# Amazon Fine Food Reviews Analysis

Data Source: https://www.kaggle.com/snap/amazon-fine-food-reviews

EDA: https://nycdatascience.com/blog/student-works/amazon-fine-foods-visualization/

The Amazon Fine Food Reviews dataset consists of reviews of fine foods from Amazon.

Number of reviews: 568,454
Number of users: 256,059
Number of products: 74,258
Timespan: Oct 1999 - Oct 2012
Number of Attributes/Columns in data: 10

Attribute Information:

1. Id
2. ProductId - unique identifier for the product
3. UserId - unqiue 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

**Objective:**

Given a review, determine whether the review is positive (rating of 4 or 5) or negative (rating of 1 or 2).

[Q] How to determine if a review is positive or negative?

[Ans] We could use Score/Rating. A rating of 4 or 5 can be cosnidered as a positive review. A rating of 1 or 2 can be considered as negative one. A review of rating 3 is considered nuetral and such reviews are ignored from our analysis. This is an approximate and proxy way of determining the polarity (positivity/negativity) of a review.

# [1]. Reading Data

## [1.1] Loading the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

In order to load the data, We have used the SQLITE dataset as it is easier to query the data and visualise the data efficiently.

Here as we only want to get the global sentiment of the recommendations (positive or negative), we will purposefully ignore all Scores equal to 3. If the score is above 3, then the recommendation wil be set to "positive". Otherwise, it will be set to "negative".

In [1]:
```python
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")


import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
```

```python
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle

from tqdm import tqdm
import os
```

In [0]:
```python
# using SQLite Table to read data.
con = sqlite3.connect('database.sqlite')

# filtering only positive and negative reviews i.e.
# not taking into consideration those reviews with Score=3
# SELECT * FROM Reviews WHERE Score != 3 LIMIT 500000, will give top 50
0000 data points
# you can change the number to any other number based on your computing
 power

# filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Sco
re != 3 LIMIT 500000""", con)
# for tsne assignment you can take 5k data points
```

```
filtered_data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score
 != 3 LIMIT 5000""", con)

# Give reviews with Score>3 a positive rating(1), and reviews with a sc
ore<3 a negative rating(0).
def partition(x):
    if x < 3:
        return 0
    return 1

#changing reviews with score less than 3 to be positive and vice-versa
actualScore = filtered_data['Score']
positiveNegative = actualScore.map(partition)
filtered_data['Score'] = positiveNegative
print("Number of data points in our data", filtered_data.shape)
filtered_data.head(3)
```

Number of data points in our data (5000, 10)

Out[0]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfulnes |
|---|---|---|---|---|---|---|
| 0 | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 |
| 1 | 2 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0 | 0 |

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfulnes |
|---|---|---|---|---|---|---|
| 2 | 3 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1 | 1 |

In [0]:
```python
display = pd.read_sql_query("""
SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
FROM Reviews
GROUP BY UserId
HAVING COUNT(*)>1
""", con)
```

In [0]:
```python
print(display.shape)
display.head()
```

(80668, 7)

Out[0]:

| | UserId | ProductId | ProfileName | Time | Score | Text | COUI |
|---|---|---|---|---|---|---|---|
| 0 | #oc-R115TNMSPFT9I7 | B007Y59HVM | Breyton | 1331510400 | 2 | Overall its just OK when considering the price... | 2 |
| 1 | #oc-R11D9D7SHXIJB9 | B005HG9ET0 | Louis E. Emory "hoppy" | 1342396800 | 5 | My wife has recurring extreme muscle spasms, u... | 3 |

| | UserId | ProductId | ProfileName | Time | Score | Text | COUI |
|---|---|---|---|---|---|---|---|
| 2 | #oc-R11DNU2NBKQ23Z | B007Y59HVM | Kim Cieszykowski | 1348531200 | 1 | This coffee is horrible and unfortunately not ... | 2 |
| 3 | #oc-R11O5J5ZVQE25C | B005HG9ET0 | Penguin Chick | 1346889600 | 5 | This will be the bottle that you grab from the... | 3 |
| 4 | #oc-R12KPBODL2B5ZD | B007OSBE1U | Christopher P. Presta | 1348617600 | 1 | I didnt like this coffee. Instead of telling y... | 2 |

In [0]: `display[display['UserId']=='AZY10LLTJ71NX']`

Out[0]:

| | UserId | ProductId | ProfileName | Time | Score | Text | |
|---|---|---|---|---|---|---|---|
| 80638 | AZY10LLTJ71NX | B006P7E5ZI | undertheshrine "undertheshrine" | 1334707200 | 5 | I was recommended to try green tea extract to ... | 5 |

In [0]: `display['COUNT(*)'].sum()`

Out[0]: 393063

# [2] Exploratory Data Analysis

## [2.1] Data Cleaning: Deduplication

It is observed (as shown in the table below) that the reviews data had many duplicate entries. Hence it was necessary to remove duplicates in order to get unbiased results for the analysis of the data. Following is an example:

```
In [0]: display= pd.read_sql_query("""
SELECT *
FROM Reviews
WHERE Score != 3 AND UserId="AR5J8UI46CURR"
ORDER BY ProductID
""", con)
display.head()
```

Out[0]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|---|---|---|---|---|---|
| 0 | 78445 | B000HDL1RQ | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 1 | 138317 | B000HDOPYC | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|---|---|---|---|---|---|
| 2 | 138277 | B000HDOPYM | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 3 | 73791 | B000HDOPZG | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |
| 4 | 155049 | B000PAQ75C | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 |

As it can be seen above that same user has multiple reviews with same values for HelpfulnessNumerator, HelpfulnessDenominator, Score, Time, Summary and Text and on doing analysis it was found that

ProductId=B000HDOPZG was Loacker Quadratini Vanilla Wafer Cookies, 8.82-Ounce Packages (Pack of 8)

ProductId=B000HDL1RQ was Loacker Quadratini Lemon Wafer Cookies, 8.82-Ounce Packages (Pack of 8) and so on

It was inferred after analysis that reviews with same parameters other than ProductId belonged to the same product just having different flavour or quantity. Hence in order to reduce redundancy it was decided to eliminate the rows having same parameters.

The method used for the same was that we first sort the data according to ProductId and then just keep the first similar product review and delelte the others. for eg. in the above just the review for ProductId=B000HDL1RQ remains. This method ensures that there is only one representative for each product and deduplication without sorting would lead to possibility of different representatives still existing for the same product.

```
In [0]:  #Sorting data according to ProductId in ascending order
         sorted_data=filtered_data.sort_values('ProductId', axis=0, ascending=Tr
         ue, inplace=False, kind='quicksort', na_position='last')
```

```
In [0]:  #Deduplication of entries
         final=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time"
         ,"Text"}, keep='first', inplace=False)
         final.shape
```

Out[0]: (4986, 10)

```
In [0]:  #Checking to see how much % of data still remains
         (final['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

Out[0]: 99.72

**Observation:-** It was also seen that in two rows given below the value of HelpfulnessNumerator is greater than HelpfulnessDenominator which is not practically possible hence these two rows too are removed from calcualtions

```
In [0]:  display= pd.read_sql_query("""
         SELECT *
         FROM Reviews
         WHERE Score != 3 AND Id=44737 OR Id=64422
         ORDER BY ProductID
         """, con)

         display.head()
```

Out[0]:

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | Helpfuln |
|---|---|---|---|---|---|---|
| 0 | 64422 | B000MIDROQ | A161DK06JJMCYF | J. E. Stephens "Jeanne" | 3 | 1 |
| 1 | 44737 | B001EQ55RW | A2V0I904FH7ABY | Ram | 3 | 2 |

```
In [0]:  final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
```

```
In [0]:  #Before starting the next phase of preprocessing lets see the number of
          entries left
         print(final.shape)

         #How many positive and negative reviews are present in our dataset?
         final['Score'].value_counts()
```

```
(4986, 10)
```

```
Out[0]:  1    4178
         0     808
         Name: Score, dtype: int64
```

## [3] Preprocessing

## [3.1]. Preprocessing Review Text

Now that we have finished deduplication our data requires some preprocessing before we go on further with analysis and making the prediction model.

Hence in the Preprocessing phase we do the following in the order below:-

1. Begin by removing the html tags
2. Remove any punctuations or limited set of special characters like , or . or # etc.
3. Check if the word is made up of english letters and is not alpha-numeric
4. Check to see if the length of the word is greater than 2 (as it was researched that there is no adjective in 2-letters)
5. Convert the word to lowercase
6. Remove Stopwords
7. Finally Snowball Stemming the word (it was obsereved to be better than Porter Stemming)

After which we collect the words used to describe positive and negative reviews

```
In [0]:  # printing some random reviews
         sent_0 = final['Text'].values[0]
         print(sent_0)
         print("="*50)

         sent_1000 = final['Text'].values[1000]
         print(sent_1000)
         print("="*50)

         sent_1500 = final['Text'].values[1500]
         print(sent_1500)
         print("="*50)

         sent_4900 = final['Text'].values[4900]
         print(sent_4900)
         print("="*50)
```

Why is this $[...] when the same product is available for $[...] here?<
br />http://www.amazon.com/VICTOR-FLY-MAGNET-BATT-REFILL/dp/B00004RBDY<

br /><http://www.amazon.com/VICTOR-FLY-MAGNET-BAIT-REFILL/dp/B0000HRDD~

br /><br />The Victor M380 and M502 traps are unreal, of course -- total fly genocide. Pretty stinky, but only right nearby.
====================================================
I recently tried this flavor/brand and was surprised at how delicious these chips are.  The best thing was that there were a lot of "brown" chips in the bsg (my favorite), so I bought some more through amazon and shared with family and friends.  I am a little disappointed that there are not, so far, very many brown chips in these bags, but the flavor is still very good.  I like them better than the yogurt and green onion flavor because they do not seem to be as salty, and the onion flavor is better.  If you haven't eaten Kettle chips before, I recommend that you try a bag before buying bulk.  They are thicker and crunchier than Lays but just as fresh out of the bag.
====================================================
Wow.  So far, two two-star reviews.  One obviously had no idea what they were ordering; the other wants crispy cookies.  Hey, I'm sorry; but these reviews do nobody any good beyond reminding us to look  before ordering.<br /><br />These are chocolate-oatmeal cookies.  If you don't like that combination, don't order this type of cookie.  I find the combo quite nice, really.  The oatmeal sort of "calms" the rich chocolate flavor and gives the cookie sort of a coconut-type consistency.  Now let's also remember that tastes differ; so, I've given my opinion.<br /><br />Then, these are soft, chewy cookies -- as advertised.  They are not "crispy" cookies, or the blurb would say "crispy," rather than "chewy."
I happen to like raw cookie dough; however, I don't see where these taste like raw cookie dough.  Both are soft, however, so is this the confusion?  And, yes, they stick together.  Soft cookies tend to do that.  They aren't individually wrapped, which would add to the cost.  Oh yeah, chocolate chip cookies tend to be somewhat sweet.<br /><br />So, if you want something hard and crisp, I suggest Nabiso's Ginger Snaps.  If you want a cookie that's soft, chewy and tastes like a combination of chocolate and oatmeal, give these a try.  I'm here to place my second order.
====================================================
love to order my coffee on amazon.  easy and shows up quickly.<br />This k cup is great coffee.  dcaf is very good as well
====================================================

In [0]: # remove urls from text python: https://stackoverflow.com/a/40823105/40
84039

```
sent_0 = re.sub(r"http\S+", "", sent_0)
sent_1000 = re.sub(r"http\S+", "", sent_1000)
sent_150 = re.sub(r"http\S+", "", sent_1500)
sent_4900 = re.sub(r"http\S+", "", sent_4900)

print(sent_0)
```

Why is this $[...] when the same product is available for $[...] here?<
br /> /><br />The Victor M380 and M502 traps are unreal, of course -- t
otal fly genocide. Pretty stinky, but only right nearby.

In [0]:
```
# https://stackoverflow.com/questions/16206380/python-beautifulsoup-how
-to-remove-all-tags-from-an-element
from bs4 import BeautifulSoup

soup = BeautifulSoup(sent_0, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_1000, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_1500, 'lxml')
text = soup.get_text()
print(text)
print("="*50)

soup = BeautifulSoup(sent_4900, 'lxml')
text = soup.get_text()
print(text)
```

Why is this $[...] when the same product is available for $[...] here?
/>The Victor M380 and M502 traps are unreal, of course -- total fly gen
ocide. Pretty stinky, but only right nearby.
==================================================
I recently tried this flavor/brand and was surprised at how delicious t
hese chips are.  The best thing was that there were a lot of "brown" ch

ips in the bsg (my favorite), so I bought some more through amazon and shared with family and friends.  I am a little disappointed that there are not, so far, very many brown chips in these bags, but the flavor is still very good.  I like them better than the yogurt and green onion fl avor because they do not seem to be as salty, and the onion flavor is b etter.  If you haven't eaten Kettle chips before, I recommend that you try a bag before buying bulk.  They are thicker and crunchier than Lays but just as fresh out of the bag.
==================================================
Wow.  So far, two two-star reviews.  One obviously had no idea what the y were ordering; the other wants crispy cookies.  Hey, I'm sorry; but t hese reviews do nobody any good beyond reminding us to look  before ord ering.These are chocolate-oatmeal cookies.  If you don't like that comb ination, don't order this type of cookie.  I find the combo quite nice, really.  The oatmeal sort of "calms" the rich chocolate flavor and give s the cookie sort of a coconut-type consistency.  Now let's also rememb er that tastes differ; so, I've given my opinion.Then, these are soft, chewy cookies -- as advertised.  They are not "crispy" cookies, or the blurb would say "crispy," rather than "chewy."  I happen to like raw co okie dough; however, I don't see where these taste like raw cookie doug h.  Both are soft, however, so is this the confusion?  And, yes, they s tick together.  Soft cookies tend to do that.  They aren't individually wrapped, which would add to the cost.  Oh yeah, chocolate chip cookies tend to be somewhat sweet.So, if you want something hard and crisp, I s uggest Nabiso's Ginger Snaps.  If you want a cookie that's soft, chewy and tastes like a combination of chocolate and oatmeal, give these a tr y.  I'm here to place my second order.
==================================================
love to order my coffee on amazon.  easy and shows up quickly.This k cu p is great coffee.  dcaf is very good as well

In [0]:
```python
# https://stackoverflow.com/a/47091490/4084039
import re

def decontracted(phrase):
    # specific
    phrase = re.sub(r"won't", "will not", phrase)
    phrase = re.sub(r"can\'t", "can not", phrase)
```

