

# Statistical models for natural language parsing

Mark Johnson

Microsoft Research / Brown University

July 2007

# Outline

## Introduction

Non-local dependencies and the PCFG MLE

Generative statistical parsers

Exponential (a.k.a. Maximum Entropy) parsing models

Coarse to fine reranking

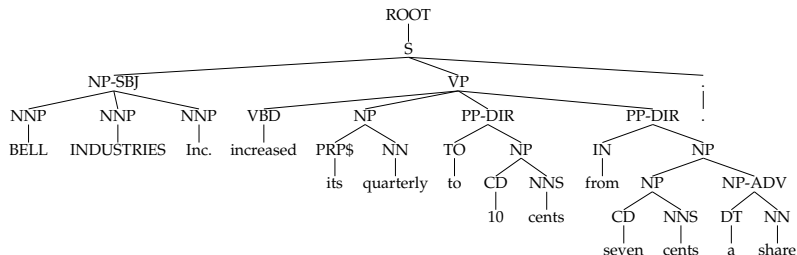
Self-training of the reranking parser

Sample parser errors

## Preview of natural language parsing

- Non-local dependencies cause PCFG Maximum Likelihood Estimator (MLE) to produce sub-optimal grammars
- State-splitting or decorating with features can make non-local dependencies local
- Exponential (a.k.a. Maximum Entropy) models aren't as adversely affected by non-local dependencies as PCFGs
- But MLE seems difficult to compute  $\Rightarrow$  Maximum Conditional Likelihood Estimation (MCLE)
- MCLE also seems better suited to parsing tasks if PCFG doesn't accurately describe distribution of strings
- Coarse-to-fine reranking combines PCFG and exponential models to produce the most accurate parsers we have today

# Treebank corpora



- The Penn treebank contains hand-annotated parse trees for ~ 50,000 sentences
- Treebanks also exist for the Brown corpus, the Switchboard corpus (spontaneous telephone conversations) and Chinese and Arabic corpora

# Outline

Introduction

Non-local dependencies and the PCFG MLE

Generative statistical parsers

Exponential (a.k.a. Maximum Entropy) parsing models

Coarse to fine reranking

Self-training of the reranking parser

Sample parser errors

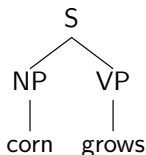
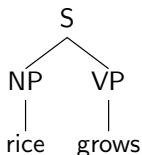
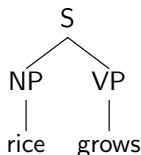
# Estimating a grammar from a treebank

- *Maximum likelihood principle*: Choose the grammar and rule probabilities that make the trees in the corpus *as likely as possible*
  - ▶ read the rules off the trees
  - ▶ for PCFGs, set rule probabilities to the *relative frequency* of each rule in the treebank

$$P(\text{VP} \rightarrow \text{V NP}) = \frac{\text{Number of times VP} \rightarrow \text{V NP occurs}}{\text{Number of times VP occurs}}$$

- *If the language is generated by a PCFG and the treebank trees are its derivation trees, the estimated grammar converges to the true grammar.*

# Estimating PCFGs from visible data

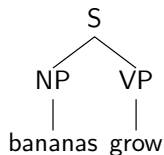
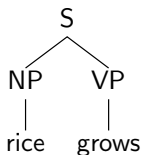
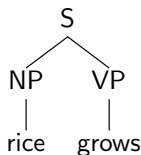


Rule	Count	Rel Freq
$S \rightarrow NP VP$	3	1
$NP \rightarrow \text{rice}$	2	$2/3$
$NP \rightarrow \text{corn}$	1	$1/3$
$VP \rightarrow \text{grows}$	3	1

$$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{rice} \quad \text{grows} \end{array} \right) = 2/3$$

$$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{corn} \quad \text{grows} \end{array} \right) = 1/3$$

# Non-local dependencies and PCFG MLE



Rule	Count	Rel Freq
S $\rightarrow$ NP VP	3	1
NP $\rightarrow$ rice	2	2/3
NP $\rightarrow$ bananas	1	1/3
VP $\rightarrow$ grows	2	2/3
VP $\rightarrow$ grow	1	1/3

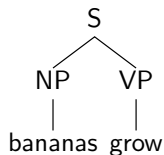
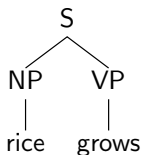
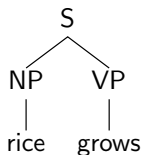
$$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ rice \quad grows \end{array} \right) = 4/9$$

$$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ bananas \quad grow \end{array} \right) = 1/9$$

*partition function*  $Z = 5/9$



# Dividing by partition function $Z$



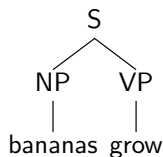
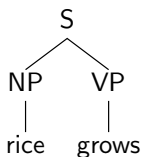
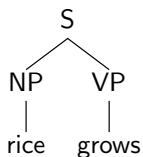
Rule	Count	Rel Freq
$S \rightarrow NP VP$	3	1
$NP \rightarrow \text{rice}$	2	$2/3$
$NP \rightarrow \text{bananas}$	1	$1/3$
$VP \rightarrow \text{grows}$	2	$2/3$
$VP \rightarrow \text{grow}$	1	$1/3$

$$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{rice} \quad \text{grows} \end{array} \right) = 4/9 \quad 4/5$$

$$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{bananas} \quad \text{grow} \end{array} \right) = 1/9 \quad 1/5$$

$$Z = \underline{5/9}$$

# Other values do better!



Rule	Count	$p$
$S \rightarrow NP VP$	3	1
$NP \rightarrow \text{rice}$	2	2/3
$NP \rightarrow \text{bananas}$	1	1/3
$VP \rightarrow \text{grows}$	2	1/2
$VP \rightarrow \text{grow}$	1	1/2

