

PPOL 628: Text as Data – Computational Linguistics for Social Scientists

Class 4: Parts of Speech and Phrases

Today

- Lecture: key points from readings
- Reading discussion
- Theory assignment questions
- Lab: partsofspeech.R
- Website: github.com/matthewjdenny/PPOL_628_Text_As_Data

Parts of Speech

- Grammatical Categories for words \leftrightarrow parts of speech.
- Nouns: typically refer to people, animals, concepts and things.
 - Cat, dog, watermelon, car, boat, chair, Matt, Susan, space, time.
- Verbs: typically express the action in a sentence.
 - Run, throw, treat, show, write, walk, climb
- Adjectives: Describe the properties of nouns:
 - Fast, green, interesting, small, hard, frequent
- “Substitution test” – determine which words belong in the same class.

Coarse tag	PTB tags	Examples
N: Nouns	NN: Common nouns (singular) NNS: Common nouns (plural) NNP: Proper nouns (singular) NNPS: Proper nouns (plural) FW: Foreign words CD: Numbers	paragraph, adoption, member, extension, exploration barriers, additions, objects, policies, negotiations Advanced, Assessment, Notice, Contents, Injury AUTHORIZATIONS, Limitations, Indians, Presidents novo, parte, pima, de, tempore, officio, pro, bona 48, 632, 1834, 2009, 1129, 1302, 381, 586, 810
A: Adjectives	JJ: Adjectives and ordinals JJR: Adjectives (comparative) JJS: Adjectives (superlative) VBG: Gerunds, present participles RB: Adverbs	renewable, following, other, scientific, subsequent, last younger, higher, Higher, More, less, smaller, earlier latest, highest, greatest, Best, least, largest setting, resulting, being, working, operating, beginning respectively, generally, forth, previously, so, no, fully
D: Determiners	DT	this, either, any (<i>Most common:</i> the, a, this)
P: Prepositions	IN: Most prepositions TO: The word “to”	<i>Most common words:</i> of, in, for, by, under, as, with to, To, TO
V: Verbs	VB: Verbs (base form) VBD: Past tense VBN: Past participle VBP: Present tense (non-3rd sing) VBZ: Present tense (3rd sing)	itemize, supply, TERMINATE, guarantee, concentrate overestimated, trained, expired, GENERATED, switched eliminated, intercepted, owed, advertised, Incorporated mitigate, nullify, Benefit, insert, fulfill, produce, seize distributes, announces, directs, respects, upholds, uses
M: Verb Modifiers	RB: Adverbs (base form) RBR: Comparative adverbs RBS: Superlative adverbs RP: Particle adverbs MD: Modal auxiliary verbs	extremely, hard, rapidly, after, now better, faster, slower, easier, shorter best, worst, fastest, slowest, easiest about, off, on, up can, should, might, musn't
C: Coord. Conj.	CC: Coordinating conjunctions	and, or, but

Part of Speech Tagging

- Process of assigning a part of speech (POS) “tag” to each word in a document.
- Canonical training set of POS tags (in English) is the Penn Treebank (1999) which is maintained by the Linguistic Data Consortium:
 - <https://catalog.ldc.upenn.edu/LDC99T42>
- In practice, this means employing either hand coding, a heuristic approach, a maximum likelihood model, a neural net, or some combination to assign POS tags to terms.
- Quality of POS taggers is dependent on training data
 - Challenges in other languages, Twitter, etc.

Part of Speech Tagging

Should a Federal agency seek to restrict photography of its installations or personnel, it shall obtain a court order that outlines the national security or other reasons for the restriction.



Should/M a/D **Federal/N agency/N** seek/V to/T **restrict/V**
photography/N of/P its/PR installations/N or/C
personnel/N, it/PR shall/M **obtain/V** a/D **court/N order/N**
that/d outlines/V the/D **national/A security/N** or/C
other/A reasons/N for/P the/D restriction/N.