```python
        # general
        phrase = re.sub(r"n\'t", " not", phrase)
        phrase = re.sub(r"\'re", " are", phrase)
        phrase = re.sub(r"\'s", " is", phrase)
        phrase = re.sub(r"\'d", " would", phrase)
        phrase = re.sub(r"\'ll", " will", phrase)
        phrase = re.sub(r"\'t", " not", phrase)
        phrase = re.sub(r"\'ve", " have", phrase)
        phrase = re.sub(r"\'m", " am", phrase)
        return phrase
```

In [0]:
```python
sent_1500 = decontracted(sent_1500)
print(sent_1500)
print("="*50)
```

```
Wow.  So far, two two-star reviews.  One obviously had no idea what the
y were ordering; the other wants crispy cookies.  Hey, I am sorry; but
these reviews do nobody any good beyond reminding us to look  before or
dering.<br /><br />These are chocolate-oatmeal cookies.  If you do not
like that combination, do not order this type of cookie.  I find the co
mbo quite nice, really.  The oatmeal sort of "calms" the rich chocolate
flavor and gives the cookie sort of a coconut-type consistency.  Now le
t is also remember that tastes differ; so, I have given my opinion.<br
/><br />Then, these are soft, chewy cookies -- as advertised.  They are
not "crispy" cookies, or the blurb would say "crispy," rather than "che
wy."  I happen to like raw cookie dough; however, I do not see where th
ese taste like raw cookie dough.  Both are soft, however, so is this th
e confusion?  And, yes, they stick together.  Soft cookies tend to do t
hat.  They are not individually wrapped, which would add to the cost.
Oh yeah, chocolate chip cookies tend to be somewhat sweet.<br /><br />S
o, if you want something hard and crisp, I suggest Nabiso is Ginger Sna
ps.  If you want a cookie that is soft, chewy and tastes like a combina
tion of chocolate and oatmeal, give these a try.  I am here to place my
second order.
==================================================
```

In [0]:
```python
#remove words with numbers python: https://stackoverflow.com/a/1808237
0/4084039
```

```python
sent_0 = re.sub("\S*\d\S*", "", sent_0).strip()
print(sent_0)
```

Why is this $[...] when the same product is available for $[...] here?<br /> /><br />The Victor  and  traps are unreal, of course -- total fly genocide. Pretty stinky, but only right nearby.

In [0]:
```python
#remove spacial character: https://stackoverflow.com/a/5843547/4084039
sent_1500 = re.sub('[^A-Za-z0-9]+', ' ', sent_1500)
print(sent_1500)
```

Wow So far two two star reviews One obviously had no idea what they were ordering the other wants crispy cookies Hey I am sorry but these reviews do nobody any good beyond reminding us to look before ordering br br These are chocolate oatmeal cookies If you do not like that combination do not order this type of cookie I find the combo quite nice really The oatmeal sort of calms the rich chocolate flavor and gives the cookie sort of a coconut type consistency Now let is also remember that tastes differ so I have given my opinion br br Then these are soft chewy cookies as advertised They are not crispy cookies or the blurb would say crispy rather than chewy I happen to like raw cookie dough however I do not see where these taste like raw cookie dough Both are soft however so is this the confusion And yes they stick together Soft cookies tend to do that They are not individually wrapped which would add to the cost Oh yeah chocolate chip cookies tend to be somewhat sweet br br So if you want something hard and crisp I suggest Nabiso is Ginger Snaps If you want a cookie that is soft chewy and tastes like a combination of chocolate and oatmeal give these a try I am here to place my second order

In [0]:
```python
# https://gist.github.com/sebleier/554280
# we are removing the words from the stop words list: 'no', 'nor', 'not'
# <br /><br /> ==> after the above steps, we are getting "br br"
# we are including them into stop words list
# instead of <br /> if we have <br/> these tags would have revmoved in
 the 1st step

stopwords= set(['br', 'the', 'i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've",\
```

```
                "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselve
s', 'he', 'him', 'his', 'himself', \
                'she', "she's", 'her', 'hers', 'herself', 'it', "it's", 'it
s', 'itself', 'they', 'them', 'their',\
                'theirs', 'themselves', 'what', 'which', 'who', 'whom', 'th
is', 'that', "that'll", 'these', 'those', \
                'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'h
ave', 'has', 'had', 'having', 'do', 'does', \
                'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or',
 'because', 'as', 'until', 'while', 'of', \
                'at', 'by', 'for', 'with', 'about', 'against', 'between',
'into', 'through', 'during', 'before', 'after',\
                'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out',
'on', 'off', 'over', 'under', 'again', 'further',\
                'then', 'once', 'here', 'there', 'when', 'where', 'why', 'h
ow', 'all', 'any', 'both', 'each', 'few', 'more',\
                'most', 'other', 'some', 'such', 'only', 'own', 'same', 's
o', 'than', 'too', 'very', \
                's', 't', 'can', 'will', 'just', 'don', "don't", 'should',
"should've", 'now', 'd', 'll', 'm', 'o', 're', \
                've', 'y', 'ain', 'aren', "aren't", 'couldn', "couldn't",
'didn', "didn't", 'doesn', "doesn't", 'hadn',\
                "hadn't", 'hasn', "hasn't", 'haven', "haven't", 'isn', "is
n't", 'ma', 'mightn', "mightn't", 'mustn',\
                "mustn't", 'needn', "needn't", 'shan', "shan't", 'shouldn',
 "shouldn't", 'wasn', "wasn't", 'weren', "weren't", \
                'won', "won't", 'wouldn', "wouldn't"])
```

In [0]:
```python
# Combining all the above stundents
from tqdm import tqdm
preprocessed_reviews = []
# tqdm is for printing the status bar
for sentance in tqdm(final['Text'].values):
    sentance = re.sub(r"http\S+", "", sentance)
    sentance = BeautifulSoup(sentance, 'lxml').get_text()
    sentance = decontracted(sentance)
    sentance = re.sub("\S*\d\S*", "", sentance).strip()
    sentance = re.sub('[^A-Za-z]+', ' ', sentance)
    # https://gist.github.com/sebleier/554280
```

```python
    sentance = ' '.join(e.lower() for e in sentance.split() if e.lower
() not in stopwords)
    preprocessed_reviews.append(sentance.strip())
```

```
100%|████████████████████████████████████████████████████████
████████| 4986/4986 [00:01<00:00, 3137.37it/s]
```

In [0]: `preprocessed_reviews[1500]`

Out[0]: `'wow far two two star reviews one obviously no idea ordering wants crispy cookies hey sorry reviews nobody good beyond reminding us look ordering chocolate oatmeal cookies not like combination not order type cookie find combo quite nice really oatmeal sort calms rich chocolate flavor gives cookie sort coconut type consistency let also remember tastes differ given opinion soft chewy cookies advertised not crispy cookies blurb would say crispy rather chewy happen like raw cookie dough however not see taste like raw cookie dough soft however confusion yes stick together soft cookies tend not individually wrapped would add cost oh yeah chocolate chip cookies tend somewhat sweet want something hard crisp suggest nabiso ginger snaps want cookie soft chewy tastes like combination chocolate oatmeal give try place second order'`

### [3.2] Preprocessing Review Summary

In [0]: `## Similartly you can do preprocessing for review summary also.`

# [4] Featurization

## [4.1] BAG OF WORDS

In [0]:
```python
#BoW
count_vect = CountVectorizer() #in scikit-learn
count_vect.fit(preprocessed_reviews)
print("some feature names ", count_vect.get_feature_names()[:10])
```

```
print('='*50)

final_counts = count_vect.transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_counts))
print("the shape of out text BOW vectorizer ",final_counts.get_shape())
print("the number of unique words ", final_counts.get_shape()[1])
```

```
some feature names  ['aa', 'aahhhs', 'aback', 'abandon', 'abates', 'abb
ott', 'abby', 'abdominal', 'abiding', 'ability']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (4986, 12997)
the number of unique words  12997
```

## [4.2] Bi-Grams and n-Grams.

In [0]:
```
#bi-gram, tri-gram and n-gram

#removing stop words like "not" should be avoided before building n-gra
ms
# count_vect = CountVectorizer(ngram_range=(1,2))
# please do read the CountVectorizer documentation http://scikit-learn.
org/stable/modules/generated/sklearn.feature_extraction.text.CountVecto
rizer.html

# you can choose these numebrs min_df=10, max_features=5000, of your ch
oice
count_vect = CountVectorizer(ngram_range=(1,2), min_df=10, max_features
=5000)
final_bigram_counts = count_vect.fit_transform(preprocessed_reviews)
print("the type of count vectorizer ",type(final_bigram_counts))
print("the shape of out text BOW vectorizer ",final_bigram_counts.get_s
hape())
print("the number of unique words including both unigrams and bigrams "
, final_bigram_counts.get_shape()[1])
```

```
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer  (4986, 3144)
the number of unique words including both unigrams and bigrams  3144
```

the number of unique words including both unigrams and bigrams    3144

## [4.3] TF-IDF

```
In [0]:  tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10)
         tf_idf_vect.fit(preprocessed_reviews)
         print("some sample features(unique words in the corpus)",tf_idf_vect.ge
         t_feature_names()[0:10])
         print('='*50)

         final_tf_idf = tf_idf_vect.transform(preprocessed_reviews)
         print("the type of count vectorizer ",type(final_tf_idf))
         print("the shape of out text TFIDF vectorizer ",final_tf_idf.get_shape
         ())
         print("the number of unique words including both unigrams and bigrams "
         , final_tf_idf.get_shape()[1])
```

```
some sample features(unique words in the corpus) ['ability', 'able', 'a
ble find', 'able get', 'absolute', 'absolutely', 'absolutely deliciou
s', 'absolutely love', 'absolutely no', 'according']
==================================================
the type of count vectorizer  <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer  (4986, 3144)
the number of unique words including both unigrams and bigrams  3144
```

## [4.4] Word2Vec

```
In [0]:  # Train your own Word2Vec model using your own text corpus
         i=0
         list_of_sentance=[]
         for sentance in preprocessed_reviews:
             list_of_sentance.append(sentance.split())
```

```
In [0]:  # Using Google News Word2Vectors

         # in this project we are using a pretrained model by google
```

```python
# its 3.3G file, once you load this into your memory
# it occupies ~9Gb, so please do this step only if you have >12G of ram
# we will provide a pickle file wich contains a dict ,
# and it contains all our courpus words as keys and  model[word] as val
ues
# To use this code-snippet, download "GoogleNews-vectors-negative300.bi
n"
# from https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS21pQmM/edi
t
# it's 1.9GB in size.


