

# SENTIMENT ANALYSIS USING LOGISTIC REGRESSION

In [1]:	<pre>import pandas as pd</pre>																																																																																																																																		
In [2]:	<pre>df = pd.read_csv('Reviews.csv') df.head(10)</pre>																																																																																																																																		
Out[2]:	<table><tr><th></th><th></th><th>Id</th><th>ProductId</th><th>UserId</th><th>ProfileName</th><th>HelpfulnessNumerator</th><th>HelpfulnessDenominator</th><th>Score</th><th>Time</th><th>Summary</th></tr><tr><td></td><td>0</td><td>1</td><td>B001E4KFG0</td><td>A3SGXH7AUHU8GW</td><td>delmartian</td><td>1</td><td>1</td><td>5</td><td>1303862400</td><td>Good Quality Dog Food</td></tr><tr><td></td><td>1</td><td>2</td><td>B00813GRG4</td><td>A1D87F6ZCVE5NK</td><td>dll pa</td><td>0</td><td>0</td><td>1</td><td>1346976000</td><td>Not as Advertised</td></tr><tr><td></td><td>2</td><td>3</td><td>B000LQOCH0</td><td>ABXLMWJIXAIN</td><td>Natalia Corres "Natalia Corres"</td><td>1</td><td>1</td><td>4</td><td>1219017600</td><td>"Delight" says it all</td></tr><tr><td></td><td>3</td><td>4</td><td>B000UA0QIQ</td><td>A395BORC6FGVXV</td><td>Karl</td><td>3</td><td>3</td><td>2</td><td>1307923200</td><td>Cough Medicine</td></tr><tr><td></td><td>4</td><td>5</td><td>B006K2ZZ7K</td><td>A1UQRSCLF8GW1T</td><td>Michael D. Bigham "M. Wassir"</td><td>0</td><td>0</td><td>5</td><td>1350777600</td><td>Great taffy</td></tr><tr><td></td><td>5</td><td>6</td><td>B006K2ZZ7K</td><td>ADT0SRK1MGOEU</td><td>Twoapennything</td><td>0</td><td>0</td><td>4</td><td>1342051200</td><td>Nice Taffy</td></tr><tr><td></td><td>6</td><td>7</td><td>B006K2ZZ7K</td><td>A1SP2KV6FXXRU1</td><td>David C. Sullivan</td><td>0</td><td>0</td><td>5</td><td>1340150400</td><td>Great! Just as good as the expensive brands!</td></tr><tr><td></td><td>7</td><td>8</td><td>B006K2ZZ7K</td><td>A3JRGQVEQN31IQ</td><td>Pamela G. Williams</td><td>0</td><td>0</td><td>5</td><td>1336003200</td><td>Wonderful, tasty taffy</td></tr><tr><td></td><td>8</td><td>9</td><td>B000E7L2R4</td><td>A1MZY09TZK08BI</td><td>R. James</td><td>1</td><td>1</td><td>5</td><td>1322006400</td><td>Yay Barley</td></tr><tr><td></td><td>9</td><td>10</td><td>B00171APVA</td><td>A21BT40VZCCYT4</td><td>Carol A. Reed</td><td>0</td><td>0</td><td>5</td><td>1351209600</td><td>Healthy Dog Food</td></tr></table>												Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary		0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	5	1303862400	Good Quality Dog Food		1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	1	1346976000	Not as Advertised		2	3	B000LQOCH0	ABXLMWJIXAIN	Natalia Corres "Natalia Corres"	1	1	4	1219017600	"Delight" says it all		3	4	B000UA0QIQ	A395BORC6FGVXV	Karl	3	3	2	1307923200	Cough Medicine		4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0	0	5	1350777600	Great taffy		5	6	B006K2ZZ7K	ADT0SRK1MGOEU	Twoapennything	0	0	4	1342051200	Nice Taffy		6	7	B006K2ZZ7K	A1SP2KV6FXXRU1	David C. Sullivan	0	0	5	1340150400	Great! Just as good as the expensive brands!		7	8	B006K2ZZ7K	A3JRGQVEQN31IQ	Pamela G. Williams	0	0	5	1336003200	Wonderful, tasty taffy		8	9	B000E7L2R4	A1MZY09TZK08BI	R. James	1	1	5	1322006400	Yay Barley		9	10	B00171APVA	A21BT40VZCCYT4	Carol A. Reed	0	0	5	1351209600	Healthy Dog Food
		Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary																																																																																																																									
	0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	5	1303862400	Good Quality Dog Food																																																																																																																									
	1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	1	1346976000	Not as Advertised																																																																																																																									
	2	3	B000LQOCH0	ABXLMWJIXAIN	Natalia Corres "Natalia Corres"	1	1	4	1219017600	"Delight" says it all																																																																																																																									
	3	4	B000UA0QIQ	A395BORC6FGVXV	Karl	3	3	2	1307923200	Cough Medicine																																																																																																																									
	4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0	0	5	1350777600	Great taffy																																																																																																																									
	5	6	B006K2ZZ7K	ADT0SRK1MGOEU	Twoapennything	0	0	4	1342051200	Nice Taffy																																																																																																																									
	6	7	B006K2ZZ7K	A1SP2KV6FXXRU1	David C. Sullivan	0	0	5	1340150400	Great! Just as good as the expensive brands!																																																																																																																									
	7	8	B006K2ZZ7K	A3JRGQVEQN31IQ	Pamela G. Williams	0	0	5	1336003200	Wonderful, tasty taffy																																																																																																																									
	8	9	B000E7L2R4	A1MZY09TZK08BI	R. James	1	1	5	1322006400	Yay Barley																																																																																																																									
	9	10	B00171APVA	A21BT40VZCCYT4	Carol A. Reed	0	0	5	1351209600	Healthy Dog Food																																																																																																																									

