

▼ DECISION TREE ALGORITHM on Amazon Fine Food Reviews

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 - 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

Objective:

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

▼ Loading,Cleaning & Preprocessing 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 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

3. If the score is above 3, then the recommendation will be set to "positive". Otherwise, it will be set to "negative".

```
1 %matplotlib inline
2 import warnings
3 import sqlite3
4 import pandas as pd
5 import numpy as np
6 import nltk
7 import string
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10 from sklearn.feature_extraction.text import TfidfTransformer
11 from sklearn.feature_extraction.text import TfidfVectorizer
12 from sklearn.feature_extraction.text import CountVectorizer
13 from sklearn.metrics import confusion_matrix
14 from sklearn import metrics
15 from sklearn.metrics import roc_curve, auc
16 from nltk.stem.porter import PorterStemmer
17
```

```

18 import re
19 import string
20 from nltk.corpus import stopwords
21 from nltk.stem import PorterStemmer
22 from nltk.stem.wordnet import WordNetLemmatizer
23
24 from gensim.models import Word2Vec
25 from gensim.models import KeyedVectors
26 import pickle
27
28 from tqdm import tqdm
29 import os
30
31 warnings.filterwarnings("ignore")

```

⌘ /usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning: pandas.util.testing is deprecated; use pandas.util._testing instead

```

1 from google.colab import drive
2 drive.mount('/content/drive')

```

⌘ Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pf

Enter your authorization code:

 Mounted at /content/drive

```

1 con = sqlite3.connect("/content/drive/My Drive/Colab Notebooks/database.sqlite")
2
3 filtered_data=pd.read_sql_query("""SELECT * FROM Reviews WHERE Score != 3""",con);
4 filtered_data.head(3)

```

⌘

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

```

1 filtered_data.shape
2
3 def partition(x):
4     if x < 3 :
5         return 'negative'
6     return 'positive'
7
8 actualScore=filtered_data['Score']
9 positive_negative=actualScore.map(partition)
10 filtered_data['Score']=positive_negative
11 print("Number of datapoints",filtered_data.shape)
12 filtered_data.head(3)

```

⌘

[illegible]

```
1 print(display.shape)
2 display.head(3)
```

	UserId	ProductId	ProfileName	Time	Score	
0	#oc-R115TNMSPFT9I7	B007Y59HVM	Breyton	1331510400	2	Overall its just OK when considering the price
1	#oc-R11D9D7SHXIJB9	B005HG9ET0	Louis E. Emory "hoppy"	1342396800	5	My wife has recurring extreme muscle spasms and this is the only product that has helped
2	#oc-R11DNU2NBKQ23Z	B007Y59HVM	Kim Cieszykowski	1348531200	1	This coffee is horrible and unfortunately I have to drink it every day

```
1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND UserId="AR5J8UI46CURR"
5 ORDER BY ProductID
6 """, con)
7 display.head()
```

	ID	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	
0	78445	B000HDL1RQ	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199
1	138317	B000HDOPYC	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199
2	138277	B000HDOPYM	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199
3	73791	B000HDOPZG	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199
4	155049	B000PAQ75C	AR5J8UI46CURR	Geetha Krishnan	2	2	5	1199

```
1 sorted_data=filtered_data.sort_values('ProductId',axis=0,ascending=True,inplace=False,kind='quicksort',na_position='last')
2
3 final_data=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"},keep='first',inplace=False)
4 final_data.shape
```

```
↳ (364173, 10)
```

```
1 (final_data['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

```
↳ 69.25890143662969
```

```
1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND Id=44737 OR Id=64422
5 ORDER BY ProductID
6 """, con)
7
8 display.head()
```

```
↳
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score
0	64422	B000MIDROQ	A161DK06JJMCYF	J. E. Stephens "Jeanne"	3	1	5
1	44737	B001EQ55RW	A2V0I904FH7ABY	Ram	3	2	4

```
1 final_data=final_data[final_data.HelpfulnessNumerator<=final_data.HelpfulnessDenominator]
```

```
1 nltk.download('stopwords')
```

```
↳ [nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
True
```

```
1 stopping_words = set(stopwords.words('english'))
2 print(stopping_words)
```

```
↳ {'before', 'ourselves', "couldn't", 'didn', 'the', 'some', "aren't", 'itself', 'can', 'that', 'them', 'did', 'sh
```

```
1 def clean_html(text):
2     clean_r = re.compile('<,*?>')
3     clean_text = re.sub(clean_r, '', text)
4     return clean_text
5
6 def Clean_punc(text):
7     clean_sentence = re.sub(r'[?|!|\'|\"|#]', r' ', text)
8     clean_data = re.sub(r'[,|,|)(|\\|/)]', r' ', clean_sentence)
9     return clean_data
```

```
1 from tqdm import tqdm
2 import os
3 import pdb
4 import pickle
5
6 from tqdm import tqdm
7 import os
8 import pdb
9 import pickle
10
11 stem_no = nltk.stem.SnowballStemmer('english')
12
13 if not os.path.isfile('final_data.sqlite'):
14     final_string=[]
```

```

15 all_positive_words=[]
16 all_negative_words=[]
17 for i,sentence in enumerate(tqdm(final_data['Text'].values)):
18     filtered_sentence=[]
19     sent_without_html_tags=clean_html(sentence)
20     #pdb.set_trace()
21     for w in sent_without_html_tags.split():
22         for cleaned_words in Clean_punc(w).split():
23             if ((cleaned_words.isalpha()) & (len(cleaned_words) > 2)):
24                 if(cleaned_words.lower() not in stopping_words) :
25                     stemming=(stem_no.stem(cleaned_words.lower())).encode('utf8')
26                     filtered_sentence.append(stemming)
27                     if(final_data['Score'].values)[i]=='positive':
28                         all_positive_words.append(stemming)
29                     if(final_data['Score'].values)[i]=='negative':
30                         all_negative_words.append(stemming)
31     str1 = b" ".join(filtered_sentence)
32     final_string.append(str1)
33
34 final_data['Cleaned_text']=final_string
35 final_data['Cleaned_text']=final_data['Cleaned_text'].str.decode("utf-8")
36
37 conn = sqlite3.connect('final_data.sqlite')
38 cursor=conn.cursor()
39 conn.text_factory = str
40 final_data.to_sql('Reviews',conn,schema=None,if_exists='replace',index=True,index_label=None,chunksize=None,dt
41 conn.close()
42
43
44 with open('positive_words.pkl','wb') as f :
45     pickle.dump(all_positive_words,f)
46 with open('negative_words.pkl','wb') as f :
47     pickle.dump(all_negative_words,f)