# http://kavita-ganesan.com/gensim-word2vec-tutorial-starter-code/#.W17
SRFAzZPY
# you can comment this whole cell
# or change these varible according to your need

is_your_ram_gt_16g=False
want_to_use_google_w2v = False
want_to_train_w2v = True

if want_to_train_w2v:
    # min_count = 5 considers only words that occured atleast 5 times
    w2v_model=Word2Vec(list_of_sentance,min_count=5,size=50, workers=4)
    print(w2v_model.wv.most_similar('great'))
    print('='*50)
    print(w2v_model.wv.most_similar('worst'))

elif want_to_use_google_w2v and is_your_ram_gt_16g:
    if os.path.isfile('GoogleNews-vectors-negative300.bin'):
        w2v_model=KeyedVectors.load_word2vec_format('GoogleNews-vectors
-negative300.bin', binary=True)
        print(w2v_model.wv.most_similar('great'))
        print(w2v_model.wv.most_similar('worst'))
    else:
        print("you don't have gogole's word2vec file, keep want_to_trai
n_w2v = True, to train your own w2v ")
```

```
[('snack', 0.9951335191726685), ('calorie', 0.9946465492248535), ('wond
```

```
erful', 0.9946032166481018), ('excellent', 0.9944332838058472), ('espec
ially', 0.9941144585609436), ('baked', 0.9940600395202637), ('salted',
0.994047224521637), ('alternative', 0.9937226176261902), ('tasty', 0.99
36816692352295), ('healthy', 0.9936649799346924)]
==================================================
[('varieties', 0.9994194507598877), ('become', 0.9992934465408325), ('p
opcorn', 0.9992750883102417), ('de', 0.9992610216140747), ('miss', 0.99
92451071739197), ('melitta', 0.999218761920929), ('choice', 0.999210238
4567261), ('american', 0.9991837739944458), ('beef', 0.999178051948547
4), ('finish', 0.9991567134857178)]
```

In [0]:
```python
w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
print("sample words ", w2v_words[0:50])
```

```
number of words that occured minimum 5 times  3817
sample words  ['product', 'available', 'course', 'total', 'pretty', 'st
inky', 'right', 'nearby', 'used', 'ca', 'not', 'beat', 'great', 'receiv
ed', 'shipment', 'could', 'hardly', 'wait', 'try', 'love', 'call', 'ins
tead', 'removed', 'easily', 'daughter', 'designed', 'printed', 'use',
'car', 'windows', 'beautifully', 'shop', 'program', 'going', 'lot', 'fu
n', 'everywhere', 'like', 'tv', 'computer', 'really', 'good', 'idea',
'final', 'outstanding', 'window', 'everybody', 'asks', 'bought', 'mad
e']
```

## [4.4.1] Converting text into vectors using Avg W2V, TFIDF-W2V

**[4.4.1.1] Avg W2v**

In [0]:
```python
# average Word2Vec
# compute average word2vec for each review.
sent_vectors = []; # the avg-w2v for each sentence/review is stored in
 this list
for sent in tqdm(list_of_sentance): # for each review/sentence
    sent_vec = np.zeros(50) # as word vectors are of zero length 50, yo
```

```
u might need to change this to 300 if you use google's w2v
    cnt_words =0; # num of words with a valid vector in the sentence/re
view
    for word in sent: # for each word in a review/sentence
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    sent_vectors.append(sent_vec)
print(len(sent_vectors))
print(len(sent_vectors[0]))
```

100%|████████████████████████████████████████████████████████
████████████| 4986/4986 [00:03<00:00, 1330.47it/s]

4986
50


**[4.4.1.2] TFIDF weighted W2v**

In [0]:
```
# S = ["abc def pqr", "def def def abc", "pqr pqr def"]
model = TfidfVectorizer()
tf_idf_matrix = model.fit_transform(preprocessed_reviews)
# we are converting a dictionary with word as a key, and the idf as a v
alue
dictionary = dict(zip(model.get_feature_names(), list(model.idf_)))
```

In [0]:
```
# TF-IDF weighted Word2Vec
tfidf_feat = model.get_feature_names() # tfidf words/col-names
# final_tf_idf is the sparse matrix with row= sentence, col=word and ce
ll_val = tfidf

tfidf_sent_vectors = []; # the tfidf-w2v for each sentence/review is st
ored in this list
row=0;
for sent in tqdm(list_of_sentance): # for each review/sentence
```

```python
        sent_vec = np.zeros(50) # as word vectors are of zero length
        weight_sum =0; # num of words with a valid vector in the sentence/r
eview
        for word in sent: # for each word in a review/sentence
            if word in w2v_words and word in tfidf_feat:
                vec = w2v_model.wv[word]
#                tf_idf = tf_idf_matrix[row, tfidf_feat.index(word)]
                # to reduce the computation we are
                # dictionary[word] = idf value of word in whole courpus
                # sent.count(word) = tf valeus of word in this review
                tf_idf = dictionary[word]*(sent.count(word)/len(sent))
                sent_vec += (vec * tf_idf)
                weight_sum += tf_idf
        if weight_sum != 0:
            sent_vec /= weight_sum
        tfidf_sent_vectors.append(sent_vec)
        row += 1
```

```
100%|████████████████████████████████████████████████████████
███████████| 4986/4986 [00:20<00:00, 245.63it/s]
```

# [5] Assignment 8: Decision Trees

1. **Apply Decision Trees on these feature sets**

   - SET 1:Review text, preprocessed one converted into vectors using (BOW)
   - SET 2:Review text, preprocessed one converted into vectors using (TFIDF)
   - SET 3:Review text, preprocessed one converted into vectors using (AVG W2v)
   - SET 4:Review text, preprocessed one converted into vectors using (TFIDF W2v)

2. **The hyper paramter tuning (best `depth` in range [1, 5, 10, 50, 100, 500, 100], and the best `min_samples_split` in range [5, 10, 100, 500])**

   - Find the best hyper parameter which will give the maximum AUC value

- Find the best hyper paramter using k-fold cross validation or simple cross validation data
- Use gridsearch cv or randomsearch cv or you can also write your own for loops to do this task of hyperparameter tuning

3. **Graphviz**

   - Visualize your decision tree with Graphviz. It helps you to understand how a decision is being made, given a new vector.
   - Since feature names are not obtained from word2vec related models, visualize only BOW & TFIDF decision trees using Graphviz
   - Make sure to print the words in each node of the decision tree instead of printing its index.
   - Just for visualization purpose, limit max_depth to 2 or 3 and either embed the generated images of graphviz in your notebook, or directly upload them as .png files.

4. **Feature importance**

   - Find the top 20 important features from both feature sets Set 1 and Set 2 using `feature_importances_` method of Decision Tree Classifier and print their corresponding feature names

5. **Feature engineering**

   - To increase the performance of your model, you can also experiment with with feature engineering like :
     - Taking length of reviews as another feature.
     - Considering some features from review summary as well.

6. **Representation of results**

   - You need to plot the performance of model both on train data and cross validation data for each hyper parameter, like shown in the figure.

Once after you found the best hyper parameter, you need to train your model with it, and find the AUC on test data and plot the ROC curve on both train and test.

Along with plotting ROC curve, you need to print the confusion matrix with predicted and original labels of test data points. Please visualize your confusion matrices using seaborn heatmaps.



7. **Conclusion**

- You need to summarize the results at the end of the notebook, summarize it in the table format. To print out a table please refer to this prettytable library link



**Note: Data Leakage**

1. There will be an issue of data-leakage if you vectorize the entire data and then split it into train/cv/test.
2. To avoid the issue of data-leakag, make sure to split your data first and then vectorize it.
3. While vectorizing your data, apply the method fit_transform() on you train data, and apply the method transform() on cv/test data.
4. For more details please go through this link.

```
In [2]: # using SQLite Table to read data.
        con = sqlite3.connect('final.sqlite')
        final = pd.read_sql_query(""" SELECT * FROM Reviews""", con)
        final = final[:50000]
```

```
In [3]: final.head()
```

Out[3]:

| | index | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator |
|---|---|---|---|---|---|---|

| | index | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator |
|---|---|---|---|---|---|---|
| 0 | 138706 | 150524 | 0006641040 | ACITT7DI6IDDL | shari zychinski | 0 |
| 1 | 138688 | 150506 | 0006641040 | A2IW4PEEKO2R0U | Tracy | 1 |
| 2 | 138689 | 150507 | 0006641040 | A1S4A3IQ2MU7V4 | sally sue "sally sue" | 1 |
| 3 | 138690 | 150508 | 0006641040 | AZGXZ2UUK6X | Catherine Hallberg " (Kate)" | 1 |

| | index | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator |
|---|---|---|---|---|---|---|
| 4 | 138691 | 150509 | 0006641040 | A3CMRKGE0P909G | Teresa | 3 |

```
In [4]: from sklearn.model_selection import train_test_split
        ##Sorting data according to Time in ascending order for Time Based Splitting
        time_sorted_data = final.sort_values('Time', axis=0, ascending=True, inplace=False, kind='quicksort', na_position='last')

        x = time_sorted_data['CleanedText'].values
        y = time_sorted_data['Score']

        # split the data set into train and test
        X_train, X_test, Y_train, Y_test = train_test_split(x, y, test_size=0.3, random_state=0,shuffle=False)
```

## Applying Decision Trees

### [5.1] Applying Decision Trees on BOW, SET 1

```
In [5]: # Please write all the code with proper documentation
        #BoW
        count_vect = CountVectorizer(min_df = 500)
        X_train_vec = count_vect.fit_transform(X_train)
```

```python
X_test_vec = count_vect.transform(X_test)
print("the type of count vectorizer :",type(X_train_vec))
print("the shape of out text BOW vectorizer : ",X_train_vec.get_shape
())
print("the number of unique words :", X_train_vec.get_shape()[1])
```

```
the type of count vectorizer : <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text BOW vectorizer :  (35000, 450)
the number of unique words : 450
```

In [6]:
```python
from sklearn.preprocessing import StandardScaler
sc = StandardScaler(with_mean=False)
X_train_vec_standardized = sc.fit_transform(X_train_vec)
X_test_vec_standardized = sc.transform(X_test_vec)
```

In [7]:
```python
# Importing libraries
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score,confusion_matrix,f1_score,pr
ecision_score,recall_score

param_grid = {'max_depth': [1, 5, 10, 50, 100, 500, 1000],'min_samples_
split':[5,10,100,500]}
model = GridSearchCV(DecisionTreeClassifier(class_weight ='balanced'),
param_grid, scoring = 'roc_auc',cv=3 , n_jobs = -1,pre_dispatch=2)
model.fit(X_train_vec_standardized, Y_train)
print("Model with best parameters :\n",model.best_estimator_)
print("Accuracy of the model : ",model.score(X_test_vec_standardized, Y
_test))
```

```
Model with best parameters :
 DecisionTreeClassifier(class_weight='balanced', criterion='gini',
            max_depth=50, max_features=None, max_leaf_nodes=None,
            min_impurity_decrease=0.0, min_impurity_split=None,
            min_samples_leaf=1, min_samples_split=500,
            min_weight_fraction_leaf=0.0, presort=False, random_state=N
one,
            splitter='best')
Accuracy of the model :  0.7980626979648984
```

```python
In [8]:  # Cross-Validation errors
         cv_errors = [(1-i)*100 for i in model.cv_results_['mean_test_score']]
         training_scores=[(1-i)*100 for i in model.cv_results_['mean_train_scor
         e']]

         # Optimal value of depth
         optimal_depth = model.best_estimator_.max_depth
         print("The optimal value of depth is : ",optimal_depth)

         optimal_split = model.best_estimator_.min_samples_split
         print("The optimal number of base learners is : ",optimal_split)
```

```
The optimal value of depth is :  50
The optimal number of base learners is :  500
```

```python
In [19]:  max_depths = [1, 5, 10, 50, 100, 500, 1000]

          train_results = []
          test_results = []
          for max_depth in max_depths:
              dt = DecisionTreeClassifier(max_depth=max_depth)
              dt.fit(X_train_vec_standardized,Y_train)

              train_pred = dt.predict(X_train_vec_standardized)

              false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
          rain, train_pred)
              roc_auc = auc(false_positive_rate, true_positive_rate)
              # Add auc score to previous train results
              train_results.append(roc_auc)

              Y_pred = dt.predict(X_test_vec_standardized)

              false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
          est, Y_pred)
              roc_auc = auc(false_positive_rate, true_positive_rate)
              # Add auc score to previous test results
              test_results.append(roc_auc)
```

```python
from matplotlib.legend_handler import HandlerLine2D

line1, = plt.plot(max_depths, train_results, 'b', label='Train AUC')
line2, = plt.plot(max_depths, test_results, 'r', label='Test AUC')

plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})

plt.ylabel('AUC score')
plt.xlabel('Tree depth')
plt.show()
```



```python
In [20]: min_samples_splits = [5,10,100,500]

train_results1 = []
test_results1 = []
for min_samples_split in min_samples_splits:
    dt = DecisionTreeClassifier(min_samples_split=min_samples_split)
    dt.fit(X_train_vec_standardized,Y_train)

    train_pred = dt.predict(X_train_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
```

```python
rain, train_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous train results
    train_results1.append(roc_auc)

    Y_pred = dt.predict(X_test_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
est, Y_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous test results
    test_results1.append(roc_auc)

from matplotlib.legend_handler import HandlerLine2D

line1, = plt.plot(min_samples_splits, train_results1, 'b', label='Train
 AUC')
line2, = plt.plot(min_samples_splits, test_results1, 'r', label='Test A
UC')

plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})

plt.ylabel('AUC score')
plt.xlabel('min samle splits')
plt.show()
```

In [21]:
```python
# DecisionTreeClassifier with Optimal value of depth
dt = DecisionTreeClassifier(max_depth=optimal_depth,min_samples_split=o
ptimal_split)
dt.fit(X_train_vec_standardized,Y_train)
predictions = dt.predict(X_test_vec_standardized)
predictions1 = dt.predict(X_train_vec_standardized)

# Variables that will be used for  making table in Conclusion part of t
his assignment
bow_depth = optimal_depth
bow_split = optimal_split
bow_train_acc = model.score(X_test_vec_standardized, Y_test)*100
bow_test_acc = accuracy_score(Y_test, predictions) * 100
```

In [22]:
```python
train_fpr, train_tpr, thresholds = roc_curve(Y_train, dt.predict_proba(
X_train_vec_standardized)[:,1])
test_fpr, test_tpr, thresholds = roc_curve(Y_test, dt.predict_proba(X_t
est_vec_standardized)[:,1])

plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
rain_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
```

```
tpr)))
plt.legend()
plt.xlabel("depth: hyperparameter")
plt.ylabel("AUC")
plt.title("ERROR PLOTS")
plt.show()
```

ERROR PLOTS



In [23]:
```
# evaluate accuracy on test data
acc = accuracy_score(Y_test, predictions) * 100
print('\nThe Test Accuracy of the DecisionTreeClassifier for depth = %d
 is %f%%' % (optimal_depth, acc))

# evaluate precision
acc = precision_score(Y_test, predictions, pos_label = 1)
print('\nThe Test Precision of the DecisionTreeClassifier for depth = %
d is %f' % (optimal_depth, acc))

