

# LING/C SC 581:

## Advanced Computational Linguistics

Lecture 23

Prof. Sandiway Fong

# Today's Topic

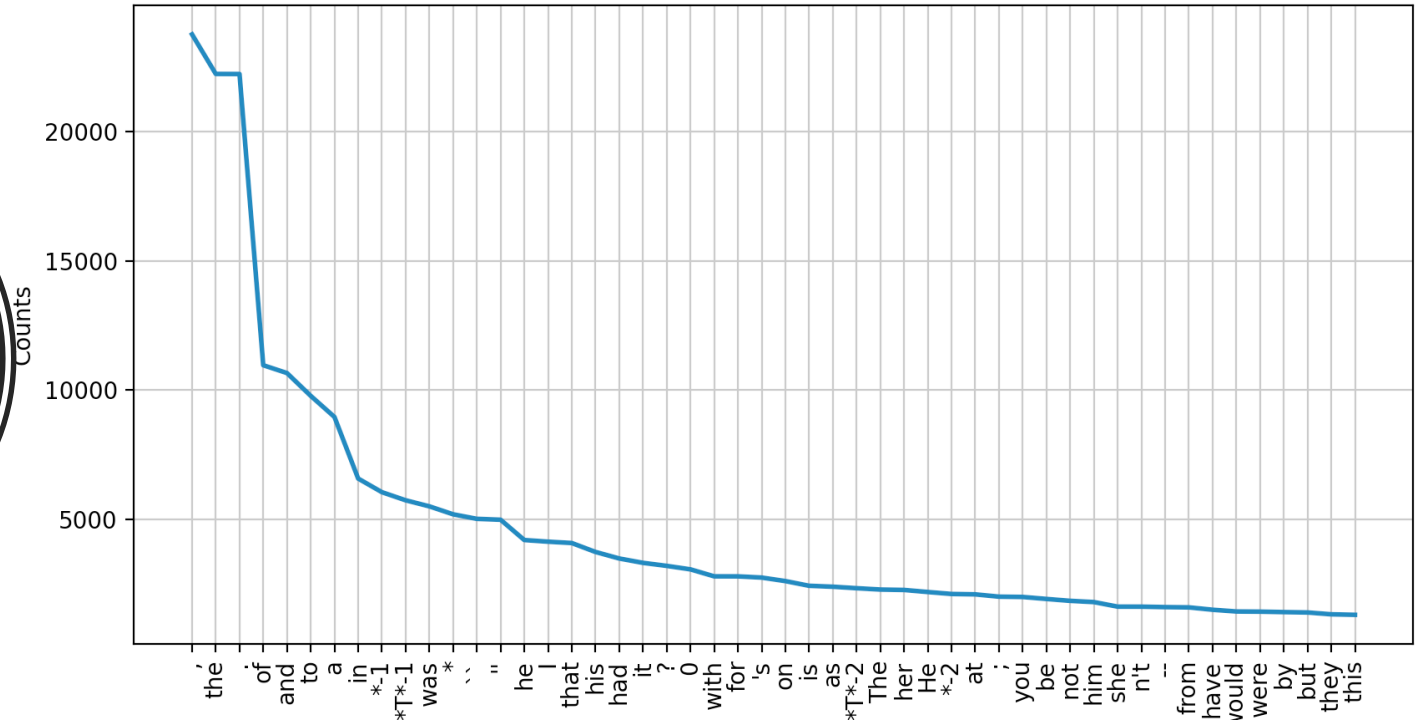
- More on nltk and ptb
  - **Zipf's Law**
  - extracting the grammar rules (called **productions**)
  - looking for words that have multiple POS tags



# Brown corpus: FreqDist

```
>>> 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)
```

# Brown corpus: FreqDist



# `nltk.FreqDist(corpus)`

`N()`

[source]

Return the total number of sample outcomes that have been recorded by this `FreqDist`. For the number of unique sample values (or bins) with counts greater than zero, use `FreqDist.B()`.

Return type

`int`

`update(*args, **kwargs)`

[source]

Override `Counter.update()` to invalidate the cached N

## `nltk.FreqDist.update`

`FreqDist.update(*args, **kwargs)`

Like `dict.update()` but add counts instead of replacing them.

