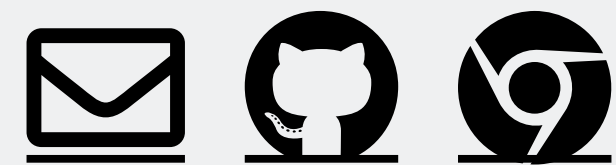


---

# LG 467 Computers in Linguistics

## [1-2021] Topic 5: POS tagging

Sakol Suethanapornkul



# Previously...

We can count frequencies of individual items with `FreqDist`:

```
import nltk

f = open('text1.txt')
txt = f.read()
f.close()

tokens = nltk.word_tokenize(txt)
freq = nltk.FreqDist(tokens)

freq.most_common(15)
freq.plot()
freq.plot(cumulative = True)
```

Code 6.2

# Previously...

You can define your own functions to automate tasks. Functions can take multiple parameters plus a docstring:

```
def ttr(lst, digits):  
    """Compute type-token ratio on a list of strings,  
    accepts one list and digits to round."""  
    result = len(set(lst))/len(lst)  
    return round(result, digits)
```

```
ttr(text1, 3)  
ttr(text2, 4)  
help(ttr)
```

Code 6.11

# Previously...

Language is not random. Things pattern together.

- We can ask what the next word might be given context
- Simple and efficient way of modeling context is using *n-grams*
- N-grams are units of n sequences
  - Characters: 'custard'
    - *(c, u, s, t, a, r, d)*     *(cu, us, st, ta, ar, rd)*     *(cus, ust, sta, tar, ard)*
  - Words: 'He is eating fried rice'
    - *(he, is, eating, fried, rice)*     *(he is, is eating, eating fried, fried rice)*

# Previously...

N-grams to model language data

- "Turn in your \_\_\_\_\_"  $\rightarrow p(\text{_____} | \text{"turn in your"})$
- decide on a useful context size (2, 3, 4, or ?)
- save analyses not of individual words but of words given the previous n words

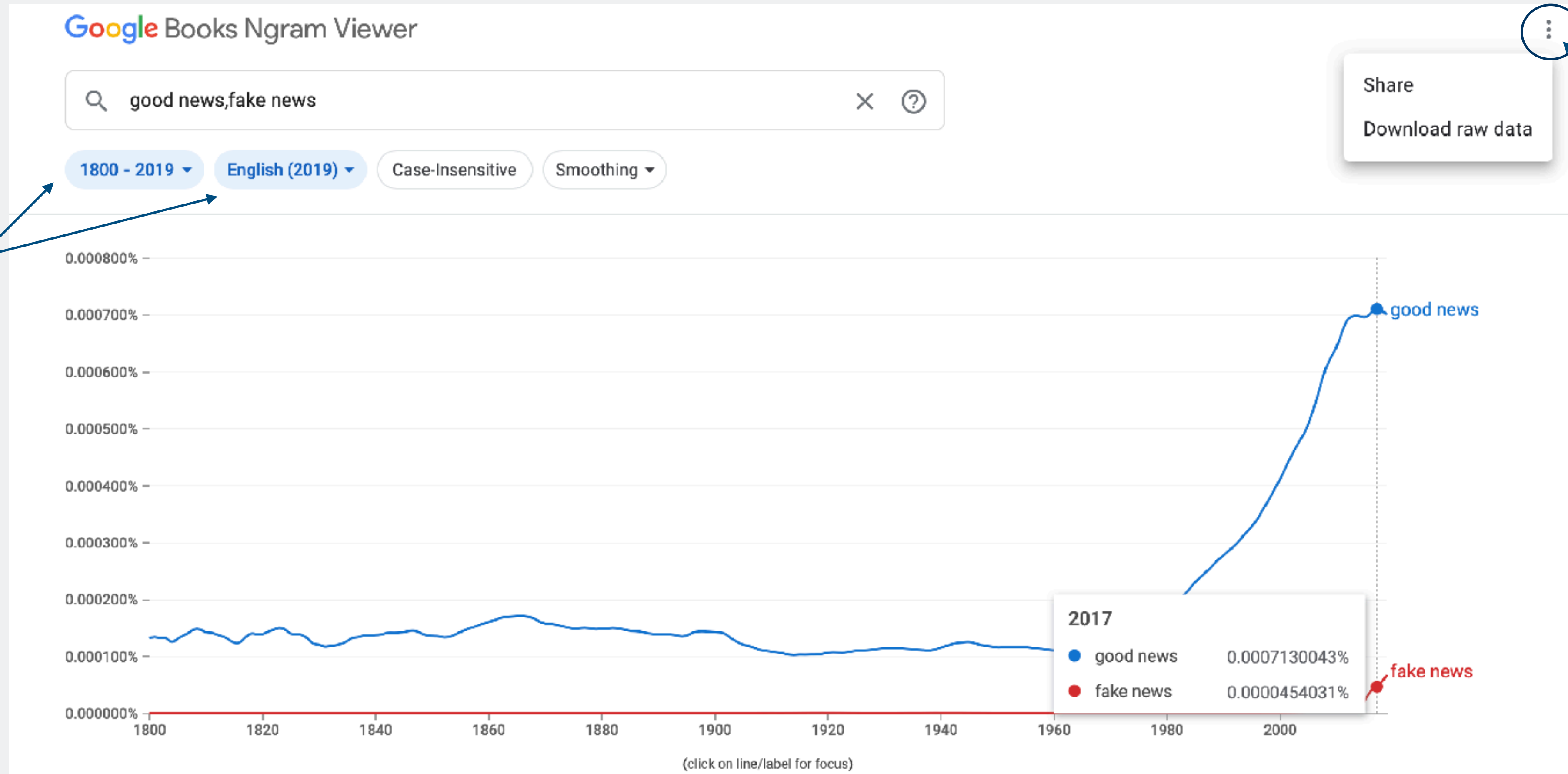
# Previously...