# evaluate recall
acc = recall_score(Y_test, predictions, pos_label = 1)
print('\nThe Test Recall of the DecisionTreeClassifier for depth = %d i
s %f' % (optimal_depth, acc))
```

```python
# evaluate f1-score
acc = f1_score(Y_test, predictions, pos_label = 1)
print('\nThe Test F1-Score of the DecisionTreeClassifier for depth = %d
 is %f' % (optimal_depth, acc))
```

The Test Accuracy of the DecisionTreeClassifier for depth = 50 is 82.94
6667%

The Test Precision of the DecisionTreeClassifier for depth = 50 is 0.85
2037

The Test Recall of the DecisionTreeClassifier for depth = 50 is 0.95786
5

The Test F1-Score of the DecisionTreeClassifier for depth = 50 is 0.901
857

In [24]:
```python
# Code for drawing seaborn heatmaps on test data
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(Y_test, predictions), index=
class_names, columns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```

## Confusion Matrix

In [25]:
```python
# evaluate accuracy on train data
acc = accuracy_score(Y_train, predictions1) * 100
print('\nThe Train Accuracy of the DecisionTreeClassifier for depth = %
d is %f%%' % (optimal_depth, acc))

# evaluate precision
acc = precision_score(Y_train, predictions1, pos_label = 1)
print('\nThe Train Precision of the DecisionTreeClassifier for depth =
%d is %f' % (optimal_depth, acc))

# evaluate recall
```

```python
acc = recall_score(Y_train, predictions1, pos_label = 1)
print('\nThe Train Recall of the DecisionTreeClassifier for depth = %d
 is %f' % (optimal_depth, acc))

# evaluate f1-score
acc = f1_score(Y_train, predictions1, pos_label = 1)
print('\nThe Train F1-Score of the DecisionTreeClassifier for depth = %
d is %f' % (optimal_depth, acc))
```

The Train Accuracy of the DecisionTreeClassifier for depth = 50 is 88.0
31429%

The Train Precision of the DecisionTreeClassifier for depth = 50 is 0.8
96766

The Train Recall of the DecisionTreeClassifier for depth = 50 is 0.9737
03

The Train F1-Score of the DecisionTreeClassifier for depth = 50 is 0.93
3652

In [26]:
```python
# Code for drawing seaborn heatmaps
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(Y_train, predictions1), inde
x=class_names, columns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```

## Confusion Matrix



**[5.1.1] Top 20 important features from <span style="color:red">SET 1</span>**

```
In [27]:  # Calculate feature importances from decision trees
          importances = dt.feature_importances_

          # Sort feature importances in descending order
          indices = np.argsort(importances)[::-1][:20]

          # Rearrange feature names so they match the sorted feature importances
```

```python
names = count_vect.get_feature_names()

sns.set(rc={'figure.figsize':(11.7,8.27)})

# Create plot
plt.figure()

# Create plot title
plt.title("Feature Importance")

# Add bars
plt.bar(range(20), importances[indices])

# Add feature names as x-axis labels
names = np.array(names)
plt.xticks(range(20), names[indices], rotation=90)

# Show plot
plt.show()
# uni_gram.get_feature_names()
```

Feature Importance

[5.1.2] Graphviz visualization of Decision Tree on BOW, SET 1

```python
In [30]:   # Please write all the code with proper documentation
           from sklearn import tree
           import pydotplus
           from IPython.display import Image
           from IPython.display import SVG
           from graphviz import Source
           from IPython.display import display
```

```
target = ['negative','positive']
# Create DOT data
data = tree.export_graphviz(dt,out_file=None,max_depth=2,feature_names=
names,class_names=target,filled=True,rounded=True,special_characters=Tr
ue)

# Draw graph
graph = pydotplus.graph_from_dot_data(data)

# Show graph
Image(graph.create_png())
```

Out[30]:



**[5.2] Applying Decision Trees on TFIDF, SET 2**

```python
In [31]: tf_idf_vect = TfidfVectorizer(min_df=1000)
         X_train_vec = tf_idf_vect.fit_transform(X_train)
         X_test_vec = tf_idf_vect.transform(X_test)
         print("the type of count vectorizer :",type(X_train_vec))
         print("the shape of out text TFIDF vectorizer : ",X_train_vec.get_shape
         ())
         print("the number of unique words :", X_train_vec.get_shape()[1])

         # Data-preprocessing: Standardizing the data
         sc = StandardScaler(with_mean=False)
         X_train_vec_standardized = sc.fit_transform(X_train_vec)
         X_test_vec_standardized = sc.transform(X_test_vec)
```

```
the type of count vectorizer : <class 'scipy.sparse.csr.csr_matrix'>
the shape of out text TFIDF vectorizer :  (35000, 216)
the number of unique words : 216
```

```python
In [32]: from sklearn.preprocessing import StandardScaler
         sc = StandardScaler(with_mean=False)
         X_train_vec_standardized = sc.fit_transform(X_train_vec)
         X_test_vec_standardized = sc.transform(X_test_vec)
```

```python
In [34]: # Importing libraries
         from sklearn.tree import DecisionTreeClassifier
         from sklearn.model_selection import GridSearchCV
         from sklearn.metrics import accuracy_score,confusion_matrix,f1_score,pr
         ecision_score,recall_score

         param_grid = {'max_depth': [1, 5, 10, 50, 100, 500, 1000],'min_samples_
         split':[5,10,100,500]}
         model = GridSearchCV(DecisionTreeClassifier(), param_grid, scoring = 'r
         oc_auc', cv=3 , n_jobs = -1,pre_dispatch=2)
         model.fit(X_train_vec_standardized, Y_train)
         print("Model with best parameters :\n",model.best_estimator_)
         print("Accuracy of the model : ",model.score(X_test_vec_standardized, Y
         _test))

         # Cross-Validation errors
```

```python
cv_errors = [1-i for i in model.cv_results_['mean_test_score']]
training_scores=[1-i for i in model.cv_results_['mean_train_score']]

# Optimal value of depth
optimal_depth = model.best_estimator_.max_depth
print("The optimal value of depth is : ",optimal_depth)


optimal_split = model.best_estimator_.min_samples_split
print("The optimal number of base learners is : ",optimal_split)
```

```
Model with best parameters :
 DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=
50,
            max_features=None, max_leaf_nodes=None,
            min_impurity_decrease=0.0, min_impurity_split=None,
            min_samples_leaf=1, min_samples_split=500,
            min_weight_fraction_leaf=0.0, presort=False, random_state=N
one,
            splitter='best')
Accuracy of the model :  0.7788415116532476
The optimal value of depth is :  50
The optimal number of base learners is :  500
```

In [35]:
```python
max_depths = [1, 5, 10, 50, 100, 500, 1000]

train_results = []
test_results = []
for max_depth in max_depths:
    dt = DecisionTreeClassifier(max_depth=max_depth)
    dt.fit(X_train_vec_standardized,Y_train)

    train_pred = dt.predict(X_train_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
rain, train_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous train results
    train_results.append(roc_auc)
```

```python
    Y_pred = dt.predict(X_test_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
est, Y_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous test results
    test_results.append(roc_auc)

from matplotlib.legend_handler import HandlerLine2D

line1, = plt.plot(max_depths, train_results, 'b', label='Train AUC')
line2, = plt.plot(max_depths, test_results, 'r', label='Test AUC')

plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})

plt.ylabel('AUC score')
plt.xlabel('Tree depth')
plt.show()
```

In [36]:
```python
min_samples_splits = [5,10,100,500]

train_results1 = []
test_results1 = []
for min_samples_split in min_samples_splits:
    dt = DecisionTreeClassifier(min_samples_split=min_samples_split)
    dt.fit(X_train_vec_standardized,Y_train)

    train_pred = dt.predict(X_train_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_train, train_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
```

```python
    # Add auc score to previous train results
    train_results1.append(roc_auc)

    Y_pred = dt.predict(X_test_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
est, Y_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous test results
    test_results1.append(roc_auc)

from matplotlib.legend_handler import HandlerLine2D

line1, = plt.plot(min_samples_splits, train_results1, 'b', label='Train
 AUC')
line2, = plt.plot(min_samples_splits, test_results1, 'r', label='Test A
UC')

plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})

plt.ylabel('AUC score')
plt.xlabel('min samle splits')
plt.show()
```

In [37]:
```python
# DecisionTreeClassifier with Optimal value of depth
dt = DecisionTreeClassifier(max_depth=optimal_depth,min_samples_split=optimal_split)
dt.fit(X_train_vec_standardized,Y_train)
predictions = dt.predict(X_test_vec_standardized)
predictions1 = dt.predict(X_train_vec_standardized)

# Variables that will be used for  making table in Conclusion part of this assignment
tfidf_depth = optimal_depth
tfidf_split = optimal_split
tfidf_train_acc = model.score(X_test_vec_standardized, Y_test)*100
tfidf_test_acc = accuracy_score(Y_test, predictions) * 100
```

```
In [38]: train_fpr, train_tpr, thresholds = roc_curve(Y_train, dt.predict_proba(
         X_train_vec_standardized)[:,1])
         test_fpr, test_tpr, thresholds = roc_curve(Y_test, dt.predict_proba(X_t
         est_vec_standardized)[:,1])

         plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
         rain_tpr)))
         plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
         tpr)))
         plt.legend()
         plt.xlabel("depth: hyperparameter")
         plt.ylabel("AUC")
         plt.title("ERROR PLOTS")
         plt.show()
```

## ERROR PLOTS



Legend:
- train AUC = 0.8847911397842253
- test AUC = 0.7795747542324558

Y-axis: AUC
X-axis: depth: hyperparameter

In [39]:
```python
# evaluate accuracy on test data
acc = accuracy_score(Y_test, predictions) * 100
print('\nThe Test Accuracy of the DecisionTreeClassifier for depth = %d
 is %f%%' % (optimal_depth, acc))

# evaluate precision
acc = precision_score(Y_test, predictions, pos_label = 1)
print('\nThe Test Precision of the DecisionTreeClassifier for depth = %
d is %f' % (optimal_depth, acc))

# evaluate recall
acc = recall_score(Y_test, predictions, pos_label = 1)
```

```python
print('\nThe Test Recall of the DecisionTreeClassifier for depth = %d i
s %f' % (optimal_depth, acc))

# evaluate f1-score
acc = f1_score(Y_test, predictions, pos_label = 1)
print('\nThe Test F1-Score of the DecisionTreeClassifier for depth = %d
 is %f' % (optimal_depth, acc))
```

The Test Accuracy of the DecisionTreeClassifier for depth = 50 is 82.90
0000%

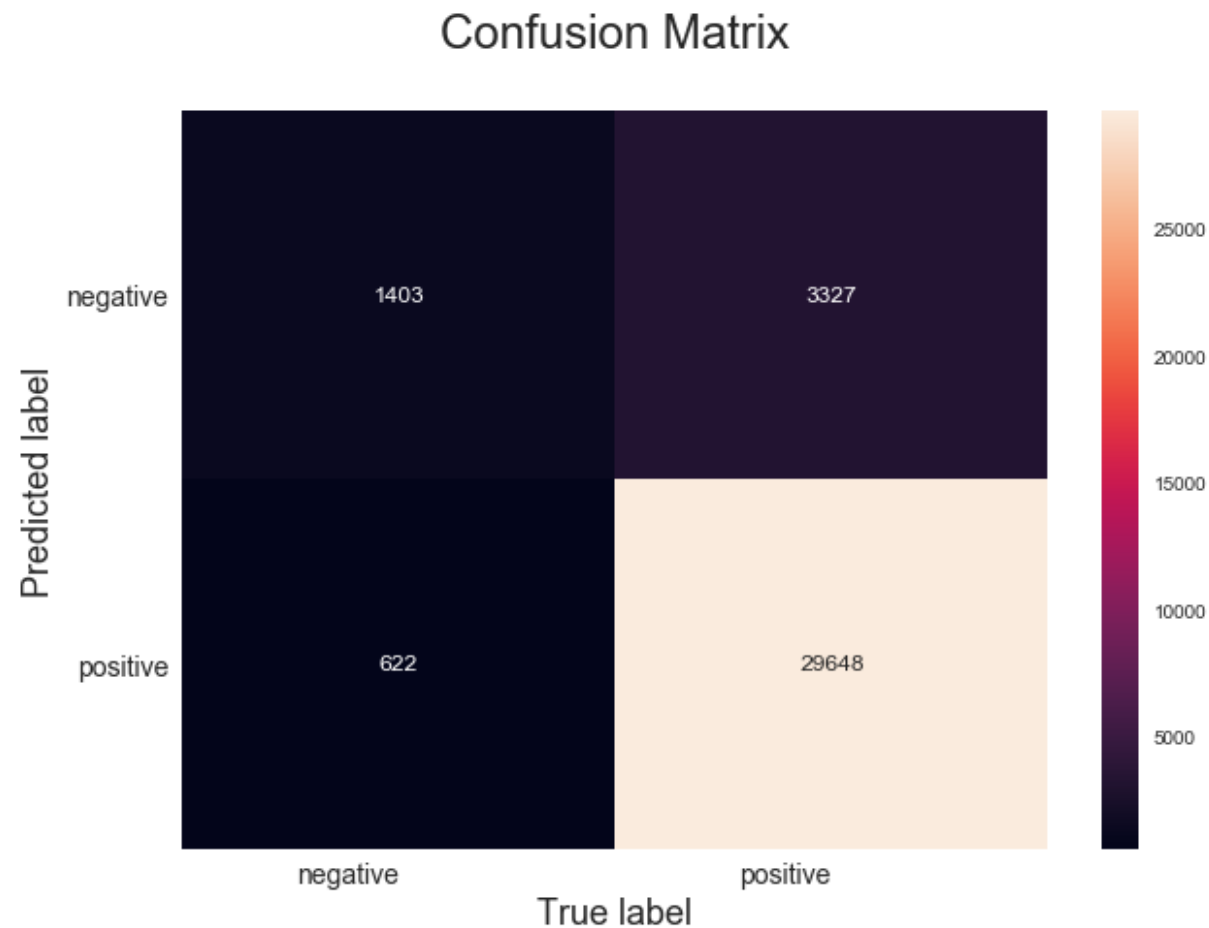The Test Precision of the DecisionTreeClassifier for depth = 50 is 0.84
7128

The Test Recall of the DecisionTreeClassifier for depth = 50 is 0.96511
8

The Test F1-Score of the DecisionTreeClassifier for depth = 50 is 0.902
282

In [40]:
```python
# Code for drawing seaborn heatmaps on test data
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(Y_test, predictions), index=
class_names, columns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```

