

# LING/C SC 581:

## Advanced Computational Linguistics

Lecture 12

# Today's Topics

*Did everyone manage to install the structural probes program?*

- Using the structural probes program
- Case study: variety of passives
- Homework 6
  - due next Monday midnight

ptb

```
>>> import nltk
>>> from nltk.corpus import ptb
>>> ptb.fileids()
```

BROWN/CF/CF01.MRG	BROWN/CF/CF02.MRG	BROWN/CF/CF03.MRG	BROWN/CF/CF04.MRG
BROWN/CF/CF05.MRG	BROWN/CF/CF06.MRG	BROWN/CF/CF07.MRG	BROWN/CF/CF08.MRG
BROWN/CF/CF09.MRG	BROWN/CF/CF10.MRG	BROWN/CF/CF11.MRG	BROWN/CF/CF12.MRG
BROWN/CF/CF13.MRG	BROWN/CF/CF14.MRG	BROWN/CF/CF15.MRG	BROWN/CF/CF16.MRG
BROWN/CF/CF17.MRG	BROWN/CF/CF18.MRG	BROWN/CF/CF19.MRG	BROWN/CF/CF20.MRG
BROWN/CF/CF21.MRG	BROWN/CF/CF22.MRG	BROWN/CF/CF23.MRG	BROWN/CF/CF24.MRG
BROWN/CF/CF25.MRG	BROWN/CF/CF26.MRG	BROWN/CF/CF27.MRG	BROWN/CF/CF28.MRG
BROWN/CF/CF29.MRG	BROWN/CF/CF30.MRG	BROWN/CF/CF31.MRG	BROWN/CF/CF32.MRG
BROWN/CG/C001.MRG	BROWN/CG/C002.MRG	BROWN/CG/C003.MRG	BROWN/CG/C004.MRG
BROWN/CG/C005.MRG	BROWN/CG/C006.MRG	BROWN/CG/C007.MRG	BROWN/CG/C008.MRG
BROWN/CG/C009.MRG	BROWN/CG/C010.MRG	BROWN/CG/C011.MRG	BROWN/CG/C012.MRG
BROWN/CG/C013.MRG	BROWN/CG/C014.MRG	BROWN/CG/C015.MRG	BROWN/CG/C016.MRG
BROWN/CG/C017.MRG	BROWN/CG/C018.MRG	BROWN/CG/C019.MRG	BROWN/CG/C020.MRG
BROWN/CG/C021.MRG	BROWN/CG/C022.MRG	BROWN/CG/C023.MRG	BROWN/CG/C024.MRG
BROWN/CG/C025.MRG	BROWN/CG/C026.MRG	BROWN/CG/C027.MRG	BROWN/CG/C028.MRG
BROWN/CG/C029.MRG	BROWN/CG/C030.MRG	BROWN/CG/C031.MRG	BROWN/CG/C032.MRG
BROWN/CK/K001.MRG	BROWN/CK/K002.MRG	BROWN/CK/K003.MRG	BROWN/CK/K004.MRG
BROWN/CK/K005.MRG	BROWN/CK/K006.MRG	BROWN/CK/K007.MRG	BROWN/CK/K008.MRG
BROWN/CK/K009.MRG	BROWN/CK/K010.MRG	BROWN/CK/K011.MRG	BROWN/CK/K012.MRG
BROWN/CK/K013.MRG	BROWN/CK/K014.MRG	BROWN/CK/K015.MRG	BROWN/CK/K016.MRG
BROWN/CK/K017.MRG	BROWN/CK/K018.MRG	BROWN/CK/K019.MRG	BROWN/CK/K020.MRG
BROWN/CL/C001.MRG	BROWN/CL/C002.MRG	BROWN/CL/C003.MRG	BROWN/CL/C004.MRG
BROWN/CL/C005.MRG	BROWN/CL/C006.MRG	BROWN/CL/C007.MRG	BROWN/CL/C008.MRG
BROWN/CL/C009.MRG	BROWN/CL/C010.MRG	BROWN/CL/C011.MRG	BROWN/CL/C012.MRG
BROWN/CL/C013.MRG	BROWN/CL/C014.MRG	BROWN/CL/C015.MRG	BROWN/CL/C016.MRG
BROWN/CL/C017.MRG	BROWN/CL/C018.MRG	BROWN/CL/C019.MRG	BROWN/CL/C020.MRG
BROWN/CN/C001.MRG	BROWN/CN/C002.MRG	BROWN/CN/C003.MRG	BROWN/CN/C004.MRG
BROWN/CN/C005.MRG	BROWN/CN/C006.MRG	BROWN/CN/C007.MRG	BROWN/CN/C008.MRG
BROWN/CN/C009.MRG	BROWN/CN/C010.MRG	BROWN/CN/C011.MRG	BROWN/CN/C012.MRG
BROWN/CN/C013.MRG	BROWN/CN/C014.MRG	BROWN/CN/C015.MRG	BROWN/CN/C016.MRG
BROWN/CN/C017.MRG	BROWN/CN/C018.MRG	BROWN/CN/C019.MRG	BROWN/CN/C020.MRG
BROWN/CN/C021.MRG	BROWN/CN/C022.MRG	BROWN/CN/C023.MRG	BROWN/CN/C024.MRG
BROWN/CN/C025.MRG	BROWN/CN/C026.MRG	BROWN/CN/C027.MRG	BROWN/CN/C028.MRG
BROWN/CN/C029.MRG	BROWN/CN/C030.MRG	BROWN/CN/C031.MRG	BROWN/CN/C032.MRG
BROWN/CP/CP001.MRG	BROWN/CP/CP002.MRG	BROWN/CP/CP003.MRG	BROWN/CP/CP004.MRG
BROWN/CP/CP005.MRG	BROWN/CP/CP006.MRG	BROWN/CP/CP007.MRG	BROWN/CP/CP008.MRG
BROWN/CP/CP009.MRG	BROWN/CP/CP010.MRG	BROWN/CP/CP011.MRG	BROWN/CP/CP012.MRG
BROWN/CP/CP013.MRG	BROWN/CP/CP014.MRG	BROWN/CP/CP015.MRG	BROWN/CP/CP016.MRG
BROWN/CP/CP017.MRG	BROWN/CP/CP018.MRG	BROWN/CP/CP019.MRG	BROWN/CP/CP020.MRG
BROWN/CP/CP021.MRG	BROWN/CP/CP022.MRG	BROWN/CP/CP023.MRG	BROWN/CP/CP024.MRG
BROWN/CR/CR001.MRG	BROWN/CR/CR002.MRG	BROWN/CR/CR003.MRG	BROWN/CR/CR004.MRG
BROWN/CR/CR005.MRG	BROWN/CR/CR006.MRG	BROWN/CR/CR007.MRG	BROWN/CR/CR008.MRG

