

Natural Language Processing (CSE 447/547M): Phrase Structure Parsing

Noah Smith

© 2019

University of Washington
`nasmith@cs.washington.edu`

February 20, 2019

Administrivia

- ▶ Please double-check your tarballs before you turn them in!
- ▶ Assignment 2: watch for updated version!
- ▶ Highlights from mid-quarter evaluation:
 - ▶ You are mostly happy with lectures, TAs, office hours, the switch to Piazza and TA responsiveness there, slides, question management in class, section slides, textbook.
 - ▶ You want more formal instruction on AllenNLP and more explanation of the connection between lecture and assignments.

Context-Free Grammar

A **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

(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?

Related problem: **parsing**:

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

(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.

Related problem: **parsing**:

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

(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?

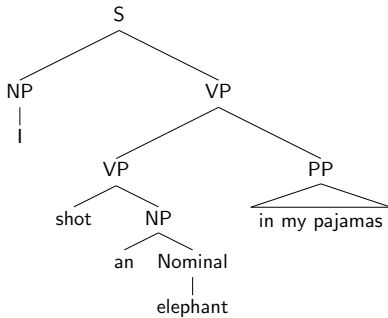
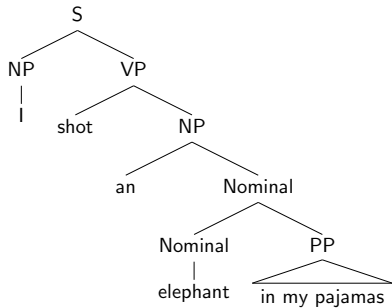
The proof is a derivation.

Related problem: **parsing**:

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

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

Ambiguity

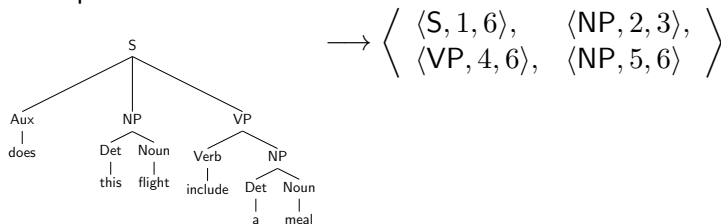


Parser Evaluation

Represent a parse tree as a collection of tuples $\langle \ell_1, i_1, j_1 \rangle, \langle \ell_2, i_2, j_2 \rangle, \dots, \langle \ell_n, i_n, j_n \rangle$, where

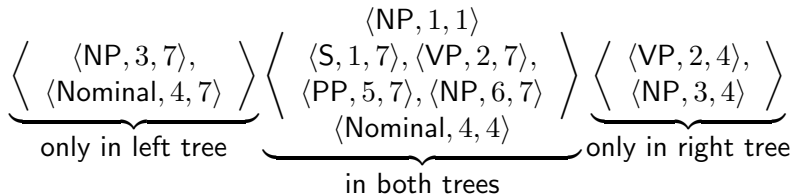
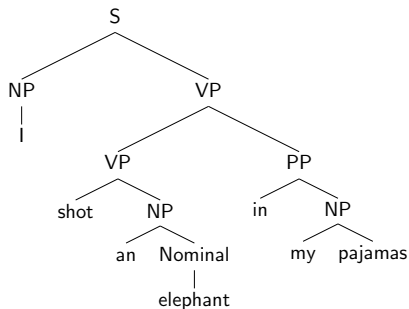
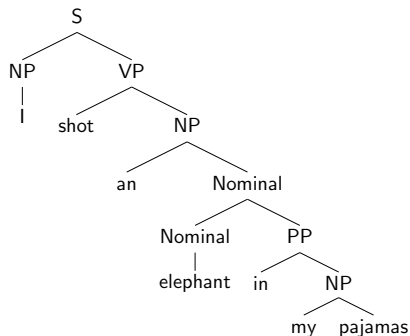
- ▶ ℓ_k is the nonterminal labeling the k th phrase
- ▶ i_k is the index of the first word in the k th phrase
- ▶ j_k is the index of the last word in the k th phrase

Example:



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

Tree Comparison Example



Two Views of Parsing

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.

