

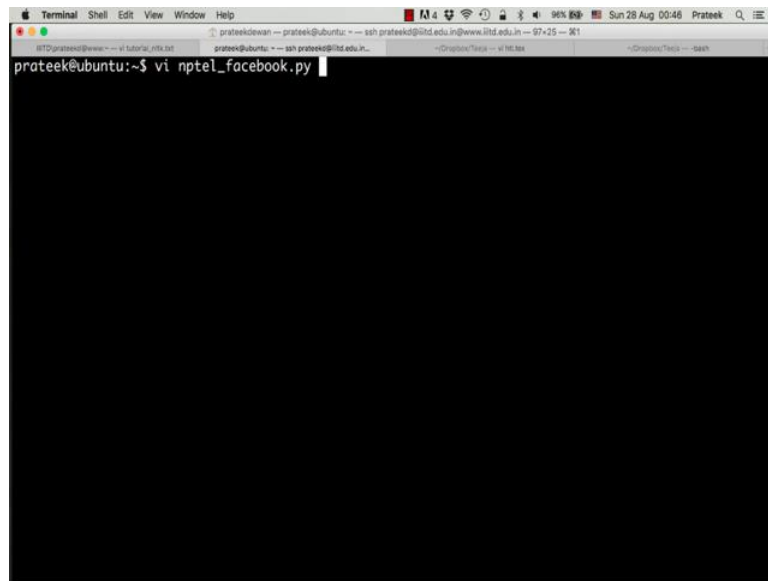
Privacy and Security in Online Social Networks
Department of Computer Science and Engineering
Indian Institute of Technology, Madras

Lecture - 25
Tutorial 5: Analyzing text using Python NLTK

Hi everyone, welcome to another tutorial for the course Privacy and Security in Online Social Media. In the last few tutorials, we have seen how to collect data from Facebook and Twitter and we have seen some basics of python and Linux - Ubuntu and we have also seen some parts of social network analysis using Gephi.

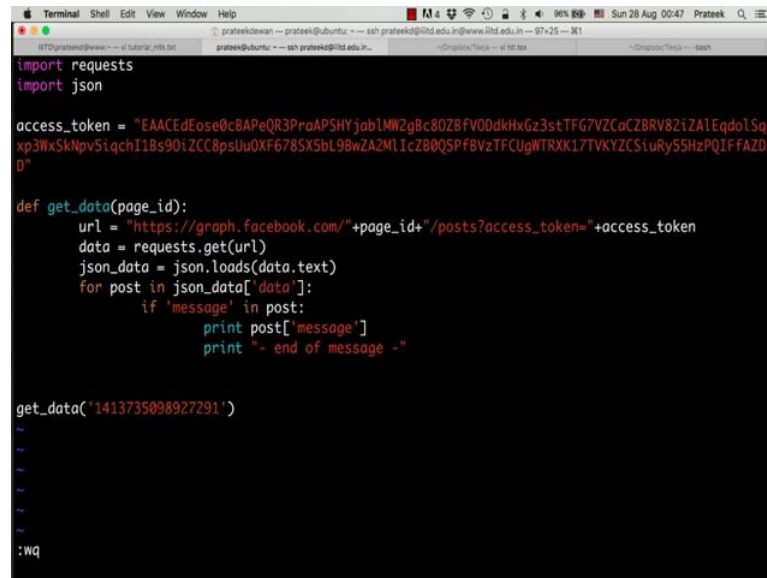
In today's tutorial, what we will learn is how to process the data that you have collected from Facebook and Twitter. **Specifically, today** we will be looking at the basics of working with textual data that you have obtained.

(Refer Slide Time: 00:50)

A screenshot of a terminal window on a Linux system. The terminal title bar shows 'Terminal Shell Edit View Window Help' and the user 'Prateek'. The prompt is 'prateek@ubuntu:~\$' and the command entered is 'vi nptel_facebook.py'. The terminal content is mostly black, indicating the file is being opened in a text editor.

Now, let us go back to the Facebook data collection script.

(Refer Slide Time: 00:54)



```
Terminal Shell Edit View Window Help
prateek@ubuntu:~$ vi nptel_facebook.py
prateek@ubuntu:~$ python nptel_facebook.py

import requests
import json

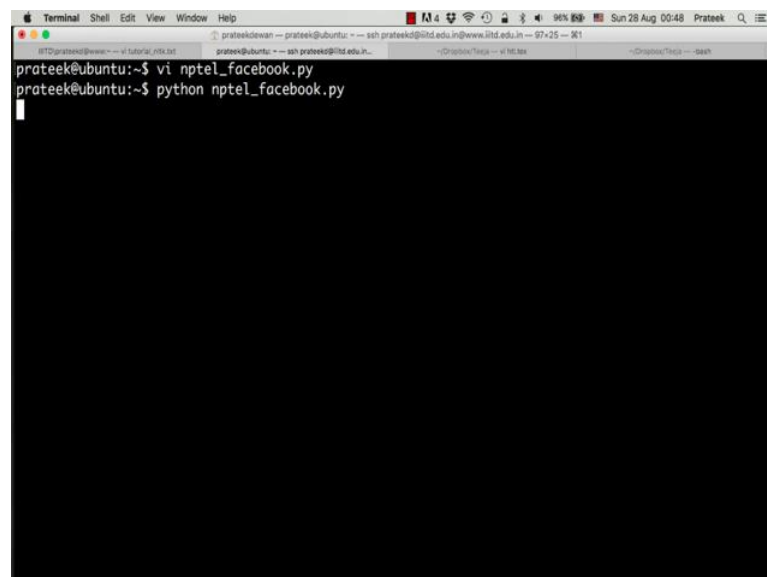
access_token = "EAACEdEose0c8APEQR3PraAPSHYjob1MW2gBc80ZBFV00dkHxGz3stTFG7VZCoCZBRV82iZAIEdqo1Soc
xp3WxSkNpv5iqchI18s90iZCC8pslu0XF6785X5bL9BwZA2M1IcZB0QSPfBVzTFCUgWTRXK17TVKYZCSluRy55HzPQIFFAZ0Z
0"

def get_data(page_id):
    url = "https://graph.facebook.com/"+page_id+"/posts?access_token="+access_token
    data = requests.get(url)
    json_data = json.loads(data.text)
    for post in json_data['data']:
        if 'message' in post:
            print post['message']
            print "- end of message -"

get_data('1413735098927291')
```

And if you remember, this is what we did to collect data from Facebook, we created a simple function. We defined a URL which was pointing to the graph API and we used a requests library to extract data from Facebook using this URL and if you remember, we saw how to print the message field present in Facebook posts. So, let us run this code again just to revise what we did in the first couple of weeks.

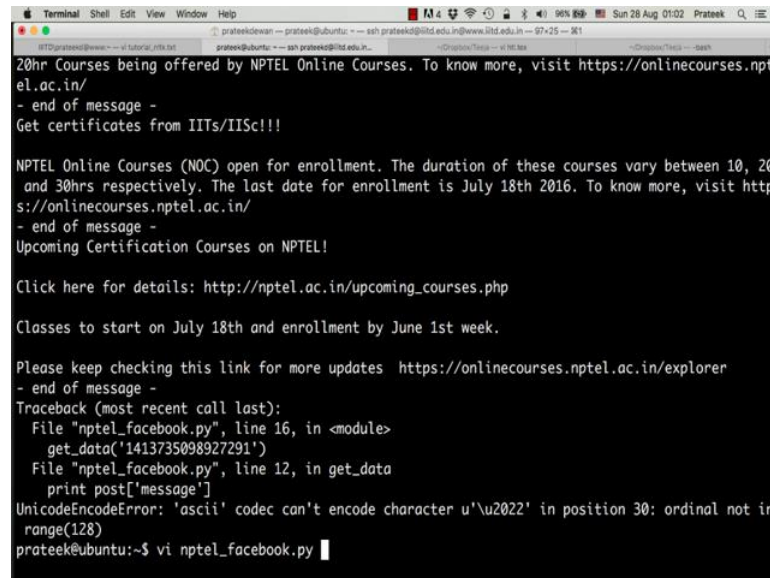
(Refer Slide Time: 01:23)



```
Terminal Shell Edit View Window Help
prateek@ubuntu:~$ vi nptel_facebook.py
prateek@ubuntu:~$ python nptel_facebook.py
```

So, we say python space nptel underscore facebook.

(Refer Slide Time: 01:30)



```
Terminal Shell Edit View Window Help
prateek@prateek:~$ python nptel_facebook.py
20hr Courses being offered by NPTEL Online Courses. To know more, visit https://onlinecourses.nptel.ac.in/
- end of message -
Get certificates from IITs/IISc!!!

NPTEL Online Courses (NOC) open for enrollment. The duration of these courses vary between 10, 20
and 30hrs respectively. The last date for enrollment is July 18th 2016. To know more, visit http
s://onlinecourses.nptel.ac.in/
- end of message -
Upcoming Certification Courses on NPTEL!

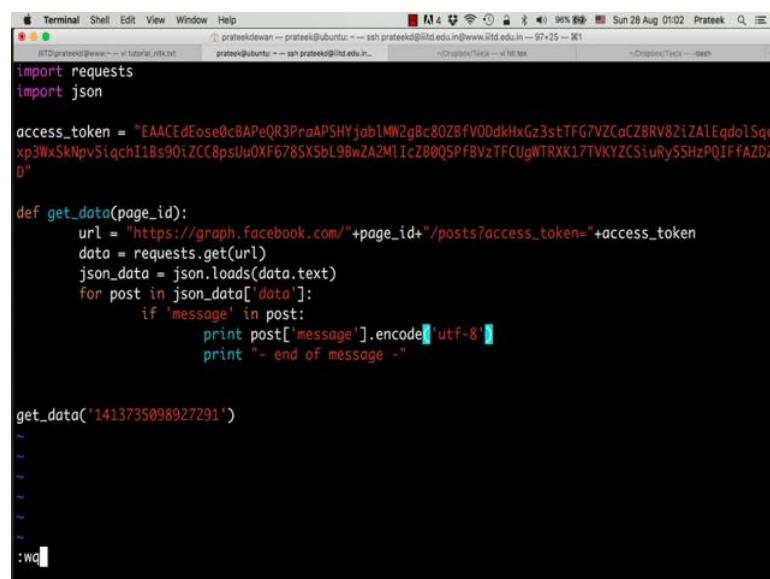
Click here for details: http://nptel.ac.in/upcoming_courses.php

Classes to start on July 18th and enrollment by June 1st week.

Please keep checking this link for more updates https://onlinecourses.nptel.ac.in/explorer
- end of message -
Traceback (most recent call last):
  File "nptel_facebook.py", line 16, in <module>
    get_data('1413735098927291')
  File "nptel_facebook.py", line 12, in get_data
    print post['message']
UnicodeEncodeError: 'ascii' codec can't encode character u'\u2022' in position 30: ordinal not in
range(128)
prateek@ubuntu:~$ vi nptel_facebook.py
```

And it will print the message part of the posts that the NPTEL page has done. So, there is an error here. So, what happens is **by** default, when you are printing something on your terminal, if the data is not in the standard ASCII format and if there is a Unicode data; by default python uses ASCII to print character. Now, if it encounters a Unicode character it is not able to print it and it usually throws an error. So, if you read this error it says that ASCII codec cannot encode character u and some code in position 30. So, any string which starts with a u in the beginning is usually a Unicode string. Now, to handle such strings what we can do is.

