

## TextPreprocessing+Word2Vec+ Machine Learning Classifier

### Import Libraries

In [27]:

```
### For computational and random seed purpose
import numpy as np
np.random.seed(42)
#to read csv file
import pandas as pd
#To split into train and cv data
from sklearn.model_selection import train_test_split
#To compute AUROC
from sklearn.metrics import auc, roc_auc_score
#for AUROC graph
import matplotlib.pyplot as plt
#for oversampling technique
from imblearn.over_sampling import SMOTE # (https://imbalanced-learn.org/stable/references/generated/imblearn.ove
r_sampling.SMOTE.html)
#Data is imbalanced, we need calibrated model
from sklearn.calibration import CalibratedClassifierCV
#for hyperparameter tuning and Cross-validation fold
from sklearn.model_selection import GridSearchCV, StratifiedKFold, RepeatedStratifiedKFold
#to ignore the error message
import warnings
warnings.filterwarnings("ignore")
#for heatmap and other plotting technique
import seaborn as sns
#to strandize the real value data
from sklearn.preprocessing import StandardScaler, MinMaxScaler
from sklearn.preprocessing import LabelEncoder
#To create Knn model on datasets
from sklearn.neighbors import KNeighborsClassifier
#for accuracy
from sklearn.metrics import accuracy_score, confusion_matrix, f1_score
# For applying model on datasets
from sklearn.naive_bayes import GaussianNB
from sklearn import svm
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.ensemble import StackingClassifier

#using simple network
from keras.models import Sequential
from keras import layers
from keras.wrappers.scikit_learn import KerasClassifier
import joblib
import sys
sys.modules['sklearn.externals.joblib'] = joblib
from mlxtend.feature_selection import SequentialFeatureSelector
from sklearn.linear_model import LogisticRegression
from sklearn.feature_selection import RFE

import warnings
warnings.filterwarnings('ignore')
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.preprocessing import RobustScaler
```

In [28]:

```
data = pd.read_csv("Modified_SQL_Dataset.csv")
```

In [29]:

```
data.drop_duplicates(inplace=True)
data.head()
```

Out[29]:

	Query	Label
0	" or pg_sleep ( __TIME__ ) --	1
1	create user name identified by pass123 tempora...	1
2	AND 1 = utl_inaddr.get_host_address ( ...	1
3	select * from users where id = '1' or @ @1 ...	1
4	select * from users where id = 1 or 1#" ( ...	1

In [30]:

```
data.shape
```

Out[30]:

(30907, 2)

In [31]:

```
import re
import nltk
nltk.download('punkt')
nltk.download('wordnet')
import string
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve,auc
from nltk.stem.porter import PorterStemmer
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer
from gensim.models import word2vec
from gensim.models import KeyedVectors
import pickle
```

```
from tqdm import tqdm
import os
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
```

In [32]:

```
#printing some random variable
sent_0 = data['Query'].values[0]
print(sent_0)
print("="*50)
```

```
sent_2000 = data['Query'].values[2000]
print(sent_2000)
print("="*50)
```

```
sent_15000 = data['Query'].values[15000]
print(sent_15000)
print("="*50)
```

```
sent_20000 = data['Query'].values[20000]
print(sent_20000)
print("="*50)
```

```
" or pg_sleep ( __TIME__ ) --
=====
-6073" ) ) ) or 9502 = 2012#
=====
wilson@autoconstruccion.fi
=====
SELECT SUBSTR ( "SQL Tutorial", -5, 5 ) AS ExtractString;
=====
```

In [33]:

```
#remove the special Character : https://stackoverflow.com/a/5843547/4084039
sent_15000 = re.sub('[^A-Za-z0-9]+',' ',sent_15000)
print(sent_15000)
```

wilson autoconstruccion fi

In [34]:

```
sent_20000 = re.sub('[^A-Za-z0-9]+',' ',sent_20000)
print(sent_20000)
```

SELECT SUBSTR SQL Tutorial 5 5 AS ExtractString

In [35]:

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

In [36]:

```
# Combining all the above students
from tqdm import tqdm
preprocessed_query = []
lemmatizer = WordNetLemmatizer()
# tqdm is for printing the status bar
for sentence in tqdm(data['Query'].values):
    sentence = re.sub('[^A-Za-z0-9]+',' ', sentence)
    sentence = re.sub(r',', ' ', sentence)
    #https://www.machinelearningplus.com/nlp/lemmatization-examples-python/
    tokenization = nltk.word_tokenize(sentence)
    sentence = ' '.join([lemmatizer.lemmatize(w) for w in tokenization])
    # https://gist.github.com/sebleier/554280
    sentence = ' '.join(e.lower() for e in sentence.split() if e.lower() not in stopwords)
    preprocessed_query.append(sentence.strip())
```

100%|██████████| 30907/30907 [00:05<00:00, 6016.61it/s]

In [37]:

```
preprocessed_query
```

Out[37]:

```
['pg sleep time',
'create user name identified pass123 temporary tablespace temp default tablespace user',
'1 utl inaddr get host address select distinct table name select distinct table name rownum limit s
ys table limit 5',
'select user id 1 1 1 union select 1 version 1',
'select user id 1 1 union select 1 version 1',
'select name syscolumns id select id sysobjects name tablename',
'select user id 1 1 1 1',
'1 load file char 47 101 116 99 47 112 97 115 115 119 100 1 1 1',
'select user id 1 1 1 union select 1 version 1',
'select user id 1 union select 1 version 1',
'1 1',
'',
'admin 1 1',
'select user id 1 1 1 1',
'1 1',
'1 utl inaddr get host address select distinct column name select distinct column name rownum limit
```

```

'1 utl_inaddr get host address select distinct column name select distinct column name rownum limit
tab column limit 5',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'distinct',
'select user id 1 1 1 1',
'1 ascii lower substring select top 1 name sysobjects xtype u 1 1 116',
'select user id 1 1 1 1',
'admin 1 1',
'select user id 1 1 1 1',
'insert',
'select user id 1 1 1 1 1 1',
'select user id 1 1 1 1',
'1 utl_inaddr get host address select distinct password select distinct password rownum limit sys u
ser limit 8',
'select user id 1 1 1 1',
'1 1',
'1 1',
'select user id 1 1 1 1',
'select user id 1 union select 1 version 1',
'select user id 1 union select version 1',
'1 1',
'27 1 1',
'select user id 1 union select 1 version 1',
'select user id 1 union select version 1',
'1 load file char 47 101 116 99 47 112 97 115',
'select user id 1 1 1 union select 1 version 1',
'select user id 1 1 1 1 1',
'select user id 1 union select 1 version 1',
'select user id 1 union select 1 1 version 1',
'select user id 1 1 union select null banner v version rownum 1 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'1',
'select user id 1 1 1 1',
'admin 1 1',
'select user id 1 1 union select null banner v version rownum 1 1',
'x27union select',
'print',
'select user id 1 union select version 1',
'1 1',
'1 pg sleep time',
'select user id 1 1 1 1',
'admin 1 1',
'2 1 3',
'select user id 1 1 1 1 1',
'select user id 1 1 union select null version 1',
'admin 1 1',
'not select system user sa waitfor delay 0 0 2',
'select user id 1 union select 1 banner v version rownum 1 1',
'0 0',
'select user id 1 union select version 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'select user id 1 union select null version 1',
'select user id 1 1 1 1',
'char 4039 41 2b 40select',
'select user id 1 1 1 1 1',
'admin 1 1',
'declare varchar 200 select 0x73656c6',
'0 0',
'select user id 1 1 1 1',
'x x',
'select user id 1 union select 1 version 1',
'sleep 50',
'select user id 1 1 union select null banner v version rownum 1 1',
'1 union select 1 2 3 4 5 6 name sysobjects xtype u',
'select user id 1 1 union select version 1',
'select user id 1 1 union select null version 1',
'select user id 1 1 1 1 1',
'1 1',
'select user id 1 1 union select 1 banner v version rownum 1 1',
'1 utl_inaddr get host address select distinct column name select distinct column name rownum limit
tab column limit 4',
'select user id 1 union select version 1',
'1 user name dbo',
'not substring select version 24 1 1 waitfor delay 0 0 2',
'select user id 1 1 union select 1 version 1',
'select user id 1 1 union select null version 1',
'select user id 1 1 1 1',

