

EMAIL SPAM DETECTION

We've all been the recipient of spam emails before. Spam mail, or junk mail, is a type of email that is sent to a massive number of users at one time, frequently containing cryptic messages, scams, or most dangerously, phishing content.

In this Project, use Python to build an email spam detector. Then, use machine learning to train the spam detector to recognize and classify emails into spam and non-spam. Let's get started!

Importing required libraries

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

Read the CSV file

```
In [8]: email = pd.read_csv('spam.csv', encoding = 'ISO-8859-1')
email
```

```
Out[8]:
```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	ham	Go until jurong point, crazy.. Available only in Sun...	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 final tk...	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
...
5567	spam	This is the 2nd time we have tried 2 contact u... We are sorry the 1st time we could not reach u...	NaN	NaN	NaN
5568	ham	Will i_b going to esplanade fr home?	NaN	NaN	NaN
5569	ham	Pity, * was in mood for that. So...any other s...	NaN	NaN	NaN
5570	ham	The guy did some bitching but I acted like i'd... nothing really happened.	NaN	NaN	NaN
5571	ham	Rofl. Its true to its name	NaN	NaN	NaN

5572 rows x 5 columns

```
In [9]: email.head()
```

```
Out[9]:
```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	ham	Go until jurong point, crazy.. Available only in Sun...	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 final tk...	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

```
In [10]: email.tail()
```

```
Out[10]:
```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
5567	spam	This is the 2nd time we have tried 2 contact u... We are sorry the 1st time we could not reach u...	NaN	NaN	NaN
5568	ham	Will i_b going to esplanade fr home?	NaN	NaN	NaN
5569	ham	Pity, * was in mood for that. So...any other s...	NaN	NaN	NaN
5570	ham	The guy did some bitching but I acted like i'd... nothing really happened.	NaN	NaN	NaN
5571	ham	Rofl. Its true to its name	NaN	NaN	NaN

```
In [11]: email.columns
```

```
Out[11]: Index(['v1', 'v2', 'Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], dtype='object')
```

```
In [14]: email.shape
```

```
Out[14]: (5572, 5)
```

```
In [12]: email.size
```

```
Out[12]: 27860
```

```
In [17]: email.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 5 columns):
#   Column          Non-Null Count  Dtype
---  -
0    v1              5572 non-null    object
1    v2              5572 non-null    object
2    Unnamed: 2      50 non-null     object
3    Unnamed: 3      12 non-null     object
4    Unnamed: 4      6 non-null      object
dtypes: object(5)
memory usage: 217.8+ KB
```

```
In [16]: email.describe()
```

```
Out[16]:
```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
count	5572	5572	50	12	6
unique	2	5169	43	10	5
top	ham	Sorry, I'll call later	bt not his girfrnd... G o o d n i g h t ...@"	GE	GNT-)"
freq	4825	30	3	2	2

Data Cleaning

```
In [ ]: email.drop(columns=['Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'],inplace=True)
```

```
In [44]: email
```

```
Out[44]:
```

	v1	v2
0	ham	Go until jurong point, crazy.. Available only in Sun...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup final tk...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...
...
5567	spam	This is the 2nd time we have tried 2 contact u... We are sorry the 1st time we could not reach u...
5568	ham	Will i_b going to esplanade fr home?
5569	ham	Pity, * was in mood for that. So...any other s...
5570	ham	The guy did some bitching but I acted like i'd... nothing really happened.
5571	ham	Rofl. Its true to its name

5572 rows x 2 columns

```
In [47]: email = email.rename(columns={'v1':'Target', 'v2':'Message'})
email
```

```
Out[47]:
```

	Target	Message
0	ham	Go until jurong point, crazy.. Available only in Sun...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup final tk...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...
...
5567	spam	This is the 2nd time we have tried 2 contact u... We are sorry the 1st time we could not reach u...
5568	ham	Will i_b going to esplanade fr home?
5569	ham	Pity, * was in mood for that. So...any other s...
5570	ham	The guy did some bitching but I acted like i'd... nothing really happened.
5571	ham	Rofl. Its true to its name

5572 rows x 2 columns

```
In [49]: email.isnull().sum()
```

```
Out[49]: Target      0
Message      0
dtype: int64
```

```
In [50]: email.duplicated().sum()
```

```
Out[50]: 403
```

```
In [51]: email.drop_duplicates(keep='first',inplace=True)
```

```
In [53]: email.duplicated().sum()
```

```
Out[53]: 0
```

```
In [54]: email.size
```

```
Out[54]: 10338
```

Label Encoding

```
In [55]: from sklearn.preprocessing import LabelEncoder
encoder=LabelEncoder()
email['Target']=encoder.fit_transform(email['Target'])
email['Target']
```

