

Probabilistic estimation of probabilistic syntactic models

Joan Andreu Sánchez

Departamento Sistemas Informáticos y Computación
Instituto Tecnológico de Informática
Universidad Politécnica Valencia

PASCAL Bootcamp 2010

URL: <http://www.dsic.upv.es/~jandreu>

e-mail: jandreu@dsic.upv.es

Index

- 1 Introduction to the problem
 - 1.1 Introduction
 - 1.2 Syntactic models and Parsing
 - 1.3 Syntactic models for Language Modeling
 - 1.4 Syntactic models for Machine Translation
2. Preliminaries
 - 2.1 Notation and definitions
 - 2.2 Basic probabilistic properties of syntactic models
 - 2.3 CKY-based parsing algorithms
3. Probabilistic estimation of PCFG
 - 3.1 Introduction
 - 3.2 Inside-Outside algorithm
 - 3.3 Viterbi algorithm
 - 3.4 Probabilistic properties of the estimated PCFG
 - 3.5 Use of PCFG for LM
4. Probabilistic estimation of SITG
 - 4.1 Introduction
 - 4.2 Parsing with SITG
 - 4.3 Use of SITG for MT
5. Advanced topics
 - 5.1 On-line learning of syntactic models
 - 5.2 Active learning of syntactic models
 - 5.3 Interactive-predictive parsing:
a framework for active learning

Index

1 Introduction to the problem

- 1.1 Introduction
- 1.2 Syntactic models and Parsing
- 1.3 Syntactic models for Language Modeling
- 1.4 Syntactic models for Machine Translation

2. Preliminaries

- 2.1 Notation and definitions
- 2.2 Basic probabilistic properties of syntactic models
- 2.3 CKY-based parsing algorithms

3. Probabilistic estimation of PCFG

- 3.1 Introduction
- 3.2 Inside-Outside algorithm
- 3.3 Viterbi algorithm
- 3.4 Probabilistic properties of the estimated PCFG
- 3.5 Use of PCFG for LM

4. Probabilistic estimation of SITG

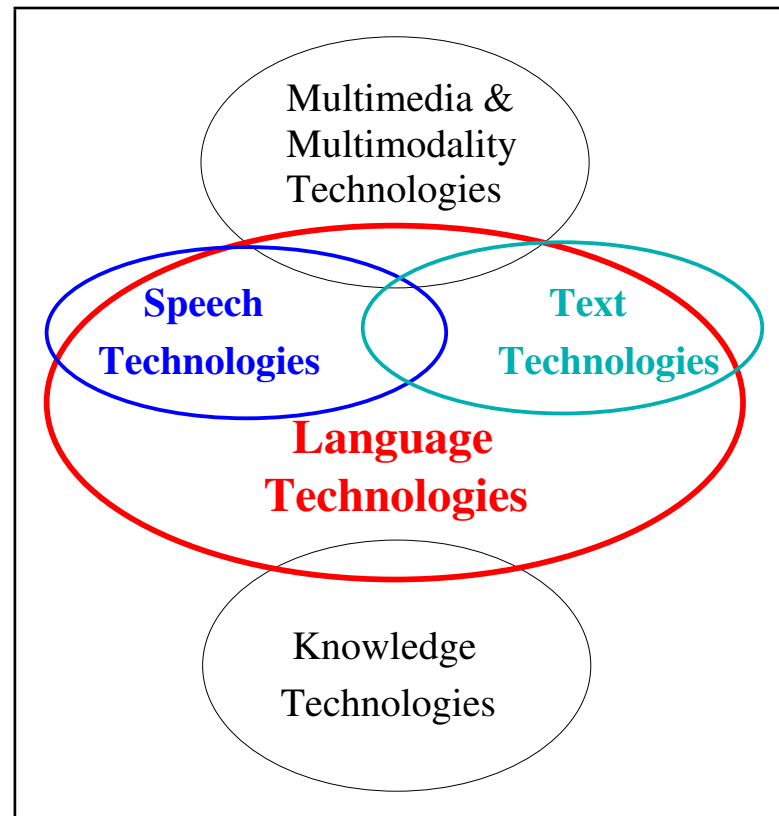
- 4.1 Introduction
- 4.2 Parsing with SITG
- 4.3 Use of SITG for MT

5. Advanced topics

- 5.1 On-line learning of syntactic models
- 5.2 Active learning of syntactic models
- 5.3 Interactive-predictive parsing:
a framework for active learning

1.1 INTRODUCTION

“Computational Linguistics” deals with the most difficult communication process:
Natural Language



1.1 INTRODUCTION

Goal:

To develop systems that are able to process, to understand and, to produce Natural Language

Motivation:

- Natural Language is the main way to represent and to transfer human knowledge
- There exist lots of information and knowledge in Natural Language
- There exist a lot of potential users that need to communicate with computers in Natural Language

Example:

Dave Bowman: Open the pod bay doors, HAL.

HAL: I'm sorry Dave. I'm afraid I can't do that.

Stanley Kubrick and Arthur C. Clarke, *2001: A Space Odyssey*.

Applications

- Systems for information extraction from text and speech

Examples:

information retrieval, information extraction, text categorization, ...

- Systems for speech/text to speech/text:

Examples:

machine translation, speech translation, speech recognition, ...

- Systems for communication with humans:

Examples:

dialog systems, query systems, ...

Probabilistic approach

- **Interpretation** by using the probabilistic decision rule
[to generate a desired interpretation (output)]
- **Modeling** the human perception with Statistical Decision techniques and Formal Language theory
[to define the statistical dependence between observations (input) and interpretation (output)]
- **Learning** knowledge from examples
[to learn the model parameters from training examples]

Main goals of the course

- **To solve difficult problems related to Natural Language with syntactic approaches**
- To study fundamentals related to Computational Linguistics
- To learn basic techniques that are necessary to develop robust systems that are able to understand text data

1.1 INTRODUCTION

Applications

- Automatic Speech Recognition
- Machine Translation
- Dialog Systems
- Automatic Summarization
- Text Classification
- Information Retrieval
- ...

Abstract tasks

- *Language Modeling*
- Part of Speech Tagging
- *Parsing*
- Lexical Disambiguation
- Semantic Analysis
- Discourse Analysis
- ...

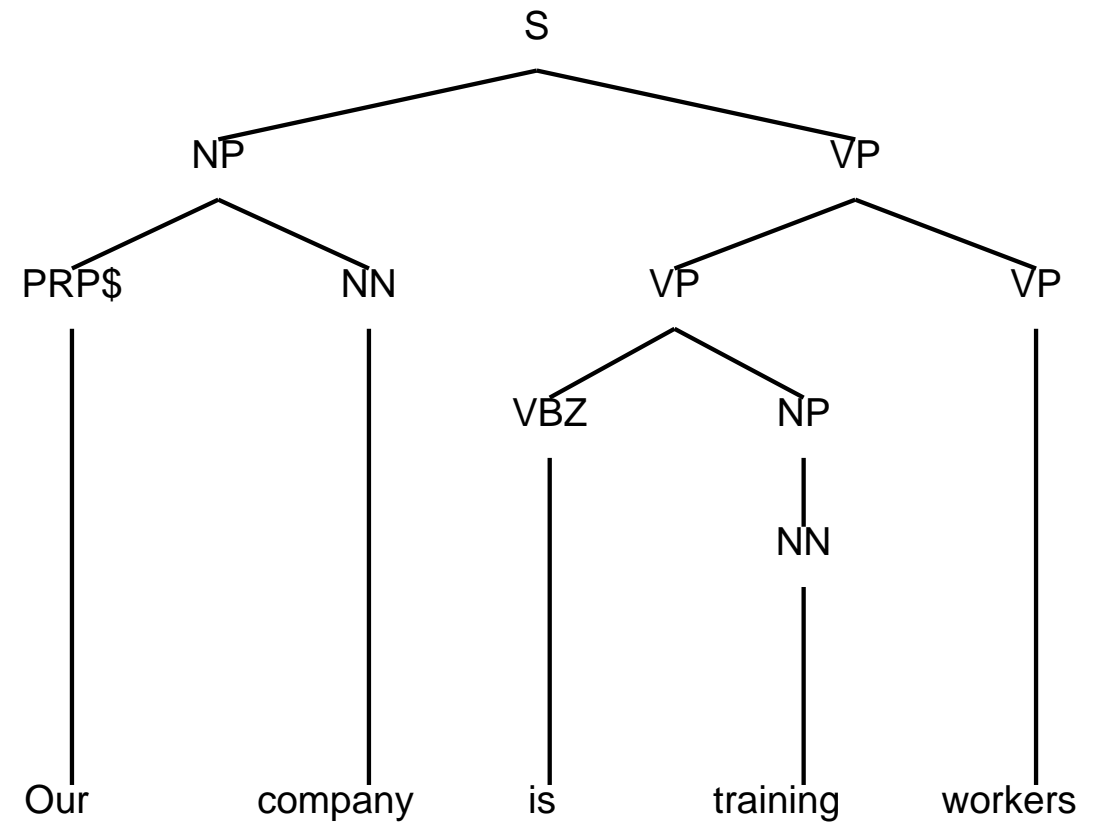
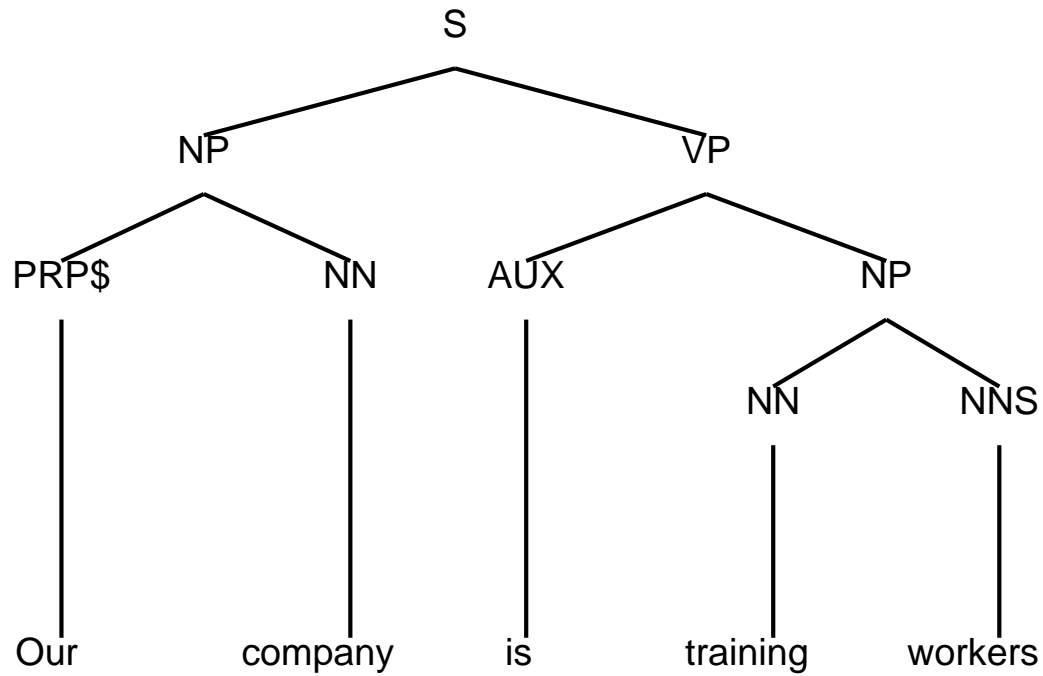
1.2 SYNTACTIC MODELS AND PARSING

Knowledge levels in Natural Language:

- Morphology: word structure
- Syntax:
 - word category — Part of Speech tagging
 - sentence structure — *Parsing*
- Semantics:
 - word semantics
 - sentence semantics
- Pragmatics: use of the language, cultural issues, environment
- Discourse: dialog structure

1.2 SYNTACTIC MODELS AND PARSING

Difficulties in Parsing: Ambiguity



1.2 SYNTACTIC MODELS AND PARSING

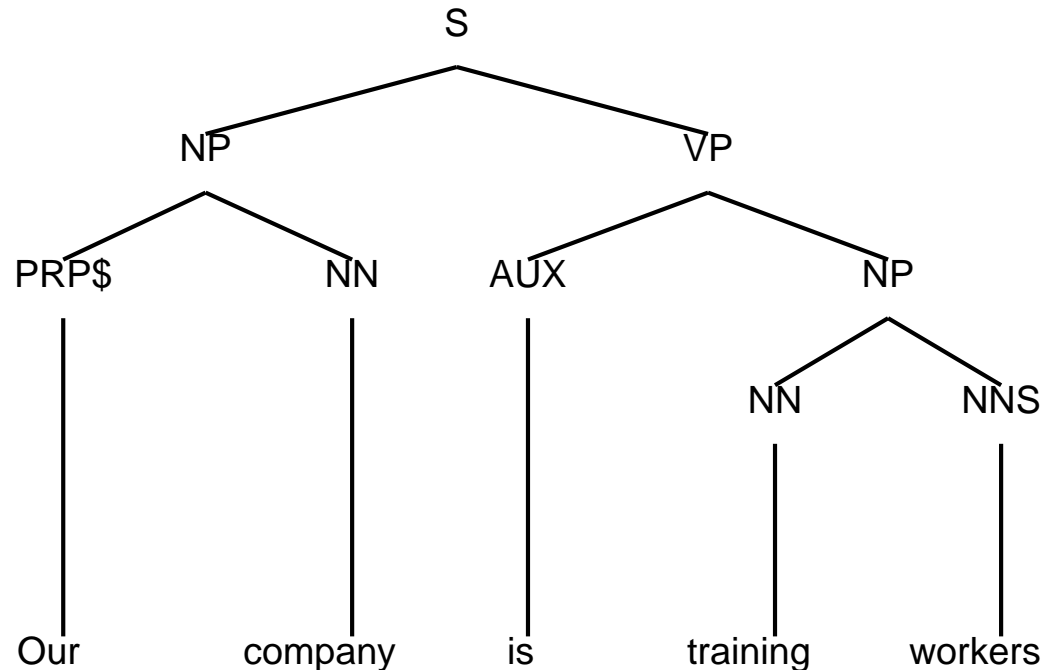
Parsing with syntactic models: (Formal) grammar

S	→	NP VP	PRP\$	→	Our
NP	→	PRP\$ NN	NN	→	company
NP	→	NN NNS	AUX	→	is
NP	→	NN	NN	→	training
VP	→	AUX NP	NNS	→	workers
VP	→	VP VP			
VP	→	VBZ NP			

1.2 SYNTACTIC MODELS AND PARSING

Parsing with syntactic models: (Formal) grammar

1.0	S	→	NP VP	1.0	PRP\$	→	Our
0.4	NP	→	PRP\$ NN	0.6	NN	→	company
0.3	NP	→	NN NNS	1.0	AUX	→	is
0.3	NP	→	NN	0.4	NN	→	training
0.5	VP	→	AUX NP	1.0	NNS	→	workers
0.3	VP	→	VP VP				
0.2	VP	→	VBZ NP				



1.3 SYNTACTIC MODELS FOR LANGUAGE MODELING

Recognition with noisy channel

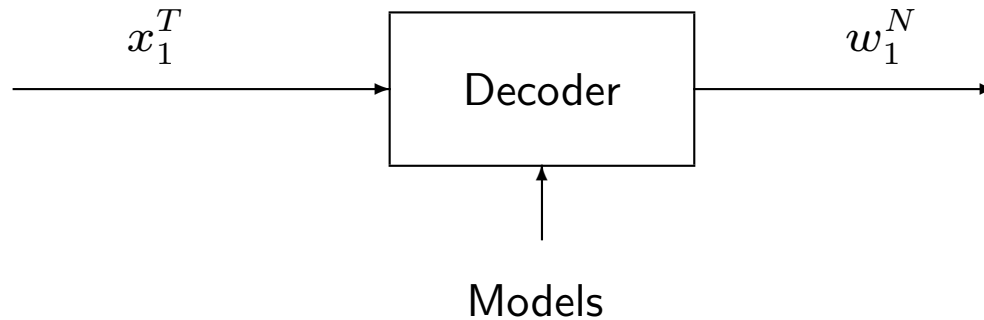


$$\hat{I} = \arg \max_I \Pr(I|O) = \arg \max_I \Pr(O|I) \Pr(I)$$

$\Pr(I)$: language model probability

$\Pr(O|I)$: channel probability

Automatic Speech Recognition



$$\widehat{w_1^N} = \arg \max_{w_1^N} \Pr(w_1^N | x_1^T) = \arg \max_{w_1^N} \Pr(x_1^T | w_1^N) \Pr(w_1^N)$$

Language Model

$$\Pr(w_1^N) = \Pr(w_1) \prod_{n=2}^N \Pr(w_n | w_1^{n-1})$$

1.3 SYNTACTIC MODELS FOR LANGUAGE MODELING

- **N-Gram models:** Restriction on the history length w_1^{n-1}

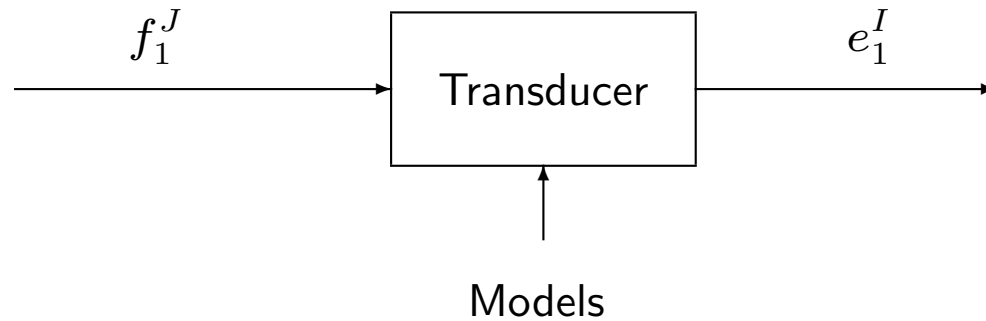
$$\Pr(w_1^N) = \Pr(w_1) \prod_{n=2}^N \Pr(w_n | w_{n-k+1}^{n-1})$$

