

Assignment 5 : Logistic Regression

- Logistic regression is a statistical method for analyzing a dataset in which there are one or more independent variables that determine an outcome. The outcome is measured with a binary variable.
- The relationship between X and Y is non_linear then we use Logistic regression.
- It is used to describe data and to explain the relationship between one dependent binary variable and one or more nominal, ordinal, interval or ratio-level independent variables.

$$\log(p(x) / (p(1-x))) = w_0 + w_1 x_1$$

$$\text{so, } p(x) = 1 / (1 + \exp(-(w_0 + w_1 x_1)))$$

- Dependent variable must be bernouli distribution.
- It does not assume linear relationship between dependent and independent variable.
- It is robust as it does not require independent variables to follow normal distribution.

```
In [5]: 1 #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 [6]: 1 #Loading the Libraries
        2 import pandas as pd
        3 import numpy as np
        4 import matplotlib.pyplot as plt
        5 import seaborn as sns
        6
        7 from sklearn.model_selection import train_test_split
        8 from sklearn.feature_extraction.text import CountVectorizer
        9 from sklearn import preprocessing
        10
        11 from sklearn.model_selection import TimeSeriesSplit
        12 from sklearn.model_selection import GridSearchCV
        13 from sklearn.linear_model import LogisticRegression
        14
        15 from sklearn.metrics import confusion_matrix
        16 from sklearn.metrics import accuracy_score
        17 from sklearn.metrics import recall_score
        18 from sklearn.metrics import precision_score
        19 from sklearn.metrics import f1_score
        20
        21 from sklearn.model_selection import RandomizedSearchCV
        22 from scipy.sparse import find
        23 from numpy import random
        24 from sklearn.feature_extraction.text import TfidfVectorizer
        25 import gensim
```

```
In [7]: 1 #Loading the dataset
        2 data_frame = openfromfile("New_Amazon_preprocess_data")
```

```
In [8]: 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)

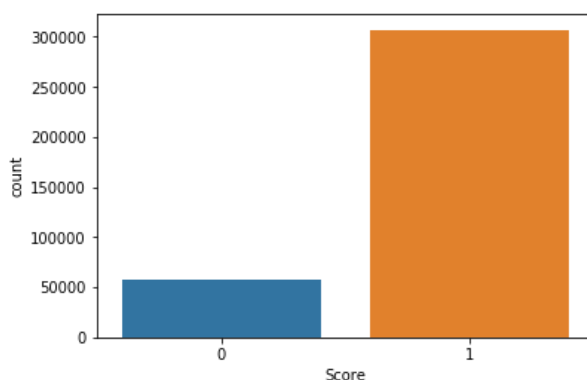
Out[8]:

		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	D
24749	24750	2734888454	A13ISQV0U9GZIC	Sandikaye		1	1	negative	1192060800	
308076	308077	2841233731	A3QD68O22M2XHQ	LABRNTN		0	0	positive	1345852800	b
150523	150524	6641040	ACITT7DI6IDDL	shari zychinski		0	0	positive	939340800	edt

```
In [9]: 1 #Storing the data_frame based on the time attribute
2 data_frame.sort_values('Time', inplace=True)
3
4 #Reseting the data_frame
5 data_frame.reset_index(drop=False, inplace=True)
```

```
In [10]: 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 [11]: 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[11]: 1    307061
0     57110
Name: Score, dtype: int64
```

```
In [12]: 1 #Storing the cleanedtext attribute into X and Score attribute into the Y
2 X = data_frame.CleanedText
3
4 y = data_frame.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: (254919,)

Shape of the test data: (109252,)

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: 8.61 s

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

```
In [17]: 1 #Normalizing the train and test data
2 bow_tr = preprocessing.normalize(bow_tr)
3 bow_test = preprocessing.normalize(bow_test)
```

```
In [18]: 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)
```

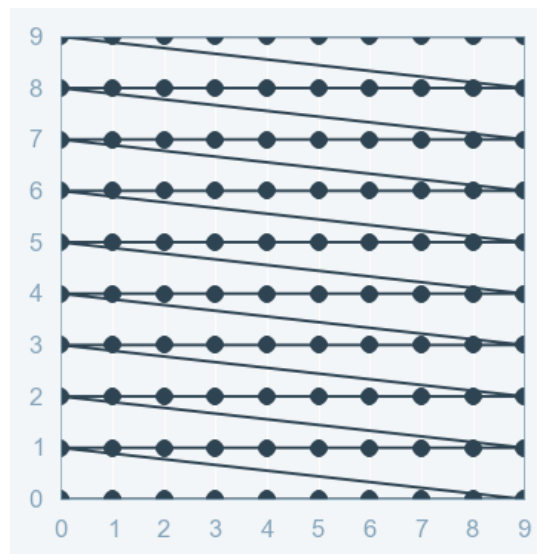
shape of train data: (254919, 78610)

shape of test data: (109252, 78610)

1. Hyperparameter tuning using grid search and random search cross validation:

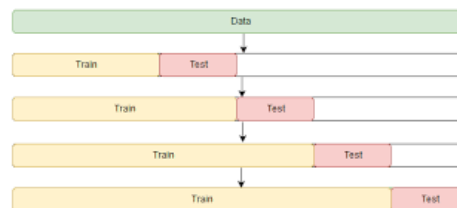
Grid Search Cross Validation:

- 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 [38]: 1 #Function for Grid Search Cross Validation
2 def grid_search(X_train, y_train):
3     #Assigning the values for hyperparameter and regularization as l1 and l2
4     parameters = {'C':[1000,100,10,5,1,0.5,0.1,0.05,0.01,0.005,0.001], 'penalty':['l1','l2']}
5
6     #splitting the data based on the time series
7     tbs = TimeSeriesSplit(n_splits=5)
8
9     log_model = LogisticRegression()
10
11     #Grid Search Cross Validation using Logistic regression
12     gsv = GridSearchCV(log_model, parameters, scoring='accuracy', n_jobs=3, cv=tbs, verbose=3)
13     gsv.fit(X_train, y_train)
14
15     #Best hyperparameter value
16     print("optimal hyperparameter:", gsv.best_params_)
17     print("Best Accuracy:", gsv.best_score_ * 100)
18
19     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 22 candidates, totalling 110 fits

```
[Parallel(n_jobs=3)]: Done 26 tasks      | elapsed: 4.9min
[Parallel(n_jobs=3)]: Done 110 out of 110 | elapsed: 7.2min finished
```

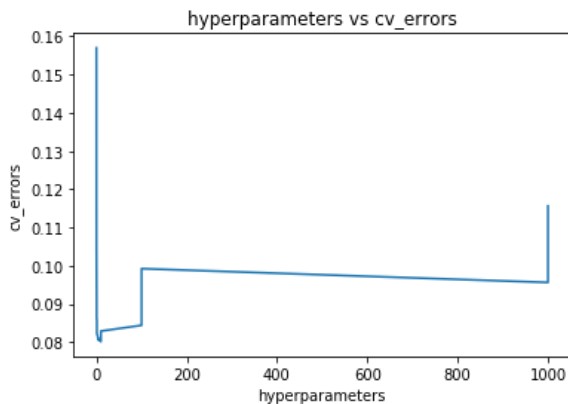
```
optimal hyperparameter: {'C': 10, 'penalty': 'l2'}
Best Accuracy: 91.9818293084781
Wall time: 7min 32s
```

```
In [20]: 1 grid_scores[:3]
```

```
Out[20]: [mean: 0.88436, std: 0.00966, params: {'C': 1000, 'penalty': 'l1'},
          mean: 0.90436, std: 0.00334, params: {'C': 1000, 'penalty': 'l2'},
          mean: 0.90075, std: 0.00446, params: {'C': 100, 'penalty': 'l1'}]
```

```
In [39]: 1 #Function for plot between hyperparameter C and cv_errors
2
3 def param_cv_error(hyperparameters, cv_errors):
4     plt.plot(hyperparameters, cv_errors)
5     plt.title("hyperparameters vs cv_errors")
6     plt.xlabel("hyperparameters")
7     plt.ylabel("cv_errors")
8     plt.show()
```

```
In [22]: 1
2 hy_params = [val[0]['C'] for val in grid_scores]
3
4 cv_errors = [1-val[1] for val in grid_scores]
5
6 #Calling the function for plot between C and cv_errors
7 param_cv_error(hy_params, cv_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
```

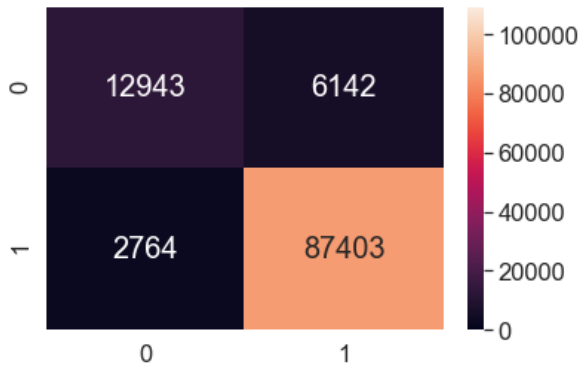
```
Out[23]: LogisticRegression(C=10, class_weight=None, dual=False, fit_intercept=True,
    intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
    penalty='l2', random_state=None, solver='liblinear', tol=0.0001,
    verbose=0, warm_start=False)
```

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

```
In [62]: 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=109252)
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()
```

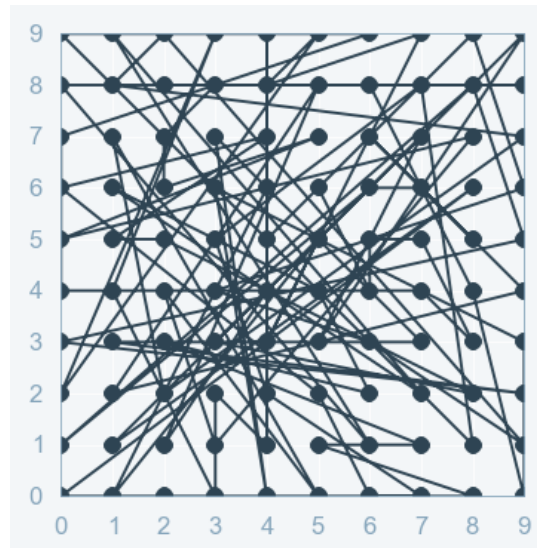
```
In [26]: 1 #Calling the function for test metrics
        2 test_metrics(y_test, y_pred)
```

