# Grammar induction in computers and people

Amy Perfors
IPAM 2007

---

## Expressivity vs. learnability

### A fundamental tradeoff

Low expressivity:
easy to learn

High expressivity:
hard to learn

---

## Expressivity vs. learnability

### Maps onto modeling choices

Low expressivity:
easy to learn

High expressivity:
hard to learn

Choice of
representation

---

## Expressivity vs. learnability

### Maps onto modeling choices

Low expressivity:
easy to learn

High expressivity:
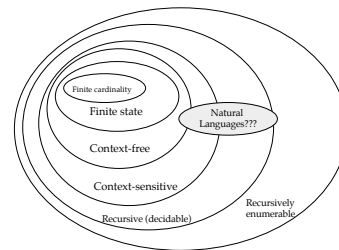hard to learn

Choice of learning
algorithm

---

## Learnability

### Not just about the search algorithm

- Complexity of the space
- Nature of the input
- What constitutes "having learned"
  … as well as the capabilities of the learner

---

## Gold's Theorem

Based on positive evidence only, it is impossible to learn a class of languages other than those of finite cardinality

Finite cardinality

Finite state

Natural Languages???

Context-free

Context-sensitive

Recursively enumerable

Recursive (decidable)

Gold (1967)

---

## Gold's Theorem

Based on positive evidence only, it is impossible to learn a class of languages other than those of finite cardinality

1. Listener "hears" a string

2. Listener decides whether that string is consistent with the grammar they currently think is being used (if it's the first trial, they just pick one randomly).

3. If not, generate a new grammar and go to #2. If so, stay with current grammar and go to #2.

You have learned the language once you never change grammars

Gold (1967)

## Gold's Theorem

Based on positive evidence only, it is impossible to learn a class of languages other than those of finite cardinality

You will eventually converge: the negative evidence will rule out incorrect grammars

Gold (1967)

## Gold's Theorem

Based on positive evidence only, it is impossible to learn a class of languages other than those of finite cardinality

You will eventually converge: the negative evidence will rule out incorrect grammars

If you guessed a grammar that is too large, you'll never realize it and never change

Gold (1967)

## Does this lead to problems?

Kids don't seem to receive (or notice) negative evidence!

Adults tend to only correct the truth of a child's utterance, not the syntax:

Child: Mama isn't boy, he a girl.
Adult: That's right.

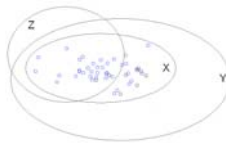When adults (rarely) try to correct a child's syntax, the kid doesn't get it

Child: Nobody don't like me.
Mother: No, say "Nobody likes me."
Child: Nobody don't like me.
   *[dialogue repeated 8 times]*
Mother: Now listen carefully, say "NOBODY LIKES ME."
Child: Oh! Nobody don't likeS me.

Source: McNeil

## Many ways to address this

(for both kids and computers)

- Change learning criterion
- Add some sort of inductive bias
  - Universal Grammar (what does this mean?)
- Change learning algorithm
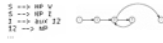  - Probabilistic learner: guaranteed in the limit (Horning, 1969)

Horning (1969), Solomonoff (1978)

## But how about *not* in the limit?

Grammar induction in the real world

## Grammar induction

### Computational question at multiple levels

✗ Grammar type

Me :: ⟵ Specific grammar structure

Mark Johnson ⎱ Parameters of specific grammar

Parsing, prediction of the dataset

## Grammar induction

### Computational question at multiple levels

Grammar type

This talk ⟵ Specific grammar structure

Johnson, Yuille ⎱ Parameters of specific grammar

Parsing, prediction of the dataset

## Plan

### Unsupervised learning of grammars, by complexity

Finite-state grammars
- N-gram models (Goldwater, Griffiths, & Johnson, 2006)
  - Word segmentation and morphology
- HMMs: Bayesian model merging (Stolcke & Omohundro, 1994)
  - Phonetic learning
- HMMs: Dirichlet prior (Goldwater & Griffiths, 2007)
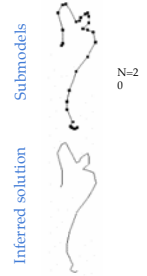  - Syntactic categories