■■■

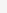
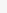

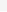
192

```
WSJ_00/WSJ_0004.MRG WSJ_00/WSJ_0005.MRG WSJ_00/WSJ_0006.MRG WSJ_00/WSJ_0007.MRG
WSJ_00/WSJ_0008.MRG WSJ_00/WSJ_0009.MRG WSJ_00/WSJ_0010.MRG WSJ_00/WSJ_0011.MRG
WSJ_00/WSJ_0012.MRG WSJ_00/WSJ_0013.MRG WSJ_00/WSJ_0014.MRG WSJ_00/WSJ_0015.MRG
WSJ_00/WSJ_0016.MRG WSJ_00/WSJ_0017.MRG WSJ_00/WSJ_0018.MRG WSJ_00/WSJ_0019.MRG
WSJ_00/WSJ_0020.MRG WSJ_00/WSJ_0021.MRG WSJ_00/WSJ_0022.MRG WSJ_00/WSJ_0023.MRG
WSJ_00/WSJ_0024.MRG WSJ_00/WSJ_0025.MRG WSJ_00/WSJ_0026.MRG WSJ_00/WSJ_0027.MRG
WSJ_00/WSJ_0032.MRG WSJ_00/WSJ_0033.MRG WSJ_00/WSJ_0034.MRG WSJ_00/WSJ_0035.MRG
WSJ_00/WSJ_0036.MRG WSJ_00/WSJ_0037.MRG WSJ_00/WSJ_0038.MRG
```

