

# Project Name:- Malignant-Comments-Classifer.

## Importing Libraries

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
from sklearn.model_selection import train_test_split
```

## Loading Dataset

```
In [6]: import pandas as pd
df=pd.read_csv('train.csv')
```

```
df.head()
```

	id	comment_text	malignant	highly_malignant	rule	threat	abuse	loathe
Out[6]:	0	000097932f770f	Explanationwhy the edits made under my usern...	0	0	0	0	0
	1	000239393b600f	Drawn He matches the background colour f'm s...	0	0	0	0	0
	2	00013970cc020e	Hey man, i'm really not trying to edit we'r...	0	0	0	0	0
	3	000234a1c15b627b	"Alkenet cant make any real suggestions on...	0	0	0	0	0
	4	000199261d45d635	You, sr, are my hero. Any chance you remember...	0	0	0	0	0

```
In [8]: df
```

	id	comment_text	malignant	highly_malignant	rule	threat	abuse	loathe
Out[8]:	0	000097932f770f	Explanationwhy the edits made under my usern...	0	0	0	0	0
	1	000239393b600f	Drawn He matches the background colour f'm s...	0	0	0	0	0
	2	00013970cc020e	Hey man, i'm really not trying to edit we'r...	0	0	0	0	0
	3	000234a1c15b627b	"Alkenet cant make any real suggestions on...	0	0	0	0	0
	4	000199261d45d635	You, sr, are my hero. Any chance you remember...	0	0	0	0	0
	...	...	...	...	...	...	...	...
	15556	8d07f799590f1e	...And for the second time of asking, when...	0	0	0	0	0
	15567	8e4a6e0c3b46d0	You should be ashamed of yourself for that h...	0	0	0	0	0
	15568	8ec3a0a5c2c1c9	Spoken withArtes, seems no actual article for...	0	0	0	0	0
	15569	8f12570f64aaf3	And it looks like it was actually you who put...	0	0	0	0	0
	15570	8f45c42a0719fa	"Arkel... i really don't think you understand...	0	0	0	0	0

159571 rows x 8 columns

```
In [88]: df.shape
```

```
Out[88]: (159571, 8)
```

```
In [77]: df['malignant'].unique()
```

```
Out[77]: array([0, 1], dtype=int64)
```

```
In [77]: df['highly_malignant'].shape
```

```
Out[77]: (159571,)
```

```
Out[9]:
```

Checking Null Values

```
In [149]: df.isnull()
```

	id	comment_text	malignant	highly_malignant	rule	threat	abuse	loathe
Out[149]:	0	False	False	False	False	False	False	False
	1	False	False	False	False	False	False	False
	2	False	False	False	False	False	False	False
	3	False	False	False	False	False	False	False
	4	False	False	False	False	False	False	False
	...	...	...	...	...	...	...	...
	15556	False	False	False	False	False	False	False
	15567	False	False	False	False	False	False	False
	15568	False	False	False	False	False	False	False
	15569	False	False	False	False	False	False	False
	15570	False	False	False	False	False	False	False

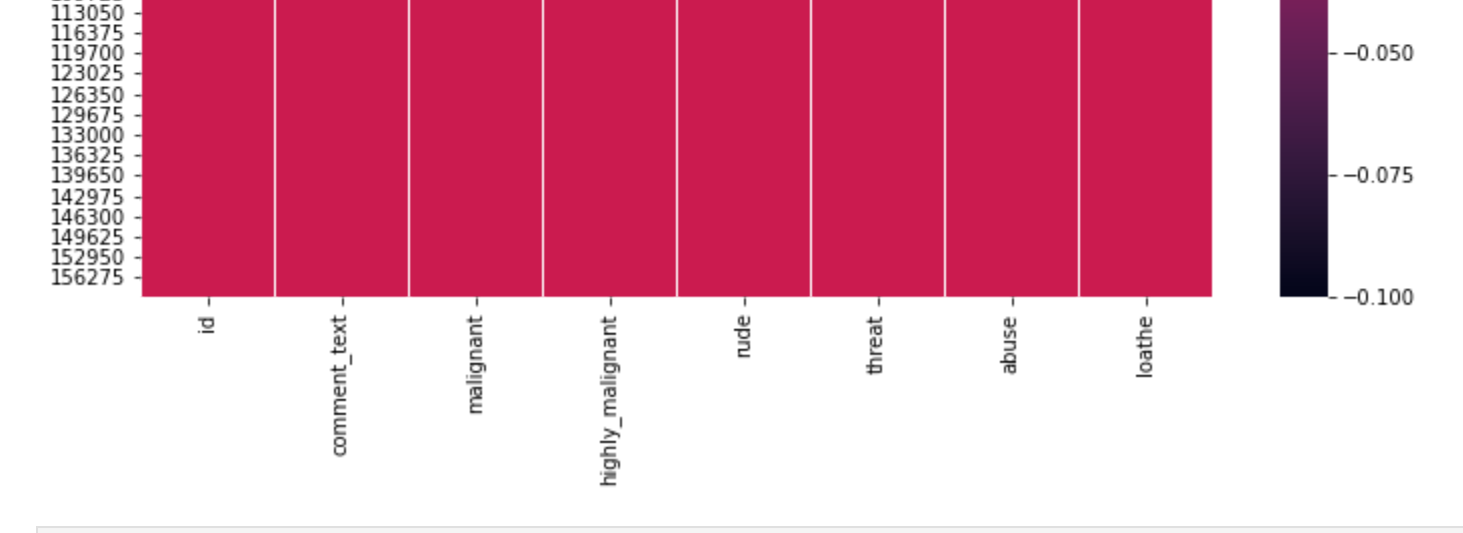
159571 rows x 8 columns