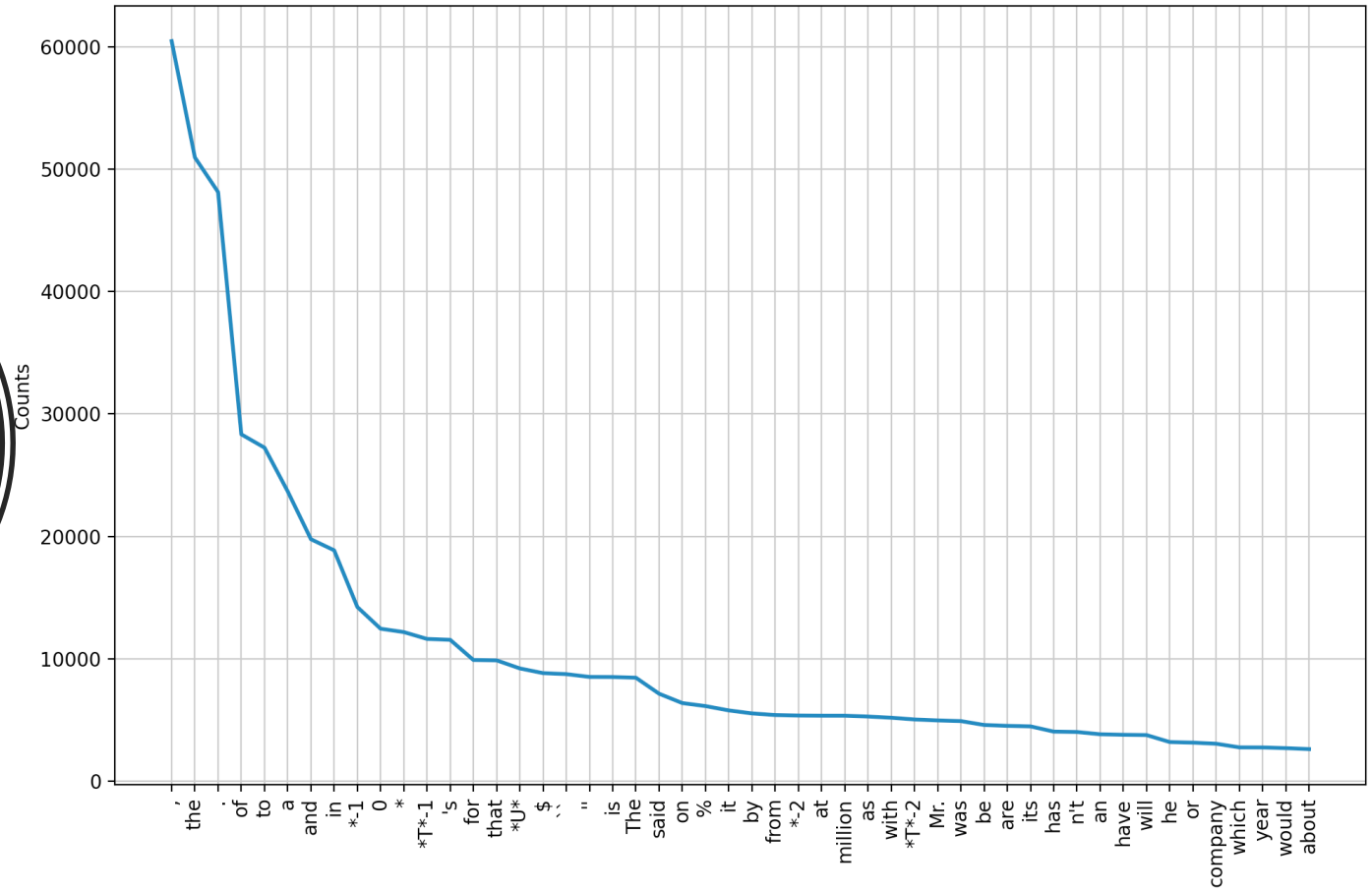
Source can be an iterable, a dictionary, or another `Counter` instance.

# WSJ corpus: FreqDist

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)
```

# WSJ corpus: FreqDist



## Brown +WSJ: FreqDist

- Merge the two FreqDists:

```
>>> fd.update(fd2)
```

```
>>> fd
```

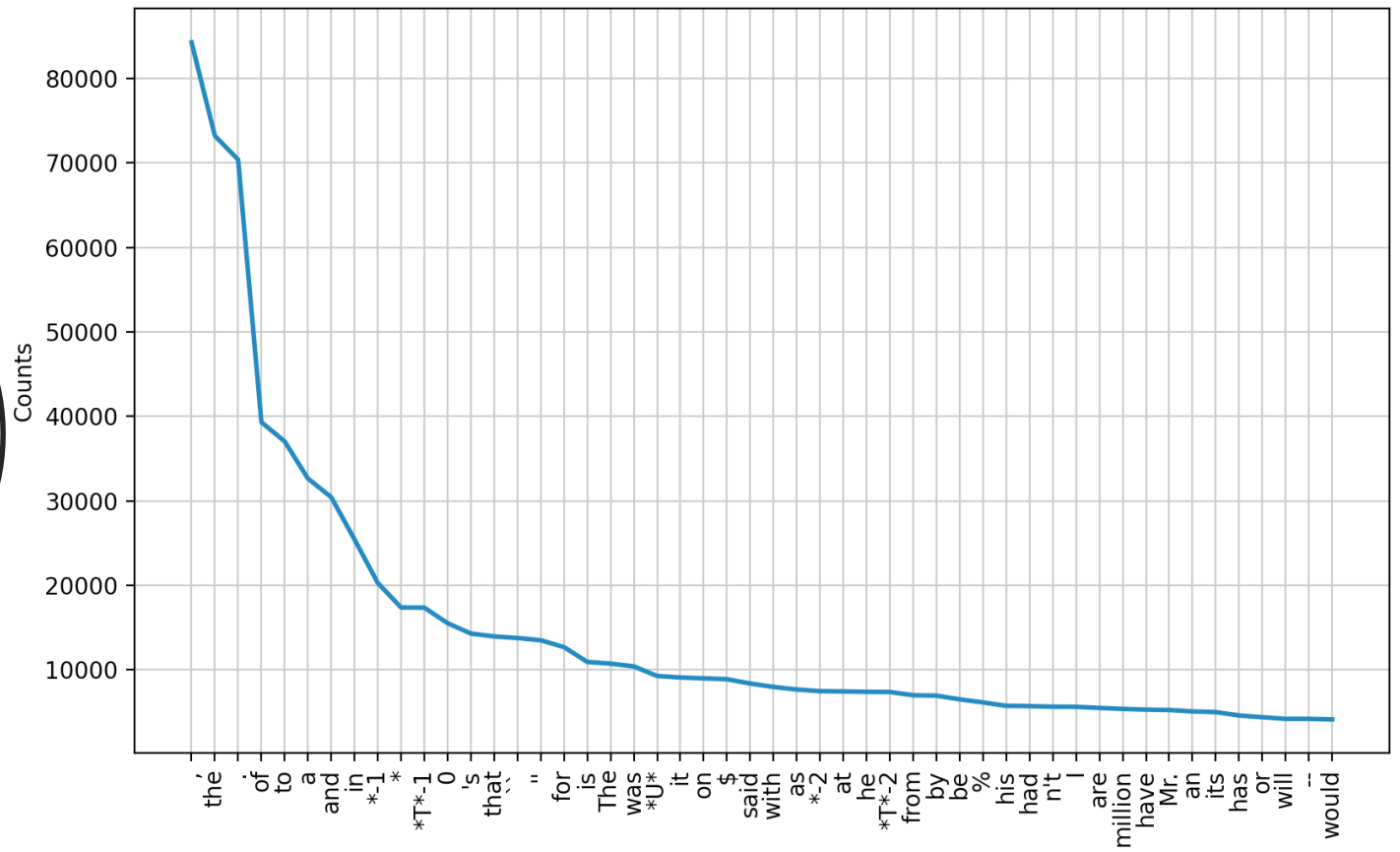
```
FreqDist({' ': 84260, 'the': 73219, '.': 70385, 'of':  
39302, 'to': 37027, 'a': 32634, 'and': 30423, 'in':  
25432, '*-1': 20268, '*': 17363, ...})
```

```
>>> fd.N()
```

```
1740895
```

```
>>> fd.plot(50)
```

# Brown +WSJ: FreqDist



# Zipf's Law

- Zipf's law was originally formulated in terms of *quantitative linguistics*, stating that given some corpus of natural language utterances, **the frequency of any word is inversely proportional to its rank** in the frequency table.
  - Thus, the **most frequent word** will occur approximately **twice** as often as the **second most frequent word**, **three times** as often as the **third most frequent word**, etc.
  - For example, in the Brown Corpus of American English text, the word "**the**" is the most frequently occurring word, and by itself accounts for nearly 7% of all word occurrences (69,971 out of slightly over 1 million).
  - True to Zipf's Law, the second-place word "of" accounts for slightly over 3.5% of words (36,411 occurrences), followed by "and" (28,852).
  - **Only 135 vocabulary items** are needed to account for **half** the Brown Corpus.

# Zipf's Law

Probability and Statistics › Descriptive Statistics ›

Wolfram

## Zipf's Law

In the English language, the probability of encountering the  $r$ th most common word is given roughly by  $P(r) = 0.1 / r$  for  $r$  up to 1000 or so. The law breaks down for less frequent words, since the [harmonic series](#) diverges. Pierce's (1980, p. 87) statement that  $\sum P(r) > 1$  for

### Equation:

- $\text{freq} = c \cdot \text{rank}^{-m}$ 
  - for positive constants  $m$  and  $c$
- $\log(\text{freq}) = -m \log(\text{rank}) + \log(c)$
- has the form of an equation of a straight line ( $y=mx+c$ )

### Code:

```
zipf.py given on the course webpage  
>>> import zipf  
>>> zipf.plot(tokens)  
tokens = list of words (a corpus)
```

# Zipf's Law: ptb

```
>>> wws = []
>>> for id in wids:
...     wws.extend(ptb.words(id))
...
>>> len(wws)
1253013
>>> bws = []
>>> for id in bids:
...     bws.extend(ptb.words(id))
...
```

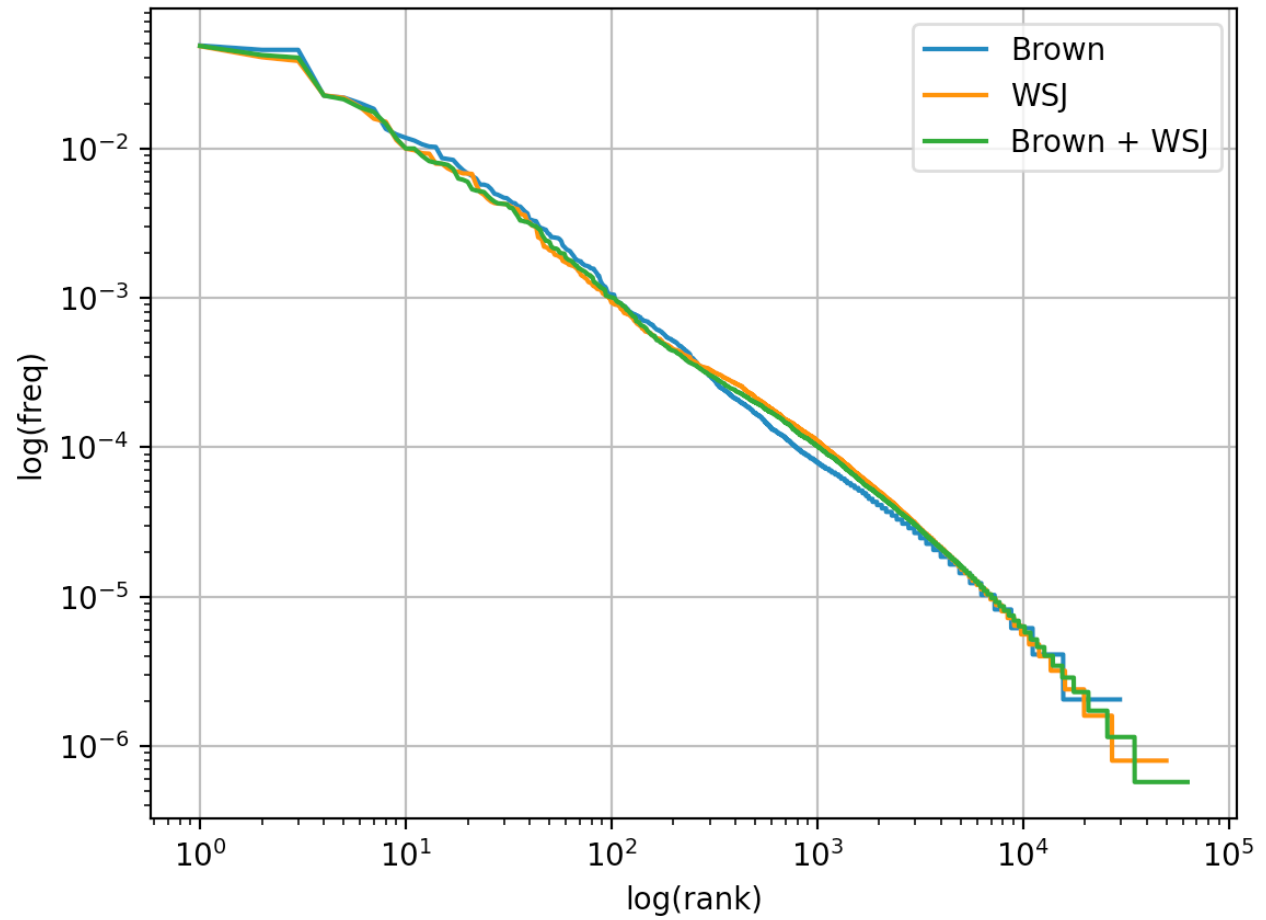
WSJ  
words

Brown  
words