## Confusion Matrix



```
In [41]:  # evaluate accuracy on train data
          acc = accuracy_score(Y_train, predictions1) * 100
          print('\nThe Train Accuracy of the DecisionTreeClassifier for depth = %
          d is %f%%' % (optimal_depth, acc))

          # evaluate precision
          acc = precision_score(Y_train, predictions1, pos_label = 1)
          print('\nThe Train Precision of the DecisionTreeClassifier for depth =
          %d is %f' % (optimal_depth, acc))
```

```python
# evaluate recall
acc = recall_score(Y_train, predictions1, pos_label = 1)
print('\nThe Train Recall of the DecisionTreeClassifier for depth = %d
 is %f' % (optimal_depth, acc))

# evaluate f1-score
acc = f1_score(Y_train, predictions1, pos_label = 1)
print('\nThe Train F1-Score of the DecisionTreeClassifier for depth = %
d is %f' % (optimal_depth, acc))
```

The Train Accuracy of the DecisionTreeClassifier for depth = 50 is 88.7
17143%

The Train Precision of the DecisionTreeClassifier for depth = 50 is 0.8
99105

The Train Recall of the DecisionTreeClassifier for depth = 50 is 0.9794
52

The Train F1-Score of the DecisionTreeClassifier for depth = 50 is 0.93
7560

In [42]:
```python
# Code for drawing seaborn heatmaps on test data
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(Y_train, predictions1), inde
x=class_names, columns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```

## Confusion Matrix



**[5.2.1] Top 20 important features from** <span style="color:red">**SET 2**</span>

```python
In [43]:  # Please write all the code with proper documentation
          # Calculate feature importances from decision trees
          importances = dt.feature_importances_

          # Sort feature importances in descending order
          indices = np.argsort(importances)[::-1][:20]
```

```python
# Rearrange feature names so they match the sorted feature importances
names = tf_idf_vect .get_feature_names()

sns.set(rc={'figure.figsize':(11.7,8.27)})

# Create plot
plt.figure()

# Create plot title
plt.title("Feature Importance")

# Add bars
plt.bar(range(20), importances[indices])

# Add feature names as x-axis labels
names = np.array(names)
plt.xticks(range(20), names[indices], rotation=90)

# Show plot
plt.show()
# uni_gram.get_feature_names()
```

Feature Importance

```
#print('other way of feature imp')
#print('feature importance')
#print('clf=DecisionTreeClassifier(max_depth=optimal_depth)\
#clf.fit(x_train_bow,y_train)\
#print('top 25 words and their IG---')\
#print('-------------------------------')\
#top=clf.feature_importances_\
#s=np.argsort(top)[-25:]\
#feature=count_vec.get_feature_names()\
#for i in range(25):\
```

```
#    index=s[i]\
#    print(feature[index],'\t\t:\t\t',round(top[index],5))')
```

**[5.2.2] Graphviz visualization of Decision Tree on TFIDF, <span style="color:red">SET 2</span>**

In [45]:
```python
# Please write all the code with proper documentation
target = ['negative','positive']
# Create DOT data
data = tree.export_graphviz(dt,max_depth=2,out_file=None,feature_names=
names,class_names=target,filled=True,rounded=True,special_characters=Tr
ue)

# Draw graph
graph = pydotplus.graph_from_dot_data(data)

# Show graph
Image(graph.create_png())
```

Out[45]:

```
In [46]:  # saving graph
          #import graphviz
          #from sklearn import tree
          #clf=DecisionTreeClassifier(max_depth=13)
          #clf.fit(x_train_bow,y_train)
          #dot_data = tree.export_graphviz(clf, out_file=None)
          #graph = graphviz.Source(dot_data)
          #graph.render("graph_bow.pdf")
```

## [5.3] Applying Decision Trees on AVG W2V, SET 3

```
In [47]:  # Please write all the code with proper documentation
```

```python
# List of sentence in X_train text
sent_of_train=[]
for sent in X_train:
    sent_of_train.append(sent.split())

# List of sentence in X_est text
sent_of_test=[]
for sent in X_test:
    sent_of_test.append(sent.split())

# Train your own Word2Vec model using your own train text corpus
# min_count = 5 considers only words that occured atleast 5 times
w2v_model=Word2Vec(sent_of_train,min_count=5,size=50, workers=4)

w2v_words = list(w2v_model.wv.vocab)
print("number of words that occured minimum 5 times ",len(w2v_words))
```

```
number of words that occured minimum 5 times  8359
```

In [48]:
```python
# compute average word2vec for each review for X_train .
train_vectors = [];
for sent in sent_of_train:
    sent_vec = np.zeros(50)
    cnt_words =0;
    for word in sent: #
        if word in w2v_words:
            vec = w2v_model.wv[word]
            sent_vec += vec
            cnt_words += 1
    if cnt_words != 0:
        sent_vec /= cnt_words
    train_vectors.append(sent_vec)

# compute average word2vec for each review for X_test .
test_vectors = [];
for sent in sent_of_test:
    sent_vec = np.zeros(50)
    cnt_words =0;
    for word in sent: #
```

```python
            if word in w2v_words:
                vec = w2v_model.wv[word]
                sent_vec += vec
                cnt_words += 1
        if cnt_words != 0:
            sent_vec /= cnt_words
        test_vectors.append(sent_vec)


import warnings
warnings.filterwarnings('ignore')
# Data-preprocessing: Standardizing the data
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train_vec_standardized = sc.fit_transform(train_vectors)
X_test_vec_standardized = sc.transform(test_vectors)
```

In [49]:
```python
# Importing libraries
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score,confusion_matrix,f1_score,pr
ecision_score,recall_score

param_grid = {'max_depth': [1, 5, 10, 50, 100, 500, 1000],'min_samples_
split':[5,10,100,500]}
model = GridSearchCV(DecisionTreeClassifier(), param_grid, scoring = 'r
oc_auc', cv=3 , n_jobs = -1,pre_dispatch=2)
model.fit(X_train_vec_standardized, Y_train)
print("Model with best parameters :\n",model.best_estimator_)
print("Accuracy of the model : ",model.score(X_test_vec_standardized, Y
_test))

# Cross-Validation errors
cv_errors = [1-i for i in model.cv_results_['mean_test_score']]
training_scores=[1-i for i in model.cv_results_['mean_train_score']]

# Optimal value of depth
optimal_depth = model.best_estimator_.max_depth
print("The optimal value of depth is : ",optimal_depth)
```

```
optimal_split = model.best_estimator_.min_samples_split
print("The optimal number of base learners is : ",optimal_split)
```

```
Model with best parameters :
 DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=
10,
            max_features=None, max_leaf_nodes=None,
            min_impurity_decrease=0.0, min_impurity_split=None,
            min_samples_leaf=1, min_samples_split=500,
            min_weight_fraction_leaf=0.0, presort=False, random_state=N
one,
            splitter='best')
Accuracy of the model :  0.7771645455875285
The optimal value of depth is :   10
The optimal number of base learners is :  500
```

In [50]:
```python
max_depths = [1, 5, 10, 50, 100, 500, 1000]

train_results = []
test_results = []
for max_depth in max_depths:
    dt = DecisionTreeClassifier(max_depth=max_depth)
    dt.fit(X_train_vec_standardized,Y_train)

    train_pred = dt.predict(X_train_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
rain, train_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous train results
    train_results.append(roc_auc)

    Y_pred = dt.predict(X_test_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
est, Y_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous test results
    test_results.append(roc_auc)
```
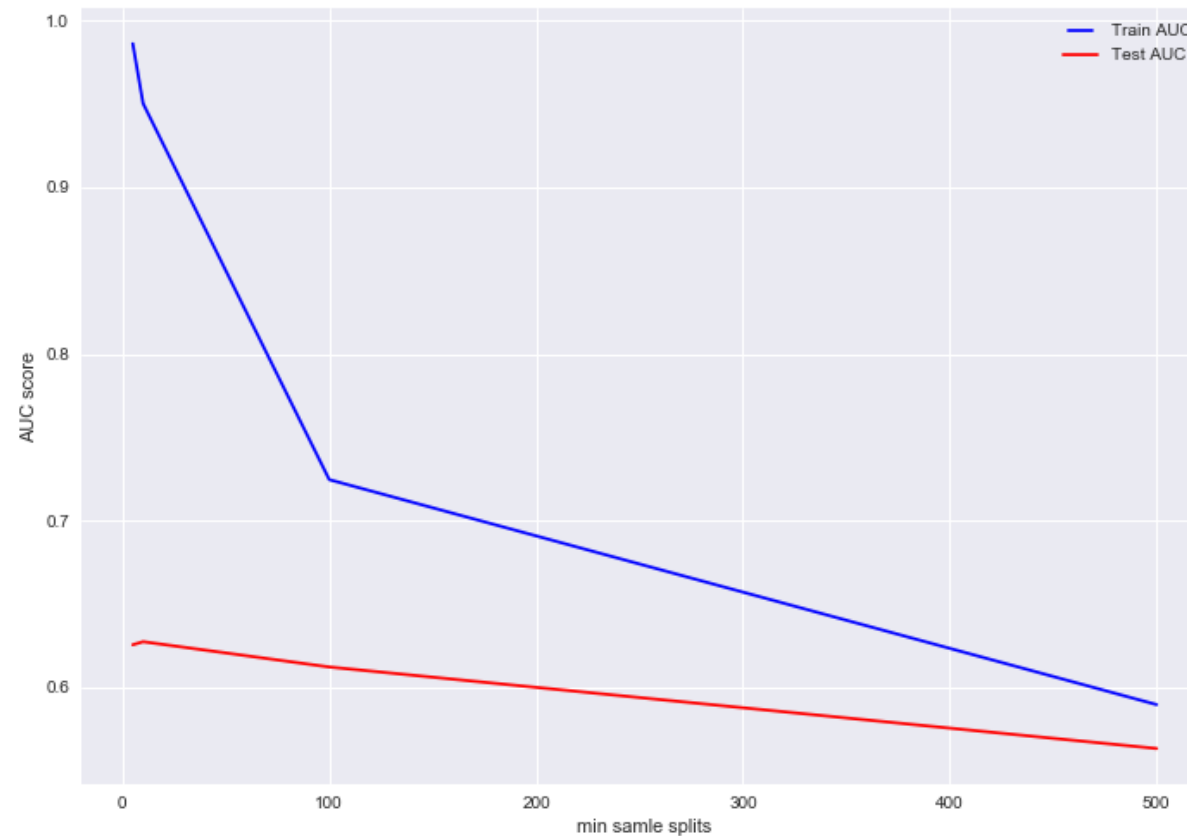
```python
from matplotlib.legend_handler import HandlerLine2D

line1, = plt.plot(max_depths, train_results, 'b', label='Train AUC')
line2, = plt.plot(max_depths, test_results, 'r', label='Test AUC')

plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})

plt.ylabel('AUC score')
plt.xlabel('Tree depth')
plt.show()
```
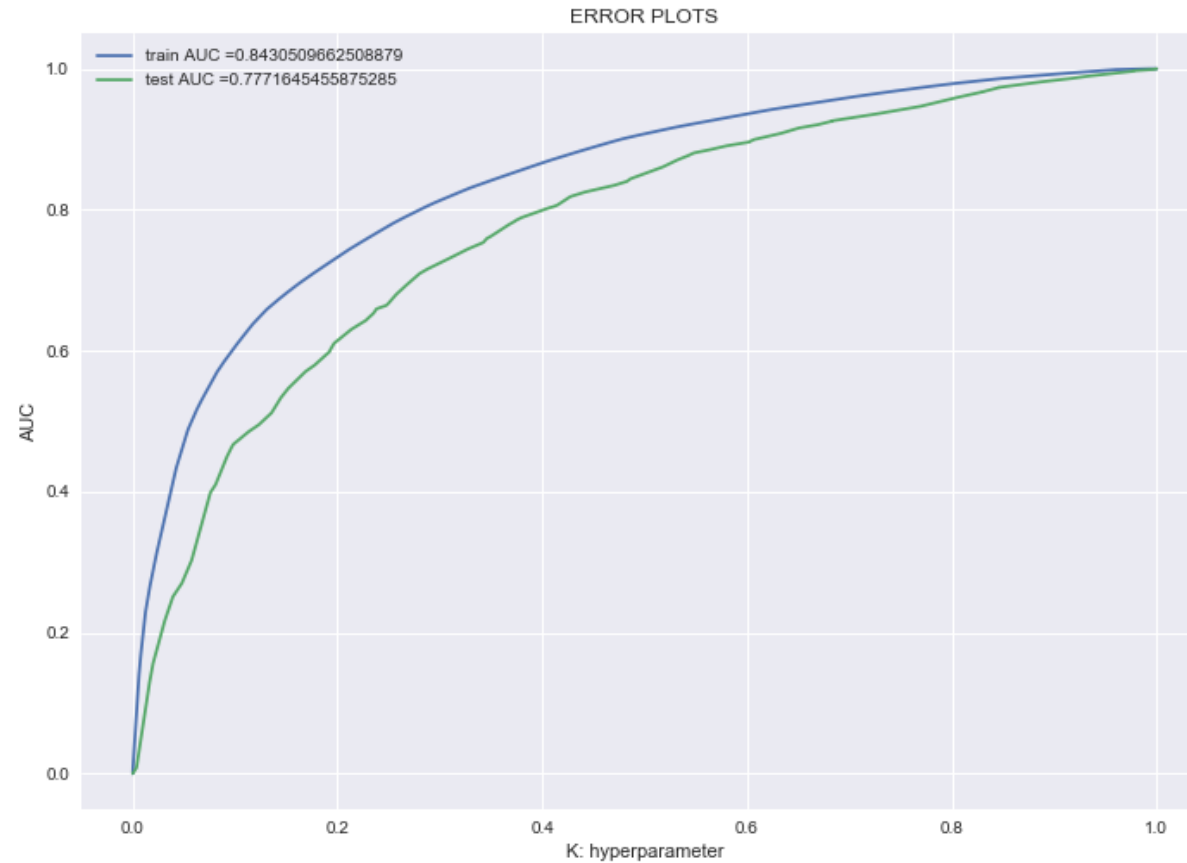