Structured grammars
- Bayesian model merging: PCFGs (Stolcke & Omohundro, 1994)
- Constituent-context model (Klein & Manning, 2002)
- Dependency grammars (Carroll & Charniak, 1992)
- Combined CCM & Dependency grammars (Klein & Manning, 2004)

## Bayesian model merging: HMMs

### Basic idea

- Constructed from subcomponents
  - With only a little data, the subcomponents are the datapoints themselves (plus slight generalization)
  - Construct more complex models by merging pairs of simple components
- Merging…
  - Try to do efficiently by sacrificing as little likelihood as possible each time
  - Put prior so know when to stop

Submodels

N=20

Inferred solution

Stolcke & Omohundro, 1994

## Bayesian model merging: HMMs

### Basic idea

- Constructed from subcomponents
  - With only a little data, the subcomponents are the datapoints themselves (plus slight generalization)
  - Construct more complex models by merging pairs of simple components
- Merging…
  - Try to do efficiently by sacrificing as little likelihood as possible each time
  - Put prior so know when to stop

Submodels

N=17

Inferred solution

Stolcke & Omohundro, 1994

## Bayesian model merging: HMMs

### Basic idea

- Constructed from subcomponents
  - With only a little data, the subcomponents are the datapoints themselves (plus slight generalization)
  - Construct more complex models by merging pairs of simple components
- Merging…
  - Try to do efficiently by sacrificing as little likelihood as possible each time
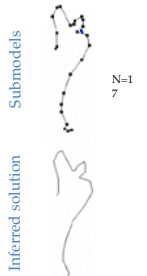  - Put prior so know when to stop

Submodels

N=14

Inferred solution

Stolcke & Omohundro, 1994

# Bayesian model merging: HMMs

## Basic idea

- Constructed from subcomponents
  - With only a little data, the submodels are the datapoints themselves (plus slight generalization)
  - Construct more complex models by merging pairs of simple components
- Merging…
  - Try to do efficiently by sacrificing as little likelihood as possible each time
  - Put prior so know when to stop
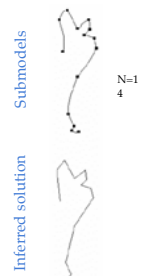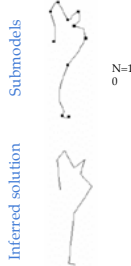
Submodels

N=10

Inferred solution

---

# Bayesian model merging: HMMs

## Notation

Set of states $Q$

Initial state $q_I$, Final state $q_F$

$$q_I \xrightarrow{\sigma_1} q_1 \xrightarrow{\sigma_1} \cdots \to q_F$$

Set of probability parameters:

Transition probabilities $p(q \to q')$: the probability that state $q'$ follows $q$

Emission probabilities $p(q \mid \sigma)$: the probability that symbol $\sigma$ is emitted when in state $q$

---

# Bayesian model merging: HMMs

## Model details

Likelihood: probability of a string $x$ is the sum of the probabilities of all paths that generate $x$

$$P(x|M) = \sum_{q_1 \ldots q_l \in Q^l} p(q_I \to q_1) p(x_1 \mid q_1) p(q_1 \to q_2) \ldots p(q_l \mid q_l) p(q_l \to q_F)$$

Prior: structural vs. parameter

Parameters: Dirichlet distribution over each set of multinomial parameters (transition and emission probabilities)

$$P(\theta_M^{(q)} \mid M_G, M_S^{(q)}) = \frac{1}{B(\alpha_t, \ldots, \alpha_t)} \prod_{i=1}^{n_t^{(q)}} \theta_{qt}^{\alpha_{qt}-1} \frac{1}{B(\alpha_e, \ldots, \alpha_e)} \prod_{j=1}^{n_e^{(q)}} \theta_{qj}^{\alpha_{qe}-1}$$