- ✗ don't capture long-term dependencies
 - ✓ efficient to compute
 - ✓ efficient methods to estimate the model parameters
- **Grammatical models:** No restriction on the history length w_1^{n-1}

$$\Pr(w_1^N) = \Pr(w_1) \prod_{n=2}^N \Pr(w_n | w_1^{n-1})$$

- ✓ capture long-term dependencies
- ✗ expensive to compute
- ✗ efficient methods to estimate the model parameters, but expensive

Machine Translation



- Inverse approach (noisy channel)

$$\hat{e}_1^I = \arg \max_{e_1^I} \Pr(e_1^I | f_1^J) = \arg \max_{e_1^I} \Pr(f_1^J | e_1^I) \Pr(e_1^I)$$

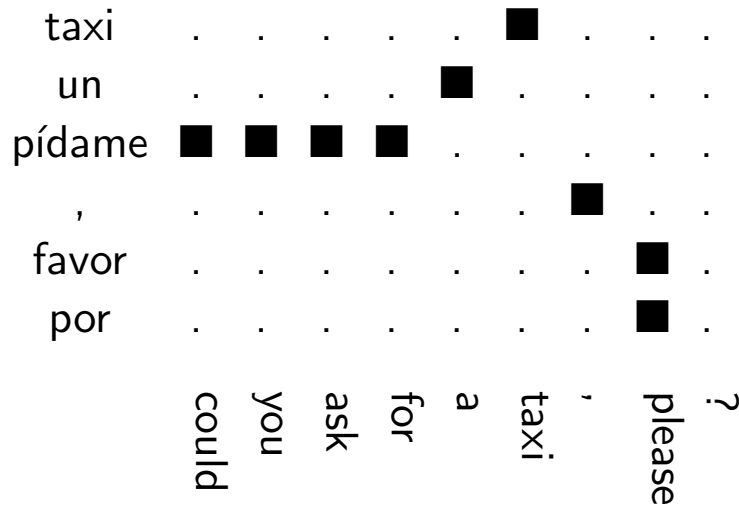
- Direct approach

$$\hat{e}_1^I = \arg \max_{e_1^I} \Pr(e_1^I | f_1^J)$$

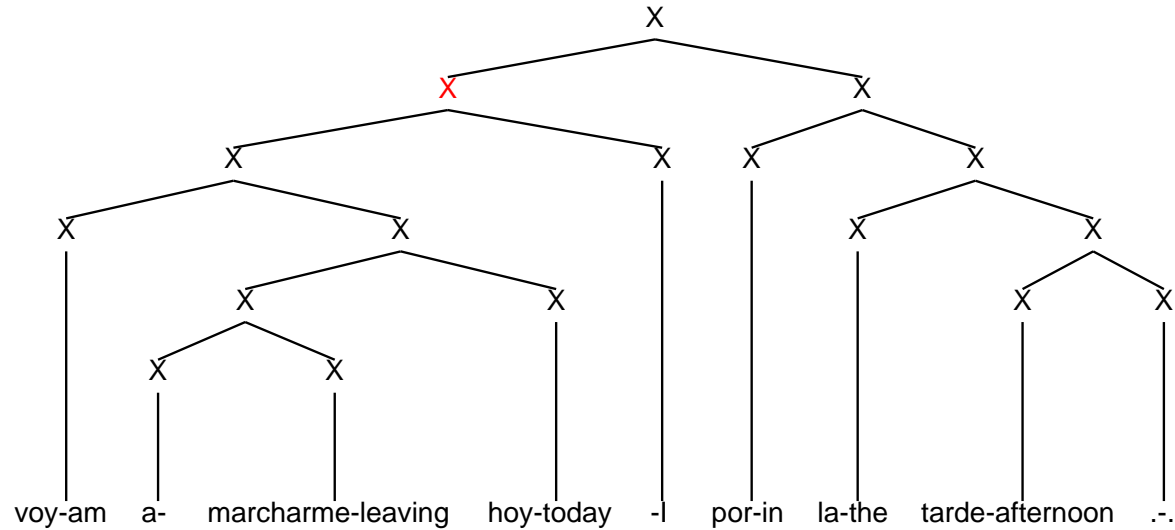
1.4 SYNTACTIC MODELS FOR MACHINE TRANSLATION

Statistical approaches to MT:

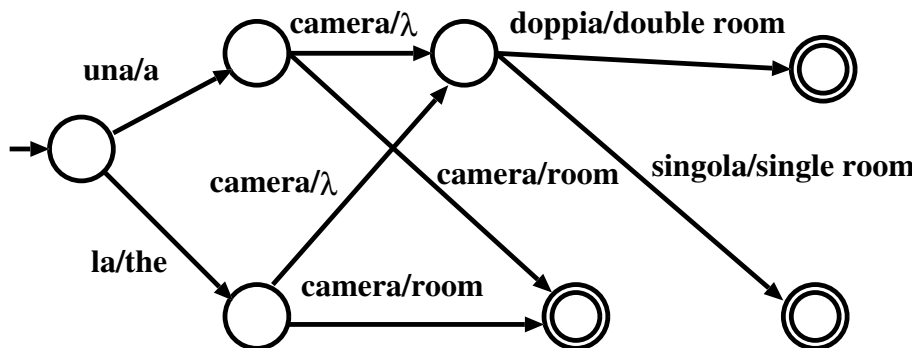
→ Word-alignment approaches



→ Syntactic approaches



→ Finite-state approaches



Index

- 1 Introduction to the problem
 - 1.1 Introduction
 - 1.2 Syntactic models and Parsing
 - 1.3 Syntactic models for Language Modeling
 - 1.4 Syntactic models for Machine Translation
2. **Preliminaries**
 - 2.1 Notation and definitions
 - 2.2 Basic probabilistic properties of syntactic models
 - 2.3 CKY-based parsing algorithms
3. Probabilistic estimation of PCFG
 - 3.1 Introduction
 - 3.2 Inside-Outside algorithm
 - 3.3 Viterbi algorithm
 - 3.4 Probabilistic properties of the estimated PCFG
 - 3.5 Use of PCFG for LM
4. Probabilistic estimation of SITG
 - 4.1 Introduction
 - 4.2 Parsing with SITG
 - 4.3 Use of SITG for MT
5. Advanced topics
 - 5.1 On-line learning of syntactic models
 - 5.2 Active learning of syntactic models
 - 5.3 Interactive-predictive parsing:
a framework for active learning

2.1 NOTATION AND DEFINITIONS

Context-free grammar: [Aho,72]

- Simple and compact models for parsing
- Formal framework well understood
- Natural Language is no regular (but almost)
- Adequate representation of long-term syntactic structures
- Adequate modeling of ambiguity

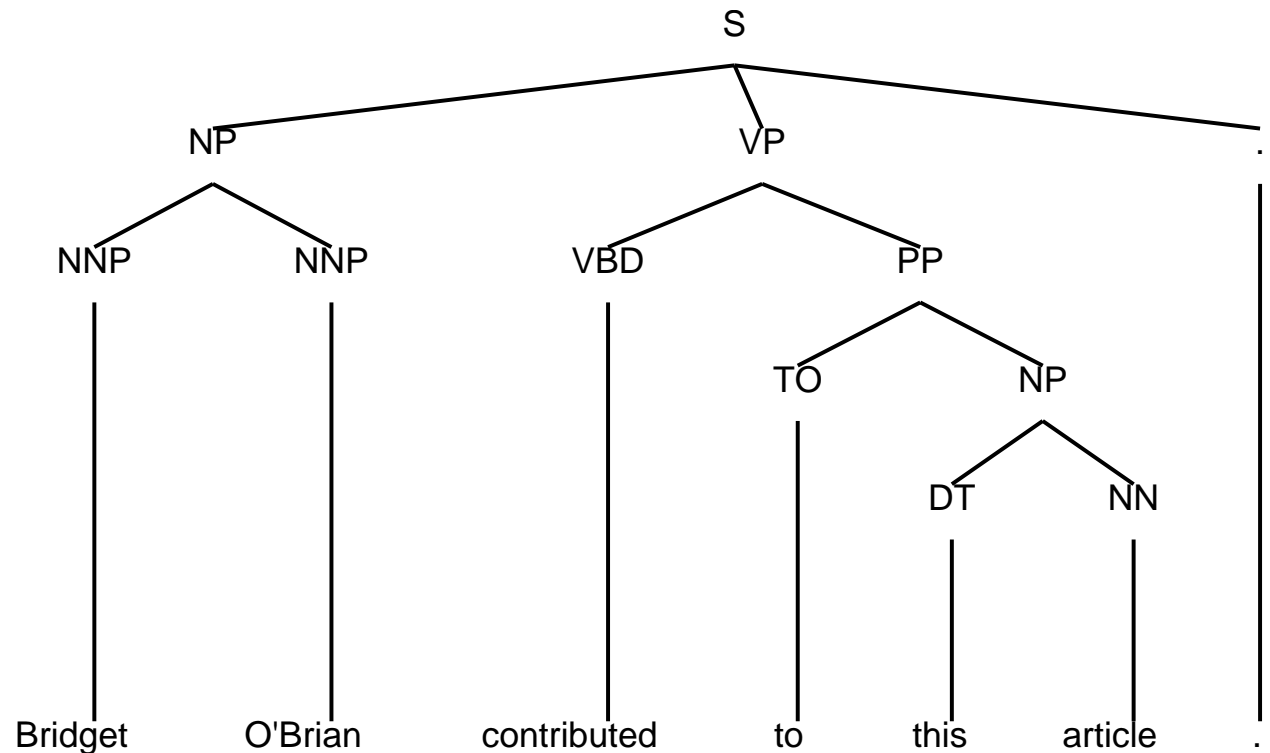
2.1 NOTATION AND DEFINITIONS

Example

- Primitives: alphabet
words, punctuation symbols, . . .
- Object representation: written sentences
“Bridget O’Brian contributed to this article”

- Pattern set
sentences

- Interpretation: syntactic analysis



2.1 NOTATION AND DEFINITIONS

- An **alphabet** T is a finite set of symbols.
- A **string** $x = a_1 \cdots a_n$ ($a_i \in T; i : 1 \dots n$), is a finite sequence of symbols of T . The length of the string is noted by $|x|$. Let x and y be two strings, $x, y \in T^*$, then the **concatenation** of x and y is the string xy . $|xy| = |x| + |y|$.
- The **empty string** ϵ , is the string with length equal to zero. For any string x , $x \in T^*$:
 $\epsilon x = x \epsilon = x$.
- The **closure** T^* is the infinite and countable set of all strings with finite length composed with symbols of T , ϵ included. The **positive closure** T^+ is defined as:
 $T^+ = T^* - \{\epsilon\}$.
- A **language** L is a set of strings composed with symbols of T ($L \subseteq T^*$).

2.1 NOTATION AND DEFINITIONS

➤ **Grammar:** $G = (N, T, P, S)$

$$V = N \cup T; N \cap T = \emptyset; S \in N; (A \rightarrow \beta) \in P; \quad A \in N; \beta \in V^*$$

➤ **Derivation:**

$$\mu A \delta \implies \mu \beta \delta \text{ iff } \exists (A \rightarrow \beta) \in P; \quad \mu, \delta \in V^*$$

➤ **Sentential Form:**

$$\alpha \in V^* \text{ is a } \textit{sentential form} \text{ of } G \text{ if } S \xRightarrow{*} \alpha$$

➤ **Language generated by G :**

$$L(G) = \{x \in T^* \mid S \xRightarrow{*} x\}$$

➤ **Grammar classification:**

- **Type 2:** *context free grammars*

$$A \rightarrow \beta$$

$$A \in N; \beta \in V^*$$

- **Type 3:** *regular grammars*

$$A \rightarrow aB, \quad A \rightarrow a$$

$$A, B \in N; a \in T$$

2.1 NOTATION AND DEFINITIONS

Approaches

- Top-Down parsing
- Down-Top parsing

Depending on time complexity

- Backtracking methods Exponential complexity
- Deterministic methods Linear complexity
 - Grammars: LL(1), SLR(1), LALR(1), LR(1), ...
- Tabular methods Cubic complexity
 - CKY algorithm**
 - Earley algorithm [Aho 72, Stolcke 95]

2.1 NOTATION AND DEFINITIONS

ALGORITHM: Cocke-Kasami-Younger

Input $G = (N, T, P, S)$ in CNF and $\mathbf{x} = x_1 \dots x_n \in T^*$

Output Parsing table $t[i, l]$ ($1 \leq i, l \leq n$)

$A \in t[i, l]$ if $A \xrightarrow{*} x_{i+1} \dots x_l$

METHOD

for all $i : 0 \dots n - 1$ **do**

$t[i, i + 1] := t[i, i + 1] \cup \{A \mid (A \rightarrow b) \in P; b = x_{i+1}\}$

for all $j : 2 \dots n$ **do**

for all $i : 0 \dots n - j$ **do**

for all $k : 1 \dots j - 1$ **do**

$t[i, i + j] := t[i, i + j] \cup \{A \mid (A \rightarrow BC) \in P;$

$B \in t[i, i + k]; C \in t[i + k, i + j]\}$

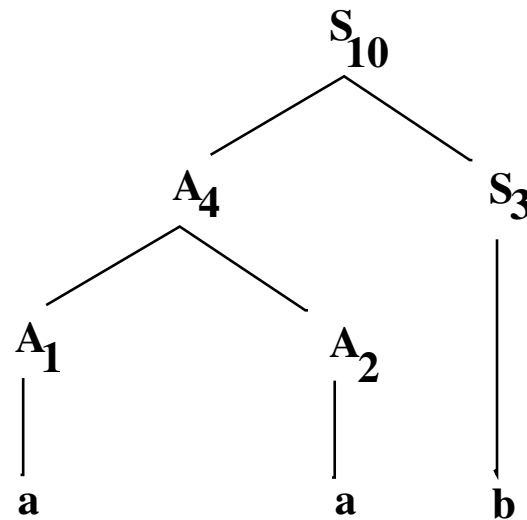
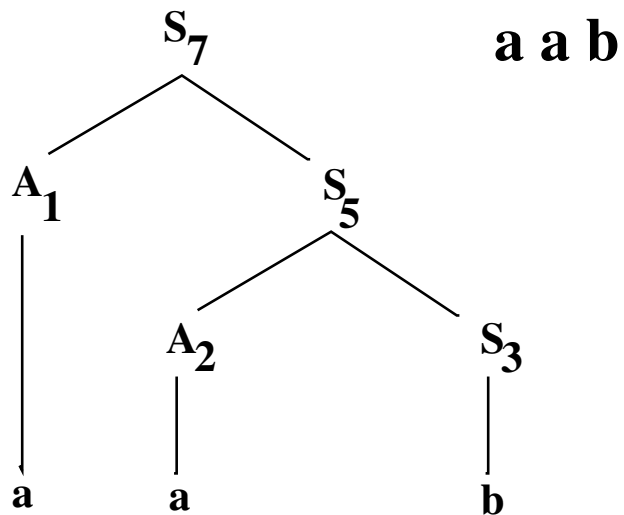
if $S \in t[0, n]$ **then** $x \in L(G)$ **else** $x \notin L(G)$

END

2.1 NOTATION AND DEFINITIONS

$S \rightarrow AS$
 $S \rightarrow b$
 $A \rightarrow AS$
 $A \rightarrow AA$
 $A \rightarrow a$

$l = 1$	$l = 2$	$l = 3$	
1: (A, 0, 0)	4: (A, 1, 2)	7: (S, 1, 5) 8: (A, 1, 5) 9: (A, 1, 6) 10: (S, 4, 3) 11: (A, 4, 3)	$i = 0$
	2: (A, 0, 0)	5: (S, 2, 3) 6: (A, 2, 3)	$i = 1$
		3: (S, 0, 0)	$i = 2$



2.1 NOTATION AND DEFINITIONS

Let $x \in T^*$ and a stochastic model M characterized by a parameter vector θ , we are interested in computing: $p_\theta(x)$

Stochastic language (L, ϕ) over T [Wetherell 80]:

- $L \subseteq T^*$ characteristic language
- $\phi : T^* \longrightarrow [0, 1]$ computable stochastic function:

- i) $x \notin L \implies \phi(x) = 0 \quad \forall x \in T^*$
- ii) $x \in L \implies 0 < \phi(x) \leq 1 \quad \forall x \in T^*$
- iii) $\sum_{x \in L} \phi(x) = 1$

Example [Booth 73]

Given the alphabet $T = \{a, b\}$, the following language is defined: $L = \{a^n b^n \mid n \geq 0\}$, where $\phi(x) = 0, \forall x \notin L$ and $\phi(a^n b^n) = \frac{1}{en!}$

$$\sum_{x \in L} \phi(x) = \sum_{0 \leq n \leq \infty} \frac{1}{en!} = \frac{1}{e} \sum_{0 \leq n \leq \infty} \frac{1}{n!} = \frac{1}{e} e = 1$$

2.1 NOTATION AND DEFINITIONS

Probabilistic context-free grammar: $G_s = (G, p)$

➤ $G = (N, T, P, S)$ characteristic grammar

➤ $p : P \rightarrow]0, 1]$ probability of the rules. $\forall A_i \in N$:

$$\sum_{1 \leq j \leq n_i} p(A_i \rightarrow \alpha_j) = 1,$$

where n_i is the number of rules with A_i in the left side of the rules.