```
>>> len(bws)
487882
>>> from zipf import *
>>> fig()
>>> plot(bws, "Brown")
>>> plot(wws, "WSJ")
>>> plot(bws+wws, "Brown + WSJ")
>>> plt.legend()
<matplotlib.legend.Legend object at
0x1278c4520>
>>> plt.show()
```

Zipf's Law:  
ptb

Log-log plot for freq vs rank



# Zipf's Law: 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')
```

# Treebanks

- are a good source of context-free phrase structure
- can compute statistics

# Trees and Productions

```
>>> len(list(ptb.parsed_sents()[0].subtrees()))
```

```
87
```

productions()

[\[source\]](#)

Generate the productions that correspond to the non-terminal nodes of the tree. For each subtree of the form (P: C1 C2 ... Cn) this produces a production of the form P -> C1 C2 ... Cn.

```
>>> len(ptb.parsed_sents()[0].productions())
```

```
87
```

# Trees and Productions

```
>>> ptb.parsed_sents()[0].productions()
[S -> S : S ., S -> PP , NP-SBJ-2 VP, PP -> IN NP, IN -> 'In', NP -> JJ
NN, JJ -> 'American', NN -> 'romance', NP-SBJ-2 -> RB NN, RB ->
'almost', NN -> 'nothing', VP -> VBZ $, 'VBZ -> 'rates', S -> NP-SBJ ADJP-
PRD, NP-SBJ -> -NONE-, -NONE- -> '*-2', ADJP-PRD -> ADJP PP, ADJP -> JJR,
JJR -> 'higher', PP -> IN SBAR-NOM, IN -> 'than', SBAR-NOM -> WHNP-1 S,
WHNP-1 -> WP, WP -> 'what', S -> NP-SBJ VP, NP-SBJ -> DT NN NNS, DT ->
'the', NN -> 'movie', NNS -> 'men', VP -> VB VP, VB -> 'have', VP -> VBN
S, VBN -> 'called', $ -> NP-SBJ, S-NOM-PRD, NP-SBJ -> -NONE-, -NONE-
-> '*T*-1', S-NOM-PRD -> NP-SBJ VP, NP-SBJ -> -NONE-, -NONE-
-> '*-1', VP -> NN NP, NN -> 'meeting', NP -> JJ, JJ -> 'cute', ' ' -> ':'
-> '-1', S -> S-ADV, NP-SBJ VP, S-ADV -> NP-SBJ VP, NP-SBJ -> DT, DT ->
'that', VP -> VBZ, VBZ -> 'is', NP-SBJ -> NN, NN -> 'boy-meets-
girl', VP -> VBZ ADJP-PRD SBAR-ADV, VBZ -> 'seems', ADJP-PRD -> RB JJ, RB
-> 'more', JJ -> 'adorable', SBAR-ADV -> IN S, IN -> 'if', S -> NP-SBJ VP,
NP-SBJ -> PRP, PRP -> 'it', VP -> VBZ RB VP, VBZ -> 'does', RB -> "n't",
VP -> VB NP PP, VB -> 'take', NP -> NN, NN -> 'place', PP -> IN NP, IN ->
'in', NP -> NP, PP -> DT NN, DT -> 'an', NN -> 'atmosphere', PP -> IN
NP, IN -> 'of', NP -> ADJP NN, ADJP -> JJ CC JJ, JJ -> 'correct', CC ->
'and', JJ -> 'acute', NN -> 'boredom', . -> '.']
```

# Trees and Productions

- <https://www.nltk.org/api/nltk.grammar.Production.html>

```
>>> ptb.parsed_sents()[0].productions()[2]
PP -> IN NP
>>> type(ptb.parsed_sents()[0].productions()[2])
<class 'nltk.grammar.Production'>
>>> ptb.parsed_sents()[0].productions()[2].lhs()
PP
>>> ptb.parsed_sents()[0].productions()[2].rhs()
(IN, NP)
>>> type(ptb.parsed_sents()[0].productions()[2].rhs())
<class 'tuple'>
```

# Trees and Productions

3<sup>rd</sup> rule is PP → IN NP:

```
>>> ptb.parsed_sents()[0].productions()[2].rhs()[0]
```

```
IN
```

```
>>> ptb.parsed_sents()[0].productions()[2].rhs()[1]
```

```
NP
```

```
>>> len(ptb.parsed_sents()[0].productions()[2].rhs())
```

```
2
```

# Trees and Productions

3<sup>rd</sup> rule is PP -> IN NP:

```
>>> ptb.parsed_sents()[0].productions()[2].is_nonlexical()
```

True

```
>>> ptb.parsed_sents()[0].productions()[2].is_lexical()
```

False

`is_nonlexical()`

Return True if the right-hand side only contains `Nonterminals`

`is_lexical()`

Return True if the right-hand contain at least one terminal token.