■■■

WSJ/23/WSJ_2398.MRG	WSJ/23/WSJ_2399.MRG	WSJ/24/WSJ_2400.MRG	WSJ/24/WSJ_2401.MRG
WSJ/24/WSJ_2402.MRG	WSJ/24/WSJ_2403.MRG	WSJ/24/WSJ_2404.MRG	WSJ/24/WSJ_2405.MRG
WSJ/24/WSJ_2406.MRG	WSJ/24/WSJ_2407.MRG	WSJ/24/WSJ_2408.MRG	WSJ/24/WSJ_2409.MRG
WSJ/24/WSJ_2410.MRG	WSJ/24/WSJ_2411.MRG	WSJ/24/WSJ_2412.MRG	WSJ/24/WSJ_2413.MRG
WSJ/24/WSJ_2414.MRG	WSJ/24/WSJ_2415.MRG	WSJ/24/WSJ_2416.MRG	WSJ/24/WSJ_2417.MRG
WSJ/24/WSJ_2418.MRG	WSJ/24/WSJ_2419.MRG	WSJ/24/WSJ_2420.MRG	WSJ/24/WSJ_2421.MRG
WSJ/24/WSJ_2422.MRG	WSJ/24/WSJ_2423.MRG	WSJ/24/WSJ_2424.MRG	WSJ/24/WSJ_2425.MRG
WSJ/24/WSJ_2426.MRG	WSJ/24/WSJ_2427.MRG	WSJ/24/WSJ_2428.MRG	WSJ/24/WSJ_2429.MRG
WSJ/24/WSJ_2430.MRG	WSJ/24/WSJ_2431.MRG	WSJ/24/WSJ_2432.MRG	WSJ/24/WSJ_2433.MRG
WSJ/24/WSJ_2434.MRG	WSJ/24/WSJ_2435.MRG	WSJ/24/WSJ_2436.MRG	WSJ/24/WSJ_2437.MRG
WSJ/24/WSJ_2438.MRG	WSJ/24/WSJ_2439.MRG	WSJ/24/WSJ_2440.MRG	WSJ/24/WSJ_2441.MRG
WSJ/24/WSJ_2442.MRG	WSJ/24/WSJ_2443.MRG	WSJ/24/WSJ_2444.MRG	WSJ/24/WSJ_2445.MRG
WSJ/24/WSJ_2446.MRG	WSJ/24/WSJ_2447.MRG	WSJ/24/WSJ_2448.MRG	WSJ/24/WSJ_2449.MRG
WSJ/24/WSJ_2450.MRG	WSJ/24/WSJ_2451.MRG	WSJ/24/WSJ_2452.MRG	WSJ/24/WSJ_2453.MRG
WSJ/24/WSJ_2454.MRG			