## REMOVING ALL THE NULL VALUES

```
Out[3]: Id 0
ProductId 0
UserId 0
ProfileName 16
HelpfulnessNumerator 0
HelpfulnessDenominator 0
Score 0
Time 0
Summary 27
Text 0
dtype: int64

In [4]: df = df.dropna()

In [5]: df.isnull().sum()

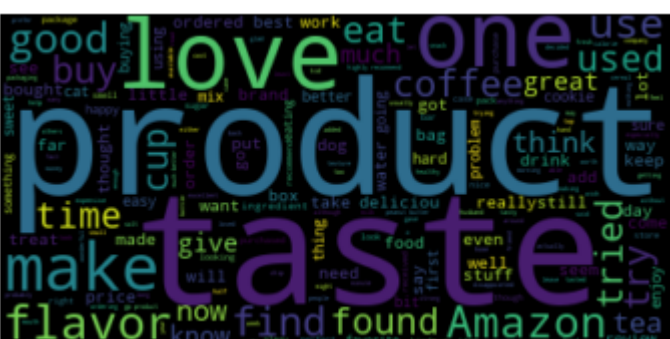
Out[5]: Id 0
ProductId 0
UserId 0
ProfileName 0
HelpfulnessNumerator 0
HelpfulnessDenominator 0
Score 0
Time 0
Summary 0
Text 0
dtype: int64
```

## CREATING WORDCLOUD TO FIND MOST FREQUENTLY USED WORDS

```
In [6]: import matplotlib.pyplot as plt
import seaborn as sns
```

## CREATING WORDCLOUD TO FIND MOST FREQUENTLY USED WORDS

```
stopwords = set(STOPWORDS)
stopwords.update(["br", "href"])
textt = " ".join(review for review in df.Text)
wordcloud = WordCloud(stopwords=stopwords).generate(textt)
plt.imshow(wordcloud, interpolation='bilinear')
plt.savefig('wordcloud01.png')
plt.axis("off")
plt.show()
```



## CLASSIFYING THE REVIEWS INTO "POS" AND "NEG" BASED ON "SCORE" & CREATING "SENTIMENT" COLUMN

```
In [8]: df = df[df['Score'] != 3]
```

