

```
In [5]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import nltk
import seaborn as sns
import string
```

```
In [6]: DF=pd.read_csv('spam.csv',encoding='latin-1')
DF.head(10)
```

```
Out[6]:
```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	ham	Go until jurong point, crazy.. Available only ...	NaN	NaN	NaN
1	ham	Ok lar... Joking wif u oni...	NaN	NaN	NaN
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	NaN	NaN	NaN
3	ham	U dun say so early hor... U c already then say...	NaN	NaN	NaN
4	ham	Nah I don't think he goes to usf, he lives aro...	NaN	NaN	NaN
5	spam	FreeMsg Hey there darling it's been 3 week's n...	NaN	NaN	NaN
6	ham	Even my brother is not like to speak with me. ...	NaN	NaN	NaN
7	ham	As per your request 'Melle Melle (Oru Minnamin...	NaN	NaN	NaN
8	spam	WINNER!! As a valued network customer you have...	NaN	NaN	NaN
9	spam	Had your mobile 11 months or more? U R entitle...	NaN	NaN	NaN

```
In [7]: DF=DF.drop(columns=['Unnamed: 2','Unnamed: 3','Unnamed: 4'])
```

```
In [8]: DF.v1.value_counts()
```

```
Out[8]: ham      4825
spam      747
Name: v1, dtype: int64
```

```
In [9]: DF.shape
```

```
Out[9]: (5572, 2)
```

```
In [10]: for i in range(0,5572):
if DF['v1'][i]=="spam":
    DF['v1'][i]=1
else:
    DF['v1'][i]=0
DF['v1']
```

```
Out[10]: 0      0
         1      0
         2      1
         3      0
         4      0
         ..
        5567    1
        5568    0
        5569    0
        5570    0
        5571    0
        Name: v1, Length: 5572, dtype: object
```

```
In [11]: DF.v1.value_counts()
```

```
Out[11]: 0      4825
         1       747
        Name: v1, dtype: int64
```

```
In [12]: DF.rename(columns={'v1': 'target', 'v2': 'text'}, inplace=True)
         DF.sample(10)
```

```
Out[12]:
```

	target	text
3403	0	Then I ask dad to pick I up lar... I wan 2 ...
2438	0	I not busy juz dun wan 2 go so early.. Hee..
3826	1	Congratulations U can claim 2 VIP row A Ticket...
393	0	Yes i think so. I am in office but my lap is i...
332	1	Call Germany for only 1 pence per minute! Call...
2586	0	If you don't respond imma assume you're still ...
4832	1	New Mobiles from 2004, MUST GO! Txt: NOKIA to ...
3773	0	Ok... But bag again..
2664	0	He remains a bro amongst bros
270	0	Come to mu, we're sorting out our narcotics si...

```
In [13]: DF.duplicated().sum()
```

```
Out[13]: 403
```

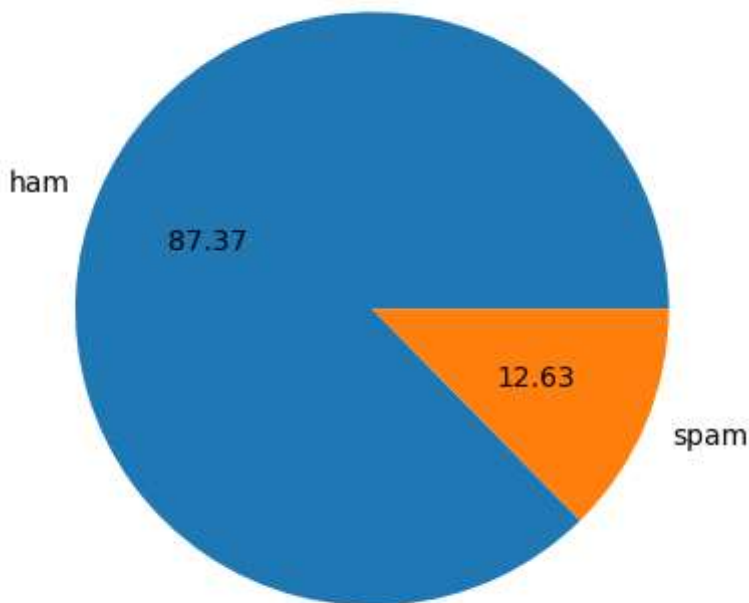
```
In [14]: DF=DF.drop_duplicates(keep='first')
         DF.duplicated().sum()
```

```
Out[14]: 0
```

```
In [15]: DF.shape
```

```
Out[15]: (5169, 2)
```

```
In [16]: plt.pie(DF['target'].value_counts(), labels=['ham', 'spam'], autopct="%0.2f")
         plt.show()
```



```
In [24]: nltk.download('punkt')
```

```
[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\91772\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
```

```
Out[24]: True
```

```
In [20]: DF['num_char']=DF['text'].apply(len)
```

```
In [21]: DF['num_word']=DF['text'].apply(lambda x:len(nltk.word_tokenize(x)))
```

```
Out[21]:
```