Stochastic derivation for PCFG

Given a sequence of stochastic events:

$$S = \alpha_0 \xrightarrow{r_1} \alpha_1 \xrightarrow{r_2} \alpha_2 \cdots \alpha_{m-1} \xrightarrow{r_m} \alpha_m = x$$

the probability of x being generated by $G_s = (G, p)$ from the rule sequence $d_x = r_1, \dots, r_m$, is:

$$\Pr_{G_s}(x, d_x) = p(r_1)p(r_2 \mid r_1) \cdots p(r_m \mid r_1 \cdots r_{m-1})$$

➤ **problem:** computation of the probabilities

➤ **restriction:** $p(r_j \mid r_1 \cdots r_{j-1}) = p(r_j)$

$$\Pr_{G_s}(x, d_x) = \prod_{j=1 \cdots m} p(r_j)$$

2.1 NOTATION AND DEFINITIONS

Probability of a derivation $d_x = r_1, \dots, r_m$

$$\Pr_{G_s}(x, d_x) = \prod_{j=1 \dots m} p(r_j) = \prod_{\forall (A \rightarrow \alpha) \in P} p(A \rightarrow \alpha)^{N(A \rightarrow \alpha, d_x)}$$

Probability of a string

$$\Pr_{G_s}(x) = \sum_{d_x \in D_x} \Pr_{G_s}(x, d_x)$$

Probability of the best derivation

$$\widehat{\Pr}_{G_s}(x) = \max_{d_x \in D_x} \Pr_{G_s}(x, d_x)$$

Probability of a string with a subset of derivations $\Delta_x \subseteq D_x$

$$\Pr_{G_s}(x, \Delta_x) = \sum_{d_x \in \Delta_x} \Pr_{G_s}(x, d_x)$$

Language generated by a PCFG

$$L(G_s) = \{x \in L(G) \mid \Pr_{G_s}(x) > 0\}$$

Consistent grammar

A PCFG $G_s = (G, p)$ is consistent iff:

$$\sum_{x \in L(G)} \Pr_{G_s}(x) = 1$$

Theorem [Booth 73]

There exist stochastic languages (L, ϕ) that can not be generated by a stochastic grammar $G_s = (G, p)$

Dem. outline Let $L = \{a^n b^n \mid n \geq 0\}$ be a stochastic language:

$$\phi(a^n b^n) = \frac{1}{en!}$$

There is not any G_s such that $\phi(x) = \Pr_{G_s}(x) \quad \forall x \in L$

2.3 CKY-BASED PARSING ALGORITHMS

Inside algorithm for PCFG [Lari 90]

- Given $x = x_1 \dots x_n \in T^*$ and $A \in N$

$$e(A \langle i, l \rangle) = \Pr_{G_s}(A \xRightarrow{*} x_i \dots x_l)$$

- Compute $\forall A \in N$:

$$e(A \langle i, i \rangle) = p(A \rightarrow b) \delta(b, x_i) \quad 1 \leq i \leq n$$

$$e(A \langle i, j \rangle) = \sum_{B, C \in N} p(A \rightarrow BC) \sum_{k=i}^{j-1} e(B \langle i, k \rangle) e(C \langle k+1, j \rangle)$$

$$1 \leq i < j \leq n$$

- $\Pr_{G_s}(x) = e(S \langle 1, n \rangle)$
- Time complexity: $O(|x|^3 |P|)$

2.3 CKY-BASED PARSING ALGORITHMS

Inside algorithm for PCFG (bracketed version [Pereira 92])

➤ Bracketed sentence:

(((Pierre Vinken) , ((61 years) old) ,) (will (join (the board) (as (a nonexecutive director)) (Nov. 29.))) .)

$$c(i, j) = \begin{cases} 1 & \text{if } (i, j) \text{ does not overlap any span in the sentence,} \\ 0 & \text{otherwise.} \end{cases}$$

➤ Compute $\forall A \in N$:

$$e(A \langle i, i \rangle) = p(A \rightarrow b) \delta(b, x_i) \quad 1 \leq i \leq n$$

$$e(A \langle i, j \rangle) = c(i, j) \sum_{B, C \in N} p(A \rightarrow BC) \sum_{k=i}^{j-1} e(B \langle i, k \rangle) e(C \langle k+1, j \rangle)$$

$$1 \leq i < j \leq n$$

➤ Linear if full bracketing

2.3 CKY-BASED PARSING ALGORITHMS

Viterbi algorithm for PCFG [Ney 91]

- Given $x = x_1 \dots x_n \in T^*$ and $A \in N$

$$\hat{e}(A \langle i, l \rangle) = \widehat{\text{Pr}}_{G_s}(A \xrightarrow{*} x_i \dots x_l)$$

- Compute $\forall A \in N$:

$$\hat{e}(A \langle i, i \rangle) = p(A \rightarrow b) \delta(b, x_i) \quad 0 \leq i \leq n$$

$$\hat{e}(A \langle i, j \rangle) = \max_{B, C \in N} p(A \rightarrow BC) \max_{k=i, \dots, j-1} \hat{e}(B \langle i, k \rangle) \hat{e}(C \langle k+1, j \rangle)$$

$$1 \leq i < j \leq n$$

- $\widehat{\text{Pr}}_{G_s}(x) = \hat{e}(S \langle 1, n \rangle)$

- Time complexity: $O(|x|^3 |P|)$ (Bracketed version: linear if full bracketing)

Outside algorithm for PCFG

- Given $x = x_1 \dots x_n \in T^*$ and $A \in N$

$$f(A \langle i, l \rangle) = \Pr_{G_s}(S \xRightarrow{*} x_1 \dots x_{i-1} A x_{l+1} \dots x_n)$$

- Compute $\forall A \in N$:

$$f(A \langle 1, n \rangle) = \delta(A, S)$$

$$f(A \langle i, j \rangle) = \sum_{B, C \in N} p(B \rightarrow CA) \sum_{k=1}^{i-1} f(B \langle k, j \rangle) e(C \langle k, i-1 \rangle)$$

$$+ \sum_{B, C \in N} p(B \rightarrow AC) \sum_{k=j+1}^n f(B \langle i, k \rangle) e(C \langle j+1, k \rangle)$$

$$1 \leq i \leq j \leq n$$

- $\Pr_{G_s}(x) = \sum_{A \in N} f(A \langle i, i \rangle) p(A \rightarrow x_i),$

$$1 \leq i \leq n$$

- Time complexity: $O(|x|^3|P|)$ (Bracketed version: linear if full bracketing)

2.3 CKY-BASED PARSING ALGORITHMS

Probability of an initial substring: LRI algorithm

$$T(A \Rightarrow B) = \sum_{\alpha} \Pr_{G_s}(A \xrightarrow{*} B\alpha)$$

$$T(A \Rightarrow BC) = p(A \rightarrow BC) + \sum_D T(A \Rightarrow D)p(D \rightarrow BC)$$

➤ Given $x = x_1 \dots x_n \in T^*$ and $A \in N$

$$e(A \ll i, l) = \Pr_{G_s}(A \xrightarrow{*} x_i \dots x_l \dots)$$

➤ Compute $\forall A \in N$:

$$e(A \ll i, i) = p(A \rightarrow x_i) + \sum_D T(A \Rightarrow D) p(D \rightarrow x_i) \quad 1 \leq i \leq n$$

$$e(A \ll i, j) = \sum_{B, C \in N} T(A \Rightarrow BC) \sum_{k=i}^{j-1} e(B \ll i, k) e(C \ll k+1, j) \quad 1 \leq i < j$$

➤ $\Pr_{G_e}(x_1 \dots x_k \dots) = e(S \ll 1, k)$

➤ Time complexity: $O(|x|^3|P|)$

Index

- 1 Introduction to the problem
 - 1.1 Introduction
 - 1.2 Syntactic models and Parsing
 - 1.3 Syntactic models for Language Modeling
 - 1.4 Syntactic models for Machine Translation
2. Preliminaries
 - 2.1 Notation and definitions
 - 2.2 Basic probabilistic properties of syntactic models
 - 2.3 CKY-based parsing algorithms
- 3. Probabilistic estimation of PCFG**
 - 3.1 Introduction
 - 3.2 Inside-Outside algorithm
 - 3.3 Viterbi algorithm
 - 3.4 Probabilistic properties of the estimated PCFG
 - 3.5 Use of PCFG for LM
4. Probabilistic estimation of SITG
 - 4.1 Introduction
 - 4.2 Parsing with SITG
 - 4.3 Use of SITG for MT
5. Advanced topics
 - 5.1 On-line learning of syntactic models
 - 5.2 Active learning of syntactic models
 - 5.3 Interactive-predictive parsing:
a framework for active learning

3.1 INTRODUCTION

➤ Supervised methods

- Maximum likelihood estimation

$$\widehat{\text{Pr}}(A \rightarrow \alpha) = \frac{C(A \rightarrow \alpha)}{C(A)}$$

- Annotated data is needed (“treebank”)

➤ Non-supervised methods

- EM algorithms
- Problem: local optimum

3.1 INTRODUCTION

Let G_s a PCFG with parameters θ and a sample $\Omega = \{x_1, x_2, \dots, x_n\}$.

$$\hat{\theta} = \arg \max_{\theta} F_{\theta}(\Omega)$$

➤ Optimization method

➤ Growth transformations

➤ Optimization function

➤ **Maximum likelihood**

➤ Corrective training

➤ Maximum mutual information

3.1 INTRODUCTION

Theorem [Baum 72]

Let $P(\Theta)$ be a homogeneous polynomial with non-negative coefficients. Let $\theta = \{\theta_{ij}\}$ be a point in the domain $D = \{\theta_{ij} \mid \theta_{ij} \geq 0; \sum_{j=1}^{q_i} \theta_{ij} = 1, i = 1, \dots, p; j = 1, \dots, q_i\}$, and let $Q(\theta)$ be a close transformation in D , that is defined as:

$$Q(\theta)_{ij} = \frac{\theta_{ij}(\partial P / \partial \Theta_{ij})_{\theta}}{\sum_{k=1}^{q_i} \theta_{ik}(\partial P / \partial \Theta_{ik})_{\theta}}$$

with the denominator different from zero. Then, $P(Q(\theta)) > P(\theta)$ except if $Q(\theta) = \theta$.

```
input  $P(\Theta)$ 
 $\theta$  = initial values
repeat
    compute  $Q(\theta)$  using  $P(\Theta)$ 
     $\theta = Q(\theta)$ 
until convergence
output  $\theta$ 
```


3.2 INSIDE-OUTSIDE ALGORITHM

Let a PCFG G_s , a sample Ω and a set of derivations Δ_x for each $x \in \Omega$

$$\Pr_{G_s}(\Omega, \Delta_\Omega) = \prod_{x \in \Omega} \Pr_{G_s}(x, \Delta_x)$$

$\forall (A \rightarrow \alpha) \in P$ (See demonstration [Benedí 05])

$$\bar{p}(A \rightarrow \alpha) = \frac{\sum_{x \in \Omega} \frac{1}{\Pr_{G_s}(x, \Delta_x)} \sum_{\forall d_x \in \Delta_x} N(A \rightarrow \alpha, d_x) \Pr_{G_s}(x, d_x)}{\sum_{x \in \Omega} \frac{1}{\Pr_{G_s}(x, \Delta_x)} \sum_{\forall d_x \in \Delta_x} N(A, d_x) \Pr_{G_s}(x, d_x)}$$

3.2 INSIDE-OUTSIDE ALGORITHM

Optimization function ($\Delta_x = D_x$)

$$\Pr_{G_s}(\Omega) = \prod_{x \in \Omega} \Pr_{G_s}(x)$$

➤ $\forall(A \rightarrow BC) \in P; \forall(A \rightarrow b) \in P$ (See demonstration)

$$\bar{p}(A \rightarrow BC) =$$

$$\frac{\sum_{x \in \Omega} \frac{p(A \rightarrow BC)}{\Pr_{G_s}(x)} \sum_{i=0}^{n-j} \sum_{j=2}^n \sum_{k=1}^{j-1} f(A \langle i, i+j \rangle) e(B \langle i, i+k \rangle) e(C \langle i+k, i+j \rangle)}{\sum_{x \in \Omega} \frac{1}{\Pr_{G_s}(x)} \sum_{i=0}^{n-j} \sum_{j=1}^n f(A \langle i, i+j \rangle) e(A \langle i, i+j \rangle)}$$

$$\bar{p}(A \rightarrow b) = \frac{\sum_{x \in \Omega} \frac{1}{\Pr_{G_s}(x)} \sum_{i=0, b=x_i}^{n-1} f(A \langle i, i \rangle) p(A \rightarrow b)}{\sum_{x \in \Omega} \frac{1}{\Pr_{G_s}(x)} \sum_{i=0}^{n-j} \sum_{j=1}^n f(A \langle i, i+j \rangle) e(A \langle i, i+j \rangle)}$$

➤ Time complexity: $O(|LT|^3|P|)$

3.3 VITERBI ALGORITHM

Optimization function ($\Delta_x = \hat{d}_x$) [Benedí 05]

$$\Pr_{G_s}(\hat{\Omega}) = \prod_{x \in \Omega} \Pr_{G_s}(x, \hat{d}_x)$$

➤ $\forall (A \rightarrow \alpha) \in P$

$$\bar{p}(A \rightarrow \alpha) = \frac{\sum_{x \in \Omega} N(A \rightarrow \alpha, \hat{d}_x)}{\sum_{x \in \Omega} N(A, \hat{d}_x)}.$$

➤ Time complexity: $O(|LT|^3|P|)$

3.3 VITERBI ALGORITHM

1. Carrying out the maximization $G_s^{(i)}$:

$$\widehat{D}_x^{(i)} = \{\widehat{d}_x^{(i)} : \widehat{d}_x^{(i)} = \arg \max_{d_x \in D_x} \Pr_{G_s^{(i)}}(x, d_x)\}$$

2. Applying the transformation: $G_s^{(i+1)}$.

The function to be optimized is defined after step 1. This function is continuous and differentiable:

$$\prod_{x \in \Omega} \Pr_{G_s^{(i)}}(x, \widehat{d}_x^{(i)}) \leq \prod_{x \in \Omega} \Pr_{G_s^{(i+1)}}(x, \widehat{d}_x^{(i)})$$

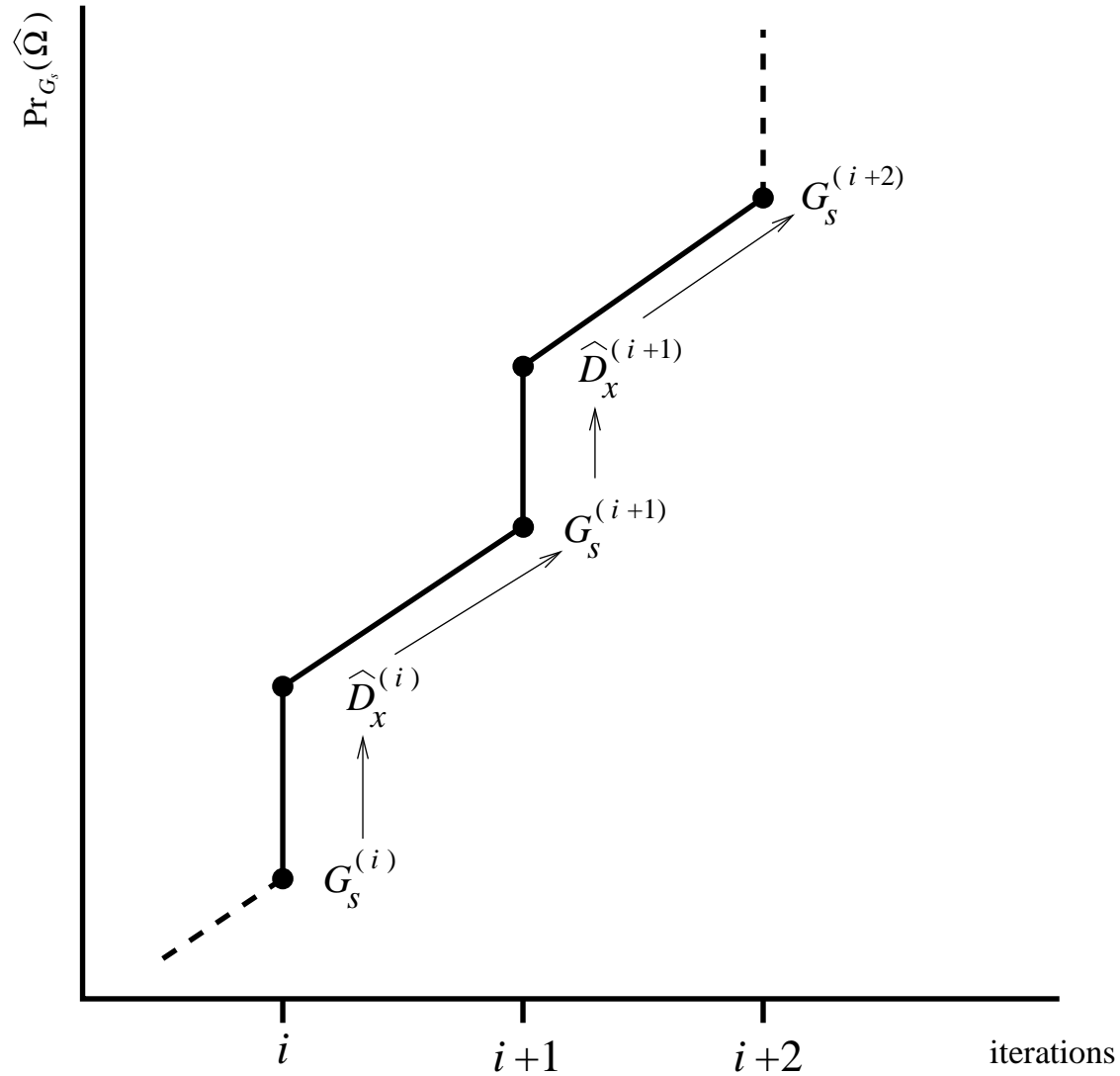
In the next step $i + 1$, the most probable tree $\widehat{D}_x^{(i+1)}$ is computed for each string x with $G_s^{(i+1)}$, and therefore:

$$\Pr_{G_s^{(i+1)}}(x, \widehat{d}_x^{(i)}) \leq \Pr_{G_s^{(i+1)}}(x, \widehat{d}_x^{(i+1)}) \quad \forall x \in \Omega,$$

and hence

$$\prod_{x \in \Omega} \Pr_{G_s^{(i+1)}}(x, \widehat{d}_x^{(i)}) \leq \prod_{x \in \Omega} \Pr_{G_s^{(i+1)}}(x, \widehat{d}_x^{(i+1)})$$

3.3 VITERBI ALGORITHM



3.4 PROBABILISTIC PROPERTIES OF THE ESTIMATED PCFG

Theorem [Booth 73] A PCFG is consistent if $\rho(E) < 1$, where $\rho(E)$ is the spectral radius (absolute value of the largest eigenvalue) of matrix E .