# Words with multiple POS tags

- Let's write a program to find words with more than one part of speech tag.
- First, let's get all the word-tag items:

```
>>> wt = [item for tree in ptb.parsed_sents() for item in tree.pos()]  
>>> len(wt)  
1740895
```

- Next, let's get the set of word-tag items, no duplicates:

```
>>> wts = set(wt) ← wts = word tag set  
>>> len(wts)  
74323
```

# Words with multiple POS tags

- Let's create a dictionary with pos tags as values:

```
>>> d = {}
>>> for item in wts:
>>>     d.setdefault(item[0], []).append(item[1])
>>>
>>> len(d)
63073
```

# Words with multiple POS tags

- *Let's look at a few examples ...*

tagguid1.pdf

```
>>> d['any']  
['DT', 'RB']  
>>> d['Any']  
['DT']
```

**any** can be a determiner (DT).

EXAMPLES: We don't have any/DT.  
Don't you want any/DT more/JJR?

However, when it precedes a comparative adverb, it is an adverb (RB).

EXAMPLES: I can't run any/RB further/**RBR**.  
I can't go on like this any/RB more/**RBR**.

# Words with multiple POS tags

**about** when used to mean “approximately” should be tagged as an adverb (RB), rather than a preposition (IN).

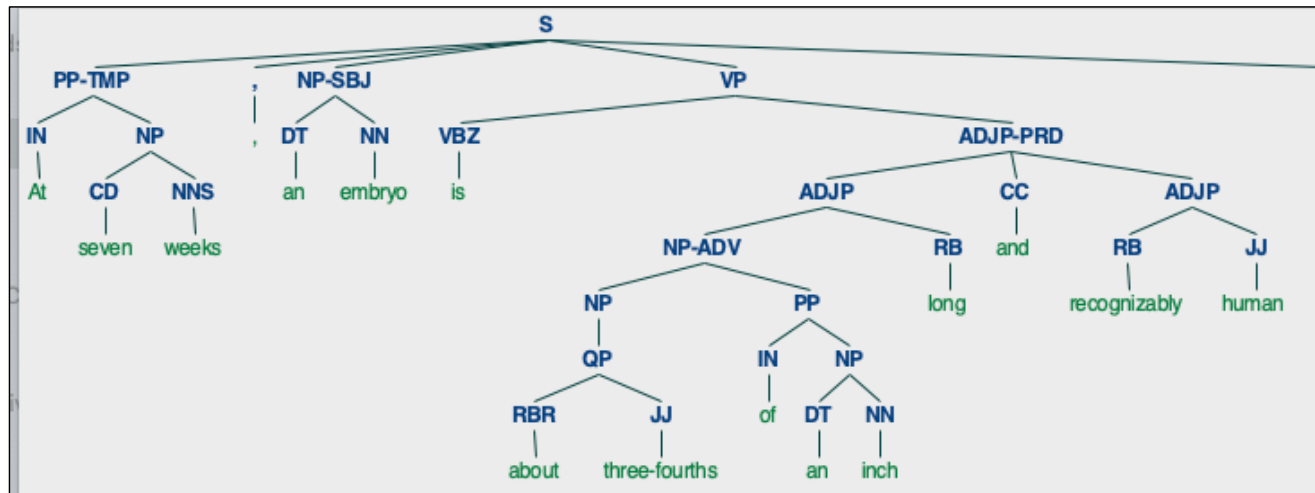
- Particle|RP
  - This category includes a number of mostly monosyllabic words that also double as prepositions.
- Adverb, comparative|RBR
  - *closer, later, less, more, further* – see previous slide

```
>>> d['about']  
['IN', 'RB', 'RP', 'RBR', 'JJ']  
>>> d['About']  
['IN', 'RB']
```

