

Assignment 8 : Decision Tree

Decision Tree :

- It is a graphical representation of all possible solutions to a decision based on certain conditions.
- some terminologies:
 1. Root node: It represents entire population or sample and this further gets divided into two or more homogenous sets.
 2. Parent/Child node: Root node is the parent node and all other nodes branched from it is known as child node.
 3. Leaf node: Node cannot be further segregated into further nodes.
- In Decision tree corresponding to every decision we have an hyperplane, so all of our hyperplanes are axis-parallel, so intuitively Decision Tree is an "set of axis parallel hyperplanes".

Splitting the decision tree can be done based on some terminologies:

Entropy: it defines the randomness in the data and it can measure the impurity.

$$\text{Entropy}(s) = -P(\text{yes}) \log(P(\text{yes})) - P(\text{no}) \log(P(\text{no}))$$

where,

s is the total sample space

$P(\text{yes})$ is probability of yes

$P(\text{no})$ is probability of no

properties: 1) For random variable Y , if probabilities are equally probable then entropy is at maximum 1.

2) if one class fully dominates then the entropy is 0.

Information Gain: it can measure the reduction in entropy and decides which attribute should be selected as

the decision node.

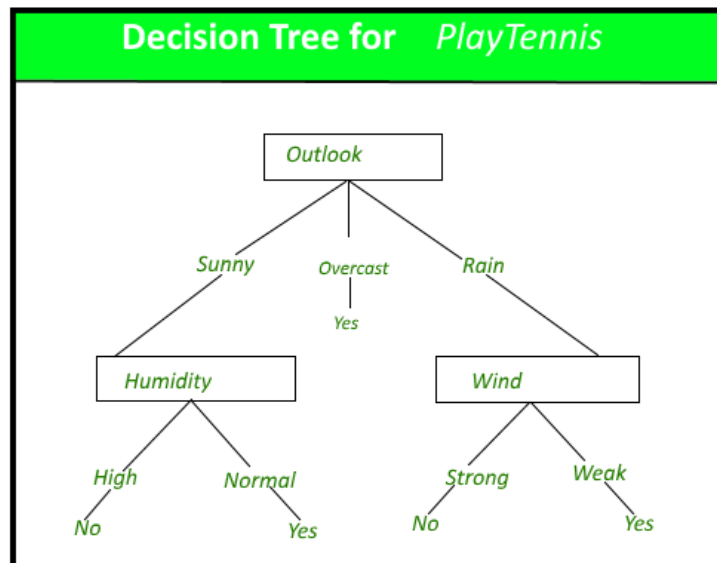
$$\text{Information Gain} = \text{Entropy}(s) - [(\text{Weighted Average}) \times \text{Entropy}(\text{each feature})]$$

properties: 1) For random variable Y , if probabilities are equally probable then information gain is 0.5.

2) if one class fully dominates then information gain is 0.

- Recursively break decision tree based on the information gain.
- When to stop tree:
 - 1) occurrences of stop nodes.
 - 2) when depth of the tree is more than at that time stop growing tree.
 - 3) if increase at every node few points are available then stop tree.
- If depth of the tree increases more then model leads to overfitting problem.
- If depth of the tree is small then model leads to underfitting problem.
- In case of decision tree no need to do feature standardization.
- If dimension is large then decision tree is not good option to use.
- As depth increases then harder to understand what happening in decision tree which means interpretability decreases.
- It can handle large data and dimension should be less.
- For low latency requirements decision tree is good.
- In case of regression we can split the decision tree of each node using the mean square error which is minimum.
- Cases:

1. Imbalanced data: Can be impacted while calculating the entropy or mean square error, so avoid that do upsampling or downsampling.
2. Large Dimension: If dimension is large then at each node we have to split, so to evaluate information gain of each feature by that time complexity to train decision tree increases more.
3. Categorical features: Do one hot encoding, if categories are more then convert into numerical feature.
4. Similarity matrix: Can't work well in case of decision tree.



```
In [4]: 1 #we use to ignore warnings
        2 import warnings
        3 warnings.filterwarnings('ignore')
```

```
In [2]: 1 import pickle
        2 def savetofile(obj,filename):
        3     pickle.dump(obj,open(filename,"wb"))
        4
        5 def openfromfile(filename):
        6     temp=pickle.load(open(filename,"rb"))
        7     return temp
```

```
In [5]: 1 #Importing the Libraries
2 import numpy as np
3 import pandas as pd
4
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 from sklearn.model_selection import train_test_split
9 from sklearn.feature_extraction.text import CountVectorizer
10
11 from sklearn.model_selection import TimeSeriesSplit
12 from sklearn.model_selection import GridSearchCV
13 from sklearn.tree import DecisionTreeClassifier
14
15 from sklearn import tree
16 from sklearn.metrics import confusion_matrix
17 from sklearn.metrics import accuracy_score
18 from sklearn.metrics import recall_score
19 from sklearn.metrics import precision_score
20 from sklearn.metrics import f1_score
21
22 from IPython.display import Image
23 import pydotplus
24 from sklearn.externals.six import StringIO
25 from sklearn.tree import export_graphviz
26 import graphviz
27 from sklearn.feature_extraction.text import TfidfVectorizer
28 import gensim
```

```
In [6]: 1 #Loading the dataset
2 data_frame = openfromfile("New_Amazon_preprocess_data")
```

```
In [7]: 1 #Shape of data
2 print("Shape of data_frame:", data_frame.shape)
3
4 #First five rows of the data_frame
5 data_frame.head()
```