Bigram counts from the Berkeley Restaurant Project:

	<b>i</b>	<b>want</b>	<b>to</b>	<b>eat</b>	<b>chinese</b>	<b>food</b>	<b>lunch</b>	<b>spend</b>
<b>i</b>	5	827	0	9	0	0	0	2
<b>want</b>	2	0	608	1	6	6	5	1
<b>to</b>	2	0	4	686	2	0	6	211
<b>eat</b>	0	0	2	0	16	2	42	0
<b>chinese</b>	1	0	0	0	0	82	1	0
<b>food</b>	15	0	15	0	1	4	0	0
<b>lunch</b>	2	0	0	0	0	1	0	0
<b>spend</b>	1	0	1	0	0	0	0	0

# Previously...

Check out: Google N-grams and Google Books ngram viewer



Click on this

Toggle to change

# Previously...

We can generate bigram (or any n-gram) counts with NLTK:

```
import nltk

nltk.bigrams(text1)

list(nltk.bigrams(text1))

nltk.ngrams(text1, 2)
list(nltk.ngrams(text1, 2))

bigrams = list(nltk.ngrams(text1, 2))
bicount = FreqDist(bigrams)
```

Code 6.12



# Beneath the surface

Our token n-gram models represent transitions between observed characters/words

- We can call them **Visible Markov Models (VMMs)**

For things that are overt, we are also interested in the probabilities of **hidden** categories

- We will need **Hidden Markov Models (HMMs)**

# Beneath the surface

Did I hear hidden categories? What did you mean?

- We may not be interested in a phrase like 'awesome news'
- Instead, we might want to know the likelihood of adjectives followed by nouns (ADJ + NOUN)
- How can we look at categories that are not in the data explicitly?

# Parts of speech (POS)

The idea that words can be classified into grammatical categories has been around for two millennia

- **part of speech**, word classes, **POS**, **POS tags**

8 parts of speech attributed to Dionysius Thrax of Alexandria (c. 1st C. BCE):

- noun, verb, pronoun, preposition, adverb, conjunction, participle, article
- These categories are relevant for NLP today.