```

```

'select user id 1 1 1 1',

'select user id 1 1 1 1 1',
'1 utl inaddr get host address select distinct username select distinct username rownum limit sys u
sers limit 6',
'1 1',
'select user id 1 1 1 1',
'text',
'1234 1 0 union select admin 81dc9bdb52d04dc20036dbd8313ed055',
'select user id 1 1 1 1',
'select user id 1 1 1 union select 1 version 1',
'admin 1 1',
'select user id 1 1 1 1 1',
'',
'1 utl inaddr get host address select distinct username select distinct username rownum limit sys u
sers limit 4',
'select information schema table',
'select user id 1 1 1 1',
'1 pg sleep time',
'select user id 1 union select version 1',
'declare q nvarchar 200 select q 0x770061',
'1 1',
'select user id 1 1 1 1 1',
'declare q nvarchar 200 select q 0x770061006900740066006f0072002000640065006c0061007900200027003000
3a0030003a0031003000270000 exec q',
'select user id 1 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'unusual unusual',
'1 select version',
'truncate',
'',
'3 3',
'select user id 1 1 union select 1 banner v version rownum 1 1',
'select user id 1 union select 1 version 1',
'admin 1 1',
'select user id 1 1 union select version 1',
'1 benchmark 10000000 md5 1',
'select user id 1 union select 1 version 1',
'1 utl inaddr get host address select sys database name dual',
'0 0',
'select user id 1 1 1 1 1',
'select user id 1 1 1 1',
'not substring select version 25 1 0 waitfor delay 0 0 2',
'1 sleep time',
'select user id 1 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'select user id 1 union select null version 1',
'select user id 1 1 1 1',
'true',
'execute immediate sel ect u er',
'hi x x',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'admin',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'union select 1 load file etc passwd 1 1 1',
'waitfor delay 0 0 time',
'unusual unusual',
'insert mysql user user host password value name localhost password pass123',
'admin 1 1',
'select user id 1 1 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'like',
'username like char 37',
'pg sleep time',
'select user id 1 1 1 1',
'3 3',
'select user id 1 1 1 1',
'select user id 1 1 1 1 1',
'exec sp',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'admin 1 1',
'waitfor delay 0 0 10',
'select user id 1 union select null version 1',
'select user id 1 union select 1 version 1',
'1 select var temp',
'hi',
'1 utl inaddr get host address select distinct column name select distinct column name rownum limit
tab column limit 1',
'admin 1 1',
..

```

```

'true',

'select user id 1 1 union select 1 version 1',
'select user id 1 1 1 1',
'1 pg sleep time',
'select user id 1 1 1 1',
'sleep time',
'select user id 1 1 1 1',
'admin 1 1',
'select user id 1 union select 1 version 1',
'exec xp regread',
'select user id 1 1 union select 1 version 1',
'select user id 1 1 1 1 1',
'select user id 1 1 1 1 1',
'select user id 1 1 1 1',
'1 1 1',
'0 0',
'union select user login char',
'select user id 1 union select 1 version 1',
'admin 1 1',
'1 select',
'select user id 1 1 1 1',
'x 1 select count tablename',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 banner v version rownum 1 1',
'1 benchmark 10000000 md5 1',
'procedure',
'',
'1 utl inaddr get host address select sys login user dual',
'1 utl inaddr get host address select distinct username select distinct username rownum limit sys u
sers limit 2',
'desc user',
'select user id 1 1 union select 1 version 1',
'select user id 1 1 1 1 1 1',
'',
'',
'select user id 1 1 union select 1 version 1',
'select user id 1 union select 1 version 1',
'select user id 1 union select null version 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'admin',
'select user id 1 1 1 1',
'utl http request',
'1 utl inaddr get host address select distinct table name select distinct table name rownum limit s
ys table limit 1',
'select user id 1 union select 1 version 1',
'benchmark 10000000 md5 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 union select 1 version 1',
'select user id 1 1 1 1',
'1 uni select',
'select user id 1 union select null version 1',
'select user id 1 1 union select null version 1',
'admin 1 1',
'utl http request http 192 168 1 1',
'7659 7659',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'1 utl inaddr get host address select distinct granted role select distinct granted role rownum lim
it dba role privs grantee sys loginuser limit 6',
'select user id 1 1 1 1',
'true',
'select user id 1 union select 1 version 1',
'select user id 1 union select version 1',
'select user id 1 1 union select 1 version 1',
'true',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'select user id 1 1 1 1',
'3 3',
'1 utl inaddr get host address select global name global name',
'select user id 1 1 union select null banner v version rownum 1 1',
'select user id 1 1 1 1 1',
'select user id 1 1 union select null version 1',
'select user id 1 1 1 union select version version 1',
'select user id 1 1 1 1 1',
'uid like',
'select user id 1 1 1 1 1',
'select user id 1 union select 1 version 1',
'23 1 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1 1',

```

```

'select user id 1 1 1 1',

'select user id 1 1 1 1',
'select user id 1 union select version 1',
'0 1 1',
'select user id 1 1 1 1',
'admin 1 1',
'select user id 1 1 1 1 1',
'select user id 1 1 1 1',
'1 utl inaddr get host address select distinct password select distinct password rownum limit sys u
ser limit 5',
'1 utl inaddr get host address select distinct column name select distinct column name rownum limit
tab column limit 7',
'apos',
'sqlvuln',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'',
'hi 1 1',
'select user id 1 union select null banner v version rownum 1 1',
'exec master xp cmdshell ping 172 10 1 255',
'select user id 1 1 1 union select 1 banner v version rownum 1 1',
'select user id 1 union select version 1',
'select information schema table',
'select user id 1 1 1 1',
'select top 1',
'1 1',
'',
'1 1',
'select user id 1 1 union select null version 1',
'select user id 1 union select version 1',
'',
'exec sp addlogin name password',
'select user id 1 union select null banner v version rownum 1 1',
'0x77616974666f722064656c61792027303a303a31302700 exec',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'group userid 1 1',
'select user id 1 1 union select null version 1',
'select user id 1 1 1 1 1',
'select user id 1 1 1 1 1',
'1 1',
'select user id 1 union select null version 1',
'select user id 1 1 1 1',
'',
'declare varchar 22 select',
'select user id 1 union select 1 version 1',
'select user id 1 union select null version 1',
'select user id 1 1 union select version version 1',
'variable',
'exec master xp cmdshell',
'select user id 1 1 1 1 1',
'select user id 1 union select 1 version 1',
'select user id 1 union select la version 1',
'select user id 1 1 union select 1 version 1',
'select user id 1 1 union select 1 version 1',
'1 select var temp',
'select user id 1 union select 1 banner v version rownum 1 1',
'select user id 1 1 1 1 1',
'select user id 1 1 1 1',
'x x',
'',
'1 1',
'select user id 11 1 union select 1 version 1',
'x x',
'select user id 1 1 union select 1 version 1',
'select user id 1 1 union select 1 version 1',
'1 utl inaddr get host address select distinct table name select distinct table name rownum limit s
ys table limit 4',
'srvrolemember sysadmin 0 waitfor delay 0 0 2',
'select user id 1 1 1 1',
'select user id 1 union select version 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 1 union select version version 1',
'1 1',
'',
'x x',
'select user id 1 1 1 1 1',
'select user id 1 union select 1 version 1',
'password',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'select user id 1 1 union select null version 1',