```

100% |██████████| 364171/364171 [06:05<00:00, 996.62it/s]

```
1 final_data['total_words'] = [len(x.split()) for x in final_data['Cleaned_text'].tolist()]
```

```
1 final_data.sort_values(by=['Time'], inplace=True, ascending=True)
```

```
1 final_data.head(3)
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score
138706	150524	0006641040	ACITT7DI6IDDL	shari zychinski	0	0	positive
138683	150501	0006641040	AJ46FKXOVC7NR	Nicholas A Mesiano	2	2	positive
417839	451856	B00004CXX9	AIUWLEQ1ADEG5	Elizabeth Medina	0	0	positive

```
1 final_data['Score'].value counts()
```

```

↳ positive    307061
   negative    57110
   Name: Score, dtype: int64

```

```
1 count_positive,count_negative=final_data['Score'].value_counts()
```

```
1 count_positive
```

```
↳ 307061
```

```

1 final_data_positive_class=final_data[final_data['Score']=='positive']
2 final_data_negative_class=final_data[final_data['Score']=='negative']

```

```
1 final_data_negative_class.shape
```

```
↳ (57110, 12)
```

▼ RANDOM DOWN SAMPLING

Note : In Down Sampling , there will be loss of information. Since we are removing the random records majority class

```

1 final_data_positive=final_data_positive_class.sample(count_negative)
2
3 fina_data_after_Sampling=pd.concat([final_data_positive,final_data_negative_class], axis=0)

```

```
1 fina_data_after_Sampling['Score'].value_counts()
```

```

↳ negative    57110
   positive    57110
   Name: Score, dtype: int64

```

```
1 fina_data_after_Sampling.shape
```

```
↳ (114220, 12)
```

```

1 final_data_100K=fina_data_after_Sampling[0:100000]
2 amazon_polarity_labels=final_data_100K['Score'].values
3 final_data_100K.head(2)

```

```
↳
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score
372042	402328	B0051S7P54	A29CNJP06GO2N1	Sandy	1	1	positiv
436815	472353	B001JG537O	AWX2C6PPAMFX0	Melissa	1	1	positiv

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.metrics import accuracy_score
4 from sklearn.model_selection import cross_val_score
5 from collections import Counter
6 from sklearn.metrics import confusion_matrix
7 from sklearn.metrics import classification_report
8
9 X_1,X_Test,Y_1,Y_Test = train_test_split(final_data_100K,amazon_polarity_labels,test_size=0.2,random_state=0)
10 X_Train,X_CV,Y_Train,Y_CV = train_test_split(X_1,Y_1,test_size=0.2)

```

APPLY BAG OF WORDS VECTORIZATION TECHNIQUE USING DECISION TREE CLASSIFIER THE BEST DEPTH

```

1 print(X_Train.shape, Y_Train.shape)
2 print(X_CV.shape, Y_CV.shape)
3 print(X_Test.shape, Y_Test.shape)
4
5 print("="*100)
6
7
8 count_vector=CountVectorizer(min_df=1)
9 X_Train_data_bow=(count_vector.fit_transform(X_Train['Cleaned_text']).values)
10 X_Test_data_bow=(count_vector.transform(X_Test['Cleaned_text']).values)
11 X_CV_data_bow=(count_vector.transform(X_CV['Cleaned_text']).values)
12
13 print("After vectorizations")
14 print(X_Train_data_bow.shape, Y_Train.shape)
15 print(X_CV_data_bow.shape, Y_CV.shape)
16 print(X_Test_data_bow.shape, Y_Test.shape)
17 print("="*100)

```



```
(64000, 12) (64000,)
(16000, 12) (16000,)
(20000, 12) (20000,)
```

```
=====
After vectorizations
```

```
(64000, 30582) (64000,)
(16000, 30582) (16000,)
(20000, 30582) (20000,)
=====
```

```
1 from sklearn.model_selection import GridSearchCV
2 from scipy.stats import randint as sp_randint
3 from sklearn.model_selection import cross_val_score
4 from sklearn.tree import DecisionTreeClassifier
5
6 def Decision_tree_Classifier(x_training_data,y_training_data):
7     grid_params = { 'max_depth' : [1,5,10,50,100,500,1000],
8                     'min_samples_split' : [5,10,100,500]
9                     }
10    Classifier_DT = DecisionTreeClassifier(random_state=None, class_weight='balanced')
11    clf=GridSearchCV(Classifier_DT,grid_params,scoring='roc_auc',return_train_score=True,cv=10)
12    clf.fit(x_training_data,y_training_data)
13    results = pd.DataFrame.from_dict(clf.cv_results_)
14    results = results.sort_values(['param_max_depth'])
15    results = results.sort_values(['param_min_samples_split'])
16    train_auc= results['mean_train_score']
17    train_auc_std= results['std_train_score']
18    cv_auc = results['mean_test_score']
19    cv_auc_std= results['std_test_score']
20    best_depth = results['param_max_depth']
21    min_sample_split = results['param_min_samples_split']
22    #log_alpha=np.log10(list(results["param_alpha"]))
23    print(clf.best_score_)
24    print(clf.best_params_)
25    plt.plot(best_depth, train_auc, label='Train AUC')
26    plt.plot(best_depth, cv_auc, label='CV AUC')
27    plt.scatter(best_depth, train_auc, label='Train AUC points')
28    plt.scatter(best_depth, cv_auc, label='CV AUC points')
29    plt.legend()
30    plt.xlabel("Best Depth : Hyperparameter")
31    plt.ylabel("AUC")
32    plt.title("Hyper parameter Vs AUC plot")
33    plt.grid()
34    plt.show()
35    return results,clf,min_sample_split,Classifier_DT

