# NATURAL LANGUAGE PROCESSING (NLP)
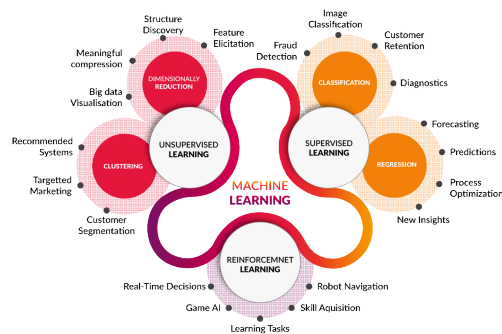
Machine Learning

Dr. Adnan Abid

---

## Natural Language Processing - NLP

Natural-language processing (NLP) is an area of computer science and artificial intelligence concerned with the interactions between computers and human (natural) languages. NLP is used to apply Machine Learning models to text and language.

Example use:

Teach machines to understand what is said in spoken and written word is the focus of Natural Language Processing. Whenever you dictate something into your iPhone / Android device that is then converted to text, that's an NLP algorithm in action.

---



---

## NLP History

The history of natural-language processing generally started in the 1950s, although work can be found from earlier periods. In 1950, Alan Turing published an article titled "Computing Machinery and Intelligence" which proposed what is now called the Turing test as a criterion of intelligence.
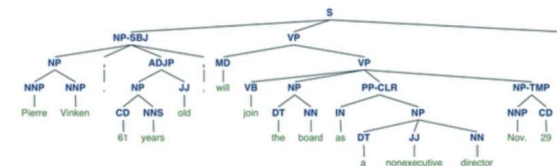
Up to the 1980s, most natural-language processing systems were based on complex sets of hand-written rules. Starting in the late 1980s, however, there was a revolution in natural-language processing with the introduction of machine learning algorithms for language processing.

https://en.wikipedia.org/wiki/Natural-language_processing

## NLP Uses

- Sentiment analysis. Identifying the mood or subjective opinions within large amounts of text, including average sentiment and opinion mining.

- Use it to predict the genre of the book.

- Question Answering

- Use NLP to build a machine translator or a speech recognition system

- Document Summarization
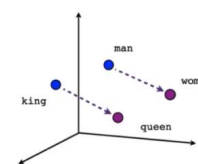
## NATURAL LANGUAGE PROCESSING - NLP



https://www.nltk.org/
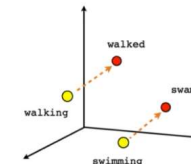
## Natural Language Processing - NLP

Main NLP library examples:

Natural Language Toolkit - NLTK
SpaCy
Stanford NLP
OpenNLP

## Natural Language Processing - NLP
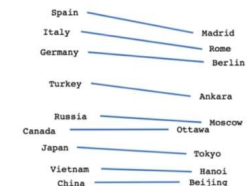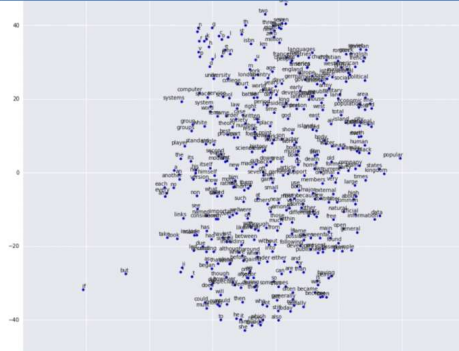


Male-Female     Verb tense     Country-Capital

https://www.tensorflow.org/tutorials/word2vec

## Natural Language Processing - NLP



## Natural Language Processing

In this part, you will understand and learn how to:

1. Clean texts to prepare them for the Machine Learning models,
2. Create a Bag of Words model,
3. Apply Machine Learning models onto this Bag of Worlds model.