# Words with multiple POS tags

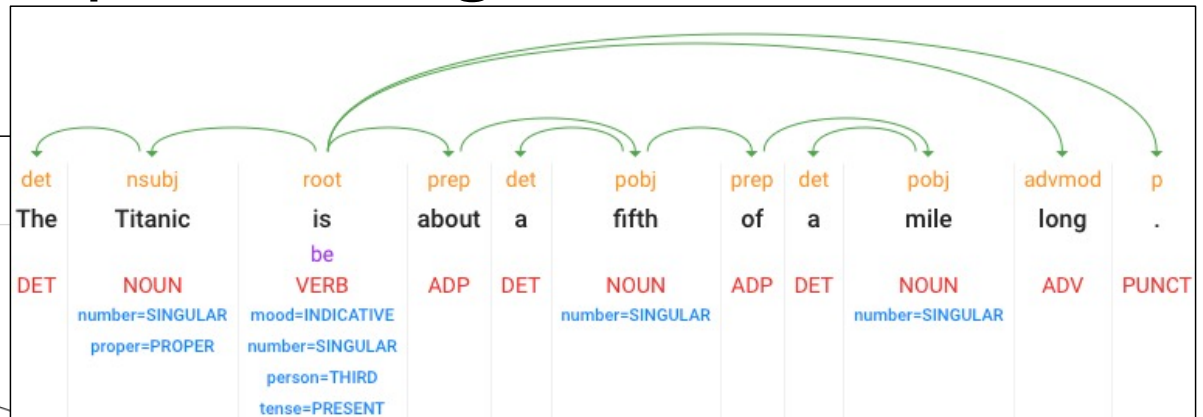
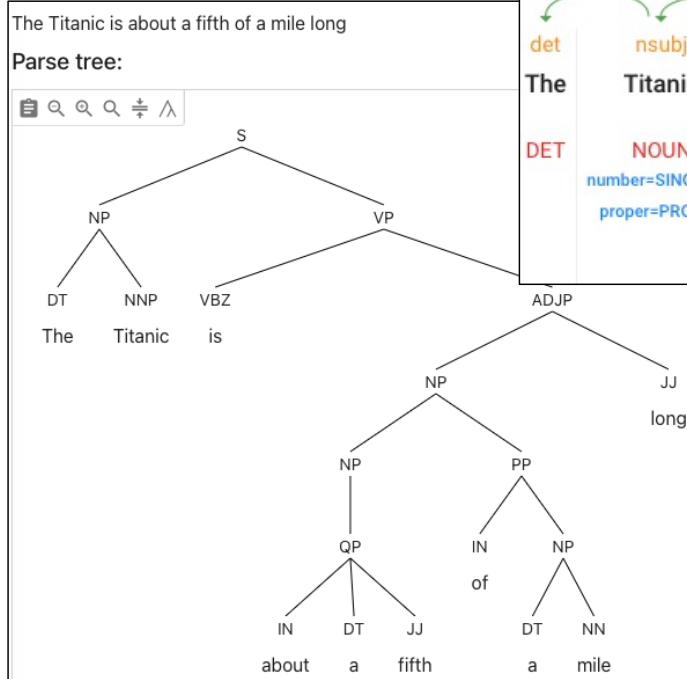
- *about*: <https://www.merriam-webster.com/dictionary/about>
  1. adverb: *about* a year ago
  2. adverb: looked about for a place to park
  3. adverb: They go about in circles.
  4. adverb: He spoke to the people standing about.
  5. adverb: the other way about
  6. preposition: People gathered about him
  7. preposition: He traveled about the country.
  8. preposition: Fish are abundant about the reefs.
  9. preposition: spoke *about* his past
  10. preposition: act as if they know what they're about
  11. adjective: is up and *about* by 7 a.m.
  12. adjective: There is a scarcity of jobs *about*.

# Words with multiple POS tags



- about = RBR?

# Words with multiple POS tags



<https://cloud.google.com/natural-language>

<https://parser.kitaev.io>

# Words with multiple POS tags

EXAMPLES: You should eat less/JJR (in terms of quantity).  
(cf. You should eat less/JJR cheese.)

You should eat less/RBR (in terms of frequency).  
(cf. You should eat rarely/RB.)

You should work less/RBR.  
(cf. You should work harder/RBR.)

*Less* should be tagged as a comparative adjective (JJR) even when it occurs without a head noun, as in *less of a problem*.

*Less* in the sense of *minus* should be tagged as a coordinating conjunction (CC).

```
>>> d['less']  
['RB', 'RBR', 'CC', 'NN', 'JJR', 'JJS']  
>>> d['Less']  
['RBR', 'NNP', 'JJR']
```

# Words with multiple POS tags

## **JJ or NN**

Nouns that are used as modifiers, whether in isolation or in sequences, should be tagged as nouns (NN, NNS) rather than as adjectives (JJ).