(Refer Slide Time: 02:17).



```
Terminal Shell Edit View Window Help
prateek@prateek:~$ python nptel_facebook.py
import requests
import json

access_token = "EAACEdEose0cBAPeQR3PraAPSHYjob1MW2gBc80ZBFV0DdkHxGz3stTFG7VZCoCZBRV82iZAIeqdo1Sac
xp3WxSkNpv5iqchI18s90iZCC8pstu0XF6785X5bL9BwZA2M1IcZB0Q5PF8VzTFCUgWTRXK17TVKYZCSiuRy55HzPQIFFAZ0Z
0"

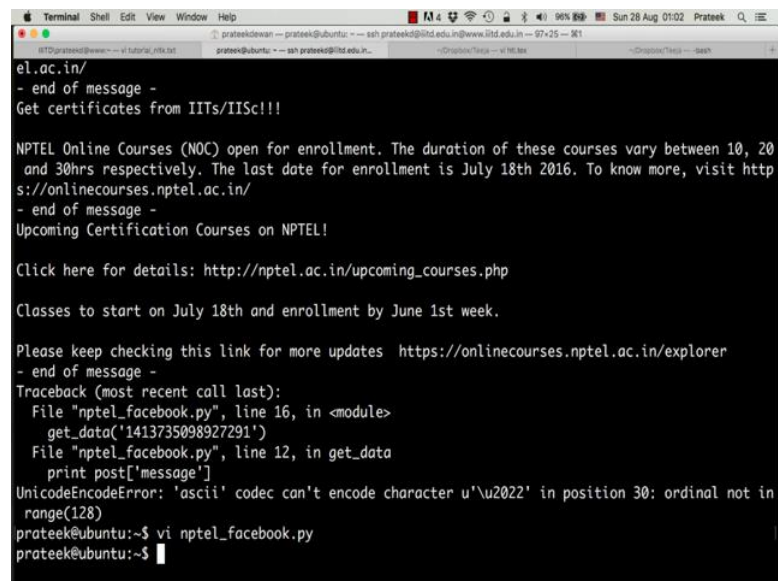
def get_data(page_id):
    url = "https://graph.facebook.com/"+page_id+"/posts?access_token="+access_token
    data = requests.get(url)
    json_data = json.loads(data.text)
    for post in json_data['data']:
        if 'message' in post:
            print post['message'].encode('utf-8')
            print "- end of message -"

get_data('1413735098927291')

:wa
```

We can just add dot encode UTF-8 in the script. So, what this will do is, this will encode all the messages or all the text that we are getting into UTF-8 format and then print it. So, if there is any Unicode character in **there** which is not the standard a to z or 0 to 9 kind of a character, this will convert it to the whatever character it is, for example, Hindi text for example, Arabic text; those are text which are not contained in the ASCII character set. So, they can be printed using UTF-8 encoding. Now, once you have done this.

(Refer Slide Time: 03:04)



```
prateek@ubuntu:~$ python nptel_facebook.py
el.ac.in/
- end of message -
Get certificates from IITs/IISc!!!

NPTEL Online Courses (NOC) open for enrollment. The duration of these courses vary between 10, 20
and 30hrs respectively. The last date for enrollment is July 18th 2016. To know more, visit http
s://onlinecourses.nptel.ac.in/
- end of message -
Upcoming Certification Courses on NPTEL!

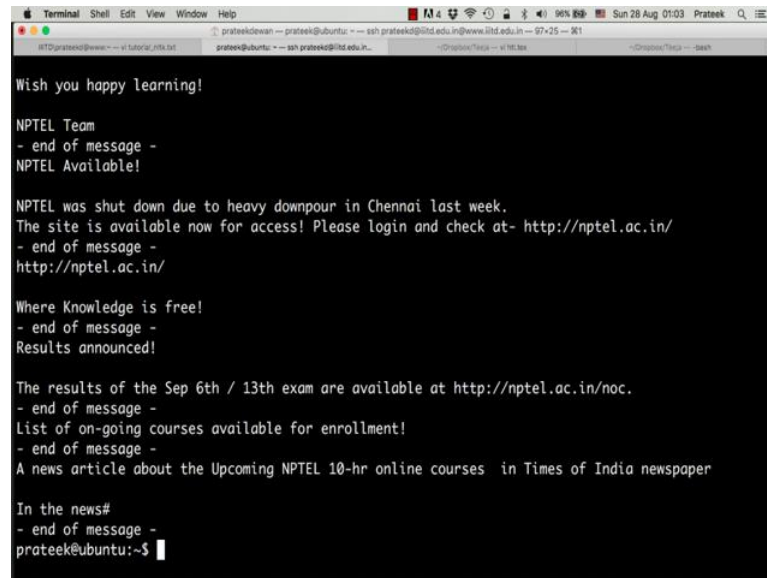
Click here for details: http://nptel.ac.in/upcoming_courses.php

Classes to start on July 18th and enrollment by June 1st week.

Please keep checking this link for more updates https://onlinecourses.nptel.ac.in/explorer
- end of message -
Traceback (most recent call last):
  File "nptel_facebook.py", line 16, in <module>
    get_data('1413735098927291')
  File "nptel_facebook.py", line 12, in get_data
    print post['message']
UnicodeEncodeError: 'ascii' codec can't encode character u'\u2022' in position 30: ordinal not in
range(128)
prateek@ubuntu:~$ vi nptel_facebook.py
prateek@ubuntu:~$
```

Now, let us try to run this script again.

(Refer Slide Time: 03:06)



```
Wish you happy learning!

NPTEL Team
- end of message -
NPTEL Available!

NPTEL was shut down due to heavy downpour in Chennai last week.
The site is available now for access! Please login and check at- http://npTEL.ac.in/
- end of message -
http://npTEL.ac.in/

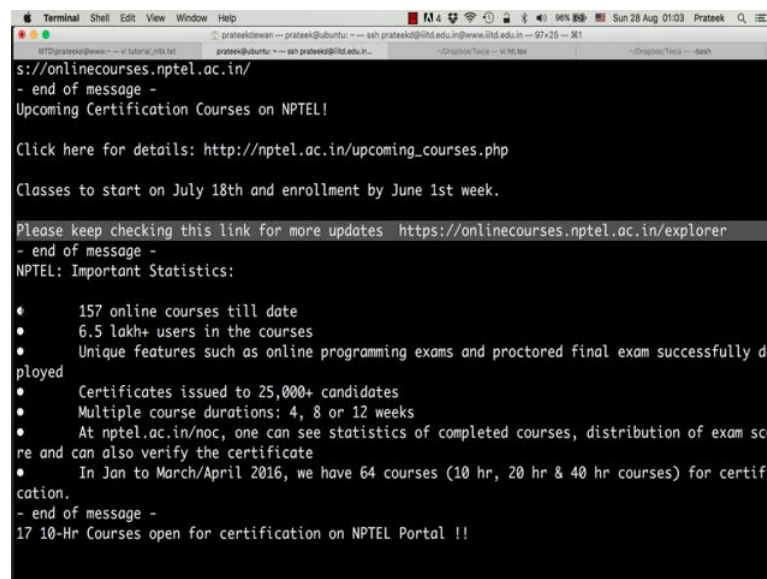
Where Knowledge is free!
- end of message -
Results announced!

The results of the Sep 6th / 13th exam are available at http://npTEL.ac.in/noc.
- end of message -
List of on-going courses available for enrollment!
- end of message -
A news article about the Upcoming NPTEL 10-hr online courses in Times of India newspaper

In the news#
- end of message -
prateek@ubuntu:~$
```

And there you go, there is no error.

(Refer Slide Time: 03:12)



```
s://onlinecourses.nptel.ac.in/
- end of message -
Upcoming Certification Courses on NPTEL!

Click here for details: http://npTEL.ac.in/upcoming\_courses.php

Classes to start on July 18th and enrollment by June 1st week.

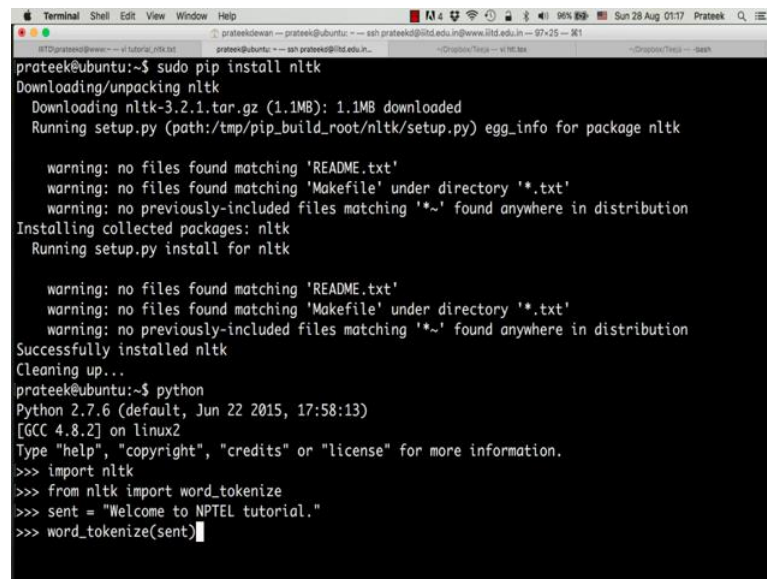
Please keep checking this link for more updates https://onlinecourses.nptel.ac.in/explorer
- end of message -
NPTEL: Important Statistics:

• 157 online courses till date
• 6.5 lakh+ users in the courses
• Unique features such as online programming exams and proctored final exam successfully deployed
• Certificates issued to 25,000+ candidates
• Multiple course durations: 4, 8 or 12 weeks
• At npTEL.ac.in/noc, one can see statistics of completed courses, distribution of exam score and can also verify the certificate
• In Jan to March/April 2016, we have 64 courses (10 hr, 20 hr & 40 hr courses) for certification.
- end of message -
17 10-Hr Courses open for certification on NPTEL Portal !!
```

And if you scroll back, we can probably look at what character it was. So, this bullet point here this is not a standard ASCII character and this is where the error was coming. If you remember, this was the last post that we received and after that there was an error message. This is probably what caused the error. **Now** to start processing text, what you need is a library which can handle text. So, in python there is a library which is called natural language tool kit or NLTK in short. Let us install that first. So, you say sudo pip

install NLTK.