Shape of data_frame: (364171, 11)

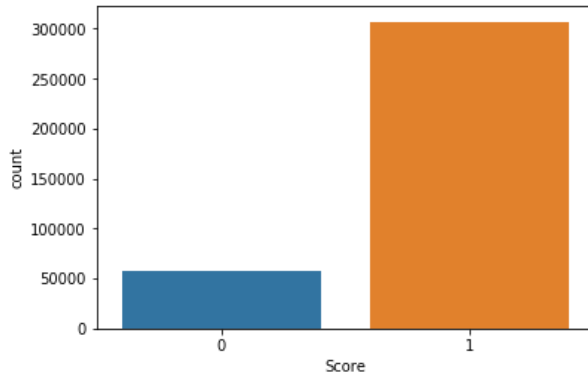
```
Out[7]:
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Si
515425	515426	141278509X	AB1A5EGHHVA9M	CHelmic	1	1	positive	1332547200	.
24750	24751	2734888454	A1C298ITT645B6	Hugh G. Pritchard	0	0	positive	1195948800	Di
24749	24750	2734888454	A13ISQV0U9GZIC	Sandikaye	1	1	negative	1192060800	
308076	308077	2841233731	A3QD68O22M2XHQ	LABRNT	0	0	positive	1345852800	b
150523	150524	6641040	ACITT7DI6IDDL	shari zychinski	0	0	positive	939340800	edi

```
In [8]: 1 #Storing the data_frame based on the time attribute
2 data_frame.sort_values('Time', inplace=True)
3
4 #Resetting the data_frame
5 data_frame.reset_index(drop=False, inplace=True)
```

```
In [9]: 1 #In the Score attribute consisting of two categories changing positive to 1 and negative to 0
2 data_frame.Score = [1 if(score == 'positive') else 0 for score in data_frame.Score]
```

```
In [10]: 1 #Count plot for score attribute
2 sns.countplot(x=data_frame.Score, data=data_frame)
3 plt.show()
4 data_frame.Score.value_counts()
```



```
Out[10]: 1    307061
0      57110
Name: Score, dtype: int64
```

```
In [11]: 1 #Taking the top 100k points
2 data_frame_100k = data_frame[0:100000]
```

```
In [12]: 1 #storing CleanedText attribute into X as a independent variable and Score attribute into y as a dependent
2 X = data_frame_100k.CleanedText
3
4 y = data_frame_100k.Score
```

```
In [13]: 1 #Splitting the data into train as 70% and test as 30%
2 X_tr, X_test, y_tr, y_test = train_test_split(X, y, test_size=0.3, shuffle=False)
```

```
In [14]: 1 #shape of train and test data
2 print("Shape of the train data:", X_tr.shape)
3
4 print("Shape of the test data:", X_test.shape)
```

```
Shape of the train data: (70000,)
Shape of the test data: (30000,)
```

BOW:

which means makes a vector for each review of length unique words from the whole dataset and makes frequency count of word.

- Bow or Bag of Words which means way of extracting features from text for use in modeling.
- A bag-of-words is a representation of text that describes the occurrence of words within a document.

It involves two things:

- 1.vocabulary of known words.
- 2.Measure of the presence of known words.

- It is called a “bag” of words, because any information about the order or structure of words in the document is discarded. The model is only concerned with whether known words occur in the document, not where in the document.

```
In [15]: 1
2 %%time
3
4 count_vec = CountVectorizer()
5 #Making the fit_transform for train data
6 bow_tr = count_vec.fit_transform(X_tr)
```

```
Wall time: 2.21 s
```

```
In [16]: 1 #Transform for test data
2 bow_test = count_vec.transform(X_test)
```

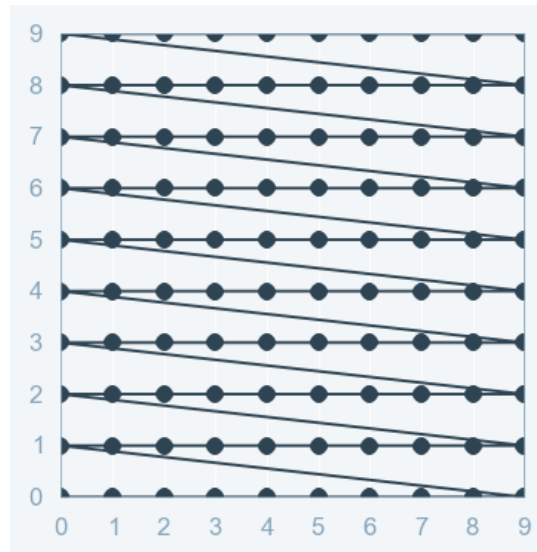
```
In [17]: 1 #Shape of train and test data after the bag of words
2 print("shape of train data:", bow_tr.shape)
3 print("shape of test data:", bow_test.shape)
```