Accuracy on test data: 91.85
Precision on test data: 93.43
Recall on test data: 96.93
F1_score on test data: 95.15



Random Search cross validation:

- implements a randomized search over parameters, where each setting is sampled from a distribution over possible parameter values.
- This has two main benefits over an exhaustive search:
 1. A budget can be chosen independent of the number of parameters and possible values
 2. Adding parameters that do not influence the performance does not decrease efficiency.
- Note: In random search, if you choose n parameters then we will have to check n combinations.



```
In [40]: 1 #Function for random Search Cross Validation
        2
        3 def random_search(X_train, y_train):
        4     #Assigning the values for hyperparameter and regularization as L1 and L2
        5     parameters = {'C': np.arange(1, 1000), 'penalty':['l1','l2']}
        6
        7     #splitting the data based on the time series
        8     tbs = TimeSeriesSplit(n_splits=5)
        9
        10    log_model = LogisticRegression()
        11    #Random search for hyperparameter tuning
        12    rsv = RandomizedSearchCV(log_model, parameters, scoring='accuracy', n_jobs=3, cv=tbs, verbose=3)
        13
        14    rsv.fit(X_train, y_train)
        15
        16    print("optimal hyperparameter:",rsv.best_params_)
        17    print("Best accuracy:",rsv.best_score_*100)
        18    return rsv.best_estimator_
```

```
In [28]: 1
2 %%time
3
4 #Calling the function for random search cross validation
5 best_estimator = random_search(bow_tr, y_tr)
```

Fitting 5 folds for each of 10 candidates, totalling 50 fits

```
[Parallel(n_jobs=3)]: Done 26 tasks | elapsed: 3.2min
[Parallel(n_jobs=3)]: Done 50 out of 50 | elapsed: 11.3min finished
```

```
optimal hyperparameter: {'penalty': 'l2', 'C': 239}
Best accuracy: 91.22864002259567
Wall time: 12min 38s
```

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

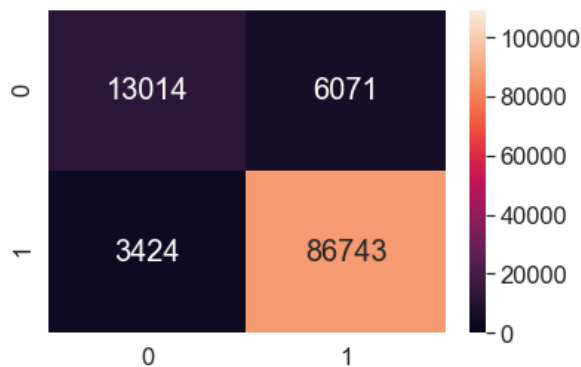
```
In [29]: 1 #Result showing the best classifier consisting of parameters
2 best_estimator
```

```
Out[29]: LogisticRegression(C=239, class_weight=None, dual=False, fit_intercept=True,
intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
penalty='l2', random_state=None, solver='liblinear', tol=0.0001,
verbose=0, warm_start=False)
```

```
In [30]: 1 #Finding the predicted values for test labels using the test data
2 y_pred = best_estimator.predict(bow_test)
```

```
In [31]: 1 #Calling the function for test metrics
2 test_metrics(y_test, y_pred)
```

```
Accuracy on test data: 91.31
Precision on test data: 93.46
Recall on test data: 96.2
F1_score on test data: 94.81
```



2. More Sparsity(Fewer elements of w^* being non-zero) as C decreases or lambda increases using L1-regularization:

```
In [41]: 1 #Function for number of non-zero elements in an vector
2 def non_zero_ele(X_train, y_train, C_value):
3     #Specifying the classifier with the hyperparameter and L1-regularization
4     clf = LogisticRegression(penalty='l1', C=C_value)
5     clf.fit(X_train, y_train)
6     #optimal weight vector
7     opt_w = clf.coef_
8     #number of non-zero elements in an optimal vector
9     print("Number of non_zero elements in optimal vector for C={} and l1_reg is {}".format(C, np.count_nonzero(opt_w)))
```

```
In [33]: 1 #To find number of non_zero elements with the C=10 and l1_reg
2 C=10
3 non_zero_ele(bow_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=10 and l1_reg is 11080:

```
In [35]: 1 #To find number of non_zero elements with the C=1 and l1_reg
2 C=1
3 non_zero_ele(bow_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=1 and l1_reg is 2229:

```
In [36]: 1 #To find number of non_zero elements with the C=0.1 and l1_reg
        2 C=0.1
        3 non_zero_ele(bow_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=0.1 and l1_reg is 522:

```
In [37]: 1 #To find number of non_zero elements with the C=0.01 and l1_reg
        2 C=0.01
        3 non_zero_ele(bow_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=0.01 and l1_reg is 75:

```
In [38]: 1 #To find number of non_zero elements with the C=0.001 and l1_reg
        2 C=0.001
        3 non_zero_ele(bow_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=0.001 and l1_reg is 2:

```
In [40]: 1 #To find number of non_zero elements with the C=0.01 and l1_reg
        2 C=0.0001
        3 non_zero_ele(bow_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=0.0001 and l1_reg is 0:

observation: By using L1_regularization with the different C values of as C decreases which means alpha increases then the number of non-zero elements in optimal vector decreases.