Probabilistic expectation matrix: $E = (e_{ij})$, expected number of times that the non-terminal A_j is derived directly from A_i :

$$e_{ij} = \sum_{(A_i \rightarrow \alpha)} p(A_i \rightarrow \alpha) N(A_j, \alpha) \quad 1 \leq i, j \leq |N|$$

Expectation matrix

$$Q = \sum_{i=0}^{\infty} E^i. \quad \text{If } G_s \text{ is consistent, then the sum converges to: } Q = (I - E)^{-1}$$

Theorem [Sánchez 97] Let $G_s = (G, p)$ be a PCFG and let Ω be a sample from $L(G)$. If $\overline{G}_s = (G, \overline{p})$ is a PCFG obtained from G_s when applying the previous growth transformation, the \overline{G}_s is consistent.

3.4 PROBABILISTIC PROPERTIES OF THE ESTIMATED PCFG

Palindrome language

$$\{ ww^R \mid w \in \{a, b\}^+; R = \text{reverse string} \}$$

➤ Original model

$$\begin{array}{llll} S \rightarrow AC & 0.4 & S \rightarrow BB & 0.1 & C \rightarrow SA & 1.0 & A \rightarrow a & 1.0 \\ S \rightarrow BD & 0.4 & S \rightarrow AA & 0.1 & D \rightarrow SB & 1.0 & B \rightarrow b & 1.0 \end{array}$$

➤ Training set: 1000 strings

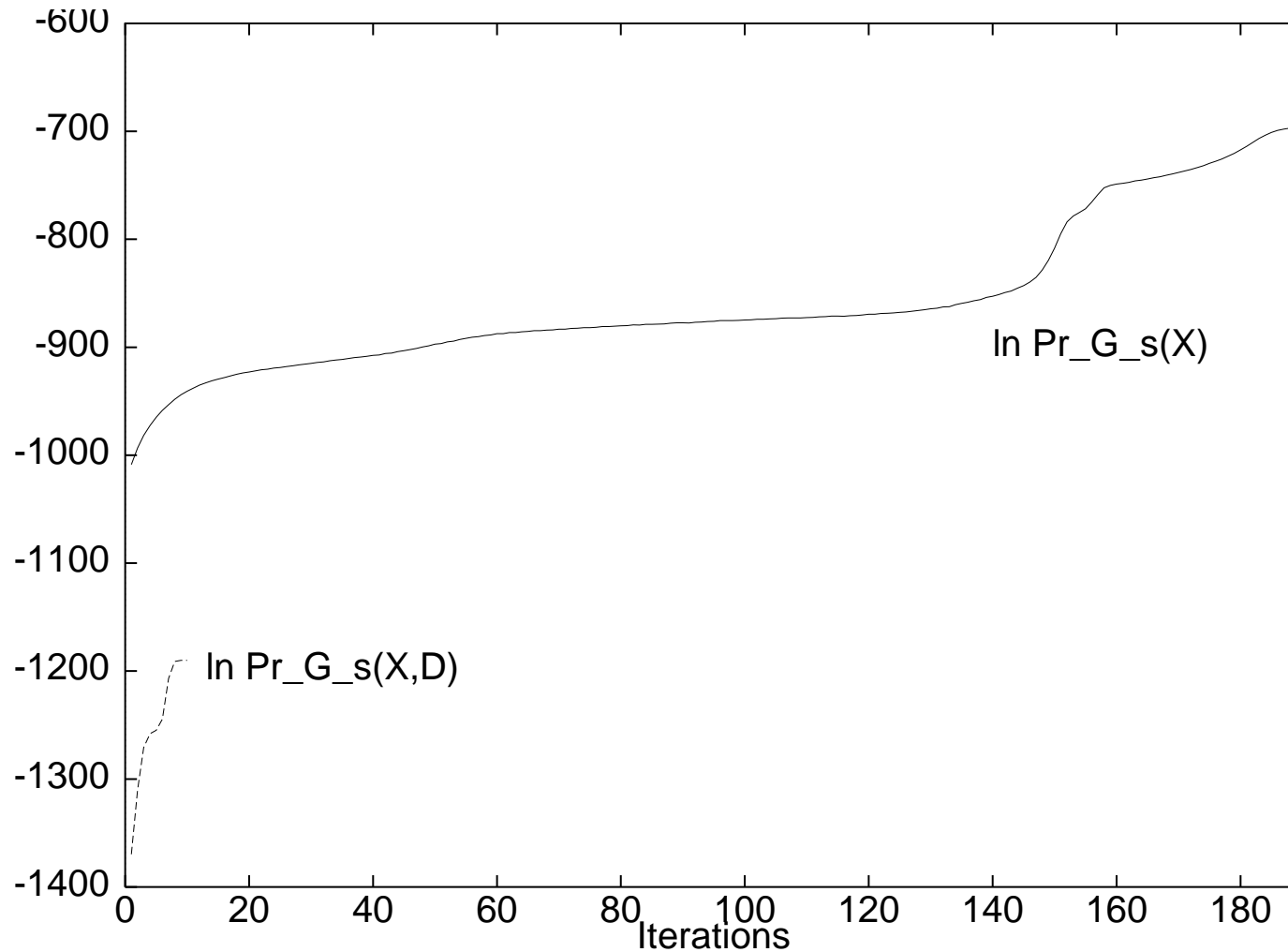
➤ Initial model to be estimated

- 5 non-terminals and 2 terminals \Rightarrow 130 rules
- Random probabilities attached to the rules

Algorithm	kld	Palindromes (%)	Non palindromes (%)
VS	6.00	1.9	98.1
IO	1.88	76.0	24.0

3.4 PROBABILISTIC PROPERTIES OF THE ESTIMATED PCFG

Functions optimized by IO algorithm and Viterbi algorithm



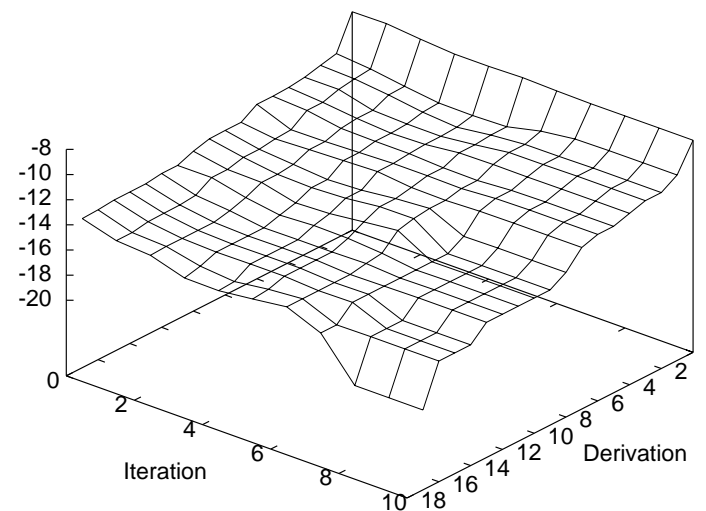
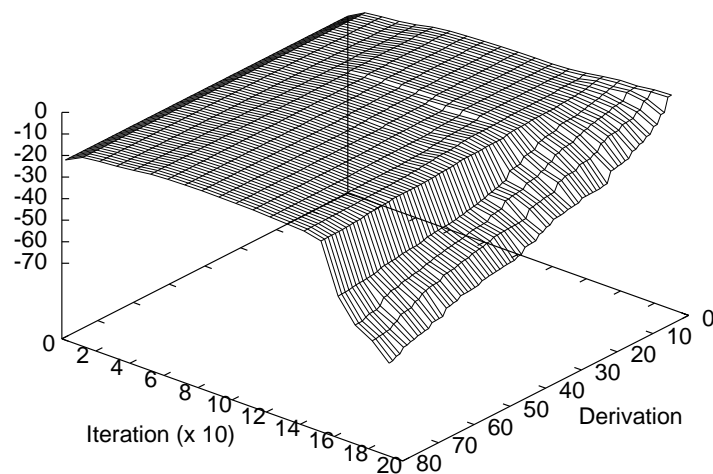
3.4 PROBABILISTIC PROPERTIES OF THE ESTIMATED PCFG

Accumulated probability mass

iteration	IO	
	apm	5bd
0	0.002	50.0%
340	0.724	97.4%

iteration	VS	
	apm	5bd
0	0.002	50.0%
10	0.024	95.7%

Evolution of the probability of a string



Combination of N-Grams and PCFG for LM [Benedi 05]

$$\Pr(w) = \Pr(w_1 \dots w_n) = \prod_{k=1}^n \Pr(w_k | w_1 \dots w_{k-1})$$

$$\Pr(w) = \prod_{k=1}^n \Pr(w_k | w_{k-n+1} \dots w_{k-1})$$

$$\Pr(w_k | w_1 \dots w_{k-1}) = \alpha \Pr_N(w_k | w_{k-n+1} \dots w_{k-1}) + (1 - \alpha) \Pr_{M_s}(w_k | w_1 \dots w_{k-1})$$

⇒ M_s : a PCFG G_c of categories (PoS tags) and a word-category distribution C_w

$$\Pr_{G_c, C_w}(w_k | w_1 \dots w_{k-1})$$

Model estimation

➤ Word N-Gram

➤ word-category distribution: C_w

$$C_w = \Pr(w|c) = \frac{N(w, c)}{\sum_{w'} N(w', c)}$$

➤ PCFG of categories: G_c

➤ (partially) unsupervised models: **IOb** and **VSb**

➤ supervised models: **“tree-bank grammars”**

Probability of the following word

$$\Pr_{G_c, C_w}(w_k | w_1 \dots w_{k-1}) = \frac{\Pr_{G_c, C_w}(w_1 \dots w_k \dots)}{\Pr_{G_c, C_w}(w_1 \dots w_{k-1} \dots)}$$

⇒ Prefix probability from G_c, C_w

$$\Pr_{G_c, C_w}(w_1 \dots w_k \dots)$$

3.5 USE OF PCFG FOR LM

➤ Definition:

$$e(A \ll i, l) = \Pr_{G_c, C_w}(A \xrightarrow{*} w_{i+1} \dots w_l \dots)$$

➤ Recursion:

$$e(A \ll i, i + j) = \sum_{B, C \in N} T(A \Rightarrow BC) \sum_{k=1}^{j-1} e(B \ll i, i + k) e(C \ll i + k, i + j)$$

➤ Base:

$$e(A \ll i, i + 1) = \sum_{\mathbf{c}} [p(A \rightarrow c) \Pr(\mathbf{w}_{i+1} | \mathbf{c}) + \sum_D T(A \Rightarrow D) p(D \rightarrow c) \Pr(\mathbf{w}_i | \mathbf{c})]$$

➤ Result: $\Pr_{G_c, C_w}(w_1 \dots w_k \dots) = e(S \ll 0, k).$

3.5 USE OF PCFG FOR LM

➤ Definition:

$$e(A \langle i, l \rangle) = \Pr_{G_c, C_w} (A \xrightarrow{*} w_{i+1} \dots w_l)$$

➤ Recursion:

$$e(A \langle i, i + j \rangle) = \sum_{B, C \in N} p(A \rightarrow BC) \sum_{k=1}^{j-1} e(B \langle i, i + k \rangle) e(C \langle i + k, i + j \rangle)$$

➤ Base:

$$e(A \langle i, i + 1 \rangle) = \sum_{\mathbf{c}} p(A \rightarrow c) \Pr(\mathbf{w}_{i+1} | \mathbf{c})$$

➤ Result:

$$\Pr_{G_c, C_w} (w_1 \dots w_n) = e(S \langle 0, n \rangle)$$

3.5 USE OF PCFG FOR LM

WSJ Experiments

➤ WSJ characteristics:

Data set	Directories	No. of senten.	No. of words
Training (full)	00-20	42,075	1,004,073
Training (≤ 50)	00-20	41,315 (98,2%)	959,390 (95,6%)
Tuning	21-22	3,371	80,156
Test	23-24	3,762	89,537

➤ **Vocabulary** (Training) 10,000 more frequent words

➤ **3-Gram model**: (*linear discounting*)

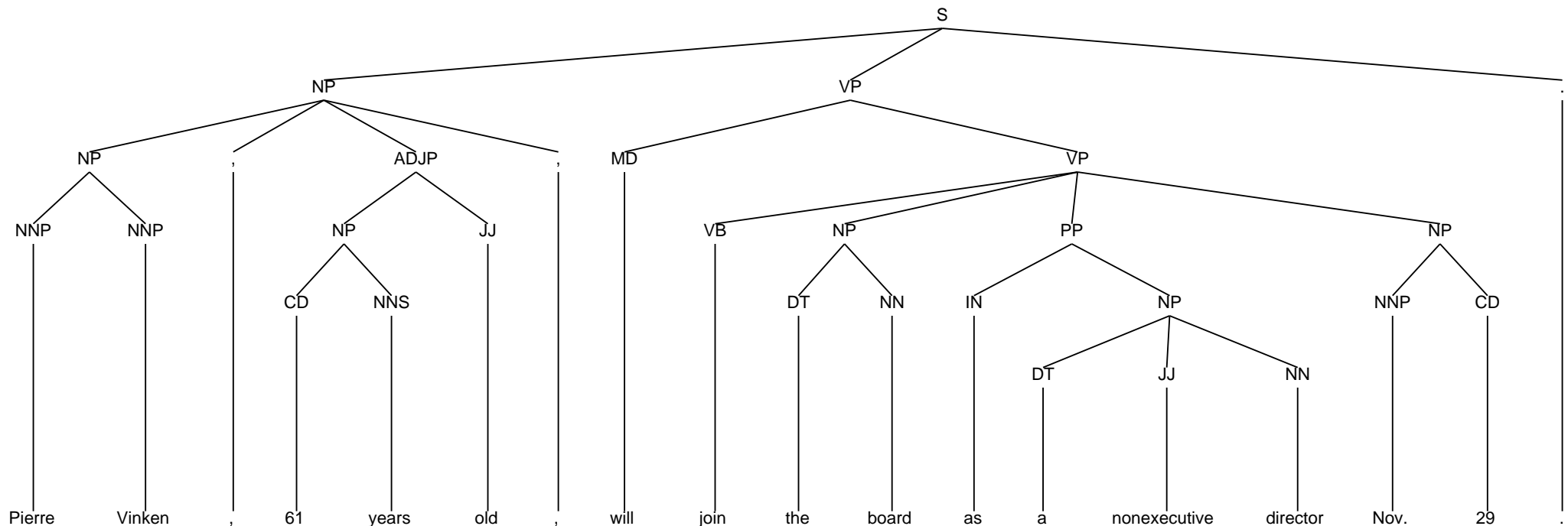
○ *Tuning set perplexity*: 160.3;

○ *Test set perplexity*: 167.3;

3.5 USE OF PCFG FOR LM

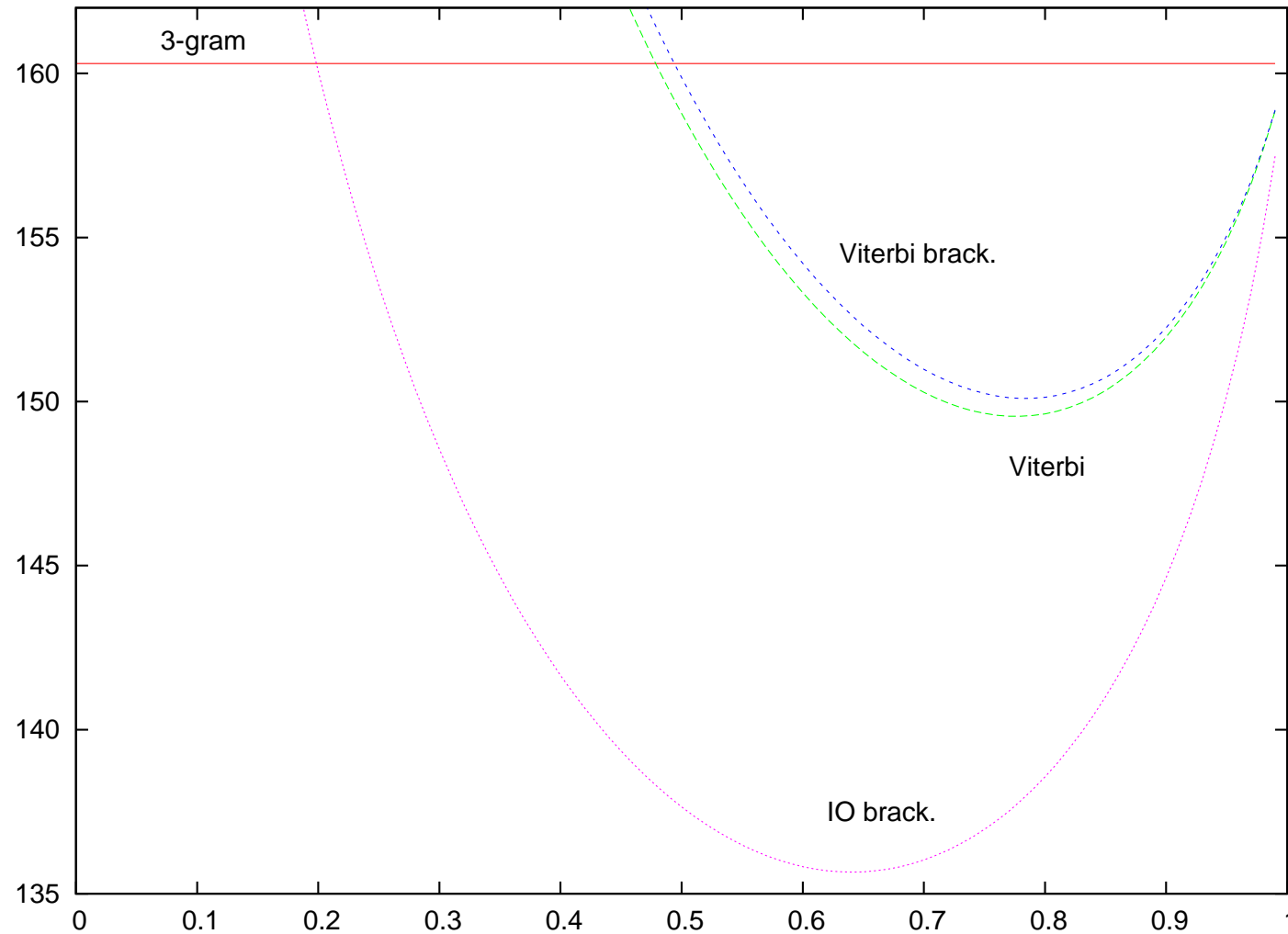
Example: *“Pierre Vinken, 61 years old, will join the board as a nonexecutive director Nov. 29.”*

((S (NP-SBJ (NP (NNP Pierre) (NNP Vinken)) (,,) (ADJP (NP (CD 61) (NNS years)) (JJ old)) (,,)) (VP (MD will) (VP (VB join) (NP (DT the) (NN board))) (PP-CLR (IN as) (NP (DT a) (JJ nonexecutive) (NN director)))) (NP-TMP (NNP Nov.) (CD 29)))) (..)))