Structural: Favor smaller number of states Q

$$P(M_S) \propto C^{-|Q|}$$

---

# Bayesian model merging: HMMs

## Algorithm

1) Construct an HMM $M_0$ that produces exactly the input strings

2) Loop:
   - Compute set of candidate merges K among the states of current HMM, $M_i$
   - For each candidate k, compute the merged model $k(M_i)$, and its posterior probability $p(k(M_i) \mid X)$
   - Let $k^*$ be the merge that maximizes $p(k \mid M_i) \mid X)$. Then let $M_{i+1} = k^*(M_i)$
   - If $P(M_{i+1} \mid X) < P(M_i \mid X)$, return $M_i$ as the induced model

Merges combine pairs of states: combined emission and transition probabilities are weighted averages of the corresponding distributions for the states that have been merged
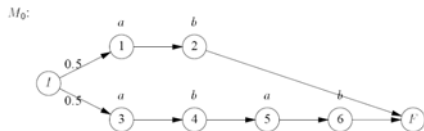
---

# Bayesian model merging: HMMs

## Example

Language: $\{ab\}^+$

Sentences: $ab$, $abab$

$M_0$:

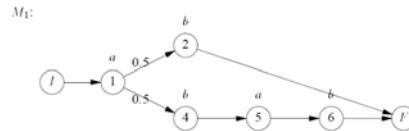---

# Bayesian model merging: HMMs
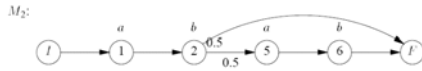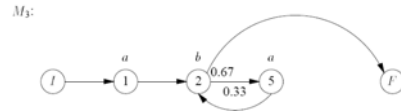
## Example

Language: $\{ab\}_+$

Sentences: $ab$, $abab$

$M_1$:

## Bayesian model merging: HMMs

### Example

Language: *{ab}+*

Sentences: *ab, abab*

$M_2$:

---

## Bayesian model merging: HMMs

### Example

Language: *{ab}+*

Sentences: *ab, abab*

$M_3$:

---

## Bayesian model merging: HMMs

### Example

Language: *{ab}+*

Sentences: *ab, abab*

$M_4$:

---

## Bayesian model merging: HMMs

### Why are simpler models favored?

Occam factor: Models with fewer states have fewer parameters, and therefore they make tighter predictions.

More effective data: When two states are merged, their probabilities are combined – so each state gets to effectively "see" more data

---

## Bayesian model merging: HMMs

### Application: Phonetic data

TIMIT: Collection of hand-labeled speech samples compiled for the purpose of training speaker-independent phonetic recognition systems

Contains acoustic data segmented by words and aligned with discrete labels from an alphabet of 62 phones

Goal is to construct a probabilistic model for each word in the database, representing its phonetic structure as accurately as possible

---

## Bayesian model merging: HMMs

### Application: Phonetic data

Initial HMM

Realizations of the word "often"

## Bayesian model merging: HMMs

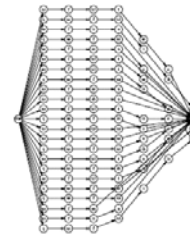### Application: Phonetic data

Learned HMM

Realizations of the word "often"



Figure: Merged HMM constructed from 37 samples of the word *often*. Merging was constrained to keep the emission on each state unique.

---

## More recent models of HMM learning

### Application: Part-of-speech (POS) tagging

---

## More recent models of HMM learning

### Application: Part-of-speech (POS) tagging

Identifies a distribution over latent variables directly, without ever fixing particular values for the model parameters:

$$P(\mathbf{t}|\mathbf{w}) = \int P(\mathbf{t}|\mathbf{w},\theta)P(\theta|\mathbf{w})d\theta$$

$w$ = the linguistic input; $\theta$ = parameters; $t$ = the hidden structure

Symmetric Dirichlet prior over the transition and output distributions

$$
\begin{aligned}
t_i | t_{i-1} = t, t_{i-2} = t', \tau^{(t,t')} &\sim \text{Mult}(\tau^{(t,t')}) \\
w_i | t_i = t, \omega^{(t)} &\sim \text{Mult}(\omega^{(t)}) \\
\tau^{(t,t')} | \alpha &\sim \text{Dirichlet}(\alpha) \\
\omega^{(t)} | \beta &\sim \text{Dirichlet}(\beta)
\end{aligned}
$$

