```
1 !pip install nltk
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: nltk in /usr/local/lib/python3.7/dist-packages (3.7)
Requirement already satisfied: click in /usr/local/lib/python3.7/dist-packages (from nltk) (7.1.2)
Requirement already satisfied: tqdm in /usr/local/lib/python3.7/dist-packages (from nltk) (4.64.0)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.7/dist-packages (from nltk) (2022.6.2)
Requirement already satisfied: joblib in /usr/local/lib/python3.7/dist-packages (from nltk) (1.1.0)
```

## ▾ - Exercici 1

Agafa un text en anglès que vulguis, i calcula'n la freqüència de les paraules.

```
1 #Loading NLTK
2 import nltk
```

```
1 nltk.download('punkt')
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
True
```

El texto sale de aquí:

https://www.nltk.org/

```
 1 from nltk.tokenize import sent_tokenize
 2
 3 text="""Natural Language Toolkit
 4 NLTK is a leading platform for building Python programs to work with human language data. It provides ea
 5
 6 Thanks to a hands-on guide introducing programming fundamentals alongside topics in computational lingui
 7
 8 NLTK has been called "a wonderful tool for teaching, and working in, computational linguistics using Pyt
 9
10 Natural Language Processing with Python provides a practical introduction to programming for language pr
11
12 tokenized_text=sent_tokenize(text)
13 print(tokenized_text)
```

```
['Natural Language Toolkit\nNLTK is a leading platform for building Python programs to work with human language data.', 'It
```

```
1 from nltk.tokenize import word_tokenize
2
3 tokenized_word=word_tokenize(text)
4 print(tokenized_word)
```

```
['Natural', 'Language', 'Toolkit', 'NLTK', 'is', 'a', 'leading', 'platform', 'for', 'building', 'Python', 'programs', 'to',
```
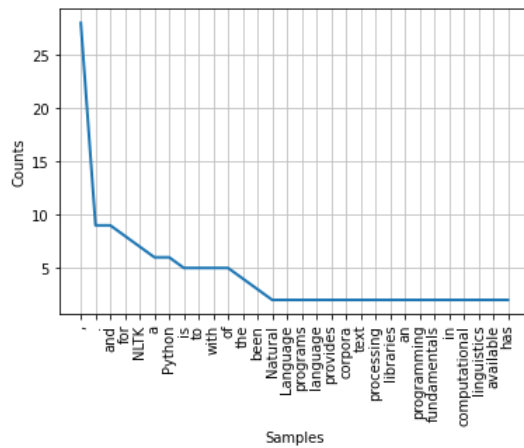
```
1 from nltk.probability import FreqDist
2 fdist = FreqDist(tokenized_word)
3 print(fdist)
```

```
<FreqDist with 137 samples and 247 outcomes>
```

```
1 fdist.most_common(5)
```

```
[(',', 28), ('.', 9), ('and', 9), ('for', 8), ('NLTK', 7)]
```

```
1 # Frequency Distribution Plot
2 import matplotlib.pyplot as plt
3 fdist.plot(30,cumulative=False)
4 plt.show()
```

## ▾ - Exercici 2

Treu les stopwords i realitza stemming al teu conjunt de dades.

```
1 nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
True
```

```
1 from nltk.corpus import stopwords
2 #nltk.download('stopwords')
3
4 stop_words=set(stopwords.words("english"))
5 print(stop_words)
6 print(len(stop_words))
```

```
{'is', "don't", 'they', 'o', 'again', 'ourselves', "shan't", 'how', "you'd", "wasn't", 'while', 'didn', 'hadn', 'you', 'was'
179
```

tokenized_text=sent_tokenize(text)

```
1 tokenized_sent = tokenized_word
2
3 filtered_sent=[]
4 for w in tokenized_sent:
5     #print(w)
6     if w not in stop_words:
7         filtered_sent.append(w)
8 print("Tokenized Sentence:",tokenized_sent)
9 print("Filterd Sentence:",filtered_sent)
```

```
Tokenized Sentence: ['Natural', 'Language', 'Toolkit', 'NLTK', 'is', 'a', 'leading', 'platform', 'for', 'building', 'Python'
Filterd Sentence: ['Natural', 'Language', 'Toolkit', 'NLTK', 'leading', 'platform', 'building', 'Python', 'programs', 'work'
```

```
1 # Stemming
2 from nltk.stem import PorterStemmer
3 from nltk.tokenize import sent_tokenize, word_tokenize
4
5 ps = PorterStemmer()
6
7 stemmed_words=[]
8 for w in filtered_sent:
9     stemmed_words.append(ps.stem(w))
10
11 print("Filtered Sentence:",filtered_sent)
12 print("Stemmed Sentence:",stemmed_words)
```

```
Filtered Sentence: ['Natural', 'Language', 'Toolkit', 'NLTK', 'leading', 'platform', 'building', 'Python', 'programs', 'work'
Stemmed Sentence: ['natur', 'languag', 'toolkit', 'nltk', 'lead', 'platform', 'build', 'python', 'program', 'work', 'human',
```

```
1 nltk.download('wordnet')
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
True
```

```
1 nltk.download('omw-1.4')
```

```
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
True
```

```python
 1 #Lexicon Normalization
 2 #performing stemming and Lemmatization
 3
 4 from nltk.stem.wordnet import WordNetLemmatizer
 5 lem = WordNetLemmatizer()
 6
 7 from nltk.stem.porter import PorterStemmer
 8 stem = PorterStemmer()
 9
10 word = "fly"
11 print("Lemmatized Word:",lem.lemmatize(word,"v"))
12 print("Stemmed Word:",stem.stem(word))
```

```
Lemmatized Word: fly
Stemmed Word: fli
```

