

DSC291: Advanced Statistical Natural Language Processing

Sequence Tagging
Parsing

Zhiting Hu

Lecture 11, May 3, 2022

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

- Sequence Tagging/Labeling
- Parsing

Sequence Labeling

[Slides adapted from UW CSE 447 by Noah Smith]

Recap: Sequence Labeling

Problem statement: given a sequence of n words \mathbf{x} , assign each a label from \mathcal{L} . Let $L = |\mathcal{L}|$.

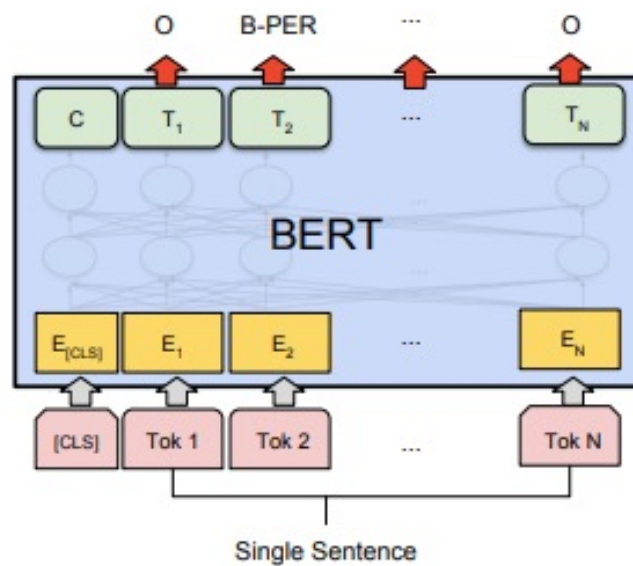
Every approach we see today will cast the problem as:

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{L}^n} \operatorname{Score}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta})$$

Naïvely, that's a classification problem where the number of possible 'labels' (output sequences) depends on the input and is $O(L^n)$ in size!

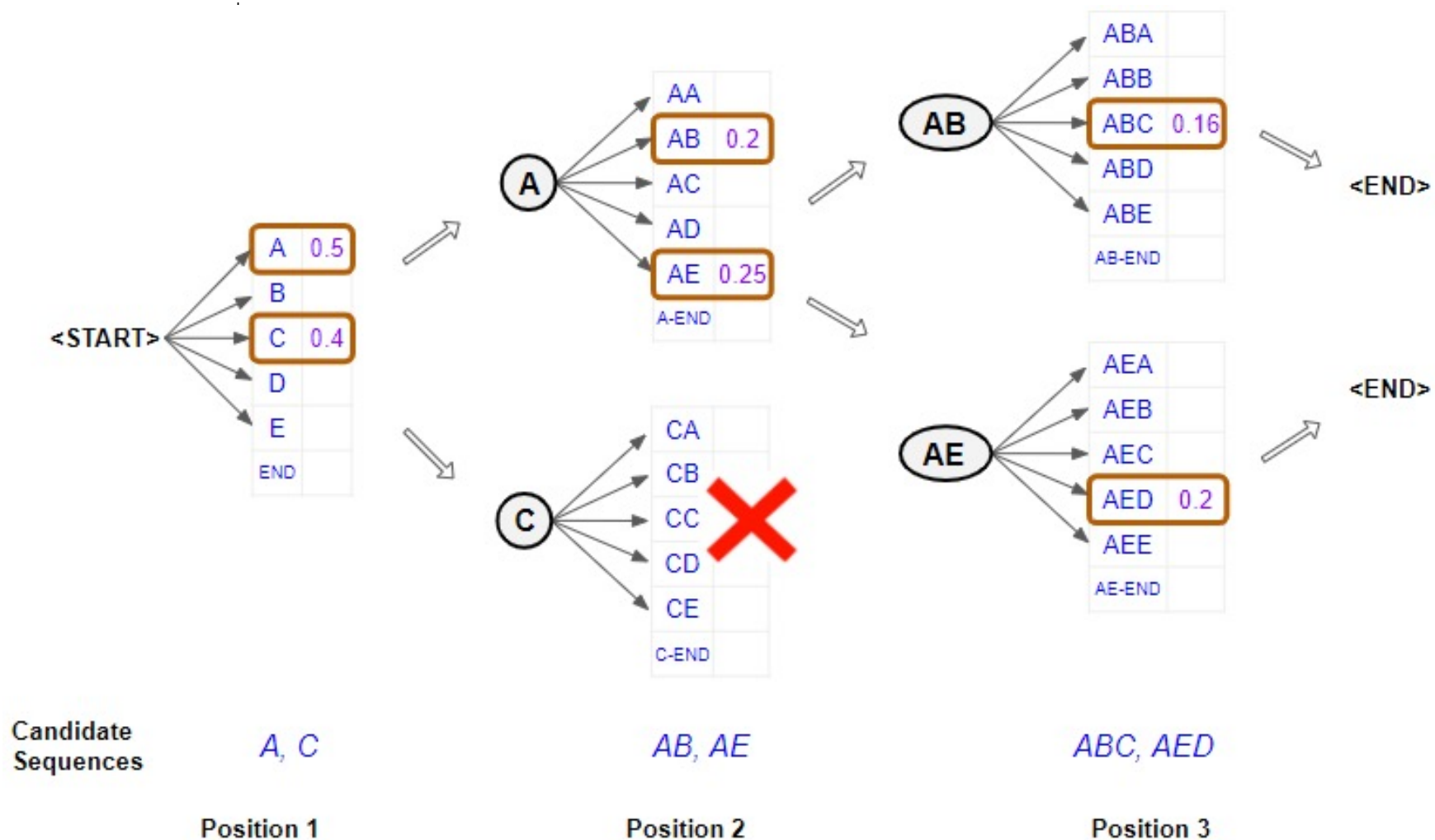
Recap: Sequence Labeling

Model	0	1	2	3	4
Score decomp.	$s(\mathbf{x}, i, y_i)$	$s(\mathbf{x}, i, \mathbf{y}_{1:i})$	emission/ transition	$s(x_i, y_i) +$ $s(y_i, y_{i+1})$	$s(\mathbf{x}, i, y_i, y_{i+1})$
learn	SGD	?	count & normalize	?	?
decode	local	beam search	Viterbi	Viterbi	Viterbi



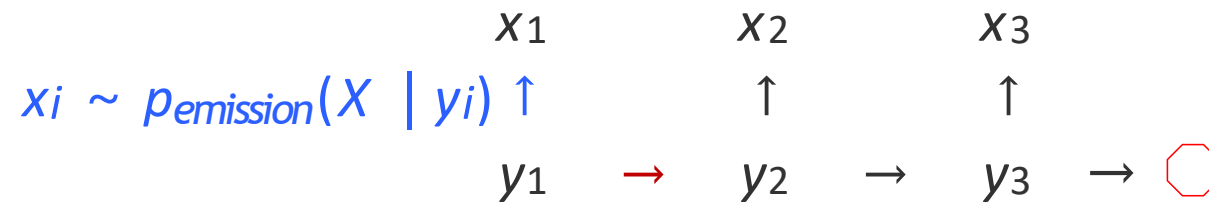
model v.0

Recap: Sequence Labeling



Recap: Sequence Labeling

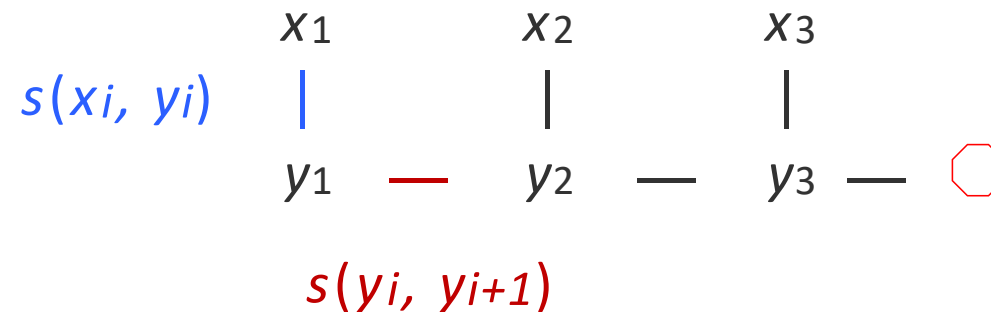
Model	0	1	2	3	4
Score decomp.	$s(\mathbf{x}, i, y_i)$	$s(\mathbf{x}, i, \mathbf{y}_{1:i})$	emission/ transition	$s(x_i, y_i) +$ $s(y_i, y_{i+1})$	$s(\mathbf{x}, i, y_i, y_{i+1})$
learn	SGD	?	count & normalize	?	?
decode	local	beam search	Viterbi	Viterbi	Viterbi



$$y_{i+1} \sim p_{\text{transition}}(Y \mid y_i)$$

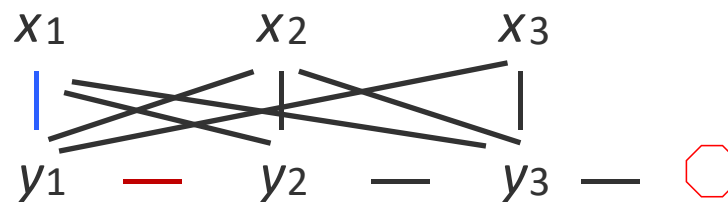
Recap: Sequence Labeling

Model	0	1	2	3	4
Score decomp.	$s(\mathbf{x}, i, y_i)$	$s(\mathbf{x}, i, \mathbf{y}_{1:i})$	emission/ transition	$s(x_i, y_i) +$ $s(y_i, y_{i+1})$	$s(\mathbf{x}, i, y_i, y_{i+1})$
learn	SGD	?	count & normalize	?	?
decode	local	beam search	Viterbi	Viterbi	Viterbi



Recap: Sequence Labeling

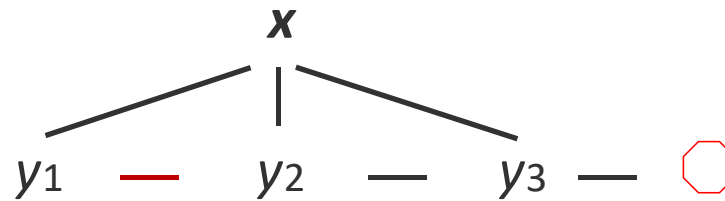
Model	0	1	2	3	4
Score decomp.	$s(\mathbf{x}, i, y_i)$	$s(\mathbf{x}, i, \mathbf{y}_{1:i})$	emission/ transition	$s(x_i, y_i) +$ $s(y_i, y_{i+1})$	$s(\mathbf{x}, i, y_i, y_{i+1})$
learn	SGD	?	count & normalize	?	?
decode	local	beam search	Viterbi	Viterbi	Viterbi



$$s(\mathbf{x}, y_i, y_{i+1})$$

Recap: Sequence Labeling

Model	0	1	2	3	4
Score decomp.	$s(\mathbf{x}, i, y_i)$	$s(\mathbf{x}, i, \mathbf{y}_{1:i})$	emission/ transition	$s(x_i, y_i) +$ $s(y_i, y_{i+1})$	$s(\mathbf{x}, i, y_i, y_{i+1})$
learn	SGD	?	count & normalize	?	?
decode	local	beam search	Viterbi	Viterbi	Viterbi



$$s(\mathbf{x}, y_i, y_{i+1})$$

Recap: Inference (Decoding) -- Viterbi Algorithm

Let $\heartsuit_i(y)$ be the score of the best label sequence for $x_{1:i}$ that ends in y . It is defined recursively:

$$\heartsuit_{n+1}(\heartsuit) = \max_{y_n \in \mathcal{L}} s(\mathbf{x}, n, y_n, \heartsuit) + \boxed{\heartsuit_n(y_n)}$$

$$\heartsuit_n(y) = \max_{y_{n-1} \in \mathcal{L}} s(\mathbf{x}, n-1, y_{n-1}, y) + \boxed{\heartsuit_{n-1}(y_{n-1})}$$

$$\heartsuit_{n-1}(y) = \max_{y_{n-2} \in \mathcal{L}} s(\mathbf{x}, n-2, y_{n-2}, y) + \boxed{\heartsuit_{n-2}(y_{n-2})}$$