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.

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.

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.
 - Today: scores are defined using the rules.

$$\text{predict}(\mathbf{x}) = \underset{\mathbf{t}=\langle r_1, \dots, r_k \rangle}{\operatorname{argmax}} \sum_{i=1}^k s(r_i) = \underset{\mathbf{t}}{\operatorname{argmax}} \sum_{r \in \mathcal{R}} c_{\mathbf{t}}(r) s(r) \quad (1)$$

where \mathbf{t} is constrained to be a grammatical tree with \mathbf{x} as the yield, and $\langle r_1, \dots, r_k \rangle$ denotes the bag of rules used in deriving the tree. Denote this set of grammatical trees with \mathbf{x} as the yield $\mathcal{T}_{\mathbf{x}}$.

Probabilistic Context-Free Grammar

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, $p(* \mid N)$.

Note: in the notation of equation 1, $s(N \rightarrow \alpha) = \log p(\alpha \mid N)$, and the total score of a tree is its *log* probability.

PCFG Example

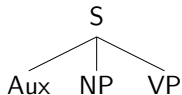
S

Write down the start symbol. Here: S

Probability:

1

PCFG Example

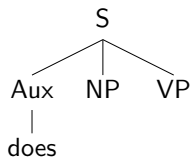


Choose a rule from the “S” distribution. Here: $S \rightarrow \text{Aux NP VP}$

Probability:

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

PCFG Example

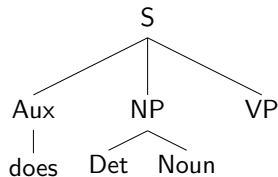


Choose a rule from the “Aux” distribution. Here: $\text{Aux} \rightarrow \text{does}$

Probability:

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

PCFG Example

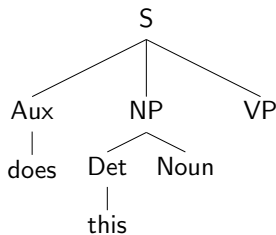


Choose a rule from the “NP” distribution. Here: $\text{NP} \rightarrow \text{Det Noun}$

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

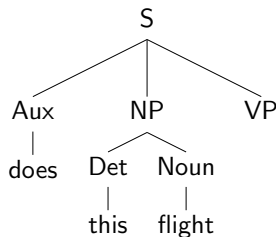


Choose a rule from the “Det” distribution. Here: $\text{Det} \rightarrow \text{this}$

Probability:

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

PCFG Example

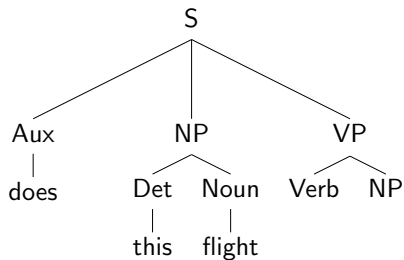


Choose a rule from the “Noun” distribution. Here: $\text{Noun} \rightarrow \text{flight}$

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

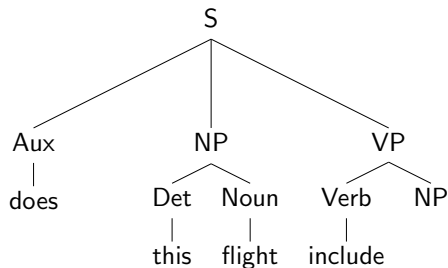


Choose a rule from the “VP” distribution. Here: $VP \rightarrow \text{Verb NP}$

Probability:

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

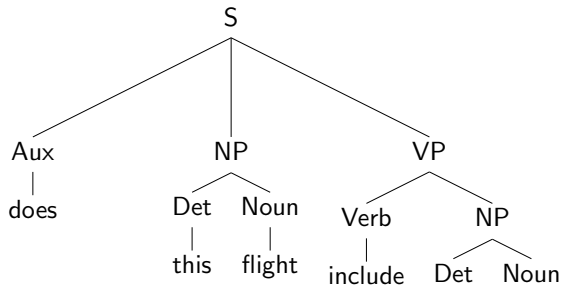
PCFG Example



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

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

PCFG Example

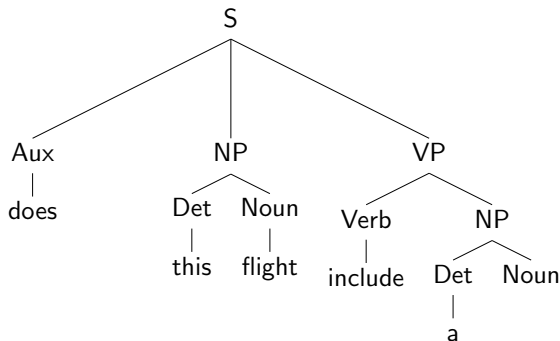


Choose a rule from the “NP” distribution. Here: $\text{NP} \rightarrow \text{Det Noun}$

Probability:

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

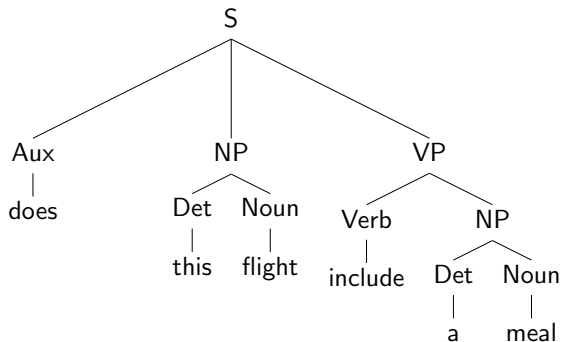
PCFG Example



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

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

PCFG Example



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

Probability:

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

PCFG as a Noisy Channel

$$\boxed{\text{source}} \longrightarrow \mathbf{T} \longrightarrow \boxed{\text{channel}} \longrightarrow \mathbf{X}$$

The PCFG defines the source model.

The channel is deterministic: it erases everything except the tree's leaves (the yield).

Decoding:

$$\begin{aligned} & \operatorname{argmax}_t p(\mathbf{t}) \cdot \begin{cases} 1 & \text{if } \mathbf{t} \in \mathcal{T}_x \\ 0 & \text{otherwise} \end{cases} \\ &= \operatorname{argmax}_{\mathbf{t} \in \mathcal{T}_x} p(\mathbf{t}) \end{aligned}$$

Probabilistic Parsing with CFGs

- ▶ How to set the probabilities $p(\text{righthand side} \mid \text{lefthand side})$?
- ▶ How to decode/parse?

Probabilistic CKY

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

Input:

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

Output: $\operatorname{argmax}_{t \in \mathcal{T}_x} \log p(t \mid x)$ (if x is in the language of the grammar)

Probabilistic CKY

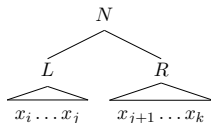
Base case: for $i \in \{1, \dots, n\}$ and for each $N \in \mathcal{N}$:

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

In general, $\heartsuit_{i:k}(N)$ is the log probability of the best tree spanning $\langle x_i, \dots, x_k \rangle$ rooted at 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:

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

Parse Chart

x_1				
	x_2			
		x_3		
			x_4	
				x_5

Parse Chart

	$\heartsuit_{1:1}(*)$				
x_1		$\heartsuit_{2:2}(*)$			
			$\heartsuit_{3:3}(*)$		
x_2				$\heartsuit_{4:4}(*)$	
					$\heartsuit_{5:5}(*)$
x_3					
x_4					
x_5					

Parse Chart

	$\heartsuit_{1:1}(\ast)$	$\heartsuit_{1:2}(\ast)$			
x_1		$\heartsuit_{2:2}(\ast)$	$\heartsuit_{2:3}(\ast)$		
			$\heartsuit_{3:3}(\ast)$	$\heartsuit_{3:4}(\ast)$	
x_2				$\heartsuit_{4:4}(\ast)$	$\heartsuit_{4:5}(\ast)$
		x_3			
			x_4		$\heartsuit_{5:5}(\ast)$
				x_5	