## ▾ - Exercici 2

Treu les stopwords i realitza stemming al teu conjunt de dades.

```python
1 sent = "Albert Einstein was born in Ulm, Germany in 1879."
2
3 tokens=nltk.word_tokenize(sent)
4 print(tokens)
```

```
['Albert', 'Einstein', 'was', 'born', 'in', 'Ulm', ',', 'Germany', 'in', '1879', '.']
```
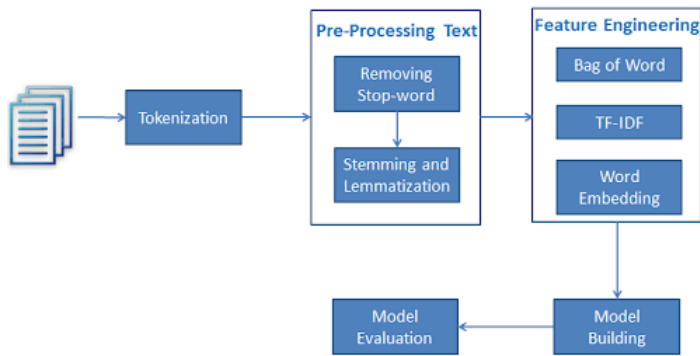
```
1 nltk.download('averaged_perceptron_tagger')
```

```
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]     /root/nltk_data...
[nltk_data]   Unzipping taggers/averaged_perceptron_tagger.zip.
True
```

```
1 nltk.pos_tag(tokens)
```

```
[('Albert', 'NNP'),
 ('Einstein', 'NNP'),
 ('was', 'VBD'),
 ('born', 'VBN'),
 ('in', 'IN'),
 ('Ulm', 'NNP'),
 (',', ','),
 ('Germany', 'NNP'),
 ('in', 'IN'),
 ('1879', 'CD'),
 ('.', '.')]
```

```
1 # Activo Google Drive
2
3 from google.colab import drive
4 drive.mount('/content/drive')
```

    Mounted at /content/drive

```
1 # Import pandas
2 import pandas as pd
3
4
5 data=pd.read_csv('/content/drive/MyDrive/01_COLAB/train.tsv', sep='\t')
```

```
1 data.head()
```

|   | PhraseId | SentenceId | Phrase | Sentiment |
|---|----------|------------|--------|-----------|
| **0** | 1 | 1 | A series of escapades demonstrating the adage ... | 1 |
| **1** | 2 | 1 | A series of escapades demonstrating the adage ... | 2 |
| **2** | 3 | 1 | A series | 2 |
| **3** | 4 | 1 | A | 2 |
| **4** | 5 | 1 | series | 2 |

```
1 data.iloc[0][2]
```

    'A series of escapades demonstrating the adage that what is good for the goose is also good for the gander , some of which
    occasionally amuses but none of which amounts to much of a story .  '

```
1 data.info()
```

    <class 'pandas.core.frame.DataFrame'>
    RangeIndex: 156060 entries, 0 to 156059
    Data columns (total 4 columns):
     #   Column      Non-Null Count   Dtype
    ---  ------      --------------   -----
     0   PhraseId    156060 non-null  int64
     1   SentenceId  156060 non-null  int64
     2   Phrase      156060 non-null  object
     3   Sentiment   156060 non-null  int64
    dtypes: int64(3), object(1)
    memory usage: 4.8+ MB

```
1 data.Sentiment.value_counts()
```