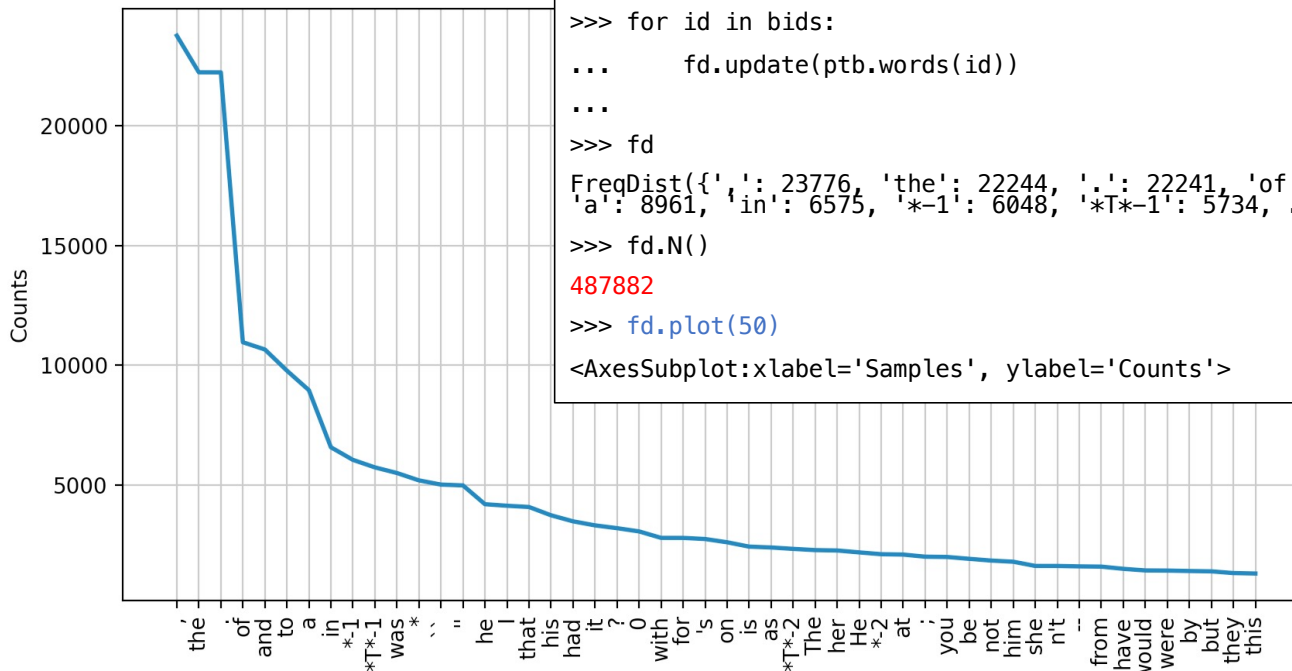
2312

- ✓  BROWN
  - >  CF
  - >  CG
  - >  CK
  - >  CL
  - >  CM
  - >  CN
  - >  CP
  - >  CR

- WSJ
  - 00
  - 01
  - 02
  - 03
  - 04
  - 05
  - 06
  - 07
  - 08
  - 09
  - 10
  - 11
  - 12
  - 13
  - 14
  - 15
  - 16
  - 17
  - 18
  - 19
  - 20
  - 21
  - 22
  - 23
  - 24

# ptb

Brown corpus file ids

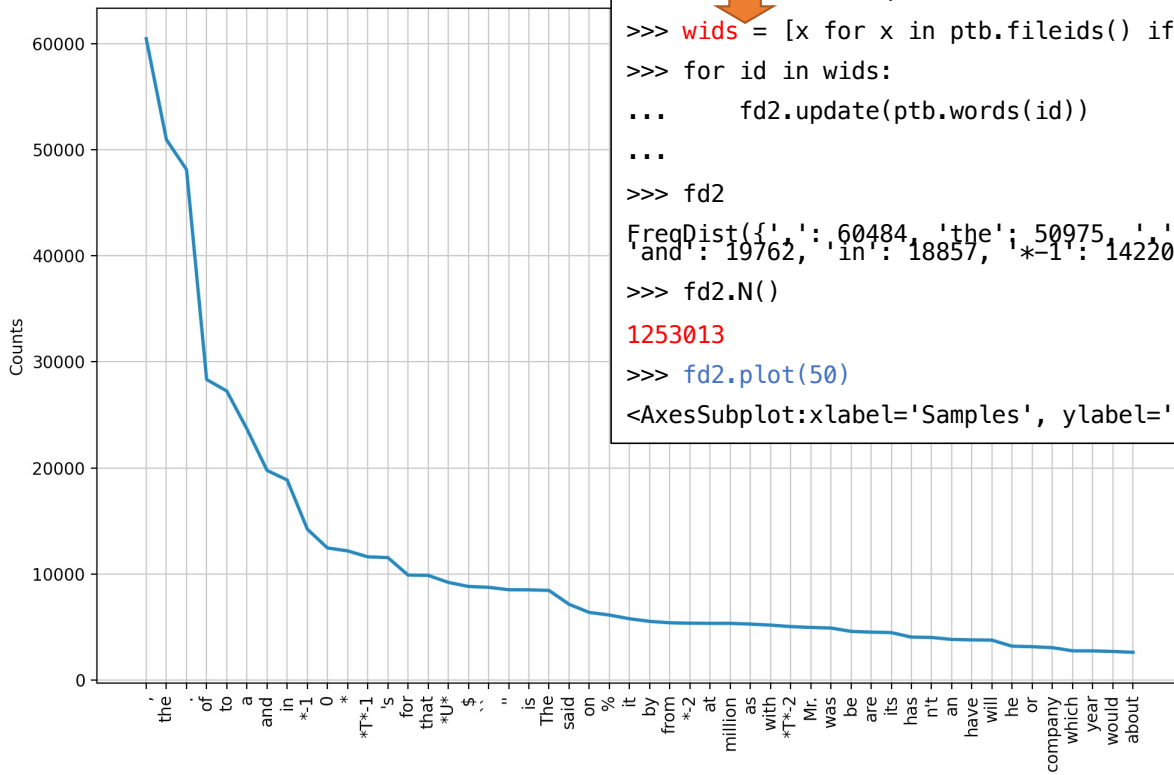


```
>>> bids = [x for x in ptb.fileids() if x.startswith('BROWN')]
>>> len(bids)
192
>>> fd = nltk.FreqDist()
>>> for id in bids:
...     fd.update(ptb.words(id))
...
>>> fd
FreqDist({' ': 23776, 'the': 22244, '': 22241, 'of': 10964, 'and': 10661, 'to': 9778,
'a': 8961, 'in': 6575, '*-1': 6048, '*T*-1': 5734, ...})
>>> fd.N()
487882
>>> fd.plot(50)
<AxesSubplot:xlabel='Samples', ylabel='Counts'>
```