⋮

$$\heartsuit_i(y) = \max_{y_{i-1} \in \mathcal{L}} s(\mathbf{x}, i-1, y_{i-1}, y) + \boxed{\heartsuit_{i-1}(y_{i-1})}$$

⋮

$$\heartsuit_1(y) = s(\mathbf{x}, 0, \heartsuit, y)$$

Viterbi Procedure (Part I: Prefix Scores)


		input sequence				
		x_1	x_2	\dots	x_n	
\mathcal{L}	l_1					
	l_2					
	\vdots					
	l_L					
						

Viterbi Procedure (Part I: Prefix Scores)

		input sequence				
		x_1	x_2	\dots	x_n	
\mathcal{L}	l_1	$\heartsuit_1(l_1)$				
	l_2	$\heartsuit_1(l_2)$				
	\vdots					
	l_L	$\heartsuit_1(l_L)$				
	\circ					


$$\heartsuit_1(y) = s(\mathbf{x}, 0, \circ, y)$$

Viterbi Procedure (Part I: Prefix Scores)

		input sequence				
		x_1	x_2	\dots	x_n	
\mathcal{L}	l_1	$\heartsuit_1(l_1)$	$\heartsuit_2(l_1)$			
	l_2	$\heartsuit_1(l_2)$	$\heartsuit_2(l_2)$			
	\vdots					
	l_L	$\heartsuit_1(l_L)$	$\heartsuit_2(l_L)$			
						

$$\heartsuit_i(y) = \max_{y_{i-1} \in \mathcal{L}} s(\mathbf{x}, i-1, y_{i-1}, y) + \boxed{\heartsuit_{i-1}(y_{i-1})}$$

Viterbi Procedure (Part I: Prefix Scores)

		input sequence				
		x_1	x_2	\dots	x_n	
\mathcal{L}	l_1	$\heartsuit_1(l_1)$	$\heartsuit_2(l_1)$		$\heartsuit_n(l_1)$	
	l_2	$\heartsuit_1(l_2)$	$\heartsuit_2(l_2)$		$\heartsuit_n(l_2)$	
	\vdots					
	l_L	$\heartsuit_1(l_L)$	$\heartsuit_2(l_L)$		$\heartsuit_n(l_L)$	
						

$$\heartsuit_n(y) = \max_{y_{n-1} \in \mathcal{L}} s(\mathbf{x}, n-1, y_{n-1}, y) + \boxed{\heartsuit_{n-1}(y_{n-1})}$$

Viterbi Procedure (Part I: Prefix Scores)

		input sequence				
		x_1	x_2	\dots	x_n	
\mathcal{L}	l_1	$\heartsuit_1(l_1)$	$\heartsuit_2(l_1)$		$\heartsuit_n(l_1)$	
	l_2	$\heartsuit_1(l_2)$	$\heartsuit_2(l_2)$		$\heartsuit_n(l_2)$	
	\vdots					
	l_L	$\heartsuit_1(l_L)$	$\heartsuit_2(l_L)$		$\heartsuit_n(l_L)$	
	\bigcirc					$\heartsuit_{n+1}(\bigcirc)$

$$\heartsuit_{n+1}(\bigcirc) = \max_{y_n \in \mathcal{L}} s(\mathbf{x}, n, y_n, \bigcirc) + \boxed{\heartsuit_n(y_n)}$$

Recap: Inference (Decoding) -- Viterbi Algorithm

Input: scores $s(\mathbf{x}, i, y, y')$, for all $i \in \{0, \dots, n\}$, $y, y' \in \mathcal{L}$

Output: \hat{y}

1. Base case: $\heartsuit_1(y) = s(\mathbf{x}, 0, \bigcirc, y)$
2. For $i \in \{2, \dots, n + 1\}$:
 - ▶ Solve for $\heartsuit_i(*)$ and $\text{bp}_i(*)$.

$$\heartsuit_i(y) = \max_{y_{i-1} \in \mathcal{L}} s(\mathbf{x}, i - 1, y_{i-1}, y) + \heartsuit_{i-1}(y_{i-1}),$$

$$\text{bp}_i(y) = \operatorname{argmax}_{y_{i-1} \in \mathcal{L}} s(\mathbf{x}, i - 1, y_{i-1}, y) + \heartsuit_{i-1}(y_{i-1})$$

(At $n + 1$ we're only interested in $y = \bigcirc$.)

3. $\hat{y}_{i+1} \leftarrow \bigcirc$
4. For $i \in \{n, \dots, 1\}$:
 - ▶ $\hat{y}_i \leftarrow \text{bp}_{i+1}(\hat{y}_{i+1})$

Recap: Inference (Decoding) -- Viterbi Algorithm

- Viterbi Asymptotics

		input sequence			
		x_1	x_2	\dots	x_n
labels in \mathcal{L}	ℓ_1				
	ℓ_2				
	\vdots				
	ℓ_L				

Space: need to store s , and fill in the cells above. $O(nL^2)$ for s (in the most general case, often less), $O(nL)$ for cells

Runtime: each cell requires an “argmax.” $O(nL^2)$

Back to s

We haven't said much about the function that scores candidate label pairs at different positions, $s(\mathbf{x}, i, y, y')$.

This function is very important; two common choices are:

- ▶ Expert-designed, task-specific features $\mathbf{f}(\mathbf{x}, i, y, y')$ and weights θ
- ▶ A neural network that encodes x_i in context, y_i , and y_{i+1} and gives back a goodness score

Either way, let θ denote the parameters of s . From now on, we'll use $s(\mathbf{x}, i, y, y'; \theta)$ and $\text{Score}(\mathbf{x}, \mathbf{y}; \theta)$ to emphasize that “ s ” is a function of parameters θ we need to estimate.

Probabilistic View of Learning

As we've done before, we start with the principle of maximum likelihood to estimate θ :

$$\begin{aligned}\theta^* &= \arg \max_{\theta \in \mathbb{R}^d} \prod_{t=1}^T p(\mathbf{Y}_t = \mathbf{y} \mid \mathbf{X}_t = \mathbf{x}; \theta) \\ &= \arg \max_{\theta \in \mathbb{R}^d} \sum_{t=1}^T \log p(\mathbf{Y}_t = \mathbf{y} \mid \mathbf{X}_t = \mathbf{x}; \theta) \\ &= \arg \min_{\theta \in \mathbb{R}^d} \sum_{t=1}^T \underbrace{-\log p(\mathbf{Y}_t = \mathbf{y} \mid \mathbf{X}_t = \mathbf{x}; \theta)}_{\text{sometimes called "log loss" or "cross entropy"}}$$

Next, we'll drill down into " $p(\mathbf{Y} = \mathbf{y}_t \mid \mathbf{X} = \mathbf{x}_t; \theta)$."

Conditional Random Fields (CRFs)

Lafferty et al. (2001)

CRFs are a tremendously influential model that generalizes multinomial logistic regression to structured outputs like sequences.

$$p_{\text{CRF}}(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta}) = \frac{\exp \text{Score}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta})}{Z(\mathbf{x}; \boldsymbol{\theta})}$$

$$Z(\mathbf{x}; \boldsymbol{\theta}) = \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} \exp \text{Score}(\mathbf{x}, \mathbf{y}'; \boldsymbol{\theta})$$

$$-\log p_{\text{CRF}}(\mathbf{y} \mid \mathbf{x}; \boldsymbol{\theta}) = \underbrace{-\text{Score}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta})}_{\text{“hope”}} + \underbrace{\log Z(\mathbf{x}; \boldsymbol{\theta})}_{\text{“fear”}}$$

So, our “CRF”:

- ▶ Uses Viterbi for decoding (our v. 4 sequence labeler)
- ▶ Trains parameters to maximize likelihood (like MLR and NNs)

Sequence-Level Log Loss

Here's the maximum likelihood learning problem (equivalently, sequence-level log loss):

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^d} \sum_{t=1}^T -\operatorname{Score}(\mathbf{x}_t, \mathbf{y}_t; \boldsymbol{\theta}) + \log Z(\mathbf{x}_t; \boldsymbol{\theta})$$

If we can calculate and differentiate (w.r.t. $\boldsymbol{\theta}$) the Score and Z functions, we can use SGD to learn.

Calculating $Z(x; \theta)$

Good news! The algorithm that gives us Z is *almost exactly like* the Viterbi algorithm.

Forward algorithm: sums the `expScore` values for all label sequences, given x , in the same asymptotic time and space as Viterbi.

Let $\alpha_i(y)$ be the sum of all (exponentiated) scores of label prefixes of length i , ending in y .

Some Algebra

Given the decomposition

$$\text{Score}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}) = \sum_{i=0}^n s(\mathbf{x}, i, y_i, y_{i+1}; \boldsymbol{\theta}),$$

it holds that

$$\exp \text{Score}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}) = \prod_{i=0}^n e^{s(\mathbf{x}, i, y_i, y_{i+1}; \boldsymbol{\theta})},$$

and therefore

$$Z(\mathbf{x}; \boldsymbol{\theta}) = \sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} \prod_{i=0}^n e^{s(\mathbf{x}, i, y'_i, y'_{i+1}; \boldsymbol{\theta})}$$

Forward Algorithm

Input: scores $s(\mathbf{x}, i, y, y'; \boldsymbol{\theta})$, for all $i \in \{0, \dots, n\}$, $y, y' \in \mathcal{L}$

Output: $Z(\mathbf{x}; \boldsymbol{\theta})$

1. Base case: $\alpha_1(y) = e^{s(\mathbf{x}, 0, \bigcirc, y; \boldsymbol{\theta})}$
2. For $i \in \{2, \dots, n + 1\}$:
 - ▶ Solve for $\alpha_i(*)$.

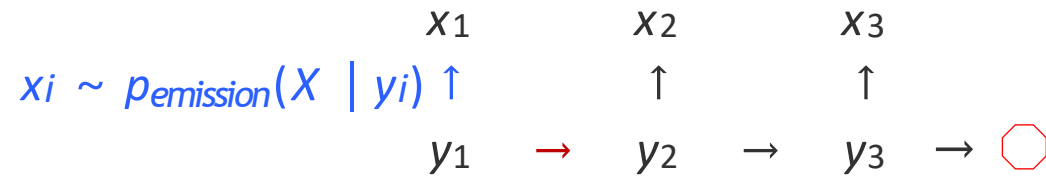
$$\alpha_i(y) = \sum_{y_{i-1} \in \mathcal{L}} e^{s(\mathbf{x}, i-1, y_{i-1}, y; \boldsymbol{\theta})} \times \alpha_{i-1}(y_{i-1})$$

(At $n + 1$ we're only interested in $y = \bigcirc$.)

3. Return $\alpha_{n+1}(\bigcirc)$, which is equal to $Z(\mathbf{x}; \boldsymbol{\theta})$.

Key Takeaways

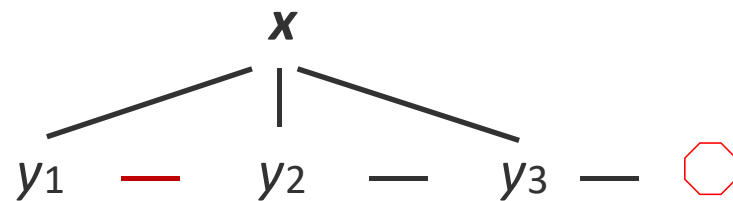
- HMM (generative model) - models joint distribution



$$p(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^N p(y_i | y_{i-1}) p(x_i | y_i)$$

$$y_{i+1} \sim p_{\text{transition}}(Y | y_i)$$

- CRF (discriminative model) - directly models conditional distribution



$$s(\mathbf{x}, y_i, y_{i+1})$$

$$p_{\text{CRF}}(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta}) = \frac{\exp \text{Score}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta})}{Z(\mathbf{x}; \boldsymbol{\theta})}$$

$$\text{Score}(\mathbf{x}, \mathbf{y}) = \sum_{i=0}^n s(\mathbf{x}, i, y_i, y_{i+1})$$

Key Takeaways

- HMM (generative model) - models joint distribution
- CRF (discriminative model) - directly models conditional distribution
- Inference: Viterbi algorithm
- Learning

Parsing

[Slides adapted from UW CSE 447 by Noah Smith; UCB Info 159/259 by David Bamman]

Motivation

As data, we tend to view natural language text as sequences (of words, wordpieces, or characters, depending on the NLP application).