```
In [11]: df.isnull().sum()
```

id	0
comment_text	0
malignant	0
highly_malignant	0
rule	0
threat	0
abuse	0
loathe	0
dtype: int64	

Checking Isnull Heatmap

```
In [14]: fig, ax = plt.subplots(figsize=(12,9))
sns.heatmap(df.isnull(),ax=ax);
```



```
In [12]: df.isnull().sum()
```

```
Out[12]: 0
```

## Describing Dataset

```
In [13]: df.describe()
```

	malignant	highly_malignant	rule	threat	abuse	loathe
count	159571.000000	128871.000000	159571.000000	159571.000000	159571.000000	159571.000000
mean	0.006944	0.009966	0.032340	0.002286	0.049364	0.008305
std	0.254079	0.090477	0.223021	0.054850	0.216827	0.093420
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
max	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000

Describing Dataset Heatmap

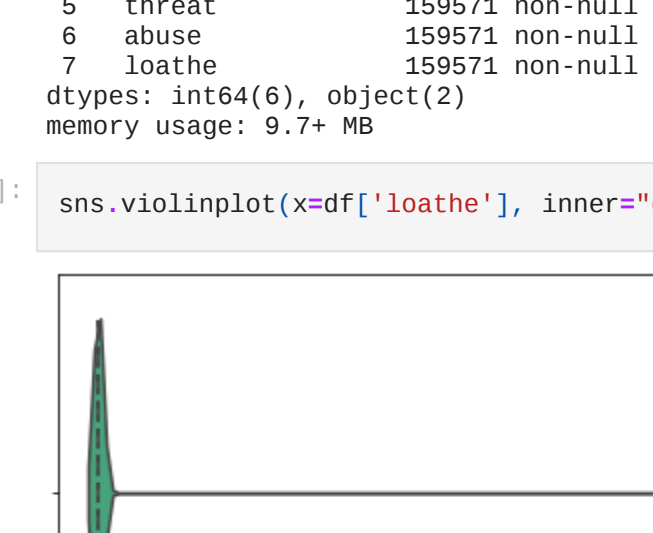
```
In [17]: import matplotlib.pyplot as plt
import seaborn as sns
fig, ax = plt.subplots(figsize=(12,9))
sns.heatmap(df.describe(),ax=ax);
```



```
In [18]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 159571 entries, 0 to 159570
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  -
0   id                    159571 non-null    object
1   comment_text         159571 non-null    object
2   malignant            159571 non-null    int64
3   highly_malignant     159571 non-null    int64
4   rule                 159571 non-null    int64
5   threat              159571 non-null    int64
6   abuse               159571 non-null    int64
7   loathe              159571 non-null    int64
dtypes: object(2), int64(6)
memory usage: 9.7+ MB
```

