

Amazon Food Reviews

Data Source:<https://www.kaggle.com/snap/amazon-fine-food-reviews>
[\(https://www.kaggle.com/snap/amazon-fine-food-reviews\)](https://www.kaggle.com/snap/amazon-fine-food-reviews)

This dataset consists of reviews of fine foods from Amazon. The data span a period of more than 10 years, including all ~500,000 reviews up to October 2012. Reviews include product and user information, ratings, and a plain text review. It also includes reviews from all other Amazon categories.



Data includes:

- Reviews from Oct 1999 - Oct 2012
- 568,454 reviews
- 256,059 users
- 74,258 products
- 260 users with > 50 reviews

Attribute Information:

1. Id
2. ProductId - unique identifier for the product
3. UserId - unique identifier for the user
4. ProfileName

5. HelpfulnessNumerator - number of users who found the review helpful
6. HelpfulnessDenominator - number of users who indicated whether they found the review helpful or not
7. Score - rating between 1 and 5
8. Time - timestamp for the review
9. Summary - brief summary of the review
10. Text - text of the review

Number of people who found the review helpful

Number of people who indicated whether or not the review was helpful

129 of 134 people found the following review helpful

 What a great TV. When the decision came down to either ...

By Cimmerian on November 20, 2014

What a great TV. When the decision came down to either sending my kids to college or buying this set, the choice was easy. Now my kids can watch this set when they come home from their McJobs and be happy like me.

1 Comment

Was this review helpful to you?

Rating

-Product ID

-Reviewer User ID

Summary

Review

Objective:- Review Polarity

Given a review, determine the review is positive or negative

1.Naive Way

Naive way to do this will be to say Score with 1 & 2 -> Negative and 4 & 5 -> positive and review with score 3 is ignored and we consider it as neutral

2. Using text review to decide the polarity

Take the summary and text of review and analyze it using NLP whether the customer feedback/review is positive or negative

In [2]: #Imports

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import sqlite3 as sql
import seaborn as sns
from time import time
import gensim
import random
import warnings

warnings.filterwarnings("ignore")

%matplotlib inline
# sets the backend of matplotlib to the 'inline' backend:
#With this backend, the output of plotting commands is displayed inline within fro
#directly below the code cell that produced it. The resulting plots will then also

#Pickle python objects to file
import pickle
def savetofile(obj,filename):
    pickle.dump(obj,open(filename+".p","wb"), protocol=4)
def openfromfile(filename):
    temp = pickle.load(open(filename+".p","rb"))
    return temp
```

In [4]: # !wget --header="Host: storage.googleapis.com" --header="User-Agent: Mozilla/5.0

First Let's do the EDA

Loading the data

In [60]: #Using sqlite3 to retrieve data from sqlite file

```
con = sql.connect("database.sqlite")#Connection object that represents the database

#Using pandas functions to query from sql table
df = pd.read_sql_query("""
SELECT * FROM Reviews
""",con)

#Reviews is the name of the table given
#Taking only the data where score != 3 as score 3 will be neutral and it won't help us
df.head()
```

Out[60]:

	Id	ProductId		UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
0	1	B001E4KFG0	A3SGXH7AUHU8GW		delmartian		1
1	2	B00813GRG4	A1D87F6ZCVE5NK		dll pa		0
2	3	B000LQOCHO	ABXLMWJIXXAIN		Natalia Corres "Natalia Corres"		1
3	4	B000UA0QIQ	A395BORC6FGVXV		Karl		3
4	5	B006K2ZZ7K	A1UQRSCLF8GW1T		Michael D. Bigham "M. Wassir"		0



In [61]: `df.describe()`

Out[61]:

	Id	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time
count	568454.000000	568454.000000	568454.000000	568454.000000	5.684540e+05
mean	284227.500000	1.743817	2.22881	4.183199	1.296257e+09
std	164098.679298	7.636513	8.28974	1.310436	4.804331e+07
min	1.000000	0.000000	0.00000	1.000000	9.393408e+08
25%	142114.250000	0.000000	0.00000	4.000000	1.271290e+09
50%	284227.500000	0.000000	1.00000	5.000000	1.311120e+09
75%	426340.750000	2.000000	2.00000	5.000000	1.332720e+09
max	568454.000000	866.000000	923.00000	5.000000	1.351210e+09

◀ ▶

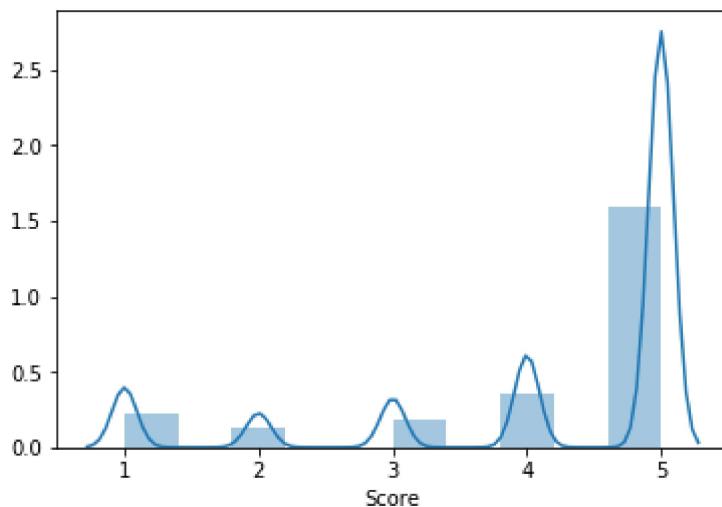
In [62]: `df.shape`

`df['Score'].size`

Out[62]: 568454

In [63]: `# sns.FacetGrid(df).\\`
`# map(plt.hist,bins=df['Score'])`
`# # plt.plot(df['Score'].value_counts())`
`# plt.show()`

In [64]: `sns.distplot(df['Score'],bins=10)`
`plt.show()`



```
In [65]: df['Score'].value_counts()
```

```
Out[65]: 5    363122  
4    80655  
1    52268  
3    42640  
2    29769  
Name: Score, dtype: int64
```

```
In [66]: #Using pandas functions to query from sql table  
df = pd.read_sql_query("""  
SELECT * FROM Reviews  
WHERE Score != 3  
""",con)
```

1. Naive Way

Score as positive or negative

```
In [67]: def polarity(x):
    if x < 3:
        return 'Negative'
    else:
        return 'Positive'
df["Score"] = df["Score"].map(polarity) #Map all the scores as the function polarity
df.head()
```

Out[67]:

	Id	ProductId		UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
0	1	B001E4KFG0	A3SGXH7AUHU8GW		delmartian		1
1	2	B00813GRG4	A1D87F6ZCVE5NK		dll pa		0
2	3	B000LQOCH0	ABXLMWJIXXAIN		Natalia Corres "Natalia Corres"		1
3	4	B000UA0QIQ	A395BORC6FGVXV		Karl		3
4	5	B006K2ZZ7K	A1UQRSCLF8GW1T		Michael D. Bigham "M. Wassir"		0

Using Score column now we can say either a Review is positive or negative

2.Using Text data and Natural Language Processing (NLP)

Firstly we need to perform some data cleaning and then text preprocessing and convert the texts as vectors so that we can train some model on those vectors and predict polarity of the review

1.Data Cleaning

(i) Data Deduplication

```
In [68]: df.duplicated(subset={"UserId", "ProfileName", "Time", "Text"}).value_counts()
```

```
Out[68]: False    364173
True     161641
dtype: int64
```

There exist a lot of duplicates wherein the different products is **reviewed by same user at the same time**

The product ID may be different but the product is similar with different variant

```
In [69]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan		2
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan		2
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan		2
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan		2
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan		2

Geeta gave the review at the same time for multiple products which is not possible ethically, the products were same but different flavours hence counted as multiple products

In [70]: #Deleting all the duplicates having the same userID, Profile, NameTime and Text a
df1 = df.drop_duplicates(subset={"UserId", "ProfileName", "Time", "Text"}, keep="first")

In [71]: size_diff = df1['Id'].size / df['Id'].size
print("%.1f %% reduction in data after deleting duplicates" % ((1 - size_diff) * 100))
print("Size of data", df1['Id'].size, " rows ")

30.7 % reduction in data after deleting duplicates
Size of data 364173 rows

(ii) Helpfullness Numerator Greater than Helpfullness Denominator

In [72]: df2 = df1[df1.HelpfulnessNumerator <= df1.HelpfulnessDenominator]
print("Size of data", df2['Id'].size, " rows ")

Size of data 364171 rows

Text Preprocessing

[1] HTML Tag Removal

In [73]: import re #Regular Expr Operations
#string = r"sdfsdfd" :- r is for raw string as Regex often uses \ backslashes(\w).

#####Function to remove html tags from data
def striphtml(data):
 p = re.compile('<.*?>')#Find this kind of pattern
 # print(p.findall(data))#List of strings which follow the regex pattern
 return p.sub('', data) #Substitute nothing at the place of strings which match.

striphml('I Want This text!<>')

Out[73]: 'I Want This text!'

[2] Punctuations Removal

In [74]: #####Function to remove All the punctuations from the text
def strippunc(data):
 p = re.compile(r'[?|!|"|#|.!,|)|(|\|/|~|%|*|]')
 return p.sub('', data)
strippunc("fsd*?~,,(sdfsdfdsvv)#")

Out[74]: 'fsd sdfsd fdsvv'

[3] Stopwords

Stop words usually refers to the most common words in a language are generally filtered out before or after processing of natural language data. Sometimes it is avoided to remove the stop words to support phrase search.

In [75]:

```
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

stop = stopwords.words('english') #All the stopwords in English Language
#excluding some useful words from stop words list as we doing sentiment analysis
excluding = ['against', 'not', 'don', "don't", 'ain', 'aren', "aren't", 'couldn', "c
    'doesn', "doesn't", 'hadn', "hadn't", 'hasn', "hasn't", 'haven', "ha
    'mightn', "mightn't", 'mustn', "mustn't", 'needn', "needn't", 'should
    "wasn't", 'weren', "weren't", 'won', "won't", 'wouldn', "wouldn't"]
stop = [words for words in stop if words not in excluding]
print(stop)
```

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're",
"you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he',
'him', 'his', 'himself', 'she', "she's", 'her', 'hers', 'herself', 'it', 'i
t's', 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 'what',
'which', 'who', 'whom', 'this', 'that', "that'll", 'these', 'those', 'am', 'i
s', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'havin
g', 'do', 'does', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or',
'because', 'as', 'until', 'while', 'of', 'at', 'by', 'for', 'with', 'about', 'b
etween', 'into', 'through', 'during', 'before', 'after', 'above', 'below', 't
o', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'again',
'further', 'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'al
l', 'any', 'both', 'each', 'few', 'more', 'most', 'other', 'some', 'such', 'n
o', 'nor', 'only', 'own', 'same', 'so', 'than', 'too', 'very', 's', 't', 'can',
'will', 'just', 'should', "should've", 'now', 'd', 'll', 'm', 'o', 're', 've',
'y', 'ma', 'shan', "shan't"]
```

[4] Stemming

Porter Stemmer: Most commonly used stemmer without a doubt, also one of the most gentle stemmers. Though it is also the most computationally intensive of the algorithms. It is also the oldest stemming algorithm by a large margin.

SnowBall Stemmer(Porter2): Nearly universally regarded as an improvement over porter, and for good reason. Porter himself in fact admits that it is better than his original algorithm. Slightly faster computation time than Porter, with a fairly large community around it.

```
In [86]: from nltk.stem import SnowballStemmer
snow = SnowballStemmer('english') #initialising the snowball stemmer
print("Stem/Root words of the some of the words using SnowBall Stemmer:")
print(snow.stem('tasty'))
print(snow.stem('tasteful'))
print(snow.stem('tastiest'))
print(snow.stem('delicious'))
print(snow.stem('amazing'))
print(snow.stem('amaze'))
print(snow.stem('initialize'))
print(snow.stem('fabulous'))
print(snow.stem('Honda City'))
print(snow.stem('unpleasant'))
```

Stem/Root words of the some of the words using SnowBall Stemmer:

tasti
tast
tastiest
delici
amaz
amaz
initi
fabul
honda c
unpleas

Stemming and Lemmatization Differences

- Both lemmatization and stemming attempt to bring a canonical form for a set of related word forms.
- Lemmatization takes the part of speech into consideration. For example, the term 'meeting' may either be returned as 'meeting' or as 'meet' depending on the part of speech.
- Lemmatization often uses a tagged vocabulary (such as Wordnet) and can perform more sophisticated normalization. E.g. transforming mice to mouse or foci to focus.
- Stemming implementations, such as the Porter's stemmer, use heuristics that truncates or transforms the end letters of the words with the goal of producing a normalized form. Since this is algorithm based, there is no requirement of a vocabulary.
- Some stemming implementations may combine a vocabulary along with the algorithm. Such an approach for example converts 'cars' to 'automobile' or even 'Honda City', 'Mercedes Benz' to a common word 'automobile'
- A stem produced by typical stemmers may not be a word that is part of a language vocabulary but lemmatizer transforms the given word forms to a valid lemma.

Preprocessing output for one review

In [87]:

```

str1=' '
final_string=[]
all_positive_words=[] # store words from +ve reviews here
all_negative_words=[] # store words from -ve reviews here.
s=''

for sent in df2['Text'][2:3].values: #Running only for 2nd review
    filtered_sentence=[]
    print(sent) #Each review
    sent=striphtml(sent)# remove HTML tags
    sent=strippunc(sent)# remove Punctuation Symbols
    print(sent.split())
    for w in sent.split():
        print("=====>",w)
        if((w.isalpha()) and (len(w)>2)):#If it is a numerical value or character
            if(w.lower() not in stop):# If it is a stopword
                s=(snow.stem(w.lower())).encode('utf8') #Stemming the word using .
                print("Selected: Stem Word->",s)
                filtered_sentence.append(s)
            else:
                print("Eliminated as it is a stopword")
                continue
        else:
            print("Eliminated as it is a numerical value or character of lenght 1")
            continue
    #     print(filtered_sentence)
    str1 = b" ".join(filtered_sentence) #final string of cleaned words

    final_string.append(str1)
print("*****")
print("Finally selected words from the review:\n",final_string)

```

This is a confection that has been around a few centuries. It is a light, pillowy citrus gelatin with nuts - in this case Filberts. And it is cut into tiny squares and then liberally coated with powdered sugar. And it is a tiny mouthful of heaven. Not too chewy, and very flavorful. I highly recommend this yummy treat. If you are familiar with the story of C.S. Lewis' "The Lion, The Witch, and The Wardrobe" - this is the treat that seduces Edmund into selling out his Brother and Sisters to the Witch.

['This', 'is', 'a', 'confection', 'that', 'has', 'been', 'around', 'a', 'few', 'centuries', 'It', 'is', 'a', 'light', 'pillowy', 'citrus', 'gelatin', 'with', 'nuts', '-', 'in', 'this', 'case', 'Filberts', 'And', 'it', 'is', 'cut', 'into', 'tiny', 'squares', 'and', 'then', 'liberally', 'coated', 'with', 'powdered', 'sugar', 'And', 'it', 'is', 'a', 'tiny', 'mouthful', 'of', 'heaven', 'Not', 'too', 'chewy', 'and', 'very', 'flavorful', 'I', 'highly', 'recommend', 'this', 'yummy', 'treat', 'If', 'you', 'are', 'familiar', 'with', 'the', 'story', 'of', 'CS', 'Lewis', 'The', 'Lion', 'The', 'Witch', 'and', 'The', 'Wardrobe', '-', 'this', 'is', 'the', 'treat', 'that', 'seduces', 'Edmund', 'into', 'selling', 'out', 'his', 'Brother', 'and', 'Sisters', 'to', 'the', 'Witch']

=====> This

Preprocessing on all the reviews

In [88]:

```
%time
# Code takes a while to run as it needs to run on around 500k sentences.
i=0
str1=' '
final_string=[]
all_positive_words=[] # store words from +ve reviews here
all_negative_words=[] # store words from -ve reviews here.
s='.'
t0=time()
for sent in df2['Text'].values:
    filtered_sentence=[]
    # print(sent) #Each review
    sent=striphtml(sent)# remove HTML tags
    sent=strippunc(sent)# remove Punctuation Symbols
    # print(sent.split())
    for w in sent.split():
        # print("=====>",w)
        if((w.isalpha() and (len(w)>2)):#If it is a numerical value or character
           if(w.lower() not in stop):# If it is a stopword
               s=(snow.stem(w.lower())).encode('utf8') #Stemming the word using .
                                         #encoding as byte-string/utf-8
        #
        # print("Selected: Stem Word->",s)
        filtered_sentence.append(s)
        if (df2['Score'].values)[i] == 'Positive':
            all_positive_words.append(s) #list of all words used to descr
        if(df2['Score'].values)[i] == 'Negative':
            all_negative_words.append(s) #list of all words used to descr
        else:
            print("Eliminated as it is a stopword")
            continue
        else:
            print("Eliminated as it is a numerical value or character of Length"
                  continue
        #
        # print(filtered_sentence)
        str1 = b" ".join(filtered_sentence) #final string of cleaned words
                                         #encoding as byte-string/utf-8

        final_string.append(str1)
        #
        # print("*****")
        # print("Finally selected words from the review:\n",final_string)
        i+=1
print("Preprocessing completed in ")
```

Preprocessing completed in
CPU times: user 9min 7s, sys: 200 ms, total: 9min 8s
Wall time: 9min 8s

Cleaned text Without Stemming for Google trained W2Vec

In [76]:

```
%time
# Code takes a while to run as it needs to run on around 500k sentences.
i=0
str1=' '
final_string_nostem=[]
all_positive_words=[] # store words from +ve reviews here
all_negative_words=[] # store words from -ve reviews here.
s=''
t0=time()
for sent in df2['Text'].values:
    filtered_sentence=[]
    sent=striphtml(sent)# remove HTML tags
    sent=strippunc(sent)# remove Punctuation Symbols
    for w in sent.split():
        if((w.isalpha()) and (len(w)>2)):#If it is a numerical value or character
            if(w.lower() not in stop):# If it is a stopword
                s=w.lower().encode('utf8') #encoding as byte-string/utf-8
            else:
                continue
        else:
            continue
    str1 = b" ".join(filtered_sentence)
    final_string_nostem.append(str1)
    i+=1
print("Preprocessing completed in ")
```

Preprocessing completed in 245.63548946380615
CPU times: user 4min 5s, sys: 164 ms, total: 4min 5s
Wall time: 4min 5s

The above code uses string as byte-string / utf-8(uses 1 byte), Python default stores string as Unicode / (utf16/utf32) {depends on how python was compiled}-(uses 2/4 byte) as our data is large 1 byte difference can save a lot of memory. Hence encoding the data as byte-string
For more info: <https://stackoverflow.com/questions/10060411/byte-string-vs-unicode-string-python>
[\(https://stackoverflow.com/questions/10060411/byte-string-vs-unicode-string-python\)](https://stackoverflow.com/questions/10060411/byte-string-vs-unicode-string-python)

Positive and Negative words in reviews

```
In [24]: from collections import Counter
print("No. of positive words:", len(all_positive_words))
print("No. of negative words:", len(all_negative_words))
# print("Sample positive words", all_positive_words[:9])
# print("Sample negative words", all_negative_words[:9])
positive = Counter(all_positive_words)
print("\nMost Common positive words", positive.most_common(10))
negative = Counter(all_negative_words)
print("\nMost Common negative words", negative.most_common(10))
```

No. of positive words: 11678044

No. of negative words: 2393854

Most Common positive words [(b'not', 145019), (b'like', 138335), (b'tast', 126024), (b'good', 109838), (b'love', 106551), (b'flavor', 106408), (b'use', 102872), (b'great', 101125), (b'one', 94396), (b'product', 88466)]

Most Common negative words [(b'not', 53634), (b'tast', 33828), (b'like', 32059), (b'product', 27411), (b'one', 20176), (b'flavor', 18898), (b'would', 17858), (b'tri', 17515), (b'use', 15148), (b'good', 14616)]

- "tast" , "like" , "flavor", "good" and "one" are some of the most common words in both negative and positive reviews
- "good" and "great" are some of the most common words in positive reviews
- "would" and "coffe" are some of the most common words in negative reviews
- tasty, good, etc are some of the words common in both **because there may be a not before it like "not tasty" , "not good"**

Storing our preprocessed data in DB

In [89]: `#Adding a column of CleanedText which displays the data after pre-processing of the reviews`
`df2['CleanedText']=final_string`
`df2['CleanedText_NoStem']=final_string_nostem`
`df2.head(3)`

Out[89]:

	Id	ProductId		UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator
0	1	B001E4KFG0	A3SGXH7AUHU8GW		delmartian		1
1	2	B00813GRG4	A1D87F6ZCVE5NK		dll pa		0
2	3	B000LQOCH0	ABXLMWJIXXAIN		Natalia Corres "Natalia Corres"		1

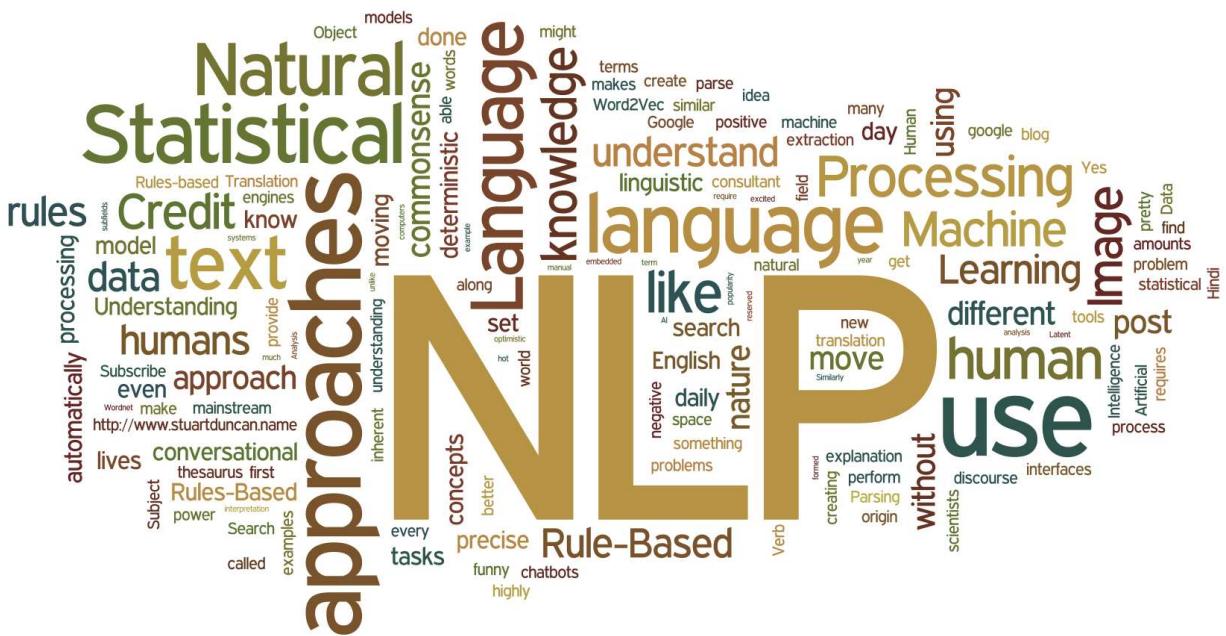


In [81]: `df2['CleanedText_NoStem'][1]`

Out[81]: `b'product arrived labeled jumbo salted peanuts the peanuts actually small sized unsalted not sure error vendor intended represent product jumbo'`

In [90]: `### Storing dataframe in sqlite3`
`import sqlite3`
`con = sqlite3.connect('final.sqlite')`
`con.text_factory = str #To store the string as byte strings only`
`df2.to_sql('Reviews', con, if_exists='replace')`

In [6]: `#Using sqlite3 to retrieve data from sqlite file`
`con = sql.connect("final.sqlite")#Loading Cleaned/ Preprocessed text that we did in step 8`
`#Using pandas functions to query from sql table`
`df2 = pd.read_sql_query("""`
`SELECT * FROM Reviews`
`""",con)`



Some Key NLP Terms:

Natural Language Processing (NLP)

A Computer Science field connected to Artificial Intelligence and Computational Linguistics which focuses on interactions between computers and human language and a machine's ability to understand, or mimic the understanding of human language. Examples of NLP applications include Siri and Google Now.

Information Extraction

The process of automatically extracting structured information from unstructured and/or semi-structured sources, such as text documents or web pages for example.

Sentiment Analysis

The use of Natural Language Processing techniques to extract subjective information from a piece of text. i.e. whether an author is being subjective or objective or even positive or negative. (can also be referred to as Opinion Mining). As in this case we doing sentiment analysis of reviews of users from Amazon.

Data Corpus or Corpora

A usually large collection of documents that can be used to infer and validate linguistic rules, as well as to do statistical analysis and hypothesis testing. e.g. The Amazon Fine Food Review dataset is a corpus.

Document

A "document" is a distinct text, you could treat an individual paragraph or even sentence as a "document".

In our case our each review is a document

Bag of Words (BoW)

A commonly used model in methods of Text Classification. As part of the BOW model, a piece of text (sentence or a document) is represented as a bag or multiset of words, disregarding grammar and even word order and the frequency or occurrence of each word is used as a feature for training a classifier.

OR

Simply,Converting a collection of text documents to a matrix of token counts

Ways to convert text to vector

1. Uni-gram BOW

```
In [8]: from sklearn.feature_extraction.text import CountVectorizer
```

```
In [29]: %%time
uni_gram = CountVectorizer() #in scikit-Learn
uni_gram_vectors = uni_gram.fit_transform(df2['CleanedText'].values)
```

CPU times: user 17.9 s, sys: 120 ms, total: 18.1 s
Wall time: 18.1 s

```
In [30]: #Saving the variable to access later without recomputing
savetofile(uni_gram_vectors,"uni_gram")
```

```
In [16]: #Loading the variable from file
uni_gram_vectors = openfromfile("uni_gram")
```

```
In [31]: uni_gram_vectors.shape[1]
```

Out[31]: 209129

```
In [11]: uni_gram_vectors[0]
```

Out[11]: <1x209129 sparse matrix of type '<class 'numpy.int64'>'
with 20 stored elements in Compressed Sparse Row format>

```
In [12]: type(uni_gram_vectors)
```

Out[12]: scipy.sparse.csr.csr_matrix

```
In [95]: %%time
from sklearn.decomposition import TruncatedSVD

tsvd_uni = TruncatedSVD(n_components=1000)#No of components as total dimensions
tsvd_uni_vec = tsvd_uni.fit_transform(uni_gram_vectors)
```

CPU times: user 30min 43s, sys: 29.7 s, total: 31min 13s
Wall time: 9min

```
In [96]: savetofile(tsvd_uni,"tsvd_uni")
savetofile(tsvd_uni_vec,"tsvd_uni_vec")
```

```
In [6]: tsvd_uni = openfromfile("tsvd_uni")
tsvd_uni_vec = openfromfile("tsvd_uni_vec")
```

```
In [27]: tsvd_uni.explained_variance_ratio_[:].sum()
```

Out[27]: 0.82439113222951488

```
In [28]: %%time
from sklearn.manifold import TSNE
from time import time
import random

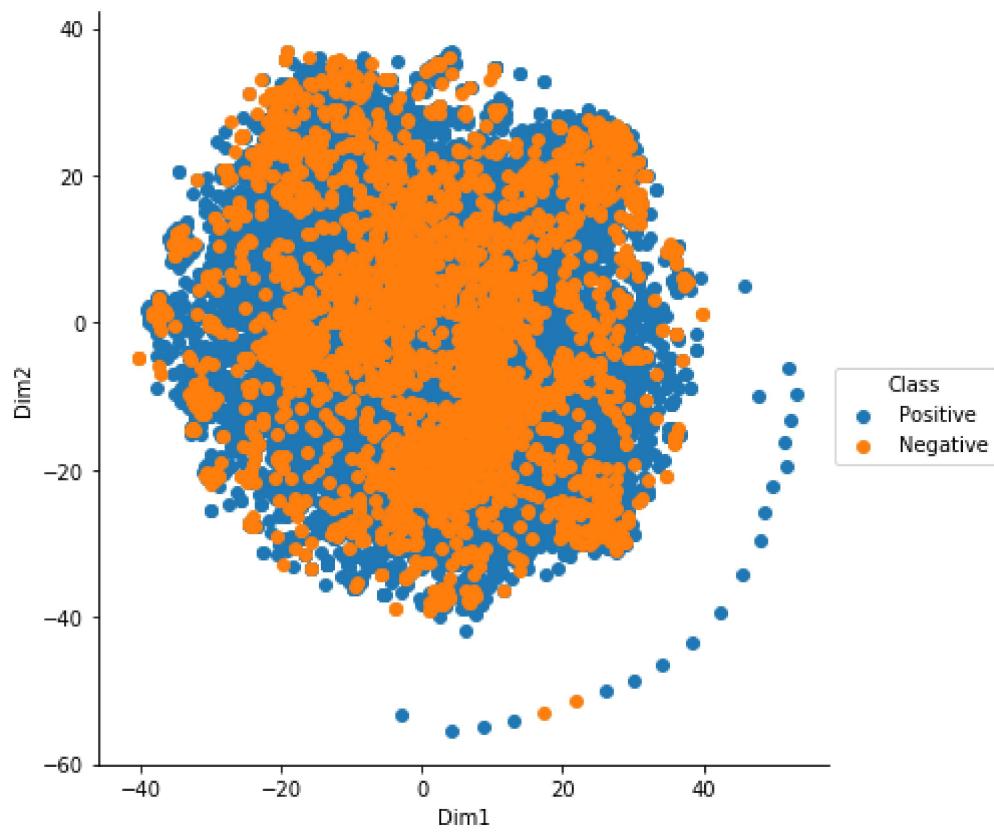
n_samples = 20000
sample_cols = random.sample(range(1, tsvd_uni_vec.shape[0]), n_samples)
sample_features = tsvd_uni_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=30)
# print(sample_features,sample_class)

t0 = time()
embedded_data = model.fit_transform(sample_features)
print("TSNE done in %0.3fs." % (time() - t0))
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])
```

(20000, 1000) (20000, 1)
TSNE done in 1716.121s.
(20000, 3)
CPU times: user 27min 30s, sys: 1min 5s, total: 28min 36s
Wall time: 28min 36s

In [29]: `#Perplexity = 30`

```
sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend  
plt.show()
```



In [32]:

```
%time
#Perplexity = 40
from sklearn.manifold import TSNE
from time import time
import random

n_samples = 20000
sample_cols = random.sample(range(1, tsvd_uni_vec.shape[0]), n_samples)
sample_features = tsvd_uni_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=40)
# print(sample_features,sample_class)