# Parts of speech (POS)

Parts of speech are defined based on:

- grammatical relationship with their neighboring words
- Morphological properties about their affixes (*-ness, -able, -ed*)

# Two classes of English words

## 1. Closed class words

- Relatively fixed membership
- Usually function words:
  - determiners: *this, that, a, an, the*
  - pronouns: *she, I, me, them, her, our, who, whose*
  - prepositions: *on, under, over, at, with, ...*
  - conjunctions: *and, but, while, because, that, ...*
  - particles: *go over, turn in, ...*

# Two classes of English words

## 2. Open class words

- Usually content words: Nouns, Verbs, Adjectives, Adverbs
  - Plus interjections: *oh, ouch, uh-huh, yes, no, hello*
- New nouns and verbs like *iPhone, facebook, mansplain, google*

## Open class ("content") words

### Nouns

#### Proper

*Janet*  
*Italy*

#### Common

*cat, cats*  
*mango*

### Verbs

#### Main

*eat*  
*went*

#### Auxiliary

*can*  
*had*

Adjectives *old green tasty*

Adverbs *slowly yesterday*

### Numbers

*122,312*  
*one*

Interjections *Ow hello*

*... more*

## Closed class ("function")

Determiners *the some*

Conjunctions *and or*

Pronouns *they its*

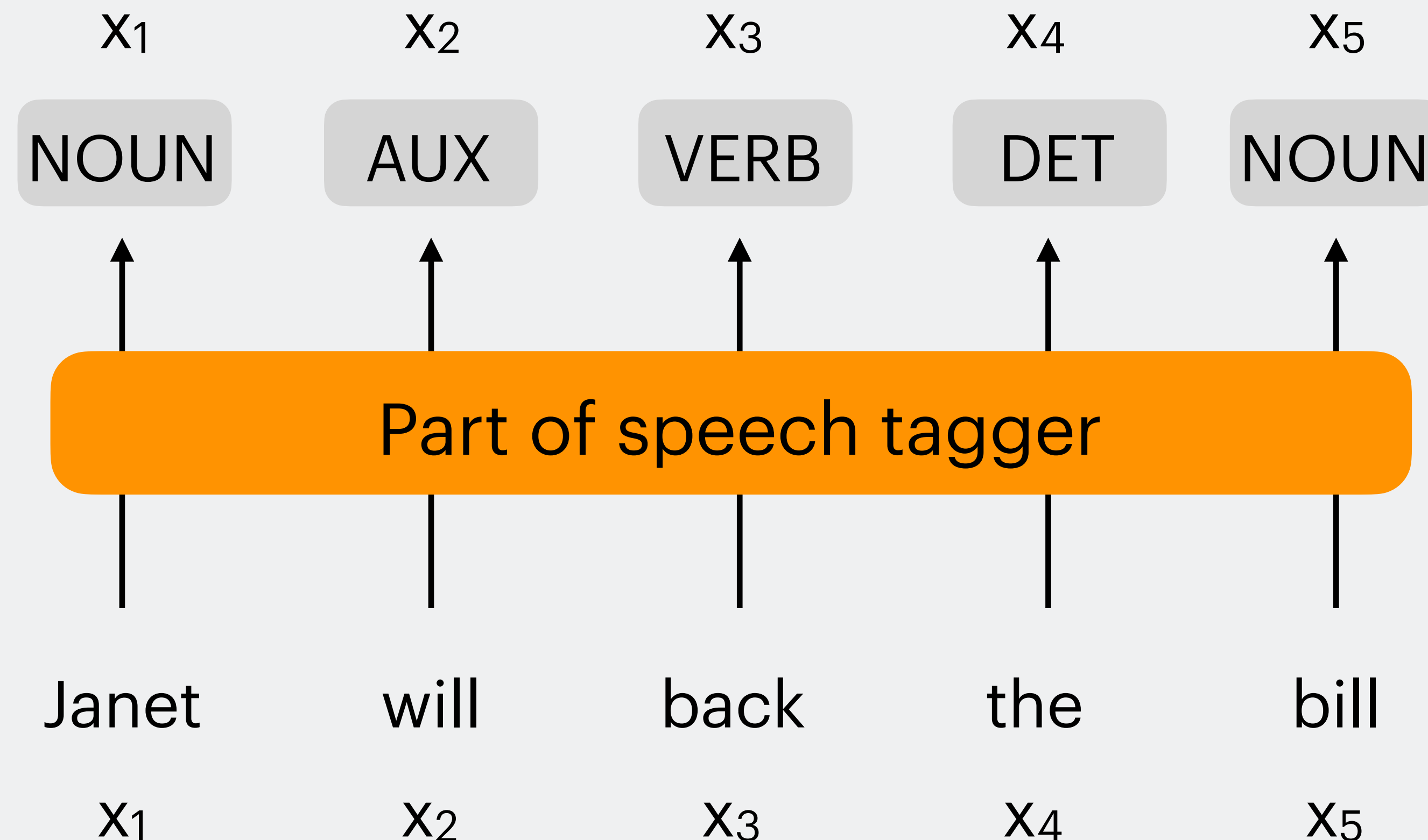
Prepositions *to with*

Particles *off up*

*... more*

# POS tagging

POS tagging = assigning a part of speech to each word in a text





# POS tagging

Tagging is a disambiguation task

- Words often have more than one POS
- The goal is to find the correct tag for the situation

Case 1: 'book'

- Book a flight
- Whose book is it?