```
In [20]: sns.violinplot(x=df['loathe'],inner='quartile',color='#322772');
```

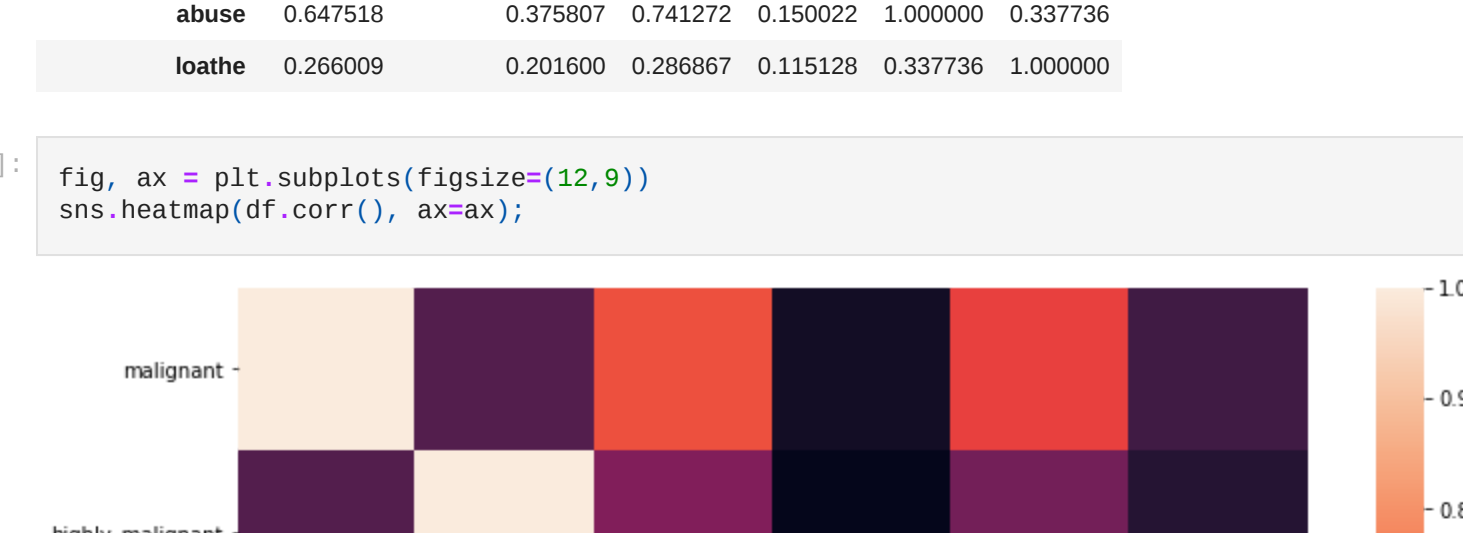


## Correlation

```
In [22]: df.corr()
```

	malignant	highly_malignant	rule	threat	abuse	loathe
malignant	1.000000	0.300129	0.470525	0.337058	0.647218	0.266406
highly_malignant	0.300129	1.000000	0.403054	0.228601	0.378807	0.205600
rule	0.470525	0.403054	1.000000	0.341179	0.741372	0.208667
threat	0.337058	0.228601	0.341179	1.000000	0.500002	0.131028
abuse	0.647218	0.378807	0.741372	0.500002	1.000000	0.337736
loathe	0.266406	0.205600	0.208667	0.131028	0.337736	1.000000

```
In [21]: fig, ax = plt.subplots(figsize=(12,9))
sns.heatmap(df.corr(),ax=ax);
```



## EXPLORATORY DISTRIBUTION ANALYSIS

### Scatterplot