3.5 USE OF PCFG FOR LM

Tuning set perplexity



3.5 USE OF PCFG FOR LM

Test set perplexity

Model	Perplexity		% improvement
	Trigram	Interpolated	
[Chelba 00]	167.1	148.9	10.9
[Roark 01]	167.0	137.3	17.8
IOb	167.3	142.3	14.9

WER

Model	Training Size	Vocabulary Size	LM Weight	WER
[Chelba 00]	20M	20K	16	13.0
[Roark 01]	1M	10K	15	15.1
Treebank trigram	1M	10K	5	16.6
No language model			0	16.8
Current model	1M	10K	6	16.0

Index

- 1 Introduction to the problem
 - 1.1 Introduction
 - 1.2 Syntactic models and Parsing
 - 1.3 Syntactic models for Language Modeling
 - 1.4 Syntactic models for Machine Translation
2. Preliminaries
 - 2.1 Notation and definitions
 - 2.2 Basic probabilistic properties of syntactic models
 - 2.3 CKY-based parsing algorithms
3. Probabilistic estimation of PCFG
 - 3.1 Introduction
 - 3.2 Inside-Outside algorithm
 - 3.3 Viterbi algorithm
 - 3.4 Probabilistic properties of the estimated PCFG
 - 3.5 Use of PCFG for LM
4. **Probabilistic estimation of SITG**
 - 4.1 Introduction
 - 4.2 Parsing with SITG
 - 4.3 Use of SITG for MT
5. Advanced topics
 - 5.1 On-line learning of syntactic models
 - 5.2 Active learning of syntactic models
 - 5.3 Interactive-predictive parsing:
a framework for active learning

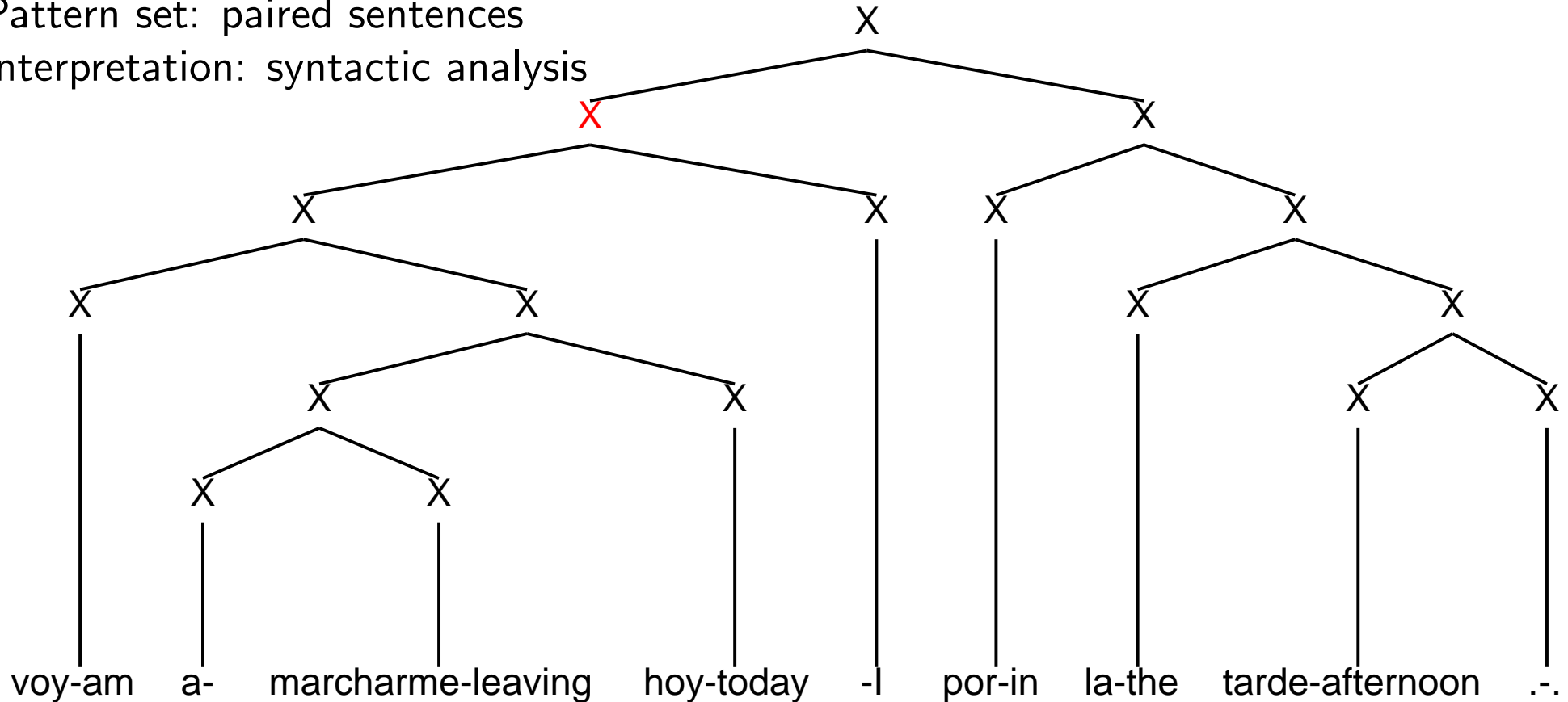
4.1 INTRODUCTION

Stochastic inversion transduction grammars [Wu 97, Maryanski 79]

- Primitives: two alphabets (words, punctuation symbols, . . .)
- Object representation: two paired written sentences

“voy a marcharme hoy por la tarde” \iff “I am leaving today in the afternoon”

- Pattern set: paired sentences
- Interpretation: syntactic analysis



4.1 INTRODUCTION

➤ **ITG:** $G = (N, W_1, W_2, R, S)$

R is a finite set of straight orientation rules $A \rightarrow [a_1 a_2 \dots a_r]$ and inverted orientation rules $A \rightarrow \langle a_1 a_2 \dots a_r \rangle$, $a_i \in N \cup X$ and $X = (W_1 \cup \{\epsilon\}) \times (W_2 \cup \{\epsilon\})$

Theorem. For any ITG G , there exists an equivalent ITG G' in which every production takes one of the following forms:

$$\begin{array}{lll} S \rightarrow \epsilon/\epsilon & A \rightarrow x/\epsilon & A \rightarrow [BC] \\ S \rightarrow x/y & A \rightarrow \epsilon/y & A \rightarrow \langle BC \rangle \end{array}$$

➤ **SITG:** $G_s = (G, p)$ where:

➤ G is an ITG

➤ p is a function that attaches a probability to each rule:

$$p : R \rightarrow]0, 1] \quad \sum_{1 \leq j \leq n_i} p(A_i \rightarrow \alpha_j) = 1, \quad \forall A_i \in N$$

4.1 INTRODUCTION

Stochastic derivation for SITG

Given a sequence of stochastic events:

$$(S, S) = (\alpha_0, \beta_0) \xrightarrow{r_1} (\alpha_1, \beta_1) \xrightarrow{r_2} (\alpha_2, \beta_2) \cdots (\alpha_{m-1}, \beta_{m-1}) \xrightarrow{r_m} (\alpha_m, \beta_m) = (x, y)$$

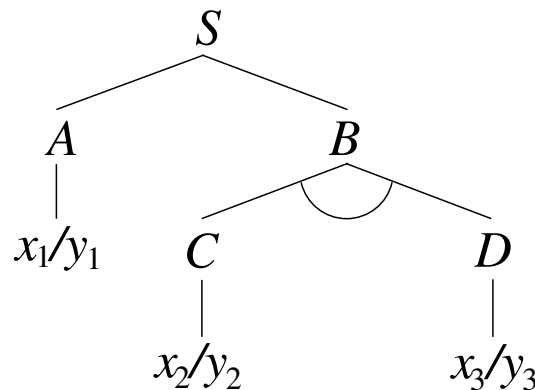
Probability of (x, y) being generated by $G_s = (G, p)$ from the rule sequence

$d_x = (r_1, \dots, r_m)$, is:

$$P_{G_s}((x, y), d_x) = \prod_{j=1 \dots m} p(r_j)$$

Example

$$\begin{aligned} S &\rightarrow [AB] \\ A &\rightarrow x_1/y_1 \\ B &\rightarrow \langle CD \rangle \\ C &\rightarrow x_2/y_2 \\ D &\rightarrow x_3/y_3 \end{aligned}$$



$$(S, S) \Rightarrow (AB, AB) \Rightarrow (x_1B, y_1B) \Rightarrow (x_1CD, y_1DC) \Rightarrow (x_1x_2D, y_1Dy_2) \Rightarrow (x_1x_2x_3, y_1y_3y_2)$$

4.1 INTRODUCTION

Probability of a string pair

$$\Pr_{G_s}(x, y) = \sum_{d_x \in D_x} \Pr_{G_s}((x, y), d_x)$$

Probability of the best derivation

$$\widehat{\Pr}_{G_s}(x, y) = \max_{d_x \in D_x} \Pr_{G_s}((x, y), d_x)$$

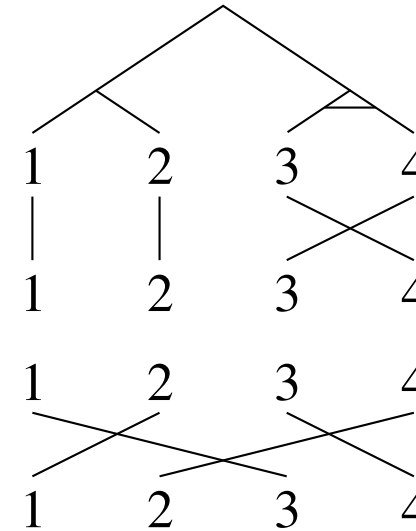
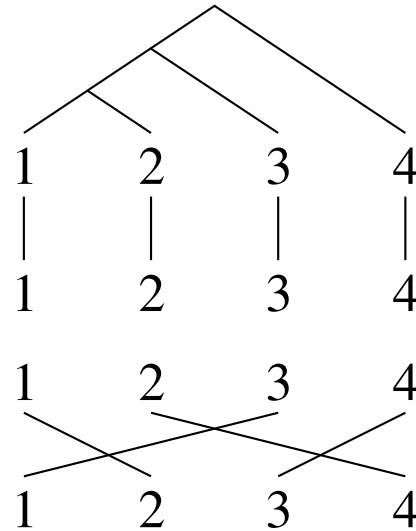
Language generated by a SITG

$$L(G_s) = \{(x, y) \mid \Pr_{G_s}(x, y) > 0\}$$

4.1 INTRODUCTION

Expressiveness of ITGs

YES



NO

r	ITG	all matchings	ratio
1	1	1	1.000
2	2	2	1.000
3	6	6	1.000
4	22	24	0.917
5	90	120	0.750

r	ITG	all matchings	ratio
6	394	720	0.547
7	1,806	5,040	0.358
8	8,558	40,320	0.212
9	41,586	362,880	0.115
10	206,098	3,628,800	0.057

4.2 PARSING WITH SITG

➤ Parsing:

- Inside algorithm
- Viterbi algorithm

➤ Learning:

- Structure learning
- Probabilistic estimation: Inside-outside estimation
Viterbi-based estimation

➤ Translation:

- Adapted Cooke-Kasami-Younger parser algorithm

Viterbi algorithm [Wu 97, Gascó 10b]

➤ Given $(x, y) \in (W_1^*, W_2^*)$ and $A \in N$

$$\delta_{i,j,k,l}(A) = \widehat{\text{Pr}}(A \xrightarrow{*} x_{i+1} \cdots x_j / y_{k+1} \cdots y_l)$$

➤ Initialization

$$\delta_{i-1,i,k-1,k}(A) = p(A \rightarrow x_i / y_k) \quad 1 \leq i \leq |x|, 1 \leq k \leq |y|$$

$$\delta_{i-1,i,k,k}(A) = p(A \rightarrow x_i / \epsilon) \quad 1 \leq i \leq |x|, 0 \leq k \leq |y|$$

$$\delta_{i,i,k-1,k}(A) = p(A \rightarrow \epsilon / y_k) \quad 0 \leq i \leq |x|, 1 \leq k \leq |y|$$

4.2 PARSING WITH SITG

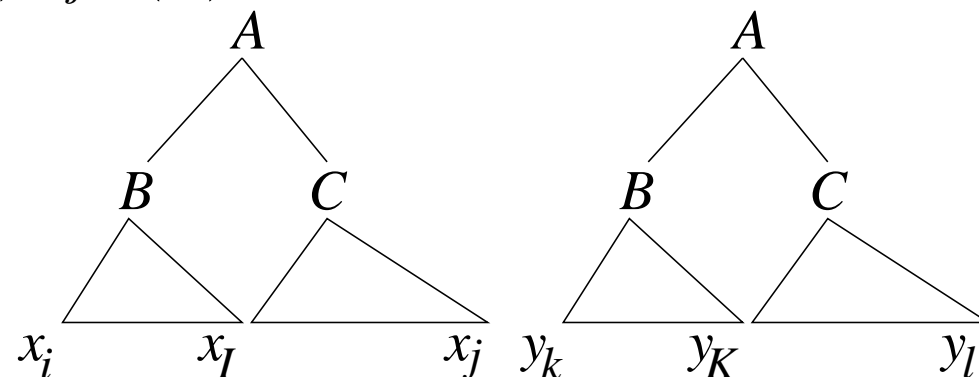
- Recursion. For all $A \in N$, and i, j, k, l such that $0 \leq i < j \leq |x|$, $0 \leq k < l \leq |y|$ and $j - i + l - k \geq 2$:

$$\delta_{ijkl}(A) = \max(\delta_{ijkl}^{\square}(A), \delta_{ijkl}^{\langle \rangle}(A))$$

$$\delta_{ijkl}^{\square}(A) = \max_{B, C \in N} p(A \rightarrow [BC]) \delta_{iIkK}(B) \delta_{IjKl}(C)$$

$$i \leq I \leq j, k \leq K \leq l$$

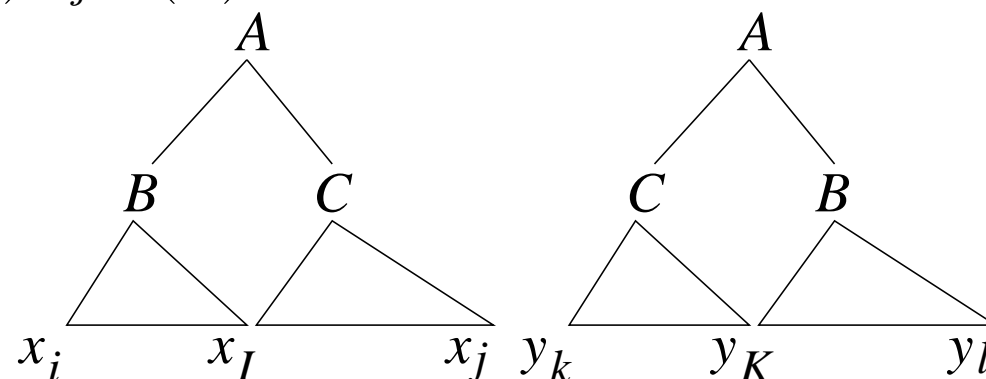
$$((j-I) + (l-K)) \times ((l-i) + (K-k)) \neq 0$$



$$\delta_{ijkl}^{\langle \rangle}(A) = \max_{B, C \in N} p(A \rightarrow \langle BC \rangle) \delta_{iIKl}(B) \delta_{IjkK}(C)$$

$$i \leq I \leq j, k \leq K \leq l$$

$$((j-I) + (K-k)) \times ((I-i) + (I-K)) \neq 0$$



Inside algorithm [Gascó 10c]

- Given $(x, y) \in (W_1^*, W_2^*)$ and $A \in N$

$$\mathcal{E}_{i,i+j,k,k+l}[A] = \Pr(A \xRightarrow{*} x_{i+1} \cdots x_{i+j} / y_{k+1} \cdots y_{k+l})$$

Substitute max by \sum in Viterbi algorithm

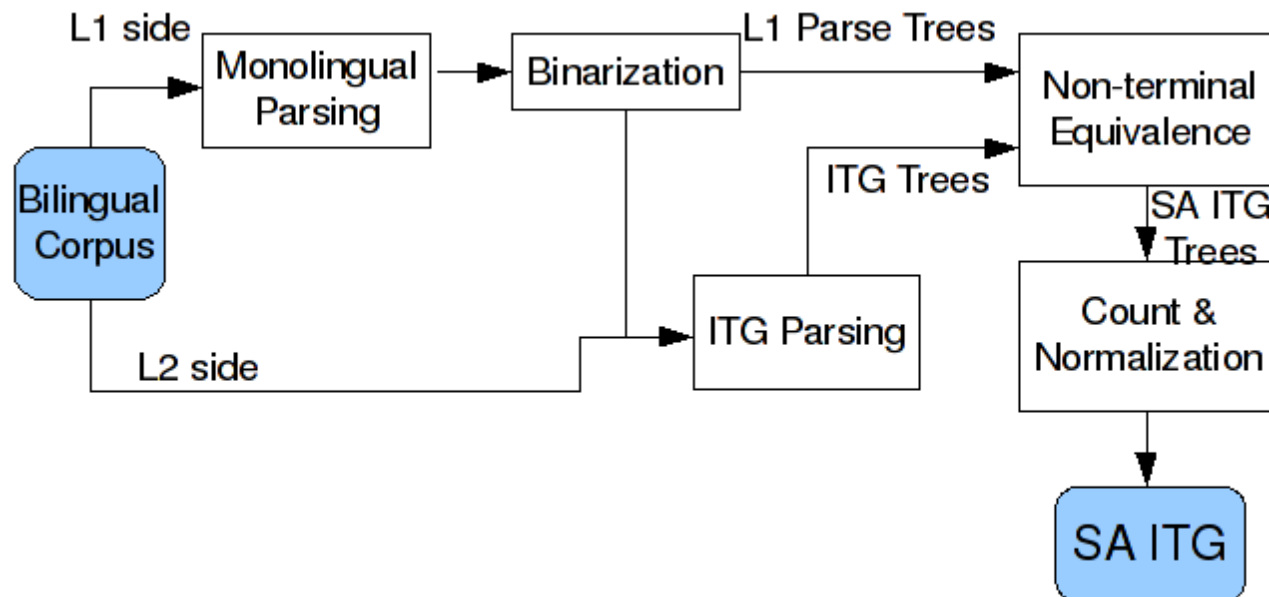
- **Theorem.** If the *inside* algorithm is applied to the substring pair $(x_{i+1} \cdots x_{i+j}, y_{k+1} \cdots y_{k+l})$ with a SITG \mathcal{G} , then the probabilistic parse matrix \mathcal{E} collects correctly the probability of this substring pair.
- **Corollary.** The probability of the pair string $(x_1 \cdots x_{|x|}, y_1 \cdots y_{|y|})$ can be computed by means of the probabilistic parse matrix \mathcal{E} as:

$$\Pr(S \xRightarrow{+} x_1 \cdots x_{|x|} / y_1 \cdots y_{|y|}) = \mathcal{E}_{0,|x|,0,|y|}[S]$$