Tag Definitions

A = adjective, N = noun, V = verb, M = modal, D = determiner,
C = conjunction, PR = pronoun, T = to

Part of Speech Tagging

Peen Treebank (45-tag corpus)		
Unambiguous (1 tag)	38,857	(81%)
Ambiguous (2-7 tags)	8,844	(19%)
Details: 2 tags	6,731	
3 tags	1,621	
4 tags	357	
5 tags	90	
6 tags	32	
7 tags	6	<i>(well, set, round, open, fit, down)</i>
8 tags	4	<i>('s, half, back, a)</i>
9 tags	3	<i>(that, more, in)</i>

A simple approach which assigns only the most common tag to each word performs with 90% accuracy!

Syntax → Phrases

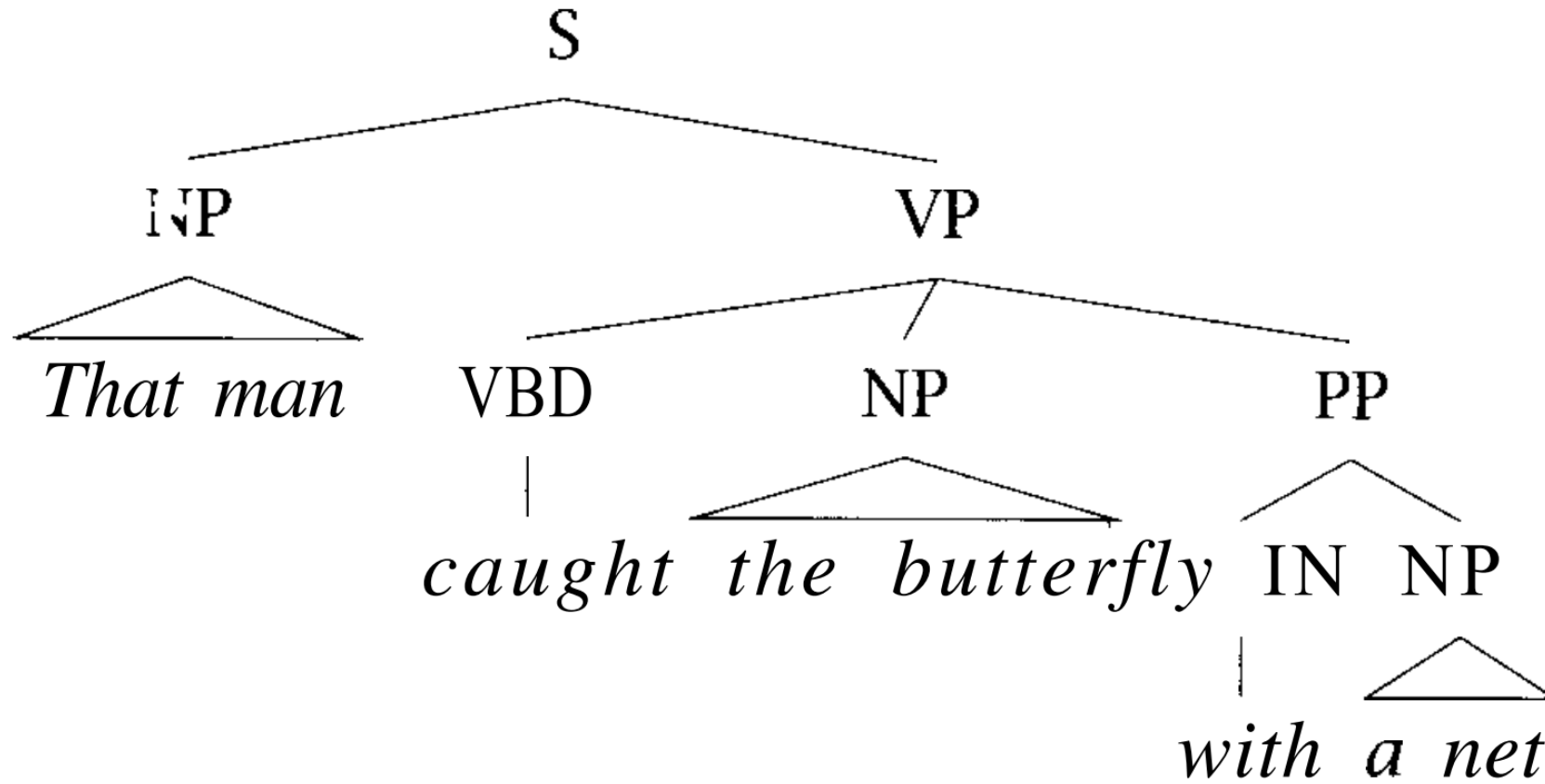
- **Syntax** is the study of the regularities and constraints of word order and phrase structure.
- **Phrases** are *syntactically coherent* groupings of words.