```
In [23]: sns.scatterplot(x='malignant',y='loathe',data=df)
```

```
Out[23]: <AxesSubplot:label='malignant', ylabel='loathe'>
```



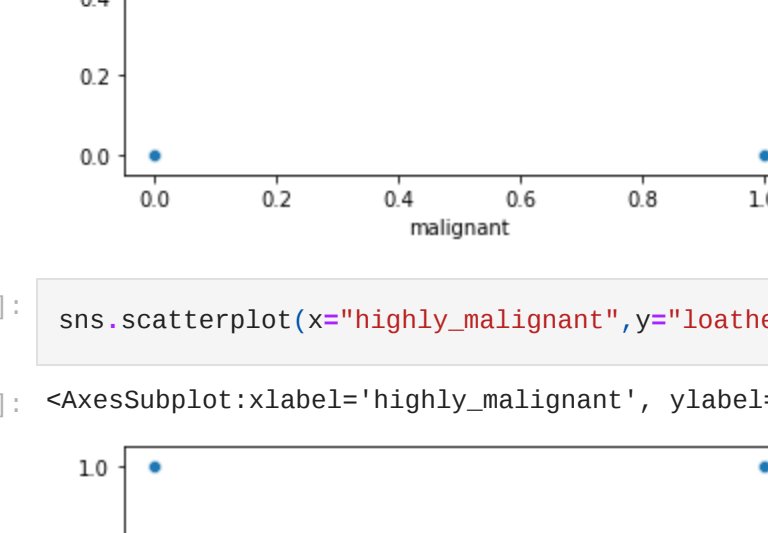
```
In [24]: sns.scatterplot(x='highly_malignant',y='loathe',data=df)
```

```
Out[24]: <AxesSubplot:label='highly_malignant', ylabel='loathe'>
```



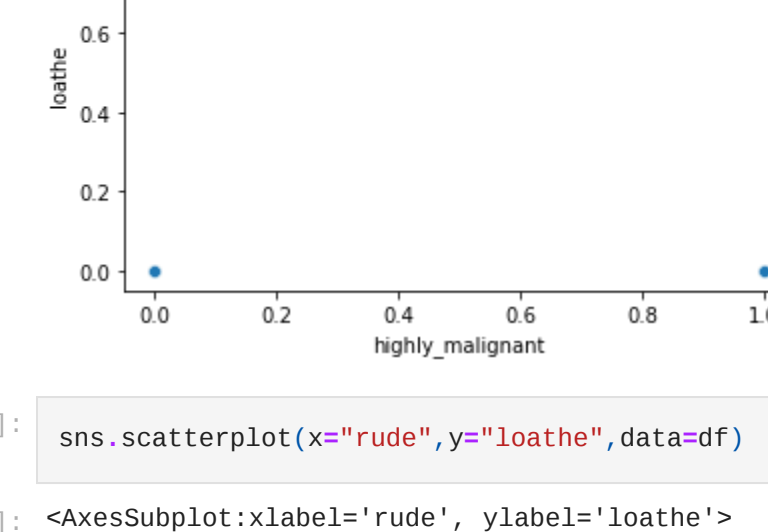
```
In [25]: sns.scatterplot(x='rule',y='loathe',data=df)
```

```
Out[25]: <AxesSubplot:label='rule', ylabel='loathe'>
```



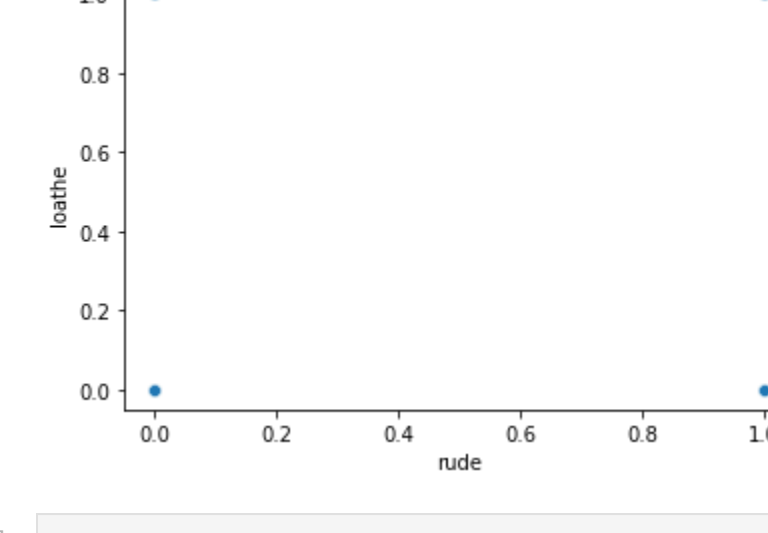
```
In [26]: sns.scatterplot(x='threat',y='loathe',data=df)
```