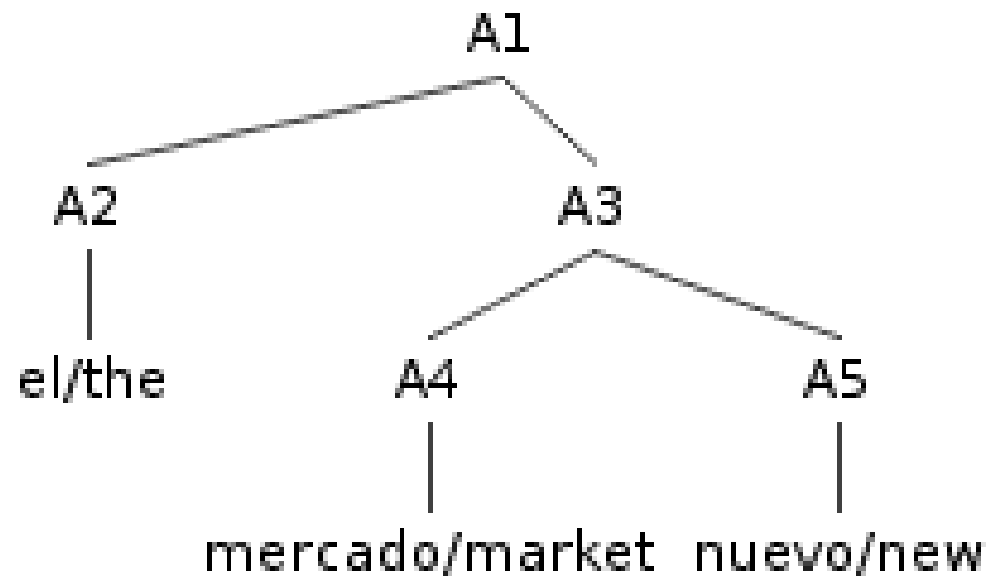
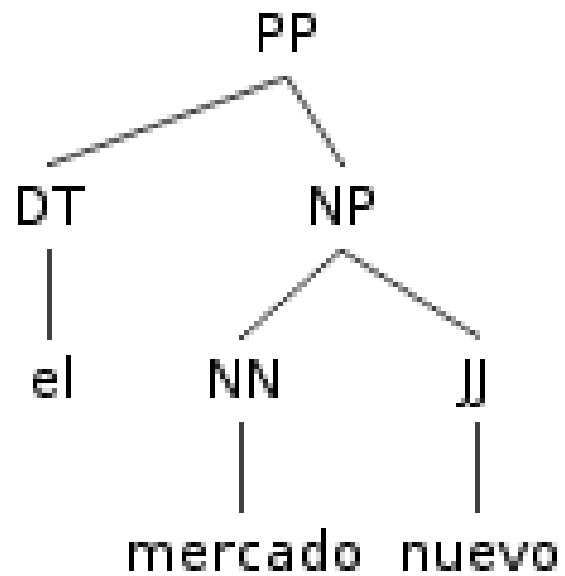
4.3 USE OF SITG FOR MT

[Gascó 10a]

1. To create an initial SITG
2. To estimate the probabilities
3. To attach linguistic information to the non-terminal symbols



4.3 USE OF SITG FOR MT



4.3 USE OF SITG FOR MT

- IWSLT 2008 (Chinese-English BTEC)
- Standard tools: GIZA++, ZMERT
- Stanford parser for Chinese
- Baseline: Moses, 5-gram

Corpus Set	Statistic	Chinese	English
Training	Sentences	42,655	
	Words	330,163	380,431
	Voc. Size	8,773	8,387
DevSet	Sentences	489	
	Words	3,169	3,861
	OOV Words	111	115
Test	Sentences	507	
	Words	3,357	-
	OOV Words	97	-

System	%BLEU
Baseline PBT	41.1
Initial ITG	41.2
Re-estimated ITG	41.8
Source SAITG	42.9
Target SAITG	43.0

Index

- 1 Introduction to the problem
 - 1.1 Introduction
 - 1.2 Syntactic models and Parsing
 - 1.3 Syntactic models for Language Modeling
 - 1.4 Syntactic models for Machine Translation
2. Preliminaries
 - 2.1 Notation and definitions
 - 2.2 Basic probabilistic properties of syntactic models
 - 2.3 CKY-based parsing algorithms
3. Probabilistic estimation of PCFG
 - 3.1 Introduction
 - 3.2 Inside-Outside algorithm
 - 3.3 Viterbi algorithm
 - 3.4 Probabilistic properties of the estimated PCFG
 - 3.5 Use of PCFG for LM
4. Probabilistic estimation of SITG
 - 4.1 Introduction
 - 4.2 Parsing with SITG
 - 4.3 Use of SITG for MT
5. **Advanced topics**
 - 5.1 On-line learning of syntactic models
 - 5.2 Active learning of syntactic models
 - 5.3 Interactive-predictive parsing:
a framework for active learning

Problem definition [Liang 09]:

- Probabilistic model: $p(\mathbf{x}, \mathbf{z}; \theta)$

Input: \mathbf{x} (a sentence)

Hidden output: \mathbf{z} (a parse tree)

Parameters: θ (rule probabilities)

- Given a set of unlabeled example $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}$, maximize the marginal log-likelihood:

$$l(\theta) = \sum_{i=1}^n \log p(\mathbf{x}^{(i)}; \theta)$$

- Evaluation of the trained model $\hat{\theta}$: accuracy

$$\text{true output } \mathbf{z}^{(i)} \leftrightarrow \arg \max_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}^{(i)}; \theta)$$

- Training algorithm: **EM algorithm** [Dempster 77, Neal 98, Cappé 09]

5.1 ON-LINE LEARNING OF SYNTACTIC MODELS

EM algorithm [Liang 09]:

Batch EM

```
 $\mu \leftarrow$  initialization  
for each iteration  $t = 1, \dots, T$ :  
   $\mu' \leftarrow 0$   
  for each example  $i = 1, \dots, n$ :  
     $s'_i \leftarrow \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}^{(i)}; \theta(\mu)) \phi(\mathbf{x}^{(i)}, \mathbf{z})$   
     $\mu' \leftarrow \mu' + s'_i$   
 $\mu \leftarrow \mu'$ 
```

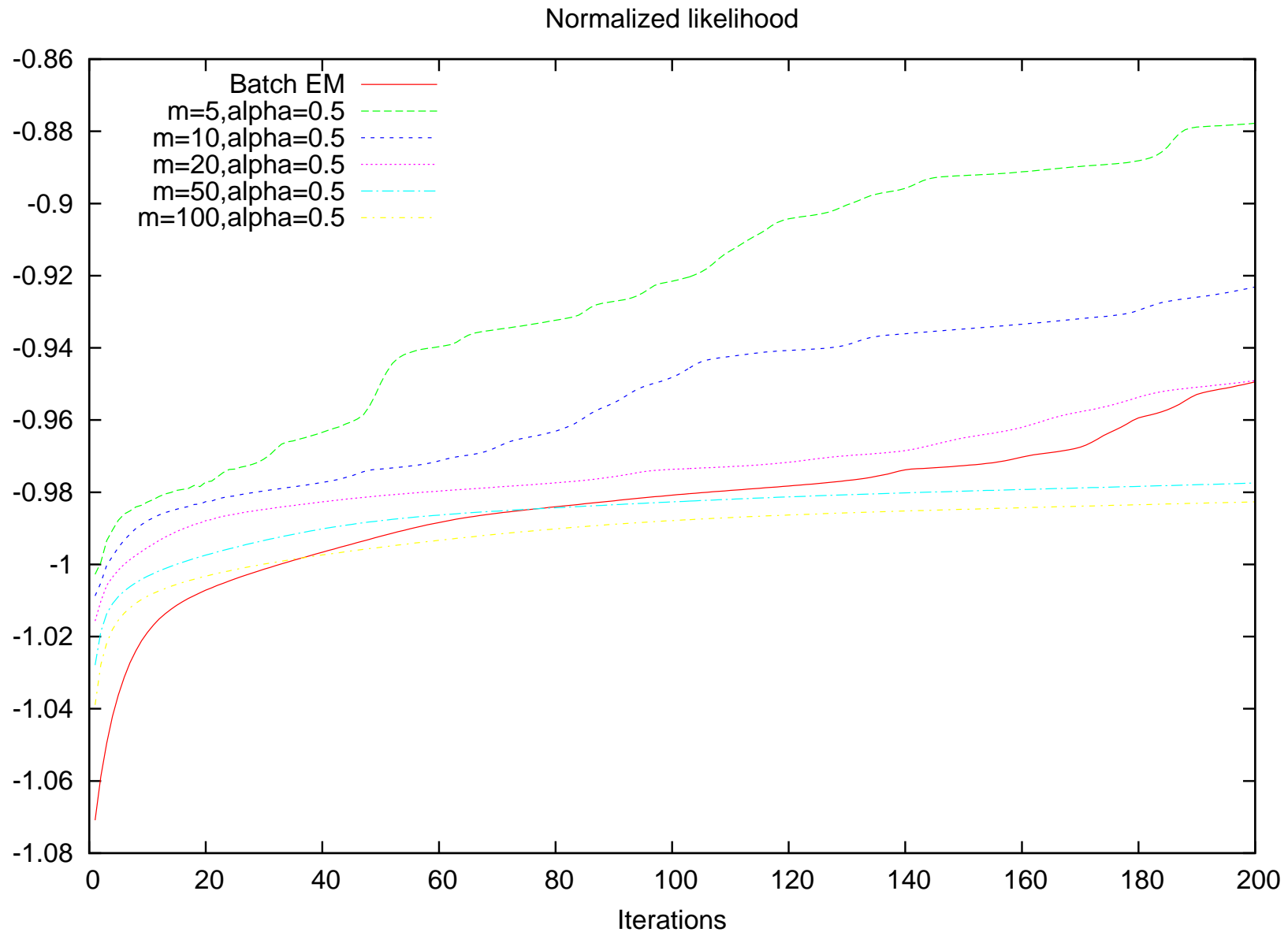
Stepwise EM

```
 $\mu \leftarrow$ ;  $k = 0$  initialization  
for each iteration  $t = 1, \dots, T$ :  
  for each example  $i = 1, \dots, n$  in  
  random order:  
     $s'_i \leftarrow \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}^{(i)}; \theta(\mu)) \phi(\mathbf{x}^{(i)}, \mathbf{z})$   
     $\mu \leftarrow (1 - \eta_k)\mu + \eta_k s'_i$   
     $k \leftarrow k + 1$ 
```

- $\phi(\mathbf{x}, \mathbf{z})$: mapping from a labelled example (\mathbf{x}, \mathbf{z}) to a vector of sufficient statistics (μ)
- $\theta(\mu)$: maximum likelihood estimate
- Stepwise EM: convergence is guaranteed if $\sum_{k=0}^{\infty} \eta_k = \infty$ and $\sum_{k=0}^{\infty} \eta_k^2 < \infty$
 - $\eta_k = (k + 2)^{-\alpha}$ with $0.5 < \alpha \leq 1$
 - Approach: take m examples at once

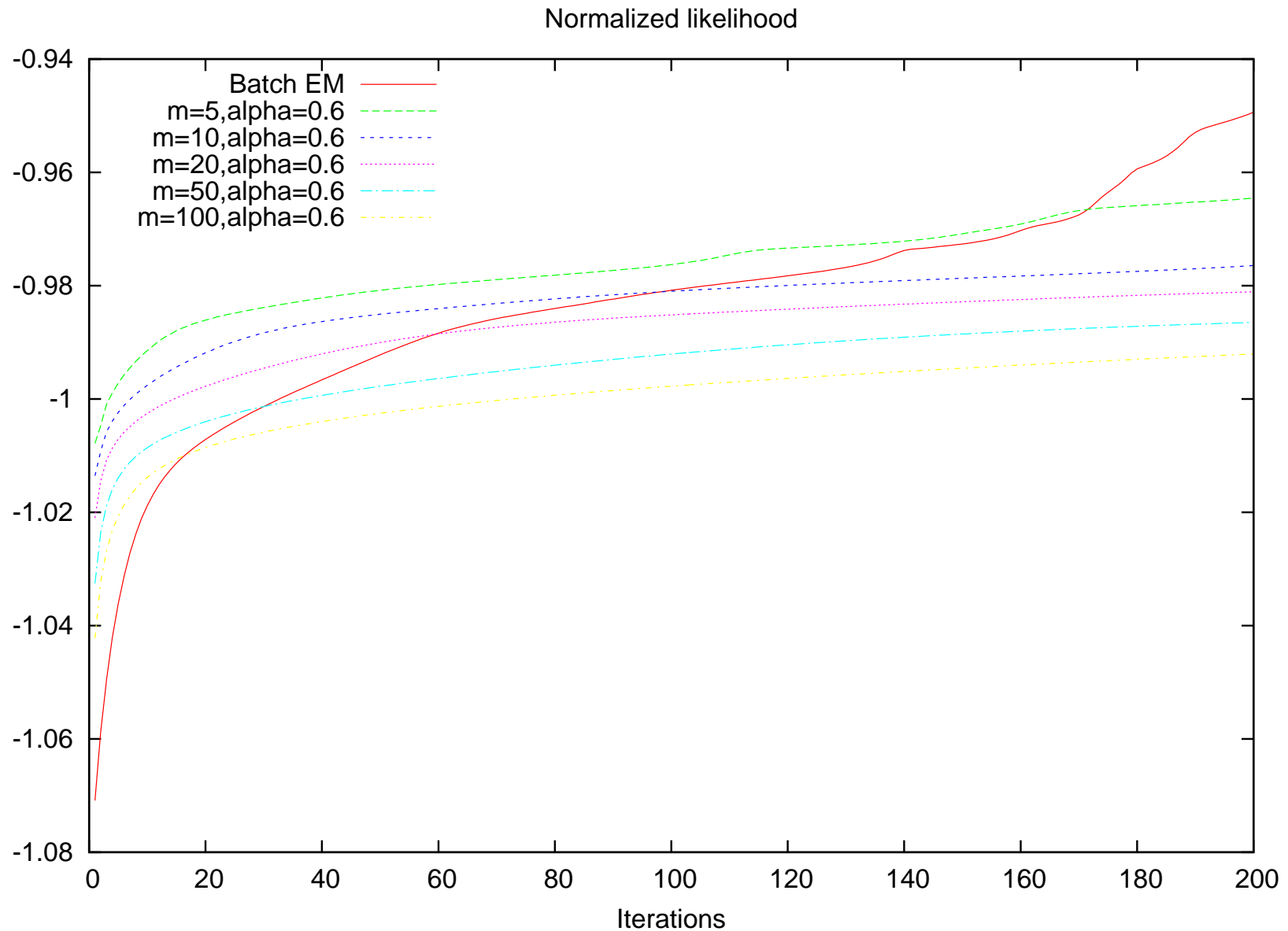
5.1 ON-LINE LEARNING OF SYNTACTIC MODELS

Palindrome language (15 random initializations, $\alpha = 0.5$)



5.1 ON-LINE LEARNING OF SYNTACTIC MODELS

Palindrome language (15 random initializations, $\alpha = 0.6$)



5.1 ON-LINE LEARNING OF SYNTACTIC MODELS

Palindrome language

➤ Batch EM

- normalized likelihood in average: -0.949 ± 0.052
- % of palindromes generated in average by the estimated grammar (50K sentences): $21.9\% \pm 10.1$

➤ Online EM

$\alpha \backslash m$	5	10	20	50	100
0.5	-0.878 ± 0.080	-0.923 ± 0.071	-0.949 ± 0.052	-0.977 ± 0.007	-0.983 ± 0.007
0.6	-0.965 ± 0.031	-0.976 ± 0.007	-0.981 ± 0.007	-0.986 ± 0.008	-0.992 ± 0.007
0.7	-0.984 ± 0.008	-0.987 ± 0.008	-0.992 ± 0.007	-0.998 ± 0.006	-1.002 ± 0.004
0.8	-0.997 ± 0.006	-1.000 ± 0.005	-1.003 ± 0.004	-1.006 ± 0.003	-1.009 ± 0.002
0.9	-1.007 ± 0.003	-1.009 ± 0.002	-1.011 ± 0.002	-1.014 ± 0.002	-1.016 ± 0.002
1.0	-1.016 ± 0.002	-1.018 ± 0.003	-1.021 ± 0.004	-1.023 ± 0.004	-1.025 ± 0.005

$\alpha \backslash m$	5	10	20	50	100
0.5	35.1 ± 21.0	29.5 ± 18.3	21.8 ± 10.2	16.6 ± 2.3	15.4 ± 2.4
0.6	20.3 ± 11.3	16.8 ± 2.4	15.8 ± 2.4	14.9 ± 2.8	14.6 ± 3.2
0.7	15.1 ± 2.7	14.9 ± 2.9	14.6 ± 3.3	14.9 ± 3.4	15.4 ± 3.1
0.8	15.0 ± 3.2	15.3 ± 3.2	15.7 ± 2.9	16.0 ± 2.4	16.2 ± 2.0
0.9	16.4 ± 2.1	16.3 ± 1.9	16.4 ± 1.7	16.5 ± 1.4	16.4 ± 1.2
1.0	16.6 ± 1.3	16.2 ± 1.2	16.0 ± 1.0	15.7 ± 1.0	15.5 ± 1.0