(Abney 1997)

$$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{rice} \quad \text{grows} \end{array} \right) = 2/6 \quad 2/3$$

$$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ NP \quad VP \\ | \quad | \\ \text{bananas} \quad \text{grow} \end{array} \right) = 1/6 \quad 1/3$$

$$Z = \underline{3/6}$$

# Make dependencies local – GPSG-style

rule	count	rel freq	
$S \rightarrow \begin{matrix} \text{NP} \\ +\text{singular} \end{matrix} \begin{matrix} \text{VP} \\ +\text{singular} \end{matrix}$	2	2/3	$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ \begin{matrix} \text{NP} \\ +\text{singular} \\   \\ \text{rice} \end{matrix} \quad \begin{matrix} \text{VP} \\ +\text{singular} \\   \\ \text{grows} \end{matrix} \end{array} \right) = 2/3$
$S \rightarrow \begin{matrix} \text{NP} \\ +\text{plural} \end{matrix} \begin{matrix} \text{VP} \\ +\text{plural} \end{matrix}$	1	1/3	
$\begin{matrix} \text{NP} \\ +\text{singular} \end{matrix} \rightarrow \text{rice}$	2	1	$P \left( \begin{array}{c} S \\ \swarrow \quad \searrow \\ \begin{matrix} \text{NP} \\ +\text{plural} \\   \\ \text{bananas} \end{matrix} \quad \begin{matrix} \text{VP} \\ +\text{plural} \\   \\ \text{grow} \end{matrix} \end{array} \right) = 1/3$
$\begin{matrix} \text{NP} \\ +\text{plural} \end{matrix} \rightarrow \text{bananas}$	1	1	
$\begin{matrix} \text{VP} \\ +\text{singular} \end{matrix} \rightarrow \text{grows}$	2	1	
$\begin{matrix} \text{VP} \\ +\text{plural} \end{matrix} \rightarrow \text{grow}$	1	1	

# Outline

Introduction

Non-local dependencies and the PCFG MLE

**Generative statistical parsers**

Exponential (a.k.a. Maximum Entropy) parsing models

Coarse to fine reranking

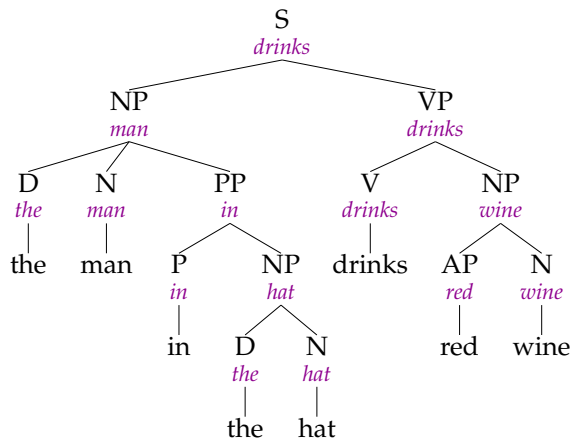
Self-training of the reranking parser

Sample parser errors

# Generative statistical parsers

- Splitting node labels (a.k.a. decorating the tree with features) enables PCFG to capture non-local dependencies
- Modern generative statistical parsers track around 7 different non-local dependencies
- These dependencies are encoded as “features” on nodes
- Most combinations of features are not observed in training data, but will occur in new sentences  
⇒ smoothing is essential!

# “Head to head” dependencies



## Rules:

$S \rightarrow NP \ VP$   
*drinks*  $\rightarrow$  *man* *drinks*

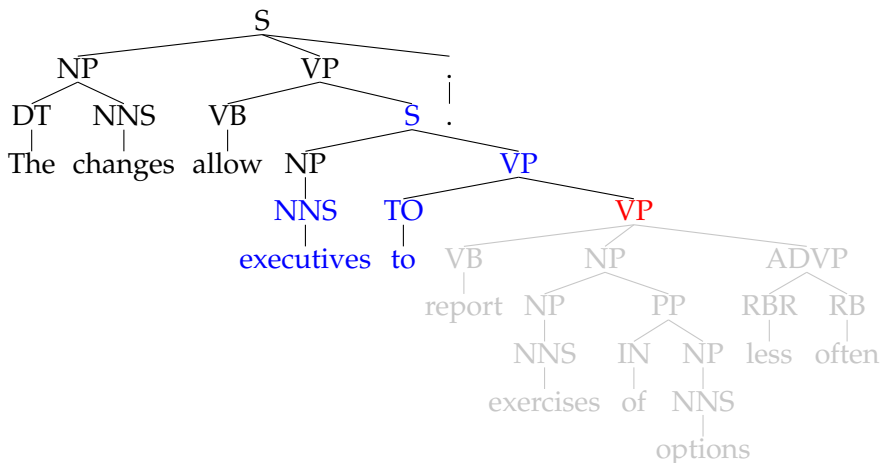
$VP \rightarrow V \ NP$   
*drinks*  $\rightarrow$  *drinks* *wine*

$NP \rightarrow AP \ N$   
*wine*  $\rightarrow$  *red* *wine*

...

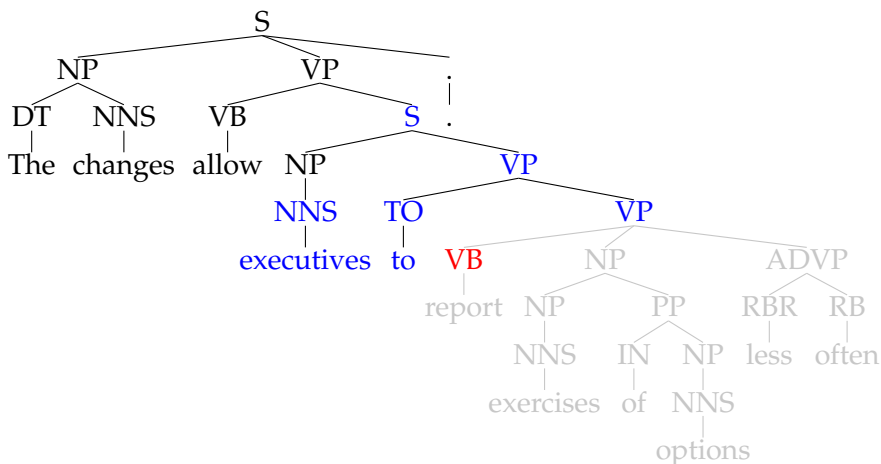
- *Lexicalization* captures syntactic and semantic dependencies
- Lexicalized structural preferences may be most important

# Generative language model (Charniak 2001)



- Predicted node is shown in red
- Conditioning nodes are shown in blue

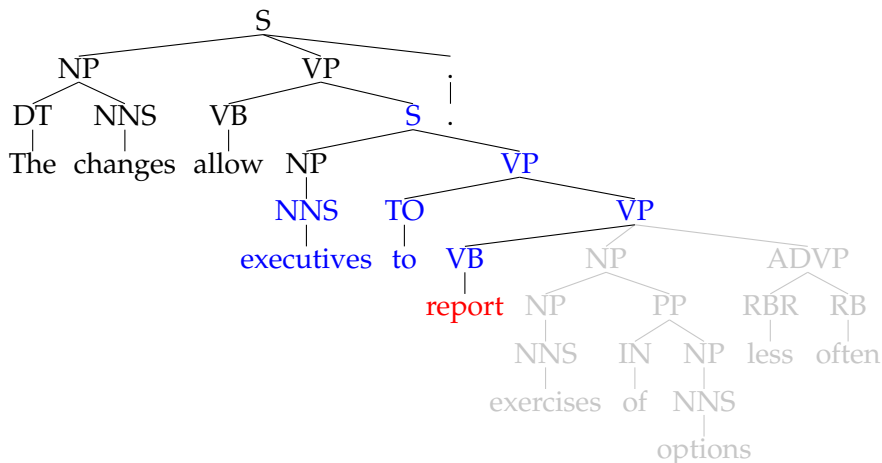
# Generative language model (Charniak 2001)



