

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

```
In [2]: span_df = pd.read_csv("C:\Users\Owner\Downloads\span.csv", encoding='latin1')
span_df

Out [2]:
```

	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 dont think he goes to ust, he lives aro...	NaN	NaN	NaN
...
5567	spam	This is the 2nd time we have tried 2 contact u...	NaN	NaN	NaN
5568	ham	Will b going to explanade fr home?	NaN	NaN	NaN
5569	ham	Pey, * 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...	NaN	NaN	NaN
5571	ham	Rofl. Its true to its name	NaN	NaN	NaN

5572 rows x 5 columns

```
In [3]: span_df.shape
Out [3]: (5572, 5)
```

DATA CLEANING

```
In [4]: span_df.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  32 non-null     object
 4   Unnamed: 4  6 non-null      object
dtypes: object(5)
memory usage: 217.8+ KB

In [5]: # Dropping unnecessary columns
span_df.drop(columns=['Unnamed: 2','Unnamed: 3','Unnamed: 4'],inplace=True)

In [6]: span_df.head()
```

```
Out [6]:
```

	v1	v2
0	ham	Go until jurong point, crazy.. Available only ...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I dont think he goes to ust, he lives aro...

```
In [7]: # Changing column names
span_df.rename(columns={'v1':'Target', 'v2':'Text'},inplace=True)

In [8]: span_df.head()
```

```
Out [8]:
```

	Target	Text
0	ham	Go until jurong point, crazy.. Available only ...
1	ham	Ok lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I dont think he goes to ust, he lives aro...

```
In [9]: # Checking for null values
span_df.isna().sum()

Out [9]:
Target      0
Text        0
dtype: int64

In [10]: # Checking Duplicate values
span_df.duplicated().sum()

Out [10]:
483

In [11]: span_df.describe()
```

```
Out [11]:
```

	Target	Text
count	5572	5572
unique	2	5169
top	ham	Sory, I'll call later
freq	4825	30

VISUALIZE DATASET

```
In [12]: # Let's get the length of the messages
span_df['length'] = span_df['Text'].apply(len)
span_df.head()
```

```
Out [12]:
```

	Target	Text	length
0	ham	Go until jurong point, crazy.. Available only ...	111
1	ham	Ok lar... Joking wif u oni...	29
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	155
3	ham	U dun say so early hor... U c already then say...	49
4	ham	Nah I dont think he goes to ust, he lives aro...	61

```
In [13]: span_df

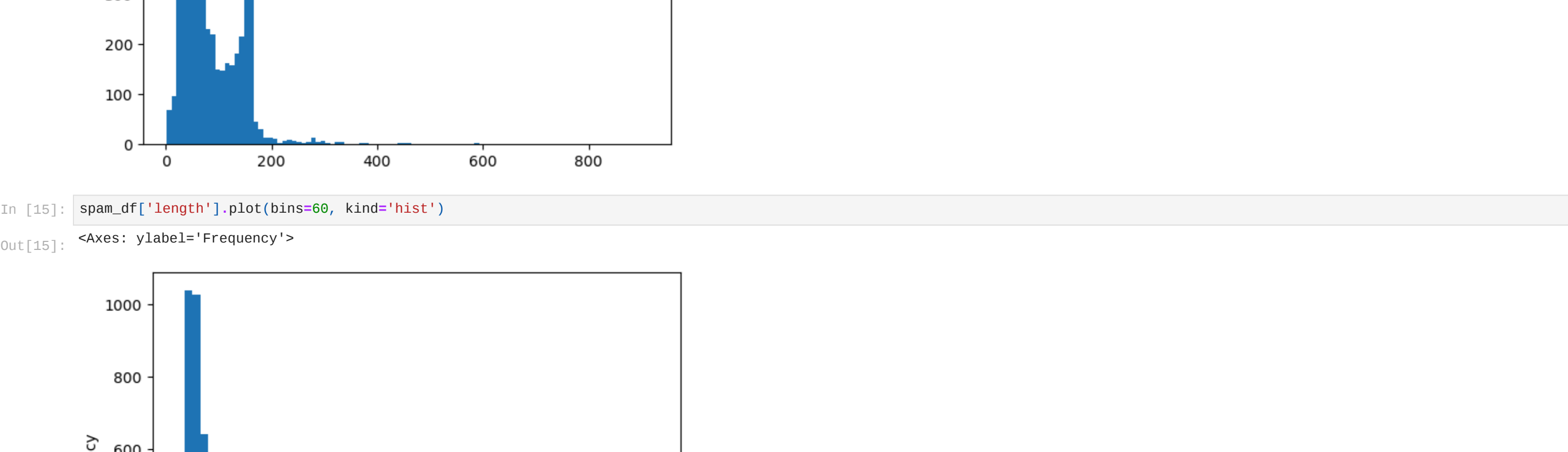
Out [13]:
```

```
Out [13]:
```