(Refer Slide Time: 03:47)



```
prateek@ubuntu:~$ sudo pip install nltk
Downloading/unpacking nltk
  Downloading nltk-3.2.1.tar.gz (1.1MB): 1.1MB downloaded
  Running setup.py (path:/tmp/pip_build_root/nltk/setup.py) egg_info for package nltk

  warning: no files found matching 'README.txt'
  warning: no files found matching 'Makefile' under directory '*.txt'
  warning: no previously-included files matching '*~' found anywhere in distribution
Installing collected packages: nltk
  Running setup.py install for nltk

  warning: no files found matching 'README.txt'
  warning: no files found matching 'Makefile' under directory '*.txt'
  warning: no previously-included files matching '*~' found anywhere in distribution
Successfully installed nltk
Cleaning up...
prateek@ubuntu:~$ python
Python 2.7.6 (default, Jun 22 2015, 17:58:13)
[GCC 4.8.2] on linux2
Type "help", "copyright", "credits" or "license" for more information.
>>> import nltk
>>> from nltk import word_tokenize
>>> sent = "Welcome to NPTEL tutorial."
>>> word_tokenize(sent)
```

And press enter. Now, nltk is installed, if you go to the python shell and say import nltk, the import is successful, but interestingly just downloading and installing nltk does not suffice to perform the basic functions that you need.

So, for example, if you want to convert a sentence into a set of words, there is something called word-**tokenizer**. So, you can import it from nltk by saying, from nltk import word underscore tokenize and let us create a string say a sentence is equal to ‘Welcome to NPTEL tutorial’, now what this word underscore tokenize is supposed to do is when we give it a string. It is supposed to give us a list of the words present in the sentence. So, if we say word underscore tokenize sent.

(Refer Slide Time: 04:59)

```
Terminal Shell Edit View Window Help
prateek@prateek:~$ python
>>> tokenize.tokenize('nize')
Traceback (most recent call last):
  File "/usr/local/lib/python2.7/dist-packages/nltk/tokenize/_init_.py", line 90, in sent_tokenize
    tokenizer = load('tokenizers/punkt/{0}.pickle'.format(language))
  File "/usr/local/lib/python2.7/dist-packages/nltk/data.py", line 801, in load
    opened_resource = _open(resource_url)
  File "/usr/local/lib/python2.7/dist-packages/nltk/data.py", line 919, in _open
    return find(path_, path + ['.']).open()
  File "/usr/local/lib/python2.7/dist-packages/nltk/data.py", line 641, in find
    raise LookupError(resource_not_found)
LookupError:
*****
Resource u'tokenizers/punkt/english.pickle' not found. Please
use the NLTK Downloader to obtain the resource: >>>
nltk.download()
Searched in:
- '/home/prateek/nltk_data'
- '/usr/share/nltk_data'
- '/usr/local/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/local/lib/nltk_data'
- u''
*****
>>> nltk.download()
```

It is supposed to give us the words in the sentence, but it throws an error instead, now why is that that is because the nltk library requires some data to perform the basic operations, that data is present in a set of corpuses.

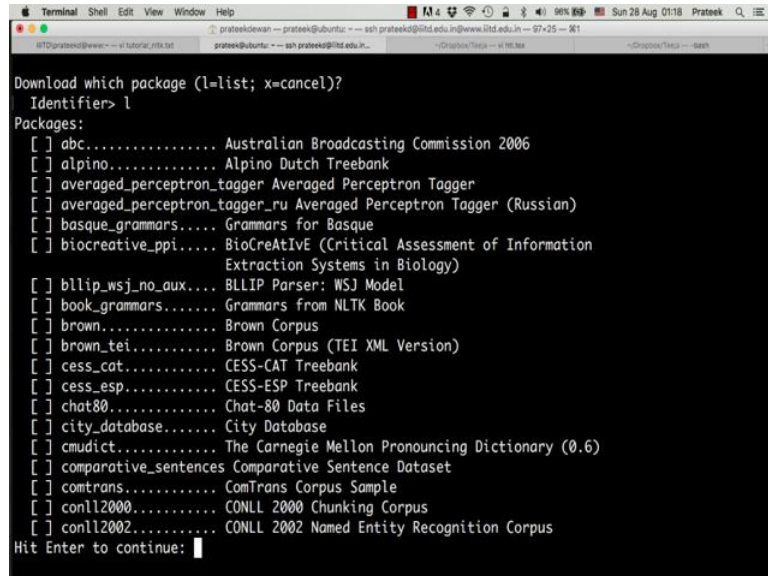
So, if you look at the error message carefully, it says that you should use this nltk dot download function to download the resources that are required to perform this operation. So, you just type nltk dot download and press enter and you will see this nltk downloader menu appear now type d and press enter.

(Refer Slide Time: 05:32)

```
Terminal Shell Edit View Window Help
prateek@prateek:~$ python
>>> tokenize.tokenize('nize')
Traceback (most recent call last):
  File "/usr/local/lib/python2.7/dist-packages/nltk/tokenize/_init_.py", line 90, in sent_tokenize
    tokenizer = load('tokenizers/punkt/{0}.pickle'.format(language))
  File "/usr/local/lib/python2.7/dist-packages/nltk/data.py", line 801, in load
    opened_resource = _open(resource_url)
  File "/usr/local/lib/python2.7/dist-packages/nltk/data.py", line 919, in _open
    return find(path_, path + ['.']).open()
  File "/usr/local/lib/python2.7/dist-packages/nltk/data.py", line 641, in find
    raise LookupError(resource_not_found)
LookupError:
*****
Resource u'tokenizers/punkt/english.pickle' not found. Please
use the NLTK Downloader to obtain the resource: >>>
nltk.download()
Searched in:
- '/home/prateek/nltk_data'
- '/usr/share/nltk_data'
- '/usr/local/share/nltk_data'
- '/usr/lib/nltk_data'
- '/usr/local/lib/nltk_data'
- u''
*****
>>> nltk.download()
NLTK Downloader
-----
d) Download l) List u) Update c) Config h) Help q) Quit
-----
Downloader> d
Download which package (l=list; x=cancel)?
Identifier> l
Packages:
|
```

So, the downloader will ask you which package you want to download, now type l and press enter to see the list of available packages.

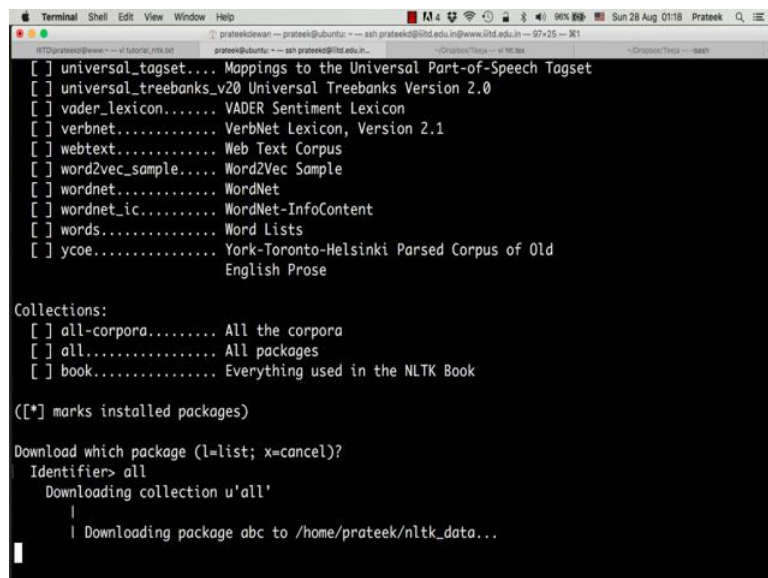
(Refer Slide Time: 05:42)



```
Terminal Shell Edit View Window Help
prateek@prateek:~$ ssh prateek@ltd.edu.in
prateek@ltd.edu.in:~$
Download which package (l=list; x=cancel)?
Identifier> l
Packages:
[ ] abc..... Australian Broadcasting Commission 2006
[ ] alpino..... Alpino Dutch Treebank
[ ] averaged_perceptron_tagger Averaged Perceptron Tagger
[ ] averaged_perceptron_tagger_ru Averaged Perceptron Tagger (Russian)
[ ] basque_grammars..... Grammars for Basque
[ ] biocreative_ppi..... BioCreAtIVE (Critical Assessment of Information
Extraction Systems in Biology)
[ ] bllip_wsj_no_aux.... BLLIP Parser: WSJ Model
[ ] book_grammars..... Grammars from NLTK Book
[ ] brown..... Brown Corpus
[ ] brown_tei..... Brown Corpus (TEI XML Version)
[ ] cess_cat..... CESS-CAT Treebank
[ ] cess_esp..... CESS-ESP Treebank
[ ] chat80..... Chat-80 Data Files
[ ] city_database..... City Database
[ ] cmudict..... The Carnegie Mellon Pronouncing Dictionary (0.6)
[ ] comparative_sentences Comparative Sentence Dataset
[ ] comtrans..... ComTrans Corpus Sample
[ ] conll2000..... CONLL 2000 Chunking Corpus
[ ] conll2002..... CONLL 2002 Named Entity Recognition Corpus
Hit Enter to continue: |
```

Now, there is a whole big list of packages that can be used with the NLTK library. So, instead of selecting **the ones** we want, what we will do is, we will just type all a double l and press enter which will allow us to download all the packages which are available.

(Refer Slide Time: 05:50)



```
Terminal Shell Edit View Window Help
prateek@prateek:~$ ssh prateek@ltd.edu.in
prateek@ltd.edu.in:~$
[ ] universal_tagset.... Mappings to the Universal Part-of-Speech Tagset
[ ] universal_treebanks_v20 Universal Treebanks Version 2.0
[ ] vader_lexicon..... VADER Sentiment Lexicon
[ ] verbnets..... VerbNet Lexicon, Version 2.1
[ ] webtext..... Web Text Corpus
[ ] word2vec_sample.... Word2Vec Sample
[ ] wordnet..... WordNet
[ ] wordnet_ic..... WordNet-InfoContent
[ ] words..... Word Lists
[ ] ycoe..... York-Toronto-Helsinki Parsed Corpus of Old
English Prose
Collections:
[ ] all-corpora..... All the corpora
[ ] all..... All packages
[ ] book..... Everything used in the NLTK Book
([*] marks installed packages)
Download which package (l=list; x=cancel)?
Identifier> all
Downloaded collection u'all'
|
| Downloading package abc to /home/prateek/nltk_data...
```



```
Terminal Shell Edit View Window Help
prateek@prateek: ~ --- ssh prateek@lhd.edu.in --- 97x25 --- 3E1
ITD@prateek@www: ~ --- ssh prateek@lhd.edu.in --- 97x25 --- 3E1
prateek@prateek: ~ --- ssh prateek@lhd.edu.in --- 97x25 --- 3E1
prateek@prateek: ~ --- ssh prateek@lhd.edu.in --- 97x25 --- 3E1

| /home/prateek/nltk_data...
| Downloading package mte_teip5 to /home/prateek/nltk_data...
| Unzipping corpora/mte_teip5.zip.
| Downloading package averaged_perceptron_tagger to
| /home/prateek/nltk_data...
| Unzipping taggers/averaged_perceptron_tagger.zip.
| Downloading package panlex_lite to /home/prateek/nltk_data...
| Unzipping corpora/panlex_lite.zip.
| Downloading package perluniprops to
| /home/prateek/nltk_data...
| Unzipping misc/perluniprops.zip.
| Downloading package nonbreaking_prefixes to
| /home/prateek/nltk_data...
| Unzipping corpora/nonbreaking_prefixes.zip.
| Downloading package vader_lexicon to
| /home/prateek/nltk_data...
| Downloading package porter_test to /home/prateek/nltk_data...
| Unzipping stemmers/porter_test.zip.
|
Done downloading collection all

-----
d) Download l) List u) Update c) Config h) Help q) Quit
-----
Downloader>
```

Now, this entire corpus is about 10 GB in size. So, if you do not **have enough bandwidth** **or** have slow internet, you can just download the packages that are required. So, for example, for this particular example, if you scroll back, you can see that you needed something called; you needed a package called punkt. Now, if you scroll down in the list that nltk showed you can locate this package here.