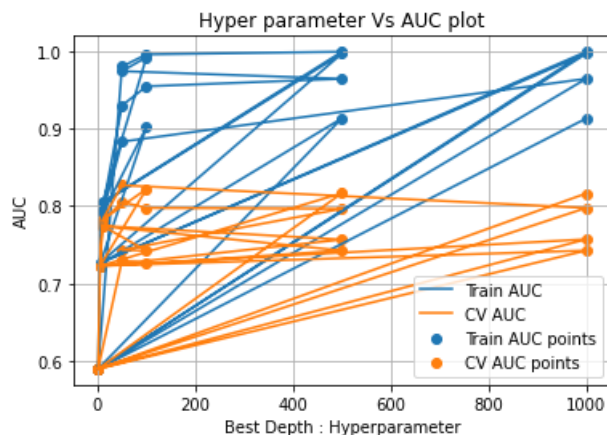
1 print ('-----BEST DEPTH USING DIFFERENT RANGE OF SAMPLE SPLIT')
2 result,best_depth,min_sample_split,decision_tree=Decision_tree_Classifier(X_Train_data_bow,Y_Train)
```



-----BEST DEPTH USING DIFFERENT RANGE OF SAMPLE SPLIT

0.8269227163421584

```
{'max_depth': 50, 'min_samples_split': 500}
```



```
1 pip install graphviz
```

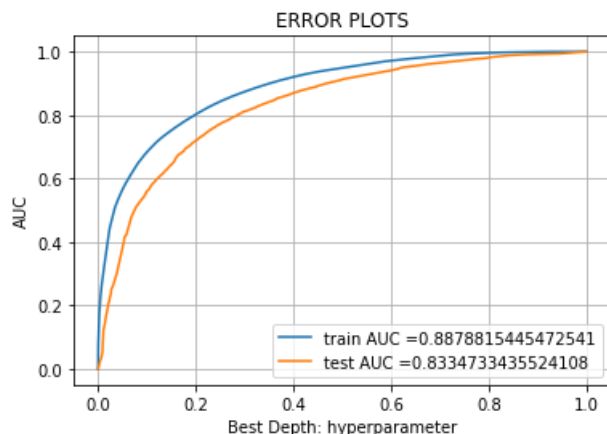
Requirement already satisfied: graphviz in /usr/local/lib/python3.6/dist-packages (0.10.1)

```
1 def Find_best_Depth(best_depth) :
2     best_depth = best_depth.best_params_
3     best_depth=best_depth.get("max_depth")
4     print(best_depth)
5     return best_depth
```

```
1 best_depth = Find_best_Depth(best_depth)
```

50

```
1 from sklearn.metrics import roc_curve, auc
2
3 decision_tree= DecisionTreeClassifier(max_depth=50,random_state=None, class_weight = 'balanced',min_samples_split=5
4 clf=decision_tree.fit(X_Train_data_bow,Y_Train)
5 pred_test_data=decision_tree.predict(X_Test_data_bow)
6 y_train_predicted_prob = decision_tree.predict_proba(X_Train_data_bow)[:,-1]
7 y_test_predicted_prob=decision_tree.predict_proba(X_Test_data_bow)[:,-1]
8 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
9 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
10 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("Best Depth: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("ERROR PLOTS")
16 plt.grid()
17 plt.show()
```



▼ FIND THE TOP 10 POSITIVE & NEGATIVE FEATURE IMPORTANCE

```
1 feature_names = count_vector.get_feature_names()
2 coef = decision_tree.feature_importances_
3 coefs_with_fns = sorted(zip(coef, feature_names))
4 top_features = zip(coefs_with_fns[:20], coefs_with_fns[:-(20 + 1):-1])
5 list(top_features)
```

```
[((0.0, 'aaa'), (0.1236684175187976, 'great')),
 ((0.0, 'aaaaaah'), (0.07609134152081516, 'disappoint')),
 ((0.0, 'aaaaaahhhhyaaaaaa'), (0.07222343378029503, 'love')),
 ((0.0, 'aaaahhhhhhhhhhh'), (0.05697057925290443, 'best')),
 ((0.0, 'aaaand'), (0.04649185534002008, 'delici')),
 ((0.0, 'aaah'), (0.03780892096838088, 'bad')),
 ((0.0, 'aaahhhhh'), (0.029834880209780194, 'perfect')),
 ((0.0, 'aachen'), (0.025822645675856973, 'good')),
 ((0.0, 'aad'), (0.022588207511681016, 'thought')),
 ((0.0, 'aadult'), (0.021714440275189956, 'favorit')),
 ((0.0, 'aaf'), (0.0206032805560982, 'excel')),
 ((0.0, 'aafco'), (0.0190078941097605, 'return')),
 ((0.0, 'aah'), (0.018835473774696902, 'money')),
 ((0.0, 'aappubl'), (0.016752904967661568, 'tast')),
 ((0.0, 'aar'), (0.014240257776803395, 'nice')),
 ((0.0, 'aarp'), (0.010728983483504234, 'review')),
 ((0.0, 'aarrgh'), (0.010648035010274169, 'tasti')),
 ((0.0, 'aback'), (0.010434062444733945, 'horribl')),
 ((0.0, 'abalon'), (0.009822693535309353, 'worst')),
 ((0.0, 'abandon'), (0.009342407470265334, 'amaz'))]
```

```
1 from sklearn.metrics import roc_auc_score
2 from sklearn.metrics import classification_report, confusion_matrix
3
4 roc_auc_score(Y_Test, y_test_predicted_prob)
```