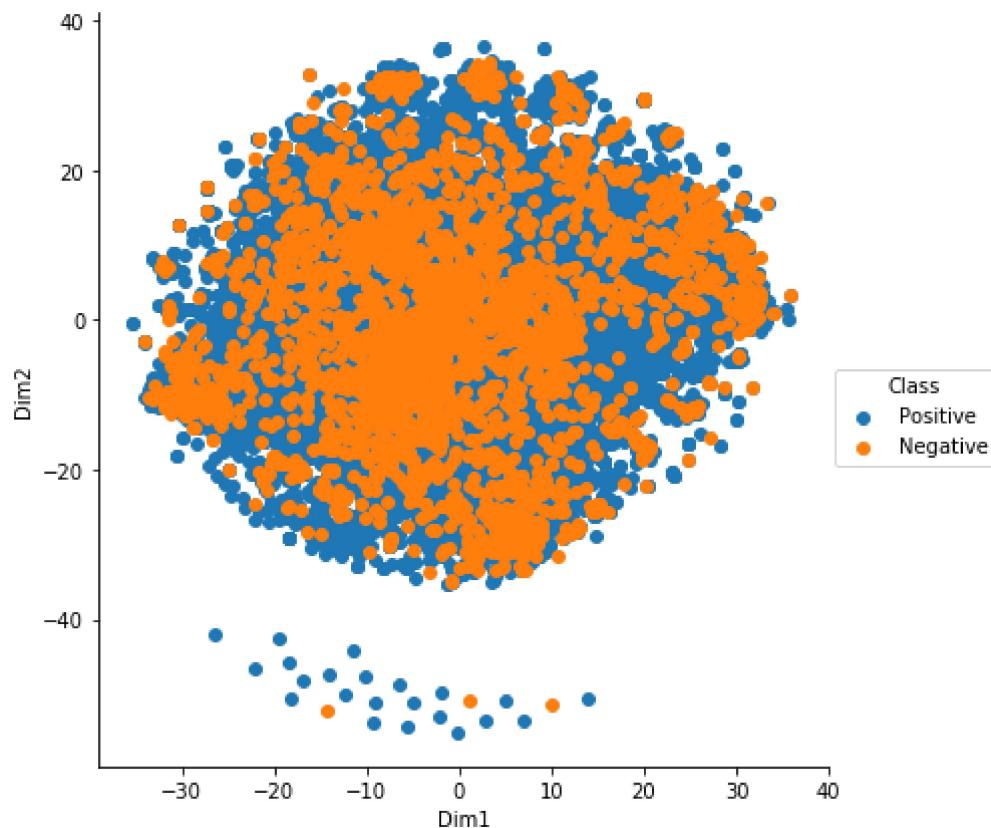
t0 = time()
embedded_data = model.fit_transform(sample_features)
print("TSNE done in %0.3fs." % (time() - t0))
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])

sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend
plt.show()
```

(20000, 1000) (20000, 1)

TSNE done in 1952.465s.

(20000, 3)



```
CPU times: user 31min 27s, sys: 1min 5s, total: 32min 33s
Wall time: 32min 33s
```

In [7]:

```
%%time
#Perplexity = 30 with 10k points
from sklearn.manifold import TSNE
from time import time
import random

n_samples = 10000
sample_cols = random.sample(range(1, tsvd_uni_vec.shape[0]), n_samples)
sample_features = tsvd_uni_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=20)
# print(sample_features,sample_class)

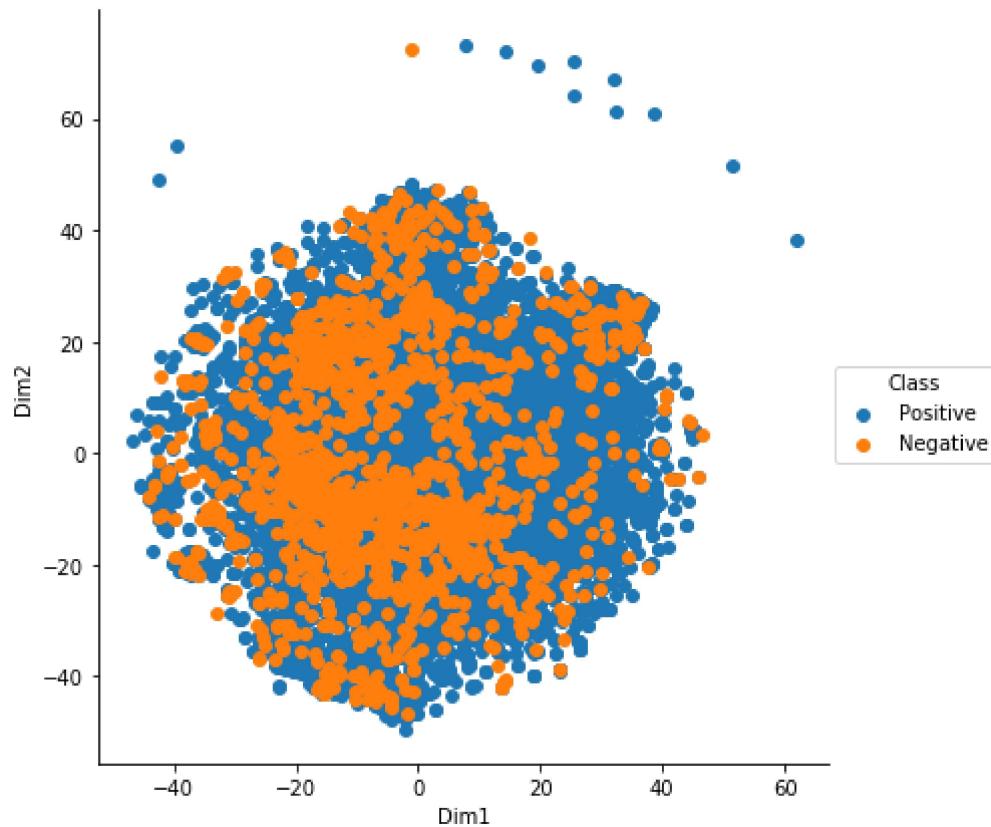
t0 = time()
embedded_data = model.fit_transform(sample_features)
print("TSNE done in %0.3fs." % (time() - t0))
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])

sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend
plt.show()
```

(10000, 1000) (10000, 1)

TSNE done in 608.475s.

(10000, 3)



CPU times: user 9min 43s, sys: 25.4 s, total: 10min 9s
 Wall time: 10min 9s

2. Bi-gram BOW

```
In [9]: %%time
#taking one words and two consecutive words together
bi_gram = CountVectorizer(ngram_range=(1,2))
bi_gram_vectors = bi_gram.fit_transform(df2['CleanedText'].values)
```

CPU times: user 58.2 s, sys: 728 ms, total: 59 s
 Wall time: 59 s

```
In [10]: #Saving the variable to access later without recomputing
savetofile(bi_gram_vectors,"bi_gram")
```

```
In [11]: #Loading the variable from file
bi_gram_vectors = openfromfile("bi_gram")
```

```
In [12]: bi_gram_vectors.shape
```

```
Out[12]: (364171, 3404647)
```

```
In [13]: bi_gram_vectors[0]
```

```
Out[13]: <1x3404647 sparse matrix of type '<class 'numpy.int64'>'  

with 42 stored elements in Compressed Sparse Row format>
```

```
In [14]: type(bi_gram_vectors)
```

```
Out[14]: scipy.sparse.csr.csr_matrix
```

```
In [17]: print("bi-gram is %.2f times more than uni-gram"%((bi_gram_vectors.shape[1]/uni_g
bi-gram is 16.28 times more than uni-gram
```

```
In [18]: %%time
from sklearn.decomposition import TruncatedSVD
sample_points = df2.sample(20000)

bi_gram = CountVectorizer(ngram_range=(1,2))
bi_gram_vectors = bi_gram.fit_transform(sample_points['CleanedText'])
tsvd_bi = TruncatedSVD(n_components=2500)#No of components as total dimensions
tsvd_bi_vec = tsvd_bi.fit_transform(bi_gram_vectors)
```

CPU times: user 1h 2min 47s, sys: 1min 4s, total: 1h 3min 51s
 Wall time: 16min 54s

```
In [21]: savetofile(tsvd_bi,"tsvd_bi")
savetofile(tsvd_bi_vec,"tsvd_bi_vec")
```

```
In [24]: tsvd_bi = openfromfile("tsvd_bi")
tsvd_bi_vec = openfromfile("tsvd_bi_vec")
```

```
In [25]: tsvd_bi.explained_variance_ratio_[:].sum()
```

```
Out[25]: 0.72117200409365134
```

In [36]:

```
%time
#Perplexity = 30 with 10k points
from sklearn.manifold import TSNE
from time import time
import random

n_samples = 10000
sample_cols = random.sample(range(1, tsvd_bi_vec.shape[0]), n_samples)
sample_features = tsvd_bi_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=30)
# print(sample_features,sample_class)

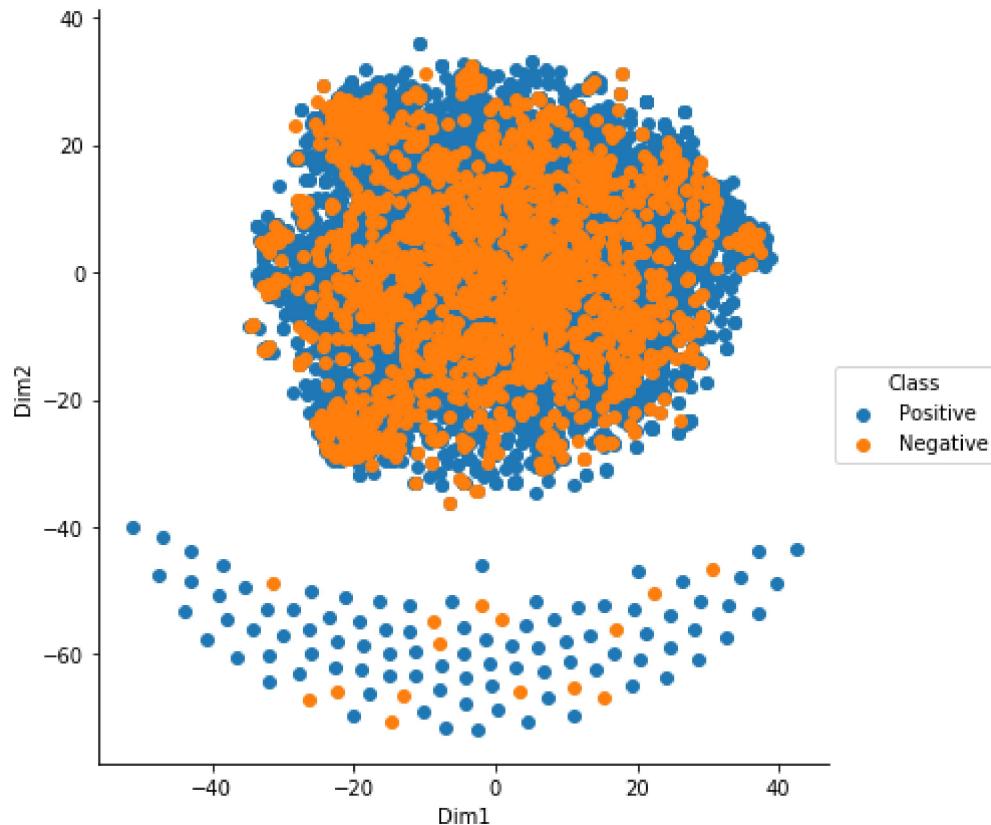
t0 = time()
embedded_data = model.fit_transform(sample_features)
print("TSNE done in %0.3fs." % (time() - t0))
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])

sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend()
plt.show()
```

(10000, 2500) (10000, 1)

TSNE done in 924.456s.

(10000, 3)



```
CPU times: user 14min 53s, sys: 31.7 s, total: 15min 25s
Wall time: 15min 25s
```

In [39]:

```

%%time
#Perplexity = 20 with 10k points
from sklearn.manifold import TSNE
from time import time
import random

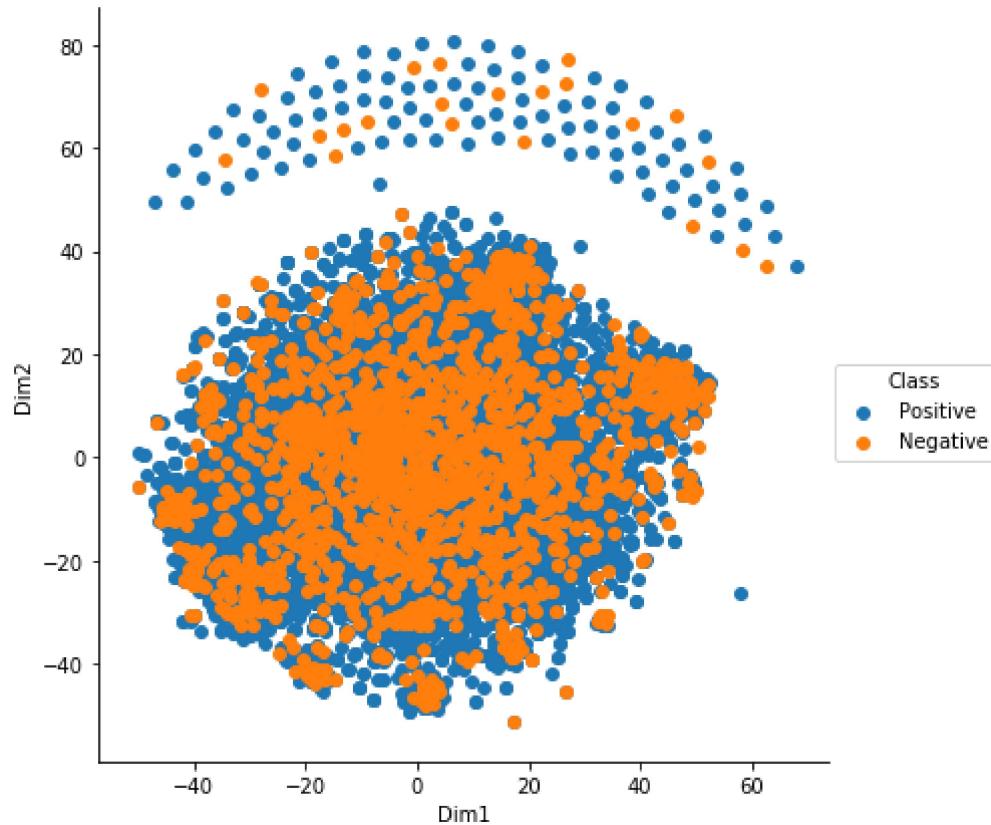
n_samples = 10000
sample_cols = random.sample(range(1, tsvd_bi_vec.shape[0]), n_samples)
sample_features = tsvd_bi_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=20)
# print(sample_features,sample_class)

t0 = time()
embedded_data = model.fit_transform(sample_features)
print("TSNE done in %0.3fs." % (time() - t0))
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])

sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend
plt.show()

```

(10000, 2500) (10000, 1)
 TSNE done in 870.144s.
 (10000, 3)



```
CPU times: user 14min, sys: 30.8 s, total: 14min 30s
Wall time: 14min 30s
```

In [40]:

```

%%time
#Perplexity = 40 with 10k points
from sklearn.manifold import TSNE
from time import time
import random

n_samples = 10000
sample_cols = random.sample(range(1, tsvd_bi_vec.shape[0]), n_samples)
sample_features = tsvd_bi_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=40)
# print(sample_features,sample_class)

t0 = time()
embedded_data = model.fit_transform(sample_features)
print("TSNE done in %0.3fs." % (time() - t0))
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])