```

```

'1 utl_inaddr get host address select distinct granted role select distinct granted role rownum limit
it dba role privs grantee sys loginuser limit 3',
'1 1',
'admin 1 1',
'select user id 1 1 1 1 1',
'select user id 1 union select 1 version 1',
'select user id 1 union select null version 1',
'select user id 1 1 1 1',
'union select user login char 114 111 111 116',
'select user id 1 1 union select 1 version 1',
'1 1',
'1 utl_inaddr get host address select distinct granted role select distinct granted role rownum limit
it dba role privs grantee sys loginuser limit 5',
'pg sleep time',
'3 3',
'select user id 1 1 union select 1 version 1',
'select user id 1 1 1 1',
'union select information schema table',
'select user id 1 1 1 1',
'benchmark 10000000 md5 1',
'1 utl_inaddr get host address select count distinct column name sys tab column',
'1 utl_inaddr get host address select distinct table name select distinct table name rownum limit s
ys table limit 7',
'',
'delete',
'waitfor delay 0 0 time',
'1 1',
'1 non existant table 1',
'select user id 1 1 1 1',
'select user id 1 union select 1 version 1',
'declare varchar 200 select 0x77616974',
'x userid null',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'select user id 1 union select version 1',
'select user id 1 union select 1 version 1',
'1 utl_inaddr get host address select host name v instance',
'1 utl_inaddr get host address select distinct column name select distinct column name rownum limit
tab column limit 2',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'union select',
'select user id 1 1 1 1 1',
'select user id 1 1 1 1 1',
'select user id 1 1 union select 1 banner v version rownum 1 1',
'x full name like bob',
'x member email null',
'2 1',
'select user id 1 1 1 1 1',
'pg sleep time',
'',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'sleep time',
'sleep 50',
'select user id 1 1 1 union select 1 banner v version rownum 1 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'sleep time',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'select user id 1 union select version 1',
'select user id 1 union select version 1',
'select user id 1 union select null version 1',
'2 1',
'exec select user',
'anything x x',
'select user id 1 union select 1 version 1',
'select user id 1 union select 1 banner v version rownum 1 1',
'select user id 1 union select null version 1',
'',
'3 3',
'select user id 1 union select 1 version 1',
'select user id 1 1 union select 1 banner v version rownum 1 1',
'admin 1 1',
'union select',
'',
'waitfor delay 0 0 time',
'select user id 1 1 1 1',
'1 waitfor delay 0 0 10',
'bfilename',
'. . . . .

```



```

'admin 1 1',

'uef',
'password',
'1 sleep time',
'1 utl inaddr get host address select distinct table name select distinct table name rownum limit s
ys table limit 6',
'select user id 1 union select version 1',
'select user id 1 1 union select 1 version 1',
'1 1 select count tablenames',
'sleep time',
'1 utl inaddr get host address select distinct granted role select distinct granted role rownum lim
it dba role privs grantee sys loginuser limit 8',
'select user id 1 1 1 1 1 1',
'6',
'1',
'1 load file char 110 46 101 120 116 char 39 39 1 0',
'select user id 1 1 1 1',
'select user id 1 1 1 1 1',
'select user id 1 1 1 union select 1 version 1',
'select user id 1 1 1 1',
'1 1',
'1 utl inaddr get host address select distinct password select distinct password rownum limit sys u
ser limit 7',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 banner v version rownum 1 1',
'whatever whatever',
'declare varchar 200 select 0x73656c6563742040407665727369666e exec',
'select user id 1 1 1 1',
'',
'1 utl inaddr get host address select count distinct table name sys table',
'select user id 1 union select 1 version 1',
'admin 1 1',
'select user id 1 1 union select version 1',
'select user id 1 1 1 1',
'',
'select user id 1 1 1 1',
'1 1',
'drop table temp',
'1 utl inaddr get host address select distinct password select distinct password rownum limit sys u
ser limit 2',
'benchmark 10000000 md5 1',
'x x',
'select user id 1 1 1 1',
'select user id 1 union select 1 version 1',
'something like',
'not select serverproperty isintegratedsecurityonly 0 waitfor delay 0 0 2',
'1 utl inaddr get host address select distinct password select distinct password rownum limit sys u
ser limit 4',
'select user id 1 1 1 1',
'exec master xp cmdshell ping 10 10 1 2',
'admin 1 1',
'not select serverproperty isintegratedsecurityonly 1 waitfor delay 0 0 2',
'1 utl inaddr get host address select distinct table name select distinct table name rownum limit s
ys table limit 2',
'select user id 1 union select version version 1',
'select user id 1 1 union select null version 1',
'waitfor delay 0 0 time',
'x userid null',
'select user id 1 1 1 1 1',
'select user id 1 union select null version 1',
'select user id 1 1 1 1',
'1 1',
'select user id 1 union select 1a banner v version rownum 1 1',
'admin 1 1',
'select user id 1 1 union select 1 version 1',
'1 load file char 110 46 101 120 11',
'exec',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'1 1',
'select user id 1 1 1 1 1',
'select user id 1 union select version version 1',
'x email null',
'select user id 1 1 union select 1 version 1',
'',
'select user id 1 union select null banner v version rownum 1 1',
'select user id 11 1 union select 1 banner v version rownum 1 1',
'select user id 1 1 union select 1 banner v version rownum 1 1',
'update',
'1',
'select user id 1 1 union select 1 version 1',
'pg sleep time',
'6',

```

```

'select user id 1 union select 1 version 1',

'exists',
'',
'text n text',
'select user id 1 union select null version 1',
'select user id 1 union select version 1',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'0 0',
'select user id 1 1 1 1 1 1',
'select user id 1 1 1 1 1',
'select user id 1 1 union select version version 1',
'select user id 1 1 union select version 1',
'password 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 union select 1 version 1',
'union select null select version',
'select user id 1 1 1 1',
'select user id 1 union select 1 version 1',
'select user id 1 union select null version 1',
'select user id 1 union select 1 version 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 union select 1 version 1',
'select user id 1 union select null version 1',
'select user id 1 union select null version 1',
'select user id 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'1 utl inaddr get host address select banner v version rownum 1',
'select user id 1 1 1 1',
'1 utl inaddr get host address select distinct column name select distinct column name rownum limit
tab column limit 3',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'union select',
'asc',
'myappadmin adduser admin newpass',
'exec sp addsrvrolemember name sysadmin',
'select user id 1 1 union select 1 version 1',
'true',
'select user id 1 1 union select 1 banner v version rownum 1 1',
'1 utl inaddr get host address select distinct password select distinct password rownum limit sys u
ser limit 3',
'userid like',
'union select null null select version',
'',
'select user id 1 1 1 1',
'admin',
'select user id 1 union select 1 version 1',
'waitfor delay 0 0 time',
'select user id 1 union select version 1',
'select user id 1 1 1 1',
'admin 1 1',
'union select version',
'',
'sleep time',
'union select null null null null null select version',
'select user id 1 union select 1 version 1',
'select user id 1 union select version 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'select user id 1 1 union select null version 1',
'select user id 1 1 1 1',
'hi',
'select name syscolumns id sele',
'select user id 1 1 1 1 1',
'not substring select version 25 1 5 waitfor delay 0 0 2',
'select user id 1 union select 1 version 1',
'2 1 3',
'union select version',
'1 utl inaddr get host address select distinct column name select distinct column name rownum limit
tab column limit 8',
'exec master xp cmdshell nslookup www google com',
'isnull 1 0',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'declare q nvarchar 200 0x730065006c00650063',
'',
'union select version',
'select user id 1 1 1 union select 1 version 1',