- Predicted node is shown in red
- Conditioning nodes are shown in blue

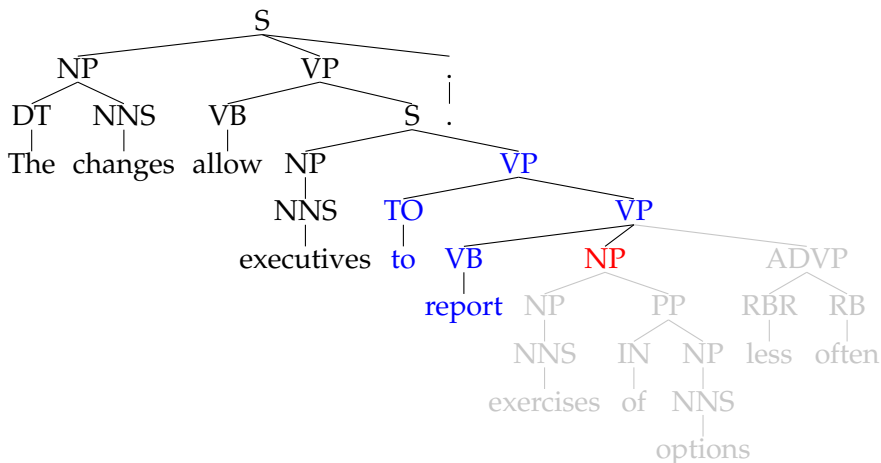


# Generative language model (Charniak 2001)



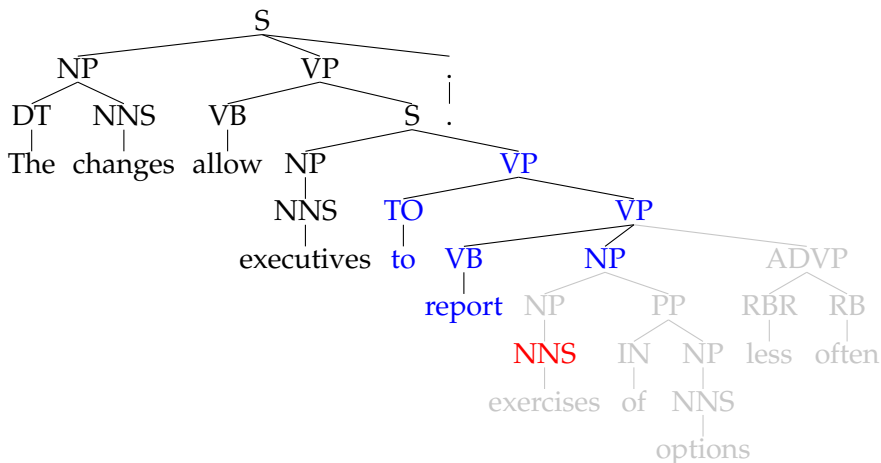
- Predicted node is shown in red
- Conditioning nodes are shown in blue

# Generative language model (Charniak 2001)



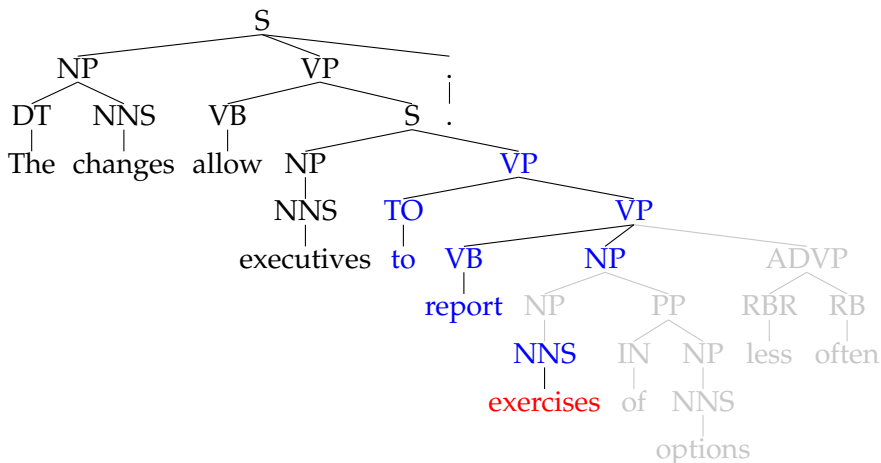
- Predicted node is shown in red
- Conditioning nodes are shown in blue

# Generative language model (Charniak 2001)



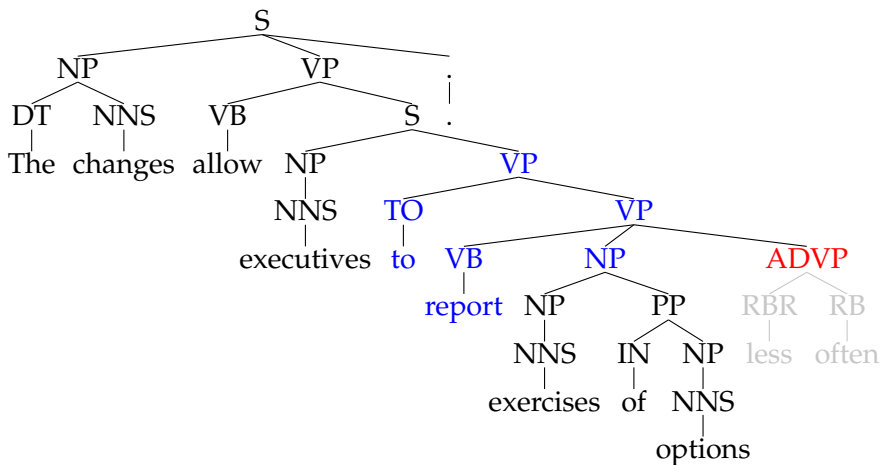
- Predicted node is shown in red
- Conditioning nodes are shown in blue

# Generative language model (Charniak 2001)



- Predicted node is shown in red
- Conditioning nodes are shown in blue

# Generative language model (Charniak 2001)



- Predicted node is shown in red
- Conditioning nodes are shown in blue

## Summary so far

- Maximum likelihood is a good way of estimating a grammar
  - Maximum likelihood estimation of a PCFG from a treebank is easy, and works well *if the trees are accurate*
  - But real language has many more dependencies than treebank grammar describes
- ⇒ relative frequency estimator not MLE
- ▶ Make non-local dependencies local by splitting categories
  - ⇒ Astronomical number of possible categories
- Find some way of accurately estimating models in the presence of unmodeled dependencies
  - ⇒ exponential models

# Outline

Introduction

Non-local dependencies and the PCFG MLE

Generative statistical parsers

Exponential (a.k.a. Maximum Entropy) parsing models

Coarse to fine reranking

Self-training of the reranking parser

Sample parser errors

# Exponential models

Exponential models are defined in terms of features, where a *feature* is any real-valued function on  $\Psi_G$ .