```
shape of train data: (70000, 37189)
shape of test data: (30000, 37189)
```

Hyperparameter tuning using grid search cross validation:

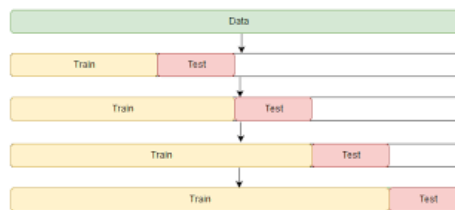
Grid Search Cross Validation:

- working through multiple combinations of parameter tunes, cross validate each and determine which one gives the best performance.
- Note: In grid search, if you choose n parameters then we will have to check 2^n combinations.



Time based splitting:

- Provides train/test indices to split time series data samples that are observed at fixed time intervals, in train/test sets. In each split, test indices must be higher than before, and thus shuffling in cross validator is inappropriate.



```
In [18]: 1 #Function for Grid Search Cross Validation
2 def grid_search(X_train, y_train):
3
4     #(1 <= depth <= 20)
5     values = [i for i in range(1,21)]
6
7     #Giving some set of parameters as input to grid search cross validation
8     parameters = {'max_depth':values}
9
10    #splitting the data based on time of 5 folds
11    tbs = TimeSeriesSplit(n_splits=5)
12
13    #measuring the quality of split using the entropy and making the data should be balanced
14    clf = DecisionTreeClassifier(criterion='entropy', class_weight='balanced')
15
16    gsv = GridSearchCV(clf, parameters, cv=tbs, verbose=3, n_jobs=3, scoring='f1')
17
18    gsv.fit(X_train,y_train)
19
20    print("Max depth:",gsv.best_params_)
21    print("Best F1-score:",gsv.best_score_*100)
22
23    return gsv.grid_scores_, gsv.best_estimator_
```

```
In [19]: 1
2 %%time
3
4 #Calling the function for Grid Search Cross Validation
5 grid_scores, best_estimator = grid_search(bow_tr, y_tr)
```

Fitting 5 folds for each of 20 candidates, totalling 100 fits

```
[Parallel(n_jobs=3)]: Done 26 tasks | elapsed: 9.0s
[Parallel(n_jobs=3)]: Done 100 out of 100 | elapsed: 1.4min finished
```

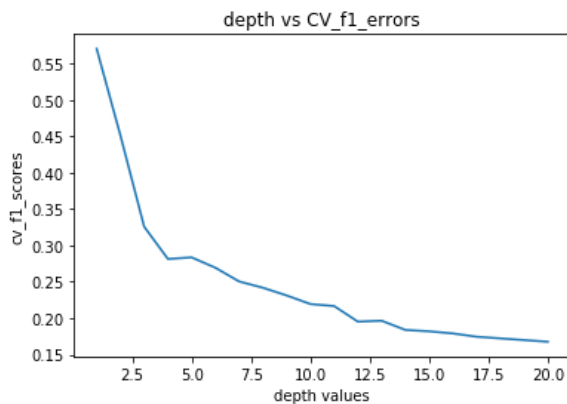
```
Max depth: {'max_depth': 20}
Best F1-score: 83.24871820830502
Wall time: 1min 35s
```

```
In [20]: 1 grid_scores[:2]
```

```
Out[20]: [mean: 0.42954, std: 0.00503, params: {'max_depth': 1},
mean: 0.54870, std: 0.00487, params: {'max_depth': 2}]
```

```
In [21]: 1 #Function for plot between F1_cv_errors and depth values
2
3 def depth_scores(depths, cv_f1_errors):
4     plt.plot(depths,cv_f1_errors)
5     plt.xlabel("depth values")
6     plt.ylabel("cv_f1_scores")
7     plt.title("depth vs CV_f1_errors")
8     plt.show()
```

```
In [22]: 1 #From grid_scores storing max_depths into depth
2 depths = [val[0]['max_depth'] for val in grid_scores]
3
4 #From grid_scores storing mean f1_cross validation scores into cv_f1_scores
5 cv_f1_errors = [1-val[1] for val in grid_scores]
6
7 #plot for depths and cv_f1_errors
8 depth_scores(depths, cv_f1_errors)
```



Testing the model from best_estimator which can be return by the grid search cross validation.