(Refer Slide Time: 06:20)

```
Terminal Shell Edit View Window Help
prateek@prateek: ~ --- ssh prateek@lhd.edu.in --- 97x25 --- 3E1
ITD@prateek@www: ~ --- ssh prateek@lhd.edu.in --- 97x25 --- 3E1
prateek@prateek: ~ --- ssh prateek@lhd.edu.in --- 97x25 --- 3E1
prateek@prateek: ~ --- ssh prateek@lhd.edu.in --- 97x25 --- 3E1

[ ] paradigms..... Paradigm Corpus
[ ] pe08..... Cross-Framework and Cross-Domain Parser
                  Evaluation Shared Task
Hit Enter to continue:
[ ] perluniprops..... perluniprops: Index of Unicode Version 7.0.0
                      character properties in Perl
[ ] pil..... The Patient Information Leaflet (PIL) Corpus
[ ] pl196x..... Polish language of the XX century sixties
[ ] porter_test..... Porter Stemmer Test Files
[ ] ppattach..... Prepositional Phrase Attachment Corpus
[ ] problem_reports.... Problem Report Corpus
[ ] product_reviews_1... Product Reviews (5 Products)
[ ] product_reviews_2... Product Reviews (9 Products)
[ ] propbank..... Proposition Bank Corpus 1.0
[ ] pros_cons..... Pros and Cons
[ ] ptb..... Penn Treebank
[ ] punkt..... Punkt Tokenizer Models
[ ] qc..... Experimental Data for Question Classification
[ ] reuters..... The Reuters-21578 benchmark corpus, AptMod
                  version
[ ] rslp..... RSLP Stemmer (Removedor de Sufixos da Lingua
              Portuguesa)
[ ] rte..... PASCAL RTE Challenges 1, 2, and 3
[ ] sample_grammars.... Sample Grammars
[ ] semcor..... SemCor 3.0
```

And you can type punkt instead of all to download only this model, this will just be faster, and this will save a lot of your bandwidth.

(Refer Slide Time: 06:30)

```
Terminal Shell Edit View Window Help
prateek@prateek:~$ python3 -m nltk.downloader
Downloader> l

Packages:
[*] abc..... Australian Broadcasting Commission 2006
[*] alpine..... Alpino Dutch Treebank
[*] averaged_perceptron_tagger Averaged Perceptron Tagger
[] averaged_perceptron_tagger_ru Averaged Perceptron Tagger (Russian)
[*] basque_grammars..... Grammars for Basque
[*] biocreative_ppi..... BioCreAtIvE (Critical Assessment of Information
Extraction Systems in Biology)
[*] bllip_wsj_no_aux.... BLLIP Parser: WSJ Model
[*] book_grammars..... Grammars from NLTK Book
[*] brown..... Brown Corpus
[*] brown_tei..... Brown Corpus (TEI XML Version)
[*] cess_cat..... CESS-CAT Treebank
[*] cess_esp..... CESS-ESP Treebank
[*] chat80..... Chat-80 Data Files
[*] city_database..... City Database
[*] cmudict..... The Carnegie Mellon Pronouncing Dictionary (0.6)
[*] comparative_sentences Comparative Sentence Dataset
[*] comtrans..... ComTrans Corpus Sample
[*] conll2000..... CONLL 2000 Chunking Corpus
[*] conll2002..... CONLL 2002 Named Entity Recognition Corpus
Hit Enter to continue: █
```

So, if you type l you can see that a star is marked in front of all the packages that we already have because we downloaded all the packages, but what you can do is just press d and type punkt to download only the tokenizer.

(Refer Slide Time: 06:36)

```
Terminal Shell Edit View Window Help
prateek@prateek:~$ python3 -m nltk.downloader
Downloader> d

Download which package (l=list; x=cancel)?
Identifier> punkt
Downloading package punkt to /home/prateek/nltk_data...
Package punkt is already up-to-date!

d) Download l) List u) Update c) Config h) Help q) Quit

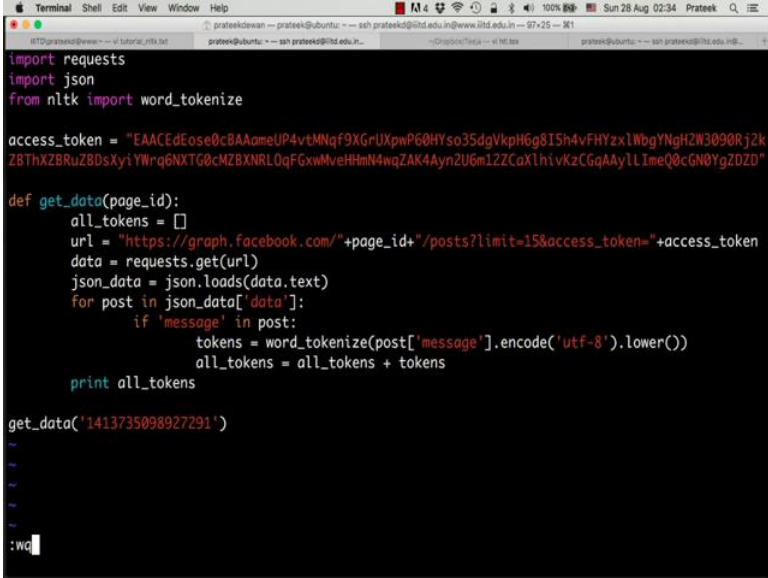
Downloader> q
True
>>> word_tokenize(sent)
['Welcome', 'to', 'NPTEL', 'tutorial', '.']
>>> exit()
prateek@ubuntu:~$ vi nptel_facebook.py █
```

So, this says that package punkt is already up to date. Now, once you are done with downloading the corpus, just type q and press enter to exit; now if you execute that statement again, you will see that word underscore tokenize has tokenized the sentence into words and we have a list of words. Now, we will use this word tokenizer to analyze

the post that we downloaded from the NPTEL Facebook page. Type exit and exit the python shell.

Now, let us go back to our script and try to use this nltk to find the most commonly used words by the NPTEL page for their posts. So, what we will essentially do is we will go through some of the recent posts that the NPTEL page has done and we will try to find out which are the most commonly appearing words in their posts. Now, let us go back to our script.

(Refer Slide Time: 07:39)



```
Terminal Shell Edit View Window Help
prateek@prateek:~$ python3 script.py
import requests
import json
from nltk import word_tokenize

access_token = "EAACEdEose0cBAAameUP4vtMNqf9XGrUXpwP60HYso35dgVkpH6g8I5h4vFHYzx1WbgYNgH2W3090RjzkZBThXZBRuZ8DsXyiYWrq6NXTG0cMZBXNRL0qFGxwMveHhmM4wqZAK4Ayn2U6m1ZzCaXlhivKzCGgAAyLLImeQ0cGN0YgZDZD"

def get_data(page_id):
    all_tokens = []
    url = "https://graph.facebook.com/" + page_id + "/posts?limit=15&access_token=" + access_token
    data = requests.get(url)
    json_data = json.loads(data.text)
    for post in json_data['data']:
        if 'message' in post:
            tokens = word_tokenize(post['message'].encode('utf-8').lower())
            all_tokens = all_tokens + tokens
    print all_tokens

get_data('1413735098927291')
```

Type from nltk import word underscore tokenize and for every message that we are iterating, what we will do is we will say tokens is equal to word underscore tokenize post message dot encode UTF-8. So, what we are essentially doing is we are converting all the posts that we have, into tokens into a list of words, now before doing that what we should do is we should also probably convert this string into lower case letters. So, we say dot lower. So, that we avoid any repetition because when we are counting word frequencies we do not want to count a word in capitals and in smalls as two different words.

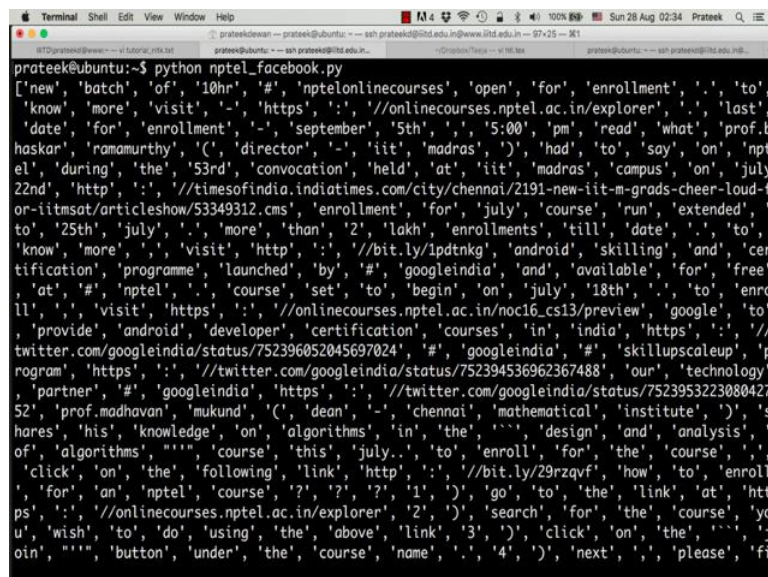
We have to maintain uniformity we will convert everything into lower case and then try to find the frequencies. Now, we have a list of words in the tokens variable for each of the posts; now we need to collect all these tokens in a single list to be able to count the frequencies. So, we will create all underscore tokens as an empty list in the beginning of

the function and every time we calculate the tokens, what we will do is we will say all underscore tokens is equal to all underscore tokens plus tokens.

