer to the training data specifying a minimum document frequency
vect = TfidfVectorizer(min_df=5).fit(X_train)
len(vect.get_feature_names())
```

```
Out[14]: 5442
```

```
In [15]: X_train_vectorized = vect.transform(X_train)
```

```
model = LogisticRegression()
model.fit(X_train_vectorized, y_train)

predictions = model.predict(vect.transform(X_test))

print('AUC: ', roc_auc_score(y_test, predictions))
```

```
AUC: 0.889951006492
```

```
In [16]: feature_names = np.array(vect.get_feature_names())
```

```
sorted_tfidf_index = X_train_vectorized.max(0).toarray()[0].argsort()

print('Smallest tfidf:\n{}\n'.format(feature_names[sorted_tfidf_index[:10]]))
print('Largest tfidf: \n{}\n'.format(feature_names[sorted_tfidf_index[-11:-1]]))
```

Smallest tfidf:

```
['61' 'printer' 'approach' 'adjustment' 'consequences' 'length' 'emailing'
 'degrees' 'handsfree' 'chipset']
```

Largest tfidf:

```
['unlocked' 'handy' 'useless' 'cheat' 'up' 'original' 'exelent' 'exelente'
 'exellent' 'satisfied']
```

```
In [17]: sorted_coef_index = model.coef_[0].argsort()
```

```
print('Smallest Coefs:\n{}\n'.format(feature_names[sorted_coef_index[:10]]))
print('Largest Coefs: \n{}\n'.format(feature_names[sorted_coef_index[-11:-1]]))
```

Smallest Coefs:

```
['not' 'slow' 'disappointed' 'worst' 'terrible' 'never' 'return' 'doesn'
 'horrible' 'waste']
```

Largest Coefs:

```
['great' 'love' 'excellent' 'good' 'best' 'perfect' 'price' 'awesome' 'far'
 'perfectly']
```

```
In [18]: # These reviews are treated the same by our current model
print(model.predict(vect.transform(['not an issue, phone is working',
                                     'an issue, phone is not working'])))
```

[0 0]

4 n-grams

```
In [19]: # Fit the CountVectorizer to the training data specifying a minimum
# document frequency of 5 and extracting 1-grams and 2-grams
vect = CountVectorizer(min_df=5, ngram_range=(1,2)).fit(X_train)
```

```
X_train_vectorized = vect.transform(X_train)
```

```
len(vect.get_feature_names())
```

Out[19]: 29072

```
In [20]: model = LogisticRegression()
model.fit(X_train_vectorized, y_train)

predictions = model.predict(vect.transform(X_test))

print('AUC: ', roc_auc_score(y_test, predictions))
```

AUC: 0.91106617946

```
In [21]: feature_names = np.array(vect.get_feature_names())

sorted_coef_index = model.coef_[0].argsort()

print('Smallest Coefs:\n{}\n'.format(feature_names[sorted_coef_index[:10]]))
print('Largest Coefs: \n{}\n'.format(feature_names[sorted_coef_index[-11:-1]]))
```

Smallest Coefs:

```
['no good' 'junk' 'poor' 'slow' 'worst' 'broken' 'not good' 'terrible'
 'defective' 'horrible']
```

Largest Coefs:

```
['excellent' 'excelente' 'excelent' 'perfect' 'great' 'love' 'awesome'
 'no problems' 'good' 'best']
```

```
In [22]: # These reviews are now correctly identified
print(model.predict(vect.transform(['not an issue, phone is working',
                                     'an issue, phone is not working'])))
```

```
[1 0]
```

```
In [ ]:
```