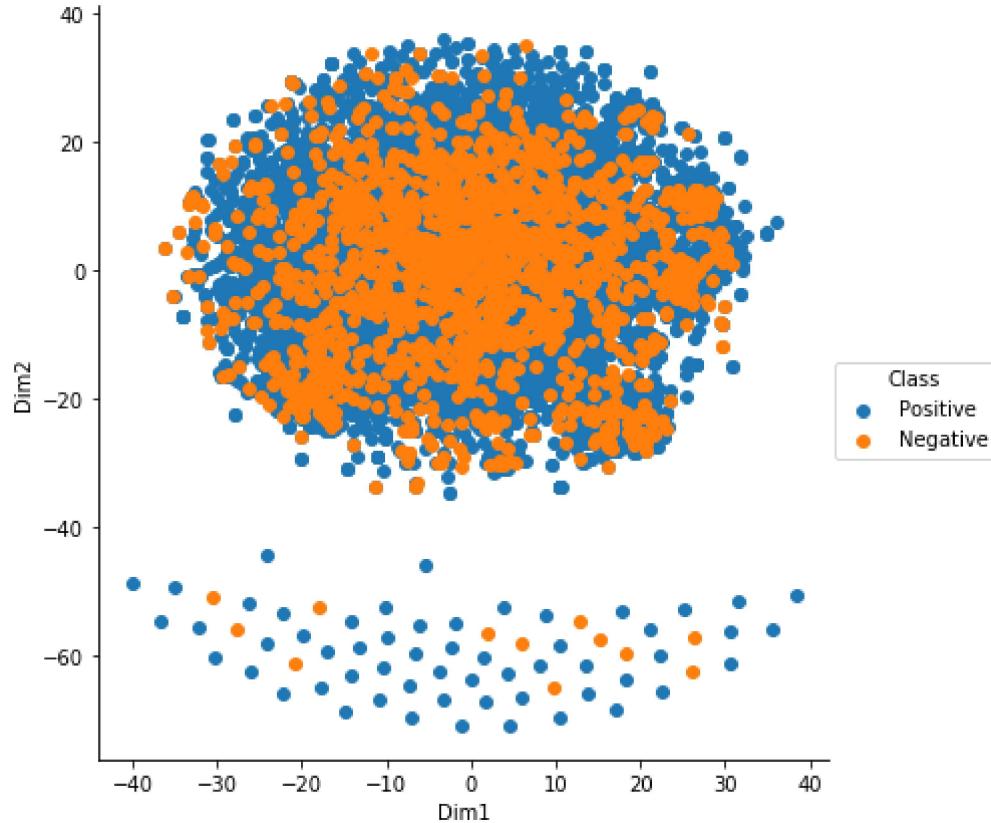
sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend
plt.show()

```

(10000, 2500) (10000, 1)

TSNE done in 1505.767s.

(10000, 3)



CPU times: user 24min 31s, sys: 34.6 s, total: 25min 6s
 Wall time: 25min 6s

3. tf-idf

$$\text{TFIDF} = \text{TF} \times \text{IDF}$$

Term Frequency: This summarizes how often a given word appears within a document.

Inverse Document Frequency: This downscals words that appear a lot across documents in the corpus.

In information retrieval, tf-idf or TFIDF, short for term frequency-inverse document frequency, is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus. The tf-idf value increases proportionally to the number of times a word appears in the document and is offset by the frequency of the word in the corpus, which helps to adjust for the fact that some words appear more frequently in general. It is often used as a weighting factor in searches of information retrieval, text mining, and user modeling. Tf-idf is one of the most popular term-weighting schemes today; 83% of text-based recommender systems in digital libraries use tf-idf.

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right)$$

$tf_{i,j}$ = number of occurrences of i in j
 df_i = number of documents containing i
 N = total number of documents

```
In [10]: %%time
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf = TfidfVectorizer(ngram_range=(1,2)) #Using bi-grams
tfidf_vec = tfidf.fit_transform(df2['CleanedText'])
```

CPU times: user 1min, sys: 788 ms, total: 1min
 Wall time: 1min

```
In [11]: #Saving the variable to access later without recomputing
savetofile(tfidf_vec, "tfidf")
```

```
In [12]: #Loading the variable from file
tfidf_vec = openfromfile("tfidf")
```

```
In [13]: tfidf_vec.shape
```

```
Out[13]: (364171, 3404647)
```

tf-idf came up with 2.9 million features for the data corpus

In [14]: `print(tfidf_vec[2])`

```
(0, 1995684) 0.0230577620001
(0, 634637) 0.0962976136656
(0, 148485) 0.0491621357676
(0, 483299) 0.095130835222
(0, 1680739) 0.0511254495732
(0, 2212861) 0.131631500755
(0, 556420) 0.0781934548627
(0, 1232676) 0.0865863526019
(0, 2018759) 0.0553755411213
(0, 465022) 0.0506518591503
(0, 1093493) 0.119601669622
(0, 737841) 0.0569454130211
(0, 3064216) 0.128558300715
(0, 2798119) 0.0741999555212
(0, 1673827) 0.0960763467562
(0, 575573) 0.0613536698434
(0, 2271789) 0.0542982037779
(0, 2885378) 0.0416697262506
(0, 1920053) 0.0562561162586
(0, 1397444) 0.070424363201
(0, 516953) 0.0606048623527
(0, 1125576) 0.027927961642
(0, 1414824) 0.040096342878
(0, 2423531) 0.0377442598615
(0, 3396938) 0.0571169892128
:
:
(0, 576437) 0.132967109652
(0, 2273524) 0.0936229482509
(0, 2888619) 0.126621307635
(0, 3064745) 0.124419662051
(0, 1920733) 0.129355981494
(0, 1397906) 0.117978967687
(0, 1997120) 0.09759836373
(0, 517278) 0.108210127217
(0, 1128708) 0.0911014244865
(0, 1416273) 0.0498816191539
(0, 2427136) 0.125107119153
(0, 3398129) 0.102164881838
(0, 3112449) 0.144089290953
(0, 1052035) 0.144089290953
(0, 2851106) 0.148687090651
(0, 1673525) 0.148687090651
(0, 1703206) 0.132967109652
(0, 3330091) 0.134481298135
(0, 3256500) 0.148687090651
(0, 3114021) 0.148687090651
(0, 2595216) 0.148687090651
(0, 933619) 0.148687090651
(0, 2607918) 0.148687090651
(0, 379436) 0.113896691731
(0, 2684850) 0.148687090651
```

Returns all the features which is non-zero for a particular review from the sparse matrix

```
In [15]: features = tfidf.get_feature_names()  
features[190000:190010]
```

```
Out[15]: ['babi health',  
          'babi healthi',  
          'babi healthier',  
          'babi healthiest',  
          'babi healthyp',  
          'babi healtyplus',  
          'babi heart',  
          'babi heat',  
          'babi heimlich',  
          'babi help']
```

Some of the feature of the tf-idf

```
In [16]: def top_tfidf_features(row, features, top_n=25):
    ''' Get top n tfidf values in row and return them with their corresponding fe
topn_ind = np.argsort(row)[::-1][:top_n]
#Sorting and getting the indexes using argsort and reversing to get descending
top_feats = [(features[i], row[i]) for i in topn_ind]
df = pd.DataFrame(top_feats,columns = ['feature', 'tfidf'])
return df
top_tfidfs = top_tfidf_features(tfidf_vec[3000,:].toarray()[0],features,20)#top 2
top_tfidfs
```

Out[16]:

	feature	tfidf
0	stretch strong	0.399925
1	cup stretch	0.378784
2	delici mocha	0.357643
3	delici hot	0.268651
4	chocol best	0.262176
5	best cup	0.256560
6	coffe delici	0.251700
7	stretch	0.229628
8	delici	0.220825
9	strong coffe	0.208862
10	mocha	0.199555
11	hot chocol	0.192724
12	strong	0.128024
13	hot	0.122785
14	chocol	0.115835
15	cup	0.114143
16	coffe	0.097874
17	best	0.094951
18	flavorless dark	0.000000
19	flavorless cup	0.000000

Top 20 tfidf features of 3000th review in the data corpus

```
In [ ]: %time
from sklearn.decomposition import TruncatedSVD

tsvd_tfidf = TruncatedSVD(n_components=100)#No of components as total dimensions
tsvd_tfidf_vec = tsvd_tfidf.fit_transform(tfidf_vec)
```

```
In [ ]: savetofile(tsvd_tfidf,"tsvd_tfidf")
savetofile(tsvd_tfidf_vec,"tsvd_tfidf_vec")
```

```
In [27]: tsvd_tfidf_vec = openfromfile("tsvd_tfidf_vec")
tsvd_tfidf = openfromfile("tsvd_tfidf")
```

```
In [41]: tsvd_tfidf.explained_variance_ratio_[:].sum()
```

```
Out[41]: 0.0030303842146799662
```

In [43]:

```

%%time
#Perplexity = 20 with 10k points
from sklearn.manifold import TSNE
from time import time
import random

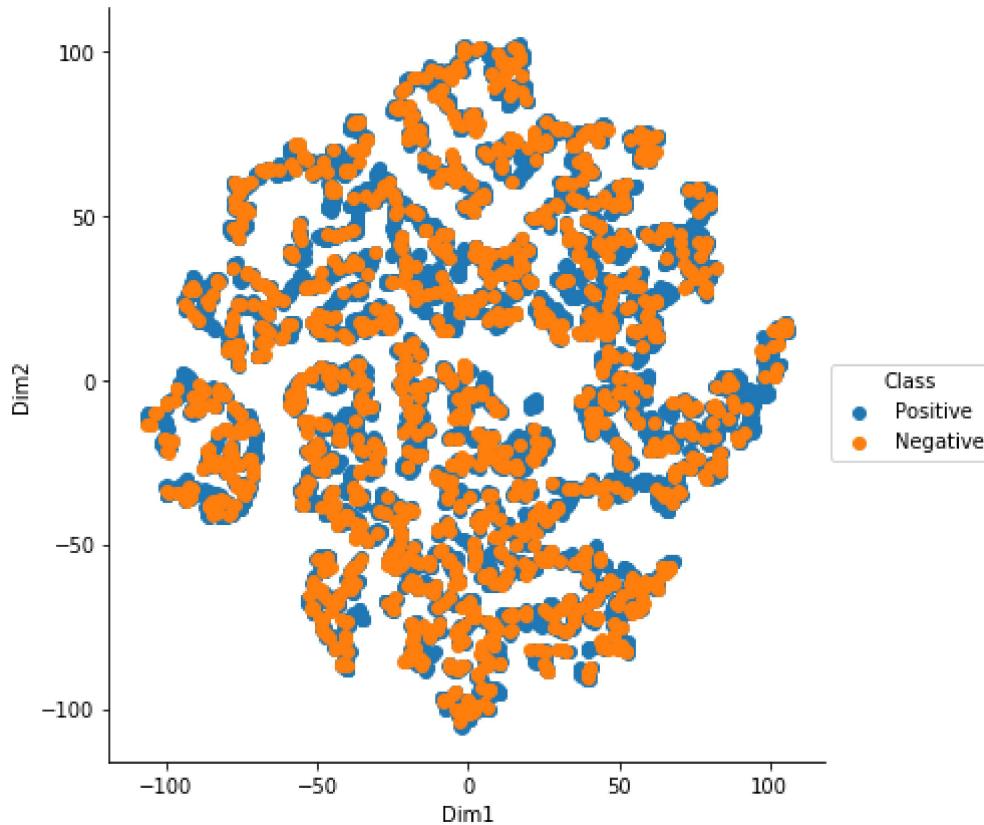
n_samples = 10000
sample_cols = random.sample(range(1, tsvd_tfidf_vec.shape[0]), n_samples)
sample_features = tsvd_tfidf_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=20)
# print(sample_features,sample_class)

t0 = time()
embedded_data = model.fit_transform(sample_features)
print("TSNE done in %0.3fs." % (time() - t0))
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])

sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend()
plt.show()

```

(10000, 2) (10000, 1)
 TSNE done in 248.691s.
 (10000, 3)



```
CPU times: user 3min 37s, sys: 32 s, total: 4min 9s
Wall time: 4min 9s
```

In [44]:

```

%%time
#Perplexity = 30 with 10k points
from sklearn.manifold import TSNE
from time import time
import random

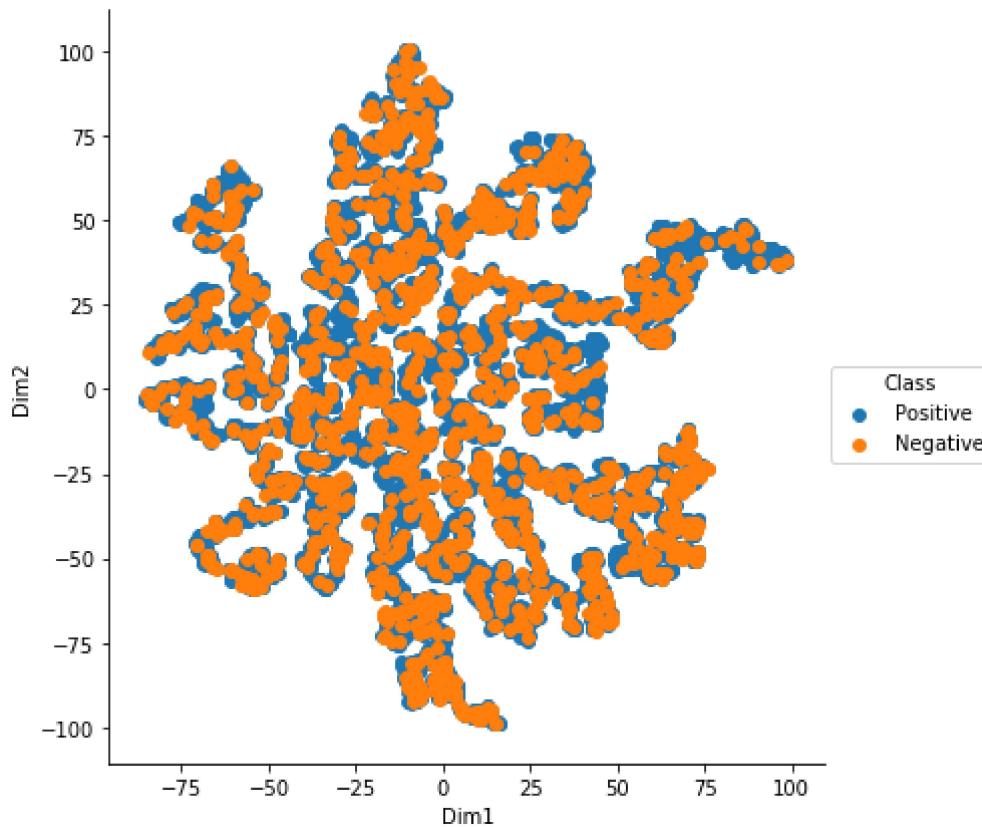
n_samples = 10000
sample_cols = random.sample(range(1, tsvd_tfidf_vec.shape[0]), n_samples)
sample_features = tsvd_tfidf_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=30)
# print(sample_features,sample_class)

t0 = time()
embedded_data = model.fit_transform(sample_features)
print("TSNE done in %0.3fs." % (time() - t0))
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])

sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend()
plt.show()

```

(10000, 2) (10000, 1)
 TSNE done in 283.880s.
 (10000, 3)



```
CPU times: user 4min 13s, sys: 31.2 s, total: 4min 44s
Wall time: 4min 44s
```



In [45]:

```

%%time
#Perplexity = 40 with 10k points
from sklearn.manifold import TSNE
from time import time
import random