	Target	Text	length
0	ham	Go until jurong point, crazy.. Available only ...	111
1	ham	Ok lar... Joking wif u oni...	29
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	155
3	ham	U dun say so early hor... U c already then say...	49
4	ham	Nah I dont think he goes to ust, he lives aro...	61
...
5567	spam	This is the 2nd time we have tried 2 contact u...	161
5568	ham	Will b going to explanade fr home?	37
5569	ham	Pey, * was in mood for that. So...any other s...	57
5570	ham	The guy did some bitching but I acted like i'd...	125
5571	ham	Rofl. Its true to its name	26

5572 rows x 3 columns

```
In [14]: span_df['length'].plot(bins=100, kind='hist')
Out [14]: <Axes: ylabel='Frequency'>
```



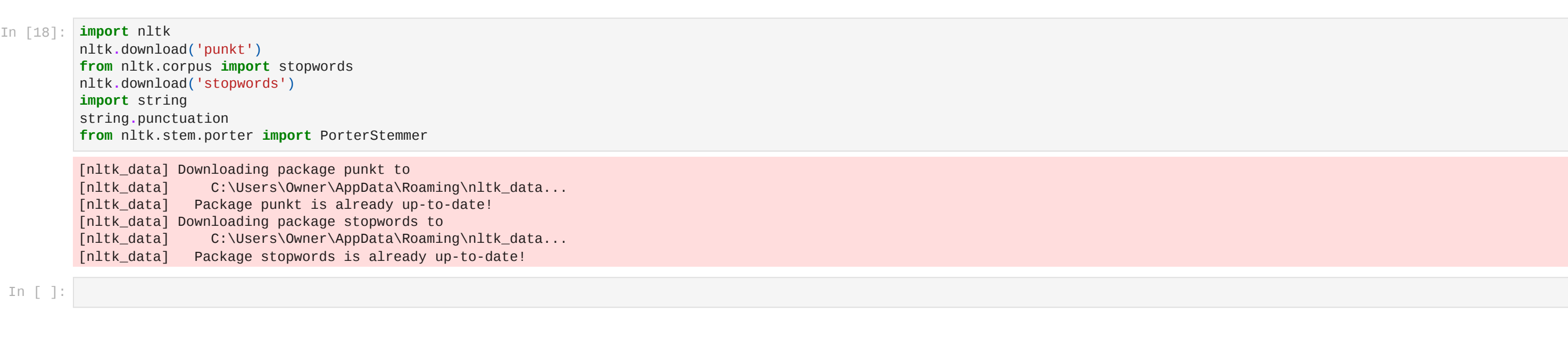
```
In [15]: span_df['length'].plot(bins=60, kind='hist')
Out [15]: <Axes: ylabel='Frequency'>
```



```
In [16]: # Ham and Spam count
span_df['Target'].value_counts()
```

```
Out [16]:
ham      4825
spam      747
Name: Target, dtype: int64
```

```
In [17]: # Check the Target value using pie chart
plt.pie(span_df['Target'].value_counts(), labels=['Ham','Spam'],autopct='%2f')
plt.show()
```



```
In [18]: import nltk
nltk.download('punkt')
from nltk.corpus import stopwords
nltk.download('stopwords')
import string
from nltk.stem.porter import PorterStemmer
```

```
[nltk_data] downloading package punkt to
[nltk_data] c:\Users\Owner\AppData\Roaming\nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] downloading package stopwords to
[nltk_data] c:\Users\Owner\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

```
In [ ]:
```

CREATE TESTING AND TRAINING DATASET/DATA CLEANING

```
In [22]: import string
string.punctuation

Out [22]: '!"#$%&'()*+,-./:;<=>?@[\]^_`{|}~'-'
```

```
In [23]: Test = 'Hello Mr. Future, I am so happy to be learning AI now!!!'
```

```
In [24]: Test_punc_removed = [char for char in Test if char not in string.punctuation]
Test_punc_removed

Out [24]:
```

```
['H',
'e',
'l',
'l',
'o',
',',
'M',
'r',
',',
'F',
'u',
't',
'u',
'r',
'e',
',',
'I',
'm',
'a',
'm',
's',
'o',
'h',
'a',
'p',
'y',
't',
'o',
'b',
'e',
'l',
'e',
'a',
'r',
'n',
'i',
'n',
',',
'g',
'o',
'o',
'A',
'I',
'H',
'e',
'l',
'l',
'o',
'o',
'w']

In [25]: # Join the characters again to form the string.
Test_punc_removed_join = ''.join(Test_punc_removed)
Test_punc_removed_join

Out [25]: 'Hello Mr Future I am so happy to be learning AI now'
```

REMOVE STOPWORDS

```
In [26]: # Download stopwords Package to execute this command
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
stopwords.words('english')
```

