

## Advanced NLP Processing

Beyond tokens...

1

Stanford NLP

## NLP Text Processing Pipeline

nlTK covers POS tagging, phrase chunking  
Stanford NLP toolkit provides parsing, coreference, NER

- Document → Sections and Paragraphs
- Paragraphs → Sentences (sentence segmentation / extraction)
- Sentences → Tokens
- Tokens → Lemmas or Morphological Variants / Stems
- Tokens → Part-of-speech (POS) Tags
- Tokens, POS Tags → Phrase Chunks (Named entities and Keyphrases)
- Tokens, POS Tags → Parse Trees
  - Augment above with coreference, entailment, sentiment, ...

2



## Part-of-speech tagging

A simple but useful form of linguistic analysis

Christopher Manning

3

Stanford NLP

## Parts of Speech

- Perhaps starting with Aristotle in the West (384–322 BCE), there was the idea of having parts of speech
  - a.k.a lexical categories, word classes, “tags”, POS
- It comes from Dionysius Thrax of Alexandria (c. 100 BCE) the idea that is still with us that there are 8 parts of speech
  - But actually his 8 aren’t exactly the ones we are taught today
    - Thrax: noun, verb, article, adverb, preposition, conjunction, participle, pronoun
    - School grammar: noun, verb, adjective, adverb, preposition, conjunction, pronoun, interjection

4

Open class (lexical) words			
Nouns		Verbs	Adjectives <i>old older oldest</i>
Proper <i>IBM Italy</i>	Common <i>cat / cats snow</i>	Main <i>see registered</i>	Adverbs <i>slowly</i>
			Numbers <i>122,312 one</i> ... more
Closed class (functional)		Modals	Prepositions <i>to with</i>
Determiners <i>the some</i>		<i>can had</i>	Particles <i>off up</i> ... more
Conjunctions <i>and or</i>			Interjections <i>Ow Eh</i>
Pronouns <i>he its</i>			

Stanford NLP

## Open vs. Closed classes

- Open vs. Closed classes
  - Closed:
    - determiners: *a, an, the*
    - pronouns: *she, he, I*
    - prepositions: *on, under, over, near, by, ...*
    - Why “closed”?
  - Open:
    - Nouns, Verbs, Adjectives, Adverbs.

6





## Phrase Chunking

- Find all non-recursive noun phrases (NPs) and verb phrases (VPs) in a sentence.
  - [NP I] [VP ate] [NP the spaghetti] [PP with] [NP meatballs].
  - [NP He] [VP reckons] [NP the current account deficit] [VP will narrow] [PP to] [NP only # 1.8 billion] [PP in] [NP September]

13



## Phrase Chunking as Sequence Labeling

- Tag individual words with one of 3 tags
  - B (Begin) word starts new target phrase
  - I (Inside) word is part of target phrase but not the first word
  - O (Other) word is not part of target phrase
- Sample for NP chunking
  - He reckons the current account deficit will narrow to only # 1.8 billion in September.

Begin Inside Other

14



## Named Entity Recognition (NER)

- A special class of **Proper Noun Phrases**
- People:** Scott Sanner, President Obama, Madonna
- Places:** New York, Madison Square Garden, Millenium Park
- Organizations:** New York Times, University of Toronto

15



## Keyphrases

- Useful noun phrases, but not necessarily Proper Nouns, e.g.,
  - "machine learning"
  - "support vector machines"
  - "genetically modified organisms"
- A subset of frequent noun phrases (harder to extract than NEs)
  - This paper has the best method I've found so far: "Automatic Recognition of Multi-Word Terms: the C-value/NC-value Method Katerina Frantziy, Sophia Ananiadouy, Hideki Mima" IJODL 2000.

<http://personalpages.manchester.ac.uk/staff/sophia.ananiadou/ijodl2000.pdf>

16



## Statistical Natural Language Parsing

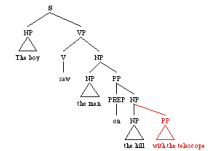
Parsing: Two views of syntactic structure

17



## Why parsing?

- "The boy saw the man on the hill with the telescope."
  - Who had the telescope?
- Depends on whether you attach "with the telescope" to "I" or "man on the hill"
- How do you determine attachments? Parsing.
  - Some sentences are inherently ambiguous: attachment ambiguity.



18



## For fun

- Who polices the police?
- Police police police police.
- Who polices the police police? ☺
- Point: we need more than word order / POS to interpret sentences... we need structure.

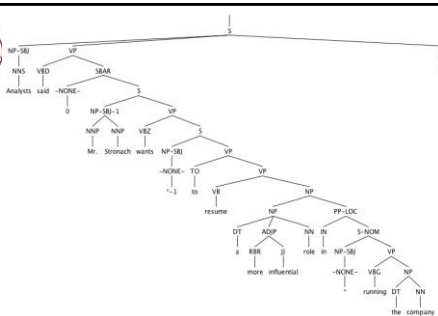
19



## Two views of linguistic structure: 1. Constituency (phrase structure)

- Phrase structure organizes words into nested constituents.
- What is a **constituent**?
- Constituent behaves as unit that can appear in different places:
  - John talked [to the children] [about drugs].
  - John talked [about drugs] [to the children].
- Substitution/expansion/pro-forms:
  - I sat [on the box/right on top of the box/there].
- Coordination, regular internal structure, no intrusion, fragments, semantics, ...

20



21



## Grammars for Parse Tree Production

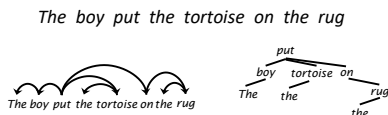
- Parent  $\rightarrow$  Child1 Child2 | Child3 Child4 ... | ...
- $S \rightarrow NP VP \mid \dots$
- $NP \rightarrow \dots NN^* \dots$
- $VP \rightarrow \dots VB^* \dots$
- $ADJP \rightarrow \dots JJ^* \dots$
- $ADVP \rightarrow \dots RB^* \dots$
- ...

22



## Two views of linguistic structure: 2. Dependency structure

- Dependency structure shows which words depend on (modify or are arguments of) which other words.



23



## Two views of linguistic structure: 2. Dependency structure

- Dependency structure shows which words depend on (modify or are arguments of) which other words.
- Can derive dependency tree from parse tree
- What about reverse?

24

## 26

## 27

## 28