So, what this will do is, this will append the tokens that we got from the current post into another list which will maintain the list of all the tokens that we have generated. So, we can get rid of these print messages and what we will now do is and let us just print all underscore tokens to see a list of all the **words** that are present in these posts.

So, let us also add a limit for the number of posts that we are analyzing. So, we will say limit is equal to say 15 and access token is equal to access token. So, what we are doing is we are printing all the words that appeared in the last fifteen posts posted by the NPTEL page.

(Refer Slide Time: 10:06)



```
prateek@ubuntu:~$ python npTEL_facebook.py
['new', 'batch', 'of', '10hr', '#', 'npTELonlinecourses', 'open', 'for', 'enrollment', 'to',
'know', 'more', 'visit', '-', 'https', ':', '//onlinecourses.nptel.ac.in/explorer', '.', 'last',
'date', 'for', 'enrollment', '-', 'september', '5th', ',', '5:00', 'pm', 'read', 'what', 'prof.b
haskar', 'ramamurthy', '(', 'director', '-', 'iit', 'madrass', ')', 'had', 'to', 'say', 'on', 'npt
el', 'during', 'the', '53rd', 'convocation', 'held', 'at', 'iit', 'madrass', 'campus', 'on', 'july
22nd', 'http', ':', '//timesofindia.indiatimes.com/city/chennai/2191-new-iit-m-grads-cheer-loud-f
or-iitmsat/articleshow/53349312.cms', 'enrollment', 'for', 'july', 'course', 'run', 'extended',
to', '25th', 'july', ',', 'more', 'than', '2', 'lakh', 'enrollments', 'till', 'date', ',', 'to',
'know', 'more', ',', 'visit', 'http', ':', '//bit.ly/1pdtngk', 'android', 'skilling', 'and', 'cer
tification', 'programme', 'launched', 'by', '#', 'googleindia', 'and', 'available', 'for', 'free',
'at', '#', 'npTEL', ',', 'course', 'set', 'to', 'begin', 'on', 'july', '18th', ',', 'to', 'enro
ll', ',', 'visit', 'https', ':', '//onlinecourses.nptel.ac.in/noc16_cs13/preview', 'google', 'to',
', 'provide', 'android', 'developer', 'certification', 'courses', 'in', 'india', 'https', ':', '//
twitter.com/googleindia/status/752396052045697024', '#', 'googleindia', '#', 'skillupscaleup', 'p
rogram', 'https', ':', '//twitter.com/googleindia/status/752394536962367488', 'our', 'technology',
', 'partner', '#', 'googleindia', 'https', ':', '//twitter.com/googleindia/status/7523953223080427
52', 'prof.madhavan', 'mukund', '(', 'dean', '-', 'chennai', 'mathematical', 'institute', ')', 's
hares', 'his', 'knowledge', 'on', 'algorithms', 'in', 'the', 'design', 'and', 'analysis',
of', 'algorithms', 'course', 'this', 'july', 'to', 'enroll', 'for', 'the', 'course', ',',
'click', 'on', 'the', 'following', 'link', 'http', ':', '//bit.ly/29rzqvf', 'how', 'to', 'enroll
', 'for', 'an', 'npTEL', 'course', '?', '?', '?', '1', ')', 'go', 'to', 'the', 'link', 'at', 'htt
ps', ':', '//onlinecourses.nptel.ac.in/explorer', '2', ')', 'search', 'for', 'the', 'course', 'yo
u', 'wish', 'to', 'do', 'using', 'the', 'above', 'link', '3', ')', 'click', 'on', 'the', 'j
oin', 'button', 'under', 'the', 'course', 'name', '4', ')', 'next', 'please', 'fi
```

So, there you see a list of all the words that were contained in the posts that the NPTEL page did. Now, to calculate the frequency of these words, what we will do is we will use the counter function from python.

(Refer Slide Time: 10:24)

```
Terminal Shell Edit View Window Help
prateek@ubuntu:~$ cat nptel_facebook.py
import requests
import json
from nltk import word_tokenize
from collections import Counter

access_token = "EAACEdEose0cBAAameUP4vtMNqf9XGrUXpwP60HYso3SdgVkpH6g8I5h4vFHYzx1WbgYNgHZW3090RjZk
ZBTHXZBRuZBDsXyiYWrq6NXTG0cMZXNRLQqFGxwMveHmN4wqZAK4Ayn2U6m1Z2CoX1hivKzCGgAay1LlImeQ0cGN0YgZDZD"

def get_data(page_id):
    all_tokens = []
    url = "https://graph.facebook.com/" + page_id + "/posts?limit=15&access_token="+access_token
    data = requests.get(url)
    json_data = json.loads(data.text)
    for post in json_data['data']:
        if 'message' in post:
            tokens = word_tokenize(post['message'].encode('utf-8').lower())
            all_tokens = all_tokens + tokens
    frequencies = Counter(all_tokens)
    for token, count in frequencies.most_common(25):
        print token, count

get_data('1413735098927291')
```

Now, say from collections import counter, what this counter does is given a list it will count the number of instances of all the unique elements present in this list. So, what we will do is, we will say **frequencies** is equal to counter all underscore tokens. Now, to print the token and **its** frequency for the most commonly occurring words, we can say, for token comma count in frequencies dot most underscore common.

So, let us look at the most common 25 words and we say print token comma count. Exit the editor and say python space NPTEL Facebook.

(Refer Slide Time: 11:22)

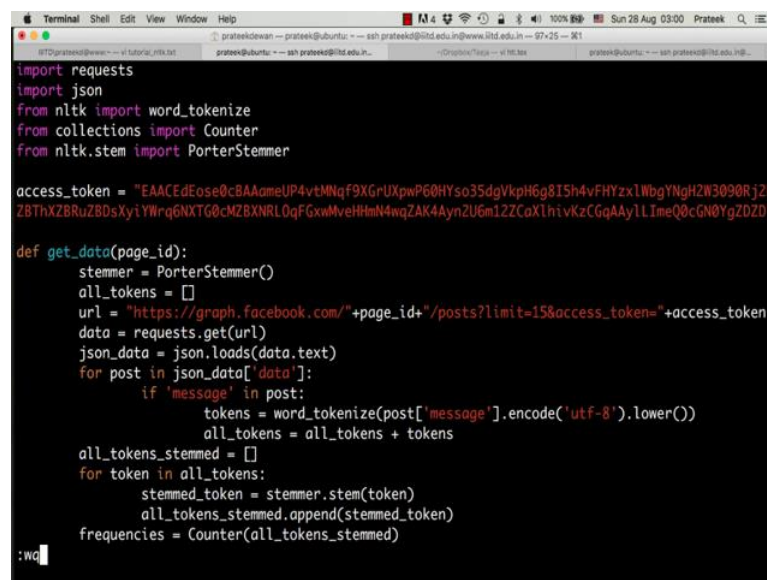
```
Terminal Shell Edit View Window Help
prateek@ubuntu:~$ vi nptel_facebook.py
prateek@ubuntu:~$ python nptel_facebook.py
for', 'enrollment', '.', 'the', 'duration', 'of', 'these', 'courses', 'vary', 'between', '10',
', '20', 'and', '30hrs', 'respectively', '.', 'the', 'last', 'date', 'for', 'enrollment', 'is',
july', '18th', '2016.', 'to', 'know', 'more', ',', 'visit', 'https', ':', '//onlinecourses.nptel
ac.in/']
prateek@ubuntu:~$ vi nptel_facebook.py
prateek@ubuntu:~$ python nptel_facebook.py
the 24
to 22
. 16
: 14
for 12
, 10
course 9
nptel 8
https 8
you 7
# 7
on 7
visit 7
more 7
) 7
know 6
? 6
and 6
! 6
```

So, there you see, if you see the word "the" has been used 24 times. Similarly, "to" has occurred 22 times, there is this full stop, which has occurred 16 times; for, course, NPTEL and all of these. Now, you have an idea of what are the most common words NPTEL uses in their posts. So, as evident, the word course and NPTEL is appearing a lot of times, similarly https which means they have posted URLs containing https which are occurring a lot. If you see, there is another problem here, what we are seeing is we are seeing course and courses appear together.

Now, essentially if you notice, course and courses represent the same base word which is course. So, for example, read, reading, reader, read and all of these essentially are derived from the same word which is read; now the process of reaching this base word from any word is called stemming. We can use stemming to eliminate repetitions like this, to avoid course and courses appear separately and instead just consider them as the base word course.

Now, NLTK provides some stemming functions which you can use to convert your words into the base stem word.

(Refer Slide Time: 12:50)



```
import requests
import json
from nltk import word_tokenize
from collections import Counter
from nltk.stem import PorterStemmer

access_token = "EAACEdEose0cBAAameUP4vtMNqF9XGrUXpwP60HYso35dgVkpH6g8I5h4vFHYzx1WbgYNgH2W3090RjZkZBThXZBRuZBDsXyiYWrq6NXTG0cMZBXNRL0qFGxwMveHmM4wqZAK4Ayn2U6m1Z2CoX1hiVzKzCGgAAY1LlmeQ0cGN0YgZDZD"

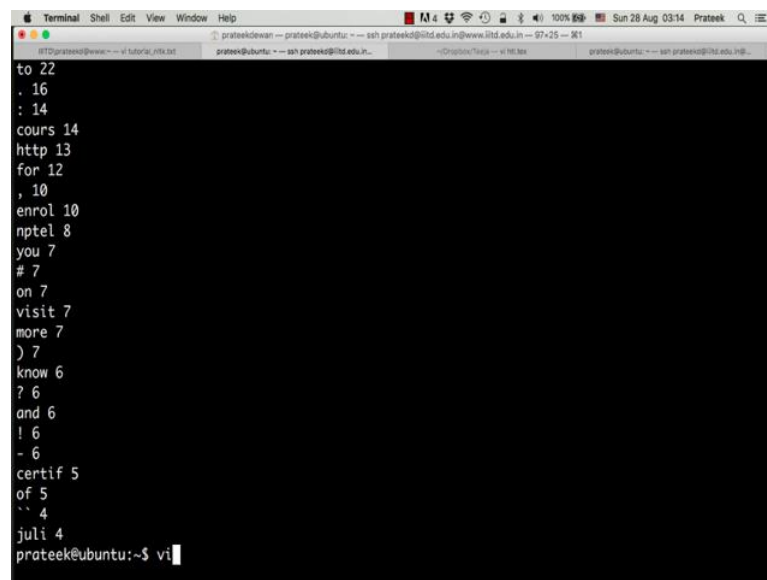
def get_data(page_id):
    stemmer = PorterStemmer()
    all_tokens = []
    url = "https://graph.facebook.com/" + page_id + "/posts?limit=15&access_token="+access_token
    data = requests.get(url)
    json_data = json.loads(data.text)
    for post in json_data['data']:
        if 'message' in post:
            tokens = word_tokenize(post['message'].encode('utf-8').lower())
            all_tokens = all_tokens + tokens
    all_tokens_stemmed = []
    for token in all_tokens:
        stemmed_token = stemmer.stem(token)
        all_tokens_stemmed.append(stemmed_token)
    frequencies = Counter(all_tokens_stemmed)
```

So, let us go back to our script and say from nltk dot stem import porter stemmer there are multiple types of stemmers; the one we will use is called the porter stemmer and we will create a stemmer object by saying stemmer is equal to porter stemmer and we will use this stemmer to stem all the tokens in our list. So, what we will do is let us create

another list called all tokens stemmed as a blank list and for token in all underscore tokens we say stemmed underscore token is equal to stemmer dot stem token.