## CLASSIFYING THE REVIEWS INTO "POS" AND "NEG" BASED ON "SCORE" & CREATING "SENTIMENT" COLUMN

In [8]:	<pre>df = df[df['Score'] != 3] df['sentiment'] = df['Score'].apply(lambda rating : +1 if rating &gt; 3 else -1) df.head(10)</pre>																																																																																																																																		
Out[8]:	<table><tr><th></th><th></th><th>Id</th><th>ProductId</th><th>UserId</th><th>ProfileName</th><th>HelpfulnessNumerator</th><th>HelpfulnessDenominator</th><th>Score</th><th>Time</th><th>Summary</th></tr><tr><td></td><td>0</td><td>1</td><td>B001E4KFG0</td><td>A3SGXH7AUHU8GW</td><td>delmartian</td><td>1</td><td>1</td><td>5</td><td>1303862400</td><td>Good Quality Dog Food</td></tr><tr><td></td><td>1</td><td>2</td><td>B00813GRG4</td><td>A1D87F6ZCVE5NK</td><td>dll pa</td><td>0</td><td>0</td><td>1</td><td>1346976000</td><td>Not as Advertised</td></tr><tr><td></td><td>2</td><td>3</td><td>B000LQOCH0</td><td>ABXLMWJIXAIN</td><td>Natalia Corres "Natalia Corres"</td><td>1</td><td>1</td><td>4</td><td>1219017600</td><td>"Delight" says it all</td></tr><tr><td></td><td>3</td><td>4</td><td>B000UA0QIQ</td><td>A395BORC6FGVXV</td><td>Karl</td><td>3</td><td>3</td><td>2</td><td>1307923200</td><td>Cough Medicine</td></tr><tr><td></td><td>4</td><td>5</td><td>B006K2ZZ7K</td><td>A1UQRSCLF8GW1T</td><td>Michael D. Bigham "M. Wassir"</td><td>0</td><td>0</td><td>5</td><td>1350777600</td><td>Great taffy</td></tr><tr><td></td><td>5</td><td>6</td><td>B006K2ZZ7K</td><td>ADT0SRK1MGOEU</td><td>Twoapennything</td><td>0</td><td>0</td><td>4</td><td>1342051200</td><td>Nice Taffy</td></tr><tr><td></td><td>6</td><td>7</td><td>B006K2ZZ7K</td><td>A1SP2KV6FXXRU1</td><td>David C. Sullivan</td><td>0</td><td>0</td><td>5</td><td>1340150400</td><td>Great! Just as good as the expensive brands!</td></tr><tr><td></td><td>7</td><td>8</td><td>B006K2ZZ7K</td><td>A3JRGQVEQN31IQ</td><td>Pamela G. Williams</td><td>0</td><td>0</td><td>5</td><td>1336003200</td><td>Wonderful, tasty taffy</td></tr><tr><td></td><td>8</td><td>9</td><td>B000E7L2R4</td><td>A1MZY09TZK08BI</td><td>R. James</td><td>1</td><td>1</td><td>5</td><td>1322006400</td><td>Yay Barley</td></tr><tr><td></td><td>9</td><td>10</td><td>B00171APVA</td><td>A21BT40VZCCYT4</td><td>Carol A. Reed</td><td>0</td><td>0</td><td>5</td><td>1351209600</td><td>Healthy Dog Food</td></tr></table>												Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary		0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	5	1303862400	Good Quality Dog Food		1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	1	1346976000	Not as Advertised		2	3	B000LQOCH0	ABXLMWJIXAIN	Natalia Corres "Natalia Corres"	1	1	4	1219017600	"Delight" says it all		3	4	B000UA0QIQ	A395BORC6FGVXV	Karl	3	3	2	1307923200	Cough Medicine		4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0	0	5	1350777600	Great taffy		5	6	B006K2ZZ7K	ADT0SRK1MGOEU	Twoapennything	0	0	4	1342051200	Nice Taffy		6	7	B006K2ZZ7K	A1SP2KV6FXXRU1	David C. Sullivan	0	0	5	1340150400	Great! Just as good as the expensive brands!		7	8	B006K2ZZ7K	A3JRGQVEQN31IQ	Pamela G. Williams	0	0	5	1336003200	Wonderful, tasty taffy		8	9	B000E7L2R4	A1MZY09TZK08BI	R. James	1	1	5	1322006400	Yay Barley		9	10	B00171APVA	A21BT40VZCCYT4	Carol A. Reed	0	0	5	1351209600	Healthy Dog Food
		Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary																																																																																																																									
	0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	5	1303862400	Good Quality Dog Food																																																																																																																									
	1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	1	1346976000	Not as Advertised																																																																																																																									
	2	3	B000LQOCH0	ABXLMWJIXAIN	Natalia Corres "Natalia Corres"	1	1	4	1219017600	"Delight" says it all																																																																																																																									
	3	4	B000UA0QIQ	A395BORC6FGVXV	Karl	3	3	2	1307923200	Cough Medicine																																																																																																																									
	4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0	0	5	1350777600	Great taffy																																																																																																																									
	5	6	B006K2ZZ7K	ADT0SRK1MGOEU	Twoapennything	0	0	4	1342051200	Nice Taffy																																																																																																																									
	6	7	B006K2ZZ7K	A1SP2KV6FXXRU1	David C. Sullivan	0	0	5	1340150400	Great! Just as good as the expensive brands!																																																																																																																									
	7	8	B006K2ZZ7K	A3JRGQVEQN31IQ	Pamela G. Williams	0	0	5	1336003200	Wonderful, tasty taffy																																																																																																																									
	8	9	B000E7L2R4	A1MZY09TZK08BI	R. James	1	1	5	1322006400	Yay Barley																																																																																																																									
	9	10	B00171APVA	A21BT40VZCCYT4	Carol A. Reed	0	0	5	1351209600	Healthy Dog Food																																																																																																																									

In [9]:	<pre>positive = df[df['sentiment']==1] negative = df[df['sentiment']==-1]</pre>									
---------	---	--	--	--	--	--	--	--	--	--