```
0.8334733435524108
```

```
1 print(classification_report(Y_Test, pred_test_data))
2 print(confusion_matrix(Y_Test, pred_test_data))
```

	precision	recall	f1-score	support
negative	0.70	0.77	0.73	8582
positive	0.81	0.75	0.78	11418
accuracy			0.76	20000
macro avg	0.76	0.76	0.76	20000
weighted avg	0.76	0.76	0.76	20000

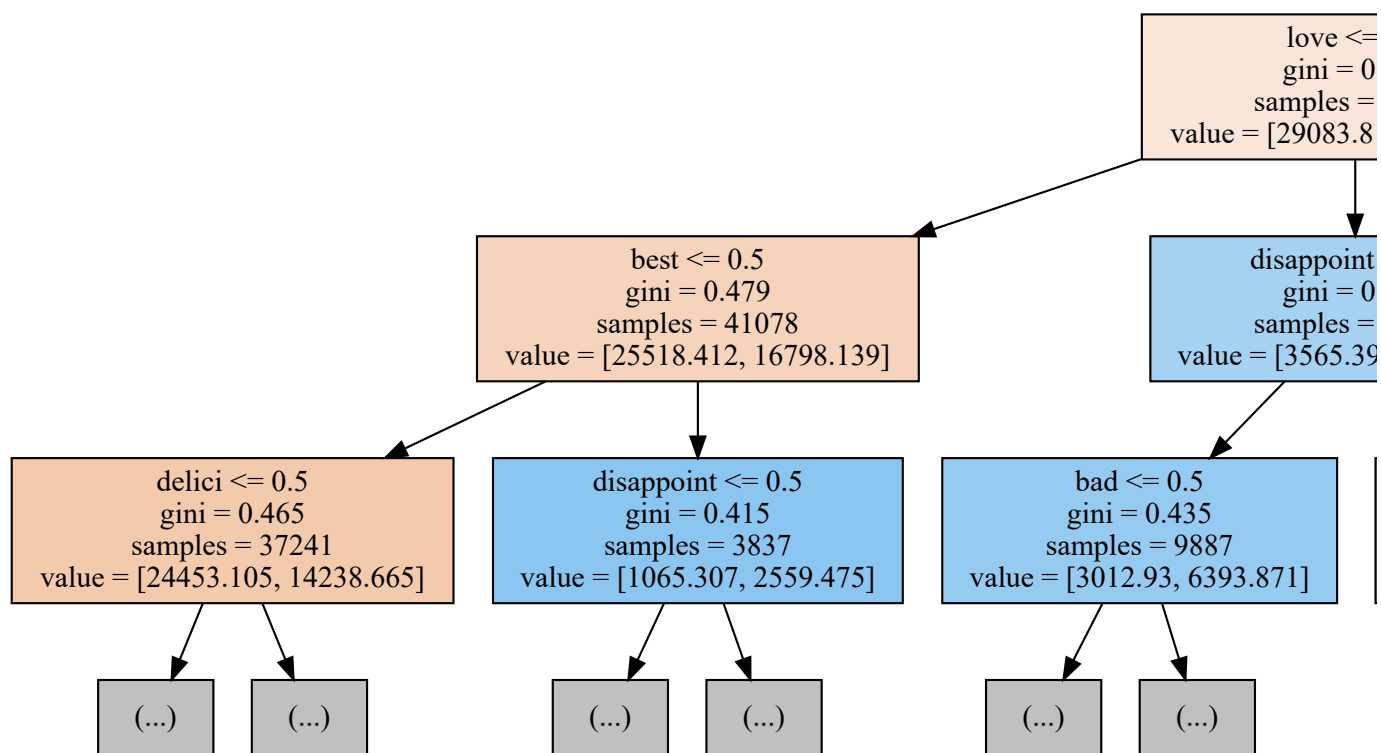
[[6640 1942]
[2868 8550]]

▼ VISUALIZE THE DECISION TREE USING GRAPHVIZ

```

1 import graphviz
2 from sklearn.tree import export_graphviz
3
4 export_graphviz(decision_tree, out_file="my_bow_tree.dot", feature_names = count_vector.get_feature_names(), max_c
5 with open("my_bow_tree.dot") as f:
6     dot_graph = f.read()
7 graphviz.Source(dot_graph)

```



▼ TF-IDF VECTORIZATION TECHNIQUE USING DECISION TREE

```

1 from sklearn.feature_extraction.text import TfidfVectorizer

```

```

2
3 print(X_Train.shape, Y_Train.shape)
4 print(X_CV.shape, Y_CV.shape)
5 print(X_Test.shape, Y_Test.shape)
6
7 print("="*100)
8
9
10 tfidf_vector=TfidfVectorizer(min_df=10)
11 X_Train_data_tfidf=(tfidf_vector.fit_transform(X_Train['Cleaned_text'].values))
12 X_Test_data_tfidf=(tfidf_vector.transform(X_Test['Cleaned_text'].values))
13 X_CV_data_tfidf=(tfidf_vector.transform(X_CV['Cleaned_text'].values))
14
15 print("After vectorizations")
16 print(X_Train_data_tfidf.shape, Y_Train.shape)
17 print(X_CV_data_tfidf.shape, Y_CV.shape)
18 print(X_Test_data_tfidf.shape, Y_Test.shape)
19 print("="*100)

```

```

↳ (64000, 12) (64000,)
   (16000, 12) (16000,)
   (20000, 12) (20000,)
=====
After vectorizations
(64000, 6982) (64000,)
(16000, 6982) (16000,)
(20000, 6982) (20000,)
=====

```

```

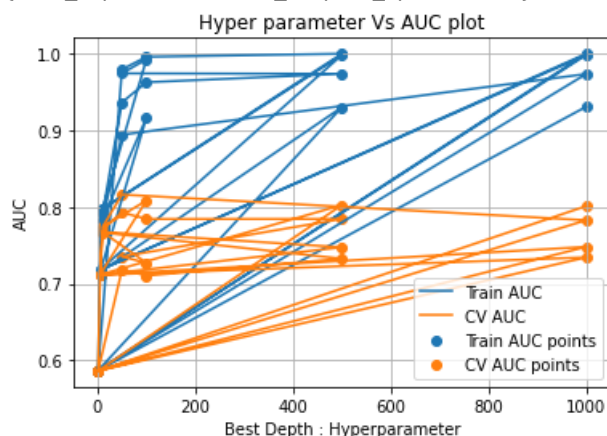
1 print ('-----BEST DEPTH USING DIFFERENT RANGE OF SAMPLE SPLIT')
2 result,best_depth,min_sample_split,decision_tree_tfidf=Decision_tree_Classifier(X_Train_data_tfidf,Y_Train)

```