Let  $f_1, \dots, f_m$  be features, and  $\lambda_1, \dots, \lambda_m$  be real-valued feature weights. An *exponential model* has the form:

$$\begin{aligned}P_\lambda(\psi) &= \frac{W_\lambda(\psi)}{Z_\lambda} \\W_\lambda(\psi) &= \exp \sum_{j=1}^m \lambda_j f_j(\psi) \\Z_\lambda &= \sum_{\psi' \in \Psi_G} W_\lambda(\psi')\end{aligned}$$

$W_\lambda(\psi)$  is the weight (unnormalized probability) of parse  $\psi$ .

$Z_\lambda$  is called the *partition function*.

Exponential models are also known as *Gibbs models*, *log-linear models* and *Maximum Entropy models*.



## PCFGs are exponential models

$\Psi$  = set of all trees generated by PCFG  $G$

$f_j(\psi)$  = number of times the  $j$ th rule is used in  $\psi$

$p(r_j)$  = probability of  $j$ th rule in  $G$

Set weight  $\lambda_j = \log p(r_j)$

$$f \left( \begin{array}{c} \text{S} \\ \swarrow \quad \searrow \\ \text{NP} \quad \text{VP} \\ | \quad | \\ \text{rice} \quad \text{grows} \end{array} \right) = [ \underbrace{1}_{\text{S} \rightarrow \text{NP VP}}, \underbrace{1}_{\text{NP} \rightarrow \text{rice}}, \underbrace{0}_{\text{NP} \rightarrow \text{bananas}}, \underbrace{1}_{\text{VP} \rightarrow \text{grows}}, \underbrace{0}_{\text{VP} \rightarrow \text{grow}} ]$$

$$P(\psi) = \prod_{j=1}^m p(r_j)^{f_j(\psi)} = \prod_{j=1}^m (\exp \lambda_j)^{f_j(\psi)} = \exp \sum_{j=1}^m \lambda_j f_j(\psi)$$

So *a PCFG is just an exponential model* with  $Z_\lambda = 1$ .

# Advantages of exponential models

- Exponential models are very flexible ...
- Features  $f$  can be *any function of parses* ...
  - ▶ whether a particular structure occurs in a parse
  - ▶ conjunctions of prosodic and syntactic structure
- Parses  $\psi$  need not be trees, but *can be anything at all*
  - ▶ Feature structures (LFG, HPSG)
- Exponential models are related to other popular models
  - ▶ Harmony theory and optimality theory
  - ▶ They are also called *Maximum Entropy* models and *log-linear* models

# Modeling dependencies

- It's usually difficult to design a PCFG model that captures a particular set of dependencies
  - ▶ probability of the tree must be broken down into a product of *conditional probability distributions*
  - ▶ non-local dependencies must be expressed in terms of GPSG-style feature passing
- It's easy to make exponential models sensitive to new dependencies
  - ▶ add a new feature functions to existing feature functions
  - ▶ *figuring out what the right dependencies are is hard, but incorporating them into an exponential model is easy*

# MLE of exponential models from visible data

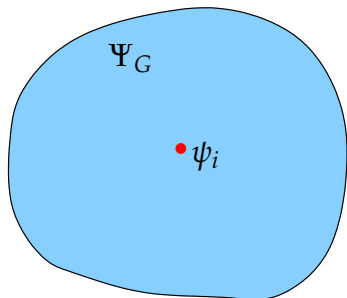
Visible training data: Parses  $\Psi = \psi_1, \dots, \psi_n$

$$\begin{aligned}\log L(\lambda) &= \sum_{i=1}^n \log P_\lambda(\psi_i) \\ &= \sum_{i=1}^n \left( \log W_\lambda(\psi_i) - \log \sum_{\psi \in \Psi_G} W_\lambda(\psi) \right)\end{aligned}$$

$$\frac{\partial \log L(\lambda)}{\partial \lambda_j} = \sum_{i=1}^n (f_j(\psi_i) - \mathbb{E}_\lambda[f_j])$$

So the likelihood is maximized when the empirical frequency of each feature equals its expected frequency.

## Maximizing likelihood of visible data is hard!



Maximizing likelihood requires summation over all of  $\Psi_G$ , even with fully visible data!

Maximizing likelihood contrasts the training data trees  $\Psi$  with  $\Psi_G$ ; i.e., select  $\lambda$  to maximize

$$\sum_{i=1}^n (\log W_\lambda(\psi_i) - \log \sum_{\psi \in \Psi_G} W_\lambda(\psi)).$$

But  $\Psi_G$  is the set of all parses of all sentences!

## Estimation by maximizing conditional likelihood

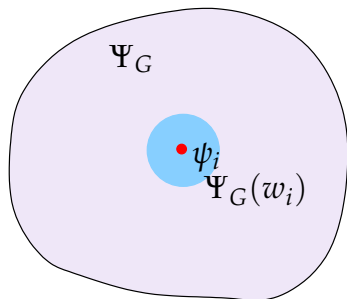
Maximize the *conditional likelihood* of the correct parses  $\Psi$  *given their yield*  $w$ .

$$\begin{aligned}\log L(\lambda) &= \sum_{i=1}^n \log P_{\lambda}(\psi_i | w_i) \\ &= \sum_{i=1}^n \left( \log W_{\lambda}(\psi_i) - \log \sum_{\psi \in \Psi_G(w_i)} W_{\lambda}(\psi) \right)\end{aligned}$$

$$\frac{\partial \log L(\lambda)}{\partial \lambda_j} = \sum_{i=1}^n (f_j(\psi_i) - E_{\lambda}[f_j | w_i])$$

So conditional likelihood is maximized when the empirical frequency of each feature equals its expected frequency *conditioned on the yields*.

