

```
import nltk
import matplotlib.pyplot as plt
import re
import unicodedata
from nltk.corpus import stopwords
```

Out[34]: True

Out[35]: True

```
Out[36]: ['a',
          'about',
          'above',
          'after',
          'again',
          'against',
          'ain',
          'all',
          'am',
          'an',
          'and',
          'any',
          'are',
          'aren',
          "aren't",
          'as',
          'at',
          'be',
          'because',
          'been']
```

```
Out[37]: ['neg/cv000_29416.txt',
'neg/cv001_19502.txt',
'neg/cv002_17424.txt',
'neg/cv003_12683.txt',
'neg/cv004_12641.txt',
'neg/cv005_29357.txt',
'neg/cv006_17022.txt',
'neg/cv007_4992.txt',
'neg/cv008_29326.txt',
'neg/cv009_29417.txt',
'neg/cv010_29063.txt',
'neg/cv011_13044.txt',
'neg/cv012_29411.txt',
'neg/cv013_10494.txt',
'neg/cv014_15600.txt',
'neg/cv015_29356.txt',
'neg/cv016_4348.txt',
'neg/cv017_23487.txt',
'neg/cv018_21672.txt',
neg/cv019_16117.txt']
```

```
Out[38]: ['neg', 'pos']
```

```
Out[39]: FreqDist({' ': 77717, 'the': 76529, '.': 65876, 'a': 38106, 'and': 35576, 'of': 34123, 'to': 31937, '"': 30585, 'is': 25195, 'i': 21822, ...})
```

```
In [40]: #Esta funcion elimina caracteres especiales y acentos
#Fue creada por CHATGP
def remove_special_characters_and_accents(text):
    text = re.sub(r'^a-zA-Z\s]', '', text)
    text = unicodedata.normalize('NFKD', text).encode('ASCII', 'ignore').decode('ASCII')
    return text
```

```
In [41]: categoriaPositiva = nltk.corpus.movie_reviews.words(categories='pos')
print(f"Longitud de categoria positiva antes de la limpieza {len(categoriaPositiva)}")
categoriaNegativas = nltk.corpus.movie_reviews.words(categories='neg')
print(f"Longitud de categoria negativa antes de la limpieza {len(categoriaNegativas)}")
```

Longitud de categoria positiva antes de la limpieza 832564  
Longitud de categoria negativa antes de la limpieza 751256

Limpieza de datos de categoria positiva

```
In [56]: def limpiezaPalabras(lista, stopWords):
    sinStopWords = [w for w in lista if not w.lower() in stopWords]
    sinCaracteresEspeciales = [remove_special_characters_and_accents(w) for w in sinStopWords]
    listaPalabras = [w for w in sinCaracteresEspeciales if w != ""]
    return listaPalabras
```

```
In [106]: #Codigo de prueba
# Test = ["Alexis", "aLexis", "Leal"]
# Len(nltk.FreqDist([w.lower() for w in Test]))
```

Out[106]: 2

```
In [50]: frecuenciaPositiva = nltk.FreqDist([w.lower() for w in limpiezaPalabras(categoriaPositiva, stop_words)])
print(f"Longitud de categoria positiva despues de la limpieza {len(frecuenciaPositiva)}")
```

Longitud de categoria positiva despues de la limpieza 29750

```
In [55]: frecuenciaNegativa = nltk.FreqDist([w.lower() for w in limpiezaPalabras(categoriaNegativas, stop_words)])
print(f"Longitud de categoria negativa despues de la limpieza {len(frecuenciaNegativa)}")
```

Longitud de categoria negativa despues de la limpieza 27841

```
In [109]: print(f"Numero de caracteres categoria positiva {frecuenciaPositiva.N()}")
print(f"Numero de caracteres catageria negativa {frecuenciaNegativa.N()}")
```

Numero de caracteres categoria positiva 372286  
Numero de caracteres catageria negativa 330613

```
In [61]: top20Positivas = frecuenciaPositiva.most_common(20)
```

```
In [100]: top20Negativas = frecuenciaNegativa.most_common(20)
```