```

↳ -----BEST DEPTH USING DIFFERENT RANGE OF SAMPLE SPLIT
0.8161551568177459
{'max_depth': 50, 'min_samples_split': 500}

```



```

1 best_depth = Find_best_Depth(best_depth)

```

```

↳ 50

```

```

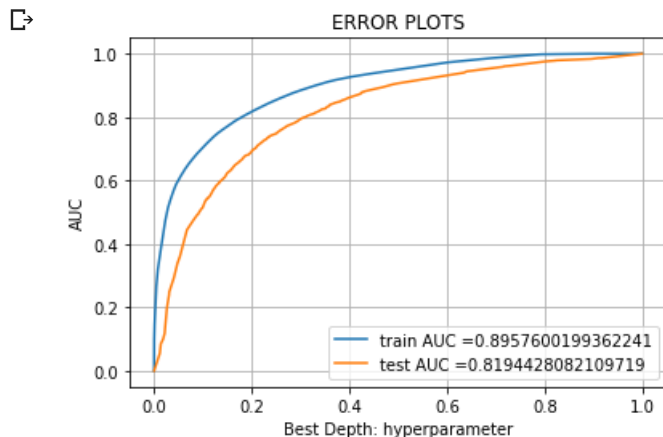
1 from sklearn.metrics import roc_curve, auc
2
3 decision_tree= DecisionTreeClassifier(max_depth=50,random_state=None, class_weight='balanced',min_samples_split=5)
4 clf=decision_tree.fit(X_Train_data_tfidf,Y_Train)
5 pred_test_data=decision_tree.predict(X_Test_data_tfidf)
6 y_train_predicted_prob = decision_tree.predict_proba(X_Train_data_tfidf)[:,-1]
7 y_test_predicted_prob=decision_tree.predict_proba(X_Test_data_tfidf)[:,-1]
8 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
9 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
10 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))

```

```

11 plt.plot(test_fpr, test_tpr, label="test AUC =" + str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("Best Depth: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("ERROR PLOTS")
16 plt.grid()
17 plt.show()

```



```

1 feature_names = tfidf_vector.get_feature_names()
2 coef = decision_tree.feature_importances_
3 coefs_with_fns = sorted(zip(coef, feature_names))
4 top_features = zip(coefs_with_fns[:20], coefs_with_fns[-(20 + 1):-1])
5 list(top_features)

```

```

[((0.0, 'abandon'), (0.119583149703108, 'great')),
 ((0.0, 'abdomin'), (0.07876461867874446, 'love')),
 ((0.0, 'abil'), (0.06551048114796722, 'disappoint')),
 ((0.0, 'abnorm'), (0.05513699661609732, 'best')),
 ((0.0, 'abomin'), (0.043514076712564805, 'delici')),
 ((0.0, 'abroad'), (0.031624697865599585, 'perfect')),
 ((0.0, 'absenc'), (0.02855670121807223, 'good')),
 ((0.0, 'absent'), (0.024225330947117193, 'bad')),
 ((0.0, 'absolut'), (0.02338306543496268, 'favorit')),
 ((0.0, 'absorb'), (0.020743958280810784, 'excel')),
 ((0.0, 'absorpt'), (0.018155354606437885, 'return')),
 ((0.0, 'absurd'), (0.016899149378417147, 'thought')),
 ((0.0, 'abund'), (0.016727685635741073, 'nice')),
 ((0.0, 'abus'), (0.013437025603691977, 'wast')),
 ((0.0, 'acai'), (0.013367511071551532, 'tast')),
 ((0.0, 'accent'), (0.012387272541635018, 'easi')),
 ((0.0, 'access'), (0.01150590456613557, 'aw')),
 ((0.0, 'accessori'), (0.010232869023075354, 'keep')),
 ((0.0, 'accid'), (0.009755070269865338, 'horribl')),
 ((0.0, 'accident'), (0.009202447902923704, 'amaz'))]

```

```

1 from sklearn.metrics import roc_auc_score
2 from sklearn.metrics import classification_report, confusion_matrix
3
4 roc_auc_score(Y_Test, y_test_predicted_prob)

```