She		him
the woman		the man
the tall woman		the short man
the very tall woman	saw	the very short man
the tall woman with sad eyes		the short man with red hair
...		...

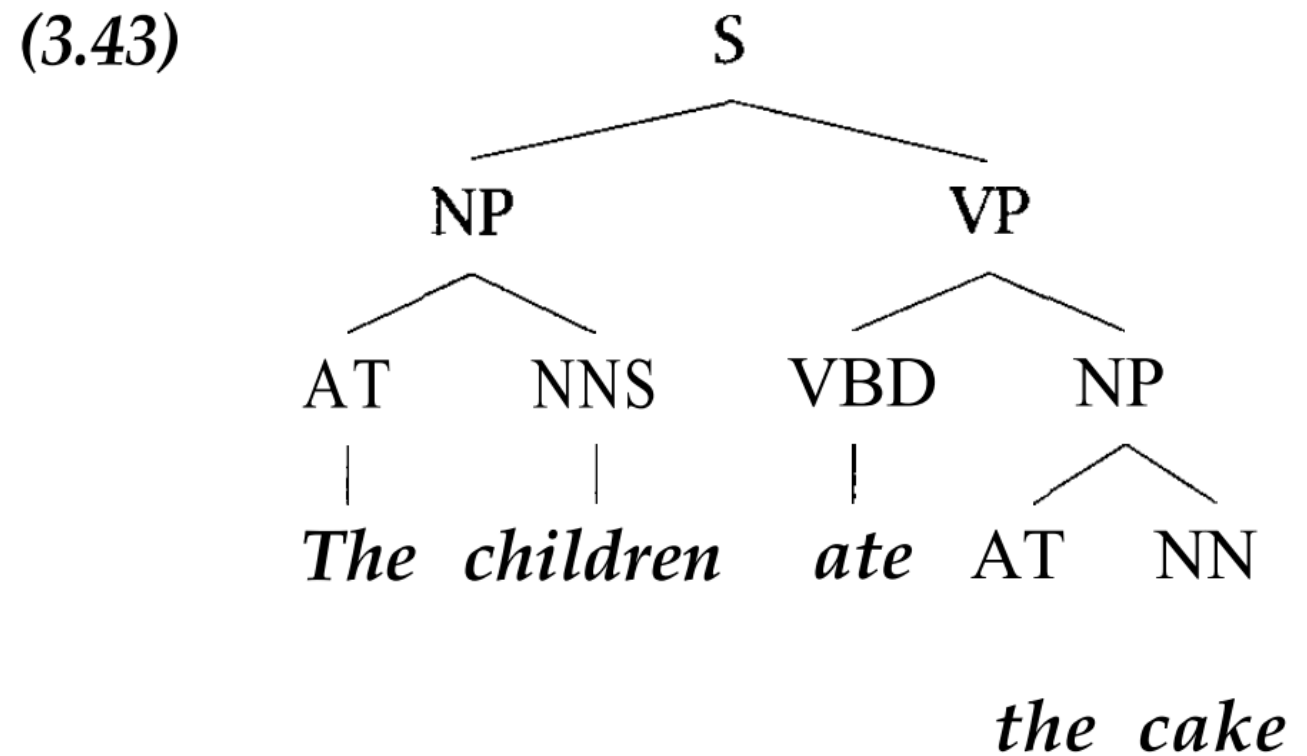
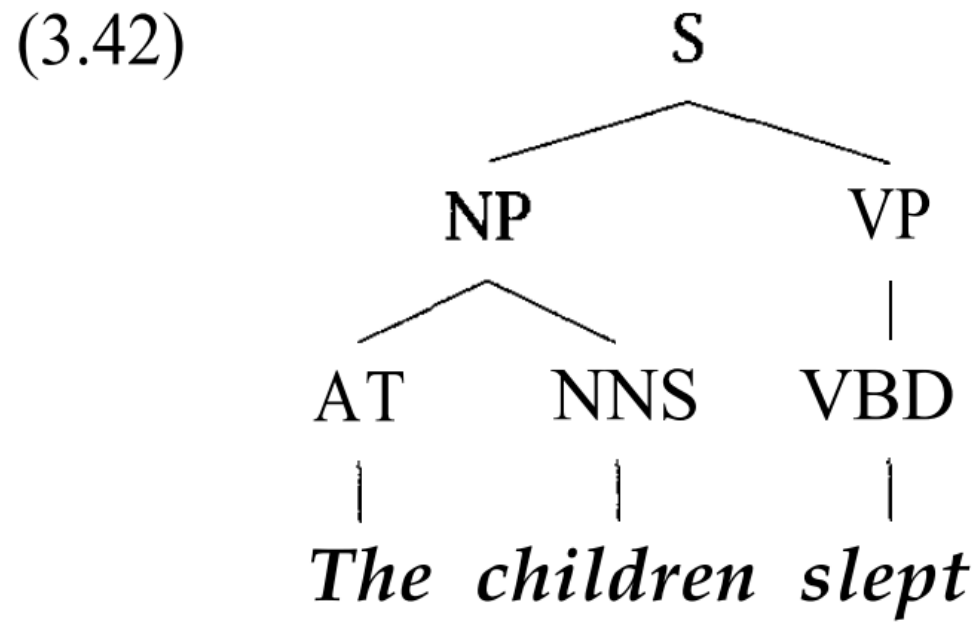
Phrases