```

```

'admin 1 1',

'exists',
'select user id 1 union select 1 version 1',
'declare varchar 8000 select 0x73656c',
'select user id 1 union select version 1',
'select user id 1 union select 1 banner v version rownum 1 1',
'select user id 1 1 1 1',
'select user id 1 union select 1 1a version 1',
'begin declare var varchar 8000 set var select var var login password user login',
'union select null null null null select version',
'hi 1 1',
'union select',
'select user id 1 1 1 1',
'select user id 1 1 union select version 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1 1',
'desc',
'select user id 1 1 1 1 1',
'anything x x',
'select user id 1 1 1 union select 1 version 1',
'select user id 1 1 1 1',
'0x770061006900740066006f0072002000640065006c00',
'select user id 1 union select version version 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'isnull 1 0',
'insert user login password level value char 0x70 char 0x65 char 0x74 char 0x65 char 0x72 char 0x70
char 0x65 char 0x74 char 0x65 char 0x72 char 0x64',
'select user id 1 1 1 1 1',
'select user id 1 1 1 1',
'admin 1 1',
'admin 1 1',
'user like',
'admin 1 1',
'sleep time',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'',
'select user id 1 union select version 1',
'admin 1 1',
'1 utl inaddr get host address select distinct table name select distinct table name rownum limit s
ys table limit 3',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'x 1 1 x',
'1 utl inaddr get host address select distinct table name select distinct table name rownum limit s
ys table limit 8',
'select user id 1 union select 1 version 1',
'x x',
'',
'select user id 1 1 union select version version 1',
'select user id 1 union select 1 version 1',
'select user id 1 union select null version 1',
'select user id 1 union select version 1',
'select user id 1 1 union select null version 1',
'select user id 1 union select null version 1',
'select user id 1 1 union select null version 1',
'sleep time',
'1 1',
'',
'select user id 1 1 1 1 1 1',
'group userid 1 1',
'select user id 1 1 1 union select 1 version 1',
'select user id 1 1 union select 1 version 1',
'0 0',
'0 0',
'pg sleep time',
'1 1',
'select user id 1 1 1 union select version 1',
'select user id 1 1 1 union select 1 version 1',
'select user id 1 1 union select 1 version 1',
'1 utl inaddr get host address select count distinct username sys users',
'sqlvuln',
'select user id 1 1 1 1',
'select user id 1 1 union select version 1',
'0x730065006c006500630074002000400076006500',
'select user id 1 union select 1 1 version 1',
'select user id 1 union select version version 1',
'declare varchar 8000 select 0x73656c6563742040407665727369666e',
'select user id 1 1 1 1',
'',

```

```

'select user id 1 1 1 1',

'1 utl inaddr get host address select distinct column name select distinct column name rownum limit
tab column limit 6',
'select user id 1 union select 1 version 1',
'username like char 37',
'select user id 1 1 union select version version 1',
'something thing',
'name',
'1 select version',
'declare q nvarchar 4000 select q',
'select user id 1 union select 1 version 1',
'1 utl inaddr get host address select distinct granted role select distinct granted role rownum lim
it dba role privs grantee sys loginuser limit 2',
'select user id 1 1 1 1',
'select user id 1 union select null version 1',
'select user id 1 1 1 1 1',
'1 utl inaddr get host address select distinct password select distinct password rownum limit sys u
ser limit 1',
'select user id 1 1 union select version version 1',
'select user id 1 union select 1a banner v version rownum 1 1',
'select user id 1 1 union select 1 version 1',
'admin 1 1',
'',
'select user id 1 1 1 1',
'1 1',
'declare q nvarchar 200 0x730065006c00650063007400200040004000760065007200730069006f006e00 exec q',
'admin',
'select user id 1 1 1 1 1',
'select user id 1 1 union select 1 version 1',
'select user id 1 1 1 1 1',
'begin declare var varchar 8000 set var',
'select user id 1 1 1 1',
'pg sleep time',
'print variable',
'select user id 1 1 1 1 1',
'',
'x x',
'select user id 1 1 1 1 1',
'1 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'',
'2a 28 7c 28mail 3d 2a 29 29',
'select user id 1 1 1 1',
'1 1',
'x 1 select count tablename',
'1 utl inaddr get host address select distinct username select distinct username rownum limit sys u
sers limit 3',
'select user id 1 union select 1 version 1',
'admin 1 1',
'',
'select user id 1 1 1 union select null version 1',
'select user id 1 union select 1 banner v version rownum 1 1',
'x x',
'union select',
'23 1 1',
'select user id 1 1 1 1 1',
'select user id 1 1 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'username like',
'elt 3 5 bin 15 ord 10 hex char 45',
'uname like',
'declare varchar 200 select 0x776169746666722064656c61792027303a303a31302700 exec',
'1',
'select user id 1 union select null version 1',
'1 utl inaddr get host address select distinct username select distinct username rownum limit sys u
sers limit 1',
'select user id 1 union select version 1',
'select user id 1 1 1 1',
'1 utl inaddr get host address select distinct granted role select distinct granted role rownum lim
it dba role privs grantee sys loginuser limit 1',
'select user id 1 1 1 1',
'admin 1 1',
'1 sleep time',
'select user id 1 1 1 1',
'sqlattempt2',
'select user id 1 union select null version 1',
'x email null',
'select user id 1 1 union select 1 banner v version rownum 1 1',
'1 utl inaddr get host address select distinct username select distinct username rownum limit sys u
sers limit 8',

```

```

'select user id 1 1 union select version version 1',

'1234 1 0 union select admin 81dc9bdb52d04dc20036dbd8313ed055',
'select user id 1 union select 1 version 1',
'var select var var temp end',
'select user id 1 union select null version 1',
'',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'select user id 1 union select version 1',
'',
'waitfor delay 0 0 time',
'',
'select user id 1 1 1 1',
'select user id 1 1 1 1',
'1 utl inaddr get host address select distinct granted role select distinct granted role rownum limit 7',
'select user id 1 union select 1 version 1',
'union select',
'select user id 1 1 1 1',
'select user id 1 1 1 1 1',
'select user id 1 1 1 union select 1 version 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 1 1 1 1',
'select user id 1 union select 1 version 1',
'exec master xp cmdshell nslookup www googl',
'benchmark 10000000 md5 1',
'union select',
'select user id 1 1 1 union select 1 version 1',
'select user id 1 1 1 1',
'',
'x member email null',
'select',
'1 1',
'',
'select user id 1 union select 1 version 1',
'select user id 1 1 union select 1 version 1',
'admin 1 1',
'select user id 1 union select 1 version 1',
'select user id 1 union select 1 banner v version rownum 1 1',
'1 utl inaddr get host address select distinct username select distinct username rownum limit sys users limit 5',
'1 utl inaddr get host address select count distinct granted role dba role privs grantee sys login user',
'select user id 1 1 union select null version 1',
'1 utl inaddr get host address select distinct password select distinct password rownum limit sys user limit 6',
'select user id 1 1 1 1',
'union select null null null select version',
'x x',
'select user id 1 union select version 1',
'1 1',
'select user id 1 1 1 1',
'3 3',
'1 benchmark 10000000 md5 1',
'1 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 union select 1 version 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1 1',
'admin 1 1',
'text',
'',
'select user id 1 1 union select version 1',
'benchmark 10000000 md5 1',
'exec xp',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1 1',
'hi',
'',
'sqlattemp1',
'admin 1 1',
'select user id 1 1 1 1',
'1 1 1',
'select user id 1 1 union select 1 banner v version rownum 1 1',
'variable',
'2a 7c',
'select user id 1 union select 1 version 1',
'mail',
'not substring select version 25 1 8 waitfor delay 0 0 2',
'select user id 1 1 union select 1 version 1',
'select user id 1 1 1 union select version version 1',
'1 utl inaddr get host address select count distinct password sys user',

```

```

'select user id 1 1 union select 1 version 1',

'order',
'select user id 1 1 1 1',
'username not null username',
'admin 1 1',
'',
'select user id 1 1 union select 1 version 1',
'select user id 1 union select version version 1',
'select user id 1 1 1 1 1',
'select user id 1 union select 1 version 1',
'',
'1 0 union',
'',
'0x770061006900740066006f0072002000640065006c00610079002000270030003a0030003a',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 union select 1 banner v version rownum 1 1',
'sqlvuln',
'text n text',
'1 1',
'0x730065006c00650063007400200040004000760065007200730069006f006e00 exec q',
'select user id 1 1 union select 1 version 1',
'select user id 1 union select version 1',
'select user id 1 1 1 1 1 1',
'select user id 1 1 1 1',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'1 exec sp exec xp',
'select user id 1 1 1 union select 1 banner v version rownum 1 1',
'3 3',
'admin',
'select top 1',
'select user id 1 1 1 1 1',
'',
'benchmark 10000000 md5 1',
'select user id 1 1 1 1',
'7659 7659',
'select user id 1 1 union select version version 1',
'select user id 1 1 1 1 1',
'1 1 1',
'',
'replace',
'select user id 1 union select 1 version 1',
'select user id 1 1 1 1',
'select user id 1 union select null version 1',
'select user id 1 union select la banner v version rownum 1 1',
'',
'objectclass',
'1 utl inaddr get host address select distinct granted role select distinct granted role rownum lim
it dba role privs grantee sys loginuser limit 4',
'limit',
'1 1',
'select user id 1 1 1 1',
'',
'1 utl inaddr get host address select distinct username select distinct username rownum limit sys u
sers limit 7',
'create user name identified pass123',
'select user id 1 union select null version 1',
'select user id 1 1 union select version 1',
'',
'select user id 1 1 1 1 1',
'',
'select user id 1 1 1 1 1',
'x full name like bob',
'select user id 1 union select version 1',
'select user id 1 1 1 1',
'select user id 1 union select version 1',
'benchmark 10000000 md5 1',
'waitfor delay 0 0 time',
'sleep time',
'wapiti',
'sleep time',
'sleep time',
'pg sleep time',
'sleep time',
'sleep time',
'pg sleep time',
'pg sleep time',
'sleep time',
'waitfor delay 0 0 time',
'1 8156 select count generate series 1 5000000',
'1 clye 7842 7842 char 109 char 79 char 70 char 90 regexp substring repeat right char 5012 0 500000
0000 null',
'4860 azyx 6901 6901 union select 6901 6901 6901 6901 6901',