## NLP - Bag of Words

Very popular NLP model - It is a model used to preprocess the texts to classify before fitting the classification algorithms on the observations containing the texts.
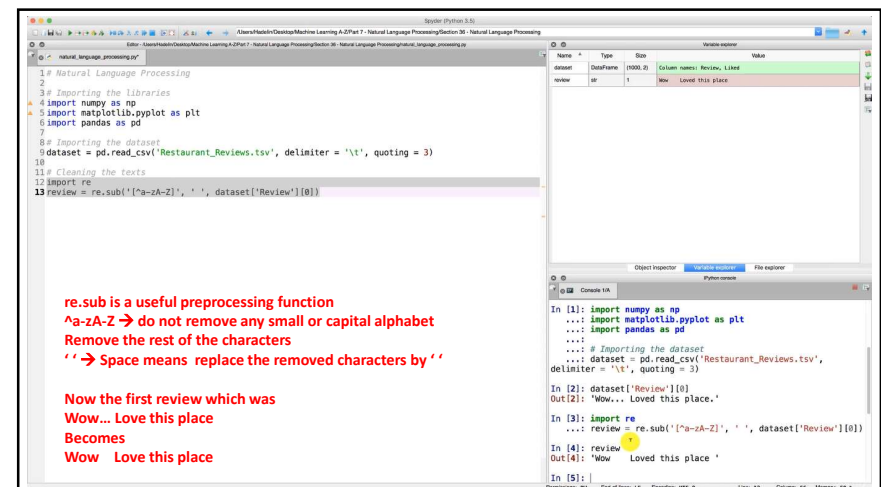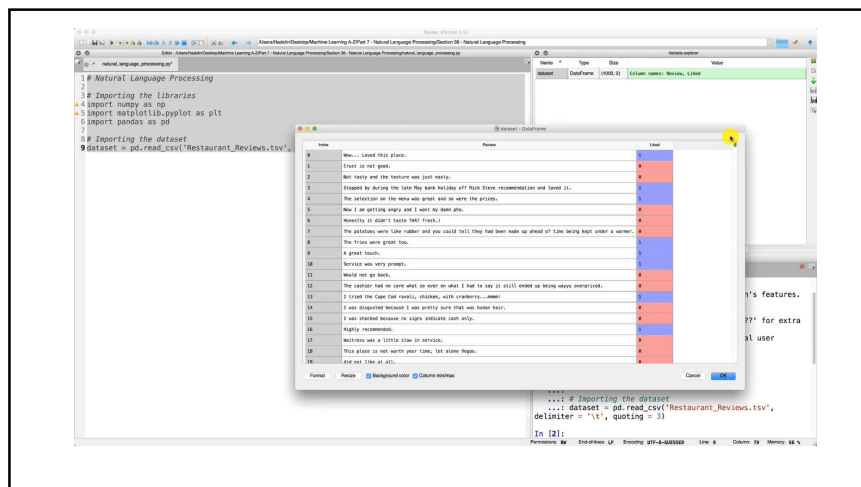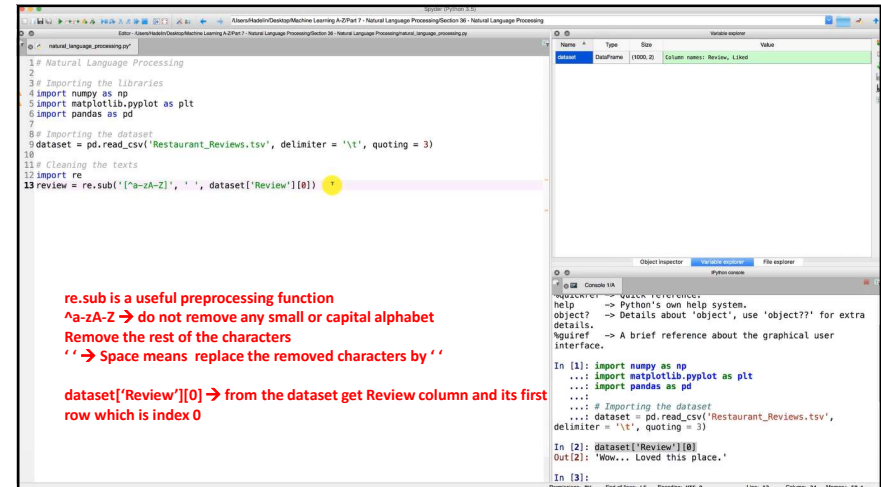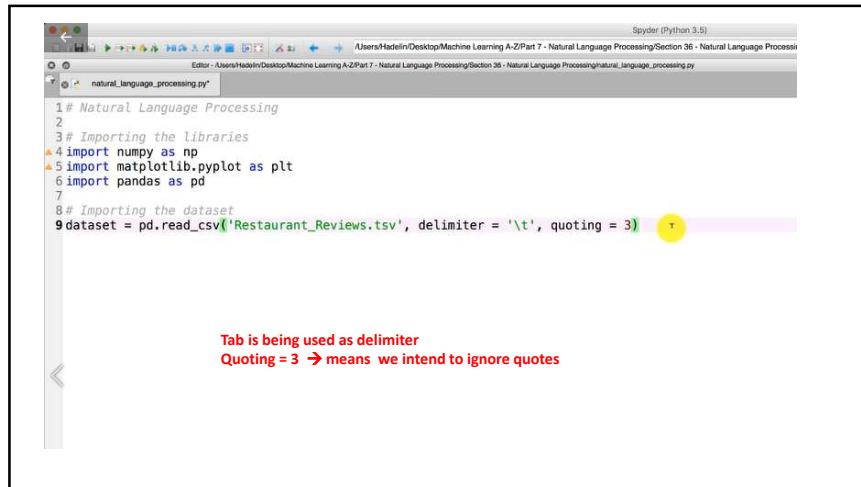
It involves two things:

1. A vocabulary of known words.
2. A measure of the presence of known words.



**CSV OR TSV?**

**Commas may be a part of the review comment... which may be misleading**
**Reviews do not contain any tabs, so TSV is suitable.**

Slide 1 (top-left):

```
# Natural Language Processing

# Importing the libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

# Importing the dataset
dataset = pd.read_csv('Restaurant_Reviews.tsv', delimiter = '\t', quoting = 3)
```

**Tab is being used as delimiter**
**Quoting = 3 → means we intend to ignore quotes**

Slide 2 (top-right):

```
# Natural Language Processing

# Importing the libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

# Importing the dataset
dataset = pd.read_csv('Restaurant_Reviews.tsv', delimiter = '\t', quoting = 3)

# Cleaning the texts
import re
review = re.sub('[^a-zA-Z]', ' ', dataset['Review'][0])
```

**re.sub is a useful preprocessing function**
**^a-zA-Z → do not remove any small or capital alphabet**
**Remove the rest of the characters**
**' ' → Space means replace the removed characters by ' '**

**dataset['Review'][0] → from the dataset get Review column and its first row which is index 0**

Slide 4 (bottom-right):

**re.sub is a useful preprocessing function**
**^a-zA-Z → do not remove any small or capital alphabet**
**Remove the rest of the characters**
**' ' → Space means replace the removed characters by ' '**

**Now the first review which was**
**Wow... Love this place**
**Becomes**
**Wow    Love this place**

## Slide 1 (top-left)

```python
1 # Natural Language Processing
2
3 # Importing the libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 # Importing the dataset
9 dataset = pd.read_csv('Restaurant_Reviews.tsv', delimiter = '\t', quoting = 3)
10
11 # Cleaning the texts
12 import re
13 review = re.sub('[^a-zA-Z]', ' ', dataset['Review'][0])
14 review = review.lower()
```

**re.sub is a useful preprocessing function**
**^a-zA-Z → do not remove any small or capital alphabet**
**Remove the rest of the characters**
**' ' → Space means replace the removed characters by ' '**

**Now the first review which was**
**Wow… Love this place**
**Becomes**
**Wow   Love this place**
**----------**
**Now convert all reviews into lower case…**
**First review becomes**
**wow   love this place**

Console:
```
...: # Importing the dataset
...: dataset = pd.read_csv('Restaurant_Reviews.tsv',
delimiter = '\t', quoting = 3)

In [2]: dataset['Review'][0]
Out[2]: 'Wow... Loved this place.'

In [3]: import re
...: review = re.sub('[^a-zA-Z]', ' ', dataset['Review'][0])

In [4]: review
Out[4]: 'Wow   Loved this place '

In [5]: review = review.lower()

In [6]: review
Out[6]: 'wow   loved this place '

In [7]:
```

## Slide 2 (top-right)

```python
1 # Natural Language Processing
2
3 # Importing the libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 # Importing the dataset
9 dataset = pd.read_csv('Restaurant_Reviews.tsv', delimiter = '\t', quoting = 3)
10
11 # Cleaning the texts
12 import re
13 import nltk
14 nltk.download('stopwords')
15 from nltk.corpus import stopwords
16 from nltk.stem.porter import PorterStemmer
17 review = re.sub('[^a-zA-Z]', ' ', dataset['Review'][0])
18 review = review.lower()
19 review = review.split()
20 ps = PorterStemmer()
21 review = [ps.stem(word) for word in review if not word in set(stopwords.words('english'))]
```

**The next step is to apply Stemming**
**Convert the words into their root form**
**Computer → compute**
**Loved → Love**
**Talking → Talk**
**It reduces the overall vocabulary and helps reducing the size of the sparse matrix.**
**Porter Stemmer is one of widely used stemmer and nltk uses it too.**
**Apply stemming after removing the stopwords so that the remain words are stemmed.**
**Now the comment [wow loved place] becomes [wow love place]**