3. Perturbation Test (Multi collinearity test):

- In case of logistic regression to find important features, firstly we have to check multi collinearity between features, if features are collinear then we should find important features using forward or backward feature selection.
- If features are not correlated then we should use optimal vector, in which consist of weight for each feature.
- **Multi collinearity:** which means very high inter correlation among the independent variables.
 - It occurs: 1. Inaccurate use of dummy variables.
 - 2. When we find one variable with help of other variables.
 - 3. Also occurs when repetition of same kind of variables.
 - 4. Variables are highly correlated to each other.
- We can also see multi collinearity by using an correlation matrix for independent variables.

```
In [19]: 1 #Training the model before adding the epsilon to training data
        2 clf = LogisticRegression(penalty='l2', C=239)
        3 clf.fit(bow_tr, y_tr)
        4
        5 #find function takes input sparse or dense matrix and it can return row indices at 0, column indices at 1
        6 #and values of non-zero elements at 2
        7 weights = find(clf.coef_[0])[2]
        8
```

```
In [20]: 1 bow_t = bow_tr
        2
        3 #Generating the random noise with N(0, 0.001), Length should be non_zero elements in bow_tr
        4 epsilon = random.normal(0, 0.1, find(bow_t)[0].size)
        5
        6 #Storing row, column and non_elements
        7 a, b, c = find(bow_t)
        8
        9 #adding epsilon to all non_zero elements
        10 bow_t[a, b] = epsilon + bow_t[a, b]
```

```
In [21]: 1 #Training the model after adding the epsilon to training data
        2 clf1 = LogisticRegression(penalty='l2', C=214)
        3 clf1.fit(bow_t, y_tr)
        4
        5 #find function takes input sparse or dense matrix and it can return row indices at 0, column indices at 1
        6 #and values of non-zero elements at 2
        7 weights1 = find(clf1.coef_[0])[2]
        8
```

```
In [22]: 1 #Calculating the difference
        2 weight_diff = (abs(weights - weights1) / weights) * 100
```



```
In [25]: 1 print(weight_diff[np.where(weight_diff > 30)].size)
```

41794

Observation: Attributes are collinear to each other.

4. Top 10 important Features for positive and negative classes:

```
In [46]: 1 #Getting the feature names from the count_vec
2 features = count_vec.get_feature_names()
3
4 #Combining coefficient values with the corresponding features
5 coefs_with_fea = sorted(zip(clf.coef_[0], features))
6
7 print("Top 10 positive important features:", [a[1] for a in coefs_with_fea[-10:]])
8 print("***100")
9 print("Top 10 negative important features:", [a[1] for a in coefs_with_fea[:10]])
```

Top 10 positive important features: ['fancier', 'goshoptnt', 'reinstat', 'borer', 'understandng', 'emeraldfor
est', 'fier', 'compait', 'brothi', 'terror']

Top 10 negative important features: ['yadayadayada', 'gind', 'plumros', 'accel', 'cocca', 'weakest', 'wrongt
h', 'holl', 'onward', 'berryblossom']

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) / In which documents a word occurs))

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

```
In [26]: 1 #Vectorizing the data
2 tfidf_vect = TfidfVectorizer(ngram_range=(1,2))
3 tfidf_tr = tfidf_vect.fit_transform(X_tr)
```

```
In [27]: 1 #Vectorizing the test data
2 tfidf_test = tfidf_vect.transform(X_test)
```

```
In [28]: 1 #Normalizing the train and test data
2 tfidf_tr = preprocessing.normalize(tfidf_tr)
3 tfidf_test = preprocessing.normalize(tfidf_test)
```

```
In [29]: 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: (254919, 2362093)

Shape of test data: (109252, 2362093)

1. Hyperparameter tuning using grid search and random search cross validation:

Grid Search Cross Validation:

```
In [51]: 1 #Calling the Grid search function
2 grid_scores, best_estimator = grid_search(tfidf_tr, y_tr)
```

Fitting 5 folds for each of 22 candidates, totalling 110 fits

[Parallel(n_jobs=3)]: Done 26 tasks | elapsed: 4.2min

[Parallel(n_jobs=3)]: Done 110 out of 110 | elapsed: 9.6min finished

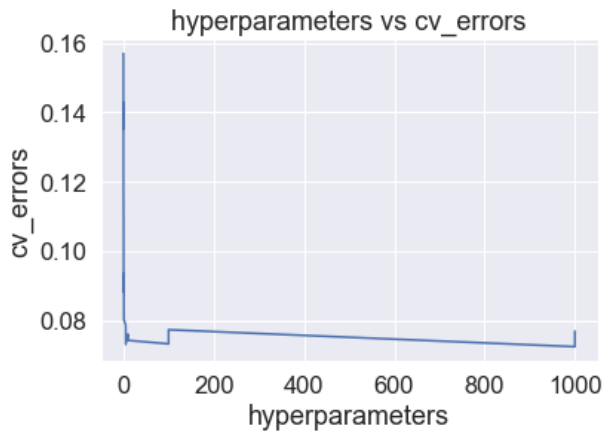
optimal hyperparameter: {'C': 1000, 'penalty': 'l2'}

Best Accuracy: 92.75761427293698