```

```

'1 union select null null null null null',

'select benchmark 5000000 md5 0x4c4d6142 6866 6866',
'1 boed 6787 6787',
'1 vdbf 7969 7969 extractvalue 1297 concat 0x5c 0x7171706a71 select elt 1297 1297 1 0x717a767a71',
'1 4281 4281',
'1 8024 3560',
'1 5452 6050 6050 ciyc like ciyc',
'1 8148 like abcdefg upper hex randomblob 500000000 2',
'1 select case 9443 9443 sleep 5 else 9443 select 9443 information schema character set end',
'1 select rawn dual 4988 4988 char 68 char 69 char 97 char 85 regexp substring repeat right char 53
89 0 5000000000 null',
'4925 union select 5686 5686 5686 5686 5686 5686 5686',
'1 dnhd 2657 2657 4240 select 4240 pg sleep 5',
'9534 3038 3038',
'1 6941 6941 6537 dbms pipe receive message chr 76 chr 116 chr 117 chr 65 5',
'select case 7978 6009 7978 else 1 select 0 end',
'1212 make set 7588 2306 2306',
'1 5466 5466 2388 benchmark 5000000 md5 0x6d457153',
'3794 union select 2485 2485 2485 2485 2485',
'1 5286 select count user t1 user t2 user t3 user t4 user t5 gnil gnil',
'1 8635 select count generate series 1 5000000',
'6400 union select 4650 4650 4650',
'1 select rttq dual 7368 7368 updatexml 1808 concat 0x2e 0x7171706a71 select elt 1808 1808 1 0x717a
767a71 8666',
'1 ztkr 1532 1532',
'1 9842 9842 union select null null null null null null null null',
'1 exp select select concat 0x7171706a71 select elt 6270 6270 1 0x717a767a71 0x78 x fbsi like fbsi'
,
'1 char 68 char 69 char 97 char 85 regexp substring repeat right char 5389 0 5000000000 null uwep u
wep',
'4291 5023 ctxsys drithsx sn 5023 chr 113 chr 113 chr 112 chr 106 chr 113 select case 5023 5023 1 e
lse 0 end dual chr 113 chr 122 chr 118 chr 122 chr 113 xyhi xyhi',
'1 boolean mode 3707 select count sysibm systables t1 sysibm systables t2 sysibm systables t3',
'1 1022 select count user t1 user t2 user t3 user t4 user t5',
'call regexp substring repeat left crypt key char 65 char 69 char 83 null 0 5000000000 null pawh paw
h',
'1589 1589 1',
'select count generate series 1 5000000 7240 7240',
'1 3715 char 113 char 113 char 112 char 106 char 113 select case 3715 3715 char 49 else char 48 end
char 113 char 122 char 118 char 122 char 113 9548 9548',
'1 6240 qqpjq select case 6240 6240 1 else 0 end rdb database qzvzq 6406 6406',
'5021 select yadq 4285 4285 order 1',
'1 4386 utl inaddr get host address chr 113 chr 113 chr 112 chr 106 chr 113 select case 4386 4386 1
else 0 end dual chr 113 chr 122 chr 118 chr 122 chr 113',
'1 6793 select 6793 pg sleep 5 baia baia',
'1 3824 benchmark 5000000 md5 0x76555642 jozh jozh',
'9508 union select 9950 9950 9950 9950 9950 9950 9950',
'1 8514 benchmark 5000000 md5 0x544d5a4c',
'1064 make set 6439 2937 2937 qojd qojd',
'1 rbpx 1264 1264',
'1 exp select select concat 0x7171706a71 select elt 6270 6270 1 0x717a767a71 0x78 x',
'8284 6171 select 8284 else drop function mbih',
'1 sleep 5',
'7868 9323 9323',
'1 ekjw 5477 5477 union select null null null null null',
'1 2462 2462 2716 select count sysusers sys1 sysusers sys2 sysusers sys3 sysusers sys4 sysusers sys
5 sysusers sys6 sysusers sys7',
'select like abcdefg upper hex randomblob 500000000 2 mfib mfib',
'1 4240 select 4240 pg sleep 5',
'1 4595 4595',
'1 oknw 8777 8777',
'1 6501 6501',
'1 char 111 char 77 char 121 char 88 regexp substring repeat left crypt key char 65 char 69 char 83
null 0 5000000000 null 8929 8929',
'7999 8422 1336',
'5742 1314 1314 5903 qqpjq select case 5903 5903 1 else 0 end rdb database qzvzq',
'1 potk 5040 5040 elt 5873 5873 sleep 5',
'4605 union select 8542 8542 8542 8542 8542 8542 8542',
'select count sysibm systables t1 sysibm systables t2 sysibm systables t3 njnr njnr',
'select case 6558 4327 1 else null end',
'7535 2724 char 113 char 113 char 112 char 106 char 113 select case 2724 2724 char 49 else char 48
end char 113 char 122 char 118 char 122 char 113',
'end rqay like rqay',
'2312 union select 5282 5282 5282 5282 5282 5282 5282',
'1 8594 select 8594 pg sleep 5',
'select case 8993 8846 8993 else 8993 select 8993 mysql db end',
'select count sysibm systables t1 sysibm systables t2 sysibm systables t3',
'select pg sleep 5',
'call regexp substring repeat right char 3702 0 5000000000 null 4142 4142',
'8153 qh1b 4948 4948 union select 4948 4948 4948 4948 4948 4948',
'7319 4493 utl inaddr get host address chr 113 chr 113 chr 112 chr 106 chr 113 select case 4493 449
3 1 else 0 end dual chr 113 chr 122 chr 118 chr 122 chr 113',
'1 bedq 8781 8781 8466 benchmark 5000000 md5 0x694a4745',
'1 8888 8888'