```
Out[26]: <AxesSubplot:label='threat', ylabel='loathe'>
```

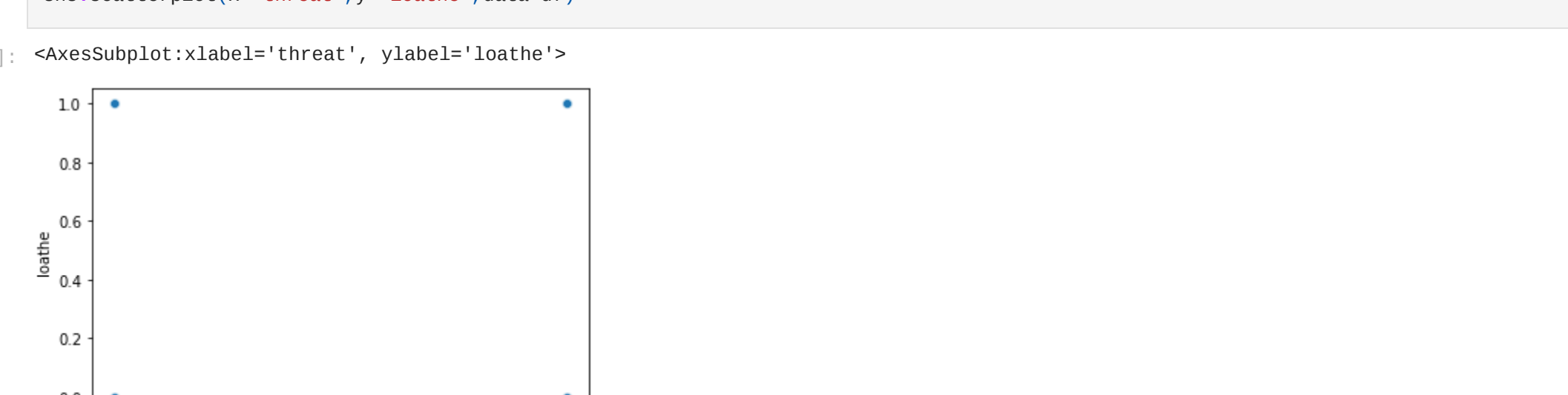


```
In [27]: sns.scatterplot(x='abuse',y='loathe',data=df)
```

```
Out[27]: <AxesSubplot:label='abuse', ylabel='loathe'>
```



```
In [29]: sns.pairplot(df)
```



```
In [38]: df
```

	id	comment_text	malignant	highly_malignant	rule	threat	abuse	loathe
Out[38]:	0	000097932f770f	Explanationwhy the edits made under my usern...	0	0	0	0	0
	1	000239393b600f	Drawn He matches the background colour f'm s...	0	0	0	0	0
	2	00013970cc020e	Hey man, i'm really not trying to edit we'r...	0	0	0	0	0
	3	000234a1c15b627b	"Alkenet cant make any real suggestions on...	0	0	0	0	0
	4	000199261d45d635	You, sr, are my hero. Any chance you remember...	0	0	0	0	0
	...	...	...	...	...	...	...	...
	15556	8d07f799590f1e	...And for the second time of asking, when...	0	0	0	0	0
	15567	8e4a6e0c3b46d0	You should be ashamed of yourself for that h...	0	0	0	0	0
	15568	8ec3a0a5c2c1c9	Spoken withArtes, seems no actual article for...	0	0	0	0	0
	15569	8f12570f64aaf3	And it looks like it was actually you who put...	0	0	0	0	0
	15570	8f45c42a0719fa	"Arkel... i really don't think you understand...	0	0	0	0	0

159571 rows x 8 columns

### skewness

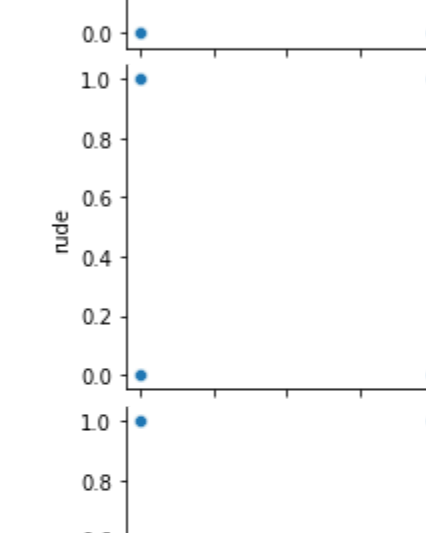
```
In [37]: df.skew()
```

malignant	2.748864
highly_malignant	9.851222
rule	3.180257
threat	0.135128
abuse	1.888088
loathe	10.535863
dtype: float64	

## Normal Distribution Curve