## CREATING WORDCLOUD WHICH CONTAINS "POSITIVE" SENTIMENTS IN REVIEWS

```
plt.show()
```



## CREATING WORDCLOUD WHICH CONTAINS "NEGATIVE" SENTIMENTS IN REVIEWS

```
In [11]: neg = " ".join(review for review in negative.Summary)
wordcloud_neg = WordCloud(stopwords=stopwords).generate(neg)
plt.imshow(wordcloud_neg, interpolation='bilinear')
```

## CREATING WORDCLOUD WHICH CONTAINS "NEGATIVE" SENTIMENTS IN REVIEWS

In [11]:	<pre>neg = " ".join(review for review in negative.Summary) wordcloud_neg = WordCloud(stopwords=stopwords).generate(neg) plt.imshow(wordcloud_neg, interpolation='bilinear') plt.savefig('wordcloud03.png') plt.show()</pre>									

In [12]:	<pre>df['sentiment'] = df['sentiment'].replace({-1 : 'negative'}) df['sentiment'] = df['sentiment'].replace({1 : 'positive'})</pre>									
----------	---	--	--	--	--	--	--	--	--	--

## REMOVING ALL THE PUNCTUATION IN "SUMMARY" COLUMN TO MAKE USE OF IT FOR THE ALGORITHM

In [13]:	<pre>def remove_punctuation(text):     final = "".join(u for u in text if u not in ("?", ".", ";", ":", "!", "'"))     return final df['Text'] = df['Text'].apply(remove_punctuation) df = df.dropna(subset=['Summary']) df['Summary'] = df['Summary'].apply(remove_punctuation)</pre>									
----------	--	--	--	--	--	--	--	--	--	--

## CREATING NEW DATAFRAME BY SPLITTING THE "SUMMARY" AND "SENTIMENT" COLUMNS IN ORIGINAL DATASET

1	Not as Advertised	negative
2	Delight says it all	positive
3	Cough Medicine	negative
4	Great taffy	positive

## SPLITTING OUR DATASET INTO "TRAIN DATA" AND "TEST DATA" USING "RANDOM" SPLIT METHOD

```
[15]: import numpy as np
index = df.index
df['random_number'] = np.random.randn(len(index))
```

## SPLITTING OUR DATASET INTO "TRAIN DATA" AND "TEST DATA" USING "RANDOM" SPLIT METHOD

In [15]:	<pre>import numpy as np index = df.index df['random_number'] = np.random.randn(len(index)) train = df[df['random_number'] &lt;= 0.8] test = df[df['random_number'] &gt; 0.8]</pre>									
----------	--	--	--	--	--	--	--	--	--	--

## CREATING "BAG OF WORDS" MODEL WHICH CONVERTS TEXT INTO BAG OF WORDS FOR THE ALGORITHM

In [16]:	<pre>from sklearn.feature_extraction.text import CountVectorizer vectorizer = CountVectorizer(token_pattern=r'\b\w+\b') train_matrix = vectorizer.fit_transform(train['Summary']) test_matrix = vectorizer.transform(test['Summary'])</pre>									
----------	---	--	--	--	--	--	--	--	--	--

## BUILDING THE LOGISTIC REGRESSION ALGORITHM FOR OUR DATASET

In [17]:	<pre>from sklearn.linear_model import LogisticRegression lr = LogisticRegression(solver='lbfgs', max_iter=1000)</pre>									
----------	---	--	--	--	--	--	--	--	--	--

In [18]:	<pre>x_train = train_matrix x_test = test_matrix y_train = train['sentiment'] y_test = test['sentiment'] lr.fit(x_train, y_train)</pre>									
----------	---	--	--	--	--	--	--	--	--	--

Out[18]:	LogisticRegression(max_iter=1000)									
----------	-----------------------------------	--	--	--	--	--	--	--	--	--

In [19]:	<pre>predictions = lr.predict(x_test)</pre>									
----------	---	--	--	--	--	--	--	--	--	--

## TESTING OUR LOGISTIC REGRESSION MODEL TO FIND "ACCURACY", "PRECISION", "RECALL"

In [20]:	<pre>from sklearn.metrics import confusion_matrix, classification_report new = np.asarray(y_test) confusion_matrix(predictions, y_test)</pre>									
----------	---	--	--	--	--	--	--	--	--	--

Out[20]:	array([[11476, 2394],        [ 5794, 91747]], dtype=int64)									
----------	--	--	--	--	--	--	--	--	--	--

In [21]:	<pre>print(classification_report(predictions, y_test))</pre>									
----------	--	--	--	--	--	--	--	--	--	--

	precision	recall	f1-score	support
negative	0.66	0.83	0.74	13870
positive	0.97	0.94	0.96	97541
accuracy			0.93	111411
macro avg	0.82	0.88	0.85	111411
weighted avg	0.94	0.93	0.93	111411