```
0.8194428082109719
```

```

1 print(classification_report(Y_Test, pred_test_data))
2 print(confusion_matrix(Y_Test, pred_test_data))

```

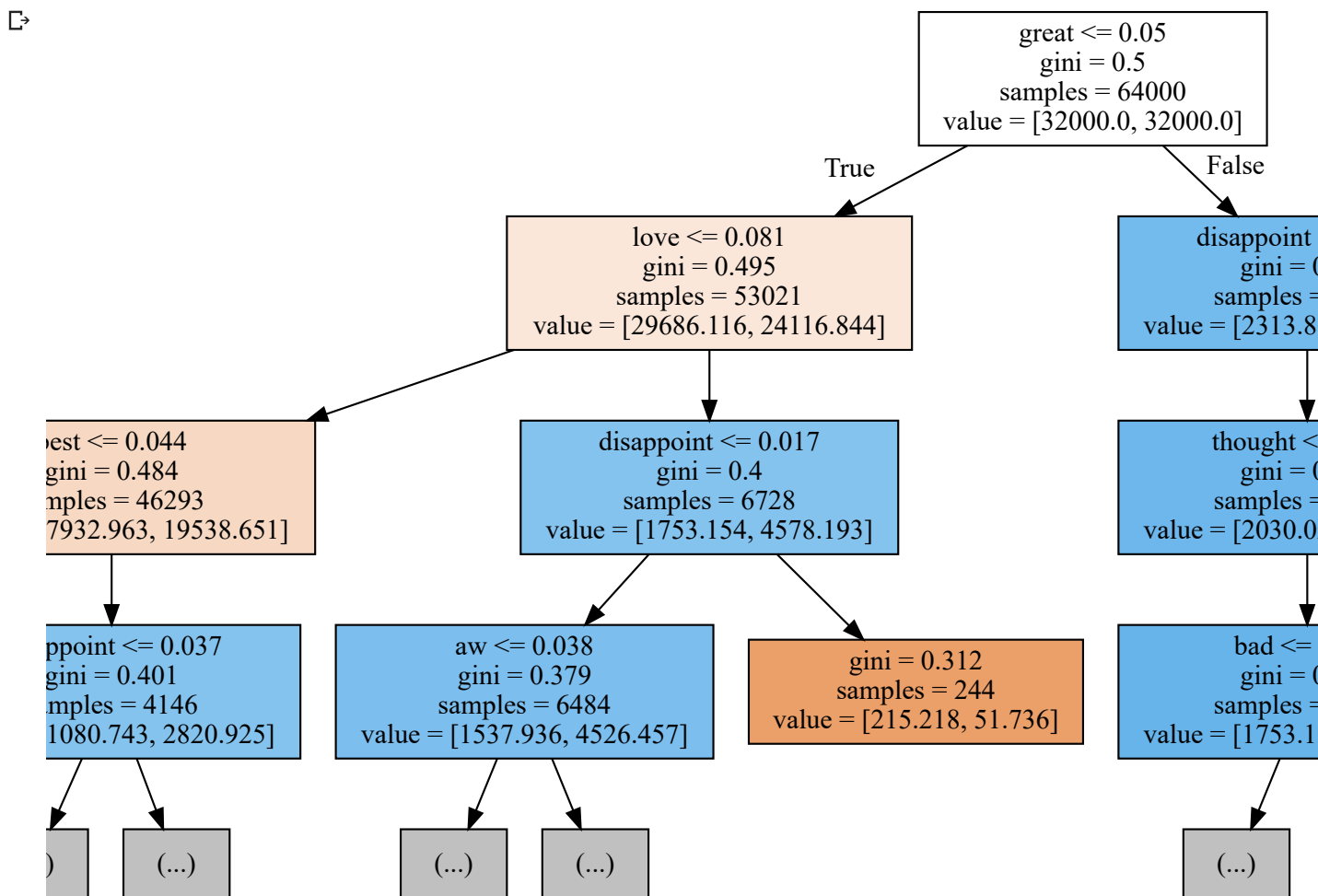
```


```

	precision	recall	f1-score	support
negative	0.69	0.76	0.72	8582
positive	0.81	0.74	0.77	11418
accuracy			0.75	20000
macro avg	0.75	0.75	0.75	20000
weighted avg	0.76	0.75	0.75	20000

```
[[6563 2019]
 [2986 8432]]
```

```
1 import graphviz
2 from sklearn.tree import export_graphviz
3
4 export_graphviz(decision_tree, out_file="my_tfidf_tree.dot", feature_names = tfidf_vector.get_feature_names(), max
5 with open("my_tfidf_tree.dot") as f:
6     dot_graph = f.read()
7 graphviz.Source(dot_graph)
```



▼ Avg Word2Vec Vectorization Technique on Decision Tree Algorithm

```
1 from gensim.models import Word2Vec
2 from gensim.models import KeyedVectors
3 import pickle
4
5 list of sent train avgw2v=[1
```

```

6 list_of_sent_test_avgw2v=[]
7 list_of_sent_cv_avgw2v=[]
8 for sent_train_avgw2v in tqdm(X_Train['Cleaned_text'].values):
9     list_of_sent_train_avgw2v.append(sent_train_avgw2v.split())

```

100%|██████████| 64000/64000 [00:00<00:00, 83969.58it/s]

```

1 for sent_test_avgw2v in tqdm(X_Test['Cleaned_text'].values):
2     list_of_sent_test_avgw2v.append(sent_test_avgw2v.split())
3
4 for sent_cv_avgw2v in tqdm(X_CV['Cleaned_text'].values):
5     list_of_sent_cv_avgw2v.append(sent_cv_avgw2v.split())

```

100%|██████████| 20000/20000 [00:00<00:00, 156933.10it/s]
100%|██████████| 16000/16000 [00:00<00:00, 176717.61it/s]

```

1 w2v_model_train = Word2Vec(list_of_sent_train_avgw2v,min_count=5,size=50,workers=4)
2 w2v_words_svm_train=list(w2v_model_train.wv.vocab)

```

```

1 w2v_model_test = Word2Vec(list_of_sent_test_avgw2v,min_count=5,size=50,workers=4)
2 w2v_words_svm_test=list(w2v_model_test.wv.vocab)

```

```

1 w2v_model_cv = Word2Vec(list_of_sent_cv_avgw2v,min_count=5,size=50,workers=4)
2 w2v_words_svm_cv=list(w2v_model_cv.wv.vocab)

```

```

1 train_vectors=[];
2 for sent in list_of_sent_train_avgw2v:
3     sent_vec=np.zeros(50)
4     cnt_words=0;
5     for word in sent:
6         if word in w2v_words_svm_train:
7             vec=w2v_model_train.wv[word]
8             sent_vec+=vec
9             cnt_words+=1
10    if cnt_words !=0:
11        sent_vec/=cnt_words
12    train_vectors.append(sent_vec)
13 print(len(train_vectors))
14 print(len(train_vectors[0]))

```

64000
50

```

1 test_vectors=[];
2 for sent in tqdm(list_of_sent_test_avgw2v):
3     sent_vec=np.zeros(50)
4     cnt_words=0;
5     for word in sent:
6         if word in w2v_words_svm_test:
7             vec=w2v_model_test.wv[word]
8             sent_vec+=vec
9             cnt_words+=1
10    if cnt_words !=0:
11        sent_vec/=cnt_words
12    test_vectors.append(sent_vec)
13 print(len(test_vectors))
14 print(len(test_vectors[0]))

```

100%|██████████| 20000/20000 [00:16<00:00, 1183.73it/s]20000
50

