## 5.4 Natural Language Processing with Python

Overview
In this tutorial you will be using a small dataset consisting of customer reviews of Video games that are available in the Amazon web store.
- Your task is to identify named entities occurring in Video Game reviews and to analyse the sentiment of the reviews
- The dataset is in a tab separated format where each line corresponds to a different Video game review
- Moreover, the tab separated file consists of three columns:
    - The first column corresponds to the review id
    - the second column to the movie id
    - the third column to the actual movie review
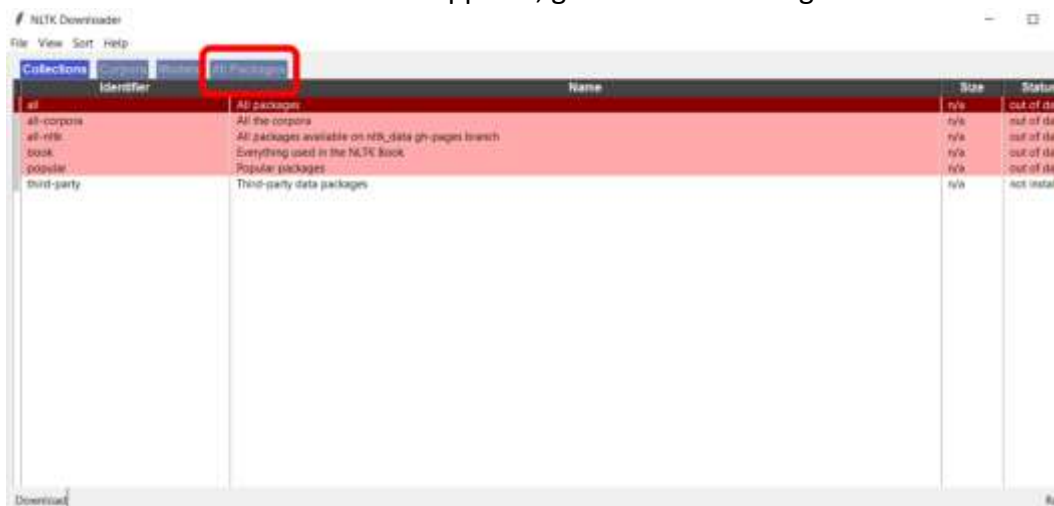
Topics Covered:
- Installing NER in NLTK
- Loading dataset using panda
- Sentiment analysis
- Named entity recognition

**Activity 1**: Installing NER in NLTK

- NLTK offers a machine learning-based named entity recogniser (NER).
- We will firstly need to install the NER module within our NLTK
- For this, create a new Python3 file in your Jupyter and enter the commands as shown below:

```
In [ ]: import nltk
        nltk.download()
```

- In the new window that it appears, go to the 'All Packages' tab

- And then download the following two packages:   1. 'maxent_ne_chunker'   2. 'words'





**Activity 2**: Loading dataset using Pandas

- We will use the read_csv function of the Pandas library to load the Video games dataset into a Pandas DataFrame
- Note that the dataset is in a tab separated format so we will need to specify that the delimiter is the tab character

- Also the dataset contains no headers so we need to set the Header argument to None

```
In [3]:   1  import pandas as pd
          2  # specify delimiter of csv file to be the tab character - '\t\
          3  # csv file contains no headers so set header=Non
          4  df = pd.read_csv('downloads/Video_games_reviews.csv',
          5                    delimiter='\t',
          6                    header=None)
          7  df
```

Out[3]:

|   | 0 | 1 | 2 |
|---|---|---|---|
| 0 | 1 | B00000JNHJ | Great game, I dont really think its a little k... |
| 1 | 2 | B000035XYC | you create things.That's about it.You also fig... |
| 2 | 3 | B0000296ZD | The improvements of the graphics are very good... |
| 3 | 4 | B00009WDLD | The graphics in this game aren't all that much... |
| 4 | 5 | B00009TW6R | This is a horrible game, average graphics, bad... |
| 5 | 6 | B00000K4MC | I never thought this game looked very good unt... |