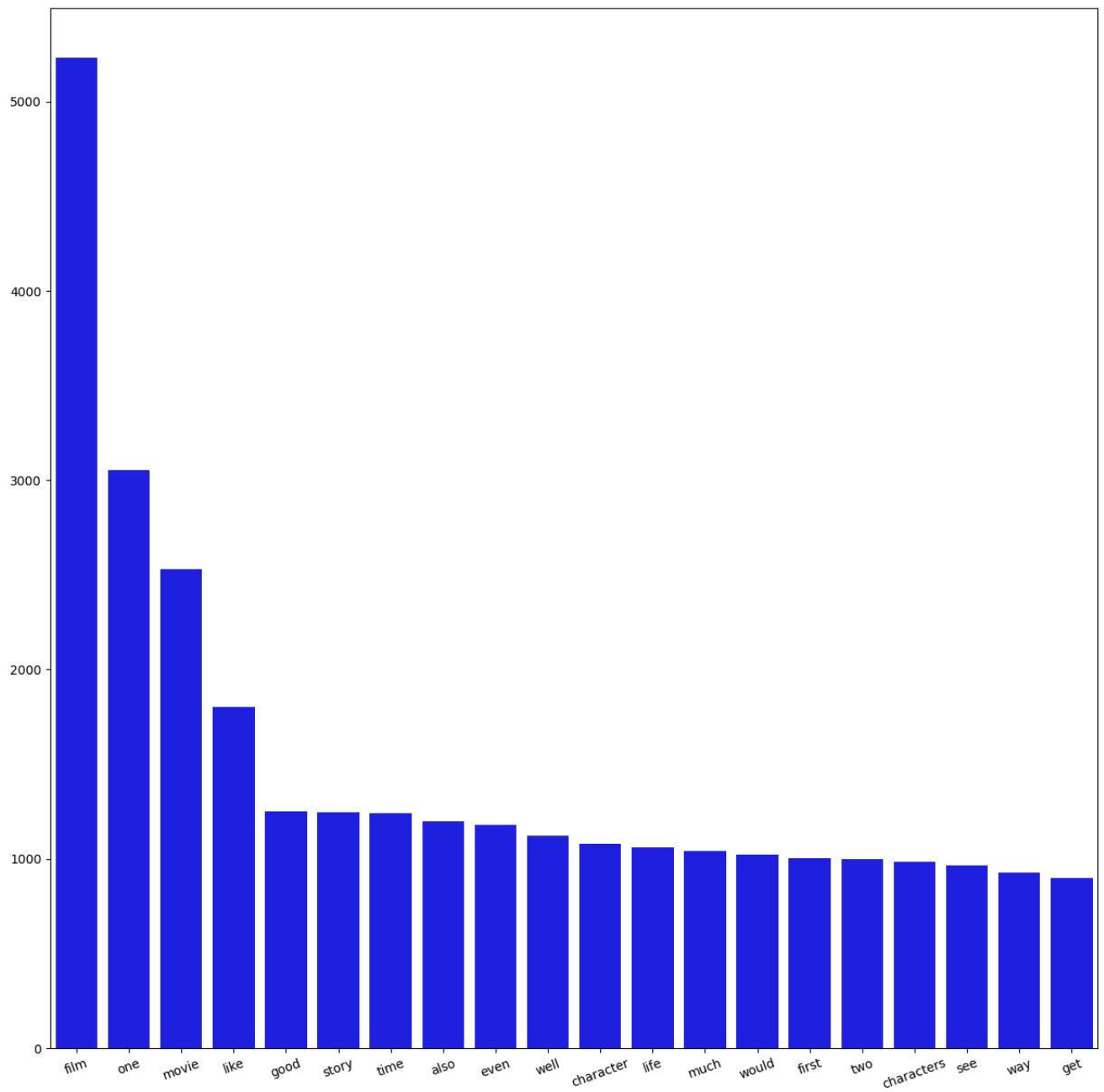
```
In [70]: #Creamos un diccionario de Las top 20 de La categoria por que este nos devuelve una tupla
dicPos = {}
for x,y in top20Positivas:
    dicPos[x] = y
```

```
In [101]: dicNeg = {}
for x,y in top20Negativas:
    dicNeg[x] = y
```

```
In [93]: import seaborn as sn
import pandas as pd
```

Graficamos el top 20 de cada categoria

```
In [110]: #Creamos una data frame con La Libreria pandas
frec_dist = pd.Series(dicPos)
fig, ax = plt.subplots(figsize=(15,15))
grafica = sn.barplot(x=frec_dist.keys(), y=frec_dist.values, ax=ax,color='blue')
plt.xticks(rotation=20);
```



```
In [103]: freq_distN = pd.Series(dicNeg)
fig, ax = plt.subplots(figsize=(15,15))
grafica = sn.barplot(x=freq_distN.keys(), y=freq_distN.values, ax=ax,color='blue')
plt.xticks(rotation=20);
```