    2    79582
    3    32927
    1    27273
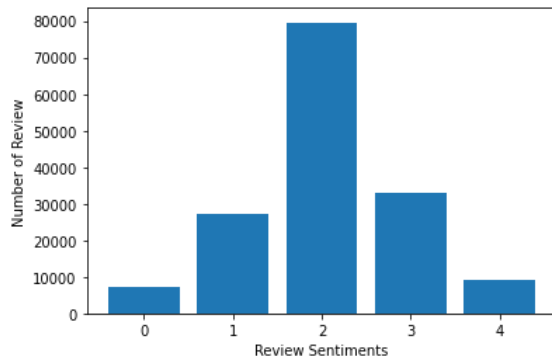    4     9206
    0     7072
    Name: Sentiment, dtype: int64

```
1 Sentiment_count=data.groupby('Sentiment').count()
2 plt.bar(Sentiment_count.index.values, Sentiment_count['Phrase'])
3 plt.xlabel('Review Sentiments')
```

```
4 plt.ylabel('Number of Review')
5 plt.show()
```



```
1 from sklearn.feature_extraction.text import CountVectorizer
2 from nltk.tokenize import RegexpTokenizer
3 #tokenizer to remove unwanted elements from out data like symbols and numbers
4 token = RegexpTokenizer(r'[a-zA-Z0-9]+')
5 cv = CountVectorizer(lowercase=True,stop_words='english',ngram_range = (1,1),tokenizer = token.tokenize)
6 text_counts= cv.fit_transform(data['Phrase'])
```

```
1 print(text_counts)
```

```
  (0, 11671)    1
  (0, 4517)     1
  (0, 3444)     1
  (0, 294)      1
  (0, 5735)     2
  (0, 5751)     1
  (0, 5512)     1
  (0, 9065)     1
  (0, 593)      1
  (0, 584)      1
  (0, 12673)    1
  (1, 11671)    1
  (1, 4517)     1
  (1, 3444)     1
  (1, 294)      1
  (1, 5735)     1
  (1, 5751)     1
  (2, 11671)    1
  (4, 11671)    1
  (5, 4517)     1
  (5, 3444)     1
  (5, 294)      1
  (5, 5735)     1
  (5, 5751)     1
  (7, 4517)     1
  :       :
  (156050, 11305)      1
  (156050, 9054)       1
  (156051, 11305)      1
  (156051, 9054)       1
  (156052, 11305)      1
  (156053, 11281)      1
  (156053, 1281)       1
  (156053, 5252)       1
  (156053, 6156)       1
  (156053, 1006)       1
  (156053, 2271)       1
  (156054, 11281)      1
  (156054, 5252)       1
  (156054, 6156)       1
  (156054, 1006)       1
  (156054, 2271)       1
  (156055, 11281)      1
  (156055, 6156)       1
  (156056, 5252)       1
  (156056, 1006)       1
  (156056, 2271)       1
  (156057, 1006)       1
  (156057, 2271)       1
  (156058, 1006)       1
  (156059, 2271)       1
```

```python
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(
3     text_counts, data['Sentiment'], test_size=0.3, random_state=1)
```

```python
1 from sklearn.naive_bayes import MultinomialNB
2 #Import scikit-learn metrics module for accuracy calculation
3 from sklearn import metrics
4 # Model Generation Using Multinomial Naive Bayes
5 clf = MultinomialNB().fit(X_train, y_train)
6 predicted= clf.predict(X_test)
7 print("MultinomialNB Accuracy:",metrics.accuracy_score(y_test, predicted))
```

```
MultinomialNB Accuracy: 0.6049169122986885
```

```python
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 tf=TfidfVectorizer()
3 text_tf= tf.fit_transform(data['Phrase'])
```

```python
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(
3     text_tf, data['Sentiment'], test_size=0.3, random_state=123)
```

```python
1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn import metrics
3 # Model Generation Using Multinomial Naive Bayes
4 clf = MultinomialNB().fit(X_train, y_train)
5 predicted= clf.predict(X_test)
6 print("MultinomialNB Accuracy:",metrics.accuracy_score(y_test, predicted))
```

```
MultinomialNB Accuracy: 0.5865265496176684
```

```python
1 import nltk
2 sentence = """At eight o'clock on Thursday morning Arthur didn't feel very good."""
3
4 tokens = nltk.word_tokenize(sentence)
5 tokens
6 tagged = nltk.pos_tag(tokens)
7 tagged[0:6]
```

```
[('At', 'IN'),
 ('eight', 'CD'),
 ("o'clock", 'NN'),
 ('on', 'IN'),
 ('Thursday', 'NNP'),
 ('morning', 'NN')]
```

```python
1 nltk.download('maxent_ne_chunker')
2 nltk.download('words')
```

```
[nltk_data] Downloading package maxent_ne_chunker to
[nltk_data]     /root/nltk_data...
[nltk_data]   Package maxent_ne_chunker is already up-to-date!
[nltk_data] Downloading package words to /root/nltk_data...
[nltk_data]   Package words is already up-to-date!
True
```

```python
1 entities = nltk.chunk.ne_chunk(tagged)
2 #entities
```

```python
1 nltk.download('treebank')
```

```
[nltk_data] Downloading package treebank to /root/nltk_data...
[nltk_data]   Unzipping corpora/treebank.zip.
True
```

```python
1 from nltk.corpus import treebank
2 t = treebank.parsed_sents('the father')[0]
```

```
3 t.draw()
```

```
---------------------------------------------------------------------------
OSError                                   Traceback (most recent call last)
<ipython-input-83-7c3566bbf033> in <module>
      1 from nltk.corpus import treebank
----> 2 t = treebank.parsed_sents('the father')[0]
      3 t.draw()

                            ▲▼ 5 frames
/usr/local/lib/python3.7/dist-packages/nltk/data.py in __init__(self, _path)
    310         _path = os.path.abspath(_path)
    311         if not os.path.exists(_path):
--> 312             raise OSError("No such file or directory: %r" % _path)
    313         self._path = _path
    314

OSError: No such file or directory: '/root/nltk_data/corpora/treebank/combined/the father'
```