## Maximizing conditional likelihood is easier



Pseudo-likelihood is consistent  
*for the conditional distribution*

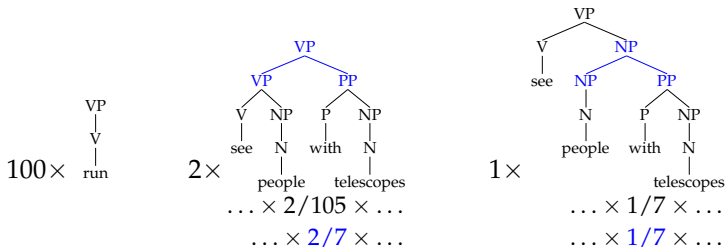
Maximizing conditional likelihood requires summing over  $\Psi_G(w_i), i = 1, \dots, n$  (obtained by parsing).

Conditional likelihood contrasts each element of training data  $\psi_i$  with the parses of  $w_i$ ; i.e., adjust  $\lambda$  to maximize  $\sum_{i=1}^n (\log W_\lambda(\psi_i) - \log \sum_{\psi' \in \Psi_G(w_i)} W_\lambda(\psi'))$ .

# Conditional likelihood is better for parsing

Parsing exploits  $P(\psi|w)$ , which MCL optimizes.

If the grammar does not generate strings accurately, ML and MCL can be quite different!



Rule	count	rel freq
VP $\rightarrow$ V	100	100/105 <b>4/7</b>
VP $\rightarrow$ V NP	3	3/105 <b>1/7</b>
VP $\rightarrow$ VP PP	2	2/105 <b>2/7</b>
NP $\rightarrow$ N	6	6/7 <b>6/7</b>
NP $\rightarrow$ NP PP	1	1/7 <b>1/7</b>



## Conditional ML estimation

$w_i$	$f(\psi_i)$	$\{f(\psi) : \psi \in \Psi_G(w_i), \psi \neq \psi_i\}$
sentence 1	(1, 3, 2)	(2, 2, 3) (3, 1, 5) (2, 6, 3)
sentence 2	(7, 2, 1)	(2, 5, 5)
sentence 3	(2, 4, 2)	(1, 1, 7) (7, 2, 1)
...	...	...

- Parser designer specifies *feature functions*  $f = (f_1, \dots, f_m)$
- A parser produces trees  $\Psi(w)$  for each sentence  
 $w \in w_1, \dots, w_n$
- Treebank tells us correct tree  $\psi_i \in \Psi(w_i)$  for sentence  $w_i$
- Feature functions  $f$  apply to each tree  $\psi \in \Psi_G(w)$ ,  
producing feature values  $f(\psi) = (f_1(\psi), \dots, f_m(\psi))$
- MCLE estimates feature weights  $\hat{\lambda}$  using a gradient-based numerical optimizer

# Regularization

- With a large number of features, exponential models can over-fit the training data
- Regularization: add *bias* term to ensure  $\hat{\lambda}$  is finite and small
- In following experiments, regularizer is a **polynomial penalty term**

$$\begin{aligned}\hat{\lambda} &= \operatorname{argmax}_{\lambda} \log \sum_{i=1}^n P_{\lambda}(\psi_i | w_i) - c \sum_{j=1}^m |\lambda_j|^p \\ &= \operatorname{argmax}_w \sum_{i=1}^n \left( \sum_{j=1}^m \lambda_j f_j(\psi_i) - \log Z_{\lambda}(w_i) \right) - c \sum_{j=1}^m |\lambda_j|^p\end{aligned}$$

- $p = 2$  gives a *Gaussian prior*.
- We maximize this expression using *numerical optimization* (Limited Memory Variable Metric)

## Conditional vs joint estimation

In this slide, let  $\psi$  be a parse tree without the terminal string  $w$

$$P(\psi, w) = P(\psi|w)P(w)$$

- ML optimizes probability of training trees  $\psi$  *and strings*  $w$
- MCLE maximizes probability of trees given strings
  - ▶ Conditional estimation uses less information from the data
  - ▶ learns nothing from distribution of strings  $P(w)$
  - ▶ learns nothing from unambiguous sentences (!)
- Joint estimation should be better (lower variance) if your model correctly relates  $P(\psi|w)$  and  $P(w)$
- Conditional estimation should be better if your model incorrectly relates  $P(\psi|w)$  and  $P(w)$

# Linguistic representations and features

- Probability of a parse  $\psi$  is completely determined by its feature vector  $(f_1(\psi), \dots, f_m(\psi))$
- The actual linguistic representation of parse  $\psi$  is *irrelevant* as long as it is rich enough to calculate features  $f(\psi)$
- Feature functions define the kinds of generalizations that the learner can extract
  - ▶ parses with the same feature values will be assigned the same probability
  - ▶ the choice of feature functions is as much a linguistic decision than the choice of representations
- Features can be arbitrary functions
  - ▶ the linguistic properties they encode *need not be directly represented in the parse*
  - ▶ very different from PCFGs, where the tree label and shape determines the generalizations extracted

# Outline

Introduction

Non-local dependencies and the PCFG MLE

Generative statistical parsers

Exponential (a.k.a. Maximum Entropy) parsing models

**Coarse to fine reranking**

Self-training of the reranking parser

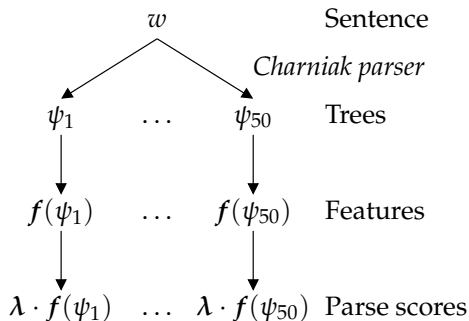
Sample parser errors

## Coarse to fine parsing

- Parsing with a grammar with a lot of features (PCFG nonterminals) is slow, even using the dynamic programming algorithms
- Coarse to fine parsing uses a sequence of grammars. The features of the coarse-grained grammars are equivalence classes of the fine-grained features.
- The parses produced by the coarse-grained grammars constrain the search with the fine-grained grammar.
- The Charniak generative parser uses a coarse-grained PCFG to identify which substrings should be parsed with the fine-grained PCFG.

## Coarse to fine reranking with exponential models

- $Z_G(w)$  is still hard to compute  $\Rightarrow$  make  $\Psi(w)$  even smaller!
- Set  $\Psi(w) =$  the 50-best parses produced by Charniak parser
- Exponential model is trained using MCLE to pick out best parse from Charniak's 50-best parses



# Features for ranking parses

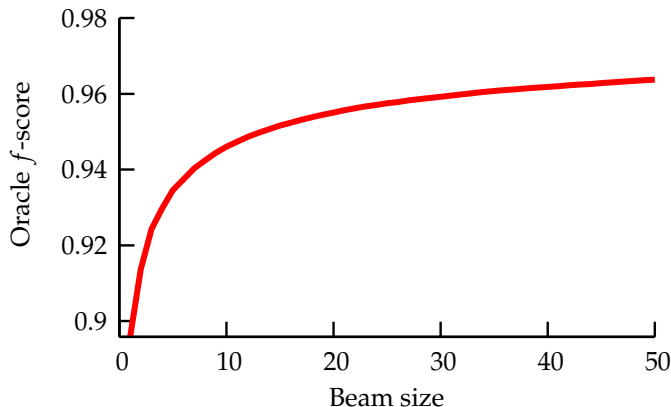
- Features can be any real-valued function of parse trees
- In these experiments the features come in two kinds:
  - ▶ The logarithm of the tree's probability estimated by the Charniak parser
  - ▶ The number of times a particular configuration appears in the parse
- *Which ones improve parsing accuracy the most?* (can you guess?)