```
In [39]: sns.distplot(df['malignant'])
```

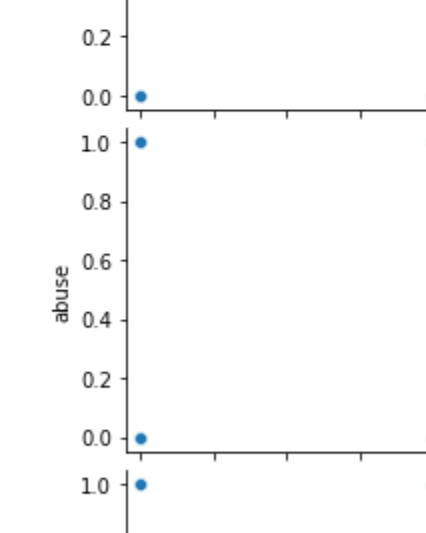
```
Out[39]: <seaborn.axisgrid.FacetGrid at 8x22c9648a70>
```



The data of the column is not normalised. The building blocks is out of the normalised curve

```
In [40]: sns.distplot(df['highly_malignant'])
```

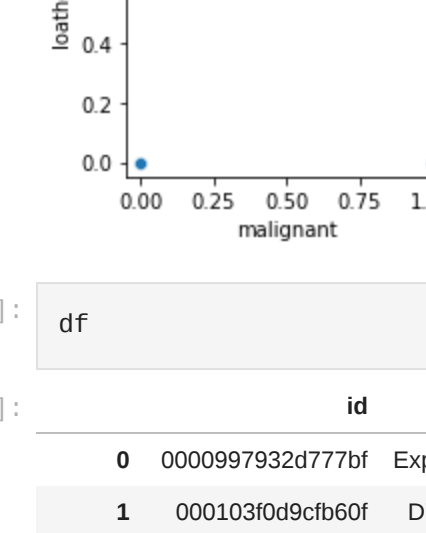
```
Out[40]: <seaborn.axisgrid.FacetGrid at 8x22c9647aacd0>
```



The data of the column is not normalised. The building blocks is out of the normalised curve

```
In [41]: sns.distplot(df['rule'])
```

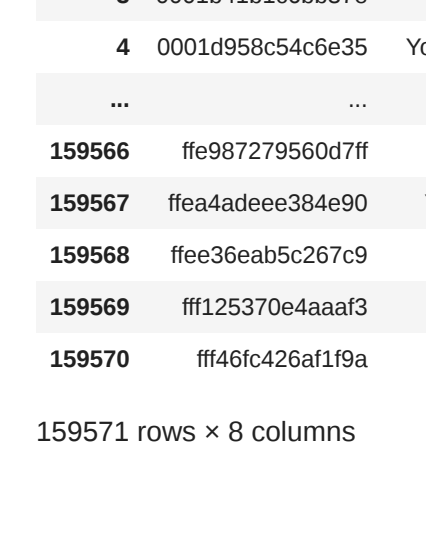
```
Out[41]: <seaborn.axisgrid.FacetGrid at 8x22c9642748b0>
```



The data of the column is not normalised. The building blocks is out of the normalised curve

```
In [42]: sns.distplot(df['threat'])
```

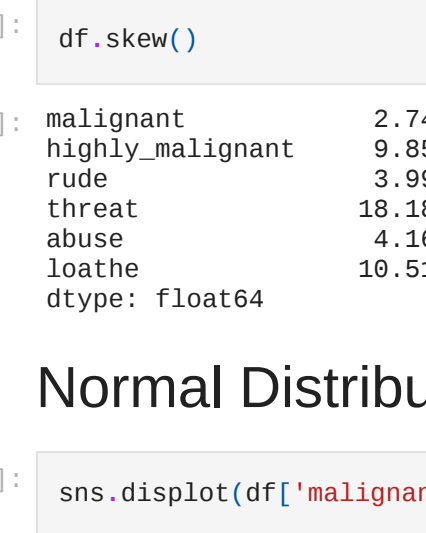
```
Out[42]: <seaborn.axisgrid.FacetGrid at 8x22c9642090b0>
```



The data of the column is not normalised. The building blocks is out of the normalised curve