By this line, what we are doing is we are converting token into **it's** stemmed version and storing it in a new variable called stemmed underscore token; now let us append this stemmed token into this new stemmed list that we created. So, we say all tokens underscore stemmed dot append stemmed underscore token and now for the frequencies, let us calculate the frequencies for the stemmed tokens. Now, this should solve our course and courses problem and consider them as the same base word course.

(Refer Slide Time: 14:26)



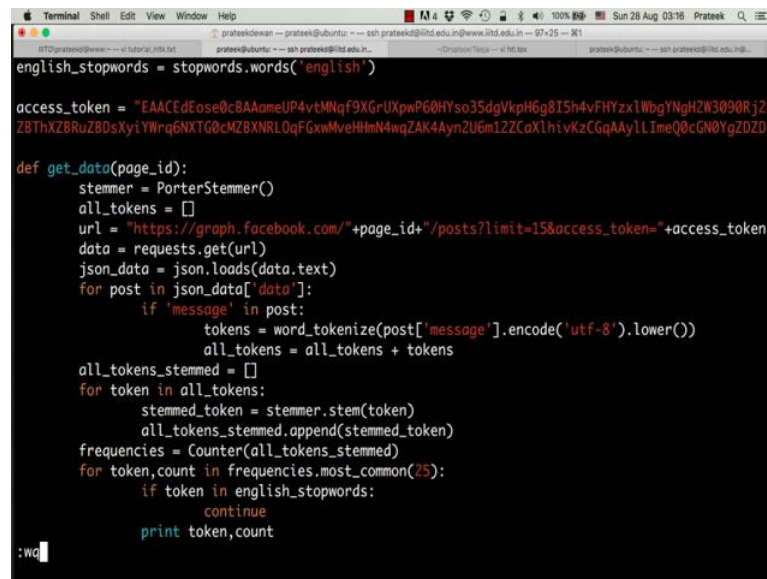
```
to 22
. 16
: 14
cours 14
http 13
for 12
, 10
enrol 10
nptel 8
you 7
# 7
on 7
visit 7
more 7
) 7
know 6
? 6
and 6
! 6
- 6
certif 5
of 5
`` 4
juli 4
prateek@ubuntu:~$ vi
```

So, you see that course and courses have been combined into one word and the frequency is now fourteen; similarly enrollment has changed to enroll and certification and certify have changed to certif., this is the stem word for the words **certify**, certification and certificates. Now, while analyzing text, words that appear commonly in the English language are not of much help for analyzing the content that is under consideration, for example, words like and or the or to or for are words which appear very frequently in the English language naturally. So, one good way to perform better analysis is to find a way to eliminate these words; now in all languages the **most** commonly appearing words in the language are known as stop words.

Now, NLTK provides stop words for the English language which we can use to identify these words and eliminate them from our analysis. So, let us go back to our script and

say from nltk dot corpus import stop words.

(Refer Slide Time: 15:30)

A terminal window with a black background and white text. The code defines a function get_data that fetches posts from a Facebook page, tokenizes the text, stems the tokens, and filters out common stopwords. The code is as follows:

```
english_stopwords = stopwords.words('english')

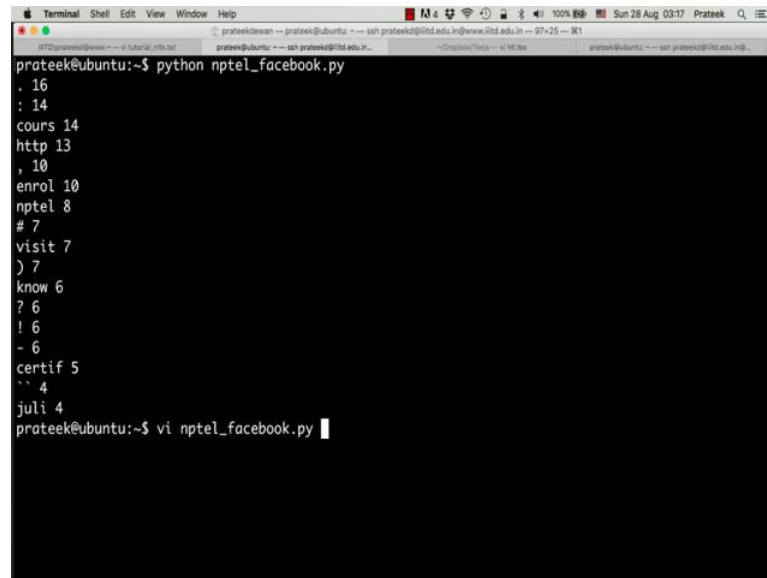
access_token = "EAACEdEose0cBAAomeUP4vtMNqf9XGrUXpwP60HYso35dgVkpH6g8T5h4vFHYzx1WbgYNgHZW3090Rj2k
ZBThXZBRuZBDsxyiYWrq6NXTG0cMZBXNRLQqFGxmVveHhmN4wqZAK4Ayn2U6m1Z2CaX1hivKzCggAAy1LlImeQ0cGN0YgZDZD"

def get_data(page_id):
    stemmer = PorterStemmer()
    all_tokens = []
    url = "https://graph.facebook.com/"+page_id+"/posts?limit=15&access_token="+access_token
    data = requests.get(url)
    json_data = json.loads(data.text)
    for post in json_data['data']:
        if 'message' in post:
            tokens = word_tokenize(post['message'].encode('utf-8').lower())
            all_tokens = all_tokens + tokens
    all_tokens_stemmed = []
    for token in all_tokens:
        stemmed_token = stemmer.stem(token)
        all_tokens_stemmed.append(stemmed_token)
    frequencies = Counter(all_tokens_stemmed)
    for token,count in frequencies.most_common(25):
        if token in english_stopwords:
            continue
        print token,count
```

Now, let us say english underscore stopwords is equal to stopwords dot words english. Now we have a variable called english underscore stopwords which contains english stopwords that we have obtained from the nltk corpus.

Now, while printing our output, let us put an if condition saying if token is present in english underscore stopwords, continue; that is if the token is present in the list of english words, this print statement will not execute and the continue statement instructs the for loop to go to the next iteration without executing the rest of the for loop body. So, essentially whenever the for loop encounters the stopword it will not print it.

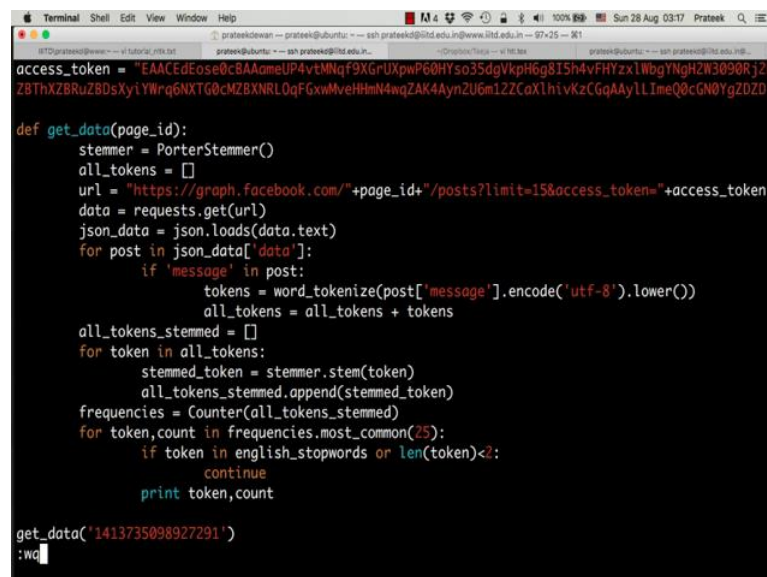
(Refer Slide Time: 16:42)



```
prateek@ubuntu:~$ python nptel_facebook.py
. 16
: 14
cours 14
http 13
, 10
enrol 10
nptel 8
# 7
visit 7
) 7
know 6
? 6
! 6
- 6
certif 5
`` 4
juli 4
prateek@ubuntu:~$ vi nptel_facebook.py
```

So, let us execute this code and see what happens. So, you see the most commonly appearing words like the, and, for and to have disappeared from this list and this gives us a better sense of the content that we are analyzing. To further refine our output we can probably eliminate these punctuation marks and special symbols as well. Let us go back to the script.

(Refer Slide Time: 17:13)



```
access_token = "EAACEdEose0c8AAameUP4vtMNqF9XGrUXpwP0HYso3SdgVkpHGg8I5h4vFHYzx1WbgYNghZWN3090RjzkZBThxZBRuzZ80sXyiYWrqNXTG0cMZBXNRL0qFGxwMveHhmM4nqZAK4AynZU6m1ZZCaX1hi vKzCGqAAy1LlmeQ0cGN0YgZDZD"

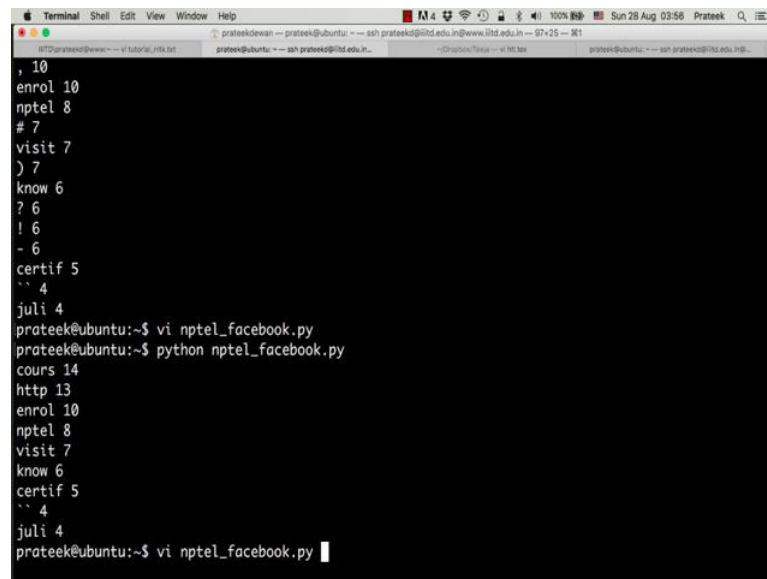
def get_data(page_id):
    stemmer = PorterStemmer()
    all_tokens = []
    url = "https://graph.facebook.com/" + page_id + "/posts?limit=15&access_token=" + access_token
    data = requests.get(url)
    json_data = json.loads(data.text)
    for post in json_data['data']:
        if 'message' in post:
            tokens = word_tokenize(post['message'].encode('utf-8').lower())
            all_tokens = all_tokens + tokens
    all_tokens_stemmed = []
    for token in all_tokens:
        stemmed_token = stemmer.stem(token)
        all_tokens_stemmed.append(stemmed_token)
    frequencies = Counter(all_tokens_stemmed)
    for token, count in frequencies.most_common(25):
        if token in english_stopwords or len(token) < 2:
            continue
        print token, count

get_data('1413735098927291')
:wa
```

Say if token in english stopwords or length of token is less than two characters which means any token which is one character will not appear in the output. So, this will

eliminate full stops, parenthesis, hash and all those kinds of tokens.