But language obeys implicit rules of grammar, and it carries meaning.

- ▶ It's helpful to consider an analogy to programming languages, which have *syntax* and *semantics*; well-formed programs can be compiled and executed to carry out a task.
- ▶ Well-formed natural language strings can be understood by others.

Computational models that analyze natural language syntax and semantics typically map into **structures** like trees, graphs, and more.

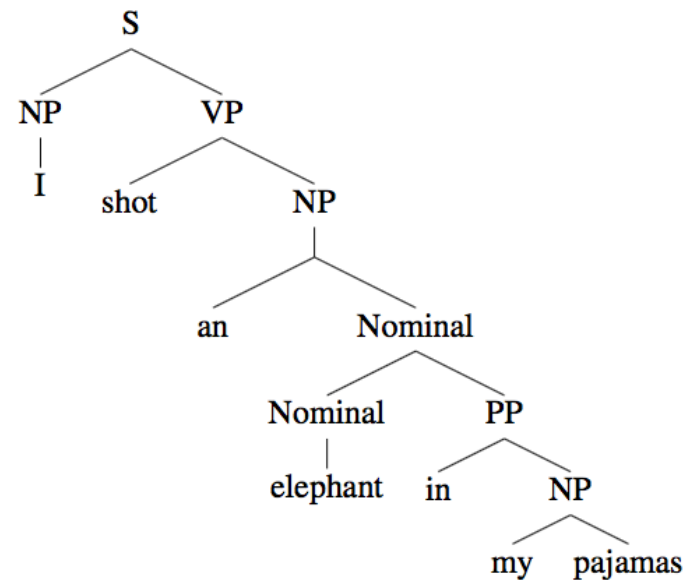
Linguistic Analysis

- ▶ **Syntax:** rules governing grammaticality or well-formedness of strings, relative to a language
- ▶ **Semantics:** how the meaning of an utterance is constructed, grounded in “the world” (or a proxy to the world)
- ▶ **Pragmatics:** the intended meaning by a speaker, in a given social context

Each has many theories, and none of them is complete!

Syntax

- With syntax, we're moving from labels for discrete items — documents (sentiment analysis), tokens (POS tagging, NER) — to the **structure** between items.



PRP VBD DT NN IN PRP\$ NNS

I shot an elephant in my pajamas

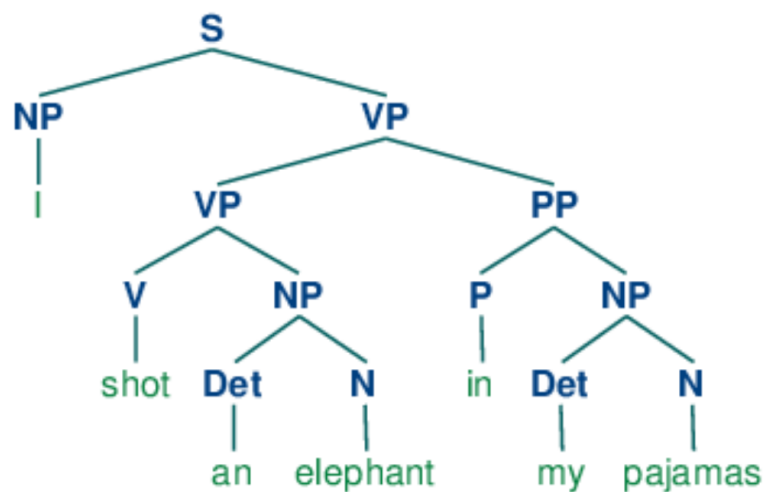
Syntax

- With syntax, we're moving from labels for discrete items — documents (sentiment analysis), tokens (POS tagging, NER) — to the **structure** between items.
- Syntax is fundamentally about the hierarchical structure of language and (in some theories) which sentences are **grammatical** in a language

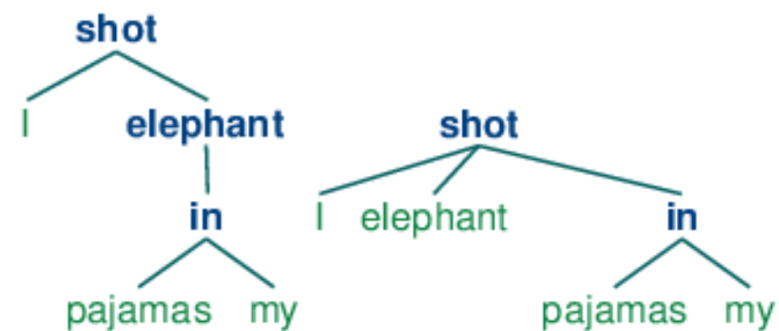
words → phrases → clauses → sentences

Formalisms

Phrase structure grammar
(Chomsky 1957)



Dependency grammar
(Mel'čuk 1988; Tesnière 1959; Pāṇini)



Constituency

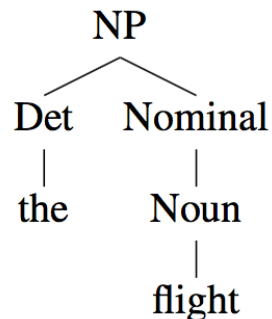
- Groups of words (“**constituents**”) behave as single units
 - Noun Phrases: groups of tokens that act Like nouns
 - ▶ Harry the Horse
 - ▶ the Broadway coppers
 - ▶ they
 - ▶ a high-class spot such as Mindy’s
 - ▶ the reason he comes into the Hot Box
 - ▶ three parties from Brooklyn

Constituency

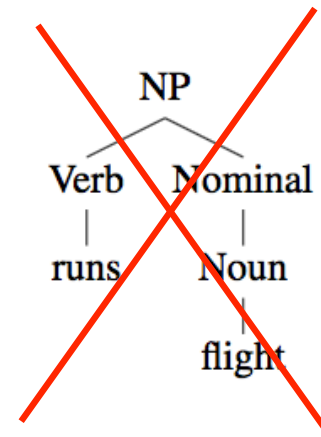
- Groups of words (“**constituents**”) behave as single units
 - Noun Phrases: groups of tokens that act like nouns
- Linguists characterize constituents in a number of ways, including:
 - ▶ where they occur (e.g., “NPs can occur before verbs”)
 - ▶ where they can *move* in variations of a sentence
 - ▶ On September 17th, I’d like to fly from Atlanta to Denver
 - ▶ I’d like to fly on September 17th from Atlanta to Denver
 - ▶ I’d like to fly from Atlanta to Denver on September 17th
 - ▶ what parts can move and what parts can’t
 - ▶ *On September I’d like to fly 17th from Atlanta to Denver
 - ▶ what they can be conjoined with
 - ▶ I’d like to fly from Atlanta to Denver on September 17th and in the morning

Context-Free Grammar (CFG)

- Take constituents to be the main building block of natural language syntax, we can attempt to formalize what makes a string grammatical in a language.
- A CFG gives a formal way to define what meaningful constituents are and exactly how a constituent is formed out of other constituents (or words). It defines **valid structure** in a language.



NP → Det Nominal



NP → Verb Nominal

Context-Free Grammar (CFG)

A **context-free grammar** consists of:

- ▶ A finite set of nonterminal symbols \mathcal{N} (sometimes called “categories”)
 - ▶ A start symbol $S \in \mathcal{N}$
- ▶ A finite alphabet Σ , called “terminal” symbols, distinct from \mathcal{N}
- ▶ Production rule set \mathcal{R} , each of the form “ $N \rightarrow \alpha$ ” where
 - ▶ The lefthand side N is a nonterminal from \mathcal{N}
 - ▶ The righthand side α is a sequence of zero or more terminals and/or nonterminals: $\alpha \in (\mathcal{N} \cup \Sigma)^*$
 - ▶ Special case: **Chomsky normal form** constrains α to be either a single terminal symbol or two nonterminals

An Example CFG for a Tiny Bit of English

From Jurafsky and Martin (forthcoming)

S → NP VP

S → Aux NP VP

S → VP

NP → Pronoun

NP → Proper-Noun

NP → Det Nominal

Nominal → Noun

Nominal → Nominal Noun

Nominal → Nominal PP

VP → Verb

VP → Verb NP

VP → Verb NP PP

VP → Verb PP

VP → VP PP

PP → Preposition NP

Det → that | this | a

Noun → book | flight | meal | money

Verb → book | include | prefer

Pronoun → I | she | me

Proper-Noun → Houston | NWA

Aux → does

Preposition → from | to | on | near

| through

“Lexicon”

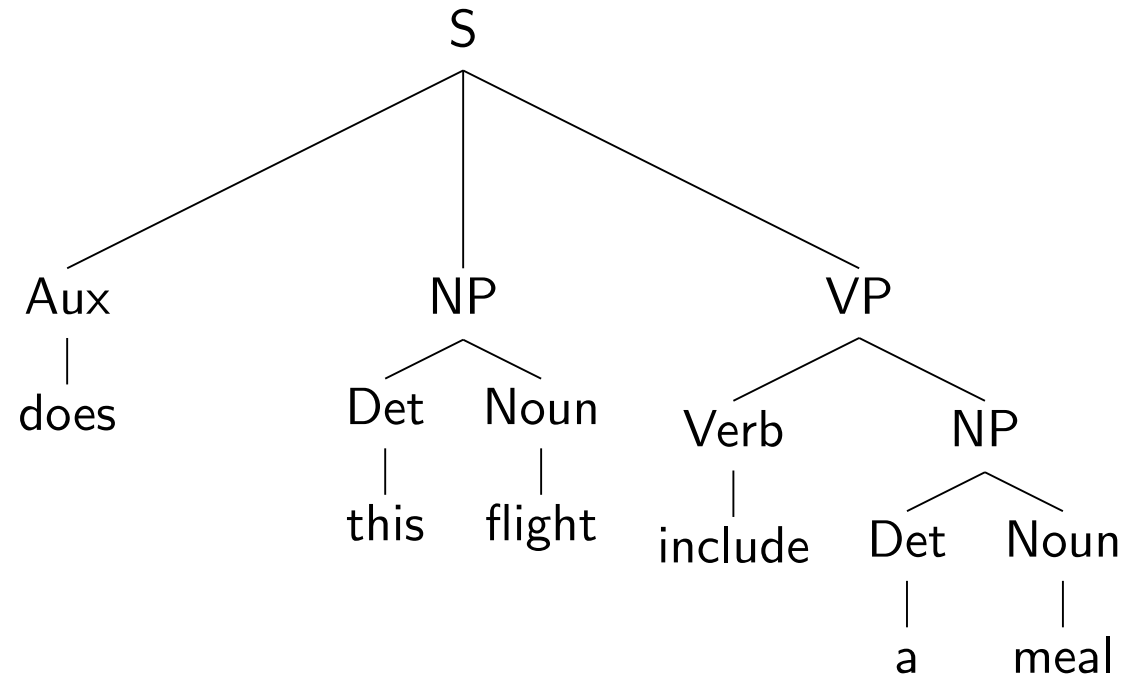
This term is used in NLP to refer to an object that associates information with words.

In a CFG, the “lexicon rules” are the rules that map a nonterminal (usually a part of speech) to a single word.