# Experimental setup

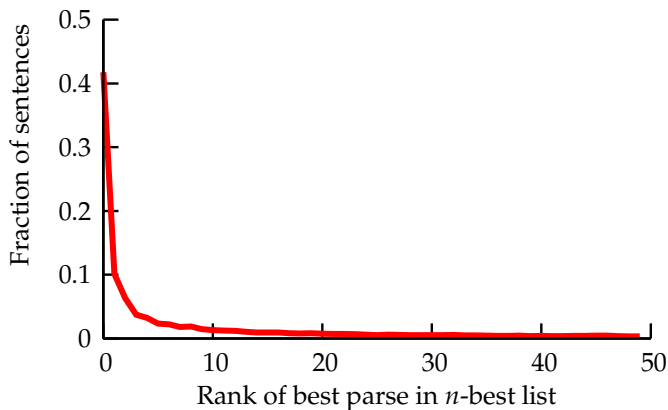
- Feature tuning experiments done using Collins' split: sections 2-19 as train, 20-21 as dev and 22 as test
- $\Psi(w)$  computed using Charniak 50-best parser
- Features which vary on less than 5 sentences pruned
- Optimization performed using LMVM optimizer from Petsc/TAO optimization package
- Regularizer constant  $c$  adjusted to maximize f-score on dev

## $f$ -score vs. $n$ -best beam size



- F-score of Charniak's most probable parse = 0.896
- Oracle f-score (f-score of best parse in beam) of Charniak's 50-best parses = 0.965 (66% redn)

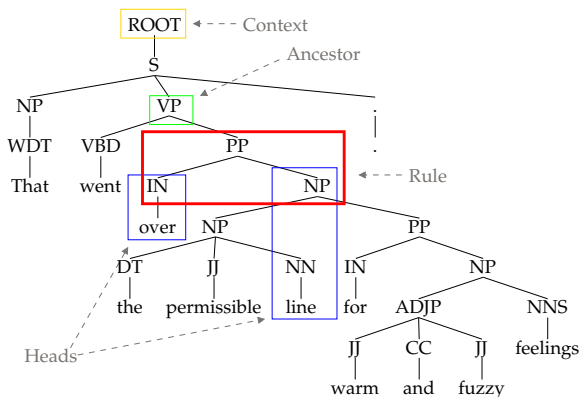
## Rank of best parse



- Charniak parser's most likely parse is the best parse 41% of the time
- Reranker picks Charniak parser's most likely parse 58% of the time

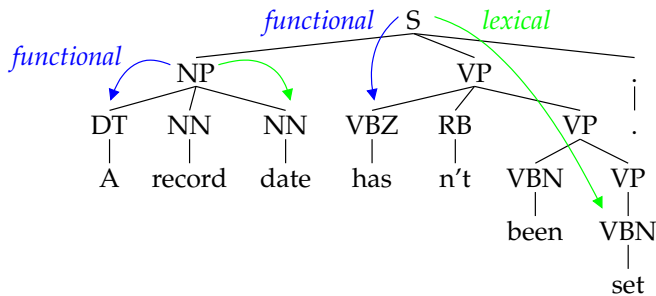
# Lexicalized and parent-annotated rules

- Rule features largely replicate features already in generative parser
- A typical Rule feature might be (PP IN NP)



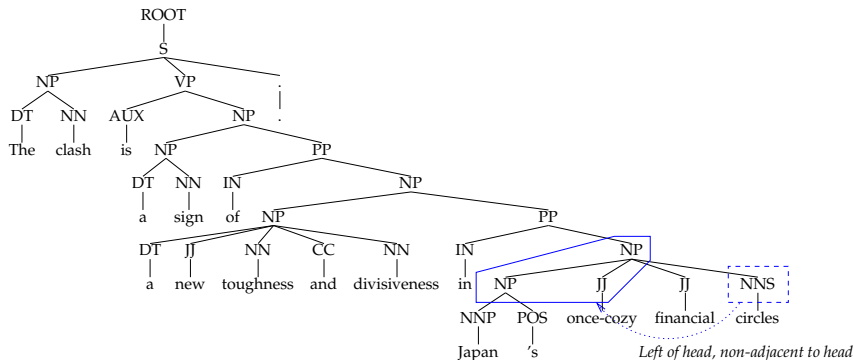
# Functional and lexical heads

- There are at least two sensible notions of head (c.f., Grimshaw)
  - ▶ *Functional heads*: determiners of NPs, auxiliary verbs of VPs, etc.
  - ▶ *Lexical heads*: rightmost Ns of NPs, main verbs in VPs, etc.
- In a log-linear model, it is easy to use both!



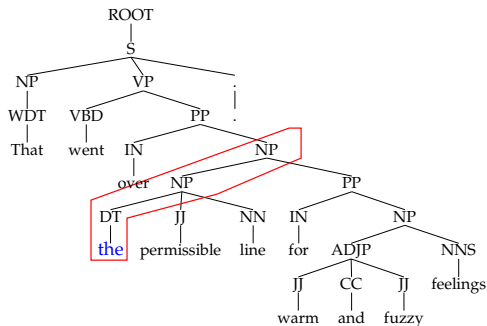
# $n$ -gram rule features generalize rules

- Breaks up long treebank constituents into shorter (phrase-like?) chunks
- Also includes *relationship to head* (e.g., adjacent? left or right?)