```
In [23]: 1 #Result showing the best classifier consisting of parameters
2 best_estimator
```

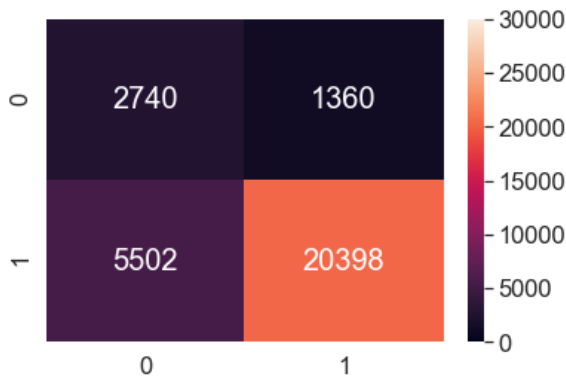
```
Out[23]: DecisionTreeClassifier(class_weight='balanced', criterion='entropy',
max_depth=20, max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best')
```

```
In [24]: 1 #Finding the predicted values for test labels using the test data
2 y_pred = best_estimator.predict(bow_test)
```

```
In [25]: 1 #Function for calculating the metrics
2 def test_metrics(y_test, y_pred):
3     cm = pd.DataFrame(confusion_matrix(y_test,y_pred),range(2),range(2))
4     sns.set(font_scale=1.5)
5     sns.heatmap(cm,annot=True,annot_kws={"size": 20}, fmt='g', vmin=0, vmax=30000)
6
7     print("Accuracy on test data:", round(accuracy_score(y_test, y_pred) * 100 , 2))
8     print("Precision on test data:", round(precision_score(y_test, y_pred) * 100 , 2))
9     print("Recall on test data:", round(recall_score(y_test, y_pred) * 100 , 2))
10    print("F1_score on test data:", round(f1_score(y_test, y_pred) * 100,2))
11
12    plt.show()
```

```
In [26]: 1 #Calling the function for test metrics
2 test_metrics(y_test, y_pred)
```

Accuracy on test data: 77.13
Precision on test data: 93.75
Recall on test data: 78.76
F1_score on test data: 85.6

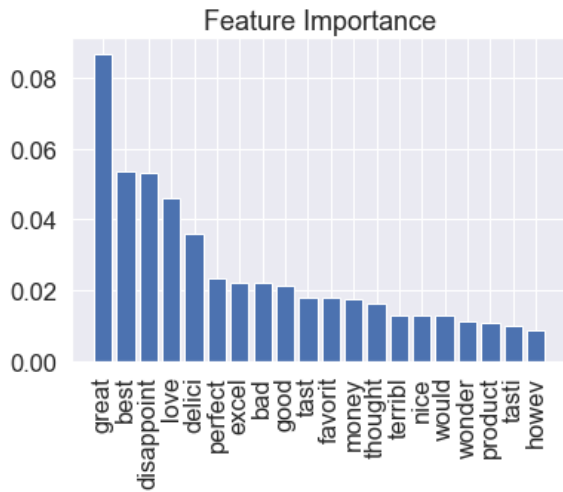


Top 10 Impotant Features :

```
In [27]: 1 #Extracting important featurrs from decision tree
2 features = best_estimator.feature_importances_
3
4 #Getting the top important feature indices
5 indices = np.argsort(features)[::-1][:20]
6
7 #Getting feature names from feature extractor object is count_vectorizer
8 fea_names = count_vec.get_feature_names()
```

```
In [28]: 1 #Function for plotting impotant features with their corresponding probability values
2 def imp_features(features, indices, fea_names):
3
4     #sns.set(rc={'figure.figsize':(11.7,8.27)})
5
6     # Create plot
7     plt.figure()
8
9     # Create plot title
10    plt.title("Feature Importance")
11
12    # Add bars
13    plt.bar(range(20), features[indices])
14
15    # Add feature names as x-axis labels
16    names = np.array(fea_names)
17    plt.xticks(range(20), names[indices], rotation=90)
18
19    # Show plot
20    plt.show()
```

```
In [29]: 1 #Bar plot for important features with their probability values
        2 imp_features(features, indices, fea_names)
```



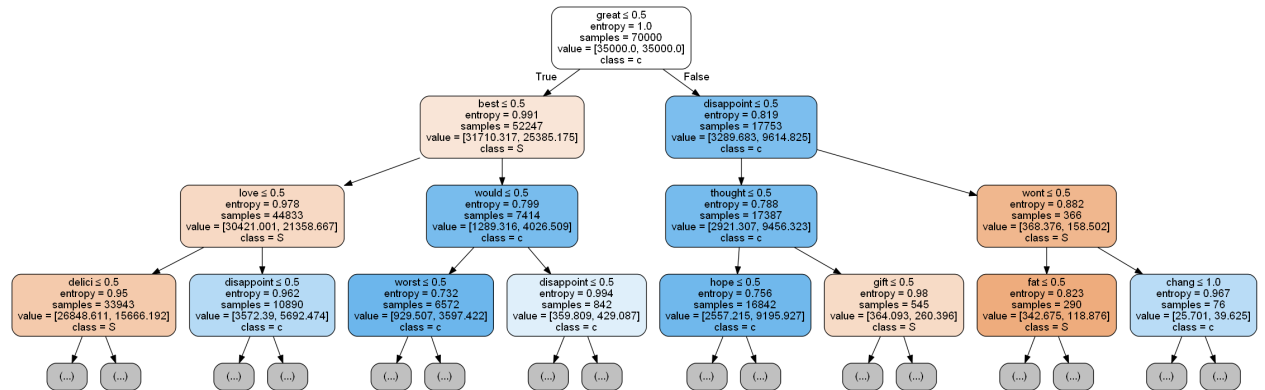
Graphviz:

Which can be used to visualize the trained decision tree, to know how model works to predict the response label.