	target	text	num_char	num_word
0	0	Go until jurong point, crazy.. Available only ...	111	24
1	0	Ok lar... Joking wif u oni...	29	8
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37
3	0	U dun say so early hor... U c already then say...	49	13
4	0	Nah I don't think he goes to usf, he lives aro...	61	15

```
In [23]: DF['num_sent']=DF['text'].apply(lambda x:len(nltk.sent_tokenize(x)))
```

```
In [25]: plt.figure(figsize=(50,10))
sns.histplot(DF[DF['target']==0]['num_char'])
sns.histplot(DF[DF['target']==1]['num_char'],color='green')

plt.figure(figsize=(50,10))
sns.histplot(DF[DF['target']==0]['num_char'])
sns.histplot(DF[DF['target']==1]['num_char'],color='orange')

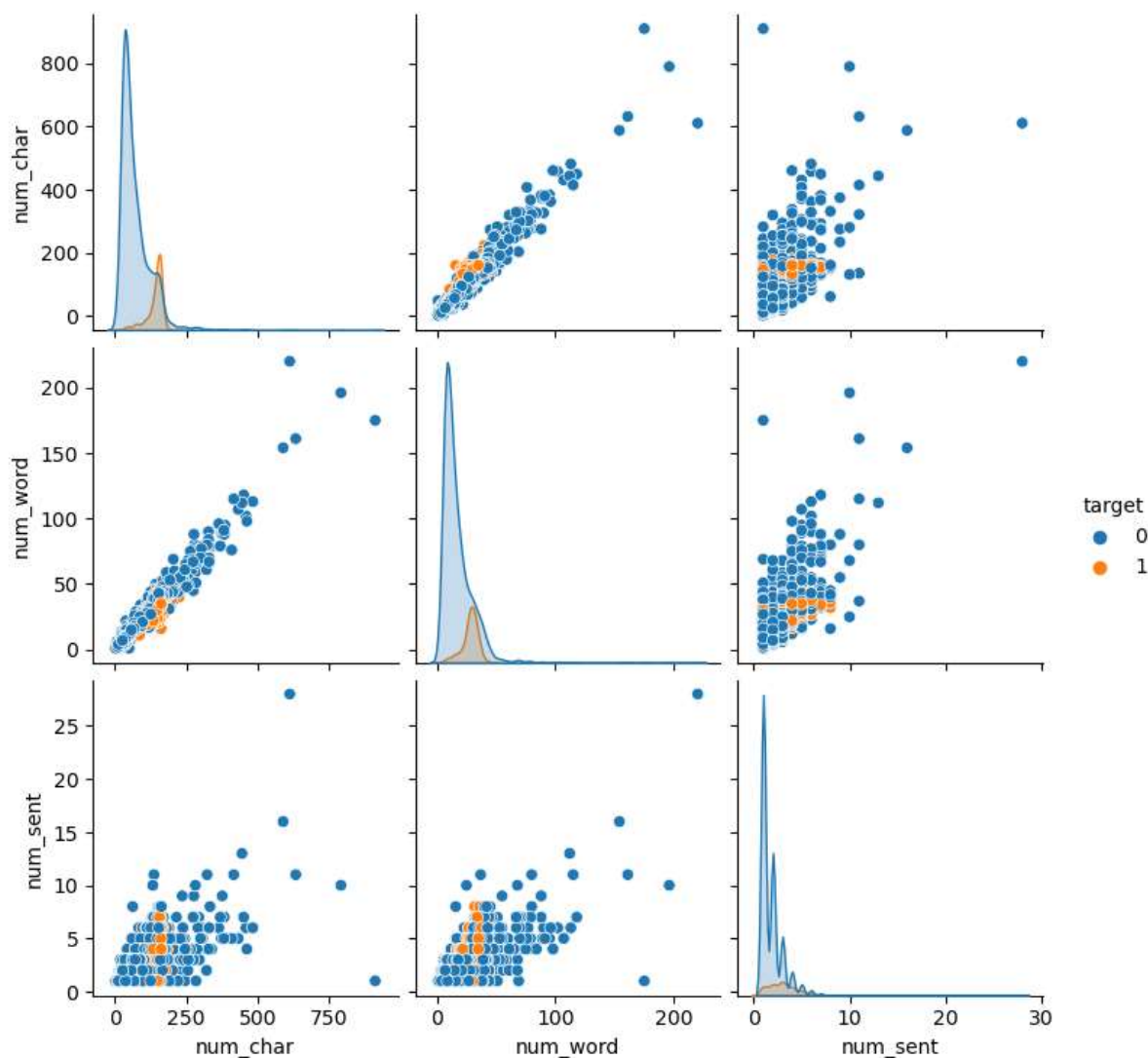
plt.figure(figsize=(50,10))
sns.histplot(DF[DF['target']==0]['num_sent'])
sns.histplot(DF[DF['target']==1]['num_sent'],color='pink')
```