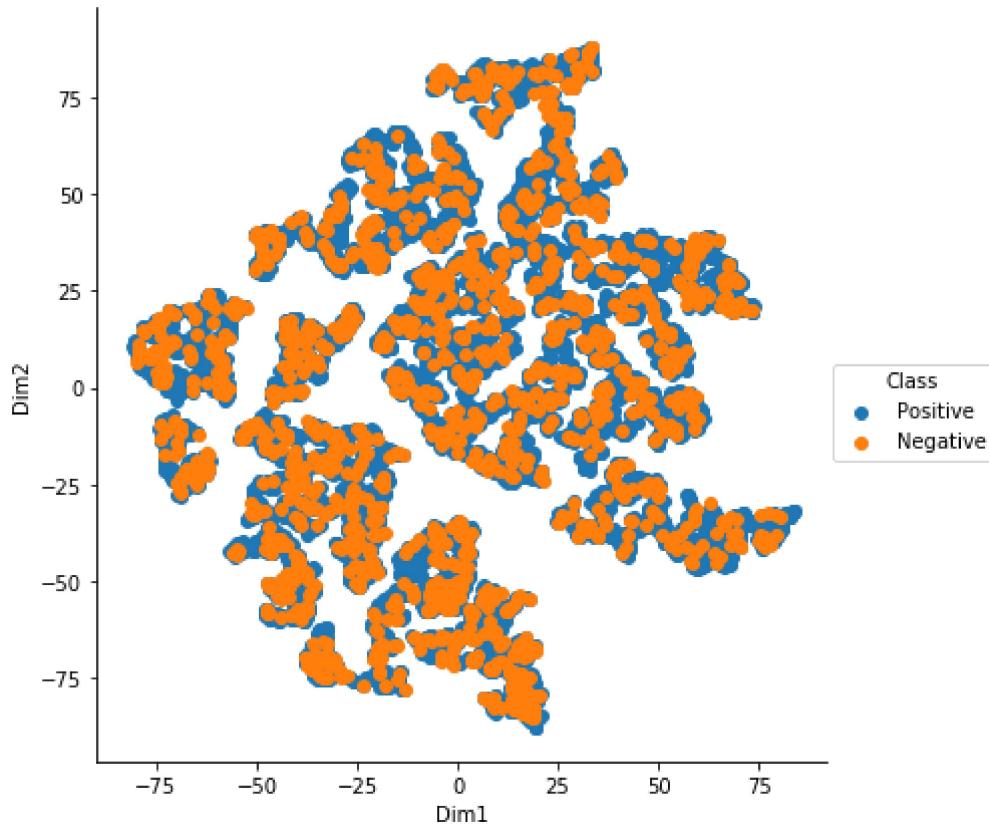
n_samples = 10000
sample_cols = random.sample(range(1, tsvd_tfidf_vec.shape[0]), n_samples)
sample_features = tsvd_tfidf_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=40)
# print(sample_features,sample_class)

t0 = time()
embedded_data = model.fit_transform(sample_features)
print("TSNE done in %0.3fs." % (time() - t0))
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])

sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend()
plt.show()

```

(10000, 2) (10000, 1)
 TSNE done in 312.785s.
 (10000, 3)



CPU times: user 4min 41s, sys: 32.2 s, total: 5min 13s
Wall time: 5min 13s

Gensim

Gensim is a robust open-source vector space modeling and topic modeling toolkit implemented in Python. It uses NumPy, SciPy and optionally Cython for performance. Gensim is specifically designed to handle large text collections, using data streaming and efficient incremental algorithms, which differentiates it from most other scientific software packages that only target batch and in-memory processing.

4. Word2Vec

[Refer Docs] :<https://radimrehurek.com/gensim/models/word2vec.html>
<https://radimrehurek.com/gensim/models/word2vec.html>

```
In [92]: final_string = []
for sent in df2['CleanedText'].values:
    sent = str(sent)
    sentence=[]
    #    print(sent)
    for word in sent.split():
        #        print(word)
        sentence.append(word)
    #        print(sentence)
    final_string.append(sentence)
```

```
In [93]: %time
# Train your own Word2Vec model using your own text corpus
import gensim

w2v_model=gensim.models.Word2Vec(final_string,min_count=5,size=50, workers=-1)
#min-count: Ignoring the words which occurs less than 5 times
#size:Creating vectors of size 50 for each word
#workers: Use these many worker threads to train the model (faster training with i
```

CPU times: user 5.28 s, sys: 0 ns, total: 5.28 s
Wall time: 5.28 s

```
In [30]: w2v_model.save('w2vmodel')#Persist/Saving the model to a file in the disk
```

```
In [31]: w2v_model = gensim.models.Word2Vec.load('w2vmodel') #Loading the model from file
```

```
In [32]: w2v_vocab = w2v_model.wv.vocab
len(w2v_vocab)
```

Out[32]: 34906

```
In [33]: w2v_model.wv.most_similar('like')
```

```
Out[33]: [('compar', 0.5214215517044067),  
 ('anywaya', 0.5147097706794739),  
 ('dryer', 0.5107567310333252),  
 ('hadiv', 0.5099674463272095),  
 ('vivani', 0.49352583289146423),  
 ('vancouv', 0.4796876013278961),  
 ('vomit', 0.47908806800842285),  
 ('mirin', 0.47721731662750244),  
 ("peet'", 0.47637227177619934),  
 ('downthi', 0.4757559895515442)]
```

```
In [34]: w2v_model.wv.most_similar('tast')
```

```
Out[34]: [('porch', 0.5215961933135986),  
 ('maisi', 0.5133664608001709),  
 ('unstick', 0.49281394481658936),  
 ('complianc', 0.48115843534469604),  
 ("b'pricey", 0.47690099477767944),  
 ('mightili', 0.47623634338378906),  
 ("b'also", 0.47145384550094604),  
 ('vow', 0.4705711007118225),  
 ('hexan', 0.46026793122291565),  
 ('ment', 0.45899584889411926)]
```

```
In [35]: w2v_model.wv.most_similar('good')
```

```
Out[35]: [('cain', 0.5610529184341431),  
 ('finea', 0.5418595671653748),  
 ('therapi', 0.5118728280067444),  
 ('sasha', 0.5085721015930176),  
 ('anton', 0.5078103542327881),  
 ("soil'", 0.5062910914421082),  
 ('discern', 0.5059114694595337),  
 ("awsom'", 0.5021228790283203),  
 ("b'newborn", 0.49307799339294434),  
 ("larg'", 0.4898505210876465)]
```

4.a Avg Word2Vec

- One of the most naive but good ways to convert a sentence into a vector
- Convert all the words to vectors and then just take the avg of the vectors the resulting vector represent the sentence

In [62]:

```
%time
avg_vec = [] #List to store all the avg w2vec's
for sent in final_string[0:1]:
    cnt = 0 #to count no of words in each reviews
    sent_vec = np.zeros(50) #Initializing with zeroes
    print("sent:",sent)
    for word in sent:
        try:
            wvec = w2v_model.wv[word] #Vector of each using w2v model
            print("wvec:",wvec)
            sent_vec += wvec #Adding the vectors
            cnt += 1
        except:
            pass #When the word is not in the dictionary then do nothing
    print("sent_vec:",sent_vec)
    a_vec =sent_vec / cnt #Taking average of vectors sum of the particular review
    print("avg_vec:",a_vec)
    avg_vec.append(a_vec) #Storing the avg w2vec's for each review
print("*****")
```

sent: ["b'bought", 'sever', 'vital', 'can', 'dog', 'food', 'product', 'foun
d', 'good', 'qualiti', 'product', 'look', 'like', 'stew', 'process', 'meat',
'smell', 'better', 'labrador', 'finicki', 'appreci', 'product', "better'"]
wvec: [-4.75264713e-03 -2.68976251e-03 3.45981750e-03 8.61182716e-03
-3.45326052e-03 -6.07945665e-04 -3.81294428e-03 7.70723587e-03
-6.99552661e-03 -2.04596203e-03 -4.63803299e-03 -7.98908077e-05
-5.44417766e-04 -6.50190702e-03 4.06994065e-03 3.06597492e-03
-7.19741359e-03 -8.84887390e-03 -8.55807401e-03 5.74907893e-03
6.90606283e-03 3.70534114e-03 4.26333630e-03 -9.67295747e-03
-5.33873122e-03 2.21293815e-03 5.71956998e-03 4.67022089e-03
-5.25551289e-03 7.90084433e-03 -8.64620041e-03 1.24186510e-03
6.92145852e-03 7.15341698e-03 1.75182312e-03 2.55550840e-03
6.01556059e-03 1.02293317e-03 -7.80893781e-04 7.74320262e-03
-5.13905776e-04 -1.03073404e-03 3.13923252e-03 1.53249697e-04
2.41127494e-03 1.59304694e-03 -1.53104786e-03 -1.76926993e-03
1.10834360e-03 -3.19925486e-03]
wvec: [-0.00556279 0.00752297 0.00460804 -0.00926204 0.0089932 -0.0055966
5
0.00586318 0.00810189 -0.00954244 0.00968478 -0.00038469 0.00798686
^ ^ ^ ^ ^ ^

```
In [63]: %%time
np.seterr(divide='ignore', invalid='ignore')
avg_vec = [] #List to store all the avg w2vec's
for sent in final_string:
    cnt = 0 #to count no of words in each reviews
    sent_vec = np.zeros(50) #Initializing with zeroes
    for word in sent:
        try:
            wvec = w2v_model.wv[word] #Vector of each using w2v model
            sent_vec += wvec #Adding the vectors
            cnt += 1
        except:
            pass #When the word is not in the dictionary then do nothing
    sent_vec /= cnt #Taking average of vectors sum of the particular review
    avg_vec.append(sent_vec) #Storing the avg w2vec's for each review
#print("*****")
# Average Word2Vec
```

CPU times: user 1min 26s, sys: 0 ns, total: 1min 26s
Wall time: 1min 26s

```
In [ ]: #Saving the variable to access later without recomputing
savetofile(avg_vec, "avg_w2v_vec")
```

```
In [4]: #Loading the variable from file
avg_vec = openfromfile("avg_w2v_vec")
```

```
In [40]: avg_vec = np.array(avg_vec)
avg_vec.shape
```

Out[40]: (364171, 50)

```
In [33]: from sklearn import preprocessing
avg_vec_norm = preprocessing.normalize(avg_vec)
```

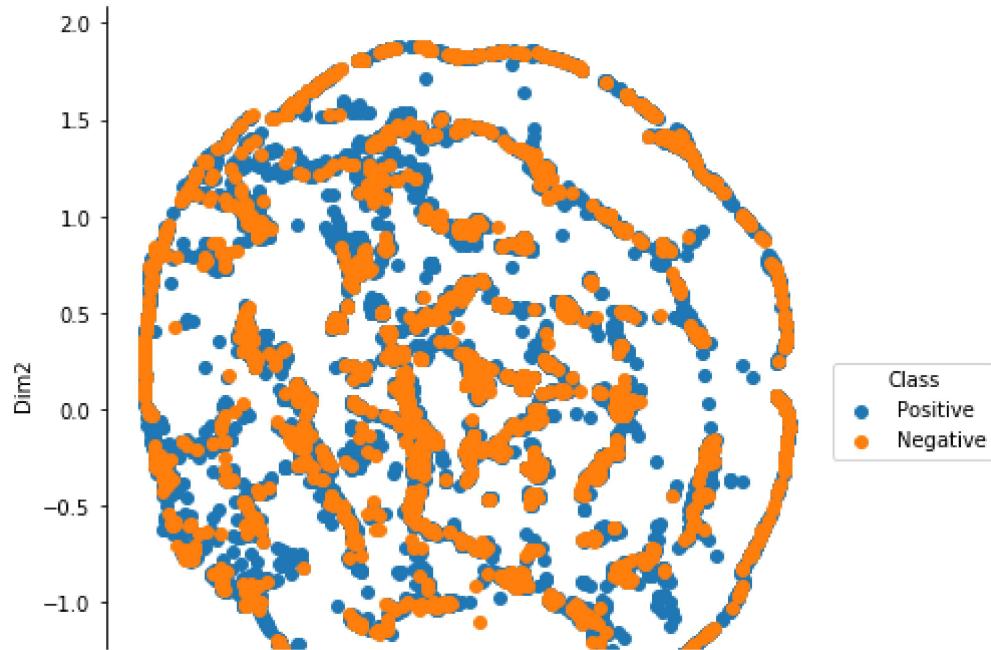
```
In [21]: %%time
from sklearn.manifold import TSNE
from time import time
import random

n_samples = 20000
sample_cols = random.sample(range(1, avg_vec.shape[0]), n_samples)
sample_features = avg_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=30)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
```

(20000, 50) (20000, 1)
 TSNE done in 767.050s.
 (20000, 3)
 CPU times: user 11min 46s, sys: 1min, total: 12min 47s
 Wall time: 12min 47s

```
In [22]: sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```



In [30]:

```
%%time
from sklearn.manifold import TSNE
from time import time
import random

n_samples = 40000
sample_cols = random.sample(range(1, avg_vec.shape[0]), n_samples)
sample_features = avg_vec[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=30)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])
```

(40000, 50) (40000, 1)

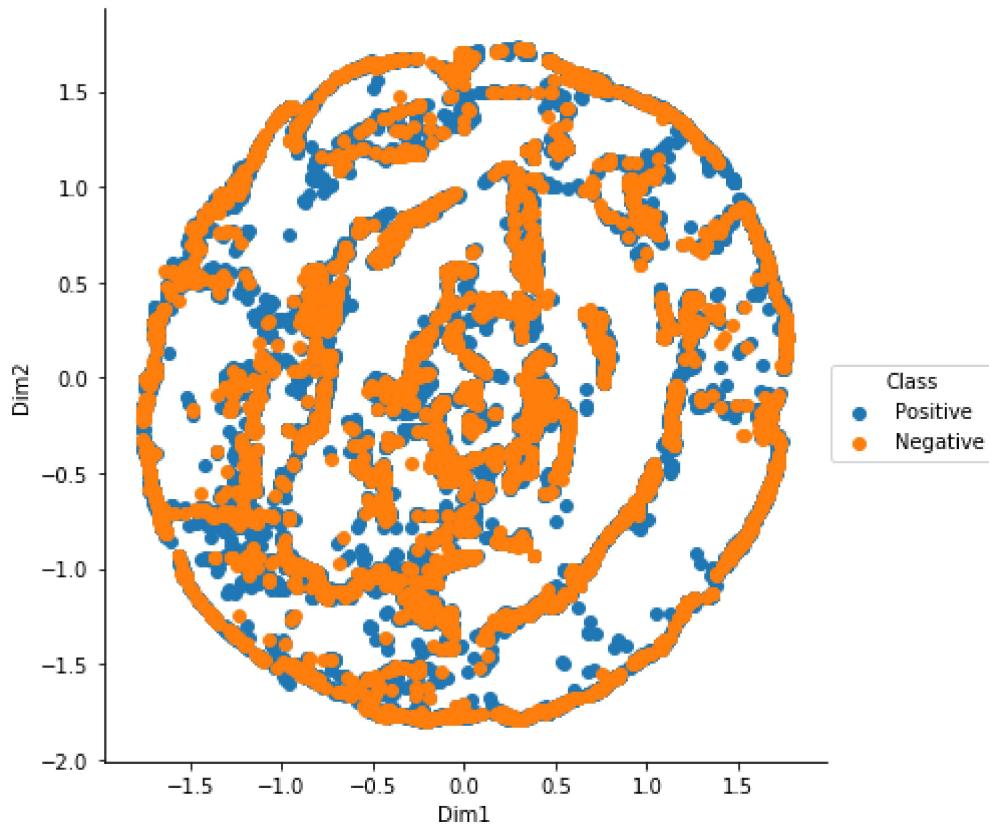
(40000, 3)

CPU times: user 26min 9s, sys: 1min 50s, total: 27min 59s

Wall time: 28min

In [31]:

```
sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend()
plt.show()
```