```
In [52]: 1 grid_scores[:2]
```

```
Out[52]: [mean: 0.92310, std: 0.00205, params: {'C': 1000, 'penalty': 'l1'},  
          mean: 0.92758, std: 0.00356, params: {'C': 1000, 'penalty': 'l2'}]
```

```
In [53]: 1  
2 hy_params = [val[0]['C'] for val in grid_scores]  
3  
4 cv_errors = [1-val[1] for val in grid_scores]  
5  
6 #Calling the function for plot between C and cv_errors  
7 param_cv_error(hy_params, cv_errors)
```



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

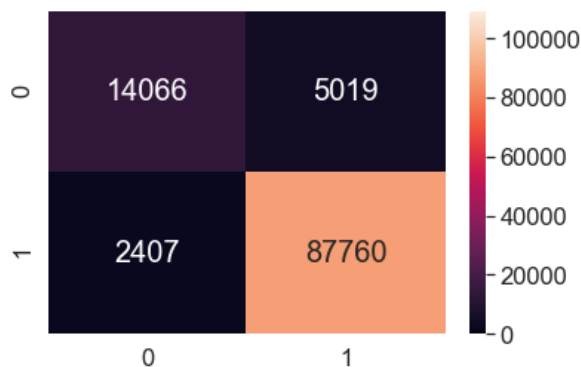
```
In [54]: 1 #Result showing the best classifier consisting of parameters  
2 best_estimator
```

```
Out[54]: LogisticRegression(C=1000, class_weight=None, dual=False, fit_intercept=True,  
                             intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,  
                             penalty='l2', random_state=None, solver='liblinear', tol=0.0001,  
                             verbose=0, warm_start=False)
```

```
In [55]: 1 #Finding the predicted values for test labels using the test data  
2 y_pred = best_estimator.predict(tfidf_test)
```

```
In [56]: 1 #Calling the function for test metrics  
2 test_metrics(y_test, y_pred)
```

Accuracy on test data: 93.2
Precision on test data: 94.59
Recall on test data: 97.33
F1_score on test data: 95.94



Random Search Cross Validation:

```
In [57]: 1 #Calling the function for random search cross validation
        2 best_estimator = random_search(tfidf_tr, y_tr)
```

Fitting 5 folds for each of 10 candidates, totalling 50 fits

[Parallel(n_jobs=3)]: Done 26 tasks | elapsed: 3.6min

[Parallel(n_jobs=3)]: Done 50 out of 50 | elapsed: 7.0min finished

optimal hyperparameter: {'penalty': 'l2', 'C': 250}

Best accuracy: 92.71807183542813

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

```
In [58]: 1 #Result showing the best classifier consisting of parameters
        2 best_estimator
```

```
Out[58]: LogisticRegression(C=250, class_weight=None, dual=False, fit_intercept=True,
                             intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
                             penalty='l2', random_state=None, solver='liblinear', tol=0.0001,
                             verbose=0, warm_start=False)
```

```
In [59]: 1 #Finding the predicted labels for test data
        2 y_pred = best_estimator.predict(tfidf_test)
```

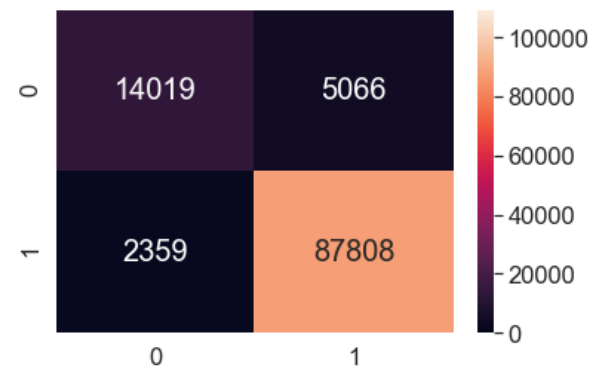
```
In [60]: 1 #Calling the function for test metrics
        2 test_metrics(y_test, y_pred)
```

Accuracy on test data: 93.2

Precision on test data: 94.55

Recall on test data: 97.38

F1_score on test data: 95.94



2. More Sparsity(Fewer elements of w^* being non-zero) as C decreases or lambda increases using L1-regularization:

```
In [61]: 1 #To find number of non_zero elements with the C=10 and l1_reg
        2 C=10
        3 non_zero_ele(tfidf_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=10 and l1_reg is 32358:

```
In [62]: 1 #To find number of non_zero elements with the C=10 and l1_reg
        2 C=5
        3 non_zero_ele(tfidf_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=5 and l1_reg is 21163:

```
In [63]: 1 #To find number of non_zero elements with the C=10 and l1_reg
        2 C=1
        3 non_zero_ele(tfidf_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=1 and l1_reg is 2656:

```
In [64]: 1 #To find number of non_zero elements with the C=10 and l1_reg
        2 C=0.1
        3 non_zero_ele(tfidf_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=0.1 and l1_reg is 339:

```
In [65]: 1 #To find number of non_zero elements with the C=10 and l1_reg
2 C=0.01
3 non_zero_ele(tfidf_tr, y_tr, C)
```

Number of non_zero elements in optimal vector for C=0.01 and l1_reg is 20:

3. Pertubation Test (Multi collinearity test):