```

```
'1 2006 2006',

'5139 union select 3373 3373 3373 3373 3373 3373 3373 3373',
'1 4411 select count sysusers sys1 sysusers sys2 sysusers sys3 sysusers sys4 sysusers sys5 sysusers
sys6 sysusers sys7 zcyc',
'1 rvch 1863 1863 8384 like abcdefg upper hex randblob 500000000 2',
'1 8514 select count domain domain t1 domain column t2 domain table t3 loao',
'9446 wmrq 3705 3705 union select 3705',
'iif 7889 5114 1 1 0',
'select dbms pipe receive message chr 66 chr 67 chr 79 chr 101 5 dual ztmd ztmd',
'1 union select null null null null null null null null',
'1 char 109 char 79 char 70 char 90 regexp substring repeat right char 5012 0 5000000000 null',
'1 elt 4249 4249 7259',
'1 4867 4867 rlike select case 7689 7689 1 else 0x28 end',
'8858 5680 select 8858 else drop function pyuo',
'1 procedure analyse extractvalue 9255 concat 0x5c benchmark 5000000 md5 0x52515a50 1 1748 1748',
'1 klie 2840 2840 8514 benchmark 5000000 md5 0x544d5a4c',
'3859 3440 cast chr 113 chr 113 chr 112 chr 106 chr 113 select case 3440 3440 1 else 0 end text chr
113 chr 122 chr 118 chr 122 chr 113 numeric 5846 5846',
'end qkkn like qkkn',
'1 pzoo 8036 8036 6793 select 6793 pg sleep 5',
'2120 8734 8844',
'1 lomw 9257 9257 union select null null null null null null',
'1 6784 6784 elt 3114 3114 sleep 5',
'1 select syrz 7699 7699 union select null null null',
'1 4595 4595',
'4984 union select 6980 6980 6980 6980 6980 6980 6980 6980',
'5224 1962 1962 union select 1962 1962 1962 1962 1962 1962 1962 1962',
'1 8514 select count domain domain t1 domain column t2 domain table t3',
'1 select gboi 4191 4191 8514 select count domain domain t1 domain column t2 domain table t3',
'1 union select null null null null null',
'1 2633 dbms pipe receive message chr 112 chr 65 chr 65 chr 103 5 xmnd xmnd',
'1 row 6237 7469 select count concat 0x7171706a71 select elt 6237 6237 1 0x717a767a71 floor rand 0
2 x select 5192 union select 3785 union select 3931 union select 7158 group x ejul ejul',
'3721 union select 9050 9050',
'1 paai 4089 4089 3707 select count sysibm systables t1 sysibm systables t2 sysibm systables t3',
'1 elt 3114 3114 sleep 5',
'1 order 1',
'9281 8363 8363 make set 8220 5127 5127',
'1 6671 6671 char 119 char 100 char 99 char 121 regexp substring repeat right char 1441 0 5000000000
0 null',
'end vwbx vwbx',
'1 boolean mode 8189 select count sysibm systables t1 sysibm systables t2 sysibm systables t3',
'1 rlike sleep 5 iwct iwct',
'9087 order 1',
'7562 8571 8571',
'3707 8571 8571',
'1 select twyt 3376 3376 7756 dbms utility sqlid sqlhash chr 113 chr 113 chr 112 chr 106 chr 113 se
lect case 7756 7756 1 else 0 end dual chr 113 chr 122 chr 118 chr 122 chr 113',
'call regexp substring repeat right char 3702 0 5000000000 null eevk eevk',
...]
```

In [38]:

```
data['Query'] = preprocessed_query
```

In [39]:

```
X = data.drop(['Label'],axis=1)
y = data['Label']
```

### Splitting data into Train and cross validation(or test) : Stratified Sampling

In [40]:

```
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.2,stratify=y) # We split the train data and test
data of both X(input features)
# and Y(class_lable)
```

In [41]:

```
print(X_train.shape)
print(y_train.shape)
print(X_test.shape)
print(y_test.shape)
```

```
(24725, 1)
(24725,)
(6182, 1)
(6182,)
```



## TF-IDF Weighted Word2Vec

In [42]:

```
import pickle
import numpy as np
with open('glove_vectors','rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())
```

In [43]:

```
train_preprocessed_query = X_train['Query']
```

In [44]:

```
tfidf_model = TfidfVectorizer()
tfidf_model.fit(train_preprocessed_query)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf_model.get_feature_names_out(), list(tfidf_model.idf_)))
tfidf_words = set(tfidf_model.get_feature_names_out())
```

In [45]:

```
# average Word2Vec
# compute average word2vec for each review.
train_preprocessed_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
row=0;
for sentence in tqdm(train_preprocessed_query):
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight =0; # num of words with a valid vector in the sentence/review
    for word in sentence.split():#for each word in review/sentance
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    train_preprocessed_w2v_vectors.append(vector)

print(len(train_preprocessed_w2v_vectors))
print(len(train_preprocessed_w2v_vectors[0]))
```

100%|██████████| 24725/24725 [00:01<00:00, 24018.73it/s]

24725  
300

In [46]:

```
print(len(train_preprocessed_w2v_vectors))
```

24725

In [47]:

```
test_preprocessed_query = X_test['Query']
```

In [48]:

```
# average Word2Vec
# compute average word2vec for each review.
test_preprocessed_w2v_vectors = []; # the avg-w2v for each sentence/review is stored in this list
for sentence in tqdm(test_preprocessed_query):
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in sentence.split(): # for each word in review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf value((sentence.count(word)/len(sentence.s
            tf_idf = dictionary[word]*(sentence.count(word)/len(sentence.split())) # getting the tfidf value for each w
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    test_preprocessed_w2v_vectors.append(vector)
print(len(test_preprocessed_w2v_vectors))
print(len(test_preprocessed_w2v_vectors[0]))
```

100%|██████████| 6182/6182 [00:00<00:00, 28710.87it/s]

6182  
300

In [49]:

```
print(np.shape(train_preprocessed_w2v_vectors))
```

(24725, 300)

In [50]:

```
print(np.shape(y_train))
```

(24725,)

In [51]:

```
from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()
X_train = scaler.fit_transform(train_preprocessed_w2v_vectors)
X_test = scaler.transform(test_preprocessed_w2v_vectors)
```

In [52]:

```
X_train
```

Out[52]:

```
array([[0.59104619, 0.68811971, 0.327026 , ..., 0.3408724 , 0.53712702,
        0.35250338],
       [0.44126443, 0.61634371, 0.353036 , ..., 0.41080645, 0.57888029,
        0.61498472],
       [0.34851322, 0.49209983, 0.49063358, ..., 0.43538778, 0.5648758 ,
        0.56375076],
       ...,
       [0.51908424, 0.45200808, 0.34155414, ..., 0.79024537, 0.46528552,
        0.37306871],
       [0.54192352, 0.39263181, 0.42231304, ..., 0.57039383, 0.6177662 ,
        0.40744118],
       [0.44481232, 0.52646882, 0.44104558, ..., 0.42116899, 0.51243162,
        0.44570459]])
```

In [ ]:

```
#kNN (See Docs: https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html)
from sklearn.model_selection import GridSearchCV, RepeatedStratifiedKFold
from sklearn.neighbors import KNeighborsClassifier

parameters = {'n_neighbors':[1,3,5,7,10,11,15,10,50,75,80,85,90,95,100]}
rkf = StratifiedKFold(n_splits=11,random_state=42,shuffle=True)
clf = KNeighborsClassifier()

grid = GridSearchCV(estimator = clf, param_grid = parameters , scoring = 'roc_auc', verbose = 1, n_jobs = -1,cv=rkf)
grid.fit(X_train,y_train)

print("Best Score:" + str(grid.best_score_))
print("Best Parameters: " + str(grid.best_params_))
```

Fitting 11 folds for each of 15 candidates, totalling 165 fits  
Best Score:0.9712641731496885  
Best Parameters: {'n\_neighbors': 7}

In [ ]:

```
from sklearn.neighbors import KNeighborsClassifier

clf = KNeighborsClassifier(n_neighbors = 7)
clf.fit(X_train,y_train)
y_pred = clf.predict(X_train)
knn_train_auc = roc_auc_score(y_train,y_pred)
print(knn_train_auc)
```

0.9106126322468551

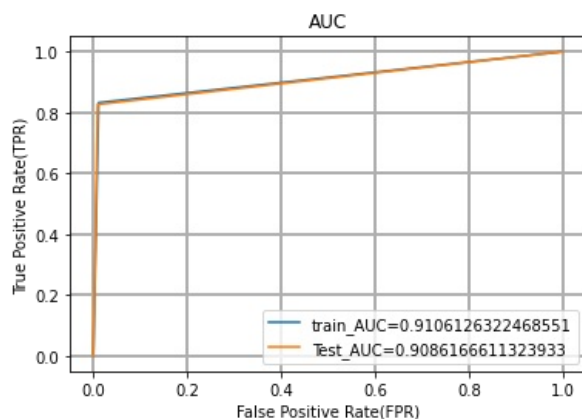
In [ ]:

```
y_test_pred = clf.predict(X_test)
knn_test_auc = roc_auc_score(y_test,y_test_pred)
print(knn_test_auc)
```

0.9086166611323933

In [ ]:

```
train_fpr,train_tpr,tr_thresholds = roc_curve(y_train,y_pred)
test_fpr,test_tpr,te_thresholds = roc_curve(y_test,y_test_pred)
plt.plot(train_fpr,train_tpr,label = "train_AUC="+str(auc(train_fpr,train_tpr)))
plt.plot(test_fpr,test_tpr,label = "Test_AUC="+str(auc(test_fpr,test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("AUC")
plt.grid(linestyle='-', linewidth=2)
```