```
[nltk_data] downloading package stopwords to
[nltk_data] c:\Users\Owner\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

```
Out [26]: ['I',
'me',
'my',
'myself',
'we',
'our',
'ours',
'ourselves',
'you',
'you're',
'you've',
'you'll',
'you'd',
'your',
'yourself',
'yourselfs',
'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',
'sam',
'is',
'isn',
'was',
'were',
'be',
'been',
'being',
'have',
'has',
'had',
'having',
'do',
'does',
'did',
'doing',
'a',
'an',
'the',
'and',
'but',
'if',
'or',
'because',
'as',
'until',
'while',
'or',
'at',
'by',
'for',
'with',
'about',
'against',
'between',
'into',
'through',
'during',
'before',
'after',
'above',
'below',
'to',
'from',
'up',
'down',
'it',
'out',
'on',
'off',
'over',
'under',
'again',
'further',
'then',
'once',
'here',
'there',
'when',
'where',
'why',
'how',
'all',
'any',
'both',
'each',
'few',
'more',
'less',
'most',
'other',
'some',
'such',
'no',
'not',
'only',
'own',
'same',
'so',
'than',
'too',
'very',
's',
't',
'can',
'will',
'just',
'don',
'don't',
'should',
'shouldn',
'should've',
'now',
'd',
'now',
'o',
'o',
're',
'w',
'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 [27]: Test_punc_removed_join

Out [27]: 'Hello Mr Future I am so happy to be learning AI now'
```

```
In [28]: Test_punc_removed_join_clean = [word for word in Test_punc_removed_join.split() if word.lower() not in stopwords.words('english')]

In [29]: Test_punc_removed_join_clean

Out [29]: 'Hello', 'Mr', 'Future', 'happy', 'learning', 'AI']
```

```
In [30]: # Define a pipeline to clean up all the messages
# The pipeline performs the following: (1) remove punctuation, (2) remove stopwords
```

```
def message_cleaning(message):
    Test_punc_removed = [char for char in message if char not in string.punctuation]
    Test_punc_removed_join = ''.join(Test_punc_removed)
    Test_punc_removed_join_clean = [word for word in Test_punc_removed_join.split() if word.lower() not in stopwords.words('english')]
    return Test_punc_removed_join_clean

In [31]: # Test the newly added function
span_df_clean = span_df['Text'].apply(message_cleaning)

In [32]: print(span_df_clean[0])

Out [32]: ['Go', 'jurong', 'point', 'crazy', 'Available', 'bugis', 'n', 'great', 'world', 'la', 'e', 'buffet', 'cine', 'got', 'amore', 'wat']

In [33]: print(span_df['Text'][0])

Out [33]: Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got amore wat...
```

```
In [34]: from sklearn.feature_extraction.text import CountVectorizer
# Define the cleaning pipeline we defined earlier
vectorizer = CountVectorizer(analyzer = message_cleaning)
spanham_countvectorizer = vectorizer.fit_transform(span_df['Text'])

In [35]: print(spanham_countvectorizer.toarray())

Out [35]: [[0 0 0 ... 0 0 0]
[0 0 0 ... 0 0 0]
[0 0 0 ... 0 0 0]
...
[0 0 0 ... 0 0 0]
[0 0 0 ... 0 0 0]
[0 0 0 ... 0 0 0]]
```

```
In [36]: spanham_countvectorizer.shape
Out [36]: (5572, 11304)
```

TRAINING THE MODEL WITH ALL DATASET

```
In [37]: from sklearn.naive.bayes import MultinomialNB
NB_classifier = MultinomialNB()
label = span_df['Text'].values
NB_classifier.fit(spanham_countvectorizer, label)
```

```
Out [37]: #MultinomialNB
MultinomialNB()
```

```
In [38]: testing_sample = ['Free money!!!!', 'Hi Kin, Please let me know if you need any further information. Thanks']
testing_sample_countvectorizer = vectorizer.transform(testing_sample)
```

```
In [39]: test_predict = NB_classifier.predict(testing_sample_countvectorizer)
test_predict

Out [39]: array(['Sorry, I'll call later',
      'Hi! You just spoke to MANEESHA V. We'd like to know if you were satisfied with the experience. Reply Toll Free with Yes or No'],
      dtype=<U910>)
```

```
In [40]: # Mini Challenge!
testing_sample = ['Hello, I am Ryan, I would like to book a hotel in Bali by January 24th', 'money viagra!!!!!!']
```

```
In [41]: testing_sample = ['money viagra!!!!!!', 'Hello, I am Ryan, I would like to book a hotel in SF by January 24th']
testing_sample_countvectorizer = vectorizer.transform(testing_sample)
test_predict = NB_classifier.predict(testing_sample_countvectorizer)
```

```
Out [41]: array(['Sorry, I'll call later', 'Sorry, I'll call later'], dtype=<U910>)
```

```
In [42]: X = spanham_countvectorizer
y = label
```

```
In [43]: X.shape
Out [43]: (5572, 11304)
```

```
In [44]: y.shape
Out [44]: (5572,)
```

```
In [45]: 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)
```

```
In [46]: from sklearn.naive.bayes import MultinomialNB
NB_classifier = MultinomialNB()
NB_classifier.fit(X_train, y_train)
```

```
Out [46]: #MultinomialNB
MultinomialNB()
```

```
In [ ]:
```