```
In [31]: 1 #Training the model before adding the epsilon to training data
2 clf = LogisticRegression(penalty='l2', C=1000)
3 clf.fit(tfidf_tr, y_tr)
4
5 #find function takes input sparse or dense matrix and it can return row indices at 0, column indices at 1
6 #and values of non-zero elements at 2
7 weights = find(clf.coef_[0])[2]
8
```

```
In [32]: 1 tfidf_t = tfidf_tr
2
3 #Generating the random noise with N(0, 0.001), Length should be non_zero elements in bow_tr
4 epsilon = random.normal(0, 0.1, find(tfidf_t)[0].size)
5
6 #Storing row, column and non_elements
7 a, b, c = find(tfidf_t)
8
9 #adding epsilon to all non_zero elements
10 tfidf_t[a, b] = epsilon + tfidf_t[a, b]
```

```
In [33]: 1 #Training the model after adding the epsilon to training data
2 clf1 = LogisticRegression(penalty='l2', C=1000)
3 clf1.fit(tfidf_t, y_tr)
4
5 #find function takes input sparse or dense matrix and it can return row indices at 0, column indices at 1
6 #and values of non-zero elements at 2
7 weights1 = find(clf1.coef_[0])[2]
8
```

```
In [34]: 1 weight_diff = (abs(weights - weights1)/weights) * 100
```

```
In [36]: 1 print(weight_diff[np.where(weight_diff > 50)].size)
```

1322430

4. Top 10 important Features for positive and negative classes:

```
In [71]: 1 #Getting the feature names from the count_vec
2 features = tfidf_vect.get_feature_names()
3
4 #Combining coefficient values with the corresponding features
5 coefs_with_fea = sorted(zip(clf.coef_[0], features))
6
7 print("Top 10 positive important features:", [a[1] for a in coefs_with_fea[-10:]])
8 print("***100)
9 print("Top 10 negative important features:", [a[1] for a in coefs_with_fea[:10]])
```

Top 10 positive important features: ['awesom', 'amaz', 'wont disappoint', 'excel', 'love', 'perfect', 'best', 'high recommend', 'delici', 'great']

Top 10 negative important features: ['worst', 'two star', 'disappoint', 'terribl', 'aw', 'horribl', 'threw', 'disgust', 'return', 'bland']

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 $(W2V(\text{word1}) + W2V(\text{word2}) + \dots + W2V(\text{wordn})) / (\text{total no.of words})$.

```

In [43]: 1 #Forming the list_of_words for 50k reviews
          2 sent_words = []
          3 for sent in X:
          4     sent_words.append(sent.split())

In [44]: 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 [45]: 1 #Word to vectors for train data
          2 w2v = gensim.models.Word2Vec(X_tr_w2v,min_count=5,size=50)

In [46]: 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 19107
['witti', 'littl', 'book', 'make', 'son', 'laugh', 'loud', 'recit', 'car', 'drive']

In [47]: 1 #Function for Avg_w2v
          2 def avg_w2v(data, w2v, w2v_words):
          3
          4     avg_vectors = [] #creating an empty list
          5     row = 0
          6     for sent in data:
          7
          8         sent_vec = np.zeros(50) #creating an vector which size should be 50 and all cells have zero's
          9         cnt_words = 0
          10        for word in sent: #From each sentence taking word
          11            if word in w2v_words:
          12                vec = w2v.wv[word] #Creating vector for word
          13                sent_vec += vec #Combining all word vectors to create sentence vector
          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 [48]: 1
          2 %%time
          3
          4 #Avg w2v for train data
          5 X_tr_avg_w2v = avg_w2v(X_tr_w2v, w2v, w2v_words)

192902
227729
Wall time: 4min 36s

In [51]: 1 #Dropping the Labels of corresponding test reviews whose words not match with train vocabulary
          2 y_tr_w2v.drop(labels=[192901, 227728], inplace=True)

In [53]: 1 #Avg w2v for test data
          2 X_test_avg_w2v = avg_w2v(X_test_w2v, w2v, w2v_words)

3296
101568

In [54]: 1 #Storing y_test_w2v indices because corresponding vectors are empty
          2 r1 = 254918 + 3296
          3 r2 = 254918 + 101568

In [55]: 1 #Dropping the Labels of corresponding test reviews whose words not match with train vocabulary
          2 y_test_w2v.drop(labels=[r1, r2], inplace=True)

```

Hyperparameter tuning using grid search and random search cross validation:

Grid Search Cross Validation:

```
In [56]: 1 #Calling the Grid search function
        2 grid_scores, best_estimator = grid_search(X_tr_avg_w2v, y_tr_w2v)
```

Fitting 5 folds for each of 22 candidates, totalling 110 fits

[Parallel(n_jobs=3)]: Done 26 tasks | elapsed: 6.6min

[Parallel(n_jobs=3)]: Done 110 out of 110 | elapsed: 16.1min finished

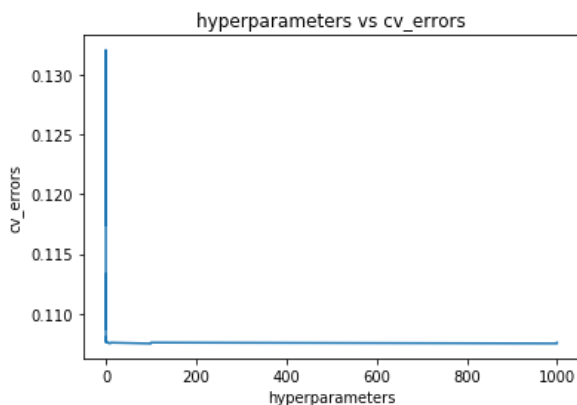
optimal hyperparameter: {'C': 100, 'penalty': 'l2'}

Best Accuracy: 89.250576660547