**2.1** Extracting the Review Column from the Data Frame
- The DataFrame consists of three columns corresponding to the three columns of the Video game dataset
- However, we are only interested in analysing the third column which is the actual review of a Video game
- It is thus useful to assign the third column for the Dataframe into a separate variable

```
In [4]:   1  import pandas as pd
          2
          3  df = pd.read_csv('downloads/Video_games_reviews.csv',
          4                    delimiter='\t',
          5                    header=None)
          6
          7  # assign the third column of the Dataframe, i.e., the actual reviews,
          8  # into a separate variable
          9  video_review_texts = df[2]
         10  video_review_texts
```

```
Out[4]: 0    Great game, I dont really think its a little k...
        1    you create things.That's about it.You also fig...
        2    The improvements of the graphics are very good...
        3    The graphics in this game aren't all that much...
        4    This is a horrible game, average graphics, bad...
        5    I never thought this game looked very good unt...
        6    The original Grand Theft Auto was a classic fo...
        7    Sonic Adventure DX is supposed to be a nearly ...
        8    I have played many many many video games in my...
```

**Activity 3**: Sentiment Analysis

- We are now ready to analyse the sentiment of the Video game reviews using the textblob library
- For this we will iterate through the reviews in the video_review_texts and for each review we will call textblob to analyse the sentiment of that review
- Textblob will assign a sentiment score to each sentence of a review
- You will need to think how you can aggregate the sentence scores to derive an overall sentiment score for each review

```
 1  import pandas as pd
 2  df = pd.read_csv('downloads/Video_games_reviews.csv',
 3                    delimiter='\t',
 4                    header=None)
 5  video_review_texts = df[2]
 6  # import textblobk
 7  from textblob import TextBlob
 8  '''
 9  Iterate through the video reviews and for each review
10  call textblob to analyse the sentiment of that review.
11  '''
12  for index, review_text in enumerate(video_review_texts):
13      # call text blob for current review_text
14      blob = TextBlob(review_text)
15      print('Analysing review\t', review_text)
16      # iterate through the sentences of the current review_text
17      # and print the sentiment scores of each sentence
18      for sentence in blob.sentences:
19          print('-------SENTIMENT OF SENTENCE--------')
20          print(sentence, '\t', sentence.sentiment.polarity)
21          print('-------END--------')
22
```

- After aggregating the sentiment scores of the constituent sentences of a review you can use a simple threshold to decide whether a review is positive or negative (e.g., if aggregate_score > 0 then review is positive else review is negative)
- Once you have decided whether a review is positive or negative you can store the sentiment classification labels on a separate column of your DataFrame
- The code below randomly decides whether a review is positive/negative and it then stores the sentiment classification labels in the dataframe

```
 1  import pandas as pd
 2  df = pd.read_csv('downloads/Video_games_reviews.csv',
 3                    delimiter='\t',
 4                    header=None)
 5  video_review_texts = df[2]
 6
 7  from textblob import TextBlob
 8  import random
 9
10  # list that stores the sentiment classification labels
11  sentiment_classification_labels = []
12  for index, review_text in enumerate(video_review_texts):
13      blob = TextBlob(review_text)
14  #     for sentence in blob.sentences:
15  #         print(sentence, '\t', sentence.sentiment.polarity)
16      # Here we randomly assign a sentiment classification label
17      # to each review but for CW2 you will need to firstly agggregate
18      # the sentence scores and then decide whether the review is positive/negative
19      # based on the aggregated score
20      sentiment_label_for_current_review = random.randint(0,1)
21      sentiment_classification_labels.append(sentiment_label_for_current_review)
22
23  # Finally append the list of sentiment_classification_labels into the DataFrame
24  df['Sentiment_Classification_Labels'] = sentiment_classification_labels
25
```