- **Noun phrases.** A noun is usually embedded in a phrase a syntactic unit of the sentence in which information about the noun is gathered. The noun is the head of the noun phrase, the central constituent that determines the syntactic character of the phrase. Noun phrases are usually the arguments of verbs, the participants in the action, activity or state described by the verb.
- **Verb phrases.** Analogous to the way nouns head noun phrases, the verb is the head of the verb phrase (VP). In general, the verb phrase organizes all elements of the sentence that depend syntactically on the verb

Tree Structure of Sentences



Tree Structure of Sentences



Justeson and Katz (1995) – TERMS

- Technical terminology (the kind of stuff you would put in an index) tends to consist of noun phrases).
- Simple heuristic algorithm to find these noun phrases efficiently in a document.

AN: linear function; lexical ambiguity; mobile phase

NN: regression coefficients; word sense; surface area

AAN: Gaussian random variable; lexical conceptual paradigm; aqueous mobile phase

ANN: cumulative distribution function; lexical ambiguity resolution; accessible surface area

NAN: mean squared error; domain independent set; silica based packing

NNN: class probability function; text analysis system; gradient elution chromatography

NPN: degrees of freedom; [*no example*]; energy of adsorption

POS Tag Patterns → Noun Phrases

Tag Pattern	Example
AN	<i>equal employment</i>
NN	<i>research project</i>
AAN	<i>local educational agency</i>
ANN	<i>recreational land resource</i>
NAN	<i>health related service</i>
NNN	<i>health care provider</i>
NPN	<i>election by majority</i>

(Justeson and Katz, 1995)

POS Tag Patterns → Verb Phrases

Tag Pattern	Example
VN	<i>reduce funding</i>
VAN	<i>encourage dissenting members</i>
VNN	<i>restrict government agencies</i>
VPN	<i>prescribe in paragraph</i>
ANV	<i>eligible employee means</i>
VDN	<i>establish a commission</i>

Regular Expressions over POS Tags

- **Noun Phrases:** $(A|N)^*N(PD^*(A|N)^*N)^*$
 - Zero or more adjectives or nouns, followed by a noun, followed (optionally) by zero or more groups of terms containing a preposition and zero or more determiners, then zero or more adjectives or nouns, and ending in a noun.
- **Verb Groups:** $(M(CM)^*|V)^*V(M(CM)^*|V)^*$
 - Modifier followed by zero or more coordinating conjunction-modifier pairs, or a verb, all repeated zero or more times, then a verb, then a modifier followed by zero or more coordinating conjunction-modifier pairs, or a verb, all repeated zero or more times.

