



```
df['Processed Text']=df['Processed Text'].apply(lambda x: tokenization(x))
df['Processed Text']=df['Processed Text'].apply(lambda x: remove_punctuation(x))
df['Processed Text']=df['Processed Text'].apply(lambda x: remove_stopwords(x))

df['Processed Text']=df['Processed Text'].apply(lambda x: ' '.join(x))
```

```
df.head(5)
```

	CVE ID Number	Vulnerability Type	Summary Text	Processed Text
0	CVE-2022-43766	Denial of Service	Apache IoTDB version 0.12.2 to 0.12.6, 0.13.0 ...	apache iotdb version vulnerable denial ser...
1	CVE-2022-43365	Denial of Service	IP-COM EW9 V15.11.0.14(9732) was discovered to...	ipcom ew v discovered contain buffer overflow ...
	CVE-			

```
df.shape
```

```
(175850, 4)
```

```
df.sample(10)
```

	CVE ID Number	Vulnerability Type	Summary Text	Processed Text
159305	CVE-2009-2332	Gain Information	CMS Chainuk 1.2 and earlier allows remote attackers	cms chainuk earlier allows remote attackers o...
36660	CVE-2020-7154	Code Execute	A ifviewselectpage expression language injecti...	ifviewselectpage expression language injection...
47975	CVE-2016-1027	Code Execute	Adobe Flash Player before 18.0.0.343 and 19.x ...	adobe flash player x x windows os x linux a...
26560	CVE-2003-0926	Denial of Service	Ethereal 0.9.15 and earlier, and Tethereal, al...	ethereal earlier tethereal allows remote atta...
146399	CVE-2009-	Bypass	Mozilla Firefox executes DOM calls in response	mozilla firefox executes dom calls response is

## ▼ TF-IDF

TF-IDF->-> Term Frequency\*Inverse Document Frequency  
for finding out the importance of a token in a document

```
import heapq
from collections import defaultdict
freq=defaultdict(int)
```

```
#data is a dictionary with all vulnerability types as keys and it contains information on the occurrence of each token
for vul in data.keys():
    items=heapq.nlargest(3,data[vul][0],key=data[vul][0].get)
    d=defaultdict(int)
    for item in items:
        d[item]=data[vul][0][item]
    freq[vul]=d
```

```
freq['Overflow']
```

```
for item in freq.keys():
    display(item)
    display(freq[item])
```

```
import os
```

[illegible]



samp=10 #10 hard limit

#p=int((samp\*n)//100)

new=df.sample(10000)

len(new)

10000

new.head(5)

	CVE ID Number	Vulnerability Type	Summary Text	Processed Text
54561	CVE-2011-5102	Code Execute	The Investigative Reports web interface in the...	investigative reports web interface triton man...
164341	CVE-2010-3387	gain privilege	<b>** DISPUTED **</b> vdrleaktest in Video Disk Reco...	disputed vdrleaktest video disk recorder vd...
115098	CVE-2011-2071	Cross Site Scripting	Cross-site scripting vulnerability in Message	crosssite scripting vulnerability message

```
newvect=TfidfVectorizer()  
Xnew=newvect.fit_transform(new['Processed Text'])  
ynew=new['Vulnerability Type']
```

```
from sklearn.linear_model import LogisticRegression  
lr = LogisticRegression()  
lr.fit(Xnew,ynew)
```

```
from sklearn.naive_bayes import GaussianNB  
gnb = GaussianNB()  
gnb.fit(Xnew.toarray(),ynew)
```

```
from sklearn.tree import DecisionTreeClassifier  
dtt = DecisionTreeClassifier()  
dtt.fit(Xnew,ynew)
```

```
from sklearn.ensemble import RandomForestClassifier  
rfc = RandomForestClassifier()  
rfc.fit(Xnew,ynew)
```

```
from sklearn.neighbors import KNeighborsClassifier  
knn = KNeighborsClassifier()  
knn.fit(Xnew,ynew)  
print(' ')
```

```
from sklearn.metrics import accuracy_score  
from collections import defaultdict  
dicnew=defaultdict(int)
```

```
lrnew=lr.predict(Xnew)  
dicnew['Logistic Regression Sample TF-IDF']=accuracy_score(ynew, lrnew)
```

```
gnbnew=gnb.predict(Xnew.toarray())  
dicnew['Gaussian NB Sample TF-IDF']=accuracy_score(ynew, gnbnew)
```

```
dttnew=dtt.predict(Xnew)  
dicnew['Decision Tree Sample TF-IDF']=accuracy_score(ynew, dttnew)
```

```
rfcnew=rfc.predict(Xnew)  
dicnew['Random Forest Sample TF-IDF']=accuracy_score(ynew, rfcnew)
```

```
knnnew=knn.predict(Xnew)  
dicnew['K Nearest Neighbour Sample TF-IDF']=accuracy_score(ynew, knnnew)
```

```
dicnew #TF-IDF vectorizer is used as it assigns value to a token based on its importance in the corpus(dataset).This is suppr  
defaultdict(int,  
            {'Logistic Regression Sample TF-IDF': 0.8395,  
             'Gaussian NB Sample TF-IDF': 0.8666,
```

```
'Decision Tree Sample TF-IDF': 0.8395,
'Random Forest Sample TF-IDF': 0.9633,
'K Nearest Neighbour Sample TF-IDF': 0.8395})
```

```
new['TF-IDF Predicted']=rfcnew
```

```
new.sample(10)
```