**EXAMPLES:** wool/NN sweater (vs. woollen/JJ sweater)  
terminal/NN type (vs. terminal/JJ cancer)  
life/NN insurance/NN company

Hyphenated modifiers, on the other hand, should always be tagged as adjectives (JJ). Thus, we have different part-of-speech assignments in examples like the following—depending on the orthographic conventions used:

**EXAMPLES:** income-tax/JJ return; income/NN tax/NN return  
value-added/JJ tax; value/NN added/VBN tax

# Words with multiple POS tags

```
>>> d['wool']
```

```
['NN']
```

```
>>> d['terminal']
```

```
['JJ', 'NN']
```

```
>>> d['woollen']
```

```
Traceback (most recent call last):
```

```
  File "<stdin>", line 1, in <module>
```

```
KeyError: 'woollen'
```

```
>>> d['Wool']
```

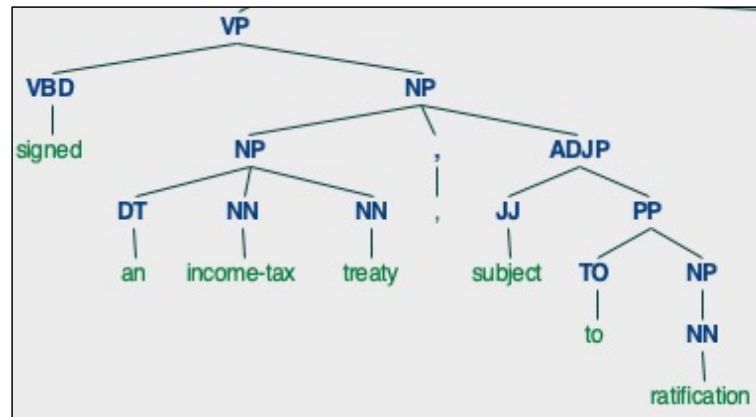
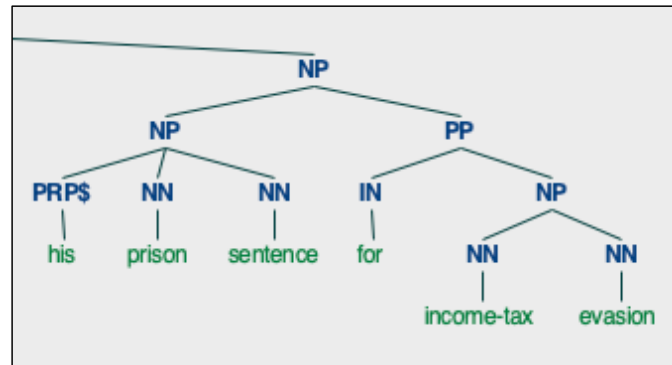
```
Traceback (most recent call last):
```

```
  File "<stdin>", line 1, in <module>
```

```
KeyError: 'Wool'
```

# Words with multiple POS tags

```
>>> d['income-tax']  
['JJ', 'NN']  
>>> d['income']  
['NN']  
>>> d['tax']  
['NN', 'VB']
```



# Words with multiple POS tags

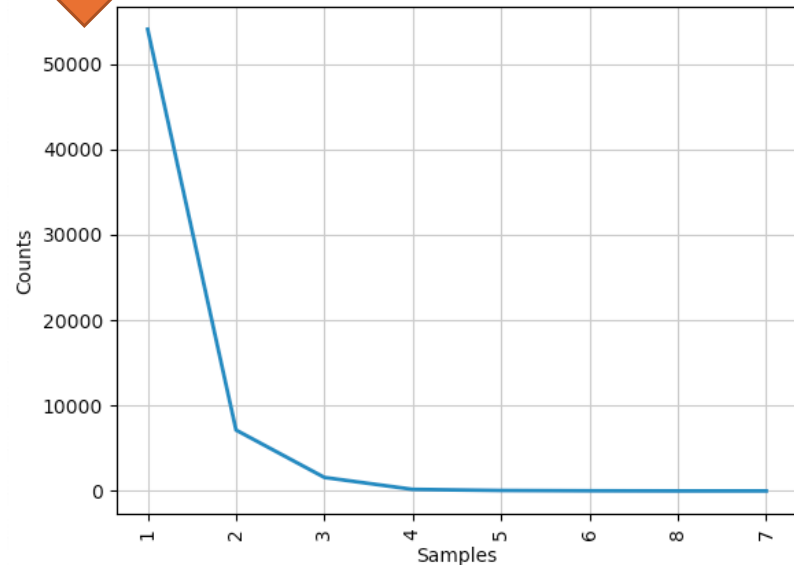
- Let's look at the frequency distribution by # of pos tags:

- recall our dictionary d maps words to pos tags*

```
>>> fd = nltk.FreqDist([len(d[k])  
for k in d])  
>>> fd.most_common()  
[(1, 54075), (2, 7137), (3, 1588),  
(4, 188), (5, 60), (6, 20), (8, 3),  
(7, 2)]  
>>> fd.plot()
```

(#pos tags, #words)

# vocab  
items



# pos tags