```

1 cv_vectors=[];
2 for sent in tqdm(list_of_sent_cv_avgw2v):
3     sent_vec=np.zeros(50)
4     cnt_words=0;
5     for word in sent:
6         if word in w2v_words_svm_cv:
7             vec=w2v_model_cv.wv[word]
8             sent_vec+=vec
9             cnt_words+=1
10    if cnt_words !=0:
11        sent_vec/=cnt_words
12    cv_vectors.append(sent_vec)
13 print(len(cv_vectors))
14 print(len(cv_vectors[0]))

```

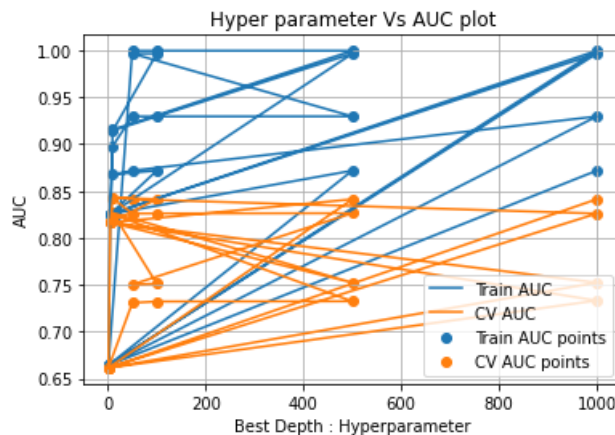
100%|██████████| 16000/16000 [00:13<00:00, 1203.18it/s]16000
50

```

1 print ('-----BEST DEPTH USING DIFFERENT RANGE OF SAMPLE SPLIT')
2 result,best_depth,min_sample_split,decision_tree=Decision_tree_Classifier(train_vectors,Y_Train)

```

-----BEST DEPTH USING DIFFERENT RANGE OF SAMPLE SPLIT
0.8429747830570099
{'max_depth': 10, 'min_samples_split': 500}



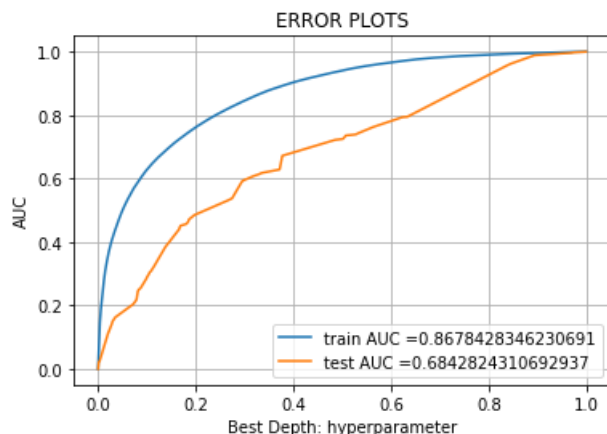
```
1 best_depth = Find_best_Depth(best_depth)
```

10

```

1 from sklearn.metrics import roc_curve, auc
2
3 decision_tree= DecisionTreeClassifier(max_depth=10,random_state=None, class_weight = 'balanced',min_samples_split=5
4 clf=decision_tree.fit(train_vectors,Y_Train)
5 pred_test_data=decision_tree.predict(test_vectors)
6 y_train_predicted_prob = decision_tree.predict_proba(train_vectors)[:,:1]
7 y_test_predicted_prob=decision_tree.predict_proba(test_vectors)[:,:1]
8 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
9 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
10 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("Best Depth: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("ERROR PLOTS")
16 plt.grid()
17 plt.show()

```

```
1 from sklearn.metrics import roc_auc_score
2 from sklearn.metrics import classification_report, confusion_matrix
3
4 roc_auc_score(Y_Test, y_test_predicted_prob)
```

0.6842824310692937

```
1 print(classification_report(Y_Test, pred_test_data))
2 print(confusion_matrix(Y_Test, pred_test_data))
```

precision recall f1-score support

negative	0.50	0.89	0.64	8582
positive	0.79	0.32	0.46	11418
accuracy			0.56	20000
macro avg	0.64	0.60	0.55	20000
weighted avg	0.66	0.56	0.53	20000

```
[[7611 971]
 [7761 3657]]
```

▼ TF-IDF Word2Vec Vectorization Technique for DECISION TREE on Amazon Fine Food Rev

```
1 model_Avgw2v = TfidfVectorizer()
2 X_Train_Avgw2v=model_Avgw2v.fit_transform(X_Train['Cleaned_text'].values)
```

```
1 X_Test_Avgw2v=model_Avgw2v.transform(X_Test['Cleaned_text'].values)
2 X_CV_Avgw2v=model_Avgw2v.transform(X_CV['Cleaned_text'].values)
```

```
1 dictionary = dict(zip(model_Avgw2v.get_feature_names(), list(model_Avgw2v.idf_)))
```

```
1 tfidf_feature=model_Avgw2v.get_feature_names()
2
3 tfidf_sent_vectors_train=[];
4 #final_tf_idf = [];
5 row=0;
6
7 for sent in tqdm(list_of_sent_train_avgw2v):
8     sent_vec=np.zeros(50)
9     weight_sum=0;
10    for word in sent :
11        if word in w2v_words_svm_train and word in tfidf_feature :
12            vec=w2v_model_train.wv[word]
```

```

13         #tf_idf=final_tf_idf[row,tfidf_feature.index(word)]
14         tf_idf=dictionary[word]*(sent.count(word)/len(sent))
15         sent_vec+=(vec*tf_idf)
16         weight_sum+=tf_idf
17
18     if weight_sum!=0:
19         sent_vec/=weight_sum
20     tfidf_sent_vectors_train.append(sent_vec)
21     row+=1

```

100% |██████████| 64000/64000 [14:31<00:00, 73.43it/s]