```
In [57]: 1 grid_scores[:2]
```

```
Out[57]: [mean: 0.89243, std: 0.00638, params: {'C': 1000, 'penalty': 'l1'},
         mean: 0.89250, std: 0.00634, params: {'C': 1000, 'penalty': 'l2'}]
```

```
In [58]: 1
        2 hy_params = [val[0]['C'] for val in grid_scores]
        3
        4 cv_errors = [1-val[1] for val in grid_scores]
        5
        6 #Calling the function for plot between C and cv_errors
        7 param_cv_error(hy_params, cv_errors)
```



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

```
In [59]: 1 #Result showing the best classifier consisting of parameters
        2 best_estimator
```

```
Out[59]: LogisticRegression(C=100, class_weight=None, dual=False, fit_intercept=True,
                             intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
                             penalty='l2', random_state=None, solver='liblinear', tol=0.0001,
                             verbose=0, warm_start=False)
```

```
In [60]: 1 y_pred = best_estimator.predict(X_test_avg_w2v)
```

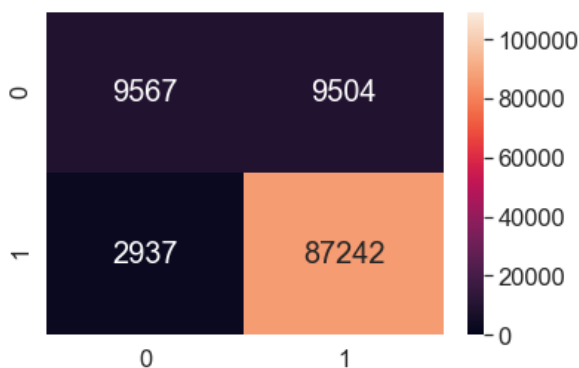
```
In [63]: 1 #Calling the function for test metrics
        2 test_metrics(y_test_w2v, y_pred)
```

Accuracy on test data: 88.61

Precision on test data: 90.18

Recall on test data: 96.74

F1_score on test data: 93.34



Random Search Cross Validation:

```
In [64]: 1 #Calling the function for random search cross validation
        2 best_estimator = random_search(X_tr_avg_w2v, y_tr_w2v)
```

Fitting 5 folds for each of 10 candidates, totalling 50 fits

[Parallel(n_jobs=3)]: Done 26 tasks | elapsed: 4.7min

[Parallel(n_jobs=3)]: Done 50 out of 50 | elapsed: 9.9min finished

optimal hyperparameter: {'penalty': 'l2', 'C': 72}

Best accuracy: 89.250576660547

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

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

```
Out[65]: LogisticRegression(C=72, class_weight=None, dual=False, fit_intercept=True,
                             intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
                             penalty='l2', random_state=None, solver='liblinear', tol=0.0001,
                             verbose=0, warm_start=False)
```

```
In [66]: 1 y_pred = best_estimator.predict(X_test_avg_w2v)
```

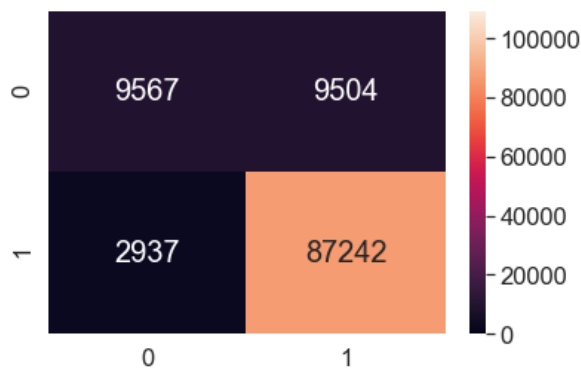
```
In [67]: 1 #Calling the function for test metrics
        2 test_metrics(y_test_w2v, y_pred)
```

Accuracy on test data: 88.61

Precision on test data: 90.18

Recall on test data: 96.74

F1_score on test data: 93.34



TFIDF-w2v:

```
In [68]: 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 [69]: 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 [70]: 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 seach and random search cross validation:

Grid Search Cross Validation:

```
In [71]: 1 #Calling the Grid search function
        2 grid_scores, best_estimator = grid_search(tfidf_w2v_tr, y_tr_w2v)
```