In [ ]:

```
print("f1-score:",f1_score(y_test,y_test_pred))
```

f1-score: 0.8972096351061292

## Logistic Regression

In [ ]:

```
#ref https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html
parameters = {'penalty':['l1', 'l2'], 'C':[0.001,0.01,0.1,1,10], 'solver':['liblinear','saga','sag']}
rkf = StratifiedKFold(n_splits=11,random_state=42,shuffle=True)
clf = LogisticRegression(random_state=42)

grid = GridSearchCV(estimator = clf, param_grid = parameters , scoring = 'roc_auc', verbose = 1, n_jobs = -1,cv=rkf)
grid.fit(X_train,y_train)

print("Best Score:" + str(grid.best_score_))
print("Best Parameters: " + str(grid.best_params_))
```

Fitting 11 folds for each of 30 candidates, totalling 330 fits  
Best Score:0.9707054340363939  
Best Parameters: {'C': 10, 'penalty': 'l2', 'solver': 'liblinear'}

In [ ]:

```
clf = LogisticRegression(C=10,penalty='l2',solver='liblinear')
clf.fit(X_train,y_train)

y_pred = clf.predict(X_train)
logistic_train_auc = roc_auc_score(y_train,y_pred)
print(logistic_train_auc)
```

0.9160849283940161

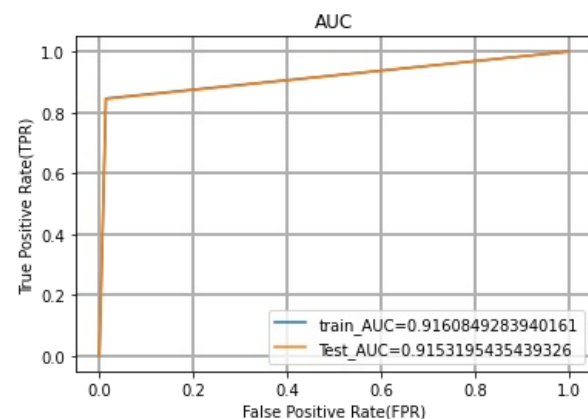
In [ ]:

```
y_test_pred = clf.predict(X_test)
logistic_test_auc = roc_auc_score(y_test,y_test_pred)
print(logistic_test_auc)
```

0.9153195435439326

In [ ]:

```
train_fpr,train_tpr,tr_thresholds = roc_curve(y_train,y_pred)
test_fpr,test_tpr,te_thresholds = roc_curve(y_test,y_test_pred)
plt.plot(train_fpr,train_tpr,label = "train_AUC="+str(auc(train_fpr,train_tpr)))
plt.plot(test_fpr,test_tpr,label = "Test_AUC="+str(auc(test_fpr,test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("AUC")
plt.grid(linestyle='-', linewidth=2)
```



In [ ]:

```
print("f1-score:",f1_score(y_test,y_test_pred))
```

f1-score: 0.9040451552210724

## Naive Bayes Classifier

In [ ]:

```
from sklearn.naive_bayes import MultinomialNB
```

In [ ]:

```
#ref https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html

parameters = {'alpha':[0.00001,0.0001,0.001,0.01,0.1,1,10,100,1000]}

rkf = StratifiedKFold(n_splits=11,random_state=42,shuffle=True)
clf = MultinomialNB()

grid = GridSearchCV(estimator = clf, param_grid = parameters, scoring = 'roc_auc', verbose = 1, n_jobs = -1,cv=rkf)
grid.fit(X_train,y_train)

print("Best Score:" + str(grid.best_score_))
print("Best Parameters: " + str(grid.best_params_))
```

Fitting 11 folds for each of 9 candidates, totalling 99 fits  
Best Score:0.8615649736325582  
Best Parameters: {'alpha': 1000}

In [ ]:

```
clf = MultinomialNB(alpha=1000)
clf.fit(X_train,y_train)

y_pred = clf.predict(X_train)
multinomial_train_auc = roc_auc_score(y_train,y_pred)
print(multinomial_train_auc)
```

0.6636754685878632

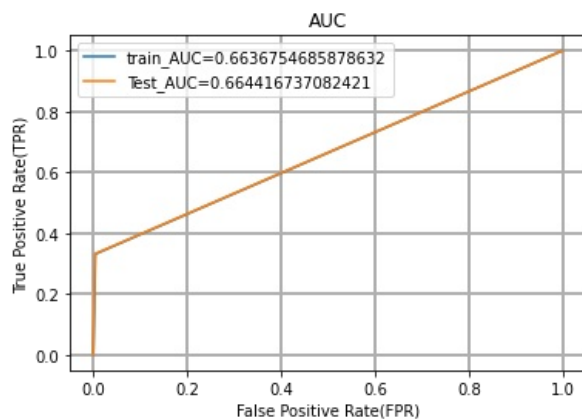
In [ ]:

```
y_test_pred = clf.predict(X_test)
multinomial_test_auc = roc_auc_score(y_test,y_test_pred)
print(multinomial_test_auc)
```

0.664416737082421

In [ ]:

```
train_fpr,train_tpr,tr_thresholds = roc_curve(y_train,y_pred)
test_fpr,test_tpr,te_thresholds = roc_curve(y_test,y_test_pred)
plt.plot(train_fpr,train_tpr,label = "train_AUC="+str(auc(train_fpr,train_tpr)))
plt.plot(test_fpr,test_tpr,label = "Test_AUC="+str(auc(test_fpr,test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("AUC")
plt.grid(linestyle='-', linewidth=2)
```



In [ ]:

```
print("f1-score:",f1_score(y_test,y_test_pred))
```

f1-score: 0.496551724137931

## Decision Tree Classifier

In [ ]:

```
#ref =https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html

params = {'max_depth':[2,3,5,7,9]}

#The instance of SVC
rkf = StratifiedKFold(n_splits=11,random_state=42,shuffle=True)
tree_model = DecisionTreeClassifier(random_state=42)
#Used GridsearchCV for Hyper-parameter
grid = GridSearchCV(estimator = tree_model, param_grid = params, scoring = 'roc_auc', verbose = 1, n_jobs = -1,cv
=rkf)
grid.fit(X_train,y_train)

print("Best Score:" + str(grid.best_score_))
print("Best Parameters: " + str(grid.best_params_))
```

Fitting 11 folds for each of 5 candidates, totalling 55 fits  
Best Score:0.9516458538537644  
Best Parameters: {'max\_depth': 9}

In [ ]:

```
tree_clf = DecisionTreeClassifier(max_depth=9)
tree_clf.fit(X_train,y_train)

y_pred = tree_clf.predict(X_train)
tree_train_auc = roc_auc_score(y_train,y_pred)
print(tree_train_auc)
```

0.9162501703755764

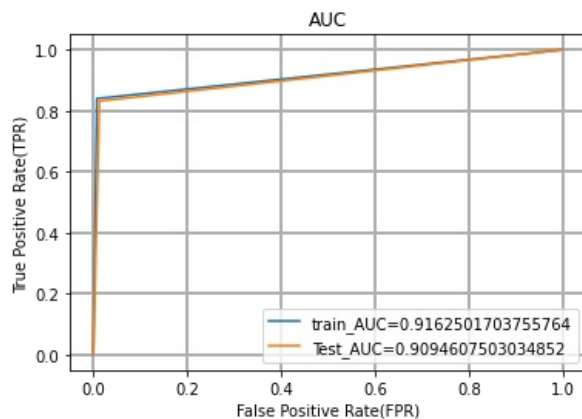
In [ ]:

```
y_test_pred = tree_clf.predict(X_test)
tree_test_auc = roc_auc_score(y_test,y_test_pred)
print(tree_test_auc)
```

0.9094607503034852

In [ ]:

```
train_fpr,train_tpr,tr_thresholds = roc_curve(y_train,y_pred)
test_fpr,test_tpr,te_thresholds = roc_curve(y_test,y_test_pred)
plt.plot(train_fpr,train_tpr,label = "train_AUC="+str(auc(train_fpr,train_tpr)))
plt.plot(test_fpr,test_tpr,label = "Test_AUC="+str(auc(test_fpr,test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("AUC")
plt.grid(linestyle='-', linewidth=2)
```