```
In [30]: 1 #Function for graphical representation of the trained decision tree
        2 def graphviz(classifier, features, response_feature):
        3     dot_data = StringIO()
        4     export_graphviz(classifier, out_file=dot_data,max_depth=3,
        5                     feature_names= features,
        6                     class_names=response_feature,
        7                     filled=True, rounded=True,
        8                     special_characters=True)
        9     graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
        10
        11     return Image(graph.create_png())
```

```
In [31]: 1 #Graphical visualization of the trained Decision Tree
        2 graphviz(best_estimator, fea_names, data_frame_100k.columns[7])
```

Out[31]:



TFIDF:

TF-IDF stands for term frequency-inverse document frequency. TF-IDF weight is a statistical measure used to evaluate how important a word is to a document in a collection or corpus. The importance increases proportionally to the number of times a word appears in the document but is offset by the frequency of the word in the corpus.

Term_frequency(TF) = (number of times word occur in document) / (Total number of words in the document).

Inverse_Document_frequency(IDF) = log((total number of documents) / ln which documents a word occurs))

So, TF-IDF(word) = TF(word) * IDF(word)


```
In [32]: 1 #Vectorizing the data
2 tfidf_vect = TfidfVectorizer(ngram_range=(1,2))
3 tfidf_tr = tfidf_vect.fit_transform(X_tr)
```

```
In [33]: 1 #Vectorizing the test data
2 tfidf_test = tfidf_vect.transform(X_test)
```

```
In [34]: 1 #Shape of the train and test data
2 print("Shape of train data:", tfidf_tr.shape)
3 print("Shape of test data:", tfidf_test.shape)
```

Shape of train data: (70000, 932177)

Shape of test data: (30000, 932177)

Hyperparameter tuning using grid search cross validation:

Grid Search Cross Validation:

```
In [35]: 1
2 %%time
3
4 #Calling the function for Grid Search Cross Validation
5 grid_scores, best_estimator = grid_search(tfidf_tr, y_tr)
```

Fitting 5 folds for each of 20 candidates, totalling 100 fits

[Parallel(n_jobs=3)]: Done 26 tasks | elapsed: 2.7min

[Parallel(n_jobs=3)]: Done 100 out of 100 | elapsed: 18.3min finished

Max depth: {'max_depth': 20}

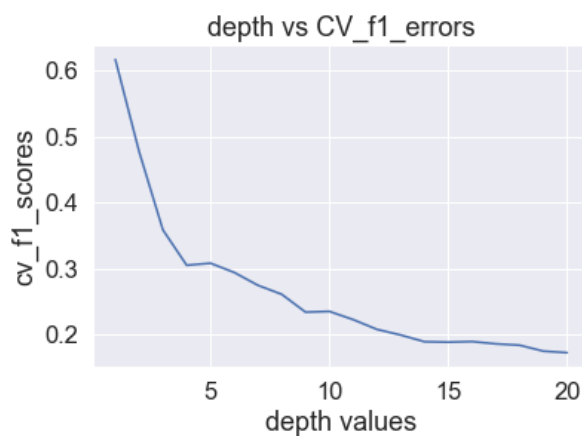
Best F1-score: 82.7953779627519

Wall time: 19min 10s

```
In [36]: 1 #Which can return the cv_scores, std_dev and max_depth for all hyper parameters
2 grid_scores[:2]
```

```
Out[36]: [mean: 0.38351, std: 0.02836, params: {'max_depth': 1},
mean: 0.52322, std: 0.02182, params: {'max_depth': 2}]
```

```
In [37]: 1 #From grid_scores storing max_depths into depth
2 depths = [val[0]['max_depth'] for val in grid_scores]
3
4 #From grid_scores storing mean f1_cross validation scores into cv_f1_scores
5 cv_f1_errors = [1-val[1] for val in grid_scores]
6
7 #plot for depths and cv_f1_errors
8 depth_scores(depths, cv_f1_errors)
```



Testing the model from best_estimator which can be return by the grid search cross validation.

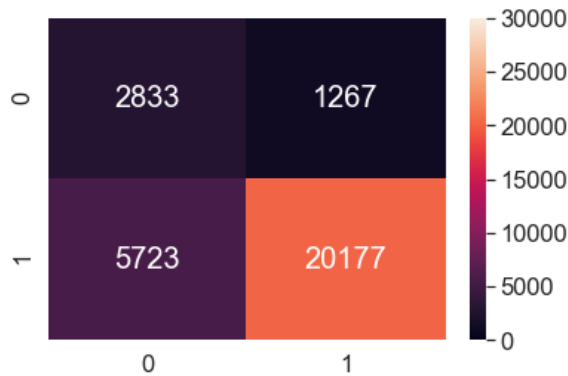
```
In [38]: 1 #Result showing the best classifier consisting of parameters
2 best_estimator
```