> SEARCH STACK OVERFLOW

```
1
```

## ▾ - Exercici 3

Realitza sentiment analysis al teu conjunt de dades.

```
1 import nltk
2 nltk.download('vader_lexicon')
3 nltk.download('punkt')
```

```
[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
True
```

```
1 import nltk
2 tokenizer = nltk.data.load('tokenizers/punkt/english.pickle')
3 sentences = tokenizer.tokenize(text)
```

```
1 from nltk.sentiment.vader import SentimentIntensityAnalyzer
2 from nltk import sentiment
3 from nltk import word_tokenize
```

```
1 analizador = SentimentIntensityAnalyzer()
```

```
 1 import statistics
 2
 3 valores=[]
 4
 5 neg=[]
 6 neu=[]
 7 pos=[]
 8 compound=[]
 9
10
11 for i, sentence in enumerate(sentences):
12     #print(sentence)
13     scores = analizador.polarity_scores(sentence)
14     print('i= ', i, scores, type(scores))
15     neu.append(scores['neu'])
16     neg.append(scores['neg'])
17     compound.append(scores['compound'])
18     pos.append(scores['pos'])
19     valores.append(scores)
20
21
22
```

```
23
24
25      '''
26      for key in scores:
27          print(key, ': ', scores[key])
28          #print()
29          '''
30 print('\n Valores: ')
31 valores[2]['neu']
32
33 print('neutros:', neu,'     -> ', round(statistics.mean(neu),2))
34 print('pos:     ', pos,'    -> ', round(statistics.mean(pos),2))
35 print('neg:     ', neg,'    -> ', round(statistics.mean(neg),2))
36 print('Compound:','     -> ', compound, round(statistics.mean(compound),2))
```

```
i=  0 {'neg': 0.0, 'neu': 0.865, 'pos': 0.135, 'compound': 0.3612} <class 'dict'>
i=  1 {'neg': 0.0, 'neu': 0.935, 'pos': 0.065, 'compound': 0.4019} <class 'dict'>
i=  2 {'neg': 0.0, 'neu': 0.846, 'pos': 0.154, 'compound': 0.5994} <class 'dict'>
i=  3 {'neg': 0.0, 'neu': 1.0, 'pos': 0.0, 'compound': 0.0} <class 'dict'>
i=  4 {'neg': 0.0, 'neu': 0.516, 'pos': 0.484, 'compound': 0.8176} <class 'dict'>
i=  5 {'neg': 0.0, 'neu': 0.689, 'pos': 0.311, 'compound': 0.9313} <class 'dict'>
i=  6 {'neg': 0.0, 'neu': 1.0, 'pos': 0.0, 'compound': 0.0} <class 'dict'>
i=  7 {'neg': 0.0, 'neu': 1.0, 'pos': 0.0, 'compound': 0.0} <class 'dict'>
i=  8 {'neg': 0.0, 'neu': 0.777, 'pos': 0.223, 'compound': 0.3182} <class 'dict'>

 Valores:
neutros: [0.865, 0.935, 0.846, 1.0, 0.516, 0.689, 1.0, 1.0, 0.777]     -> 0.85
pos:     [0.135, 0.065, 0.154, 0.0, 0.484, 0.311, 0.0, 0.0, 0.223]     -> 0.15
neg:     [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0]      -> 0.0
Compound:     -> [0.3612, 0.4019, 0.5994, 0.0, 0.8176, 0.9313, 0.0, 0.0, 0.3182] 0.38
```

✓  0 s    completado a las 8:23                                                          ● ✕