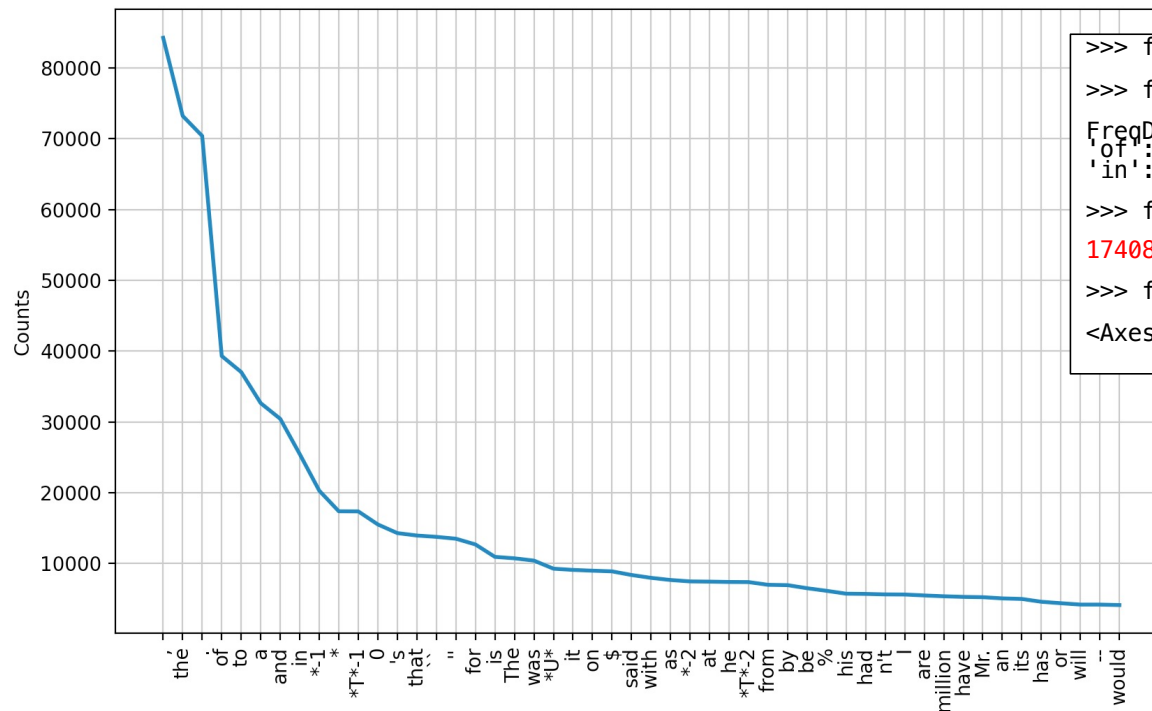
ptb

WSJ corpus file ids

```
>>> fd2 = nltk.FreqDist()
>>> wids = [x for x in ptb.fileids() if x.startswith('WSJ')]
>>> for id in wids:
...     fd2.update(ptb.words(id))
...
>>> fd2
FreqDist({'': 60484, 'the': 50975, '': 48144, 'of': 28338, 'to': 27249, 'a': 23673,
'and': 19762, 'in': 18857, '*-1': 14220, '0': 12447, ...})
>>> fd2.N()
1253013
>>> fd2.plot(50)
<AxesSubplot:xlabel='Samples', ylabel='Counts'>
```



# ptb



```
>>> fd.update(fd2)
>>> fd
FreqDist({'': 84260, 'the': 73219, '': 70385,
'of': 39302, 'to': 37027, 'a': 32634, 'and': 30423,
'in': 25432, '*-1': 20268, '*': 17363, ...})
>>> fd.N()
1740895
>>> fd.plot(50)
<AxesSubplot:xlabel='Samples', ylabel='Counts'>
```