(Refer Slide Time: 17:30)



```
prateek@ubuntu:~$ vi nptel_facebook.py
prateek@ubuntu:~$ python nptel_facebook.py
cours 14
http 13
enrol 10
nptel 8
visit 7
know 6
certif 5
` 4
juli 4
prateek@ubuntu:~$ vi nptel_facebook.py
```

Now, let us execute this code and now you see the difference in the output. So, this gives us a much better refined analysis of the content that the NPTEL page is posting. We can also look at bigrams and trigrams and other n grams to analyze this content even better. So, right now we have looked at only singular words that are appearing most commonly. Now in some cases it might be more helpful to know what are the pair of words which are occurring together or what are say 3 or 4 or x number of words which are appearing together most frequently.

For example, names. **If** you consider a full name for example, Barack Obama and you look at the most commonly appearing words in text which contains Barack Obama, we will notice that Barack and Obama will appear separately, but what if you want to know the number of times Barack Obama appeared together or for that matter, any pair of words that have appeared together. **For** that we will use an NLTK package called n grams.

(Refer Slide Time: 18:36)

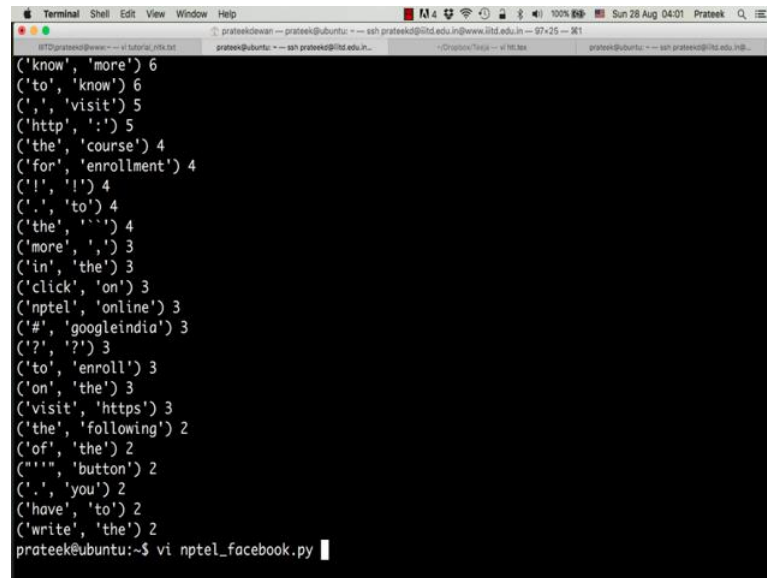
```
Terminal Shell Edit View Window Help
prateek@prateek:~$ cat get_data.py
def get_data(page_id):
    stemmer = PorterStemmer()
    all_tokens = []
    all_bigrams = []
    url = "https://graph.facebook.com/" + page_id + "/posts?limit=15&access_token=" + access_token
    data = requests.get(url)
    json_data = json.loads(data.text)
    for post in json_data['data']:
        if 'message' in post:
            tokens = word_tokenize(post['message'].encode('utf-8').lower())
            bigrams = list(ngrams(tokens, 2))
            all_tokens = all_tokens + tokens
            all_bigrams = all_bigrams + bigrams
    all_tokens_stemmed = []
    for token in all_tokens:
        stemmed_token = stemmer.stem(token)
        all_tokens_stemmed.append(stemmed_token)
    frequencies = Counter(all_bigrams)
    for token, count in frequencies.most_common(25):
        print token, count

get_data('1413735098927291')
```

So, let us go back to the code and say from nltk `import ngrams`. Now, let us create another list called all underscore bigrams where we will append all the bigrams that are appearing. Once, we have the tokens, all we need to do to get bigrams is to say bigrams is equal to `ngrams` tokens comma two this line will calculate bigrams and store them in a list called bigrams; to obtain trigrams we can simply change this number two to three and the `ngrams` function will then return trigrams instead of bigrams. We also need to convert the output of this `ngrams` function into a list. So, we say list `ngrams` tokens comma 2.

Now, as we did for tokens, we also need to append these bigrams to the all bigrams list. So, we say all underscore bigrams is equal to all underscore bigrams plus bigrams and now we directly go to the counter function to count the frequency of all underscore bigrams and we can get rid of this stopwords and token length condition.

(Refer Slide Time: 20:36)

A terminal window showing the output of a script named 'npTEL_facebook.py'. The output lists various bigrams and their frequencies. The most frequent bigram is ('know', 'more') with a count of 6. Other frequent bigrams include ('to', 'know') with 6, ('.', 'visit') with 5, ('http', ':') with 5, ('the', 'course') with 4, ('for', 'enrollment') with 4, ('!', '!') with 4, ('.', 'to') with 4, ('the', '') with 4, ('more', ',') with 3, ('in', 'the') with 3, ('click', 'on') with 3, ('npTEL', 'online') with 3, ('#', 'googleindia') with 3, ('?', '?') with 3, ('to', 'enroll') with 3, ('on', 'the') with 3, ('visit', 'https') with 3, ('the', 'following') with 2, ('of', 'the') with 2, ('"', 'button') with 2, ('.', 'you') with 2, ('have', 'to') with 2, and ('write', 'the') with 2. The terminal prompt shows the user is in the directory 'prateek@ubuntu:~\$' and has just run 'vi npTEL_facebook.py'.

```
prateek@ubuntu:~$ vi npTEL_facebook.py
('know', 'more') 6
('to', 'know') 6
('.', 'visit') 5
('http', ':') 5
('the', 'course') 4
('for', 'enrollment') 4
('!', '!') 4
('.', 'to') 4
('the', '') 4
('more', ',') 3
('in', 'the') 3
('click', 'on') 3
('npTEL', 'online') 3
('#', 'googleindia') 3
('?', '?') 3
('to', 'enroll') 3
('on', 'the') 3
('visit', 'https') 3
('the', 'following') 2
('of', 'the') 2
('"', 'button') 2
('.', 'you') 2
('have', 'to') 2
('write', 'the') 2
prateek@ubuntu:~$ vi npTEL_facebook.py
```

There you go, you see that https with a colon is the most frequently appearing bigram in this data set followed by **know more** or **to know** and **so on**. If you compare this output with the output of the uni grams, you will see considerable differences. **The** bigram output clearly tells us that NPTEL posts a lot of posts containing phrases like **to know more**, **visit URL** or 'for enrollment' or **npTEL online** or **hashtag googleindia** and so on. In the unigrams output, we were not able to identify such phrases.

Now, let see what are the most commonly appearing trigrams. So, let us go back to the code and just replace this 2 by 3 and that is all we need to do.

(Refer Slide Time: 21:03)

```
Terminal Shell Edit View Window Help
prateek@prateek:~$ cat nptel_facebook.py
import json
from nltk import word_tokenize
from collections import Counter
from nltk.stem import PorterStemmer
from nltk.corpus import stopwords
english_stopwords = stopwords.words('english')
from nltk import ngrams

access_token = "EAACEdEose0cBAAameUP4vtMNqf9XGrUXpwP60HYso35dgVkpH6g8I5h4vFHYzx1WbgYNgH2W3090Rj2kZBThXZBRuZBDsXyiYWrq6NXTG0cMZBXNRLQqF6xwMveHmM4wqZAK4Ayn2U6m12ZCaXlHivKzCGgAAyLLImeQ0cGN0YgZDZD"

def get_data(page_id):
    stemmer = PorterStemmer()
    all_tokens = []
    all_bigrams = []
    url = "https://graph.facebook.com/"+page_id+"/posts?limit=15&access_token="+access_token
    data = requests.get(url)
    json_data = json.loads(data.text)
    for post in json_data['data']:
        if 'message' in post:
            tokens = word_tokenize(post['message'].encode('utf-8').lower())
            bigrams = list(ngrams(tokens,2))
            all_tokens = all_tokens + tokens
            all_bigrams = all_bigrams + bigrams

:wa
```

Run this code again and you see the most commonly appearing trigrams.

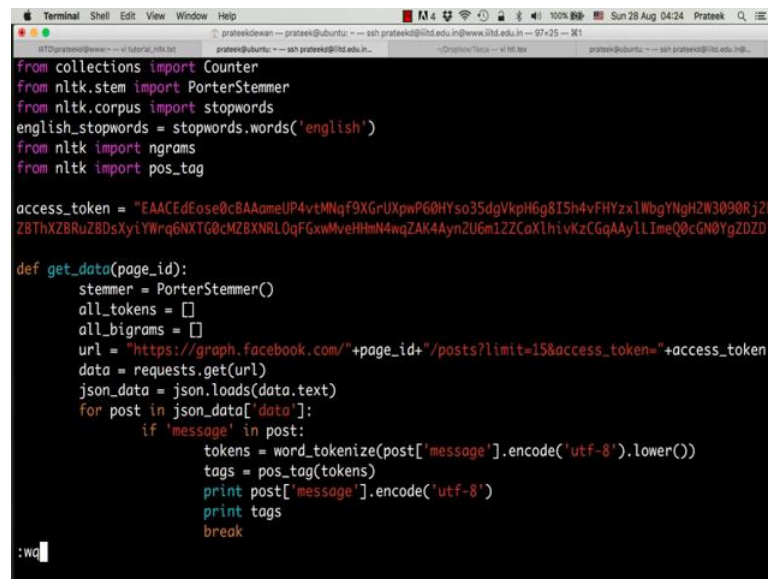
(Refer Slide Time: 21:08)

```
Terminal Shell Edit View Window Help
prateek@prateek:~$ python nptel_facebook.py
('more', ',', 'visit') 3
(('.', 'to', 'know') 3
('visit', 'https', ':') 3
('know', 'more', ',') 3
('click', 'on', 'the') 3
((',', 'visit', 'https') 3
('date', 'for', 'enrollment') 2
((',', 'visit', 'http') 2
('last', 'date', 'for') 2
('https', ':', '//onlinecourses.nptel.ac.in/explorer') 2
('https', ':', '//onlinecourses.nptel.ac.in/') 2
('on', 'the', '') 2
('for', 'the', 'course') 2
('to', 'enroll', 'for') 2
('the', 'following', 'link') 2
('have', 'to', 'pay') 2
('open', 'for', 'enrollment') 2
('nptel', 'online', 'courses') 2
('!', '!', '!') 2
('join', '"', 'button') 2
('', 'join', '') 2
('visit', 'http', ':') 2
('the', '', 'join') 2
('for', 'enrollment', '.') 2
prateek@ubuntu:~$ vi nptel_facebook.py
```

So, like we noticed - to know more comma visit and visit http click on the date for enrollment and other similar phrases which are appearing very frequently in NPTEL posts. Another interesting aspect of social media content that can be studied is the part of speech tags or the sentence structures that people use while they are posting online, for example, it might be very helpful in some cases to identify nouns, pronouns, verbs, adverbs and all such words in a sentence automatically. So, let us see how to do that

using nltk.