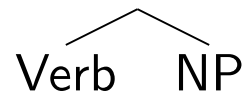
Derivation

- Given a CFG, a derivation is the sequence of productions used to generate a string of words (e.g., a sentence), often visualized as a **parse tree**.
- Language: the formal language defined by a CFG is the set of strings derivable from S (start symbol)

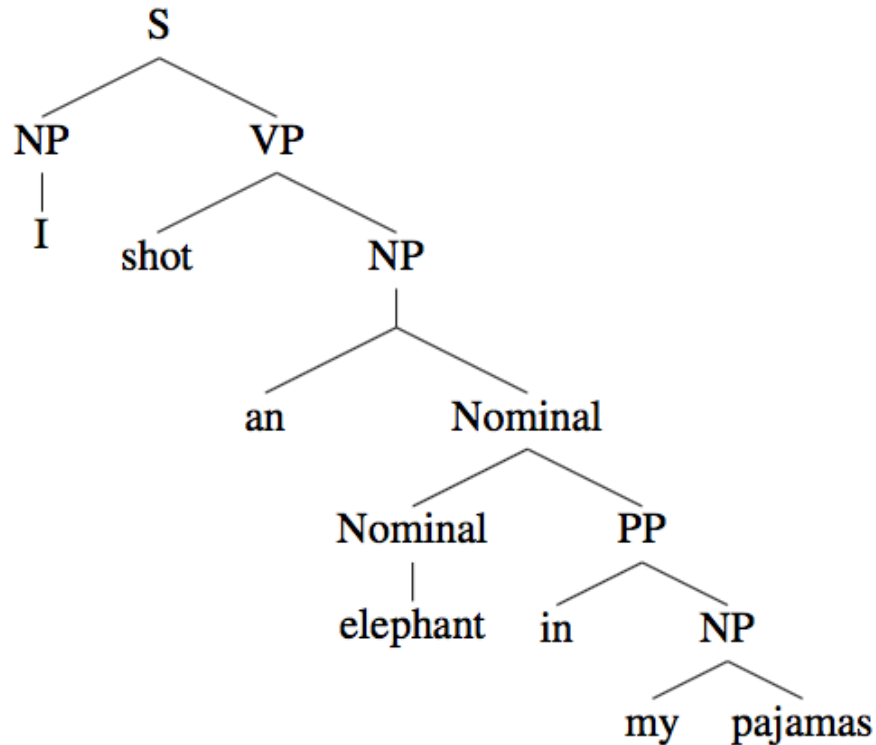
Example Derivation



The phrase-structure tree represents both the syntactic structure of the sentence and the **derivation** of the sentence under the grammar. E.g., $VP \rightarrow Verb\ NP$ corresponds to the rule $VP \rightarrow Verb\ NP$.



Example Derivation



Every internal node is a phrase

- my pajamas
- in my pajamas
- elephant in my pajamas
- an elephant in my pajamas
- shot an elephant in my pajamas
- I shot an elephant in my pajamas

Each phrase could be replaced by another of the same type of constituent

Where do natural language CFGs come from?

Building a CFG for a natural language by hand is really hard (Jurafsky and Martin, forthcoming, chapter 10).

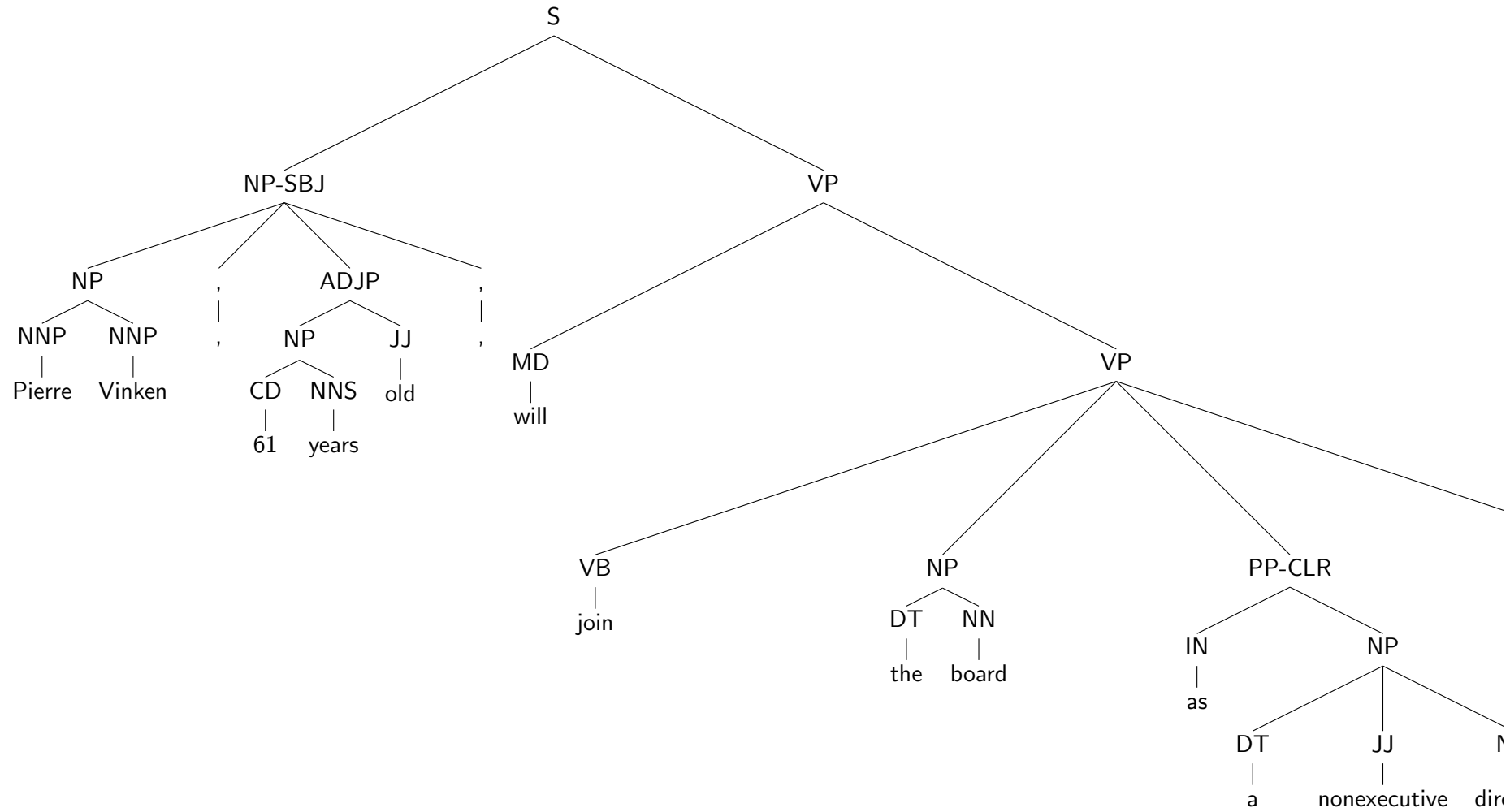
- ▶ Need lots of categories to make sure all and only grammatical sentences are included.
- ▶ Categories tend to start exploding combinatorially.
- ▶ Alternative grammar formalisms are typically used for manual grammar construction; these are often based on constraints and a powerful algorithmic tool called *unification*.

A data-driven approach:

1. Build a corpus of annotated sentences, called a **treebank**. (e.g., the Penn Treebank, Marcus et al., 1993.)
2. Extract rules from the treebank.
3. Optionally, use statistical models to generalize the rules.

collections of sentences annotated with syntactic structure

Example from the Penn Treebank



Some Penn Treebank Rules with Counts

40717 PP → IN NP	100 VP → VBD PP-PRD
33803 S → NP-SBJ VP	100 PRN → : NP :
22513 NP-SBJ → -NONE-	100 NP → DT JJS
21877 NP → NP PP	100 NP-CLR → NN
20740 NP → DT NN	99 NP-SBJ-1 → DT NNP
14153 S → NP-SBJ VP .	98 VP → VBN NP PP-DIR
12922 VP → TO VP	98 VP → VBD PP-TMP
11881 PP-LOC → IN NP	98 PP-TMP → VBG NP
11467 NP-SBJ → PRP	97 VP → VBD ADVP-TMP VP
11378 NP → -NONE-	...
11291 NP → NN	10 WHNP-1 → WRB JJ
...	10 VP → VP CC VP PP-TMP
989 VP → VBG S	10 VP → VP CC VP
985 NP-SBJ → NN	ADVP-MNR
983 PP-MNR → IN NP	10 VP → VBZ S , SBAR-ADV
983 NP-SBJ → DT	10 VP → VBZ S ADVP-TMP
969 VP → VBN VP	

Penn Treebank Rules: Statistics

32,728 rules in the training section (not including 52,257 lexicon rules)

4,021 rules in the development section

overlap: 3,128

(Phrase-Structure) Recognition and Parsing

Given a CFG $(\mathcal{N}, S, \Sigma, \mathcal{R})$ and a sentence x , the **recognition** problem is:

Is x in the language of the CFG?

The proof is a derivation of x using the rules \mathcal{R} .

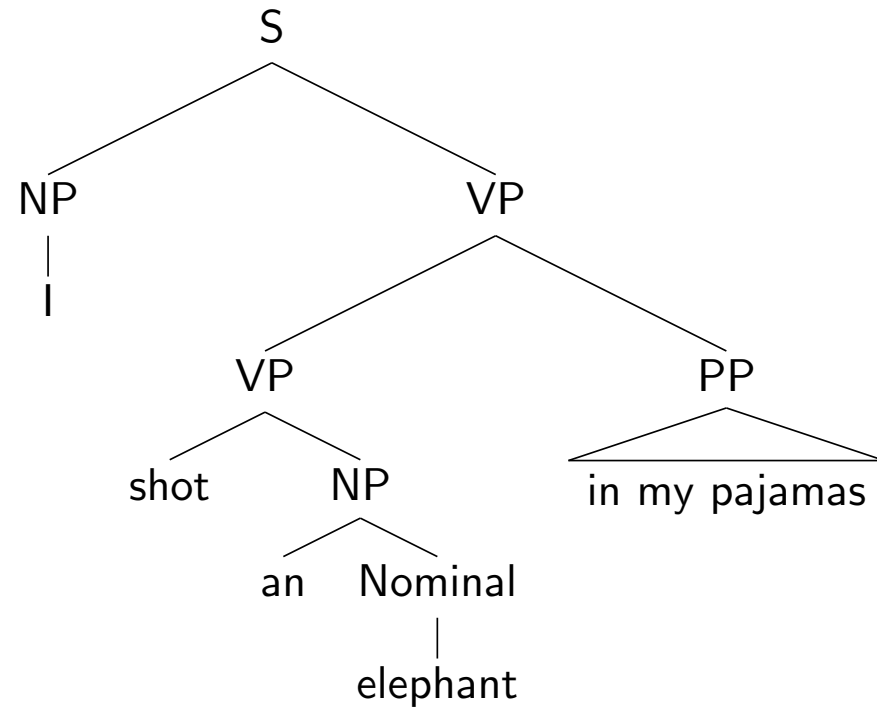
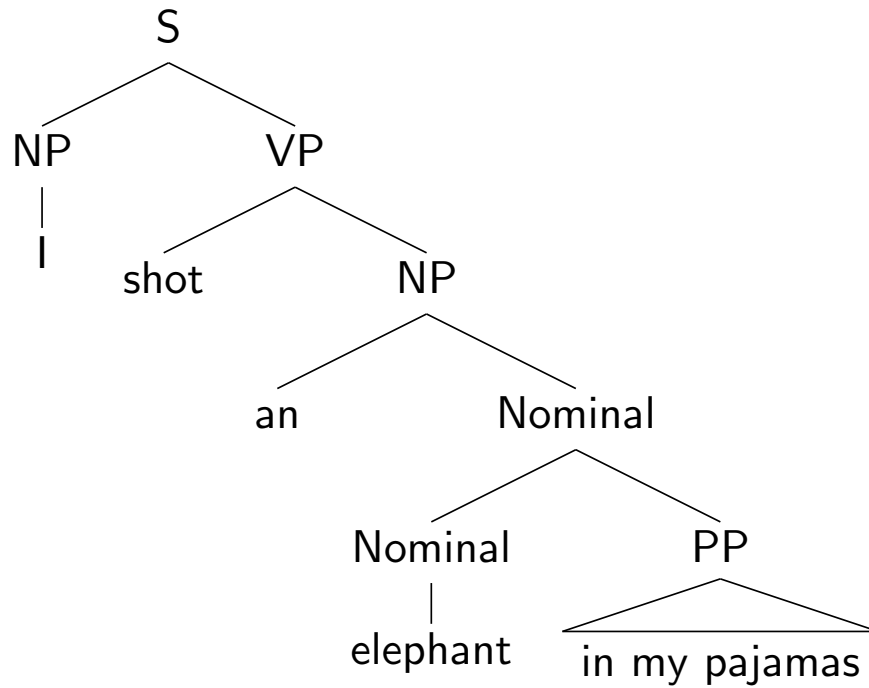
Related problem: **parsing**:

Show one or more derivations for x , using \mathcal{R} .