Console:
```
...: # Cleaning the texts
...: import re
...: import nltk
...: nltk.download('stopwords')
...: from nltk.corpus import stopwords
...: from nltk.stem.porter import PorterStemmer
...: review = re.sub('[^a-zA-Z]', ' ', dataset['Review'][0])
...: review = review.lower()
[nltk_data] Downloading package stopwords to
[nltk_data]     /Users/Hadelin/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!

In [2]: review = review.split()
...: ps = PorterStemmer()
...: review = [ps.stem(word) for word in review if not word
in set(stopwords.words('english'))]

In [3]:
```

## Slide 3 (bottom-left)

```python
1 # Natural Language Processing
2
3 # Importing the libraries
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 # Importing the dataset
9 dataset = pd.read_csv('Restaurant_Reviews.tsv', delimiter = '\t', quoting = 3)
10
11 # Cleaning the texts
12 import re
13 import nltk
14 nltk.download('stopwords')
15 from nltk.corpus import stopwords
16 review = re.sub('[^a-zA-Z]', ' ', dataset['Review'][0])
17 review = review.lower()
18 review = review.split()
19 review = [word for word in review if not word in set(stopwords.words('english'))]
```
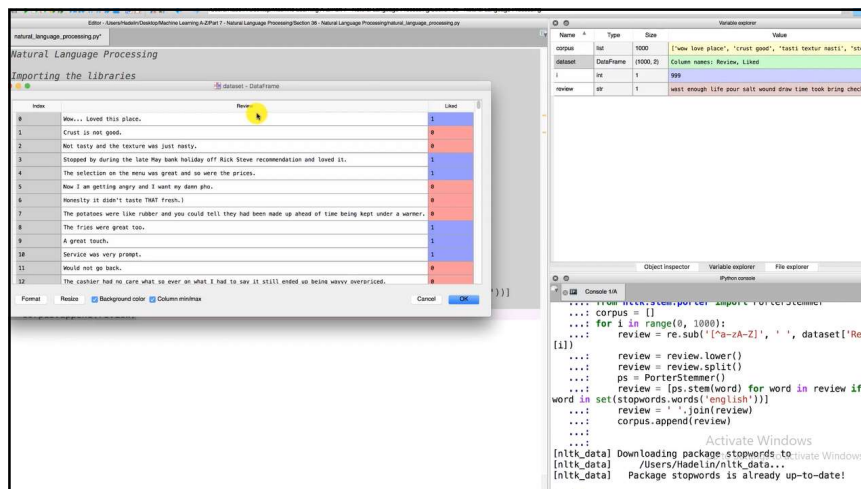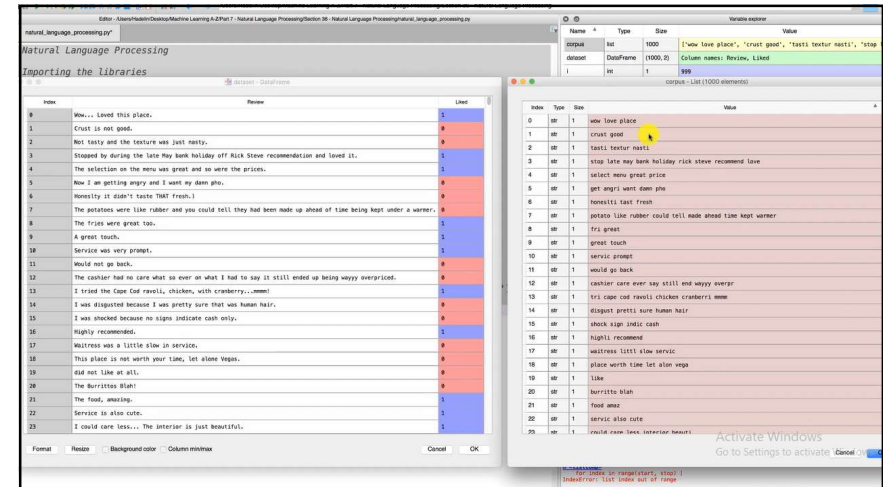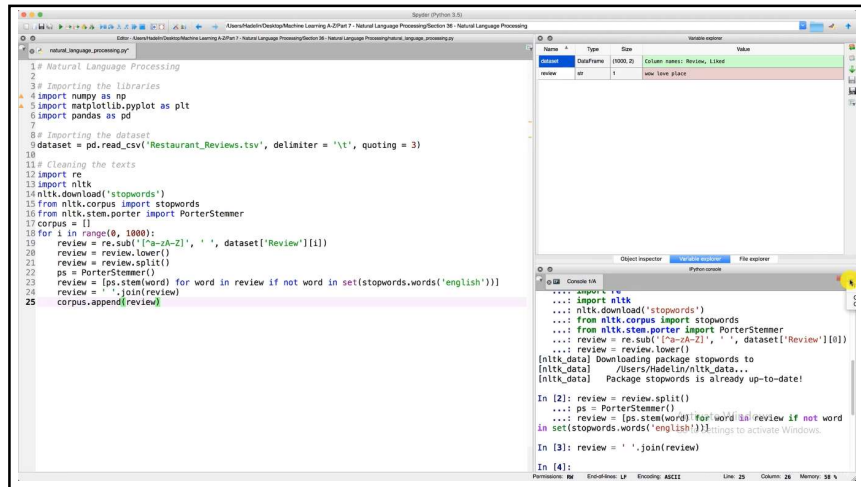
**Now we need to remove STOP WORDS**
**These are the words which do not have any meaning towards the positivity or negativity of the comment, but they are there to complete the structure of the sentences e.g. this, that, of, on etc.**