Out[25]: <Axes: xlabel='num_sent', ylabel='Count'>



```
In [26]: plt.figure(figsize=(30,10))
sns.pairplot(DF,hue='target')
```

Out[26]: <seaborn.axisgrid.PairGrid at 0x20afb9688d0>
<Figure size 3000x1000 with 0 Axes>

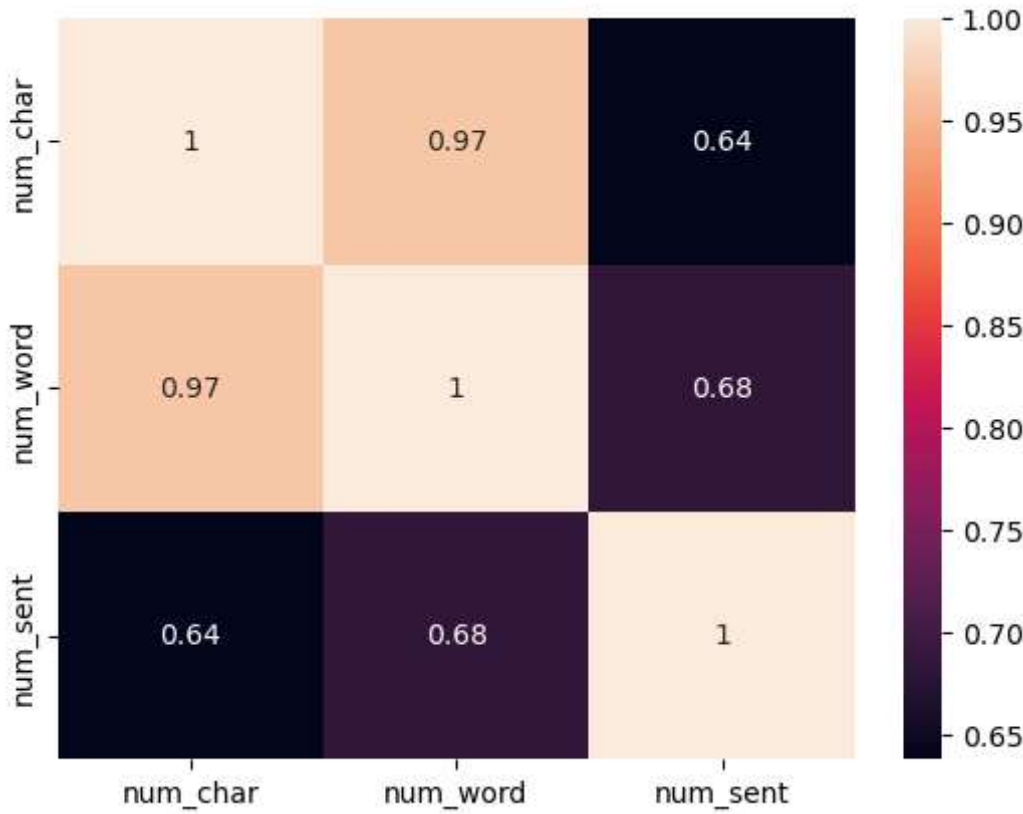


```
In [27]: sns.heatmap(DF.corr(),annot=True)
```