	CVE ID Number	Vulnerability Type	Summary Text	Processed Text	TF-IDF LR predicted
80083	CVE-2017-8748	Overflow	Internet Explorer in Microsoft Windows 7 SP1, ...	internet explorer microsoft windows sp window...	Overflow
9343	CVE-2017-2531	Denial of Service	An issue was discovered in certain Apple produ...	issue discovered certain apple products ios a...	Denial of Service
39271	CVE-2019-7797	Code Execute	Adobe Acrobat and Reader versions 2019.010.201...	adobe acrobat reader versions earlier earlie...	Code Execute
71082	CVE-1999-0935	Code Execute	classifieds.cgi allows remote attackers to exe...	classifiedscgi allows remote attackers execute...	Code Execute
	CVE-		Microsoft Publisher		

## ▼ Best-Models TF-IDF

```
from sklearn.linear_model import LogisticRegression
lr = LogisticRegression()
lr.fit(X,y1)
```

```
LogisticRegression()
```

```
y1r=lr.predict(X)
```

```
from sklearn.metrics import accuracy_score
accuracy_score(y1r, y1)
```

```
0.7195280068239978
```

```
df['TF-IDF LR predicted']=y1r
```

```
df.sample(10)
```

	CVE ID Number	Vulnerability Type	Summary Text	Processed Text	TF-IDF LR predicted
94425	CVE-2021-0399	Memory Corruption	In qtaguid_untag of xt_qtaguid.c, there is a p...	qtaguiduntag xtqtaguidc possible memory corrup...	Memory Corruption
4638	CVE-2019-11479	Denial of Service	Jonathan Looney discovered that the Linux kern...	jonathan looney discovered linux kernel defaul...	Denial of Service
61997	CVE-2008-0300	Code Execute	mapFiler.php in Mapbender 2.4 to 2.4.4 allows ...	mapfilerphp mapbender allows remote attacker...	Code Execute
76198	CVE-2019-10877	Overflow	In Teeworlds 0.7.2, there is an integer overfl...	teeworlds integer overflow cmapload enginesha...	Overflow
410248	CVE-2005-	SQL Injection	SQL injection vulnerability in	sql injection vulnerability	Code

```
from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier()
#rfc.fit(X, y1)
```

```
#y1rfc=rfc.predict(X)
```

```
from sklearn.metrics import accuracy_score
#accuracy_score(yrffc, y1)
```

```
#df['TF-IDF RFC predicted']=yrffc
```

```
# df.head(10)
```

▼ **Classifier**

```
from sklearn.feature_extraction.text import TfidfVectorizer
vec = TfidfVectorizer()
```

▼ **TF-IDF**

```
n=len(df['Processed Text'])
print(n)
```

175850

```
from sklearn.feature_extraction.text import TfidfVectorizer
nvectorizer = TfidfVectorizer()
X = nvectorizer.fit_transform(df['Processed Text'])
y1=df['Vulnerability Type']
```

```
X.shape,y1.shape
```

((175850, 138817), (175850,))

```
df.head(10)
```

	CVE ID Number	Vulnerability Type	Summary Text	Processed Text
0	CVE-2022-43766	Denial of Service	Apache IoTDB version 0.12.2 to 0.12.6, 0.13.0 ...	apache iotdb version vulnerable denial ser...
1	CVE-2022-43365	Denial of Service	IP-COM EW9 V15.11.0.14(9732) was discovered to...	ipcom ew v discovered contain buffer overflow ...
2	CVE-2022-43035	Denial of Service	An issue was discovered in Bento4 v1.6.0-639. ...	issue discovered bento v heapbufferoverflow ap...
3	CVE-2022-43033	Denial of Service	An issue was discovered in Bento4 1.6.0-639. T...	issue discovered bento bad free component aph...
4	CVE-2022-	Denial of Service	The py library through 1.11.0 for Python allow	py library python allows remote attackers con

```
train=df[df['Vulnerability Type']=='Denial of Service']
```

```
train.shape
```

(28100, 4)

```
train.head(5)
```

	CVE ID Number	Vulnerability Type	Summary Text	Processed Text
0	CVE-2022-43766	Denial of Service	Apache IoTDB version 0.12.2 to 0.12.6, 0.13.0 ...	apache iotdb version vulnerable denial ser...
1	CVE-2022-43365	Denial of Service	IP-COM EW9 V15.11.0.14(9732) was discovered to...	ipcom ew v discovered contain buffer overflow ...
	CVE-			

```
df['DoS']=df['Denial of Service'].where()
```