

In []:

```
1  #####
2  #
3  #           Problem : To Find the Frequency of Words in Webpage
4  #
5  #####
6  # Importing the Libraries for NLP
7  import nltk
8  from bs4 import BeautifulSoup
9  %matplotlib inline
10
11 # importing urllib package to read the webpage
12 import urllib.request
13
14 # opening the given Url
15 response = urllib.request.urlopen('http://php.net/')
16
17 # getting the web page with the html tags ( unstructured /non formatted )
18 html = response.read()
19
20
21 # beautifulsoup library is used for reading web pages and gives the page in the format
22 soup = BeautifulSoup(html,"html5lib")
23
24 #removes the space
25 text = soup.get_text(strip=True)
26 #print(text)
27
28 # getting the words from the text found on the page
29 tokens = [t for t in text.split()]
30
31 # calling the Frequency Distribution function , to get the frequency of words in the web page
32 freq = nltk.FreqDist(tokens)
33 maxwd = freq.max()
34 maxcnt=0
35
36 print( " The Words along with their frequency in the web page :\n",'-'*80)
37
38 # The frequency is a dictionary with word and the frequency
39 for key,val in freq.items():
40     if (key == maxwd):
41         maxcnt = val
```

```
42     print (str(key) + ':' + str(val))
43
44
```

```
In [ ]: 1 # The words with maximun frequency are
2 print("\n The word with maximun frequency in the webpage and frequency is :\n",'-'*80)
3 print(freq.max(),maxcnt)
4
```

```
In [ ]: 1 import pandas as pd
2 print("\n The Top most 10 words used in the webpage are :\n",'-'*80)
3 # Most common 10 words in the webpage are
4 words,frq = zip(*freq.most_common(10))
5 print(words)
```

```
In [ ]: 1 # Drawing a plot showing the words and their frequency distribution
2 print(" This Plot shows the Frequency Distribution of top 30 words in webpage \n")
3 freq.plot(30,cumulative=False)
4 print ('''\n From The plot , we observe PHP is used maximum on this page , hence the page appears to be
5         based on PHP ''')
6 #*****
```