With reasonable grammars, the number of parses is exponential in $|x|$.

Syntactic Ambiguity

Ambiguity is the most serious problem faced by syntactic parsers



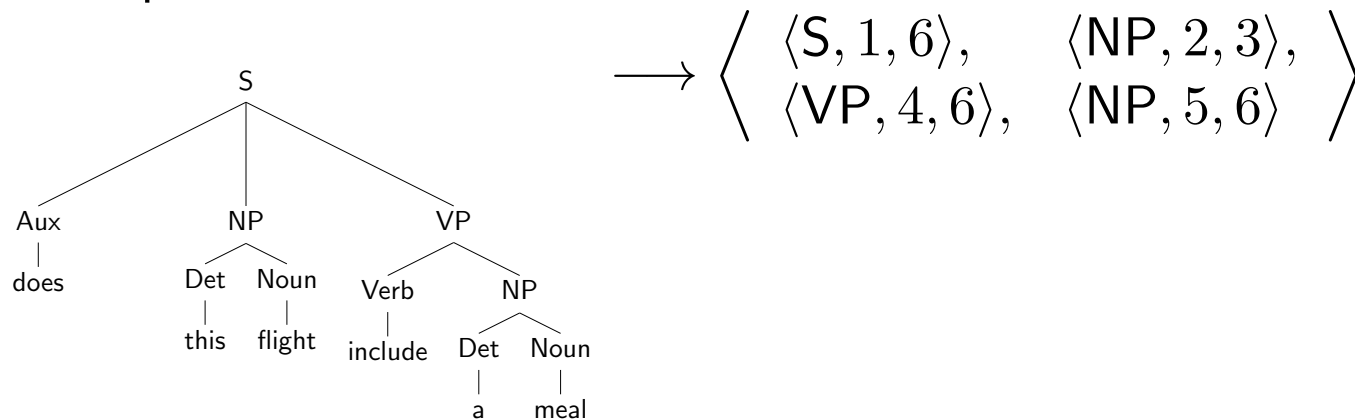
Parser Evaluation

Represent a parse tree as a collection of tuples

$\langle \langle \ell_1, i_1, j_1 \rangle, \langle \ell_2, i_2, j_2 \rangle, \dots, \langle \ell_n, i_n, j_n \rangle \rangle$, where

- ▶ ℓ_k is the nonterminal labeling the k th phrase
- ▶ i_k is the index of the first word in the k th phrase
- ▶ j_k is the index of the last word in the k th phrase

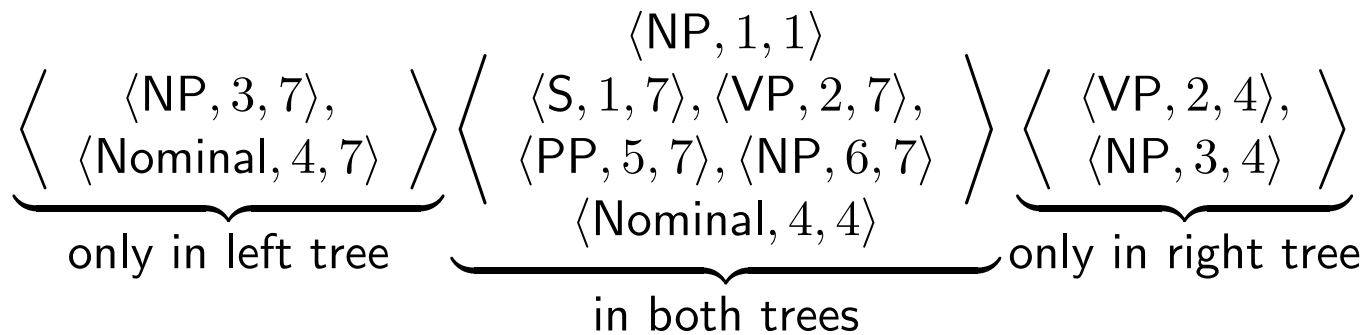
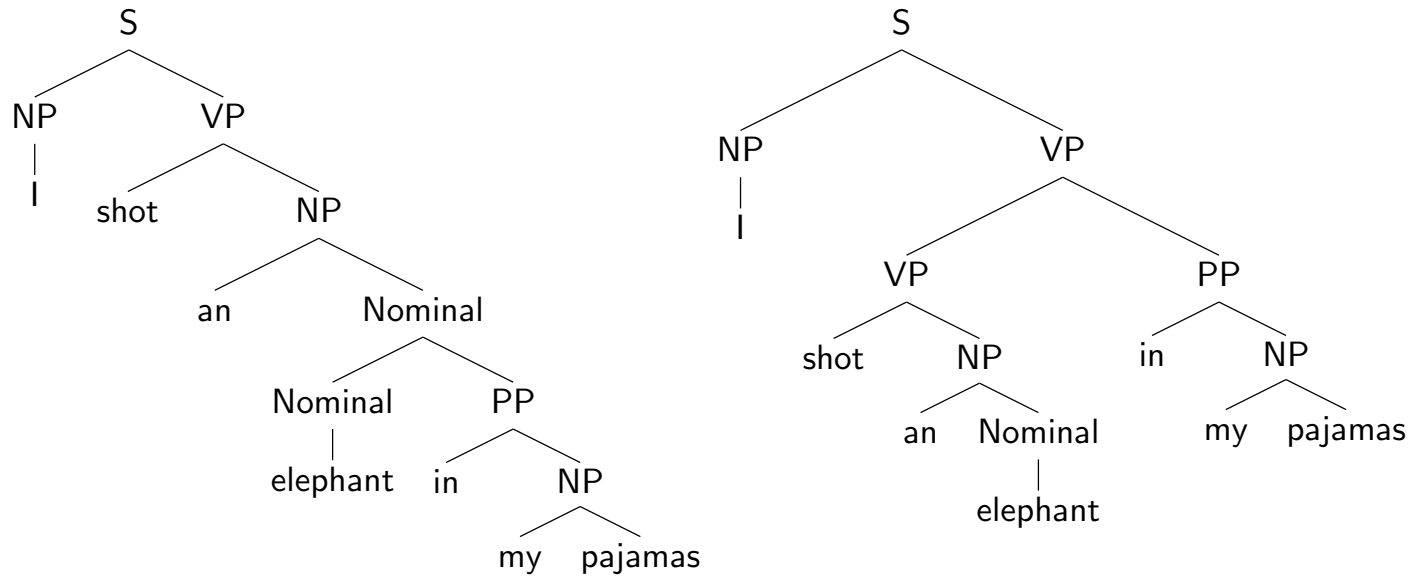
Example:



Convert gold-standard tree and system hypothesized tree into this representation, then estimate precision, recall, and F_1 .

Parser Evaluation

Tree Comparison Example



Two Views of Parsing

1. Incremental search: the state of the search is the partial structure built so far; each action incrementally extends the tree.
 - ▶ Often **greedy**, with a statistical classifier deciding what action to take in every state.
2. Discrete optimization: define a scoring function and seek the tree with the highest score.

Probabilistic Context-Free Grammar (PCFG)

- A basic CFG allows us to check whether a sentence is grammatical in the language it defines
- Binary decision: a sentence is either in the language (a series of productions yields the words we see) or it is not.
- PCFG: each production is also associated with a probability.
- This lets us calculate the probability of a parse for a given sentence
 - For a given parse tree T for sentence S comprised of n rules (each $A \rightarrow \beta$):

$$P(T, S) = \prod_i^n P(\beta \mid A)$$

Probabilistic Context-Free Grammar (PCFG)

A **probabilistic context-free grammar** consists of:

- ▶ A finite set of nonterminal symbols \mathcal{N}
 - ▶ A start symbol $S \in \mathcal{N}$
 - ▶ A finite alphabet Σ , called “terminal” symbols, distinct from \mathcal{N}
 - ▶ Production rule set \mathcal{R} , each of the form “ $N \rightarrow \alpha$ ” where
 - ▶ The lefthand side N is a nonterminal from \mathcal{N}
 - ▶ The righthand side α is a sequence of zero or more terminals and/or nonterminals: $\alpha \in (\mathcal{N} \cup \Sigma)^*$
 - ▶ Special case: **Chomsky normal form** constrains α to be either a single terminal symbol or two nonterminals
- ▶ For each $N \in \mathcal{N}$, a probability distribution over the rules where N is the lefthand side, $p(* \mid N)$.

PCFG Scores Trees

We can write the parsing problem as finding the best-scoring tree:

$$\hat{t} = \operatorname{argmax}_{t \in \mathcal{T}_x} \operatorname{Score}(t)$$

PCFGs view each tree t as a “bag of rules” (from \mathcal{R}), and define:

$$\begin{aligned} \operatorname{Score}(t) &= p(t) \\ &= \prod_{(N \rightarrow \alpha) \in \mathcal{R}} p(\alpha \mid N)^{\operatorname{count}(N \rightarrow \alpha; t)} \end{aligned}$$

PCFG Example

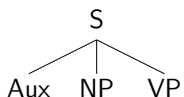
S

Write down the start symbol. Here: S

Probability:

1

PCFG Example

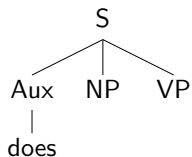


Choose a rule from the “S” distribution. Here: $S \rightarrow \text{Aux NP VP}$

Probability:

$$p(\text{Aux NP VP} \mid S)$$

PCFG Example

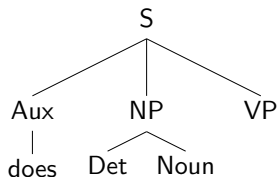


Choose a rule from the “Aux” distribution. Here: $\text{Aux} \rightarrow \text{does}$

Probability:

$$p(\text{Aux NP VP} \mid S) \cdot p(\text{does} \mid \text{Aux})$$

PCFG Example

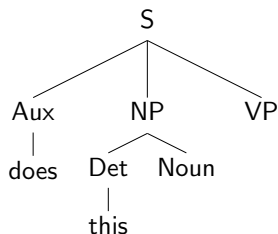


Choose a rule from the “NP” distribution. Here: NP \rightarrow Det Noun

Probability:

$$p(\text{Aux NP VP} \mid \text{S}) \cdot p(\text{does} \mid \text{Aux}) \cdot p(\text{Det Noun} \mid \text{NP})$$

PCFG Example

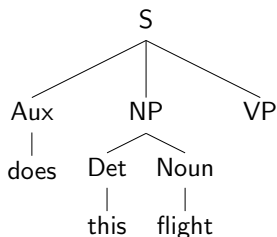


Choose a rule from the “Det” distribution. Here: Det \rightarrow this

Probability:

$$p(\text{Aux NP VP} \mid \text{S}) \cdot p(\text{does} \mid \text{Aux}) \cdot p(\text{Det Noun} \mid \text{NP}) \cdot p(\text{this} \mid \text{Det})$$

PCFG Example

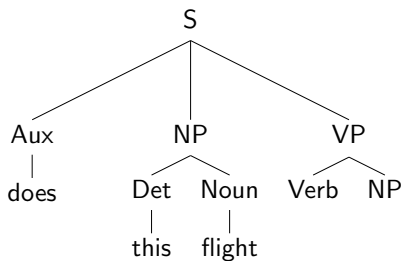


Choose a rule from the “Noun” distribution. Here: Noun \rightarrow flight

Probability:

$$p(\text{Aux NP VP} \mid \text{S}) \cdot p(\text{does} \mid \text{Aux}) \cdot p(\text{Det Noun} \mid \text{NP}) \cdot p(\text{this} \mid \text{Det}) \\ \cdot p(\text{flight} \mid \text{Noun})$$