```
In [43]: sns.distplot(df['abuse'])
```

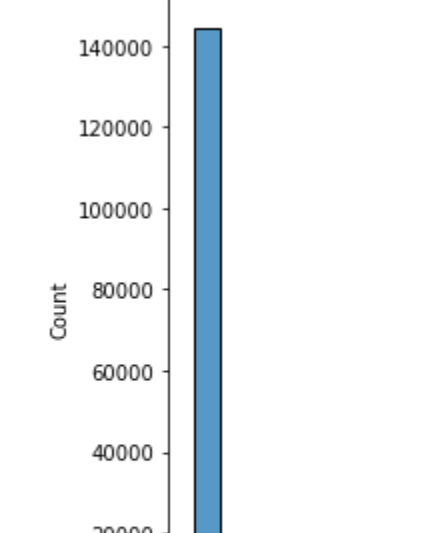
```
Out[43]: <seaborn.axisgrid.FacetGrid at 8x22c964b0d900>
```



The data of the column is not normalised. The building blocks is out of the normalised curve

```
In [44]: sns.distplot(df['loathe'])
```

```
Out[44]: <seaborn.axisgrid.FacetGrid at 8x22c9643355b0>
```



The data of the column is not normalised. The building blocks is out of the normalised curve

```
In [45]: df.corr()['loathe']
```

malignant	0.266406
highly_malignant	0.205600
rule	0.208667
threat	0.131028
abuse	0.337736
loathe	1.000000
dtype: float64	

```
In [52]: from pandas import read_csv
url = 'train.csv'
dataframe = read_csv(url, header=None)
# split into input and output columns
data = dataframe.values
X, y = data[:,1:], data[:,0]
print(X.shape, y.shape)

(159572, 7) (159572,)
```

```
In [53]: X, y = data[:,1:], data[:,0]
```

```
print(X.shape, y.shape)
```

```
(159572, 7) (159572,)
```

```
In [58]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=1)
print(X_train.shape, X_test.shape, y_train.shape, y_test.shape)
```

```
(139, 66) (59, 66) (139,) (69,)
```

```
In [59]: model = RandomForestClassifier(random_state=1)
```

```
model.fit(X_train, y_train)
```

```
Out[59]: RandomForestClassifier
```

```
In [60]: # make predictions
```

```
yhat = model.predict(X_test)
```

```
# evaluate predictions
```

```
acc = accuracy_score(y_test, yhat)
```

```
print('Accuracy: %f' % acc)
```

```
Accuracy: 0.783
```

```
In [54]: from sklearn import read_csv
```

```
url = 'train.csv'
```

```
dataframe = read_csv(url, header=None)
```

```
# summarize shape
```

```
print(dataframe.shape)
```

```
(159572, 8)
```

```
In [57]: # split into inputs and outputs
```

```
X, y = data[:,1:], data[:,0]
```

```
print(X.shape, y.shape)
```

```
(159572, 7) (159572,)
```

```
In [71]: df
```

	id	comment_text	malignant	highly_malignant	rule	threat	abuse	loathe
Out[71]:	0	000097932f770f	Explanationwhy the edits made under my usern...	0	0	0	0	0
	1	000239393b600f	Drawn He matches the background colour f'm s...	0	0	0	0	0
	2	00013970cc020e	Hey man, i'm really not trying to edit we'r...	0	0	0	0	0
	3	000234a1c15b627b	"Alkenet cant make any real suggestions on...	0	0	0	0	0
	4	000199261d45d635	You, sr, are my hero. Any chance you remember...	0	0	0	0	0
	...	...	...	...	...	...	...	...
	15556	8d07f799590f1e	...And for the second time of asking, when...	0	0	0	0	0
	15567	8e4a6e0c3b46d0	You should be ashamed of yourself for that h...	0	0	0	0	0
	15568	8ec3a0a5c2c1c9	Spoken withArtes, seems no actual article for...	0	0	0	0	0
	15569	8f12570f64aaf3	And it looks like it was actually you who put...	0	0	0	0	0
	15570	8f45c42a0719fa	"Arkel... i really don't think you understand...	0	0	0	0	0

159571 rows x 8 columns