5.1 ON-LINE LEARNING OF SYNTACTIC MODELS

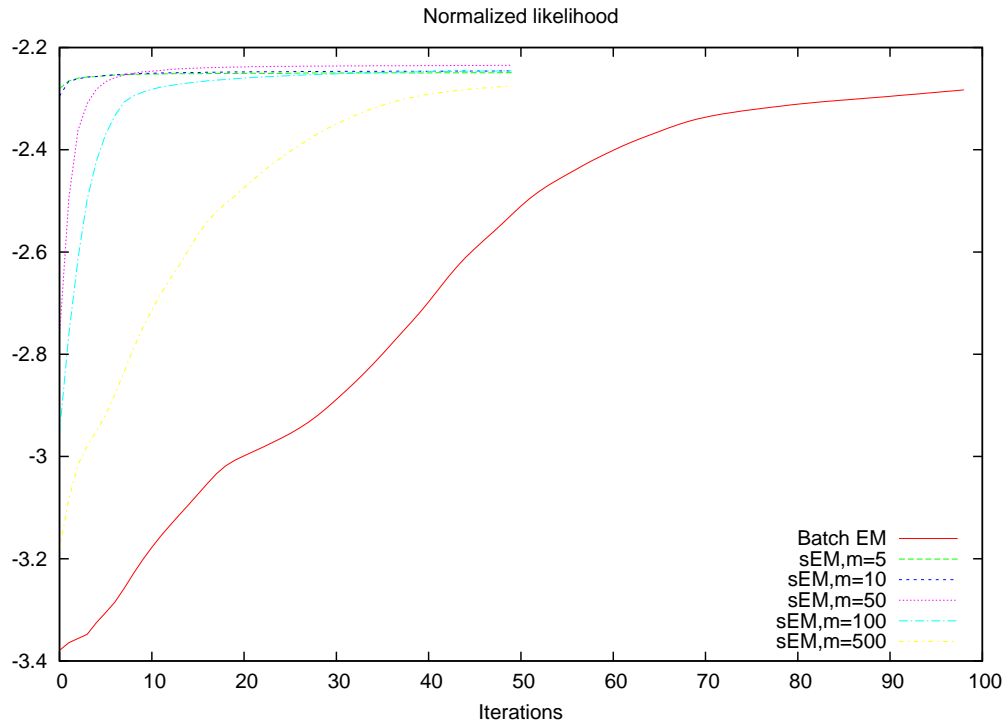
WSJ - UPenn Treebank

- Until level of preterminals: 45 PoS
- Initial PCFG: ergodic model with 45 terminals and 35 non-terminals (44,450 rules)
- Training A: length sentence ≤ 20 words
 - sections 02-21
 - 16,667 sentences
 - 229,020 running words
- Training B: length sentence ≤ 25 words
 - sections 02-21
 - 23,597 sentences
 - 388,234 running words
- Test: section 23
 - 2,416 sentences
 - 56,684 running words

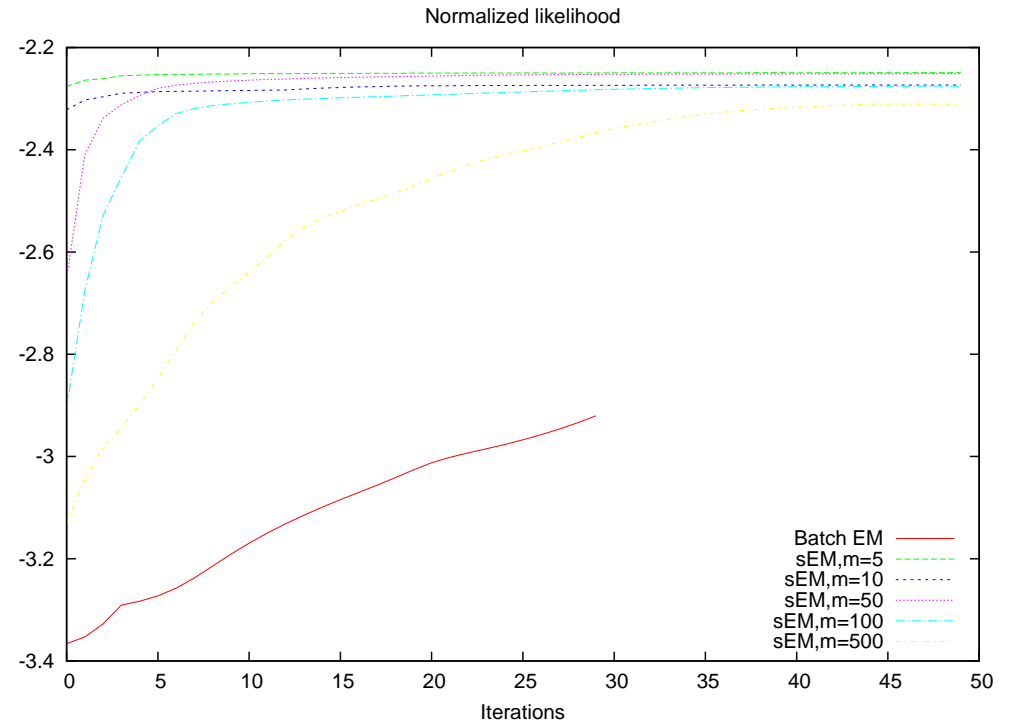
5.1 ON-LINE LEARNING OF SYNTACTIC MODELS

WSJ - UPenn Treebank

($\alpha = 0.5, \leq 20$ words)



($\alpha = 0.5, \leq 25$ words)



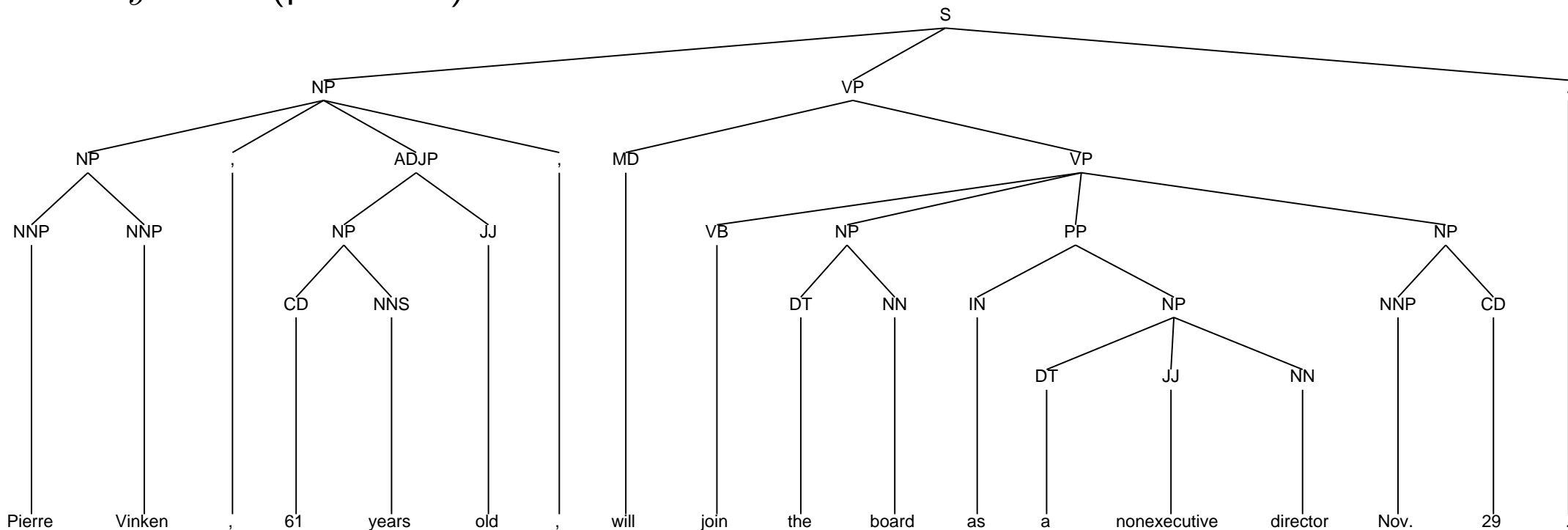
Test Bracketing F-Measure (measured with *evalb*: <http://nlp.cs.nyu.edu/evalb/>)

	BatchEM	$m = 5$	$m = 10$	$m = 50$	$m = 100$	$m = 500$
≤ 20	31	35	34	34	29	28
≤ 25		34	33	31	31	30

5.2 ACTIVE LEARNING OF SYNTACTIC MODELS

Problem definition:

- Supervised learning: (x, y)
 - x : input data (sentence)
 - y : label (parse tree)



- Problem: to annotate data is slow and expensive
- **Active learning: to annotate just the necessary data**

5.2 ACTIVE LEARNING OF SYNTACTIC MODELS

Pool-based active learning [Settles 08, Settles 10]:

```
Given: Labeled set  $\mathcal{L}$ , unlabeled pool  $\mathcal{U}$ ,  
query strategy  $\phi(\cdot)$ , query batch size  $B$   
repeat  
  // learn a model using the current  $\mathcal{L}$   
   $\theta = \text{train}(\mathcal{L})$   
  for  $b = 1$  to  $B$  do  
    // query the most informative instance  
     $\mathbf{x}_b^* = \arg \max_{\mathbf{x} \in \mathcal{U}} \phi(\mathbf{x})$   
    // move the labeled query from  $\mathcal{U}$  to  $\mathcal{L}$   
     $\mathcal{L} = \mathcal{L} \cup \langle \mathbf{x}_b^*, \text{label}(\mathbf{x}_b^*) \rangle$   
     $\mathcal{U} = \mathcal{U} - \mathbf{x}_b^*$   
  end  
until some stopping criterion
```

➤ Similar scheme for parsing in [Hwa 04]

5.2 ACTIVE LEARNING OF SYNTACTIC MODELS

Query strategies:

➤ **Uncertainty sampling:** to query the instance that is most uncertainty how to label

➤ Sequence entropy:

$$\phi^{SE}(\mathbf{x}) = - \sum_{\hat{\mathbf{y}}} P(\hat{\mathbf{y}}|\mathbf{x}; \theta) \log P(\hat{\mathbf{y}}|\mathbf{x}; \theta)$$

➤ Approach: N -best Sequence entropy:

$$\phi^{NSE}(\mathbf{x}) = - \sum_{\hat{\mathbf{y}} \in \mathcal{N}} P(\hat{\mathbf{y}}|\mathbf{x}; \theta) \log P(\hat{\mathbf{y}}|\mathbf{x}; \theta)$$

➤ **Query-By-Committee:** to query the instance for which a committee of models is in most disagreement about how to label

➤ **Information density:** to query the instance that is the most “informative” in average

$$\phi^{ID}(\mathbf{x}) = \phi^{NSE}(\mathbf{x}) \times \left(\frac{1}{U} \sum_{u=1}^U \text{sim}(\mathbf{x}, \mathbf{x}^{(u)}) \right)^{\beta}$$

5.2 ACTIVE LEARNING OF SYNTACTIC MODELS

Query strategies for parsing [Hwa 04]:

- Problem space:

- Based on novelty and frequencies of word pair co-occurrences
- Based on sentence length: f_{len}

- Performance of the hypothesis:

- Error-driven function:

$$f_{\text{err}}(\mathbf{w}, G) = 1 - P(\hat{d}_{\mathbf{w}} | \mathbf{w}, G)$$

- Normalized tree entropy (similar to $\phi^{SE}(\mathbf{x})$): f_{unc}

5.2 ACTIVE LEARNING OF SYNTACTIC MODELS

Experiments on WSJ UPenn Treebank reported in [Hwa 04]:

- Collins' model 2 parser
- Learning algorithm: statistics directly over the treebank
- Data:
 - Training: sections 02-21
 - Test: section 23
- Initial model trained on 500 sentences
- Batch size: 100
- Parsing performance: F score

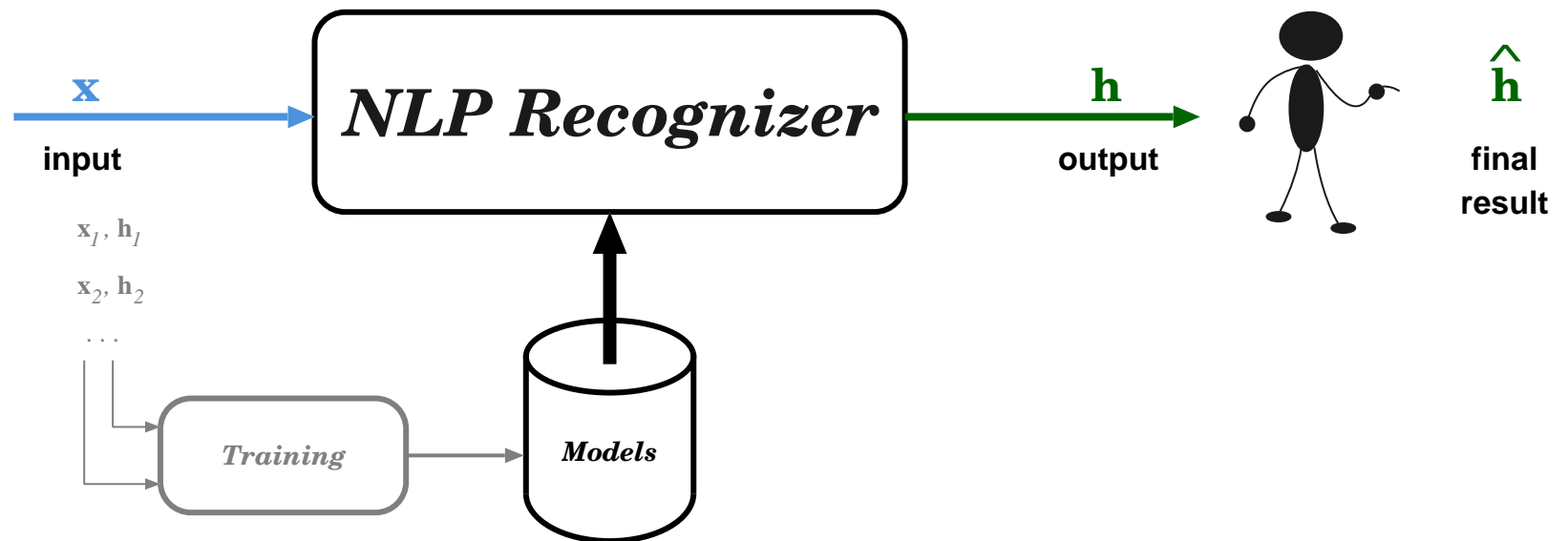
Number of labelled samples at the test performance level of 88%:

	f_{ran}	f_{len}	f_{err}	f_{unc}
# sentences	30,500	-	20,500 (33%)	17,500 (43%)
# constituents	695,000	625,000 (10%)	577,000 (17%)	505,000 (27%)

5.3 IPP: A FRAMEWORK FOR ACTIVE LEARNING

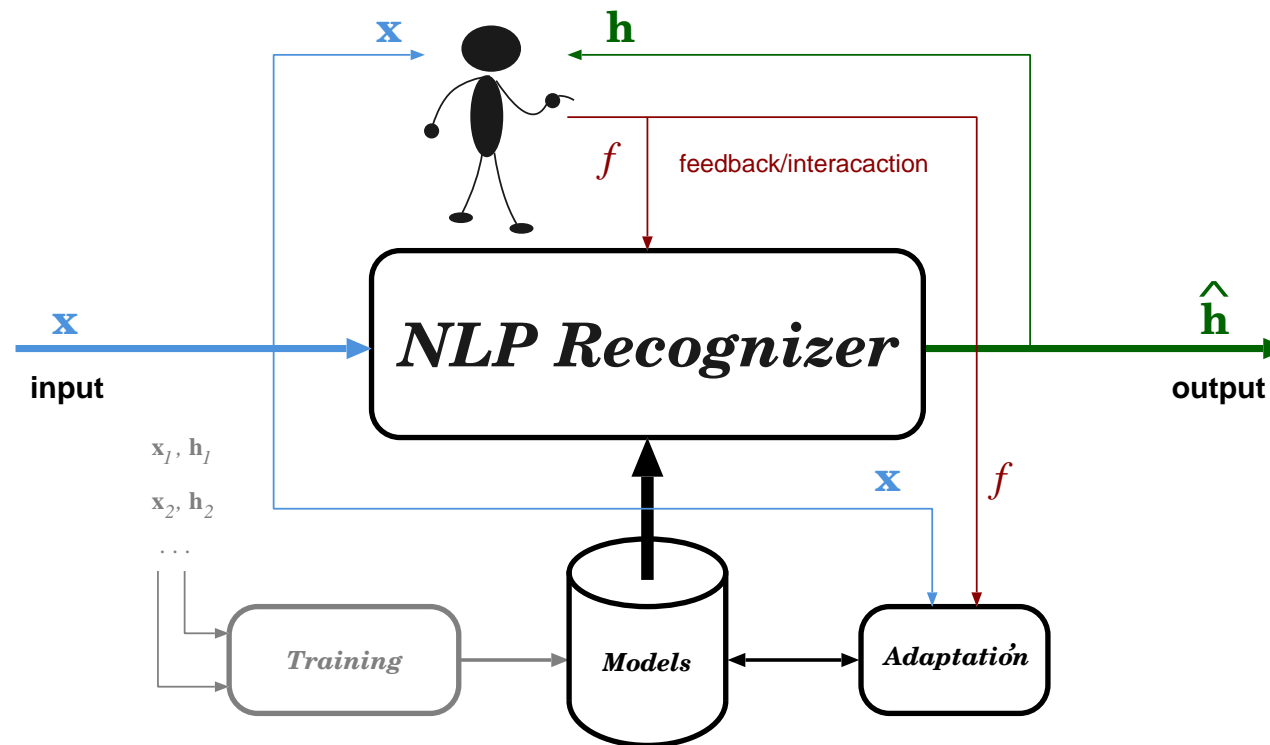
Problem definition [Sánchez 09, Sánchez 10]: Annotation parse tree is expensive and requires skilled expert humans