**NLTK provides a list of stop words.**
**Download that list and import stopwords package from nltk.corpus**

**Apply a loop on all the words of the review comment, which have been split, and remove the ones which are STOP WORDS.**

**Use set function to check for stopwords as it works faster**

Console:
```
In [6]: review
Out[6]: 'wow   loved this place '

In [7]: import nltk
...: nltk.download('stopwords')
[nltk_data] Downloading package stopwords to
[nltk_data]     /Users/Hadelin/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
Out[7]: True

In [8]: review = review.split()

In [9]: from nltk.corpus import stopwords

In [10]: review = [word for word in review if not word in
set(stopwords.words('english'))]

In [11]:
```

## Slide 4 (bottom-right)

```python
# Natural Language Processing

# Importing the libraries
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

# Importing the dataset
dataset = pd.read_csv('Restaurant_Reviews.tsv', delimiter = '\t', quoting = 3)

# Cleaning the texts
import re
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
review = re.sub('[^a-zA-Z]', ' ', dataset['Review'][0])
review = review.lower()
review = review.split()
ps = PorterStemmer()
review = [ps.stem(word) for word in review if not word in set(stopwords.words('english'))]
review = ' '.join(review)
```

**The next step is to join all the words again to convert the preprocessed review into a single string by concatenating them with a ' '.**

Console:
```
...: import nltk
...: nltk.download('stopwords')
...: from nltk.corpus import stopwords
...: from nltk.stem.porter import PorterStemmer
...: review = re.sub('[^a-zA-Z]', ' ', dataset['Review'][0])
...: review = review.lower()
[nltk_data] Downloading package stopwords to
[nltk_data]     /Users/Hadelin/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!

In [2]: review = review.split()
...: ps = PorterStemmer()
...: review = [ps.stem(word) for word in review if not w
in set(stopwords.words('english'))]

In [3]: review = ' '.join(review)
```

Variable explorer value: wow love place

## Text Cleaning

- Load Text
- CSV vs TSV
- Convert into lower case
- Split the review comment/doc into words
- Remove stop words
- Stemming
- Combining all reviews/documents

**Top-left slide:**

toarray() converts the corpus into matrix

**The BAG OF WORDS model**
**Sparse matrix with one row for each review/document**

**\*\*\*The CountVectorizer class can perform many of the preprocessing steps but it is better to explicitly perform those steps (A) for teaching purpose (B) we have the liberty to control many things e.g. in case of web scrapping the HTML tags also contain a-zA-Z letters...**

**Top-right slide:**

Set y to be the last column of the dataset.
Our data set has 2 columns 0 and 1.
Thus y = dataset.iloc[:,1] → all rows and 2nd column

**Bottom-left slide:**

Max_features = 1500 removes the words appearing in very few reviews. E.g. any proper nouns person's name etc.
It is kind of dimensional reduction
There are other methods too
Principal Component Analysis
Singular Value Decomposition

**Bottom-right slide:**

Set y to be the last column of the dataset.
Our data set has 2 columns 0 and 1.
Thus y = dataset.iloc[:,1] → all rows and 2nd column
Now, it has become a conventional ML problem

# NLP Classification

- Text Cleaning
- Bag of Words or Vectorization
- Apply Classification Models
- Compare Models
  - Confusion Matrix
  - Accuracy