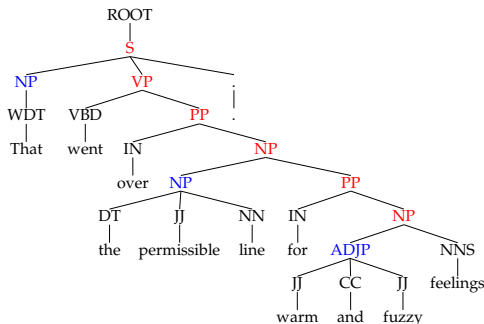
# Word and WProj features

- A Word feature is a word plus  $n$  of its parents (c.f., Klein and Manning's non-lexicalized PCFG)
- A WProj feature is a word plus all of its (maximal projection) parents, up to its governor's maximal projection



# Rightmost branch bias

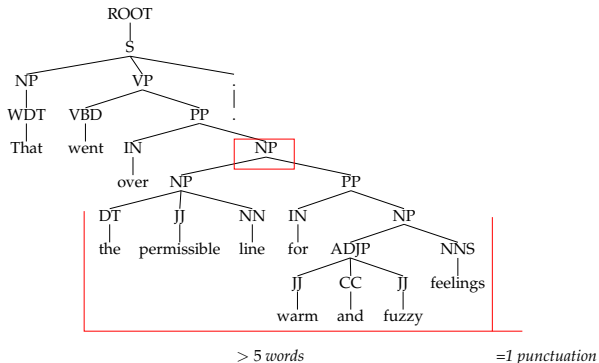
- The RightBranch feature's value is the number of nodes on the right-most branch (ignoring punctuation) (c.f., Charniak 00)
- Reflects the tendency toward right branching in English
- Only 2 different features, but very useful in final model!





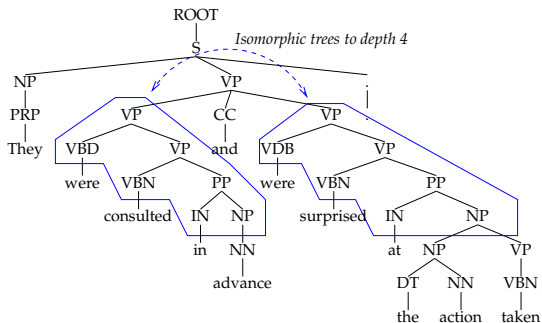
# Constituent Heaviness and location

- Heaviness measures the constituent's category, its (binned) size and (binned) closeness to the end of the sentence



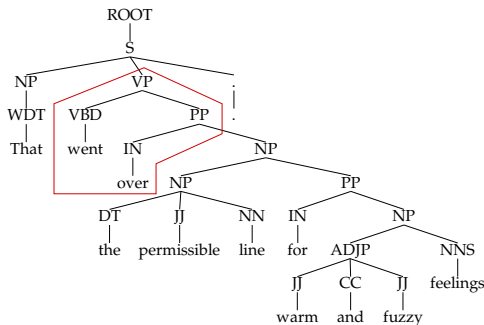
# Coordination parallelism

- A CoPar feature indicates the depth to which adjacent conjuncts are parallel



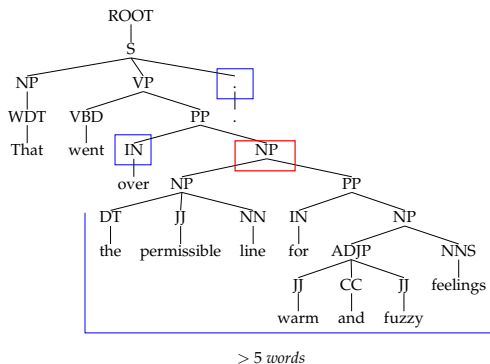
# Tree $n$ -gram

- A tree  $n$ -gram feature is a tree fragment that connect sequences of adjacent  $n$  words, for  $n = 2, 3, 4$  (c.f. Bod's DOP models)
- lexicalized and non-lexicalized variants

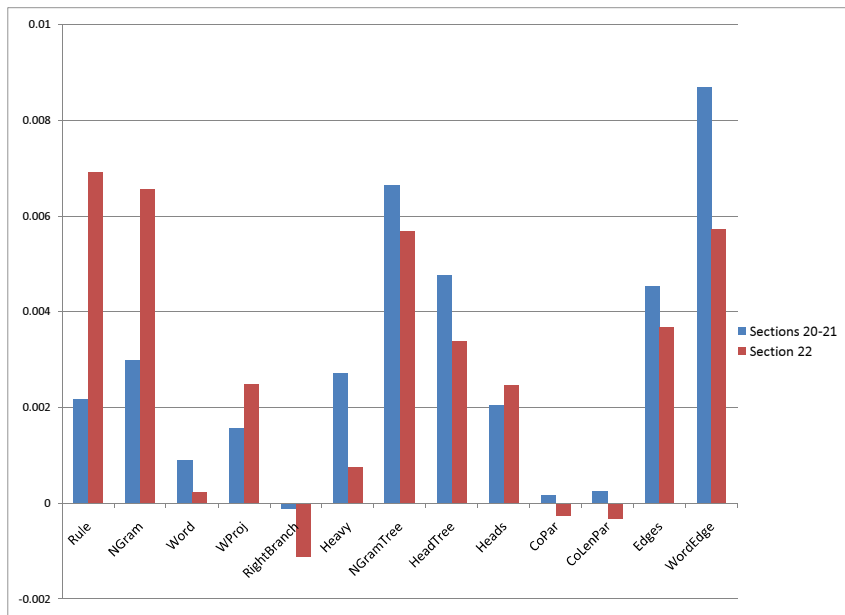


# Edges and WordEdges

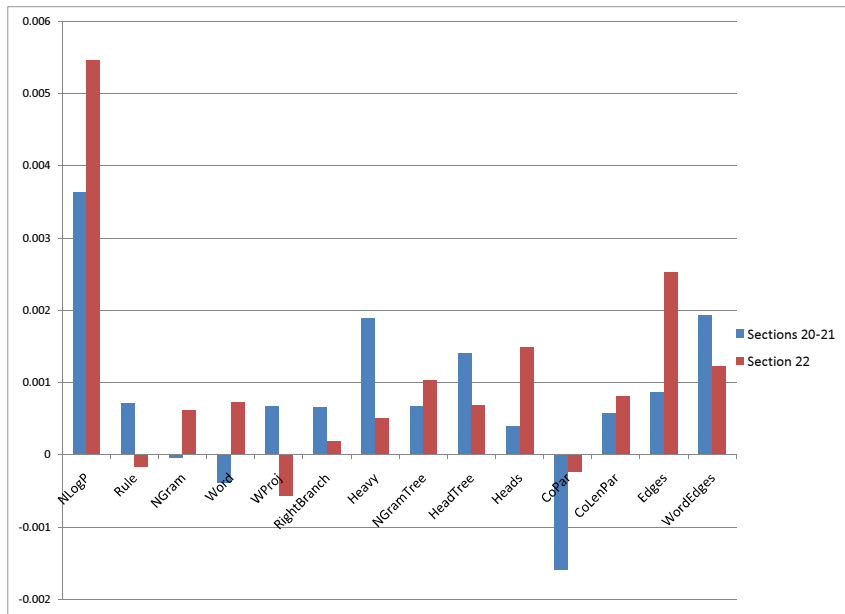
- A Neighbours feature indicates the node's category, its binned length and  $j$  left and  $k$  right lexical items and/or POS tags for  $j, k \leq 2$



