



February 26, 2018 at 1:18 PM

CIS 530 - 2/26/18

Announcements:

- ① No quiz on Wednesday. I have decided to give you that time so that you can study for midterms in your other classes.
- ② Wednesday's class will be a lecture rather than an in-class exercise. I'm thinking through what changes to make to the class based on your mid-semester feedback. Thank you for your thoughtful suggestions!
- ③ Instead of a programming assignment this week, we will have you do the first part of the midterm project. Details:

- Form a team [min size 4, max 6]
- Come up with 3 project ideas
- The goal of the term project is to design a homework assignment similar to ones you have completed in class.
- Your project will consist of these components

Shared tasks

- SEMEval
- CoNLL
- WMT

- { 1) A description of the problem, in the style of the homework write-ups
- 2) Training and evaluation data
- 3) A commented implementation of the simplest baseline that solves the problem
- 4) " " " of a baseline published in the literature
- 5) One extension per team member that



< Notes



Shared tasks

- SEMEval
- CoNLL
- WMT

NLG:
BLEU
Rouge

- Your project will consist of these components
 - 1) A description of the problem, in the style of the homework write-ups
 - 2) Training and evaluation data
 - 3) A commented implementation of the simplest baseline that solves the problem
 - 4) " " " of a baseline published in the literature
 - 5) One extension per team member that attempts to beat the baseline.
 - 6) A scoring script implementing an objective function that can be used to score submissions on the leaderboard
- These will be delivered in different milestones
 - The first milestone is due Wednesday after spring break. (March 14)
 - 1) Form a team
 - 2) Submit a short ~1-2 paragraph description of your 3 ideas.

We will use class on Monday after Spring Break to help you.

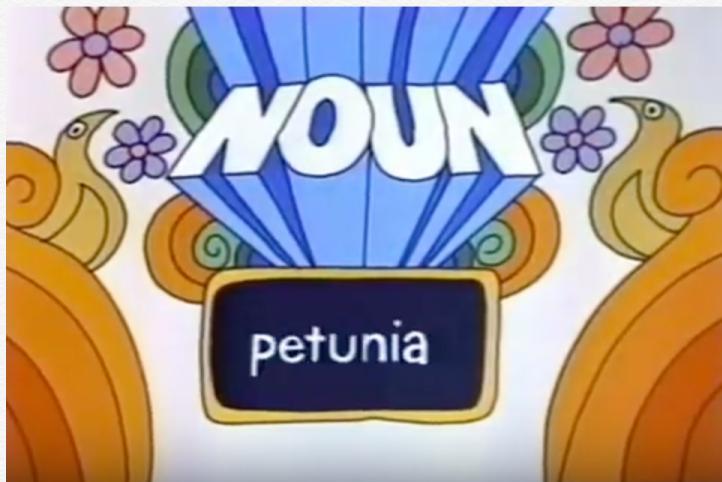
Todays lecture: Part of Speech Tagging
(Chapter 10 in the textbook)



< Notes



Today's lecture: Part of Speech Tagging
(Chapter 10 in the textbook)



Ancient Greek tag set: noun, verb, pronoun, preposition, adverb, conjunction, participle, article

Schoolhouse Rock tag set: + adjective + interjection
(c. 1970) - participle - article

Every word in Vocabulary belongs to one or more of these word classes

Assigning the classes to words in a sentence is called POS tagging.

Many words can have multiple POS tags.
Can you think of some?

< Notes



Four major classes ← open classes

Nouns

Verbs

Adjectives

Adverbs

English has all 4 but not every language does

NOUN - person, place, or thing

↳ Proper Nouns names of specific entities
or people Prof. Marcus NRA University of Pennsylvania

common nouns

the cats
a cereal
democracy

snow
salt
communism

count nouns (allow grammatical enumeration, occurring in both singular and plural ungrammatical) cat is cute — cats are cute

mass nouns (conceptualized as homogenous groups, cannot be pluralized, can appear without determiners even in singular form)

the snow is white
grammatical snow is white

Verbs - words describing actions and processes

English verbs have inflectional markers.

- + 3rd person singular
- = non -
- progressive (ing)
- past

		compute
		(it) computes + s
		(they)/you/I compute -
		computing + ing
		computed + ed

Adjectives

Scalar adjectives
temp: lukewarm < warm < hot
boiling < scalding <

- words that describe qualities

Properties or good/bad
color, age, value,
material, subjective words



< Notes



Adjectives - words that describe properties or qualities like color, age, value, material, subjective words, good/bad, amazing.

Scalar adjectives
temp: lukewarm < warm < hot
boiling < scalding < face-melting

Adverbs - modify verbs or whole verb phrases, or other words like adjectives

very
warm
→
hot

↳ locatives: here, home, uphill

very
delicious
→
scrumptious

↳ degree: very, extremely, extraordinarily, somewhat, not really, ish.

↳ Manner: slowly / quickly, softly / gently, alluringly.

↳ temporal: yesterday, Monday, last semester

Other POS tags in English

numerals: one, two, nth, first, second,

prepositions: of, on, over, under, to, from, around, "anywhere a hippo can go" Espanol el-marc.
determiners: some, a, an, the, this, that la-fam.
Indefinite los pl+
pronouns: she, he, it, they las pt+ who whom ever whatever.

closed classes

conjunctions: and or but
I really like the book that Sally wrote. when

auxiliary verbs: was hacked tense → passive voice

particles: preposition joins aspect → still going on?
↓ a verb polarity → negation mood → desirable / possible / suggested

< Notes



Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	+%, &
CD	cardinal number	<i>one, two</i>	TO	“to”	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential ‘there’	<i>there</i>	VB	verb base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VBN	verb past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, sing.	<i>IBM</i>	\$	dollar sign	\$
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	#
PDT	predeterminer	<i>all, both</i>	“	left quote	‘ or “
POS	possessive ending	<i>'s</i>	”	right quote	’ or ”
PRP	personal pronoun	<i>I, you, he</i>	(left parenthesis	[, (, {, <
PRP\$	possessive pronoun	<i>your, one's</i>)	right parenthesis],), }, >
RB	adverb	<i>quickly, never</i>	,	comma	,
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	. ! ?
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	: ; ... --
RP	particle	<i>up, off</i>			

The/DT grand/JJ jury/NN commented/VBD on/IN a/DT number/NN of/IN other/JJ topics/NNS ./.

There/EX are/VBP 70/CD children/NNS there/RB

Preliminary/JJ findings/NNS were/VBD reported/VBN in/IN today/NN 's/POS New/NNP England/NNP Journal/NNP of/IN Medicine/NNP ./.

Corpora with manual POS tags:

- Brown corpus - 1 million words of 500 written English texts from different genres
- WSJ corpus - 1 million words from the Wall Street Journal



Notes



Corpora with manual POS tags:

- Brown corpus - 1 million words of 500 written English texts from different genres
- WSJ corpus - 1 million words from the Wall Street Journal
- Switchboard corpus - 2 M words of telephone conversations

POS Tagging

words are ambiguous

so tagging must resolve / disambiguate

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

Figure 10.2 The amount of tag ambiguity for word types in the Brown and WSJ corpora, from the Treebank-3 (45-tag) tagging. These statistics include punctuation as words, and assume words are kept in their original case.

Some words have up to 6 tags