```python
In [51]: min_samples_splits = [5,10,100,500]
```

```python
train_results1 = []
test_results1 = []
for min_samples_split in min_samples_splits:
    dt = DecisionTreeClassifier(min_samples_split=min_samples_split)
    dt.fit(X_train_vec_standardized,Y_train)

    train_pred = dt.predict(X_train_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
rain, train_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous train results
    train_results1.append(roc_auc)

    Y_pred = dt.predict(X_test_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
est, Y_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous test results
    test_results1.append(roc_auc)

from matplotlib.legend_handler import HandlerLine2D

line1, = plt.plot(min_samples_splits, train_results1, 'b', label='Train
 AUC')
line2, = plt.plot(min_samples_splits, test_results1, 'r', label='Test A
UC')

plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})

plt.ylabel('AUC score')
plt.xlabel('min samle splits')
plt.show()
```

```
In [52]:  # DecisionTreeClassifier with Optimal value of depth
          dt = DecisionTreeClassifier(max_depth=optimal_depth,min_samples_split=o
          ptimal_split)
          dt.fit(X_train_vec_standardized,Y_train)
          predictions = dt.predict(X_test_vec_standardized)
          predictions1 = dt.predict(X_train_vec_standardized)

          # Variables that will be used for  making table in Conclusion part of t
          his assignment
          avg_w2v_depth = optimal_depth
          avg_w2v_split = optimal_split
          avg_w2v_train_acc = model.score(X_test_vec_standardized, Y_test)*100
          avg_w2v_test_acc = accuracy_score(Y_test, predictions) * 100
```

```python
In [53]:  train_fpr, train_tpr, thresholds = roc_curve(Y_train, dt.predict_proba(
          X_train_vec_standardized)[:,1])
          test_fpr, test_tpr, thresholds = roc_curve(Y_test, dt.predict_proba(X_t
          est_vec_standardized)[:,1])

          plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
          rain_tpr)))
          plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
          tpr)))
          plt.legend()
          plt.xlabel("K: hyperparameter")
          plt.ylabel("AUC")
          plt.title("ERROR PLOTS")
          plt.show()
```

## ERROR PLOTS



Legend:
- train AUC =0.8430509662508879
- test AUC =0.7771645455875285

Y-axis: AUC
X-axis: K: hyperparameter

In [54]:
```python
# evaluate accuracy on test data
acc = accuracy_score(Y_test, predictions) * 100
print('\nThe Test Accuracy of the DecisionTreeClassifier for depth = %d
 is %f%%' % (optimal_depth, acc))

# evaluate precision
acc = precision_score(Y_test, predictions, pos_label = 1)
print('\nThe Test Precision of the DecisionTreeClassifier for depth = %
d is %f' % (optimal_depth, acc))

# evaluate recall
acc = recall_score(Y_test, predictions, pos_label = 1)
```

```python
print('\nThe Test Recall of the DecisionTreeClassifier for depth = %d i
s %f' % (optimal_depth, acc))

# evaluate f1-score
acc = f1_score(Y_test, predictions, pos_label = 1)
print('\nThe Test F1-Score of the DecisionTreeClassifier for depth = %d
 is %f' % (optimal_depth, acc))
```

The Test Accuracy of the DecisionTreeClassifier for depth = 10 is 82.42
6667%

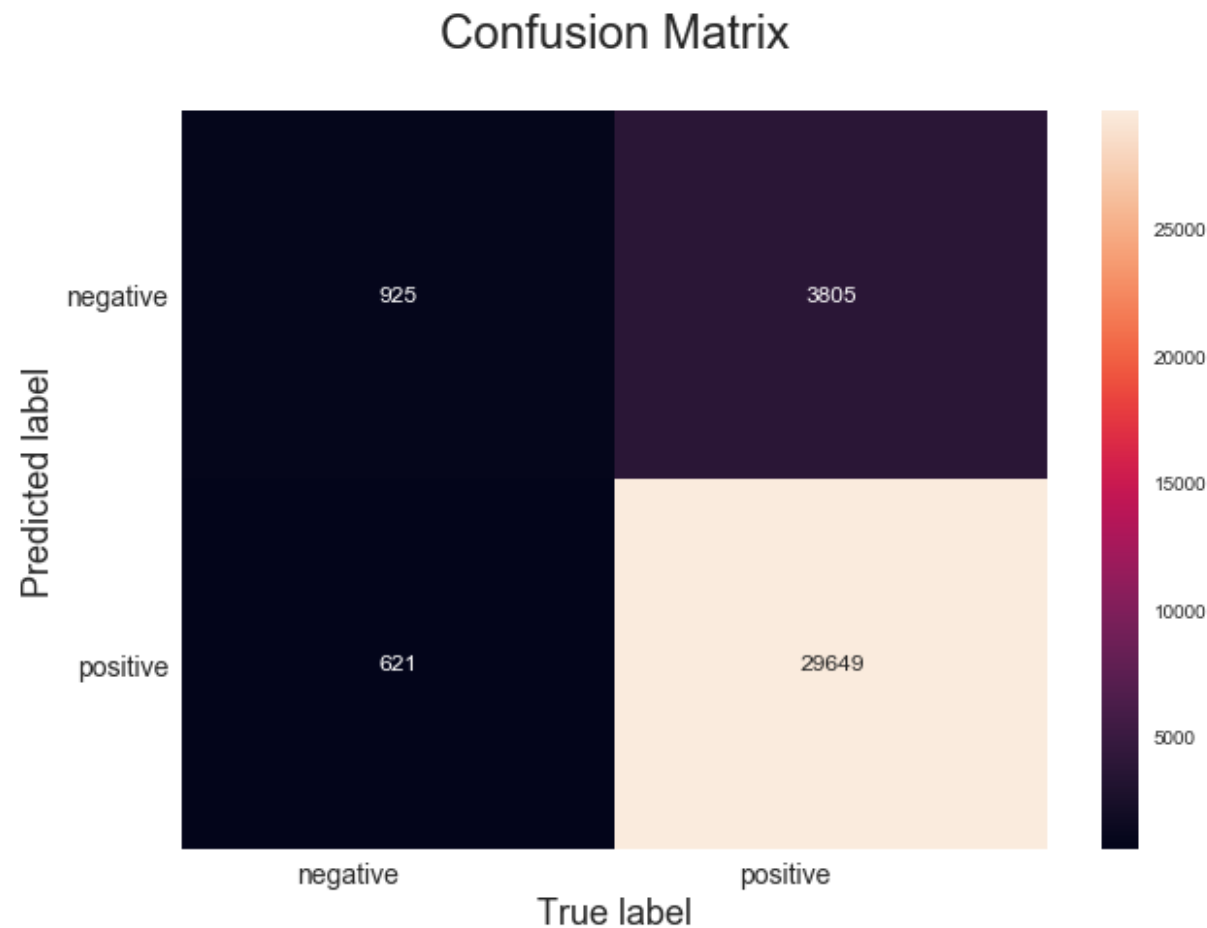The Test Precision of the DecisionTreeClassifier for depth = 10 is 0.83
7940

The Test Recall of the DecisionTreeClassifier for depth = 10 is 0.97343
1

The Test F1-Score of the DecisionTreeClassifier for depth = 10 is 0.900
618

In [55]:
```python
# Code for drawing seaborn heatmaps on test data
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(Y_test, predictions), index=
class_names, columns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```

## Confusion Matrix



```
In [56]:  # evaluate accuracy on train data
          acc = accuracy_score(Y_train, predictions1) * 100
          print('\nThe Train Accuracy of the DecisionTreeClassifier for depth = %
          d is %f%%' % (optimal_depth, acc))

          # evaluate precision
          acc = precision_score(Y_train, predictions1, pos_label = 1)
          print('\nThe Train Precision of the DecisionTreeClassifier for depth =
          %d is %f' % (optimal_depth, acc))
```

```python
# evaluate recall
acc = recall_score(Y_train, predictions1, pos_label = 1)
print('\nThe Train Recall of the DecisionTreeClassifier for depth = %d
 is %f' % (optimal_depth, acc))

# evaluate f1-score
acc = f1_score(Y_train, predictions1, pos_label = 1)
print('\nThe Train F1-Score of the DecisionTreeClassifier for depth = %
d is %f' % (optimal_depth, acc))
```

The Train Accuracy of the DecisionTreeClassifier for depth = 10 is 87.3
54286%

The Train Precision of the DecisionTreeClassifier for depth = 10 is 0.8
86262

The Train Recall of the DecisionTreeClassifier for depth = 10 is 0.9794
85

The Train F1-Score of the DecisionTreeClassifier for depth = 10 is 0.93
0544

In [57]:
```python
# Code for drawing seaborn heatmaps
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(Y_train, predictions1), inde
x=class_names, columns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```

## Confusion Matrix



## [5.4] Applying Decision Trees on TFIDF W2V, SET 4

```
In [58]:  # Please write all the code with proper documentation
          # TF-IDF weighted Word2Vec
          tf_idf_vect = TfidfVectorizer()

          # final_tf_idf1 is the sparse matrix with row= sentence, col=word and c
          ell_val = tfidf
          final_tf_idf1 = tf_idf_vect.fit_transform(X_train)
```

```python
# tfidf words/col-names
tfidf_feat = tf_idf_vect.get_feature_names()

# compute TFIDF Weighted Word2Vec for each review for X_test .
tfidf_test_vectors = [];
row=0;
for sent in sent_of_test:
    sent_vec = np.zeros(50)
    weight_sum =0;
    for word in sent:
        if word in w2v_words:
            vec = w2v_model.wv[word]
            # obtain the tf_idfidf of a word in a sentence/review
            tf_idf = final_tf_idf1[row, tfidf_feat.index(word)]
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_test_vectors.append(sent_vec)
    row += 1
```

In [59]:
```python
# compute TFIDF Weighted Word2Vec for each review for X_train .
tfidf_train_vectors = [];
row=0;
for sent in sent_of_train:
    sent_vec = np.zeros(50)
    weight_sum =0;
    for word in sent:
        if word in w2v_words:
            vec = w2v_model.wv[word]
            # obtain the tf_idfidf of a word in a sentence/review
            tf_idf = final_tf_idf1[row, tfidf_feat.index(word)]
            sent_vec += (vec * tf_idf)
            weight_sum += tf_idf
    if weight_sum != 0:
        sent_vec /= weight_sum
    tfidf_train_vectors.append(sent_vec)
    row += 1
```

```python
# Data-preprocessing: Standardizing the data
sc = StandardScaler()
X_train_vec_standardized = sc.fit_transform(tfidf_train_vectors)
X_test_vec_standardized = sc.transform(tfidf_test_vectors)
```

In [60]:
```python
# Importing libraries
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import accuracy_score,confusion_matrix,f1_score,pr
ecision_score,recall_score

param_grid = {'max_depth': [1, 5, 10, 50, 100, 500, 1000],'min_samples_
split':[5,10,100,500]}
model = GridSearchCV(DecisionTreeClassifier(), param_grid, scoring = 'r
oc_auc', cv=3 , n_jobs = -1,pre_dispatch=2)
model.fit(X_train_vec_standardized, Y_train)
print("Model with best parameters :\n",model.best_estimator_)
print("Accuracy of the model : ",model.score(X_test_vec_standardized, Y
_test))

# Cross-Validation errors
cv_errors = [1-i for i in model.cv_results_['mean_test_score']]
training_scores=[1-i for i in model.cv_results_['mean_train_score']]

# Optimal value of depth
optimal_depth = model.best_estimator_.max_depth
print("The optimal value of depth is : ",optimal_depth)


optimal_split = model.best_estimator_.min_samples_split
print("The optimal value of depth is : ",optimal_split)
```