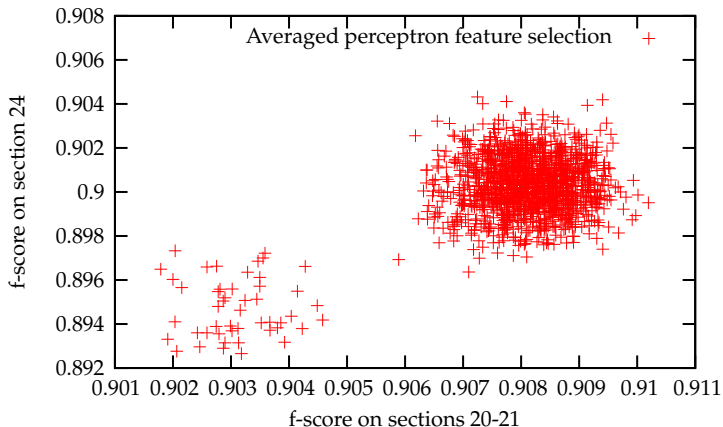
# Adding one feature class to baseline parser



# Removing one feature class from reranker



## Feature selection is hard



- Greedy feature selection using *averaged perceptron* optimizing f-score on sec 20–21
- All models also evaluated on section 24

## Results on all training data

- Features must vary on parses of at least 5 sentences in training data
- In this experiment, 1,333,863 features
- Exponential model trained on sections 2-21
- Gaussian regularization  $p = 2$ , constant selected to optimize f-score on section 22
- On section 23: *recall = 91.0, precision = 91.8, f-score = 91.4*
- Available from [www.cog.brown.edu](http://www.cog.brown.edu)



# Outline

Introduction

Non-local dependencies and the PCFG MLE

Generative statistical parsers

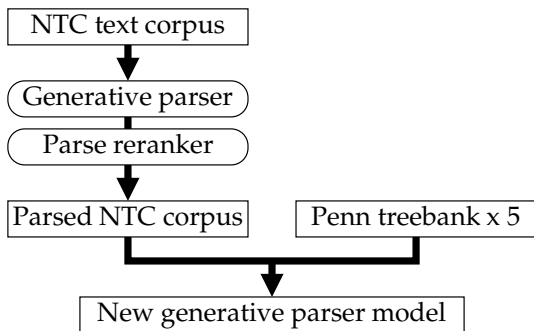
Exponential (a.k.a. Maximum Entropy) parsing models

Coarse to fine reranking

**Self-training of the reranking parser**

Sample parser errors

# Self-training for reranking parsing



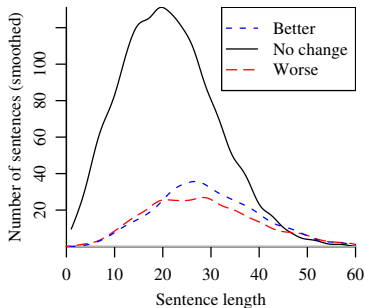
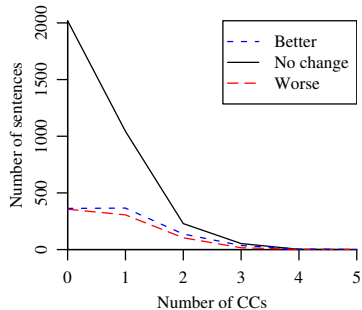
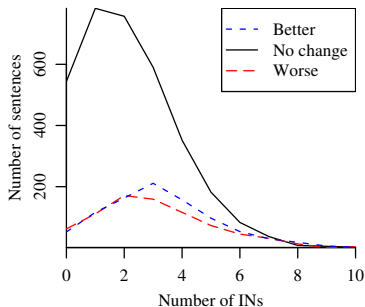
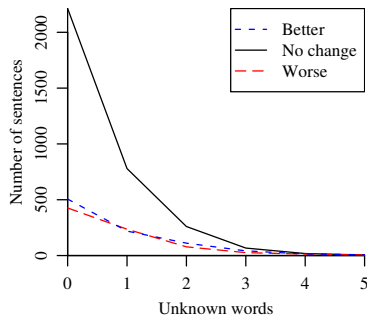
- *Improves performance* from 91.3 to **92.1 f-score**
- Self-training without the reranker does not improve performance
- Retraining the reranker on new first-stage model does not further improve performance
- Would reparsing the NTC with improved parser further improve performance?

## First-stage oracle scores

Model	1-best	10-best	50-best
Baseline	89.0	94.0	95.9
WSJ $\times$ 1 + 250k	89.8	94.6	96.2
WSJ $\times$ 5 + 1,750k	90.4	94.8	96.4

- Self-training improves first-stage generative parser's oracle scores
- First-stage parser also became more decisive: mean of  $\log_2(P(1\text{-best}) / P(50\text{th-best}))$  increased from 11.959 for the baseline parser to 14.104 for self-trained parser

# Which sentences improve?



## Self-trained WSJ parser on Brown

Sentences added	Parser	WSJ-reranker
Baseline Brown	86.4	87.4
Baseline WSJ	83.9	85.8
WSJ+50k	84.8	86.6
WSJ+250k	85.7	87.2
WSJ+1,000k	86.2	87.3
WSJ+2,500k	86.4	87.7

- Adding NTC data greatly improves performance on Brown corpus (to a lesser extent on Switchboard)

## Self-training vs in-domain training

First-stage	First stage alone	WSJ-reranker	Brown-reranker
WSJ	82.9	85.2	85.2
WSJ+NTC	87.1	87.8	87.9
Brown	86.7	88.2	88.4

- Both reranking and self-training are surprisingly domain-independent
- Self-trained NTC parser with WSJ reranker is almost as good as a parser/reranker completely trained on Brown (!)

## Summary and conclusions

- PCFG based parsers are easy to estimate, but sensitive to unmodeled dependencies
- Exponential models are difficult to estimate, but resilient to unmodeled dependencies
- Coarse to fine reranking combines both approaches
- (Re)ranking parsers can work with just about any features
- The details of linguistic representations don't matter so long as they are rich enough to compute your features from
- Self-training works with reranking parsers (why?)
- Both reranking and self-training is (surprisingly) domain-independent

# Outline

Introduction

Non-local dependencies and the PCFG MLE

Generative statistical parsers

Exponential (a.k.a. Maximum Entropy) parsing models

Coarse to fine reranking

Self-training of the reranking parser

**Sample parser errors**



# Sample parser errors