Full Generalization

$$\begin{aligned}\mathbf{NP_w_Coord} &= \mathbf{NP} \ (C \ \text{Det}^* \ \mathbf{NP})^* \\ &= (A(CA)^*I(N)^*N((P(CP)^*)+(D(CD)^*)^*(A(CA)^*I(N)^*N)^*(C(D(CD)^*)^* \\ &\quad (A(CA)^*I(N)^*N((P(CP)^*)+(D(CD)^*)^*(A(CA)^*I(N)^*N)^*)^*)^*\end{aligned}$$

$$\begin{aligned}\mathbf{Verb_Argument} &= (\mathbf{Subject_Verb} \mid \mathbf{Verb_Object} \mid \mathbf{Verb_Prep_Phrase} \mid \mathbf{NP_Verb_Phrase}) \\ &= ((A(CA)^*I(N)^*N((P(CP)^*)+(D(CD)^*)^*(A(CA)^*I(N)^*N)^*(C(D(CD)^*)^*(A(CA)^*I(N)^*N \\ &\quad ((P(CP)^*)+(D(CD)^*)^*(A(CA)^*I(N)^*N)^*)^*(P(CP)^*)^*(M(CM)^*I(V)^*V(M(CM)^*I(V)^* \\ &\quad (C(M(CM)^*I(V)^*V(M(CM)^*I(V)^*)^*I(M(CM)^*I(V)^*V(M(CM)^*I(V)^*(C(M(CM)^*I(V)^* \\ &\quad V(M(CM)^*I(V)^*)^*(D(CD)^*)^*(A(CA)^*I(N)^*N((P(CP)^*)+(D(CD)^*)^*(A(CA)^*I(N)^*N)^* \\ &\quad (C(D(CD)^*)^*(A(CA)^*I(N)^*N((P(CP)^*)+(D(CD)^*)^*(A(CA)^*I(N)^*N)^*)^*I(M(CM)^*I(V)^* \\ &\quad V(M(CM)^*I(V)^*(C(M(CM)^*I(V)^*V(M(CM)^*I(V)^*)^*((P(CP)^*)+(D(CD)^*)^*(A(CA)^*I(N)^*N)+ \\ &\quad I(A(CA)^*I(N)^*N((P(CP)^*)+(D(CD)^*)^*(A(CA)^*I(N)^*N)^*(C(D(CD)^*)^*(A(CA)^*I(N)^*N \\ &\quad ((P(CP)^*)+(D(CD)^*)^*(A(CA)^*I(N)^*N)^*)^*(P(CP)^*)^*((M(CM)^*I(V)^*V(M(CM)^*I(V)^* \\ &\quad (C(M(CM)^*I(V)^*V(M(CM)^*I(V)^*)^*(D(CD)^*)^*(A(CA)^*I(N)^*N((P(CP)^*)+(D(CD)^*)^* \\ &\quad (A(CA)^*I(N)^*N)^*(C(D(CD)^*)^*(A(CA)^*I(N)^*N((P(CP)^*)+(D(CD)^*)^*(A(CA)^*I(N)^*N)^*)^* \\ &\quad I(M(CM)^*I(V)^*V(M(CM)^*I(V)^*(C(M(CM)^*I(V)^*V(M(CM)^*I(V)^*)^*((P(CP)^*)+(D(CD)^*)^* \\ &\quad (A(CA)^*I(N)^*N)+))\end{aligned}\tag{14}$$

$$\mathbf{Phrases} = (\mathbf{NP_w_Coord} \mid \mathbf{Verb_Argument})$$

The Readings This Week

- Manning & Schtze, H. (1999). Foundations of Statistical Natural Language Processing. Chapter 3.
- Justeson, J. S., & Katz, S. M. (1995). Technical terminology: some linguistic properties and an algorithm for identification in text. Natural Language Engineering, 1(01).
- Handler, A., Denny, M. J., Wallach, H., & OConnor, B. (2016). Bag of What? Simple Noun Phrase Extraction for Text Analysis. EMNLP + CSS