PCFG Example

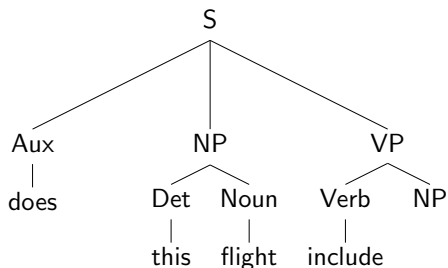


Choose a rule from the “VP” distribution. Here: $VP \rightarrow \text{Verb NP}$

Probability:

$$p(\text{Aux NP VP} \mid S) \cdot p(\text{does} \mid \text{Aux}) \cdot p(\text{Det Noun} \mid \text{NP}) \cdot p(\text{this} \mid \text{Det}) \\ \cdot p(\text{flight} \mid \text{Noun}) \cdot p(\text{Verb NP} \mid \text{VP})$$

PCFG Example

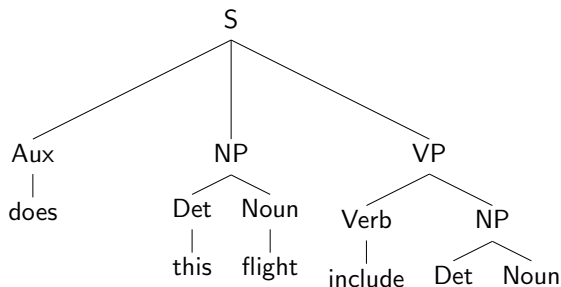


Choose a rule from the “Verb” distribution. Here: Verb \rightarrow include

Probability:

$$p(\text{Aux NP VP} \mid \text{S}) \cdot p(\text{does} \mid \text{Aux}) \cdot p(\text{Det Noun} \mid \text{NP}) \cdot p(\text{this} \mid \text{Det}) \\ \cdot p(\text{flight} \mid \text{Noun}) \cdot p(\text{Verb NP} \mid \text{VP}) \cdot p(\text{include} \mid \text{Verb})$$

PCFG Example

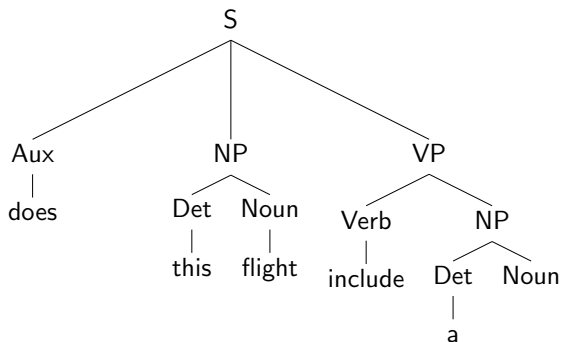


Choose a rule from the “NP” distribution. Here: NP \rightarrow Det Noun

Probability:

$$\begin{aligned} & p(\text{Aux NP VP} \mid \text{S}) \cdot p(\text{does} \mid \text{Aux}) \cdot p(\text{Det Noun} \mid \text{NP}) \cdot p(\text{this} \mid \text{Det}) \\ & \cdot p(\text{flight} \mid \text{Noun}) \cdot p(\text{Verb NP} \mid \text{VP}) \cdot p(\text{include} \mid \text{Verb}) \\ & \cdot p(\text{Det Noun} \mid \text{NP}) \end{aligned}$$

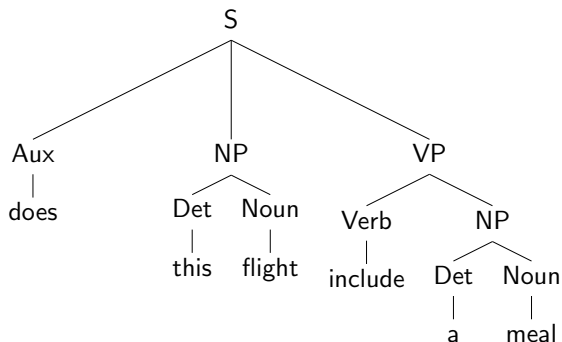
PCFG Example



Choose a rule from the “Det” distribution. Here: Det \rightarrow a
Probability:

$$p(\text{Aux NP VP} \mid \text{S}) \cdot p(\text{does} \mid \text{Aux}) \cdot p(\text{Det Noun} \mid \text{NP}) \cdot p(\text{this} \mid \text{Det}) \\ \cdot p(\text{flight} \mid \text{Noun}) \cdot p(\text{Verb NP} \mid \text{VP}) \cdot p(\text{include} \mid \text{Verb}) \\ \cdot p(\text{Det Noun} \mid \text{NP}) \cdot p(\text{a} \mid \text{Det})$$

PCFG Example



Choose a rule from the “Noun” distribution. Here: Noun \rightarrow meal
Probability:

$$p(\text{Aux NP VP} \mid \text{S}) \cdot p(\text{does} \mid \text{Aux}) \cdot p(\text{Det Noun} \mid \text{NP}) \cdot p(\text{this} \mid \text{Det}) \\ \cdot p(\text{flight} \mid \text{Noun}) \cdot p(\text{Verb NP} \mid \text{VP}) \cdot p(\text{include} \mid \text{Verb}) \\ \cdot p(\text{Det Noun} \mid \text{NP}) \cdot p(\text{a} \mid \text{Det}) \cdot p(\text{meal} \mid \text{Noun})$$

Parsing with PCFGs

- How to set the probabilities $p(\text{righthand side} \mid \text{lefthand side})$?
 - Counting / Learning (we won't discuss in this class due to time limit)

- How to decode/parse?

Probabilistic CKY (Cocke-Kasami-Younger)

(Cocke and Schwartz, 1970; Kasami, 1965; Younger, 1967)

Input:

- ▶ a PCFG $(\mathcal{N}, S, \Sigma, \mathcal{R}, p(* | *))$, in **Chomsky normal form**
- ▶ a sentence x (let n be its length)

Output: If x is in the language of the grammar.

$$\operatorname{argmax}_{t \in \mathcal{T}_x} \log p(t);$$

undefined if not.

Probabilistic CKY

Probabilistic CKY is closely related to the Viterbi algorithm; it is a dynamic programming algorithm.

The recurrence is defined around

$$\heartsuit_{i:j}(N),$$

which will store the best score (log probability) found (so far) for constructing an N -rooted constituent that spans $\langle x_i, \dots, x_j \rangle$.

In Viterbi, we used conditional independence to collapse all prefix label sequences that ended in the same label into one stored item; here we collapse all trees spanning words i to j with the same root into a single item.

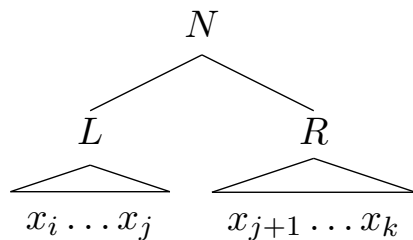
Probabilistic CKY

Base case: for $i \in \{1, \dots, n\}$ and for each $N \in \mathcal{N}$:

$$\heartsuit_{i:i}(N) = \log p(x_i \mid N)$$

For each i, k such that $1 \leq i < k \leq n$ and each $N \in \mathcal{N}$:

$$\heartsuit_{i:k}(N) = \max_{L, R \in \mathcal{N}, j \in \{i, \dots, k-1\}} \log p(L R \mid N) + \heartsuit_{i:j}(L) + \heartsuit_{(j+1):k}(R)$$

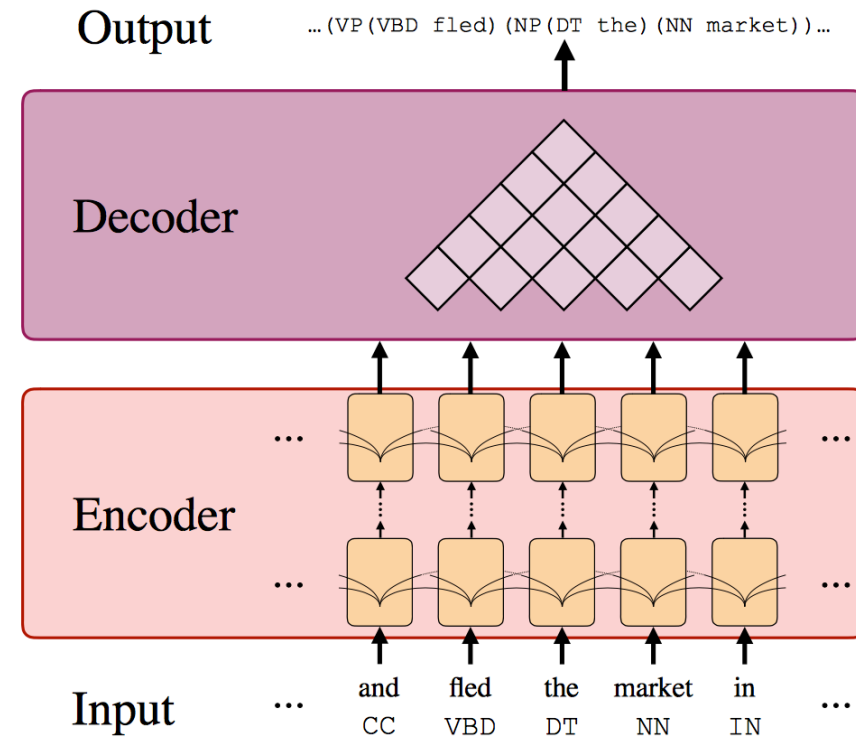


Solution:

$$\heartsuit_{1:n}(S) = \max_{t \in \mathcal{T}_x} \log p(t)$$

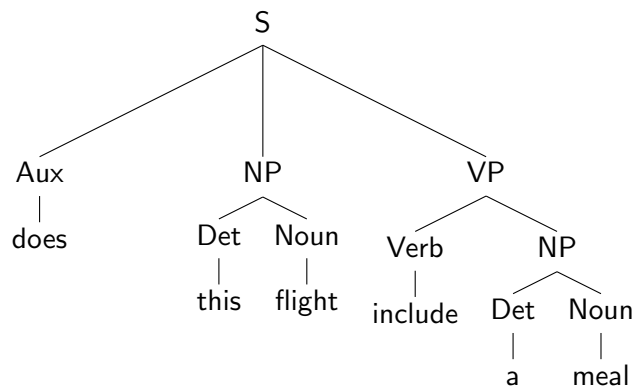
Neural Parsing

- Kitaev and Klein (2018), “Constituency Parsing with a Self-Attentive Encoder”
- Neural model (attention encoder) generates representations of each token in a sentence
- Learned scoring $s(i,j,k)$ function for each span from token i to token j with label k
- CKY for **decoding** to find the best tree through this space.



