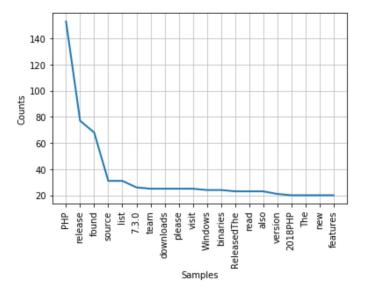
Problem Statement

In this assignment students have to find the frequency of words in a webpage. User can use urllib and BeautifulSoup to extract text from webpage.

```
In [3]: from bs4 import BeautifulSoup
import urllib.request
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
response = urllib.request.urlopen('http://php.net/')
html = response.read()
soup = BeautifulSoup(html,"html5lib")
text = soup.get_text(strip=True)
tokens = [t for t in text.split()]
clean tokens = tokens[:]
sr = stopwords.words('english')
for token in tokens:
    if token in stopwords.words('english'):
        clean tokens.remove(token)
freq = nltk.FreqDist(clean tokens)
for key,val in freq.items():
[nltk_data] Downloading package stopwords to
[nltk data] C:\Users\HP\AppData\Roaming\nltk data...
[nltk_data] Package stopwords is already up-to-date!
PHP::1
Hypertext:1
PreprocessorDownloadsDocumentationGet:1
InvolvedHelpGetting:1
StartedIntroductionA:1
simple:1
tutorialLanguage:1
ReferenceBasic:1
syntaxTypesVariablesConstantsExpressionsOperatorsControl:1
StructuresFunctionsClasses:1
ObjectsNamespacesErrorsExceptionsGeneratorsReferences:1
ExplainedPredefined:1
VariablesPredefined:1
ExceptionsPredefined:1
Interfaces:1
ClassesContext:1
```

1 of 2





2 of 2