In [ ]:

```
print("F1-score:",f1_score(y_test,y_test_pred))
```

F1-score: 0.8973690447973452

## Gradient Boosting Algorithm

In [ ]:

```
from sklearn.ensemble import GradientBoostingClassifier
params = {'n_estimators': [10,20,30,40,50,60], 'max_depth': [2,3,5,7]}
rkf = StratifiedKFold(n_splits=11, random_state=42, shuffle=True)
gbdt_model = GradientBoostingClassifier(random_state=42)
#Used GridsearchCV for Hyper-parameter
grid = GridSearchCV(estimator = gbdt_model, param_grid = params, scoring = 'roc_auc', verbose = 1, n_jobs = -1, cv
=rkf)
grid.fit(X_train,y_train)

print("Best Score:" + str(grid.best_score_))
print("Best Parameters: " + str(grid.best_params_))
```

Fitting 11 folds for each of 24 candidates, totalling 264 fits  
Best Score:0.9757435926029662  
Best Parameters: {'max\_depth': 5, 'n\_estimators': 60}

In [ ]:

```
gbdt_clf = GradientBoostingClassifier(max_depth = 5, n_estimators = 60)
gbdt_clf.fit(X_train,y_train)

y_pred = gbdt_clf.predict(X_train)
gbdt_train_auc = roc_auc_score(y_train,y_pred)
print(gbdt_train_auc)
```

0.9228726871414403

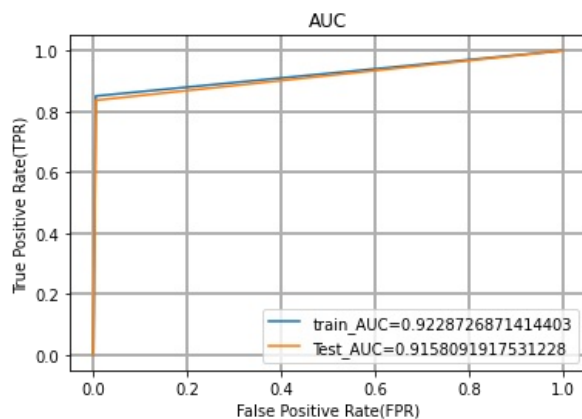
In [ ]:

```
y_test_pred = gbdt_clf.predict(X_test)
gbdt_test_auc = roc_auc_score(y_test,y_test_pred)
print(gbdt_test_auc)
```

0.9158091917531228

In [ ]:

```
train_fpr,train_tpr,tr_thresholds = roc_curve(y_train,y_pred)
test_fpr,test_tpr,te_thresholds = roc_curve(y_test,y_test_pred)
plt.plot(train_fpr,train_tpr,label = "train_AUC="+str(auc(train_fpr,train_tpr)))
plt.plot(test_fpr,test_tpr,label = "Test_AUC="+str(auc(test_fpr,test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("AUC")
plt.grid(linestyle='-', linewidth=2)
```



In [ ]:

```
print("f1-score:",f1_score(y_test,y_test_pred))
```

f1-score: 0.9067110899571633

## Stacking Classifier

In [53]:

```
import six
import sys
sys.modules['sklearn.externals.six'] = six
```

In [54]:

```
from mlxtend.classifier import StackingClassifier
from sklearn.svm import SVC
from sklearn.naive_bayes import MultinomialNB
```

In [56]:

```
#classifier 1
knn_model = KNeighborsClassifier(n_neighbors = 7)
knn_model.fit(X_train,y_train)

#Classifier 2
model = LogisticRegression(C= 10, penalty = 'l2', solver = 'liblinear')
model.fit(X_train,y_train)

#classifier 3
naive_clf = MultinomialNB(alpha=1000)
naive_clf.fit(X_train,y_train)

#classifier 4

tree_clf = DecisionTreeClassifier(max_depth = 9)
tree_clf.fit(X_train,y_train)

#classifier 5
gbdt_clf = GradientBoostingClassifier(max_depth = 5, n_estimators = 60)
gbdt_clf.fit(X_train,y_train)

#Stacking Classifier

sclf = StackingClassifier(classifiers=[knn_model,model,naive_clf,tree_clf,gbdt_clf],meta_classifier=model)

#fit the model
sclf.fit(X_train,y_train)

#predict in probabilities

y_pred = sclf.predict(X_train)
```

In [57]:

```
train_auc = roc_auc_score(y_train,y_pred)
print(train_auc)
```

0.9228726871414403

In [58]:

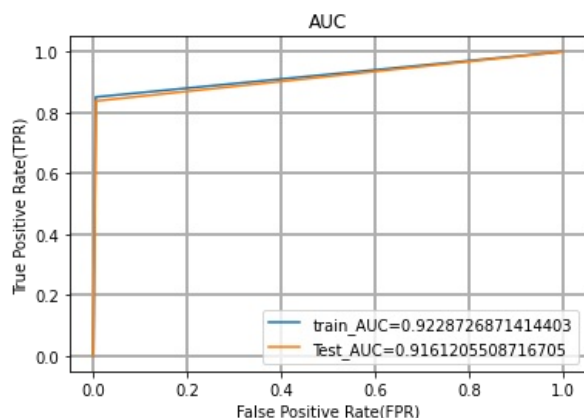
```
y_test_pred = sclf.predict(X_test)
test_auc = roc_auc_score(y_test,y_test_pred)
print(test_auc)
```

0.9161205508716705



In [59]:

```
train_fpr,train_tpr,tr_thresholds = roc_curve(y_train,y_pred)
test_fpr,test_tpr,te_thresholds = roc_curve(y_test,y_test_pred)
plt.plot(train_fpr,train_tpr,label = "train AUC="+str(auc(train_fpr,train_tpr)))
plt.plot(test_fpr,test_tpr,label = "Test AUC="+str(auc(test_fpr,test_tpr)))
plt.legend()
plt.xlabel("False Positive Rate(FPR)")
plt.ylabel("True Positive Rate(TPR)")
plt.title("AUC")
plt.grid(linestyle='--', linewidth=2)
```



In [60]:

```
print("F1 score:",f1_score(y_test,y_test_pred))
```

F1 score: 0.9070154577883471

## Summary of All Models

In [2]:

```
from texttable import Texttable
t = Texttable()
t.add_rows([['Model','Hyper-parameter','Train AUC','Test AUC','f1_score'],[ 'Knn_Model',r"{'n_neighbors': 7}",0.90,0.91,0.89],
            ['logistic Regresstion',r"{'C': 10, 'penalty': 'l2', 'solver': 'liblinear'}",0.91,0.91,0.90],
            ['Gradient Boosting Algorithm',r"{'max_depth': 5, 'n_estimators': 60}",0.92,0.91,0.90],
            ['Decision_tree Model',r"{'max_depth': 9}",0.91,0.90,0.89],
            ['Naive Bayes',r"{'alpha': 0.1}",0.66,0.66,0.49],
            ['Calibrated Model',"--",0.92,0.91,0.90]])

print(t.draw())
```

Model	Hyper-parameter	Train AUC	Test AUC	f1_score
Knn_Model	{'n_neighbors': 7}	0.900	0.910	0.890
logistic Regresstion	{'C': 10, 'penalty': 'l2', 'solver': 'liblinear'}	0.910	0.910	0.900
Gradient Boosting Algorithm	{'max_depth': 5, 'n_estimators': 60}	0.920	0.910	0.900
Decision_tree Model	{'max_depth': 9}	0.910	0.900	0.890
Naive Bayes	{'alpha': 0.1}	0.660	0.660	0.490
Calibrated Model	--	0.920	0.910	0.900

## Observation

We have done text-preprocessing technique on SQL query and apply Word2Vec to come up with vectorization.

We have applied 6 Machine learning model such as Knn\_Model, Logistic Regression, Gradient Boosting algorithm, Decision\_tree model, Naive Bayes Model and calibrated model.

Used AUC as evaluation metrics.

We found calibrated model as best model with Test AUC = 0.91 for this experiment.