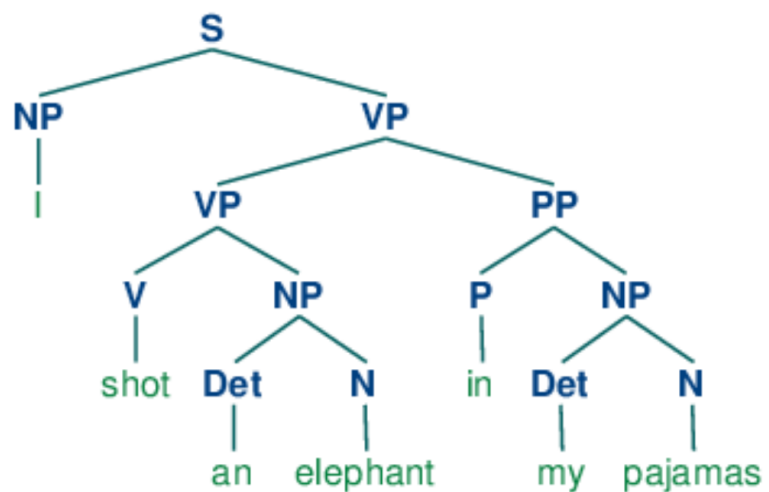
Summary so far

- Constituents: groups of words behave as single units
- Context-Free Grammar (CFG)
 - A CFG gives a formal way to define a valid structure in a language
- Probabilistic Context-Free Grammar (PCFG)
 - Each production is also associated with a probability
- Parsing:
 - Show one or more derivations for a sentence, using the grammar
 - (Probabilistic) CKY

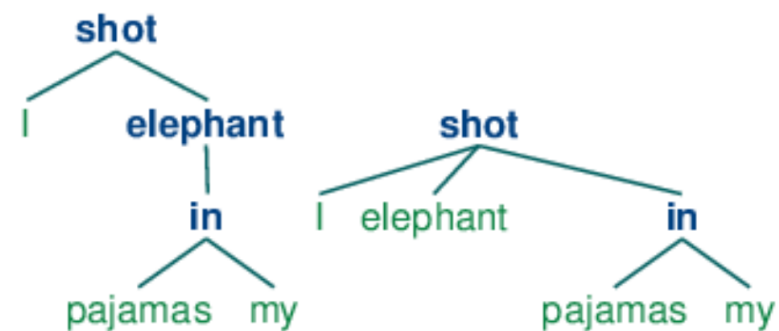


Formalisms

Phrase structure grammar
(Chomsky 1957)

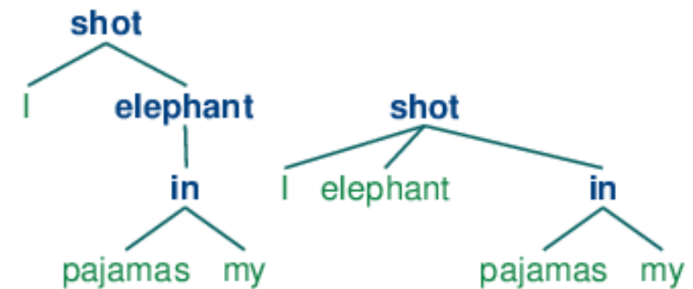
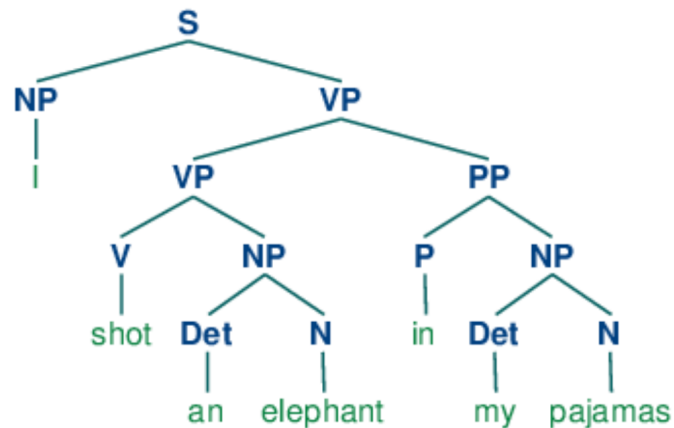


Dependency grammar
(Mel'čuk 1988; Tesnière 1959; Pāṇini)



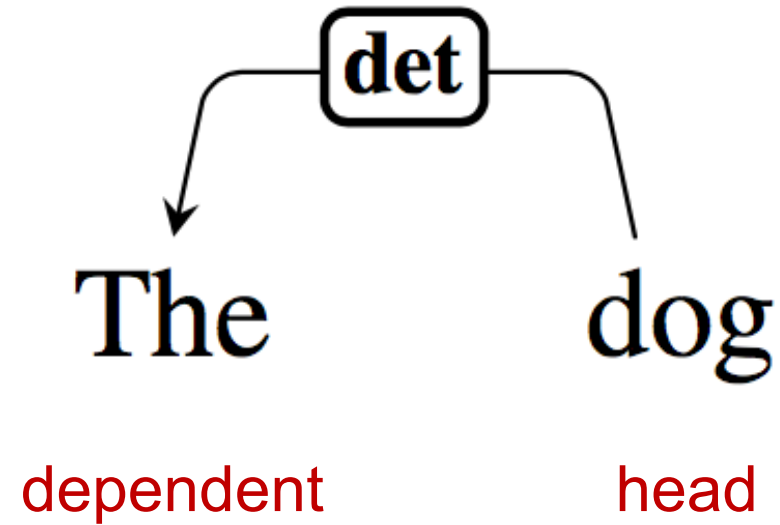
Dependency syntax

A different family of theories of syntax focuses on dependencies between words

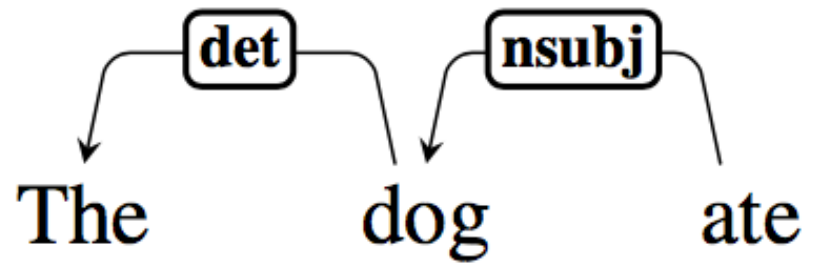


- Dependency syntax doesn't have non-terminal structure like a CFG; words are directly linked to each other.

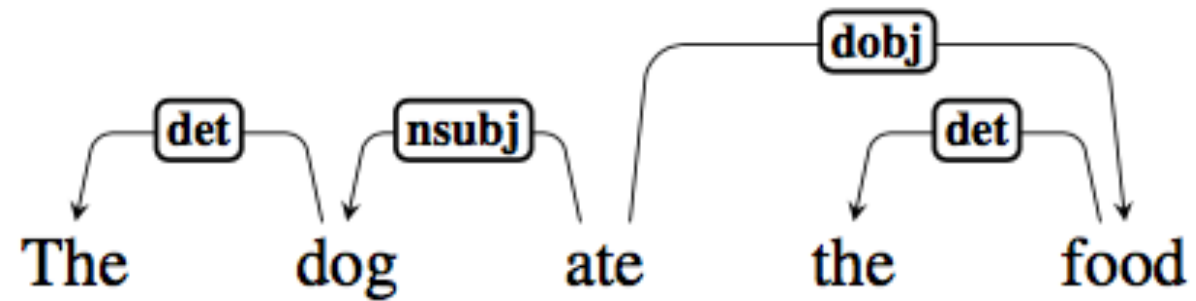
Dependency syntax



Dependency syntax

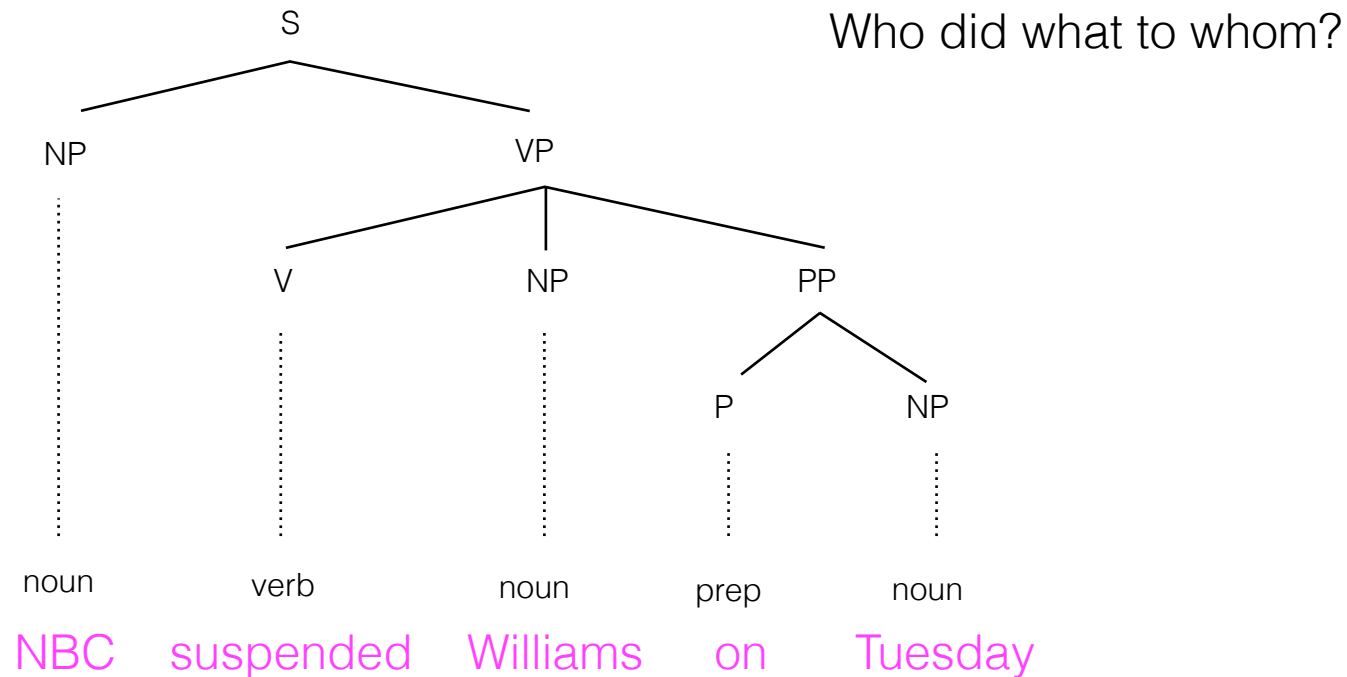


Dependency syntax



Dependencies vs constituents

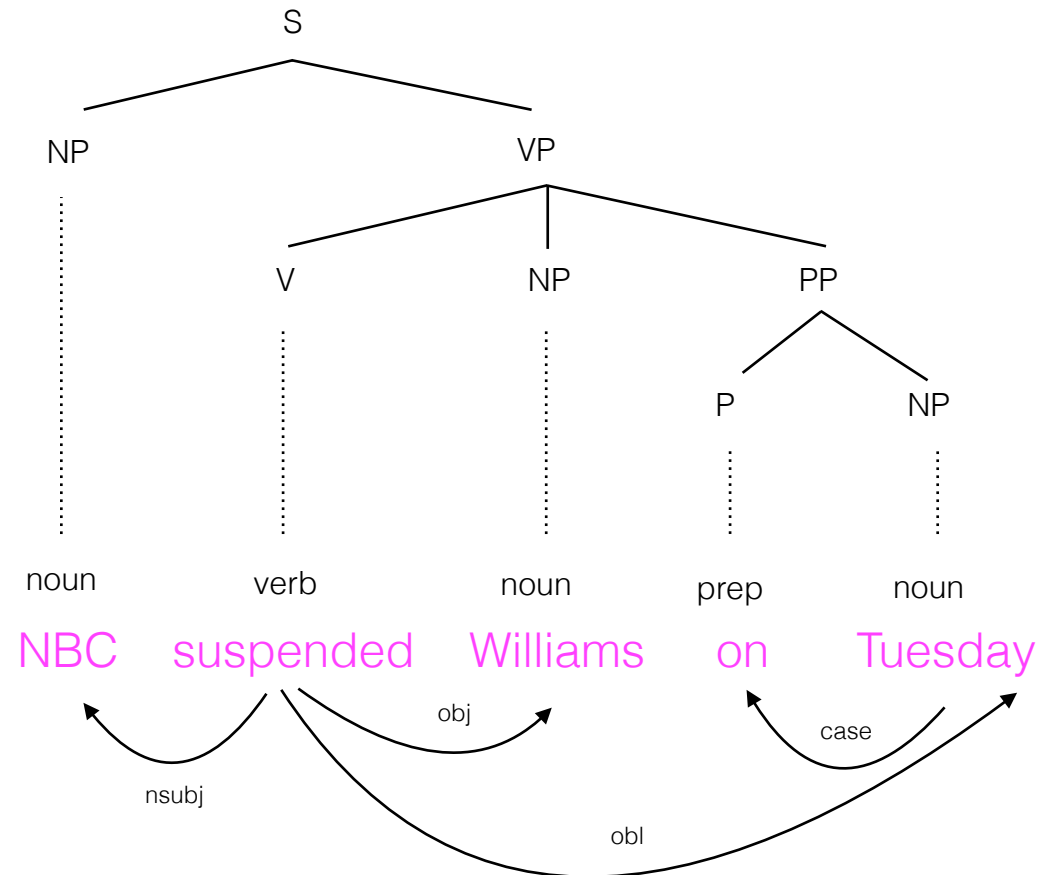
- Dependency links are closer to semantic relationships; no need to infer the relationships from the structure of a tree



subject: S → NP VP
direct object: S → NP (VP → ... NP ...)