```
Out[38]: DecisionTreeClassifier(class_weight='balanced', criterion='entropy',
max_depth=20, max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best')
```

```
In [39]: 1 #Finding the predicted values for test labels using the test data
        2 y_pred = best_estimator.predict(tfidf_test)
```

```
In [40]: 1 #Calling the function for test metrics
        2 test_metrics(y_test, y_pred)
```

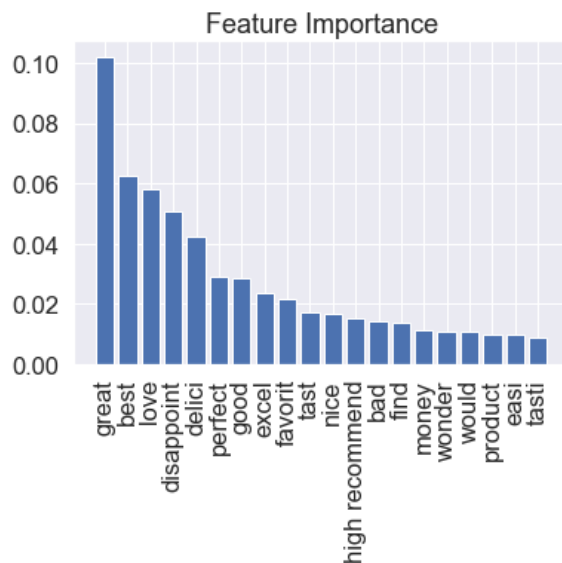
Accuracy on test data: 76.7
Precision on test data: 94.09
Recall on test data: 77.9
F1_score on test data: 85.24



Top 10 important features:

```
In [41]: 1 #Extracting important features from decision tree
        2 features = best_estimator.feature_importances_
        3
        4 #Getting the top important feature indices
        5 indices = np.argsort(features)[::-1][:20]
        6
        7 #Getting feature names from feature extractor object is count_vectorizer
        8 fea_names = tfidf_vect.get_feature_names()
```

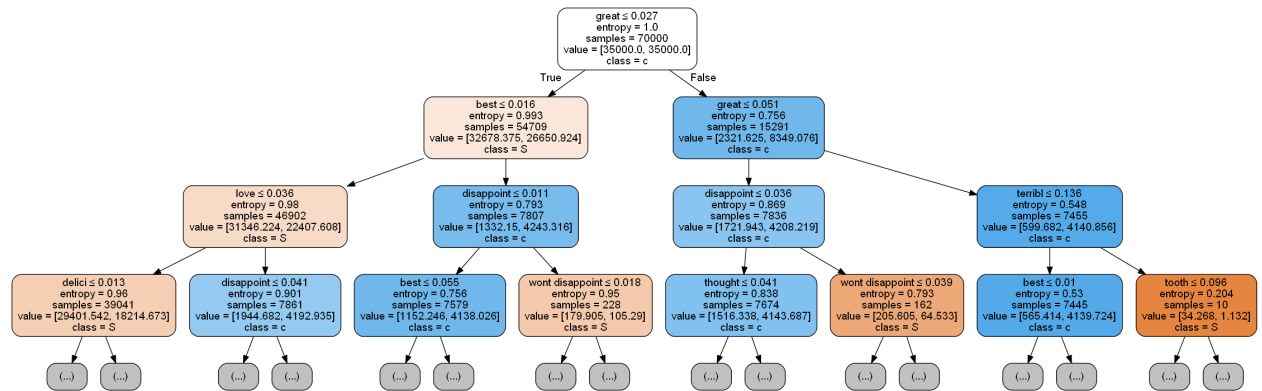
```
In [42]: 1 #Bar plot for important features with their probability values
        2 imp_features(features, indices, fea_names)
```



Visualizing the trained Decision Tree:

```
In [43]: 1 #Graphical visualization of the trained Decision Tree
2 graphviz(best_estimator, fea_names, data_frame_100k.columns[7])
```

Out[43]:



Avg_w2v:

1. W2V can take the semantic meaning of the words.
2. W2V can convert each word into an vector.
3. Avg_W2V means for each review vector should be $(W2V(word1) + W2V(word2) + \dots + W2V(wordn)) / (\text{total no. of words})$.