| | 0 | 1 | 2 | Sentiment_Classification_Labels |
|---|---|---|---|---|
| 0 | 1 | B00000JNHJ | Great game, I dont really think its a little k... | 0 |
| 1 | 2 | B000035XYC | you create things.That's about it.You also fig... | 1 |
| 2 | 3 | B0000296ZD | The improvements of the graphics are very good... | 1 |
| 3 | 4 | B00009WDLD | The graphics in this game aren't all that much... | 1 |
| 4 | 5 | B00009TW6R | This is a horrible game, average graphics, bad... | 0 |
| 5 | 6 | B00000K4MC | I never thought this game looked very good unt... | 0 |
| 6 | 7 | B00000DMAV | The original Grand Theft Auto was a classic fo... | 1 |
| 7 | 8 | B00008URUB | Sonic Adventure DX is supposed to be a nearly ... | 1 |
| 8 | 9 | B00000DMB3 | I have played many many many video games in my... | 0 |
| 9 | 10 | B00002CF9M | Lord Of TerrorWhen Diablo came out almost 4 ye... | 0 |
| 10 | 11 | B00000DMAT | Goldeneye lookes like just another shootem up ... | 1 |
| 11 | 12 | B00005NZ1G | What is the big deal about Halo. It is one of ... | 1 |
| 12 | 13 | B0000011BE | Now if you are a fan or those 2-D games and st... | 1 |

- The sentiment classification labels are now stored in the fourth column of our DataFrame
- We can apply standard Dataframe operations to retrieve only the positive or only the negative reviews

```
23
24 df['Sentiment_Classification_Labels'] = sentiment_classification_labels
25 # get only negative reviews
26 df[df.Sentiment_Classification_Labels==0]
27
28 # get only
29
```

Out[22]:

| | 0 | 1 | 2 | Sentiment_Classification_Labels |
|---|---|---|---|---|
| 3 | 4 | B00009WDLD | The graphics in this game aren't all that much... | 0 |
| 4 | 5 | B00009TW6R | This is a horrible game, average graphics, bad... | 0 |
| 6 | 7 | B00000DMAV | The original Grand Theft Auto was a classic fo... | 0 |
| 8 | 9 | B00000DMB3 | I have played many many many video games in my... | 0 |
| 9 | 10 | B00002CF9M | Lord Of TerrorWhen Diablo came out almost 4 ye... | 0 |

```
24 df['Sentiment_Classification_Labels'] = sentiment_classification_labels
25 # get only positive reviews
26 df[df.Sentiment_Classification_Labels==1]
27
28 # get only
29
```

Out[24]:

| | 0 | 1 | 2 | Sentiment_Classification_Labels |
|---|---|---|---|---|
| 0 | 1 | B00000JNHJ | Great game, I dont really think its a little k... | 1 |
| 1 | 2 | B000035XYC | you create things.That's about it.You also fig... | 1 |
| 3 | 4 | B00009WDLD | The graphics in this game aren't all that much... | 1 |
| 9 | 10 | B00002CF9M | Lord Of TerrorWhen Diablo came out almost 4 ye... | 1 |
| 12 | 13 | B0000011BE | Now if you are a fan or those 2-D games and st... | 1 |

**Activity 4**: Named Entity Recognition

- Firstly, import the NLTK libraries which will be used for NER

5

```
1  # NLTK libraries for NER
2  from nltk import word_tokenize, pos_tag, ne_chunk
3  from nltk.tree import Tree
4
```

- In NLTK we can perform tokenisation, POS tagging and NER in only one line of code
- Note that tokenisation and POS tagging are necessary pre-processing steps for NER

```
5  input_sentence = 'I study at Edgehill University which is located in Ormskirk'
6  # perform tokenisation, pos tagging and ner in only one line of code
7  chunks = ne_chunk(pos_tag(word_tokenize(input_sentence)))
```

- The ne_chunk function of NLTK will return a list of words or phrases occurring in the input sentence together with their POS tag and their named entity label
- The NLTK NER identifies the following named entities:
    o Organisation, Person, Location, Date, Time, Money, Percent, Facility and GPE (geo-political entity)
- In order to print the identified named entities that occur in the input sentence, simply iterate through the list of chunks and print only the chunks that contain a named entity (ignore chunks that contain only a POS tag)

```
10  # iterate through chunks
11  for ne in chunks:
12      # if current chunk is a named entity type
13      # then print it
14      if type(ne) == Tree:
15          print(ne)

(ORGANIZATION Edgehill/NNP University/NNP)
(GPE Ormskirk/NNP)
```

Exercise: Apply NER to the Video Game dataset