```

1 tfidf_sent_vectors_test=[];
2 #final_tf_idf = [];
3 row=0;
4
5 for sent in tqdm(list_of_sent_test_avg2v):
6     sent_vec=np.zeros(50)
7     weight_sum=0;
8     for word in sent :
9         if word in w2v_words_svm_test and word in tfidf_feature :
10             vec=w2v_model_test.wv[word]
11             #tf_idf=final_tf_idf[row,tfidf_feature.index(word)]
12             tf_idf=dictionary[word]*(sent.count(word)/len(sent))
13             sent_vec+=(vec*tf_idf)
14             weight_sum+=tf_idf
15
16     if weight_sum!=0:
17         sent_vec/=weight_sum
18     tfidf_sent_vectors_test.append(sent_vec)
19     row+=1

```

100% |██████████| 20000/20000 [04:19<00:00, 76.92it/s]

```

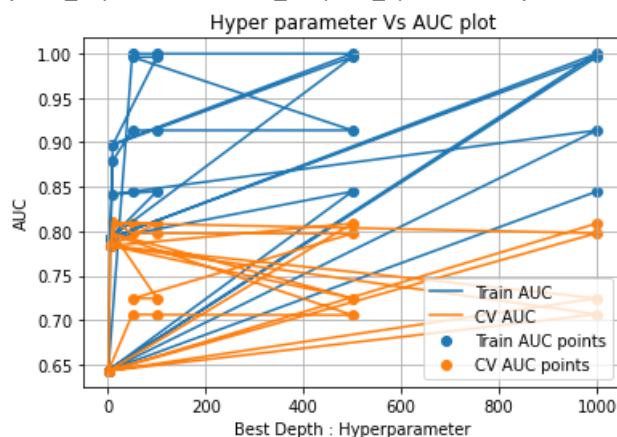
1 print ('-----BEST DEPTH USING DIFFERENT RANGE OF SAMPLE SPLIT')
2 result,best_depth,min_sample_split,decision_tree=Decision_tree_Classifier(tfidf_sent_vectors_train,Y_Train)

```

```

1 -----BEST DEPTH USING DIFFERENT RANGE OF SAMPLE SPLIT
2 0.811198148361221
3 {'max_depth': 10, 'min_samples_split': 100}

```



```
1 best_depth = Find_best_Depth(best_depth)
```

```
10
```

```

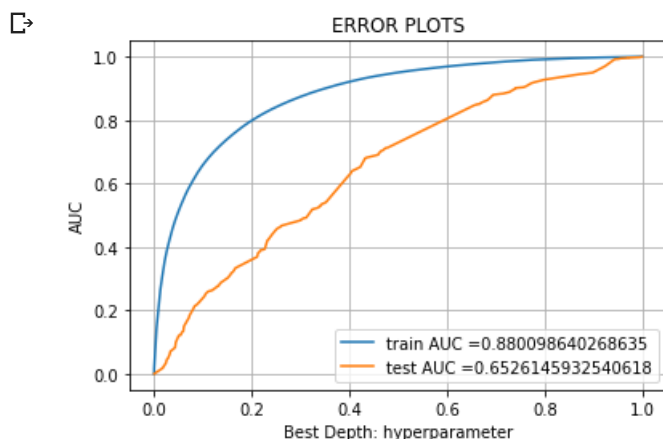
1 from sklearn.metrics import roc_curve, auc
2
3 decision_tree= DecisionTreeClassifier(max_depth=10,random_state=None, class_weight = 'balanced',min_samples_split=1

```

```

4 clf=decision_tree.fit(tfidf_sent_vectors_train,Y_Train)
5 pred_test_data=decision_tree.predict(tfidf_sent_vectors_test)
6 y_train_predicted_prob = decision_tree.predict_proba(tfidf_sent_vectors_train)[:,-1]
7 y_test_predicted_prob=decision_tree.predict_proba(tfidf_sent_vectors_test)[:,-1]
8 train_fpr, train_tpr, train_thresholds=roc_curve(Y_Train,y_train_predicted_prob,pos_label='positive')
9 test_fpr, test_tpr, test_thresholds = roc_curve(Y_Test, y_test_predicted_prob,pos_label='positive')
10 plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
11 plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
12 plt.legend()
13 plt.xlabel("Best Depth: hyperparameter")
14 plt.ylabel("AUC")
15 plt.title("ERROR PLOTS")
16 plt.grid()
17 plt.show()

```



```
1 roc_auc_score(Y_Test,y_test_predicted_prob)
```

```
0.6526145932540618
```

```

1 print(classification_report(Y_Test,pred_test_data))
2 print(confusion_matrix(Y_Test,pred_test_data))

```

```

precision    recall  f1-score   support

negative     0.50      0.77      0.60       8582
positive     0.71      0.42      0.53       11418

accuracy          0.57       20000
macro avg         0.60      0.59      0.57       20000
weighted avg      0.62      0.57      0.56       20000

[[6583 1999]
 [6600 4818]]

```

▼ PRETTY TABLE

```
1 pip install -U PTable
```

Collecting PTable

```

1 from prettytable import PrettyTable
2
3 x= PrettyTable()
4 x.field_names = ["Vectorizer" , "Hyperparameter(Best Depth)", "Minimum Sample Split","AUC"]
5 x.add_row(["Bag Of Words",50,500,0.8296])
6 x.add_row(["Tf-Idf",50,500,0.8161])
7 x.add_row(["Avg Word2Vec",10,500,0.8429])
8 x.add_row(["Tf-Idf Word2Vec",10,100,0.8111])
9 print(x)

```

```

↳ +-----+-----+-----+-----+
| Vectorizer | Hyperparameter(Best Depth) | Minimum Sample Split | AUC |
+-----+-----+-----+-----+
| Bag Of Words | 50 | 500 | 0.8296 |
| Tf-Idf | 50 | 500 | 0.8161 |
| Avg Word2Vec | 10 | 500 | 0.8429 |
| Tf-Idf Word2Vec | 10 | 100 | 0.8111 |
+-----+-----+-----+-----+

```

▼ Observation :

- 1) The AUC Score for Avg Word2Vec is 0.8429 with Hyperparameter = 10 & Minimum Sample Split = 500
- 2) To Balanced the Dataset , i have used Undersampling technique.
- 3) The Important Features does not show up the top 10 negative features

1

Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.