Parse Chart

	$\heartsuit_{1:1}(*)$	$\heartsuit_{1:2}(*)$	$\heartsuit_{1:3}(*)$		
x_1		$\heartsuit_{2:2}(*)$	$\heartsuit_{2:3}(*)$	$\heartsuit_{2:4}(*)$	
			$\heartsuit_{3:3}(*)$	$\heartsuit_{3:4}(*)$	$\heartsuit_{3:5}(*)$
x_2				$\heartsuit_{4:4}(*)$	$\heartsuit_{4:5}(*)$
					$\heartsuit_{5:5}(*)$
			x_3		
				x_4	
					x_5

Parse Chart

	$\heartsuit_{1:1}(*)$	$\heartsuit_{1:2}(*)$	$\heartsuit_{1:3}(*)$	$\heartsuit_{1:4}(*)$	
x_1		$\heartsuit_{2:2}(*)$	$\heartsuit_{2:3}(*)$	$\heartsuit_{2:4}(*)$	$\heartsuit_{2:5}(*)$
			$\heartsuit_{3:3}(*)$	$\heartsuit_{3:4}(*)$	$\heartsuit_{3:5}(*)$
x_2				$\heartsuit_{4:4}(*)$	$\heartsuit_{4:5}(*)$
					$\heartsuit_{5:5}(*)$
			x_3		
				x_4	
					x_5

Parse Chart

	$\heartsuit_{1:1}(\ast)$	$\heartsuit_{1:2}(\ast)$	$\heartsuit_{1:3}(\ast)$	$\heartsuit_{1:4}(\ast)$	$\heartsuit_{1:5}(\ast)$
x_1		$\heartsuit_{2:2}(\ast)$	$\heartsuit_{2:3}(\ast)$	$\heartsuit_{2:4}(\ast)$	$\heartsuit_{2:5}(\ast)$
x_2			$\heartsuit_{3:3}(\ast)$	$\heartsuit_{3:4}(\ast)$	$\heartsuit_{3:5}(\ast)$
				$\heartsuit_{4:4}(\ast)$	$\heartsuit_{4:5}(\ast)$
					$\heartsuit_{5:5}(\ast)$

Remarks

- ▶ Space and runtime requirements?

Remarks

- Space and runtime requirements? $O(|\mathcal{N}|n^2)$ space, $O(|\mathcal{R}|n^3)$ runtime.

Remarks

- ▶ Space and runtime requirements? $O(|\mathcal{N}|n^2)$ space, $O(|\mathcal{R}|n^3)$ runtime.
- ▶ Recovering the best tree?

Remarks

- ▶ Space and runtime requirements? $O(|\mathcal{N}|n^2)$ space, $O(|\mathcal{R}|n^3)$ runtime.
- ▶ Recovering the best tree? Backpointers.

Remarks

- ▶ Space and runtime requirements? $O(|\mathcal{N}|n^2)$ space, $O(|\mathcal{R}|n^3)$ runtime.
- ▶ Recovering the best tree? Backpointers.
- ▶ Probabilistic **Earley's** algorithm does not require the grammar to be in Chomsky normal form.

Probabilistic CKY with an Agenda

1. Initialize every item's value in the **chart** to the “default” (zero).
2. Place all initializing updates onto the **agenda**.
3. While the agenda is not empty or the goal is not reached:
 - ▶ Pop the highest-priority update from the agenda (item I with value v)
 - ▶ If $I = \text{goal}$, then return v .
 - ▶ If $v > \text{chart}(I)$:
 - ▶ $\text{chart}(I) \leftarrow v$
 - ▶ Find all combinations of I with other items in the chart, generating new possible updates; place these on the agenda.

Any priority function will work! But smart ordering will save time.

This idea can also be applied to other algorithms (e.g., Viterbi).