```
Model with best parameters :
 DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=
10,
            max_features=None, max_leaf_nodes=None,
            min_impurity_decrease=0.0, min_impurity_split=None,
            min_samples_leaf=1, min_samples_split=500,
```

```
              min_weight_fraction_leaf=0.0, presort=False, random_state=N
one,
              splitter='best')
Accuracy of the model :  0.5610866910866912
The optimal value of depth is :  10
The optimal value of depth is :  500
```

In [61]:
```python
max_depths = [1, 5, 10, 50, 100, 500, 1000]

train_results = []
test_results = []
for max_depth in max_depths:
    dt = DecisionTreeClassifier(max_depth=max_depth)
    dt.fit(X_train_vec_standardized,Y_train)

    train_pred = dt.predict(X_train_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
rain, train_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous train results
    train_results.append(roc_auc)

    Y_pred = dt.predict(X_test_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
est, Y_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous test results
    test_results.append(roc_auc)

from matplotlib.legend_handler import HandlerLine2D

line1, = plt.plot(max_depths, train_results, 'b', label='Train AUC')
line2, = plt.plot(max_depths, test_results, 'r', label='Test AUC')

plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})

plt.ylabel('AUC score')
```
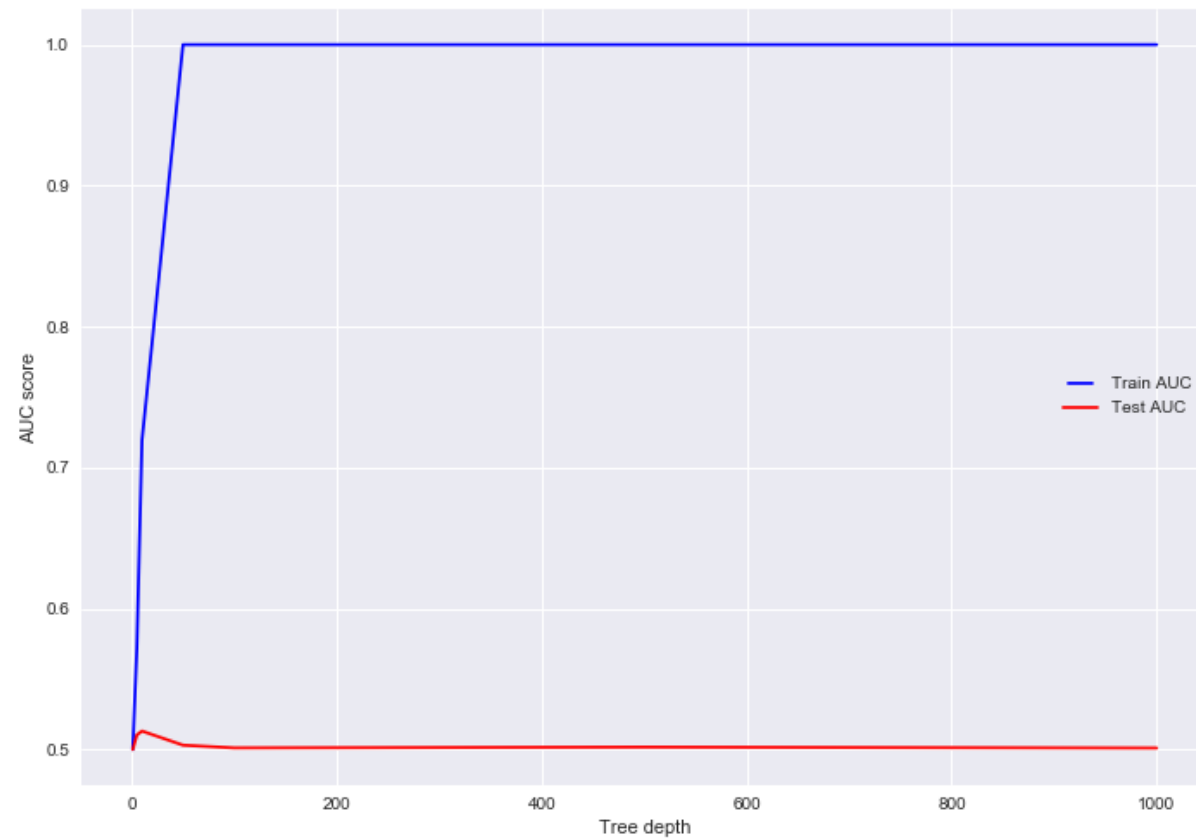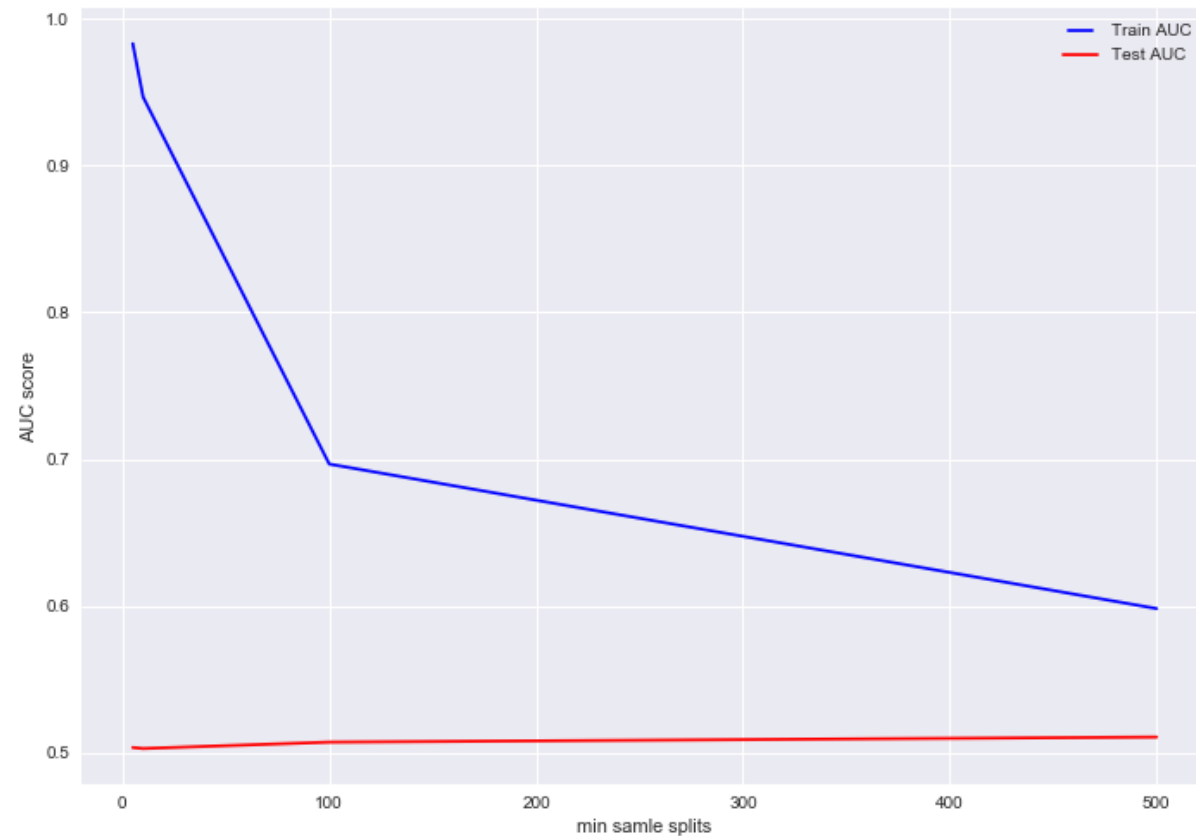
```
plt.xlabel('Tree depth')
plt.show()
```



In [62]:
```
min_samples_splits = [5,10,100,500]

train_results1 = []
test_results1 = []
for min_samples_split in min_samples_splits:
    dt = DecisionTreeClassifier(min_samples_split=min_samples_split)
    dt.fit(X_train_vec_standardized,Y_train)

    train_pred = dt.predict(X_train_vec_standardized)
```

```python
    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
rain, train_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous train results
    train_results1.append(roc_auc)

    Y_pred = dt.predict(X_test_vec_standardized)

    false_positive_rate, true_positive_rate, thresholds = roc_curve(Y_t
est, Y_pred)
    roc_auc = auc(false_positive_rate, true_positive_rate)
    # Add auc score to previous test results
    test_results1.append(roc_auc)

from matplotlib.legend_handler import HandlerLine2D

line1, = plt.plot(min_samples_splits, train_results1, 'b', label='Train
 AUC')
line2, = plt.plot(min_samples_splits, test_results1, 'r', label='Test A
UC')

plt.legend(handler_map={line1: HandlerLine2D(numpoints=2)})

plt.ylabel('AUC score')
plt.xlabel('min samle splits')
plt.show()
```
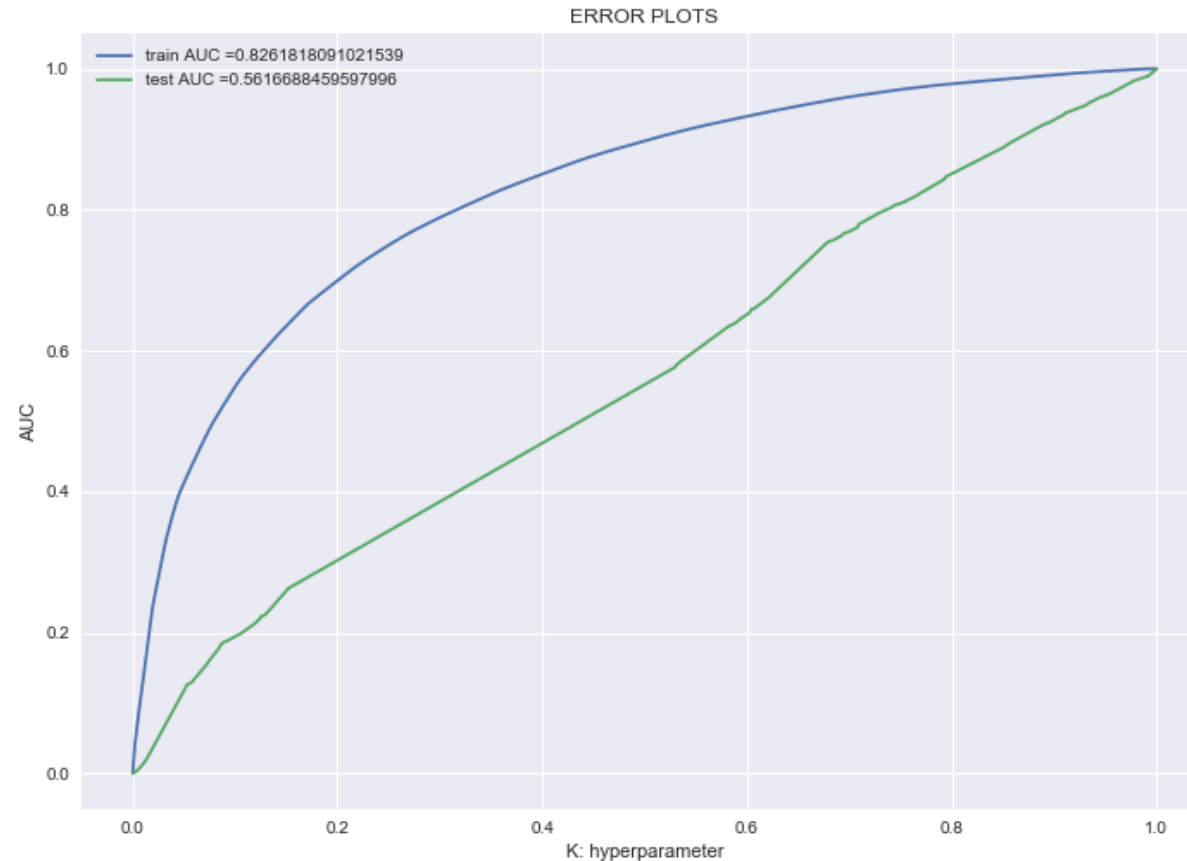
```
In [63]:  # DecisionTreeClassifier with Optimal value of depth
          dt = DecisionTreeClassifier(max_depth=optimal_depth,min_samples_split=o
          ptimal_split)
          dt.fit(X_train_vec_standardized,Y_train)
          predictions = dt.predict(X_test_vec_standardized)
          predictions1 = dt.predict(X_train_vec_standardized)

          # Variables that will be used for  making table in Conclusion part of t
          his assignment
          tfidf_w2v_depth = optimal_depth
          tfidf_w2v_split = optimal_split
          tfidf_w2v_train_acc = model.score(X_test_vec_standardized, Y_test)*100
          tfidf_w2v_test_acc = accuracy_score(Y_test, predictions) * 100
```

```python
In [64]: train_fpr, train_tpr, thresholds = roc_curve(Y_train, dt.predict_proba(
         X_train_vec_standardized)[:,1])
         test_fpr, test_tpr, thresholds = roc_curve(Y_test, dt.predict_proba(X_t
         est_vec_standardized)[:,1])

         plt.plot(train_fpr, train_tpr, label="train AUC ="+str(auc(train_fpr, t
         rain_tpr)))
         plt.plot(test_fpr, test_tpr, label="test AUC ="+str(auc(test_fpr, test_
         tpr)))
         plt.legend()
         plt.xlabel("K: hyperparameter")
         plt.ylabel("AUC")
         plt.title("ERROR PLOTS")
         plt.show()
```

ERROR PLOTS



In [65]:
```python
# evaluate accuracy on test data
acc = accuracy_score(Y_test, predictions) * 100
print('\nThe Test Accuracy of the DecisionTreeClassifier for depth = %d
 is %f%%' % (optimal_depth, acc))