Case 2: 'tired'

- He tired me out
- I'm just tired

Case 3: 'that'

- Hand me that book
- I think that you liked it

# POS tagging

Though tagging is challenging, the accuracy of POS tagging algorithms is extremely high (97% or higher)

- 55-67% of word tokens in running text are ambiguous!

<b>Types:</b>		<b>WSJ</b>	<b>Brown</b>
<b>Unambiguous</b>	(1 tag)	44,432 ( <b>86%</b> )	45,799 ( <b>85%</b> )
<b>Ambiguous</b>	(2+ tags)	7,025 ( <b>14%</b> )	8,050 ( <b>15%</b> )
<b>Tokens:</b>			
<b>Unambiguous</b>	(1 tag)	577,421 ( <b>45%</b> )	384,349 ( <b>33%</b> )
<b>Ambiguous</b>	(2+ tags)	711,780 ( <b>55%</b> )	786,646 ( <b>67%</b> )

# Tagsets for English

Common in the US:

- Penn Treebank Tagset (PTB) 45 tags
- Brown Tagset 87 tags

Common in the UK:

- CLAWS 5 (or C5) 61 tags

# Penn Treebank Tags

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coord. conj.	<i>and, but, or</i>	NNP	proper noun, sing.	<i>IBM</i>	TO	“to”	<i>to</i>
CD	cardinal number	<i>one, two</i>	NNPS	proper noun, plu.	<i>Carolinas</i>	UH	interjection	<i>ah, oops</i>
DT	determiner	<i>a, the</i>	NNS	noun, plural	<i>llamas</i>	VB	verb base	<i>eat</i>
EX	existential ‘there’	<i>there</i>	PDT	predeterminer	<i>all, both</i>	VBD	verb past tense	<i>ate</i>
FW	foreign word	<i>mea culpa</i>	POS	possessive ending	<i>'s</i>	VBG	verb gerund	<i>eating</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	PRP	personal pronoun	<i>I, you, he</i>	VBN	verb past partici- ple	<i>eaten</i>
JJ	adjective	<i>yellow</i>	PRP\$	possess. pronoun	<i>your, one's</i>	VBP	verb non-3sg-pr	<i>eat</i>
JJR	comparative adj	<i>bigger</i>	RB	adverb	<i>quickly</i>	VBZ	verb 3sg pres	<i>eats</i>
JJS	superlative adj	<i>wildest</i>	RBR	comparative adv	<i>faster</i>	WDT	wh-determ.	<i>which, that</i>
LS	list item marker	<i>1, 2, One</i>	RBS	superlatv. adv	<i>fastest</i>	WP	wh-pronoun	<i>what, who</i>
MD	modal	<i>can, should</i>	RP	particle	<i>up, off</i>	WP\$	wh-possess.	<i>whose</i>
NN	sing or mass noun	<i>llama</i>	SYM	symbol	<i>+, %, &amp;</i>	WRB	wh-adverb	<i>how, where</i>



# Practice

Try tagging the following sentences with the PTB tags:

1. The/ grand/ jury/ commented/ on/  
a/ number/ of/ other/ topics/ ./
2. Although/ preliminary/ findings/ were/  
reported/ more/ than/ a/ year/  
ago/ ./ the/ latest/ results/  
appear/ in/ today/ 's/ New/  
England/ Journal/ of/ Medicine/ ./

# Practice

Try tagging the following sentences with the PTB tags:

1. The/**DT** grand/**JJ** jury/**NN** commented/**VBD** on/**IN**  
a/**DT** number/**NN** of/**IN** other/**JJ** topics/**NNS** ./.
2. Although/**IN** preliminary/**JJ** findings/**NNS** were/**VBD**  
reported/**VBN** more/**RBR** than/**IN** a/**DT** year/**NN**  
ago/**IN** ./, the/**DT** latest/**JJS** results/**NNS**  
appear/**VBP** in/**IN** today/**NN** 's/**POS** New/**NNP**  
England/**NNP** Journal/**NNP** of/**IN** Medicine/**NNP** ./.