C:\Users\91772\AppData\Local\Temp\ipykernel_13444\1379925368.py:1: FutureWarning:
The default value of `numeric_only` in `DataFrame.corr` is deprecated. In a future version, it will default to `False`. Select only valid columns or specify the value of `numeric_only` to silence this warning.

```
sns.heatmap(DF.corr(),annot=True)
```

```
Out[27]: <Axes: >
```



```
In [28]: DF.head()
```

target			text	num_char	num_word	num_sent
0	0	Go until jurong point, crazy.. Available only ...		111	24	2
1	0	Ok lar... Joking wif u oni...		29	8	2
2	1	Free entry in 2 a wkly comp to win FA Cup fina...		155	37	2
3	0	U dun say so early hor... U c already then say...		49	13	1
4	0	Nah I don't think he goes to usf, he lives aro...		61	15	1

```
In [29]: DF.describe()
```

	num_char	num_word	num_sent
count	5169.000000	5169.000000	5169.000000
mean	78.977945	18.453279	1.947185
std	58.236293	13.324793	1.362406
min	2.000000	1.000000	1.000000
25%	36.000000	9.000000	1.000000
50%	60.000000	15.000000	1.000000
75%	117.000000	26.000000	2.000000
max	910.000000	220.000000	28.000000

```
In [30]: #Processing of data
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
```

```
ps= PorterStemmer()  
ps.stem('sleeping')
```

Out[30]: 'sleep'

```
In [31]: def transform_text(text):  
         text=text.lower()  
         text= nltk.word_tokenize(text)  
  
         y=[]  
         for i in text:  
             if i.isalnum():  
                 y.append(i)  
  
         text=y[:]  
         y.clear()  
  
         for i in text:  
             if i not in stopwords.words('English') and i not in string.punctuation:  
                 y.append(i)  
  
         text=y[:]  
         y.clear()  
  
         for i in text:  
             y.append(ps.stem(i))  
  
         return " ".join(y)
```

```
In [32]: DF.head(10)
```

Out[32]:

	target	text	num_char	num_word	num_sent
0	0	Go until jurong point, crazy.. Available only ...	111	24	2
1	0	Ok lar... Joking wif u oni...	29	8	2
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37	2
3	0	U dun say so early hor... U c already then say...	49	13	1
4	0	Nah I don't think he goes to usf, he lives aro...	61	15	1
5	1	FreeMsg Hey there darling it's been 3 week's n...	148	39	4
6	0	Even my brother is not like to speak with me. ...	77	18	2
7	0	As per your request 'Melle Melle (Oru Minnamin...	160	31	2
8	1	WINNER!! As a valued network customer you have...	158	32	5
9	1	Had your mobile 11 months or more? U R entitle...	154	31	3

```
In [ ]:
```