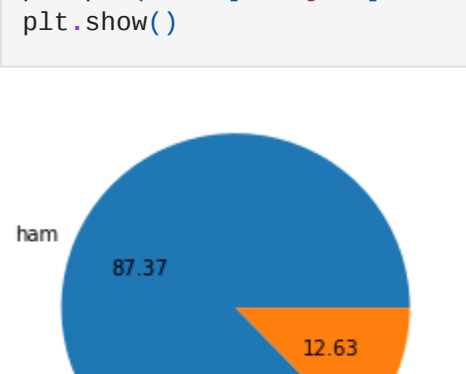
```
Out[55]: 0      0
1      0
2      1
3      0
4      0
...    .
5567   1
5568   0
5569   0
5570   0
5571   0
Name: Target, Length: 5169, dtype: int32
```

```
In [56]: email.head()
```

```
Out[56]:
```

	Target	Message
0	0	Go until jurong point, crazy.. Available only in Sun...
1	0	Ok lar... Joking wif u oni...
2	1	Free entry in 2 a wkly comp to win FA Cup final...
3	0	U dun say so early hor... U c already then say...
4	0	Nah I don't think he goes to usf, he lives aro...

```
In [57]: plt.pie(email['Target'].value_counts(), labels = ['ham', 'spam'], autopct = "%0.2f")
plt.show()
```



```
In [58]: x=email['Message']
y=email['Target']
```

```
In [59]: print(x)
```

```
0      Go until jurong point, crazy.. Available only in Sun...
1      Ok lar... Joking wif u oni...
2      Free entry in 2 a wkly comp to win FA Cup final tk...
3      U dun say so early hor... U c already then say...
4      Nah I don't think he goes to usf, he lives aro...
...
5567   This is the 2nd time we have tried 2 contact u... We are sorry the 1st time we could not reach u...
5568   Will i_b going to esplanade fr home?
5569   Pity, * was in mood for that. So...any other s...
5570   The guy did some bitching but I acted like i'd... nothing really happened.
5571   Rofl. Its true to its name
Name: Message, Length: 5169, dtype: object
```

```
In [60]: print(y)
```

```
0      0
1      0
2      1
3      0
4      0
...
5567   1
5568   0
5569   0
5570   0
5571   0
Name: Target, Length: 5169, dtype: int32
```

Training the Model

```
In [61]: 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, random_state=3)
```

```
In [62]: from sklearn.feature_extraction.text import CountVectorizer
from sklearn import svm
```

```
In [63]: cv=CountVectorizer()
```

```
In [64]: x_train_cv = cv.fit_transform(x_train)
x_test_cv = cv.transform(x_test)
```

```
In [65]: print(x_train_cv)
```

```
(0, 1879) 1
(0, 1170) 1
(0, 6840) 1
(0, 6610) 1
(0, 2779) 1
(1, 1939) 1
(1, 4467) 1
(1, 453) 1
(1, 7176) 1
(1, 7594) 1
(1, 1577) 1
(1, 203) 1
(1, 4768) 1
(1, 7175) 1
(1, 7390) 1
(1, 7590) 1
(1, 4309) 1
(1, 5157) 1
(1, 3732) 1
(1, 3015) 1
(1, 2333) 1
(1, 5219) 1
(1, 4577) 1
(1, 4731) 1
(1, 5615) 1
:
(4134, 3290) 2
(4134, 4817) 1
(4134, 1546) 1
(4134, 4195) 1
(4134, 891) 1
(4134, 1992) 1
(4134, 1261) 1
(4134, 7302) 1
(4134, 6595) 1
(4134, 1624) 1
(4134, 1977) 1
(4134, 7438) 1
(4134, 6189) 1
(4134, 6815) 1
(4134, 2357) 1
(4134, 4093) 1
(4134, 6503) 1
(4134, 5934) 1
(4134, 1661) 1
(4134, 5153) 1
(4134, 6292) 1
(4134, 3707) 1
(4134, 6172) 1
(4134, 3024) 1
(4134, 4785) 1
```

Logistic Regression

```
In [66]: from sklearn.linear_model import LogisticRegression
lr=LogisticRegression()
```

```
In [67]: lr.fit(x_train_cv,y_train)
prediction_train=lr.predict(x_train_cv)
```

```
In [68]: from sklearn.metrics import accuracy_score
print(accuracy_score(y_train, prediction_train)*100)
```

```
99.75816203143893
```

```
In [69]: prediction_test = lr.predict(x_test_cv)
```

```
In [70]: from sklearn.metrics import accuracy_score
print(accuracy_score(y_test, prediction_test)*100)
```

```
97.58220502901354
```