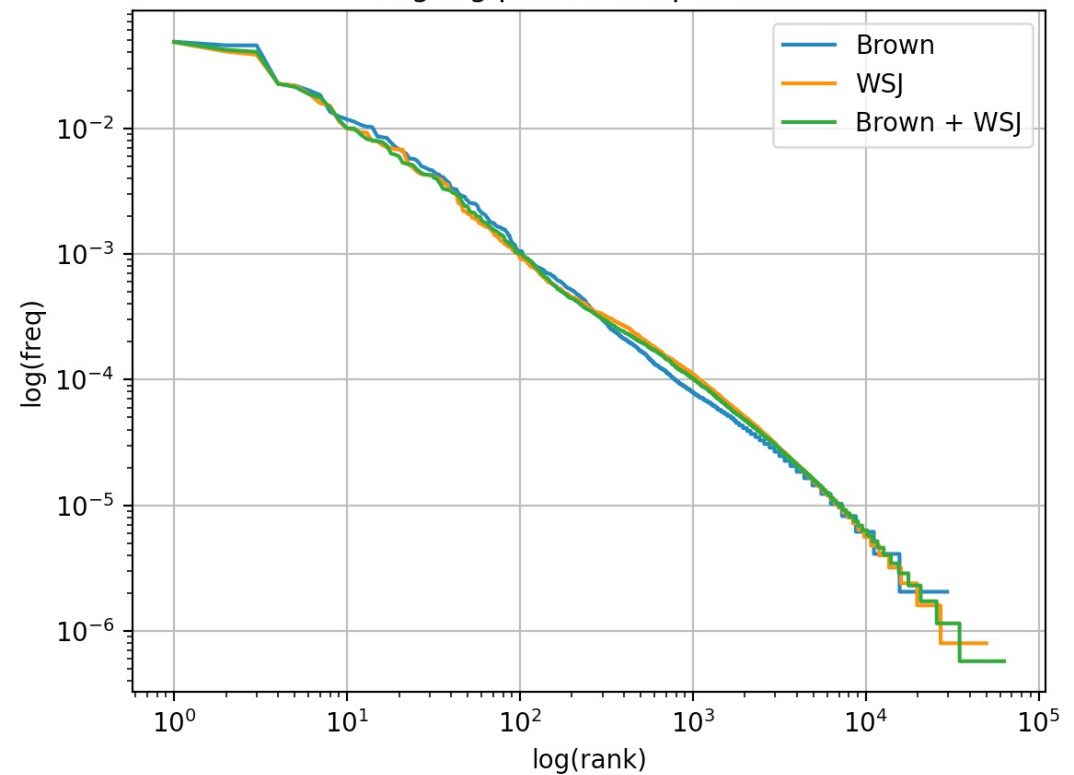
# ptb

```
>>> wws = []
>>> for id in wids:
...     wws.extend(ptb.words(id))
...
>>> bws = []
>>> for id in bids:
...     bws.extend(ptb.words(id))
...
>>> len(bws)
487882
>>> len(wws)
1253013
>>> from zipfile import *
>>> fig()
>>> plot(bws, "Brown")
>>> plot(wws, "WSJ")
>>> plot(bws+wws, "Brown + WSJ")
>>> plt.legend()
<matplotlib.legend.Legend object at 0x1278c4520>
>>> plt.show()
```

WSJ words

Brown words

Log-log plot for freq vs rank



# ptb

On course website: `zipf.py`

```
1# Sandiway Fong (c) University of Arizona 2019
2# simple function to plot Zipf's Law
3# assumes matplotlib
4from collections import Counter
5from math import log, log10
6import matplotlib.pyplot as plt
7
8def plot(tokens, text):
9    size = len(tokens)
10    c = Counter()
11    for token in tokens:
12        c[token] += 1
13    mc = c.most_common()
14    ranks = [x for x in range(1, len(mc)+1)]
15    freq = [item[1]/size for item in mc]
16    plt.plot(ranks, freq, label=text)
```

```
17
18def fig():
19    plt.figure(1)
20    plt.xscale('log')
21    plt.xlabel('log(rank)')
22    plt.yscale('log')
23    plt.ylabel('log(freq)')
24    plt.grid(True)
25    plt.title('Log-log plot for freq vs rank')
```