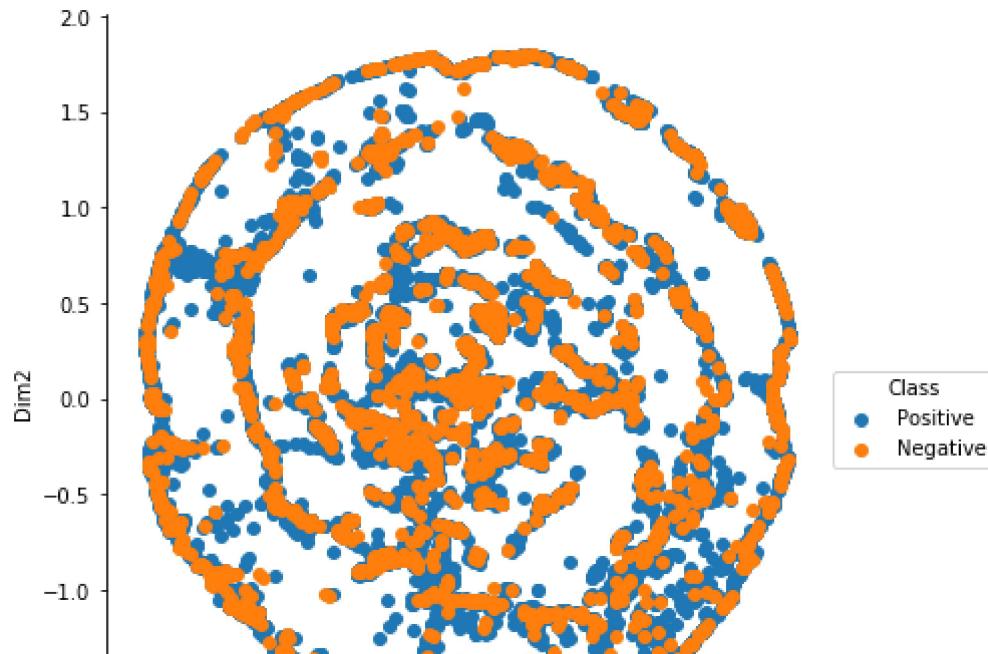
```
In [38]: %%time
from sklearn.manifold import TSNE
import random

n_samples = 20000
sample_cols = random.sample(range(1, avg_vec_norm.shape[0]), n_samples)
sample_features = avg_vec_norm[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=20)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
```

(20000, 50) (20000, 1)
(20000, 3)
CPU times: user 8min 32s, sys: 47.3 s, total: 9min 19s
Wall time: 9min 19s

```
In [39]: sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```



4.b Using Google's Trained W2Vec on Google News

```
In [3]: from gensim.models import KeyedVectors
w2v_model_google = KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative')
```

```
In [4]: w2v_vocab = w2v_model_google.wv.vocab  
len(w2v_vocab)
```

Out[4]: 3000000

```
In [5]: w2v_model_google.wv.most_similar('like')
```

```
Out[5]: [('really', 0.5752447843551636),  
 ('weird', 0.5676319599151611),  
 ('crazy', 0.5382447838783264),  
 ('kind', 0.5310239195823669),  
 ('maybe', 0.5220045447349548),  
 ('looooove', 0.5187614560127258),  
 ('anymore', 0.5177680253982544),  
 ('Kinda_reminds', 0.5151872634887695),  
 ('definitely', 0.5117843151092529),  
 ('kinda_fishy', 0.5090124607086182)]
```

```
In [39]: w2v_model_google.wv.most_similar('taste')
```

```
Out[39]: [('tastes', 0.6838272213935852),  
 ('flavor', 0.6630197763442993),  
 ('tasted', 0.6162090301513672),  
 ('Harry_Potter_butterbeer', 0.5894586443901062),  
 ('tasting', 0.5604724884033203),  
 ('tangy_taste', 0.5567916035652161),  
 ('aftertaste', 0.5558385252952576),  
 ('bitter_taste', 0.5491952300071716),  
 ('carbonated_cough_syrup', 0.5455324053764343),  
 ('taste_buds', 0.5368086695671082)]
```

```
In [50]: w2v_model_google.wv["word"].size
```

Out[50]: 300

In [107]:

```
%time
avg_vec_google = [] #List to store all the avg w2vec's
no_datapoints = 364170
sample_cols = random.sample(range(1, no_datapoints), 20001)
for sent in df2['CleanedText_NoStem'].values[sample_cols]:
    cnt = 0 #to count no of words in each reviews
    sent_vec = np.zeros(300) #Initializing with zeroes
    # print("sent:",sent)
    sent = sent.decode("utf-8")
    for word in sent.split():
        try:
            # print(word)
            wvec = w2v_model_google.wv[word] #Vector of each using w2v model
            # print("wvec:",wvec)
            sent_vec += wvec #Adding the vectors
            # print("sent_vec:",sent_vec)
            cnt += 1
        except:
            pass #When the word is not in the dictionary then do nothing
        # print(sent_vec)
    sent_vec /= cnt #Taking average of vectors sum of the particular review
    # print("avg_vec:",sent_vec)
    avg_vec_google.append(sent_vec) #Storing the avg w2vec's for each review
    # print("*****")
# print(avg_vec_google)
avg_vec_google = np.array(avg_vec_google)
```

CPU times: user 9.89 s, sys: 4 ms, total: 9.89 s
Wall time: 9.89 s

In [108]:

```
#Saving the variable to access later without recomputing
savetofile(avg_vec_google,"avg_w2v_vec_google")
```

In [109]:

```
#Loading the variable from file
avg_vec_google = openfromfile("avg_w2v_vec_google")
```

In [110]:

```
from sklearn import preprocessing

avg_vec_google_norm = preprocessing.normalize(avg_vec_google)
```

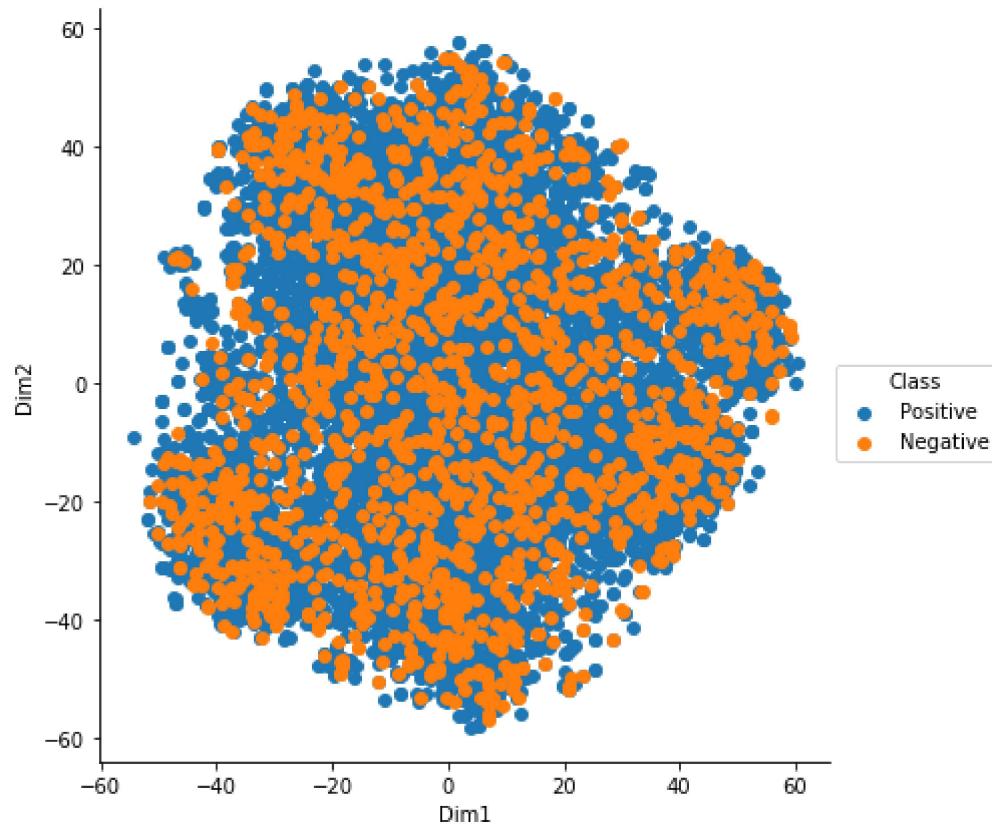
In [115]:

```
%%time
from sklearn.manifold import TSNE
import random

n_samples = 10000
sample_cols = random.sample(range(1, avg_vec_google.shape[0]), n_samples)
sample_features = avg_vec_google[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=20)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```

(10000, 300) (10000, 1)
(10000, 3)



CPU times: user 6min 32s, sys: 31.6 s, total: 7min 4s
Wall time: 7min 3s

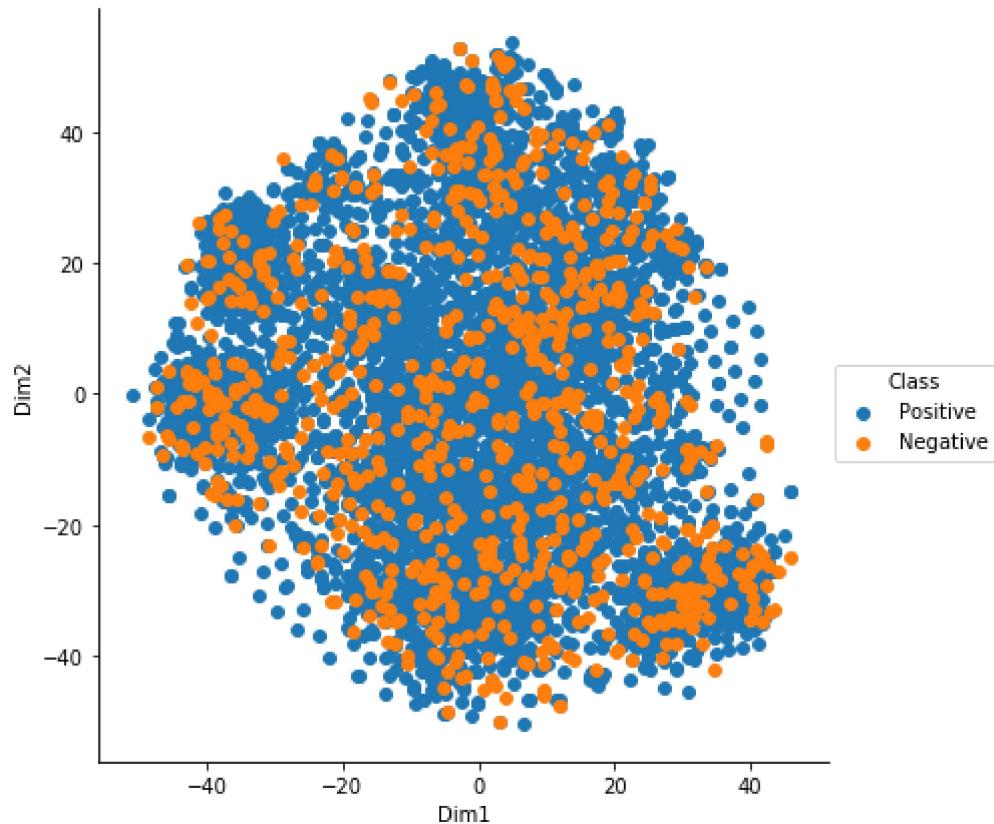
In [121]:

```
%time
from sklearn.manifold import TSNE
import random

n_samples = 5000
sample_cols = random.sample(range(1, avg_vec_google.shape[0]), n_samples)
sample_features = avg_vec_google[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=20)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])
sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend()
plt.show()
```

(5000, 300) (5000, 1)
(5000, 3)



CPU times: user 2min 43s, sys: 16.4 s, total: 2min 59s
Wall time: 2min 59s

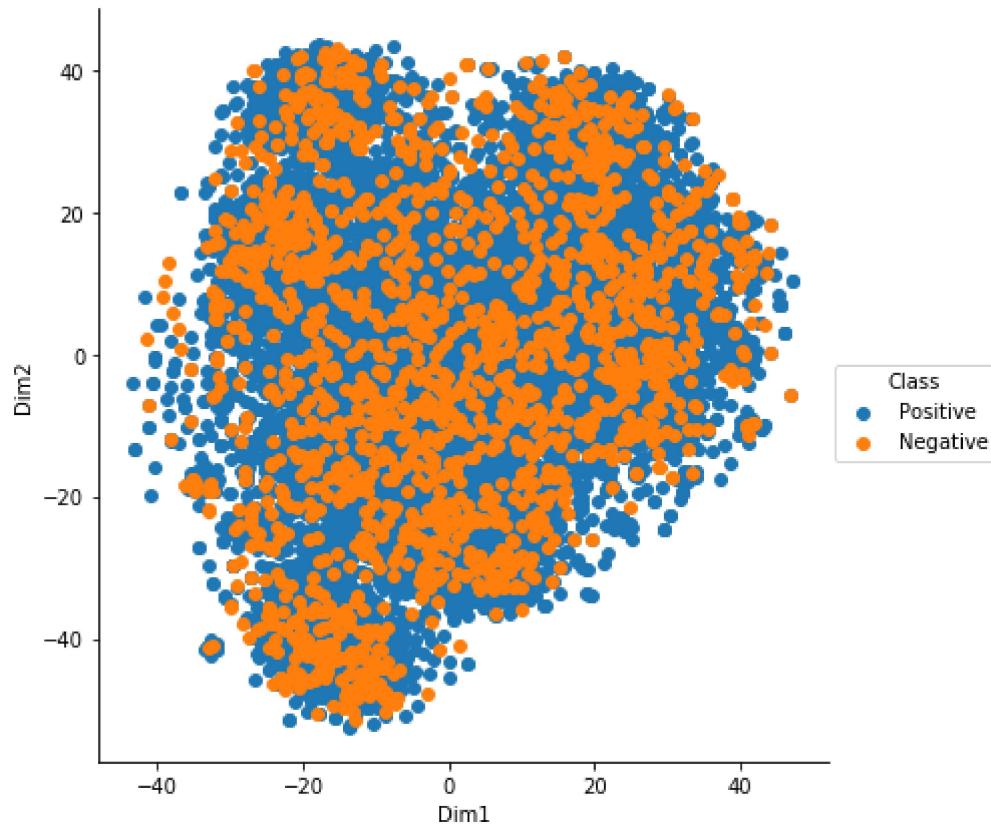
In [117]:

```
%time
from sklearn.manifold import TSNE
import random

n_samples = 10000
sample_cols = random.sample(range(1, avg_vec_google.shape[0]), n_samples)
sample_features = avg_vec_google[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=30)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```

(10000, 300) (10000, 1)
(10000, 3)



CPU times: user 7min 14s, sys: 33 s, total: 7min 47s
Wall time: 7min 47s

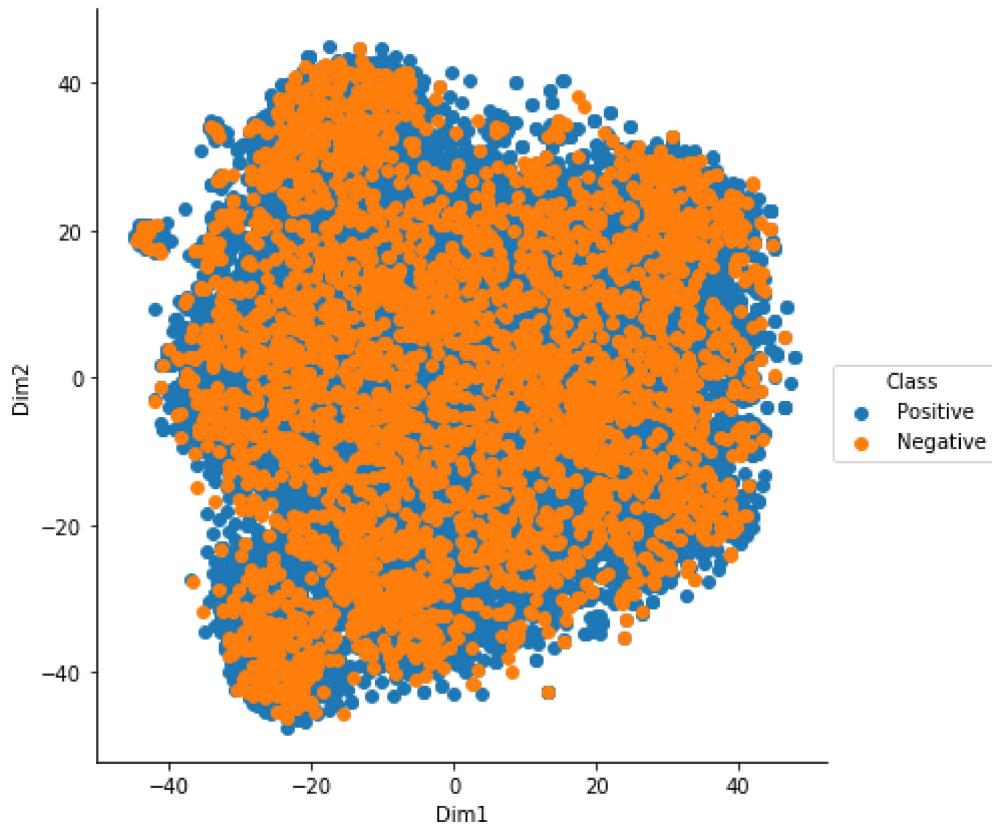
In [118]:

```
%time
from sklearn.manifold import TSNE
import random

n_samples = 20000
sample_cols = random.sample(range(1, avg_vec_google.shape[0]), n_samples)
sample_features = avg_vec_google[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=30)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```

(20000, 300) (20000, 1)
(20000, 3)



CPU times: user 27min 48s, sys: 1min 7s, total: 28min 55s
Wall time: 28min 55s

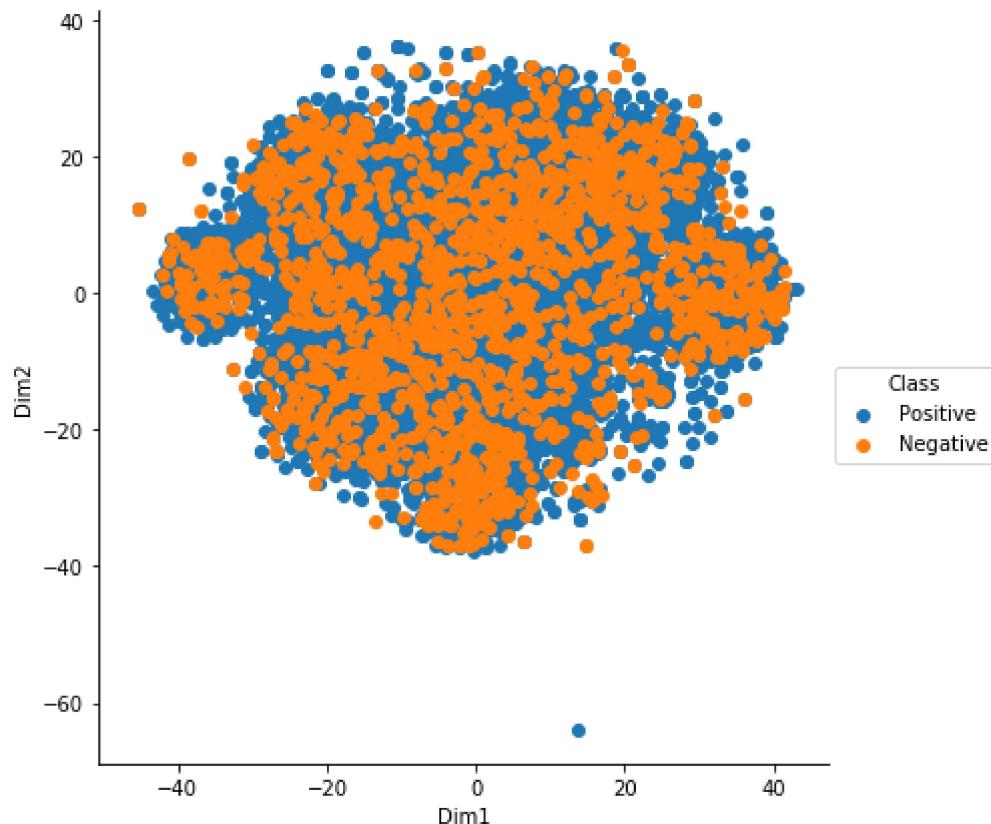
In [119]:

```
%%time
from sklearn.manifold import TSNE
import random

n_samples = 10000
sample_cols = random.sample(range(1, avg_vec_google.shape[0]), n_samples)
sample_features = avg_vec_google[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:,np.newaxis]
print(sample_features.shape,sample_class.shape)
model = TSNE(n_components=2,random_state=0,perplexity=35)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape,sample_class.shape)
final_data = np.concatenate((embedded_data,sample_class),axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data,columns=["Dim1","Dim2","Class"])
sns.FacetGrid(newdf,hue="Class",size=6).map(plt.scatter,"Dim1","Dim2").add_legend()
plt.show()
```

(10000, 300) (10000, 1)
(10000, 3)



CPU times: user 7min 44s, sys: 32.9 s, total: 8min 16s
Wall time: 8min 16s

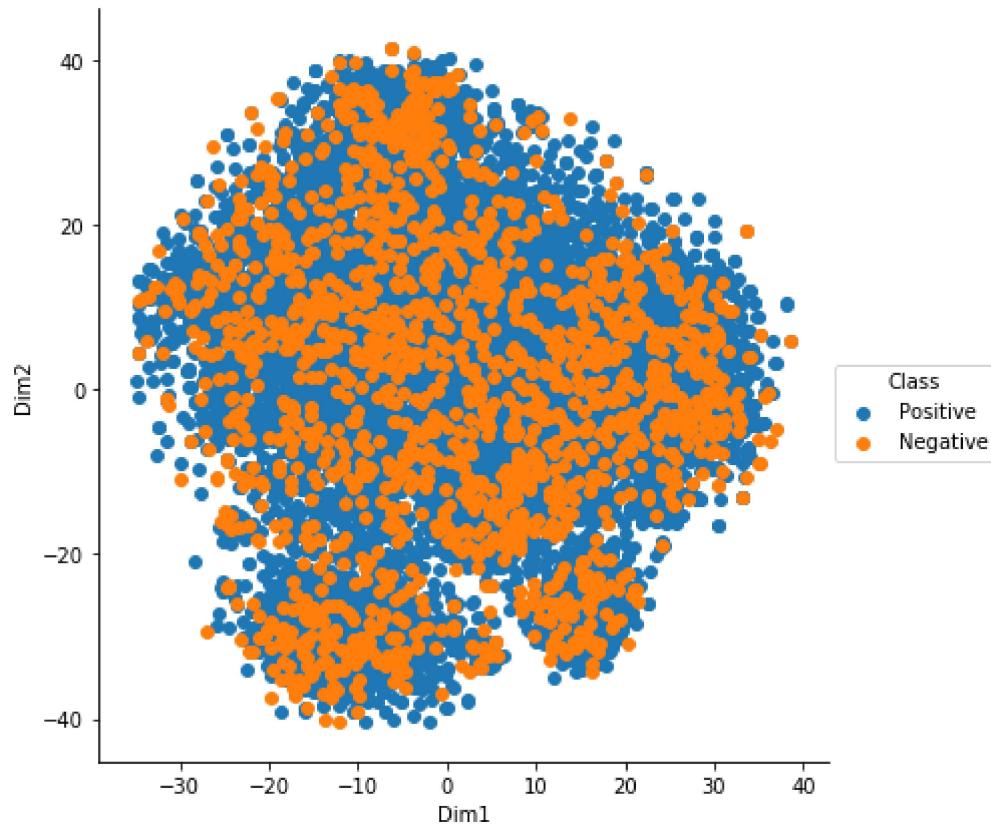
In [120]:

```
%%time
from sklearn.manifold import TSNE
import random

n_samples = 10000
sample_cols = random.sample(range(1, avg_vec_google.shape[0]), n_samples)
sample_features = avg_vec_google[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=40)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```

(10000, 300) (10000, 1)
(10000, 3)



CPU times: user 8min 13s, sys: 32 s, total: 8min 45s
Wall time: 8min 44s

5. Tf-idf W2Vec

- Another way to convert sentence into vectors

- Take weighted sum of the vectors divided by the sum of all the tfidf's
i.e. $(\text{tfidf}(\text{word}) \times \text{w2v}(\text{word})) / \sum(\text{tfidf}'s)$

In [19]: `#Taking Sample of 20k points
no_datapoints = 364170
sample_cols = random.sample(range(1, no_datapoints), 20001)`

In [20]: `%%time
###tf-idf with No Stemming
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf = TfidfVectorizer(ngram_range=(1,2)) #Using bi-grams
tfidf_vec_ns = tfidf.fit_transform(df2['CleanedText_NoStem'].values[sample_cols])

#Saving the variable to access later without recomputing
savetofile(tfidf_vec, "tfidf")

#Loading the variable from file
tfidf_vec = openfromfile("tfidf")

print(tfidf_vec_ns.shape)

tf-idf came up with 2.9 million features for the data corpus
from sklearn.decomposition import TruncatedSVD

tsvd_tfidf_ns = TruncatedSVD(n_components=300)#No of components as total dimension
tsvd_tfidf_vec_ns = tsvd_tfidf_ns.fit_transform(tfidf_vec_ns)
print(tsvd_tfidf_ns.explained_variance_ratio_[:].sum())
features = tfidf.get_feature_names()

(20001, 492431)
0.110114446613
CPU times: user 3min 56s, sys: 5.88 s, total: 4min 2s
Wall time: 58 s`

```
In [21]: %%time
tfidf_w2v_vec_google = []
review = 0

for sent in df2['CleanedText_NoStem'].values[sample_cols]:
    cnt = 0
    weighted_sum = 0
    sent_vec = np.zeros(300)
    sent = sent.decode("utf-8")
    for word in sent.split():
        try:
#            print(word)
#            wvec = w2v_model_google.wv[word] #Vector of each using w2v model
#            print("w2vec:",wvec)
#            print("tfidf:",tfidf_vec_ns[review,features.index(word)])
#            tfidf = tfidf_vec_ns[review,features.index(word)]
#            print(tfidf)
            sent_vec += (wvec * tfidf)
            weighted_sum += tfidf
        except:
            pass
    sent_vec /= weighted_sum
    tfidf_w2v_vec_google.append(sent_vec)
    review += 1
```

CPU times: user 3h 20min 51s, sys: 1.32 s, total: 3h 20min 52s
Wall time: 3h 20min 53s

```
In [22]: len(tfidf_w2v_vec_google)
```

Out[22]: 20001

```
In [23]: len(tfidf_w2v_vec_google[0])
```

Out[23]: 300

```
In [24]: tfidf_w2v_vec_google[5]
```

```
Out[24]: array([ 6.30765966e-03,  2.62348772e-02, -1.28013094e-02,
   1.66244870e-01, -4.49297504e-02,  2.66082968e-02,
   8.06051421e-02, -8.27518457e-02,  7.73820410e-03,
  1.28079768e-01,  4.09113582e-03, -1.26601291e-01,
 -3.10158627e-02,  3.02259649e-02, -1.17419000e-01,
 1.19921784e-01,  3.36843358e-02,  1.50723015e-01,
 4.26849213e-02, -3.13014550e-02, -7.91879333e-02,
 5.21383487e-02, -1.22338065e-02, -1.99465251e-03,
 7.86969992e-02, -8.04239890e-03, -4.38397773e-02,
 5.46432782e-02, -7.09601715e-03,  4.12984907e-02,
 -6.55848419e-02,  1.60784459e-03,  2.39487639e-02,
 2.04896547e-02,  3.50686609e-02,  4.94815506e-04,
 4.68977209e-02, -8.77567649e-02, -1.13872285e-02,
 1.18984474e-01,  9.81944115e-02, -9.27291092e-02,
 1.36931440e-01, -1.44343861e-02, -2.90870869e-02,
 -9.79914776e-02, -3.98462383e-02, -2.20645862e-02,
 6.52701948e-03, -2.17011543e-03, -4.47150526e-02,
 8.24235868e-03,  4.96269734e-03, -1.31958030e-02,
 -2.56885586e-04,  2.93922949e-02, -1.72273889e-02,
```

```
In [25]: savetofile(tfidf_w2v_vec_google,"tfidf_w2v_vec_google")
```

```
In [3]: tfidf_w2v_vec_google = openfromfile("tfidf_w2v_vec_google")
```

```
In [4]: from sklearn import preprocessing
tfidf_w2v_vec_google_norm = preprocessing.normalize(tfidf_w2v_vec_google)
```

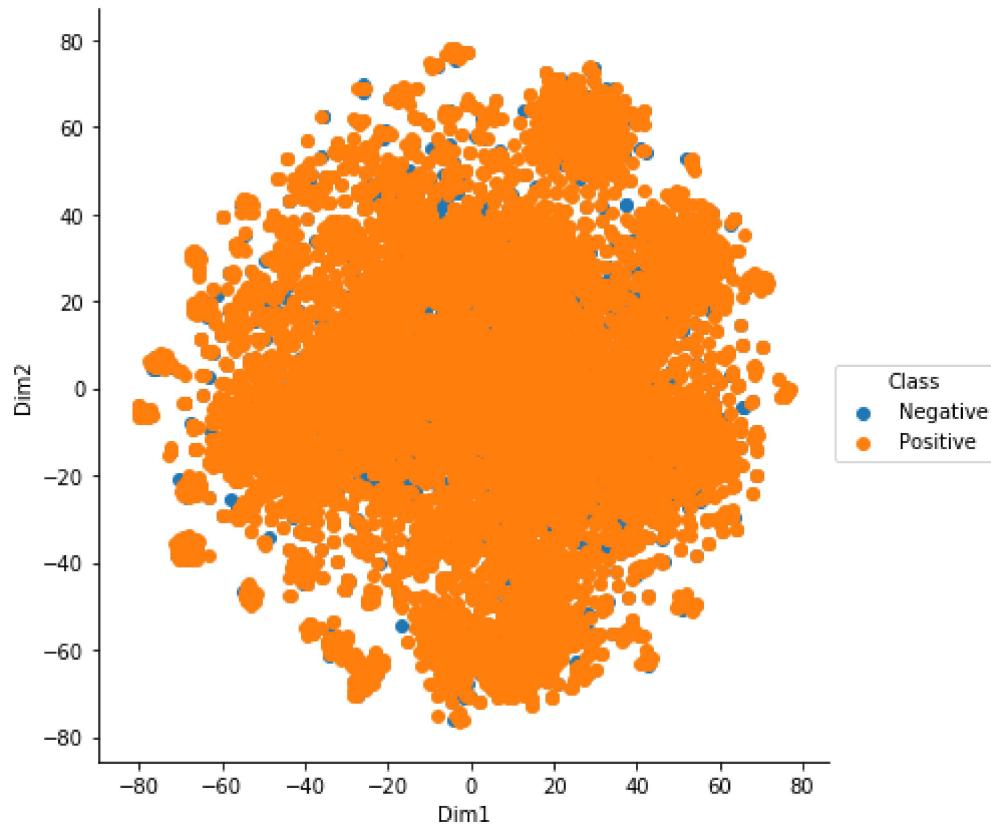
In [28]:

```
%%time
from sklearn.manifold import TSNE
import random

n_samples = 10000
sample_cols = random.sample(range(1, tfidf_w2v_vec_google_norm.shape[0]), n_samples)
sample_features = tfidf_w2v_vec_google_norm[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=20)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```

(10000, 300) (10000, 1)
(10000, 3)