- Classical two-step approach:
 - 1 Apply an automatic system
 - 2 Manually validate/correct the output



5.3 IPP: A FRAMEWORK FOR ACTIVE LEARNING

- Interactive Predictive approach:
 - Formally integrate the user into the recognition process
 - The system reacts to user feedback



- New opportunities:
 - Feedback information can be used to create efficient interactive systems
 - Each interaction step yields *ground-truth data*, which allows building *active learning systems*

5.3 IPP: A FRAMEWORK FOR ACTIVE LEARNING

Classical parsing

$$\hat{t} = \arg \max_{t \in \mathcal{T}} p_G(t|x)$$

$x \rightarrow$ input sentence

$G \rightarrow$ mode (e.g. PCFG)

$\mathcal{T} \rightarrow$ set of all possible trees for x with G

$\hat{t} \rightarrow$ obtained parse tree

Interactive predictive parsing

$$\hat{t} = \arg \max_{t \in \mathcal{T}: t_p \in t} p_G(t|x, t_p)$$

The tree prefix t_p is:

- the corrected constituent, plus
- all its ancestors, plus
- all the constituents to its left

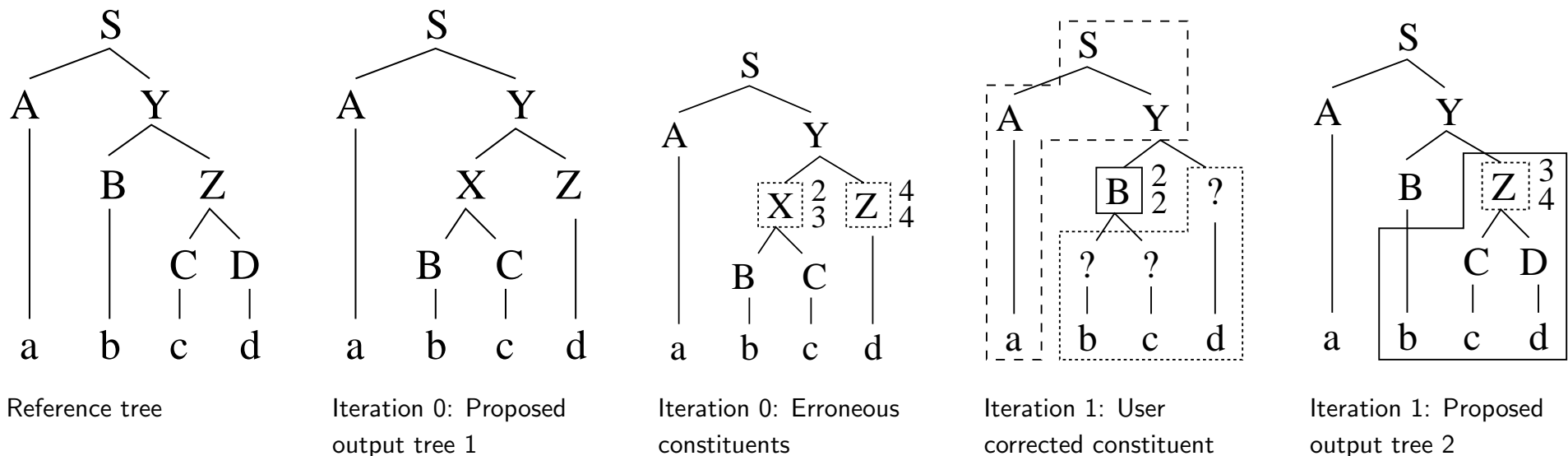
$$t_p(c'_{ij}) = \{c^B_{mn} : m \leq i, n \geq j, \text{depth}(c^B_{mn}) \leq \text{depth}(c'_{ij})\} \cup \{c^D_{pq} : p \geq 1, q < i\}$$

5.3 IPP: A FRAMEWORK FOR ACTIVE LEARNING

IPP parsing

1. The system propose a parse tree \hat{t}
2. The user finds an incorrect constituent c and corrects it, implicitly validating the prefix tree $t_p(c)$
3. The system propose a parse tree \hat{t}' taking into account the prefix tree $t_p(c)$
4. Go to step 2
- n. The user keeps iterating until an error free parse tree is achieved

Example:



Experiments [Sánchez 09]:

- Experiments were performed using the WSJ Treebank and a modified CYK parser
- Vanilla CNF PCFG obtained from sections 02-21. Test set: section 23
- The system simulates user interaction:
 1. Explore the proposed tree and find the first wrong constituent
 2. Replace it with the correct gold constituent
 3. Perform the predictive step (obtain new tree)
 - n. Repeat until the gold tree is achieved

5.3 IPP: A FRAMEWORK FOR ACTIVE LEARNING

Evaluation and results:

- Tree Constituent Error Rate (TCER): Normalized edit distance between the proposed parse tree and the gold tree
→ *User effort when manually postediting the erroneous tree*
- Tree Constituent Action Rate (TCAC): Ration of user constituent corrections performed to obtain the reference tree using the IPP system
→ *User effort when using the IPP system*

PCFG	Baseline		IPP	RelRed
	F_1	TCER	TCAC	
h=0,v=1	0.67	0.40	0.22	45%
h=0,v=2	0.68	0.39	0.21	46%
h=0,v=3	0.70	0.38	0.22	42%

5.3 IPP: A FRAMEWORK FOR ACTIVE LEARNING

IPP-ANN tool: <http://cat.iti.upv.es/ipp/>

Parser server

- Custom Viterbi implementation
- Using PCFG in CNF
- Allows requesting subtrees with
 - a root span
 - a complete root constituent

Parser client

- Light Web-client using Flash plugin
- Decodes user feedback
- Requests subtrees to the parse server based on user corrections

Communication

- Client-server communication via sockets
- Using a library specifically tailored for interactive predictive applications

REFERENCES

- [Aho 72] A.V. Aho and J.D. Ullman. *The theory of parsing, translation, and compiling. Volumen I: parsing*. Prentice-Hall, 1972.
- [Baum 72] L.E. Baum. *An inequality and Associated Maximization Technique in Statistical Estimation for Probabilistic Functions of Markov Processes*. *Inequalities*, 3:1-9, 1972.
- [Benedi 05] J.M. Benedí and J.A. Sánchez. *Estimation of stochastic context-free grammars and their use as language models*. *Computer Speech and Language*, 19(3):249-274, 2005.
- [Booth 73] T.L. Booth and R.A. Thompson *Applying Probability Measures to Abstract Languages*. *IEEE Transactions on Computers*, 22(5):442-450, May 1973.
- [Cappé 09] O. Cappé and E. Moulines. *Online Expectation-Maximization Algorithm for Latent Data Models*, *Journal of the Royal Statistics Society: Series B (Statistical Methodology)*, 71, 2009
- [Chelba 00] C. Chelba and F. Jelinek. *Structured language modeling*. *Computer Speech and Language*, 14:283-332, 2000.
- [Dempster 77] A.P. Dempster, N.M. Laird and D.B. Rubin. *Maximum Likelihood from Incomplete Data via the EM Algorithm*. *Journal of the Royal Statistical Society. Series B (Methodological)* 39 (1)::1-38, 1977.
- [Gascó 10a] G. Gascó and J.A. Sánchez. *Syntax augmented inversion transduction grammars for machine translation*. 11th International Conference on Intelligent Text Processing and Computational Linguistics (CICLING), March, 2010.

- [Gascó 10b] G. Gascó, J.A. Sánchez and J.M. Benedí. *Enlarged Search Space for SITG Parsing*, Proc. 11th Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL HLT), June, 2010, 653-656.
- [Gascó 10c] G. Gascó, J.A. Sánchez and J.M. Benedí. *Complete search space exploration for sitg inside probability*. Proceedings of the Structural, Syntactic and Statistical Pattern Recognition, Joint IAPR International Workshop, SSPR & SPR 2010, Cesme, Izmir, Turkey, August 2010.
- [Hwa 04] R. Hwa. *Sample Selection for Statistical Parsing*. Computational Linguistics, 30(3):253-276, 2004.
- [Lari 90] K. Lari and S.J. Young. *The Estimation of Stochastic Context-Free Grammars using the Inside-Outside Algorithm*. Computer Speech and Language, 4:35-56, 1990.
- [Liang 09] P. Liang and D. Klein. *Online EM for Unsupervised Models*. Proc. 10th Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL HLT), June, 2009, 611-619.
- [Maryanski 79] F.J. Maryanski and M.T. Thomason. *Properties of stochastic syntax-directed translation schemata*. Journal of Computer and Information Sciences, 8(2):89-110, 1979.
- [Neal 98] R. Neal and G. Hinton. *A view of the EM algorithm that justifies incremental, sparse, and other variants*. Learning in Graphical Models, 355-368, 1999.
- [Ney 91] H. Ney. *Dynamic Programming Parsing for Context-Free Grammars in Continuous Speech Recognition*. IEEE Trans. Signal Processing, 39(2):336-340, 1991.

- [Pereira 92] F. Pereira and Y. Schabes. *Inside-outside reestimation from partially bracketed corpora*. Proceedings of the 30th Annual Meeting of the ACL, 128-135, 1992.
- [Roark 01] B. Roark. *Probabilistic Top-Down Parsing and Language Modeling*. Computational Linguistics, 27(2):249-276, 2001.
- [Sánchez 97] J.A. Sánchez and J.M. Benedí. *Consistency of Stochastic Context-Free Grammars from Probabilistic Estimation Based on Growth Transformation*. IEEE Trans. Pattern Analysis and Machine Intelligence, 19(2):1052-1055, 1997.
- [Sánchez 09] R. Sánchez-Sáez, J.A. Sánchez and J.M. Benedí. *Interactive predictive parsing*. In Proceedings of the 11th International Conference on Parsing Technologies (IWPT'09), 222-225, Paris, France, 2009.
- [Sánchez 10a] R. Sánchez-Sáez, L. Leiva, J.A. Sánchez and J.M. Benedí. *Interactive Predictive Parsing using a Web-based Architecture*. Proceedings of the NAACL HLT 2010 Demonstration Session, 37-40, Los Angeles, California, 2010.
- [Settles 08] B. Settles and M. Craven. *An Analysis of Active Learning Strategies for Sequence Labelling Tasks*, Empirical Methods in Natural Language Processing (EMNLP), 1069-1078, 2008.
- [Settles 10] B. Settles. *Active Learning Literature Survey*. Computer Sciences Technical Report 1648, University of Wisconsin-Madison, 2010.
- [Stolcke 95] A. Stolcke. *An Efficient Probabilistic Context-Free Parsing Algorithm that Computes Prefix Probabilities*. Computational Linguistics, 21(2):165-200, 1995.

REFERENCES

- [Wetherell 80] C.S. Wetherell. *Probabilistic Languages: A Review and some Open Questions*. Computing Surveys, 12(4):361-379, 1980.
- [Wu 97] D. Wu. *Stochastic Inversion Transduction Grammars and Bilingual Parsing of Parallel Corpora*. Computational Linguistics, 23(3):377-403, 1997.

APPENDICES

A growth transformation can be defined as:

$$\bar{p}(A \rightarrow \alpha) = \frac{p(A \rightarrow \alpha) \left(\frac{\partial \Pr_{G_s}(\Omega, \Delta_\Omega)}{\partial p(A \rightarrow \alpha)} \right)_p}{\sum_{i=1}^{n_A} p(A \rightarrow \alpha_i) \left(\frac{\partial \Pr_{G_s}(\Omega, \Delta_\Omega)}{\partial p(A \rightarrow \alpha_i)} \right)_p}$$

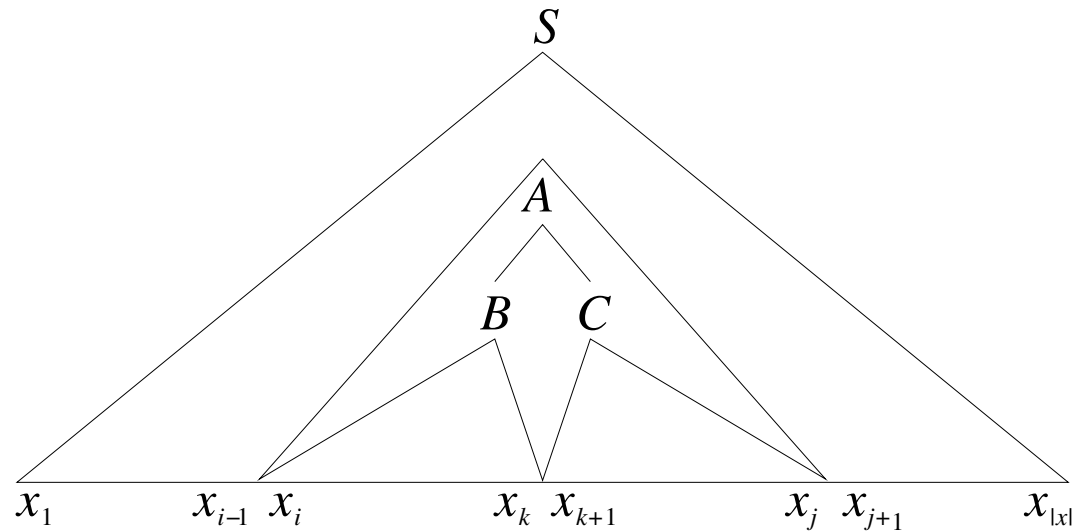
Numerator:

$$\begin{aligned} p(A \rightarrow \alpha) \left(\frac{\partial \Pr_{G_s}(\Omega, \Delta_\Omega)}{\partial p(A \rightarrow \alpha)} \right)_p &= \Pr_{G_s}(\Omega, \Delta_\Omega) \sum_{x \in \Omega} \frac{p(A \rightarrow \alpha)}{\Pr_{G_s}(x, \Delta_x)} \left(\frac{\partial \Pr_{G_s}(x, \Delta_x)}{\partial p(A \rightarrow \alpha)} \right)_p \\ &= \Pr_{G_s}(\Omega, \Delta_\Omega) \sum_{x \in \Omega} \frac{p(A \rightarrow \alpha)}{\Pr_{G_s}(x, \Delta_x)} \sum_{\forall d_x \in \Delta_x} \left(\frac{\partial \Pr_{G_s}(x, d_x)}{\partial p(A \rightarrow \alpha)} \right)_p \\ &= \Pr_{G_s}(\Omega, \Delta_\Omega) \sum_{x \in \Omega} \frac{1}{\Pr_{G_s}(x, \Delta_x)} \sum_{\forall d_x \in \Delta_x} \mathbf{N}(A \rightarrow \alpha, d_x) \Pr_{G_s}(x, d_x) \end{aligned}$$

Denominator:

$$\begin{aligned}
 & \sum_{i=1}^{n_A} p(A \rightarrow \alpha_i) \left(\frac{\partial \Pr_{G_s}(\Omega, \Delta_\Omega)}{\partial p(A \rightarrow \alpha_i)} \right)_p = \\
 &= \Pr_{G_s}(\Omega, \Delta_\Omega) \sum_{x \in \Omega} \frac{1}{\Pr_{G_s}(x, \Delta_x)} \sum_{\forall d_x \in \Delta_x} \sum_{i=1}^{n_A} \mathbf{N}(A \rightarrow \alpha_i, d_x) \Pr_{G_s}(x, d_x) \\
 &= \Pr_{G_s}(\Omega, \Delta_\Omega) \sum_{x \in \Omega} \frac{1}{\Pr_{G_s}(x, \Delta_x)} \sum_{\forall d_x \in \Delta_x} \mathbf{N}(A, d_x) \Pr_{G_s}(x, d_x).
 \end{aligned}$$

- Let $A \rightarrow BC$ in a position delimited by integers i, j, k , $1 \leq i \leq k < j \leq |x|$



- $\Delta_{x,i,j,k,A \rightarrow BC} \subseteq D_x$: subset of derivations of x in which the rule $A \rightarrow BC$ appears delimited by positions i, j, k
- $\Delta_{x,i,j,A}$: subset of derivations of x in which the non-terminal A appears delimited by positions i, j

$$\begin{aligned}
 \blacktriangleright \sum_{\forall d_x \in D_x} N(A \rightarrow BC, d_x) \Pr_{G_s}(x, d_x) &= \sum_{1 \leq i \leq k < j \leq |x|} \sum_{\forall d_x \in \Delta_{x,i,j,k,A \rightarrow BC}} \Pr_{G_s}(x, d_x) \\
 &= \sum_{1 \leq i \leq k < j \leq |x|} \Pr_{G_s}(S \xrightarrow{*} x_1 \dots x_{i-1} A x_{j+1} \dots x_{|x|}) \cdot \\
 &\quad p(A \rightarrow BC) \cdot \Pr_{G_s}(B \xrightarrow{*} x_i \dots x_k) \cdot \Pr_{G_s}(C \xrightarrow{*} x_{k+1} \dots x_j) \\
 &= \sum_{1 \leq i \leq k < j \leq |x|} f(A < i, j >) p(A \rightarrow BC) e(B < i, k >) e(C < k + 1, j >),
 \end{aligned}$$

$$\begin{aligned}
 \blacktriangleright \sum_{\forall d_x \in D_x} N(A, d_x) \Pr_{G_s}(x, d_x) &= \sum_{1 \leq i \leq j \leq |x|} \sum_{\forall d_x \in \Delta_{x,i,j,A}} \Pr_{G_s}(x, d_x) \\
 &= \sum_{1 \leq i \leq j \leq |x|} \Pr_{G_s}(S \xrightarrow{*} x_1 \dots x_{i-1} A x_{j+1} \dots x_{|x|}) \Pr_{G_s}(A \xrightarrow{*} x_i \dots x_j) \\
 &= \sum_{1 \leq i \leq j \leq |x|} f(A < i, j >) e(A < i, j >).
 \end{aligned}$$

EM algorithm [Neal 98]:

E step: Compute a distribution $\tilde{p}^{(t)}$ over the range of \mathbf{Z} such that

$$\tilde{p}^{(t)}(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}; \theta^{(t-1)})$$

M step: Set $\theta^{(t)}$ to the θ that maximizes $E_{\tilde{p}^{(t)}}[\log p(\mathbf{x}, \mathbf{z}; \theta)]$