Demo of Recent State of the Art

<https://demo.allennlp.org/constituency-parsing>

Machine Learning and Parsing

Machine Learning and Parsing

- ▶ Define arbitrary features on trees, based on linguistic knowledge; to parse, use a PCFG to generate a **k-best list** of parses (Huang and Chiang, 2005), then train a log-linear model to rerank (Charniak and Johnson, 2005). Or a neural model (Socher et al., 2013).

Machine Learning and Parsing

- ▶ Define arbitrary features on trees, based on linguistic knowledge; to parse, use a PCFG to generate a **k-best list** of parses (Huang and Chiang, 2005), then train a log-linear model to rerank (Charniak and Johnson, 2005). Or a neural model (Socher et al., 2013).
- ▶ Define rule-local features on trees and any part of the input sentence ($s(\mathbf{x}, i, k, N, L, R)$); minimize hinge or log loss. Recent state of the art is simpler, scoring labeled spans, $s(\mathbf{x}, i, k, N)$ (Stern et al., 2017).
 - ▶ These exploit dynamic programming algorithms for training.

Machine Learning and Parsing

- ▶ Define arbitrary features on trees, based on linguistic knowledge; to parse, use a PCFG to generate a **k-best list** of parses (Huang and Chiang, 2005), then train a log-linear model to rerank (Charniak and Johnson, 2005). Or a neural model (Socher et al., 2013).
- ▶ Define rule-local features on trees and any part of the input sentence ($s(\mathbf{x}, i, k, N, L, R)$); minimize hinge or log loss. Recent state of the art is simpler, scoring labeled spans, $s(\mathbf{x}, i, k, N)$ (Stern et al., 2017).
 - ▶ These exploit dynamic programming algorithms for training.
- ▶ Learn refinements on the constituents, as latent variables (Petrov et al., 2006).

Machine Learning and Parsing

- ▶ Define arbitrary features on trees, based on linguistic knowledge; to parse, use a PCFG to generate a **k-best list** of parses (Huang and Chiang, 2005), then train a log-linear model to rerank (Charniak and Johnson, 2005). Or a neural model (Socher et al., 2013).
- ▶ Define rule-local features on trees and any part of the input sentence ($s(\mathbf{x}, i, k, N, L, R)$); minimize hinge or log loss. Recent state of the art is simpler, scoring labeled spans, $s(\mathbf{x}, i, k, N)$ (Stern et al., 2017).
 - ▶ These exploit dynamic programming algorithms for training.
- ▶ Learn refinements on the constituents, as latent variables (Petrov et al., 2006).
- ▶ **Recurrent neural network grammars**, generative models like PCFGs that encode arbitrary previous derivation steps in a vector (Dyer et al., 2016). Parsing requires some tricks.

References I

- Eugene Charniak and Mark Johnson. Coarse-to-fine n-best parsing and maxent discriminative reranking. In *Proc. of ACL*, 2005.
- John Cocke and Jacob T. Schwartz. Programming languages and their compilers: Preliminary notes. Technical report, Courant Institute of Mathematical Sciences, New York University, 1970.
- Chris Dyer, Adhiguna Kuncoro, Miguel Ballesteros, and Noah A. Smith. Recurrent neural network grammars, 2016. To appear.
- Liang Huang and David Chiang. Better k -best parsing. In *Proc. of IWPT*, 2005.
- Tadao Kasami. An efficient recognition and syntax-analysis algorithm for context-free languages. Technical Report AFCRL-65-758, Air Force Cambridge Research Lab, 1965.
- Slav Petrov, Leon Barrett, Romain Thibaux, and Dan Klein. Learning accurate, compact, and interpretable tree annotation. In *Proc. of COLING-ACL*, 2006.
- Richard Socher, John Bauer, Christopher D. Manning, and Andrew Y. Ng. Parsing with compositional vector grammars. In *Proc. of ACL*, 2013.
- Mitchell Stern, Jacob Andreas, and Dan Klein. A minimal span-based neural constituency parser. In *Proc. of ACL*, 2017.
- Daniel H. Younger. Recognition and parsing of context-free languages in time n^3 . *Information and Control*, 10(2), 1967.