# evaluate precision
acc = precision_score(Y_test, predictions, pos_label = 1)
print('\nThe Test Precision of the DecisionTreeClassifier for depth = %
d is %f' % (optimal_depth, acc))

# evaluate recall
acc = recall_score(Y_test, predictions, pos_label = 1)
```

```python
print('\nThe Test Recall of the DecisionTreeClassifier for depth = %d i
s %f' % (optimal_depth, acc))

# evaluate f1-score
acc = f1_score(Y_test, predictions, pos_label = 1)
print('\nThe Test F1-Score of the DecisionTreeClassifier for depth = %d
 is %f' % (optimal_depth, acc))
```

The Test Accuracy of the DecisionTreeClassifier for depth = 10 is 78.00
6667%

The Test Precision of the DecisionTreeClassifier for depth = 10 is 0.82
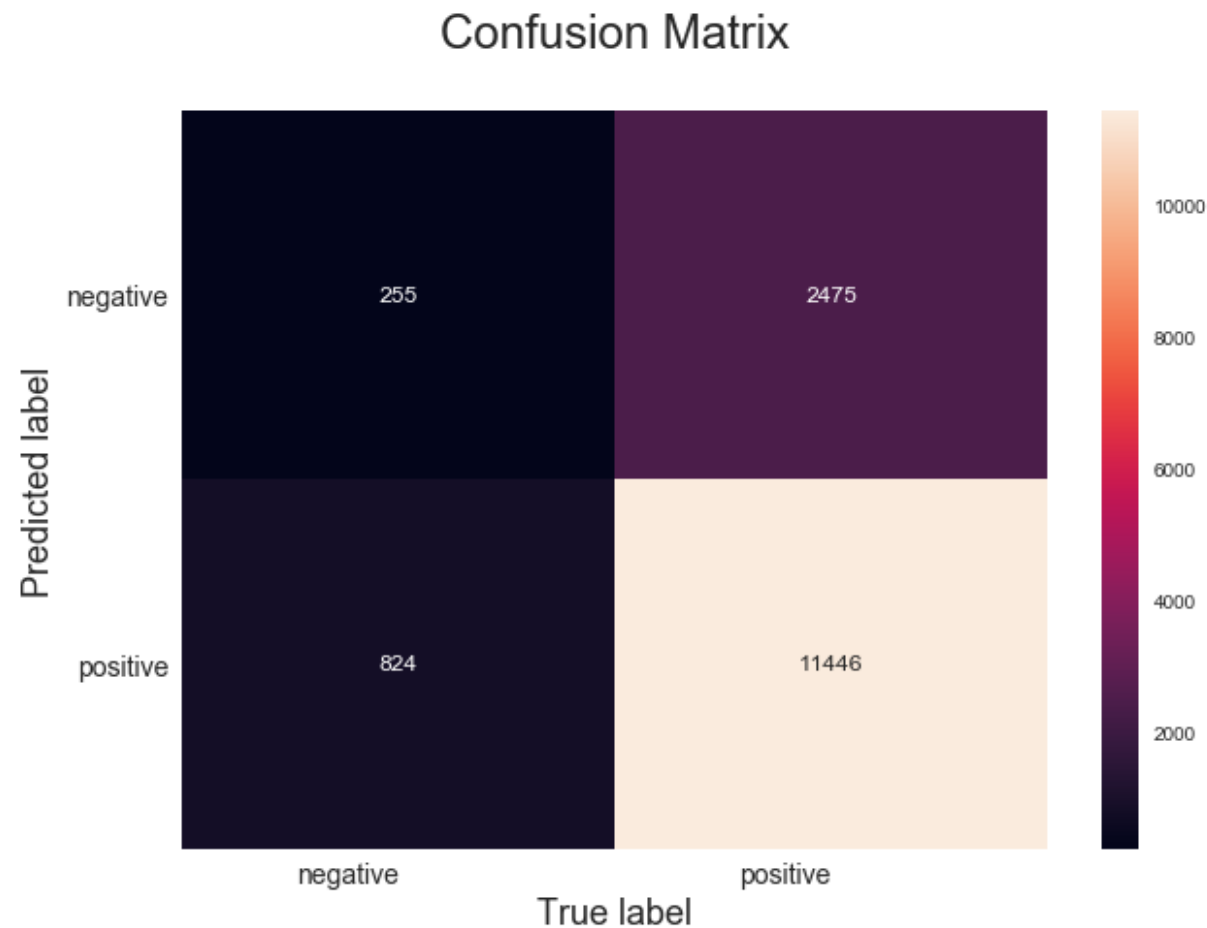2211

The Test Recall of the DecisionTreeClassifier for depth = 10 is 0.93284
4

The Test F1-Score of the DecisionTreeClassifier for depth = 10 is 0.874
041

In [66]:
```python
# Code for drawing seaborn heatmaps on test data
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(Y_test, predictions), index=
class_names, columns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```

## Confusion Matrix



```
In [67]:  # evaluate accuracy on train data
          acc = accuracy_score(Y_train, predictions1) * 100
          print('\nThe Train Accuracy of the DecisionTreeClassifier for depth = %
          d is %f%%' % (optimal_depth, acc))

          # evaluate precision
          acc = precision_score(Y_train, predictions1, pos_label = 1)
          print('\nThe Train Precision of the DecisionTreeClassifier for depth =
          %d is %f' % (optimal_depth, acc))
```

```python
# evaluate recall
acc = recall_score(Y_train, predictions1, pos_label = 1)
print('\nThe Train Recall of the DecisionTreeClassifier for depth = %d
 is %f' % (optimal_depth, acc))

# evaluate f1-score
acc = f1_score(Y_train, predictions1, pos_label = 1)
print('\nThe Train F1-Score of the DecisionTreeClassifier for depth = %
d is %f' % (optimal_depth, acc))
```

The Train Accuracy of the DecisionTreeClassifier for depth = 10 is 87.3
68571%

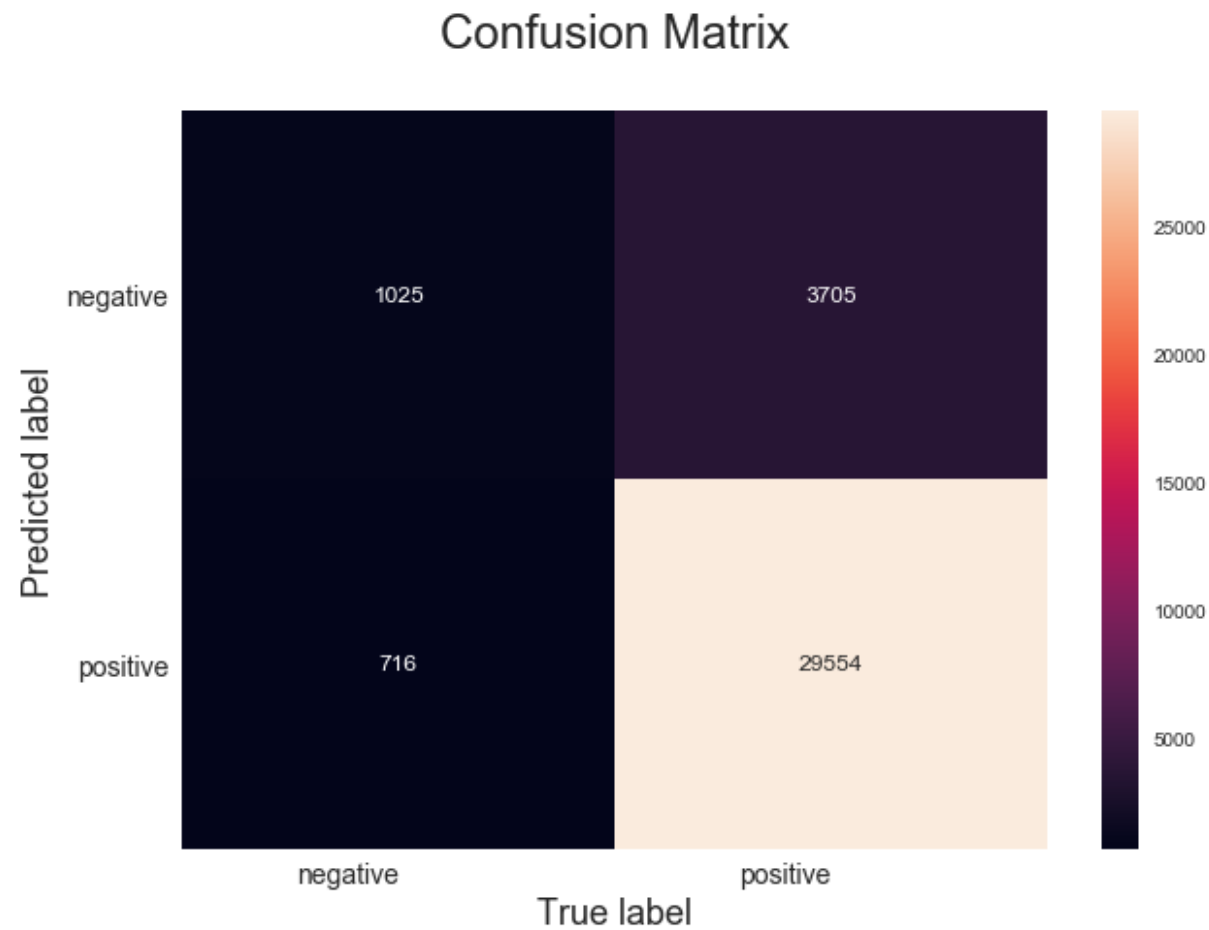The Train Precision of the DecisionTreeClassifier for depth = 10 is 0.8
88602

The Train Recall of the DecisionTreeClassifier for depth = 10 is 0.9763
46

The Train F1-Score of the DecisionTreeClassifier for depth = 10 is 0.93
0410

In [68]:
```python
# Code for drawing seaborn heatmaps
class_names = ['negative','positive']
df_heatmap = pd.DataFrame(confusion_matrix(Y_train, predictions1), inde
x=class_names, columns=class_names )
fig = plt.figure(figsize=(10,7))
heatmap = sns.heatmap(df_heatmap, annot=True, fmt="d")

# Setting tick labels for heatmap
heatmap.yaxis.set_ticklabels(heatmap.yaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
heatmap.xaxis.set_ticklabels(heatmap.xaxis.get_ticklabels(), rotation=0
, ha='right', fontsize=14)
plt.ylabel('Predicted label',size=18)
plt.xlabel('True label',size=18)
plt.title("Confusion Matrix\n",size=24)
plt.show()
```

## Confusion Matrix



## steps

- STEP 1 :- Text Preprocessing
- STEP 2:- Time-based splitting of whole dataset into train_data and test_data
- STEP 3:- Training the vectorizer on train_data and later applying same vectorizer on both train_data and test_data to transform them into vectors

- STEP 4:- Using Decision_Tree as an estimator in GridSearchCV in order to find optimal value of depth of the tree
- STEP 5:- Once , we get optimal value of depth then train Decision_Tree again with this optimal depth and make predictions on test_data
- STEP 6:- Draw Cross-Validation Error vs Depth graph,auc's
- STEP 7 :- Evaluate : Accuracy , F1-Score , Precision , Recall
- STEP 8:- Visualizing the Decision Tree using Graphviz
- STEP 9:- Draw Seaborn Heatmap for Confusion Matrix .

# [6] Conclusions

In [69]:
```python
# Please compare all your models using Prettytable library
# Creating table using PrettyTable library
from prettytable import PrettyTable

# Names of the models
names =['Decision_Tree for BoW','Decision_Tree for TFIDF','Decision_Tree for Avg_Word2Vec', 'Decision_Tree for tfidf_Word2Vec']

# Values of optimal depth
optimal_depth = [bow_depth, tfidf_depth,avg_w2v_depth, tfidf_w2v_depth]

optimal_split = [bow_split, tfidf_split,avg_w2v_split, tfidf_w2v_split]

# Training Accuracies
train_acc = [88.03,88.71,88.35,87.36]

# Test Accuracies
test_acc = [82.90,82.9,82.42,78.00]
numbering = [1,2,3,4]

# Initializing prettytable
ptable = PrettyTable()

# Adding columns
ptable.add_column("S.NO.",numbering)
```

```python
ptable.add_column("MODEL",names)
ptable.add_column("Optimal Depth",optimal_depth)
ptable.add_column("Optimal split",optimal_split)
ptable.add_column("Training Accuracy",train_acc)
ptable.add_column("Test Accuracy",test_acc)

# Printing the Table
print(ptable)
```

```
+-------+----------------------------------+---------------+-----------
----+-------------------+---------------+
| S.NO. |              MODEL               | Optimal Depth | Optimal sp
lit | Training Accuracy | Test Accuracy |
+-------+----------------------------------+---------------+-----------
----+-------------------+---------------+
|   1   |       Decision_Tree for BoW      |       50      |     500
    |       88.03       |      82.9     |
|   2   |      Decision_Tree for TFIDF     |       50      |     500
    |       88.71       |      82.9     |
|   3   |  Decision_Tree for Avg_Word2Vec  |       10      |     500
    |       88.35       |     82.42     |
|   4   | Decision_Tree for tfidf_Word2Vec |       10      |     500
    |       87.36       |      78.0     |
+-------+----------------------------------+---------------+-----------
----+-------------------+---------------+
```