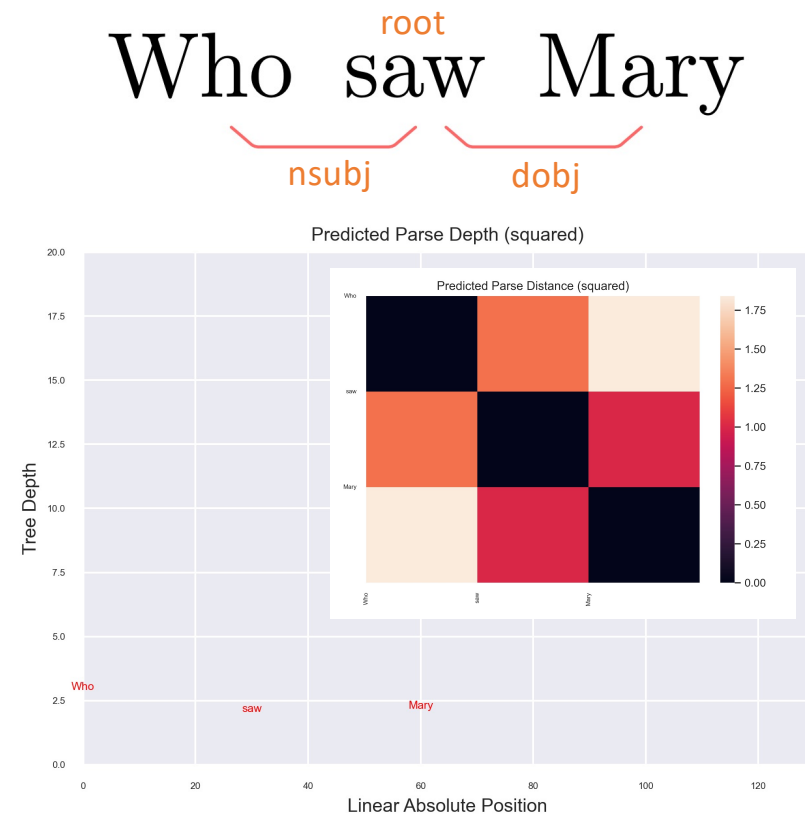


# structural probes

Run a pretrained structural probe on BERT-large quickly on the command line.

It's easy to get predictions on a sentence (or file of sentences) using our demo script and the pre-trained structural probes we release. We use `pytorch-pretrained-bert` to get BERT subword embeddings for each sentence; it should be installed during setup of the repository.

Make sure you've run `download_example.sh`; this will download two probe parameter files to `example/data/`. Also make sure you've installed all dependencies. One is a distance probe on the 16th hidden layer of BERT large, and the other is a depth probe on the same layer. The configuration file `example/demo-bert.yaml` has the right paths already plugged in; just pipe text into the demo file, as follows:



# Testing a range of sentences

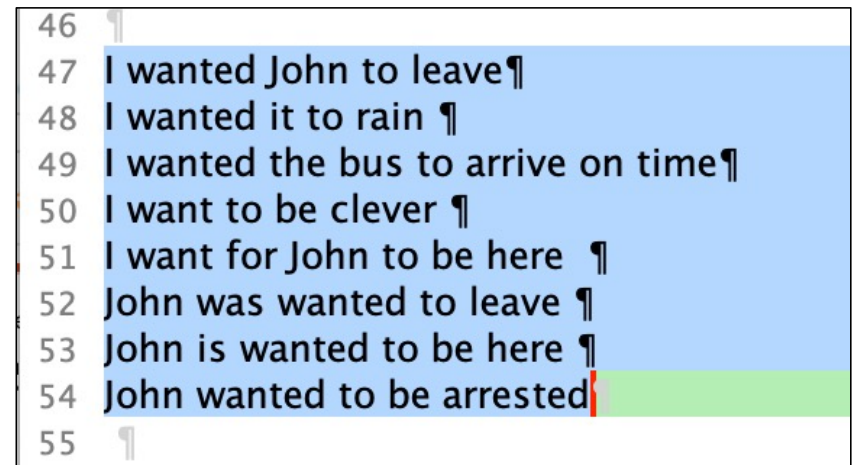
- Examples:
  - assuming you are in the (1<sup>st</sup> level) structural-probes directory on the Terminal:
    - there are further subdirectories structural-probes, example, scripts, doc-assets
    - run\_demo.py is inside structural-probes/structural-probes
    - demo-bert.yaml is inside structural-probes/example
  - 1. `echo "the bombing of the cities is a crime\nthe bombing of the cities are a crime" | python3 structural-probes/run_demo.py example/demo-bert.yaml`
  - 2. `sed -n '207,221p' l\&u1.txt | python3 structural-probes/run_demo.py example/demo-bert.yaml`