# Universal Dependencies Tagset

	Tag	Description	Example
Open Class	<b>ADJ</b>	Adjective: noun modifiers describing properties	<i>red, young, awesome</i>
	<b>ADV</b>	Adverb: verb modifiers of time, place, manner	<i>very, slowly, home, yesterday</i>
	<b>NOUN</b>	words for persons, places, things, etc.	<i>algorithm, cat, mango, beauty</i>
	<b>VERB</b>	words for actions and processes	<i>draw, provide, go</i>
	<b>PROPN</b>	Proper noun: name of a person, organization, place, etc..	<i>Regina, IBM, Colorado</i>
	<b>INTJ</b>	Interjection: exclamation, greeting, yes/no response, etc.	<i>oh, um, yes, hello</i>
Closed Class Words	<b>ADP</b>	Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation	<i>in, on, by under</i>
	<b>AUX</b>	Auxiliary: helping verb marking tense, aspect, mood, etc.,	<i>can, may, should, are</i>
	<b>CCONJ</b>	Coordinating Conjunction: joins two phrases/clauses	<i>and, or, but</i>
	<b>DET</b>	Determiner: marks noun phrase properties	<i>a, an, the, this</i>
	<b>NUM</b>	Numeral	<i>one, two, first, second</i>
	<b>PART</b>	Particle: a preposition-like form used together with a verb	<i>up, down, on, off, in, out, at, by</i>
	<b>PRON</b>	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others</i>
<b>SCONJ</b>	Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement	<i>that, which</i>	
Other	<b>PUNCT</b>	Punctuation	<i>; , ()</i>
	<b>SYM</b>	Symbols like \$ or emoji	<i>%, %</i>
	<b>X</b>	Other	<i>asdf, qwfg</i>

# Universal Dependencies Tagset for Thai

Abbreviation	Part-of-Speech tag	Examples
ADJ	Adjective	ใหม่, พิเศษ, ก่อน, มาก, สูง
ADP	Adposition	แม้, ว่า, เมื่อ, ของ, สำหรับ
ADV	Adverb	ก่อน, ก็, เล็กน้อย, เลย, สุด
AUX	Auxiliary	เป็น, ไข่, คือ, คล้าย
CCONJ	Coordinating conjunction	แต่, และ, หรือ
DET	Determiner	ที่, นี้, ซึ่ง, ทั้ง, ทุก, หลาย
INTJ	Interjection	อ๋อ, โอ๊ย
NOUN	Noun	กำมือ, พวก, สนาม, กีฬา, บัญชี
NUM	Numeral	5,000, 103.7, 2004, หนึ่ง, ร้อย
PART	Particle	มา ขึ้น ไม่ได้ เข้า
PRON	Pronoun	เรา, เขา, ตัวเอง, ใคร, เธอ
PROPN	Proper noun	โอบามา, แคปิตอลฮิล, จีไอพี, ไมเคิล
PUNCT	Punctuation	(,), “, “: ”
SCONJ	Subordinating conjunction	หาก
VERB	Verb	เปิด, ให้, ไข่, เผชิญ, อ่าน

- Universal Dependencies (UD) is one of the tagsets in PyThaiNLP
- In PyThaiNLP, this tagset is known as Parallel Universal Dependencies (PUD)



# POS tagging in NLTK

An off-the-shelf tagger is available for English:

```
from nltk import pos_tag, word_tokenize

text = "John's big idea isn't all that bad."
token = word_tokenize(text)
pos = pos_tag(token)

print(pos)
```

Code 7.1

**Question:** What tagset is this?

# POS tagging in NLTK

What can we do with pos? We can separate tags from tokens

```
from nltk import FreqDist

tok = [tok for (tok, tag) in pos]
tag = [tag for (tok, tag) in pos]

# Then, you may choose to count
FreqDist(tok)
FreqDist(tag)
```

Code 7.2



# Our plan next week...

- Part-of-speech (POS) tagging
  - More tagging!
  - Please install SpaCy
    - Windows:
      - Use Anaconda Prompt to install ([Here](#))
      - Follow the instructions on SpaCy ([Here](#))
    - Mac:
      - Use Terminal