Dependencies vs constituents

- Dependency links are closer to semantic relationships; no need to infer the relationships from the structure of a tree

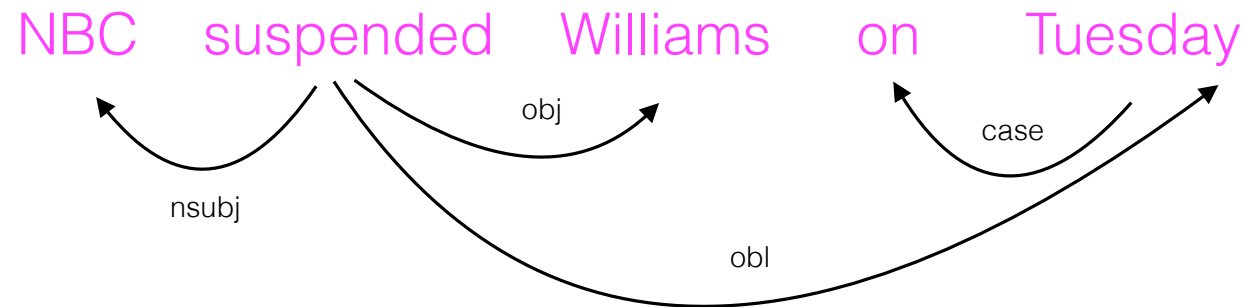


Dependencies vs constituents

- Dependency links are closer to semantic relationships; no need to infer the relationships from the structure of a tree

Captures binary relations between words

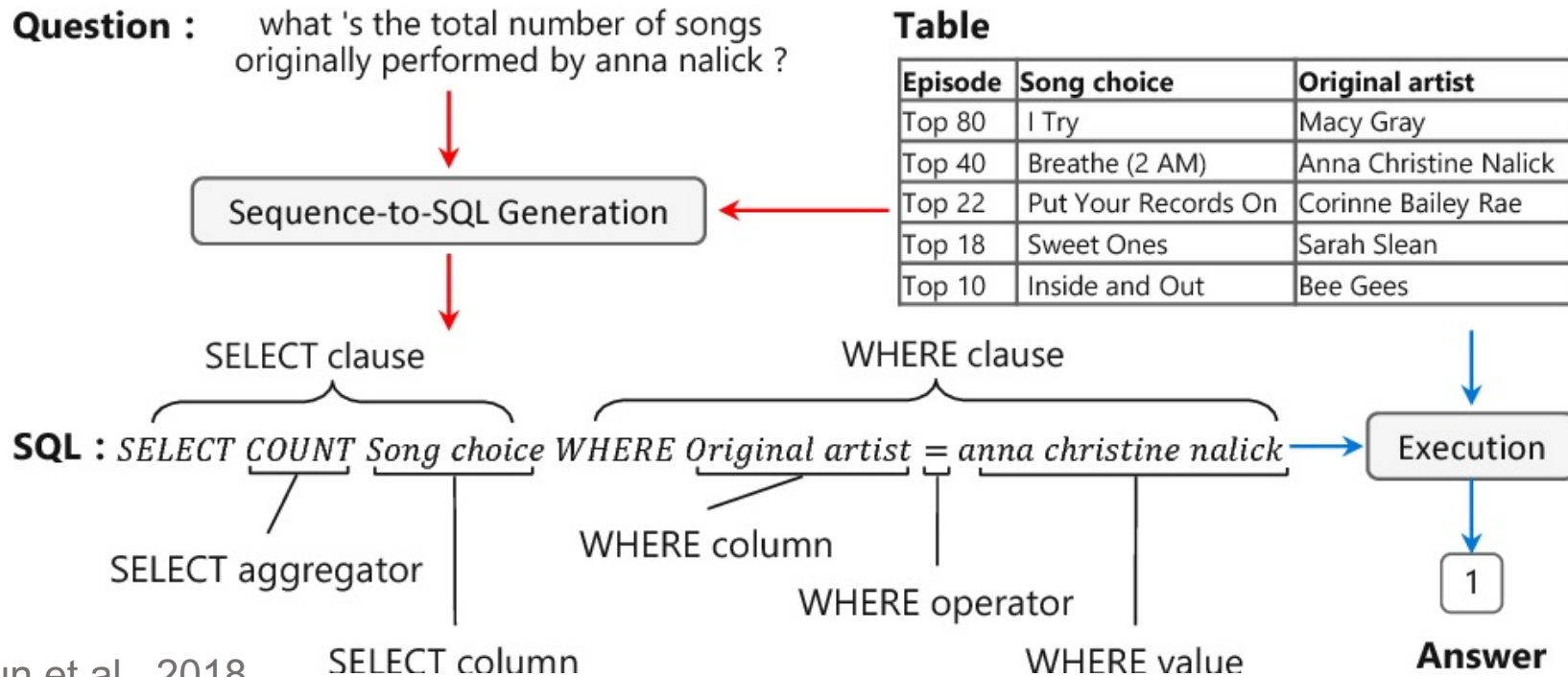
- nsubj(NBC, suspended)
- obj(Williams, suspended)



Semantic Parsing

Semantic parsing comprises a wide range of tasks where strings are mapped into meaning representation languages. Examples:

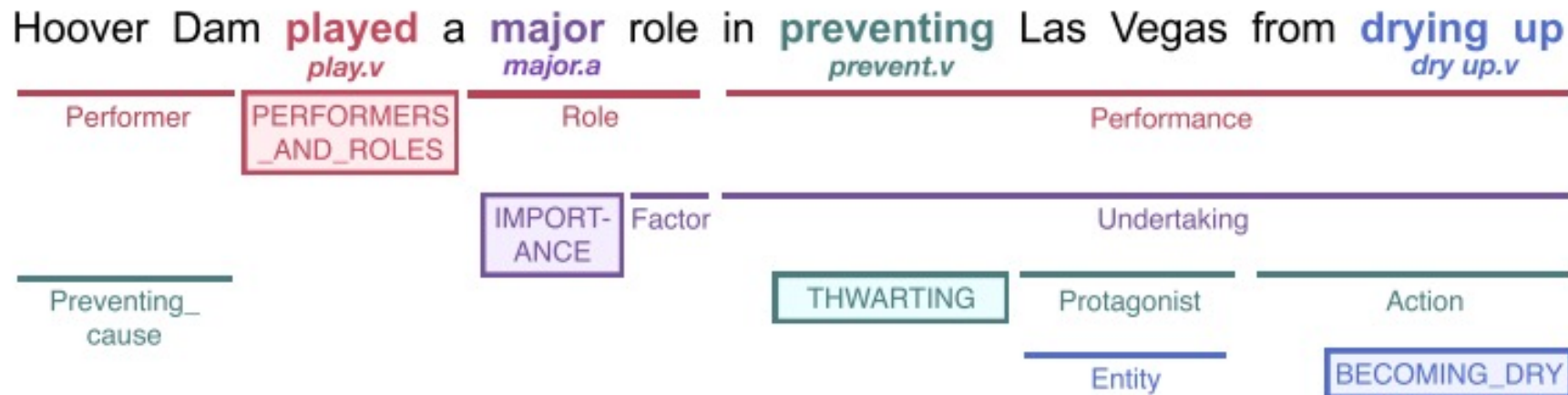
- ▶ Programming languages, especially query languages that can be used to answer questions using a database (Zettlemoyer and Collins, 2005, e.g.,)



Semantic Parsing

Semantic parsing comprises a wide range of tasks where strings are mapped into meaning representation languages. Examples:

- ▶ Programming languages, especially query languages that can be used to answer questions using a database (Zettlemoyer and Collins, 2005, e.g.,)
- ▶ Schemas designed around real-world event-types (called “frames”); trying to extract “who did what to whom?” (Baker et al., 1998; Palmer et al., 2005)



Other Examples of Linguistic Structure Prediction

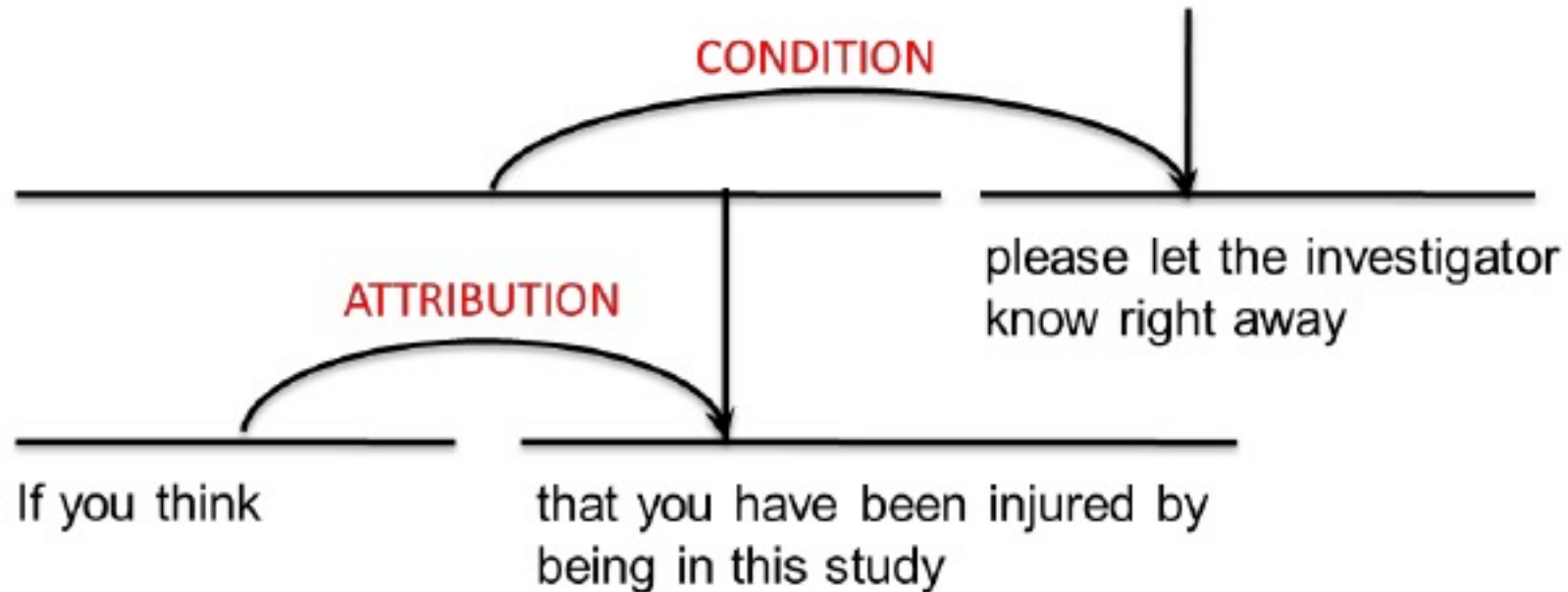
- Coreference resolution

*"I voted for Nader because he was most
aligned with my values," she said.*

The diagram shows three curved arrows indicating coreference relationships between pronouns in the text. One arrow points from 'I' to 'he', another from 'he' to 'she', and a third from 'she' to 'I', forming a cycle that identifies the speaker as the subject of the sentence.

Other Examples of Linguistic Structure Prediction

- Coreference resolution
- Discourse parsing



Key Takeaways

- Syntax
 - Constituency parsing
 - CFG, PCFG
 - Dependency parsing
- Semantic Parsing
- Coreference Resolution
- Discourse Parsing

Questions?