# Linux and Mac: sed selection

- Example:

```
sed -n '47,54p' l\u1.txt
```

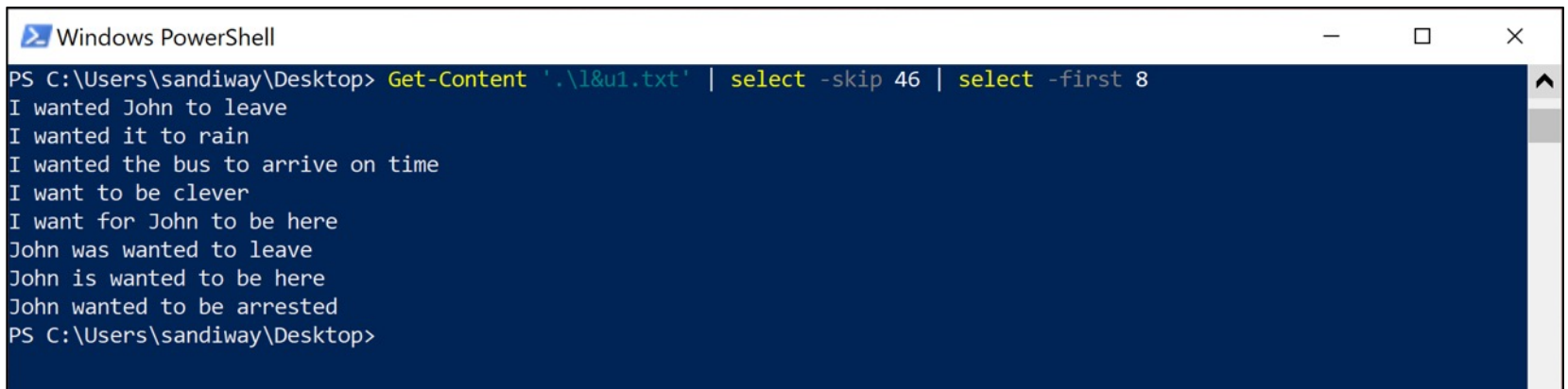
```
I wanted John to leave  
I wanted it to rain  
I wanted the bus to arrive on time  
I want to be clever  
I want for John to be here  
John was wanted to leave  
John is wanted to be here  
John wanted to be arrested
```



```
46 ¶  
47 I wanted John to leave¶  
48 I wanted it to rain ¶  
49 I wanted the bus to arrive on time¶  
50 I want to be clever ¶  
51 I want for John to be here ¶  
52 John was wanted to leave ¶  
53 John is wanted to be here ¶  
54 John wanted to be arrested¶  
55 ¶
```

# Windows 10: **sed** selection functionality

- Powershell cmdlets:
  - `Get-Content filename | Select -Skip n`
  - `Select -First n`
  - `Select -Last n`

A screenshot of a Windows PowerShell window. The title bar reads "Windows PowerShell". The command prompt shows the command: `PS C:\Users\sandiway\Desktop> Get-Content '.\l&u1.txt' | select -skip 46 | select -first 8`. The output consists of eight lines of text: "I wanted John to leave", "I wanted it to rain", "I wanted the bus to arrive on time", "I want to be clever", "I want for John to be here", "John was wanted to leave", "John is wanted to be here", and "John wanted to be arrested". The prompt then returns to `PS C:\Users\sandiway\Desktop>`.

```
Windows PowerShell
PS C:\Users\sandiway\Desktop> Get-Content '.\l&u1.txt' | select -skip 46 | select -first 8
I wanted John to leave
I wanted it to rain
I wanted the bus to arrive on time
I want to be clever
I want for John to be here
John was wanted to leave
John is wanted to be here
John wanted to be arrested
PS C:\Users\sandiway\Desktop>
```

# Passives

- **Examples:**

1. John was arrested
2. Who was arrested
3. John was believed to have been arrested      (*exceptional passivization*)
4. It is believed that John was arrested
5. That John was arrested is believed by everyone
6. John had insults shouted at him      (*indirect passive, Emonds 2007*)

# Homework 6

## Question 1

- **Raising verb *seems*:**

- John is happy
- John *seems* happy
- It *seems* John is happy
- It *seems* that John is happy
- John *seems* to be happy

- Run these examples on Google Natural Language (Google NL) and the structural probes (SP) program.
- Assume the SP unlabeled relations are labeled with the most likely dependency labels
- What is different between the SP and Google NL analyzes?

# Homework 6

## Question 2

- suppose the lower clause is:
  - John saw Mary
- and the raising predicate is *seemed*:
  - John seemed to see Mary
  - It seems (that) John saw Mary

- Does the SP program do better on these examples?

# Homework 6

## Question 3

- Consider object *wh*-question formation:
  - Who did John see
  - Who did John think (that) Mary saw
  - Who did John (that) say Bill (that) thought Mary saw
- Which does better, the SP program or Google NL?



# Homework 6

## Question 4

- Consider subject *wh*-question formation:
  - Who saw Mary
  - Who did John think saw Mary
  - Who did John say (that) Bill thought saw Mary
- Which does better, the SP program or Google NL?