(Refer Slide Time: 21:58)



```
Terminal Shell Edit View Window Help
prateek@prateek: ~ -- ssh prateek@iitd.edu.in -- ssh prateek@iitd.edu.in -- ssh prateek@iitd.edu.in -- 97x25 -- 31
prateek@prateek: ~ -- ssh prateek@iitd.edu.in -- ssh prateek@iitd.edu.in -- ssh prateek@iitd.edu.in -- 97x25 -- 31
prateek@prateek: ~ -- ssh prateek@iitd.edu.in -- ssh prateek@iitd.edu.in -- ssh prateek@iitd.edu.in -- 97x25 -- 31

from collections import Counter
from nltk.stem import PorterStemmer
from nltk.corpus import stopwords
english_stopwords = stopwords.words('english')
from nltk import ngrams
from nltk import pos_tag

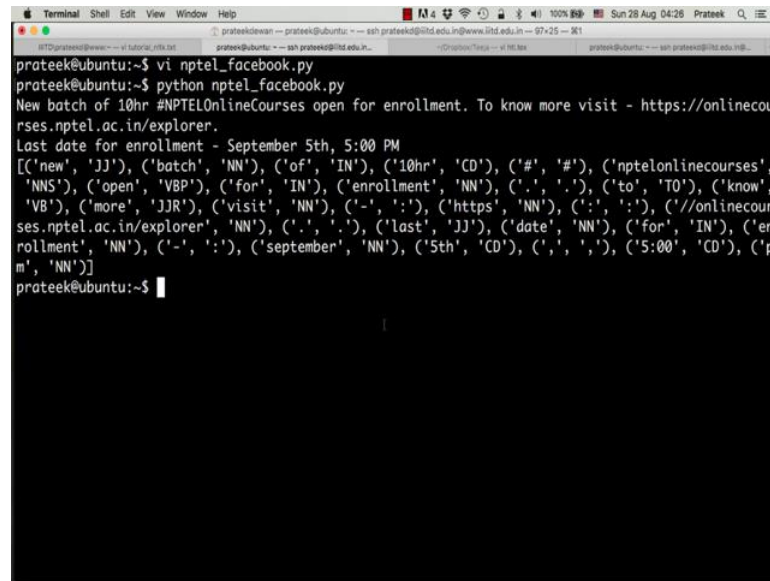
access_token = "EAACEdEose0cBAAAameUP4vtMNqf9XGrUXpwP60HYso35dgVkpH6g8I5h4vFHYzx1WbgYNgHZW3090RjZk
ZBThXZBRuZBDsXyiYWq6NXTG0cMZBXNRLQqFGxwMveHmM4wqZAK4Ayn2U6m1ZZCoX1hivKzCGqAAYlLImeQ0cGN0YgZDZD"

def get_data(page_id):
    stemmer = PorterStemmer()
    all_tokens = []
    all_bigrams = []
    url = "https://graph.facebook.com/" + page_id + "/posts?limit=15&access_token=" + access_token
    data = requests.get(url)
    json_data = json.loads(data.text)
    for post in json_data['data']:
        if 'message' in post:
            tokens = word_tokenize(post['message'].encode('utf-8').lower())
            tags = pos_tag(tokens)
            print post['message'].encode('utf-8')
            print tags
            break
:wa
```

Go back to the script and say from nltk import pos underscore tag. Pos underscore tag stands for part of speech tagging. **This** package is capable of **parsing** a sentence as tokens and marking each token with a part of speech associated with it.

So, let us see how this works, as an example we will look at the part of speech tags for only one post to make it easier for us to understand. So, after we have obtained the tokens, we say tags is equal to pos underscore tag, tokens, print tags and remove the remaining part of the code. **In** addition to tags, let us also print the message to avoid confusion to see exactly what the post is and what are the tags associated with each word in the post. So, we say print post message dot encode UTF-8 and let us break this loop here, so that we can focus on only one post for the output.

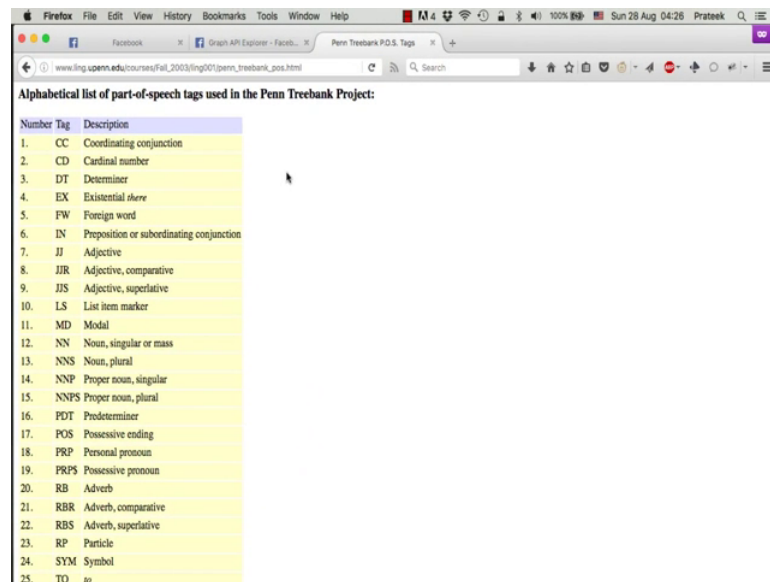
(Refer Slide Time: 23:22)



```
prateek@ubuntu:~$ vi npTEL_facebook.py
prateek@ubuntu:~$ python npTEL_facebook.py
New batch of 10hr #NPTELonlineCourses open for enrollment. To know more visit - https://onlinecou
rses.nptel.ac.in/explorer.
Last date for enrollment - September 5th, 5:00 PM
[('new', 'JJ'), ('batch', 'NN'), ('of', 'IN'), ('10hr', 'CD'), ('#', '#'), ('nptelonlinecourses',
'NNS'), ('open', 'VBP'), ('for', 'IN'), ('enrollment', 'NN'), ('.', '.'), ('to', 'TO'), ('know',
'VB'), ('more', 'JJR'), ('visit', 'NN'), ('-', '-'), ('https', 'NN'), (':', ':'), ('//onlinecour
ses.nptel.ac.in/explorer', 'NN'), ('.', '.'), ('last', 'JJ'), ('date', 'NN'), ('for', 'IN'), ('en
rollment', 'NN'), ('-', '-'), ('september', 'NN'), ('5th', 'CD'), (',', ','), ('5:00', 'CD'), ('p
m', 'NN')]
prateek@ubuntu:~$
```

So, as you see here the message reads new batch of ten our NPTEL online courses open for enrollment. Now the part of speech tagger has taken each word into consideration and assigned a part of speech tag for each of these words. The explanation for each of these tags is available online at this URL.

(Refer Slide Time: 23:47)

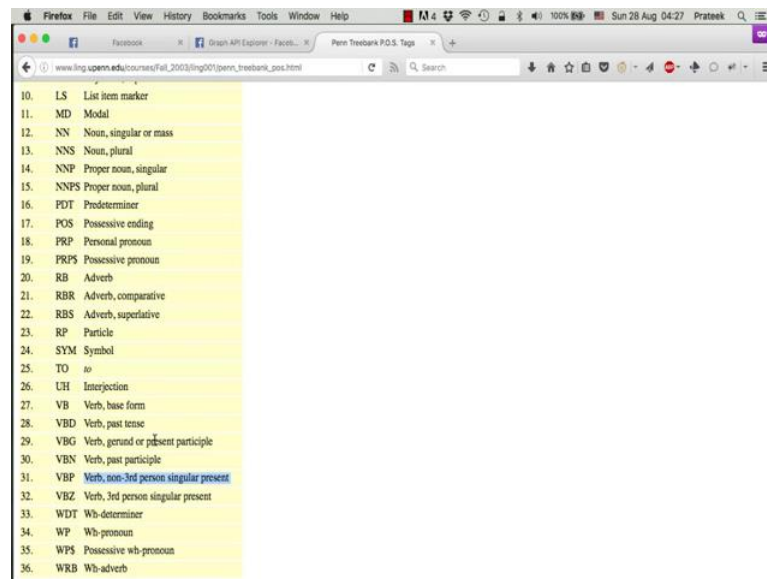


Number	Tag	Description
1.	CC	Coordinating conjunction
2.	CD	Cardinal number
3.	DT	Determiner
4.	EX	Existential <i>there</i>
5.	FW	Foreign word
6.	IN	Preposition or subordinating conjunction
7.	JJ	Adjective
8.	JJR	Adjective, comparative
9.	JJS	Adjective, superlative
10.	LS	List item marker
11.	MD	Modal
12.	NN	Noun, singular or mass
13.	NNS	Noun, plural
14.	NNP	Proper noun, singular
15.	NNPS	Proper noun, plural
16.	PDT	Predeterminer
17.	POS	Possessive ending
18.	PRP	Personal pronoun
19.	PRPS	Possessive pronoun
20.	RB	Adverb
21.	RBR	Adverb, comparative
22.	RBS	Adverb, superlative
23.	RP	Particle
24.	SYM	Symbol
25.	TO	<i>to</i>

So if you notice, the first word is marked with the tag jj, going back to the lookup table, jj refers to an adjective, which means new is an adjective. Similarly, batch is a noun. of is a preposition. 10 hour is a cardinal number. open is a non third person singular present

verb and so on.

(Refer Slide Time: 24:30)



Part of speech tags can be very helpful in understanding the structure of a sentence. They can also be used to identify whether a sentence follows correct grammar or not. The sentence structure and grammar can be used as a feature to differentiate between different kinds of online social media users. **That** is all for this tutorial.

Let us just go back and revise what all we covered in this tutorial. In terms of nltk, we looked at how to tokenize words and convert sentences into a list of words. We also saw how to calculate the frequency of these words, then we looked at stemming which converts words into their stem words. Then we looked at the nltk corpus for english stopwords and also looked at how we can calculate bigrams and trigrams using the ngrams package. Finally, we looked at part of speech tagging available as part of the natural language toolkit in python.