```
In [44]: 1 #Forming the list_of_words for 100k reviews
2 sent_words = []
3 for sent in X:
4     sent_words.append(sent.split())
```

```
In [45]: 1 #Splitting the into train and test data
2 X_tr_w2v, X_test_w2v, y_tr_w2v, y_test_w2v = train_test_split(sent_words, y, test_size=0.3, shuffle=False)
```

```
In [46]: 1 #Word to vectors for train data
2 w2v = gensim.models.Word2Vec(X_tr_w2v, min_count=5, size=50)
```

```
In [47]: 1 #storing w2v_words which can be return by w2v vocabulary
2 w2v_words = list(w2v.wv.vocab)
3 print("total words in w2v", len(w2v_words))
4 print(w2v_words[0:10])
```

total words in w2v 10701
['witti', 'littl', 'book', 'make', 'son', 'laugh', 'loud', 'car', 'drive', 'along']

```
In [48]: 1 #Function for Avg_w2v
2 def avg_w2v(data, w2v, w2v_words):
3     #creating an empty List
4     avg_vectors = []
5     row = 0
6     for sent in data:
7         #creating an vector which size should be 50 and all cells have zero's
8         sent_vec = np.zeros(50)
9         cnt_words = 0
10        for word in sent:
11            if word in w2v_words:
12                vec = w2v.wv[word]
13                sent_vec += vec
14                cnt_words += 1
15        if cnt_words != 0:
16            sent_vec /= cnt_words
17            avg_vectors.append(sent_vec)
18        row += 1
19        if cnt_words == 0:
20            print(row)
21    return avg_vectors
```

```
In [49]: 1
2 %%time
3
4 #Avg w2v for train data
5 X_tr_avg_w2v = avg_w2v(X_tr_w2v, w2v, w2v_words)
```

Wall time: 1min 7s

```
In [52]: 1 #Avg w2v for test data
2 X_test_avg_w2v = avg_w2v(X_test_w2v, w2v, w2v_words)
```

Hyperparameter tuning using grid search cross validation:

Grid Search Cross Validation:

```
In [54]: 1
2 %%time
3
4 #Calling the function for Grid Search Cross Validation
5 grid_scores, best_estimator = grid_search(X_tr_avg_w2v, y_tr)
```

Fitting 5 folds for each of 20 candidates, totalling 100 fits

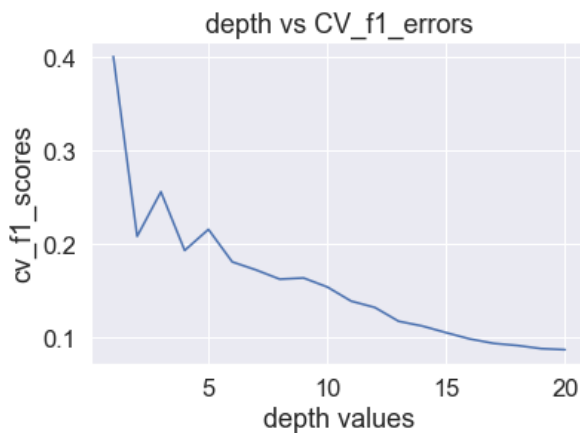
```
[Parallel(n_jobs=3)]: Done 26 tasks | elapsed: 27.2s
[Parallel(n_jobs=3)]: Done 100 out of 100 | elapsed: 1.9min finished
```

Max depth: {'max_depth': 20}
Best F1-score: 91.28832424195258
Wall time: 2min

```
In [55]: 1 #Which can return the cv_scores, std_dev and max_depth for all hyper parameters
2 grid_scores[:2]
```

```
Out[55]: [mean: 0.59932, std: 0.02166, params: {'max_depth': 1},
mean: 0.79168, std: 0.01023, params: {'max_depth': 2}]
```

```
In [56]: 1 #From grid_scores storing max_depths into depth
2 depths = [val[0]['max_depth'] for val in grid_scores]
3
4 #From grid_scores storing mean f1_cross validation scores into cv_f1_scores
5 cv_f1_errors = [1-val[1] for val in grid_scores]
6
7 #plot for depths and cv_f1_errors
8 depth_scores(depths, cv_f1_errors)
```



Testing the model from best_estimator which can be return by the grid search cross validation.

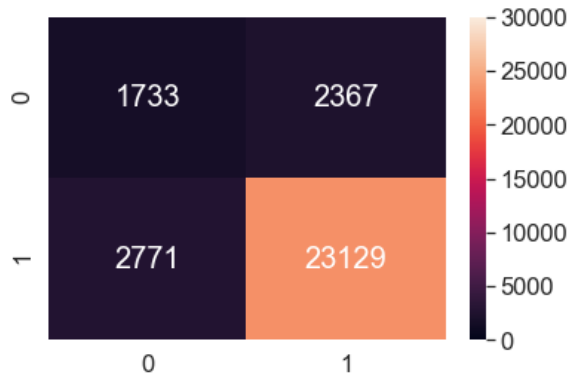
```
In [57]: 1 #Result showing the best classifier consisting of parameters
2 best_estimator
```

```
Out[57]: DecisionTreeClassifier(class_weight='balanced', criterion='entropy',
max_depth=20, max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best')
```

```
In [58]: 1 #Finding the predicted values for test labels using the test data
2 y_pred = best_estimator.predict(X_test_avg_w2v)
```

```
In [59]: 1 #Calling the function for test metrics
        2 test_metrics(y_test, y_pred)
```

Accuracy on test data: 82.87
Precision on test data: 90.72
Recall on test data: 89.3
F1_score on test data: 90.0



TF-IDF W2V:

```
In [60]: 1 #Previously done tfidf_w2v with the 50k points
        2 #Getting the train data
        3 tfidf_w2v_tr = openfromfile("tfidf_w2v_train_of_50k_pts")
        4 y_tr_w2v = openfromfile("tfidf_y_tr_w2v_of_50k_pts")
```

```
In [61]: 1 #Getting the test data
        2 tfidf_w2v_test = openfromfile("tfidf_w2v_test_of_50k_pts")
        3 y_test_w2v = openfromfile("tfidf_y_test_w2v_of_50k_pts")
```

```
In [62]: 1 #Shape of the train and test data
        2 print("Length of the train data:", len(tfidf_w2v_tr))
        3 print("Length of the test data:", len(tfidf_w2v_test))
```

Length of the train data: 35000
Length of the test data: 10164

Hyperparameter tuning using grid search cross validation:

Grid Search Cross Validation:

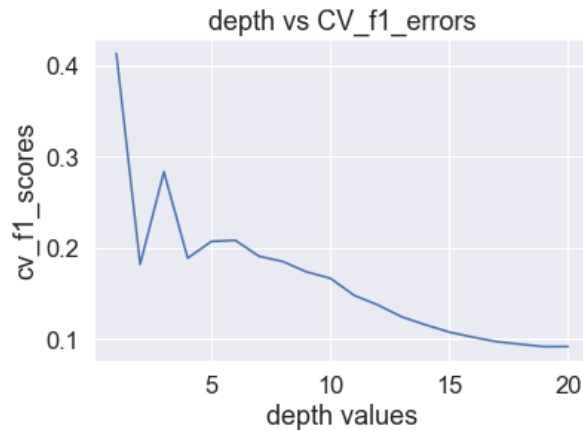
```
In [63]: 1
        2 %%time
        3
        4 #Calling the function for Grid Search Cross Validation
        5 grid_scores, best_estimator = grid_search(tfidf_w2v_tr, y_tr_w2v)
```

Fitting 5 folds for each of 20 candidates, totalling 100 fits

[Parallel(n_jobs=3)]: Done 26 tasks | elapsed: 14.1s
[Parallel(n_jobs=3)]: Done 100 out of 100 | elapsed: 54.6s finished

Max depth: {'max_depth': 19}
Best F1-score: 90.80168522657921
Wall time: 58.1 s

```
In [64]: 1 #From grid_scores storing max_depths into depth
2 depths = [val[0]['max_depth'] for val in grid_scores]
3
4 #From grid_scores storing mean f1_cross validation scores into cv_f1_scores
5 cv_f1_errors = [1-val[1] for val in grid_scores]
6
7 #plot for depths and cv_f1_errors
8 depth_scores(depths, cv_f1_errors)
```



Testing the model from best_estimator which can be return by the grid search cross validation.

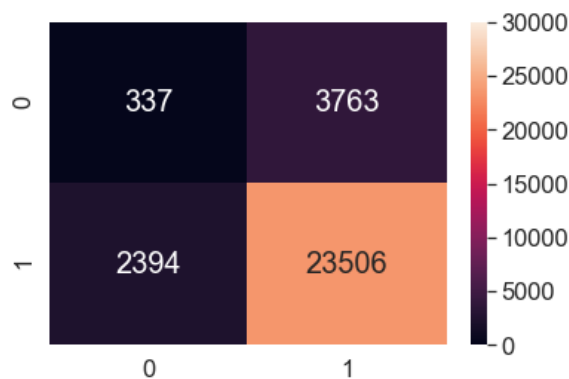
```
In [65]: 1 #Result showing the best classifier consisting of parameters
2 best_estimator
```

```
Out[65]: DecisionTreeClassifier(class_weight='balanced', criterion='entropy',
max_depth=19, max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best')
```

```
In [66]: 1 #Finding the predicted values for test labels using the test data
2 y_pred = best_estimator.predict(X_test_avg_w2v)
```

```
In [67]: 1 #Calling the function for test metrics
2 test_metrics(y_test, y_pred)
```

Accuracy on test data: 79.48
Precision on test data: 86.2
Recall on test data: 90.76
F1_score on test data: 88.42



Summary:

Performance Table:

- Given dataset is an imbalanced data which means majority class is positive and minority class is negative, so i make balanced data using class_weight is equal to balanced.
- Measurring the quality of split i used entropy.

Featurization	sample size	CV	Accuracy	F1-score	Max_depth
		Test accuracy	Test f1-score		

Featurization	sample size	CV	Accuracy	F1-score	Max_depth
BOW	100k	Grid Search	77.13%	85.60%	20
TF-IDF	100k	Grid Search	76.70%	85.24%	20
Avg-W2V	100k	Grid Search	82.87%	90.00%	20
TF-IDF W2V	50k	Grid Search	79.48%	88.42%	19

observation:

- Among all text classifications, Avg_w2v is working well for this dataset.
- In case of BOW and TFIDF most important feature is "great".