Fitting 5 folds for each of 22 candidates, totalling 110 fits

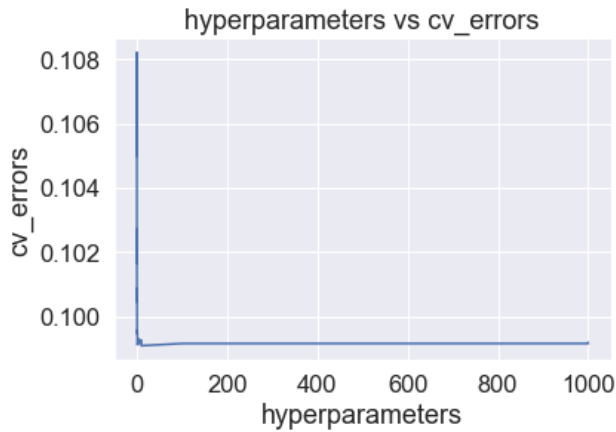
[Parallel(n_jobs=3)]: Done 26 tasks | elapsed: 17.8s

[Parallel(n_jobs=3)]: Done 110 out of 110 | elapsed: 1.0min finished

optimal hyperparameter: {'C': 10, 'penalty': 'l1'}

Best Accuracy: 90.09086233499058

```
In [72]: 1
        2 hy_params = [val[0]['C'] for val in grid_scores]
        3
        4 cv_errors = [1-val[1] for val in grid_scores]
        5
        6 #Calling the function for plot between C and cv_errors
        7 param_cv_error(hy_params, cv_errors)
```



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

```
In [73]: 1 #Result showing the best classifier consisting of parameters
        2 best_estimator
```

```
Out[73]: LogisticRegression(C=10, class_weight=None, dual=False, fit_intercept=True,
                             intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
                             penalty='l1', random_state=None, solver='liblinear', tol=0.0001,
                             verbose=0, warm_start=False)
```

```
In [74]: 1 y_pred = best_estimator.predict(tfidf_w2v_test)
```

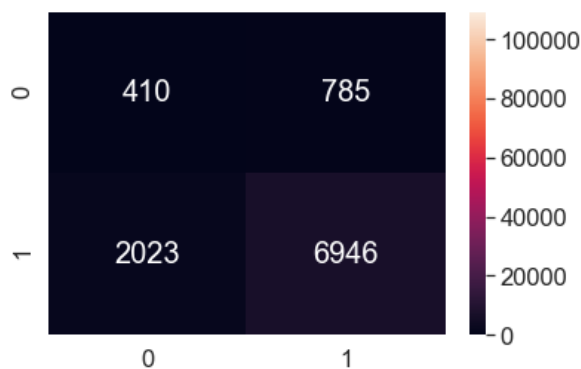
```
In [75]: 1 #Calling the function for test metrics
        2 test_metrics(y_test_w2v, y_pred)
```

Accuracy on test data: 72.37

Precision on test data: 89.85

Recall on test data: 77.44

F1_score on test data: 83.19



Random Search Cross Validation:


```
In [76]: 1 #Calling the function for random search cross validation
        2 best_estimator = random_search(tfidf_w2v_tr, y_tr_w2v)

Fitting 5 folds for each of 10 candidates, totalling 50 fits

[Parallel(n_jobs=3)]: Done 26 tasks      | elapsed: 15.6s
[Parallel(n_jobs=3)]: Done 50 out of 50 | elapsed: 29.6s finished

optimal hyperparameter: {'penalty': 'l1', 'C': 495}
Best accuracy: 90.0840048002743

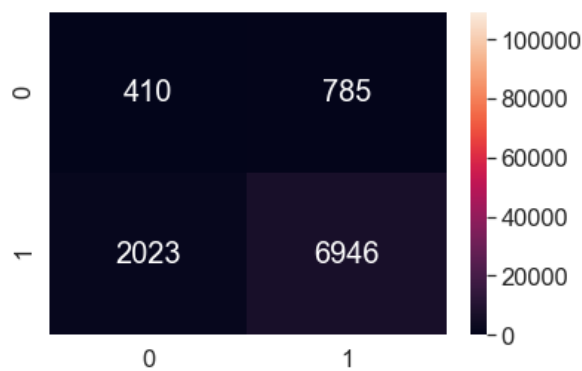
In [77]: 1 #Result showing the best classifier consisting of parameters
        2 best_estimator

Out[77]: LogisticRegression(C=495, class_weight=None, dual=False, fit_intercept=True,
                             intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
                             penalty='l1', random_state=None, solver='liblinear', tol=0.0001,
                             verbose=0, warm_start=False)

In [78]: 1 y_pred = best_estimator.predict(tfidf_w2v_test)

In [79]: 1 #Calling the function for test metrics
        2 test_metrics(y_test_w2v, y_pred)

Accuracy on test data: 72.37
Precision on test data: 89.85
Recall on test data: 77.44
F1_score on test data: 83.19
```



Summary:

Performance Table:

Featurization	sample size	CV	Accuracy	F1-score	C	Penalty
			Test accuracy	Test f1-score		
BOW	364k	Grid Search	91.85%	95.15%	10	L2
		Random Search	91.31%	94.81%	239	L2
TF-IDF	364k	Grid Search	93.20%	95.94%	1000	L2
		Random Search	93.20%	95.94%	250	L2
Avg-W2V	364k	Grid Search	88.61%	93.34%	100	L2
		Random Search	88.61%	93.34%	72	L2
TF-IDF W2V	50k	Grid Search	72.37%	83.19%	10	L1
		Random Search	72.37%	83.19%	495	L1

Observation:

- TF_IDF is working well for this dataset with logistic regression.
- Using with L1 regularization, as C decreases then number of non-zero rows decreases.
- There is an collinearity between the attributes.