CPU times: user 6min 7s, sys: 32.6 s, total: 6min 40s
Wall time: 6min 40s

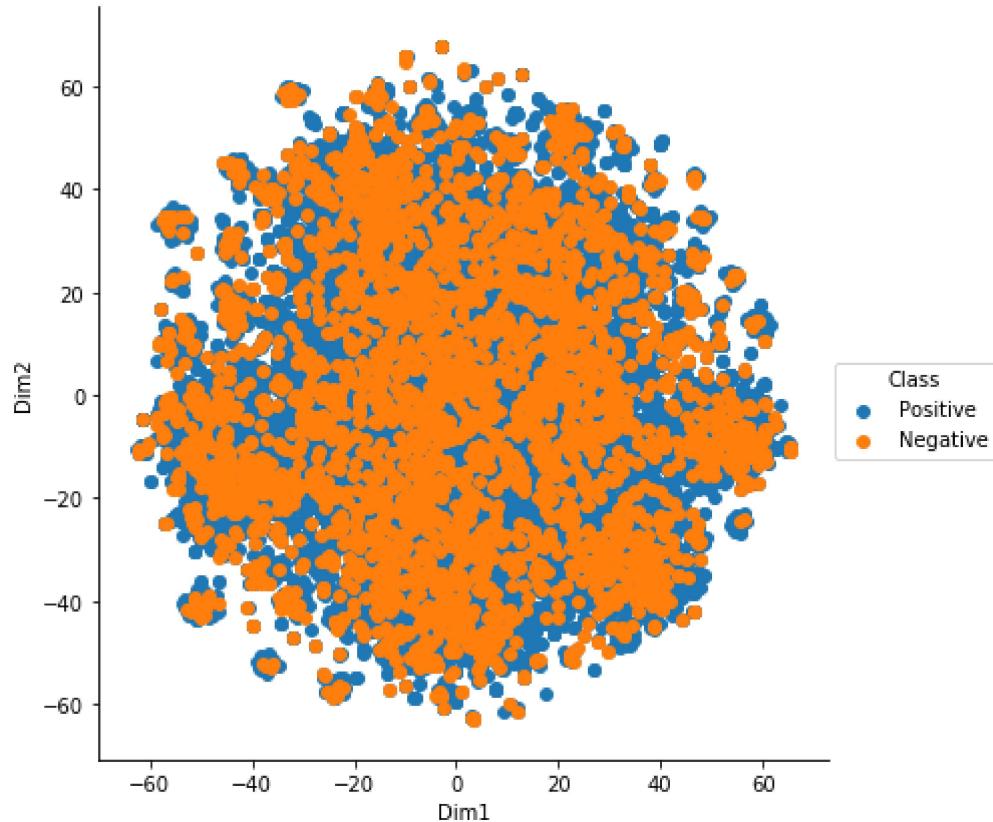
In [30]:

```
%%time
from sklearn.manifold import TSNE
import random

n_samples = 20000
sample_cols = random.sample(range(1, tfidf_w2v_vec_google_norm.shape[0]), n_samples)
sample_features = tfidf_w2v_vec_google_norm[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=35)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```

(20000, 300) (20000, 1)
(20000, 3)



CPU times: user 17min 6s, sys: 1min 3s, total: 18min 9s
Wall time: 18min 9s

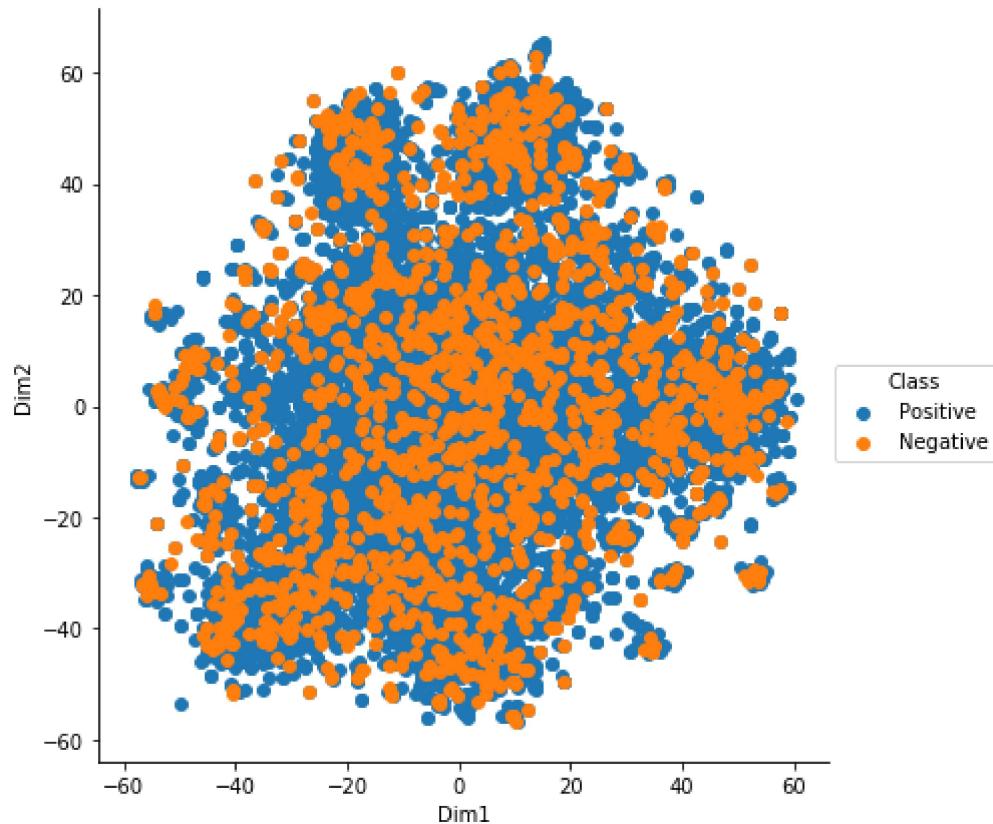
In [31]:

```
%%time
from sklearn.manifold import TSNE
import random

n_samples = 10000
sample_cols = random.sample(range(1, tfidf_w2v_vec_google_norm.shape[0]), n_samples)
sample_features = tfidf_w2v_vec_google_norm[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=40)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```

(10000, 300) (10000, 1)
(10000, 3)



CPU times: user 12min 10s, sys: 33.6 s, total: 12min 43s
Wall time: 12min 43s

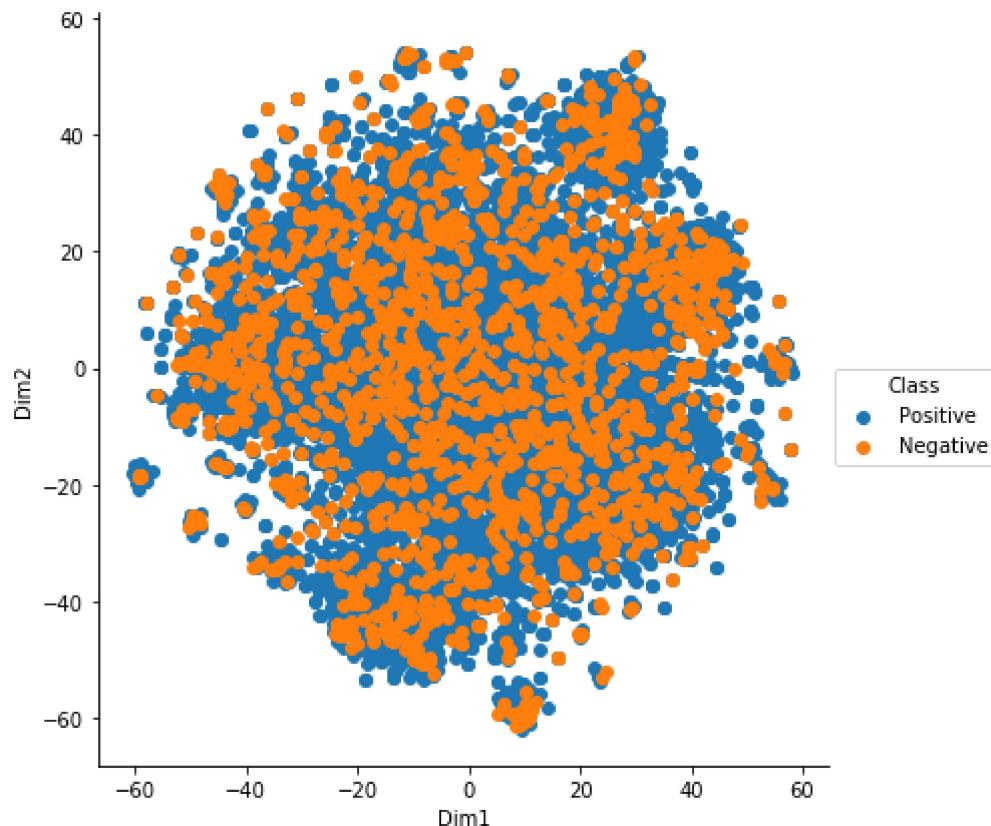
In [32]:

```
%%time
from sklearn.manifold import TSNE
import random

n_samples = 10000
sample_cols = random.sample(range(1, tfidf_w2v_vec_google_norm.shape[0]), n_samples)
sample_features = tfidf_w2v_vec_google_norm[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=45)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```

(10000, 300) (10000, 1)
(10000, 3)



CPU times: user 7min 46s, sys: 32.3 s, total: 8min 18s
Wall time: 8min 18s

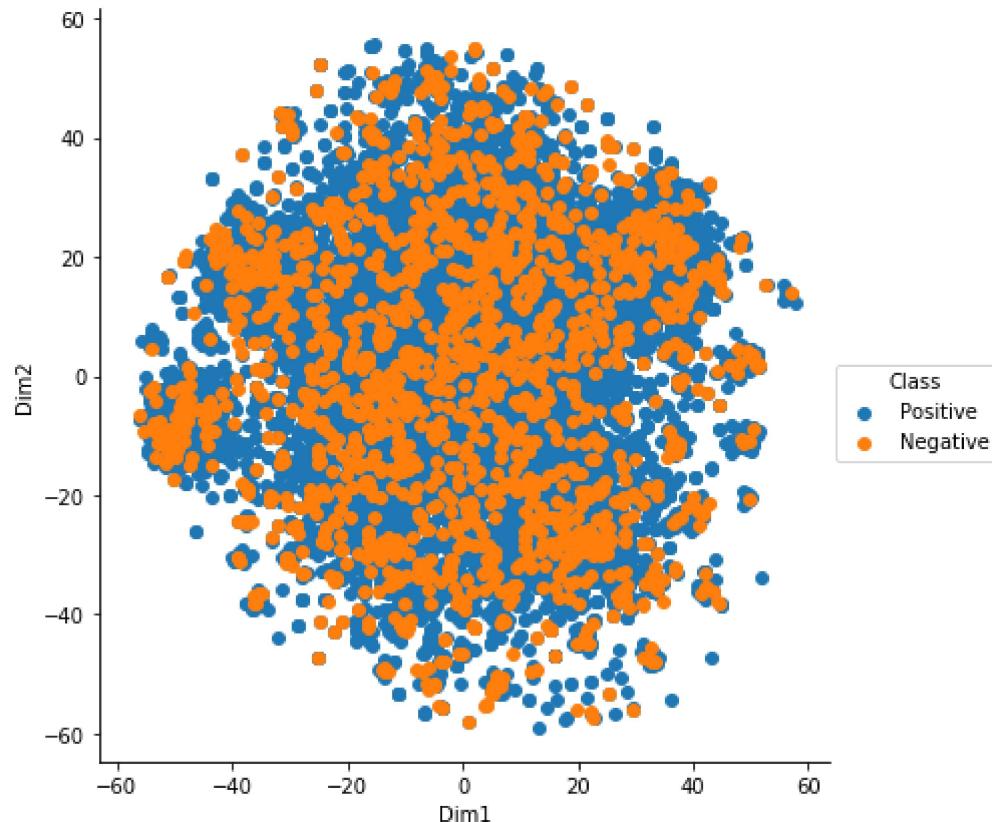
In [34]:

```
%%time
from sklearn.manifold import TSNE
import random

n_samples = 10000
sample_cols = random.sample(range(1, tfidf_w2v_vec_google_norm.shape[0]), n_samples)
sample_features = tfidf_w2v_vec_google_norm[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=50)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```

(10000, 300) (10000, 1)
(10000, 3)



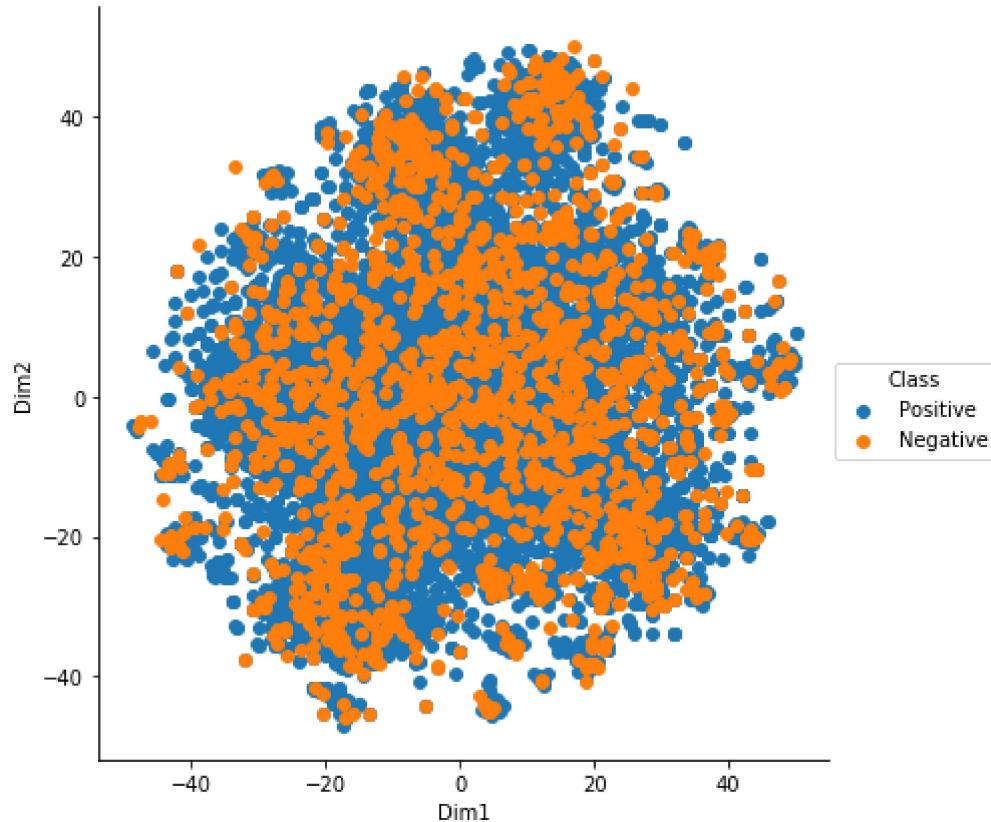
CPU times: user 8min 47s, sys: 32.5 s, total: 9min 20s
Wall time: 9min 20s

```
In [7]: %%time
from sklearn.manifold import TSNE
import random

n_samples = 10000
sample_cols = random.sample(range(1, tfidf_w2v_vec_google_norm.shape[0]), n_samples)
sample_features = tfidf_w2v_vec_google_norm[sample_cols]
# sample_features = df
sample_class = df2['Score'][sample_cols]
sample_class = sample_class[:, np.newaxis]
print(sample_features.shape, sample_class.shape)
model = TSNE(n_components=2, random_state=0, perplexity=70)

embedded_data = model.fit_transform(sample_features)
# print(embedded_data.shape, sample_class.shape)
final_data = np.concatenate((embedded_data, sample_class), axis=1)
print(final_data.shape)
newdf = pd.DataFrame(data=final_data, columns=["Dim1", "Dim2", "Class"])
sns.FacetGrid(newdf, hue="Class", size=6).map(plt.scatter, "Dim1", "Dim2").add_legend()
plt.show()
```

(10000, 300) (10000, 1)
(10000, 3)



CPU times: user 10min 1s, sys: 24.7 s, total: 10min 26s
Wall time: 10min 26s

References:

- (1) <http://blog.aylien.com/10-common-nlp-terms-explained-for-the-text/> (<http://blog.aylien.com/10-common-nlp-terms-explained-for-the-text/>)
- (2) <https://en.wikipedia.org/> (<https://en.wikipedia.org/>)
- (3) <https://buhrmann.github.io/tfidf-analysis.html> (<https://buhrmann.github.io/tfidf-analysis.html>)