---

## More recent models of HMM learning

### POS tagging: Results

**Variation of information (VI):** The VI between two clusterings C (the gold standard) and C' (the found clustering) of a set of data points is a sum of the amount of information lost in moving from C to C'
$$VI(C,C') = H(C) + H(C') - 2I(C,C')$$
(High VI indicates different clustering relative to C)

**Vary the tag information given to the model:** Contains tag information only for words that appear at least $d$ times in the training corpus (the first 24K words of the WSJ corpus)

---

## More recent models of HMM learning

### POS tagging: Results

Value of *d* → tag information

| Accuracy | 1 | 2 | 3 | 5 | 10 | ∞ |
|---|---|---|---|---|---|---|
| random | 69.6 | 56.7 | 51.0 | 45.2 | 38.6 | |
| MLHMM | 83.2 | 70.6 | 65.5 | 59.0 | 50.9 | |
| BHMM1 | 86.0 | 76.4 | 71.0 | 64.3 | 58.0 | |
| BHMM2 | 87.3 | 79.6 | 65.0 | 59.2 | 49.7 | |
| σ < | .2 | .8 | .6 | .3 | 1.4 | |
| **VI** | | | | | | |
| random | 2.65 | 3.96 | 4.38 | 4.75 | 5.13 | 7.29 |
| MLHMM | 1.13 | 2.51 | 3.00 | 3.41 | 3.89 | 6.50 |
| BHMM1 | 1.09 | 2.44 | 2.82 | 3.19 | 3.47 | 4.30 |
| BHMM2 | **1.04** | **1.78** | **2.31** | **2.49** | **2.97** | **4.04** |
| σ < | .02 | .03 | .04 | .03 | .07 | .17 |
| **Corpus stats** | | | | | | |
| % ambig. | 49.0 | 61.3 | 66.3 | 70.9 | 75.8 | 100 |
| tags/token | 1.9 | 4.4 | 5.5 | 6.8 | 8.3 | 17 |

---

## Plan

### Unsupervised learning of grammars, by complexity

Finite-state grammars
- N-gram models (Goldwater, Griffiths, & Johnson, 2006)
  - Word segmentation and morphology
- HMMs: Bayesian model merging (Stolcke & Omohundro, 1994)
  - Phonetic learning
- HMMs: Dirichlet prior (Goldwater & Griffiths, 2007)
  - Syntactic categories

Structured grammars
- Bayesian model merging: PCFGs (Stolcke & Omohundro, 1994)
- Constituent-context model (Klein & Manning, 2002)
- Dependency grammars (Carroll & Charniak, 1992)
- Combined CCM & Dependency grammars (Klein & Manning, 2004)

## Bayesian model merging: PCFGs

### Notation

Set of nonterminal symbols $N$
Set of terminal symbols
Start nonterminal $S$
Set of productions or rules $R$
Production probabilities $p(r)$ for all rules $r$

| | |
|---|---|
| 0.6 | S --> NP VP |
| 0.4 | S --> NP I |
| 0.3 | NP --> N |
| 0.4 | NP --> D N |
| 0.3 | NP --> Pro |

---

## Bayesian model merging: PCFGs

### Merging operator

Replaces two existing nonterminals $X_1$ and $X_2$ with a single new nonterminal Y.
Has a twofold effect on grammars:

1) RHS occurrences of $X_1$ and $X_2$ replaced by Y

$$Z_1 \rightarrow \lambda_1 X_1 \mu_1 \quad (c_1)$$
$$Z_1 \rightarrow \lambda_2 X_2 \mu_2 \quad (c_2)$$
$$\Downarrow \quad \text{merge}(X_1, X_2) = Y$$
$$Z_1 \rightarrow \lambda_1 Y \mu_1 \quad (c_1)$$
$$Z_2 \rightarrow \lambda_2 Y \mu_2 \quad (c_2)$$

2) The union of LHSs $X_1$ and $X_2$ replaced by Y

$$X_1 \rightarrow \lambda_1 \quad (c_1)$$
$$X_2 \rightarrow \lambda_2 \quad (c_2)$$
$$\Downarrow \quad \text{merge}(X_1, X_2) = Y$$
$$Y \rightarrow \lambda_1 \quad (c_1)$$
$$\rightarrow \lambda_2 \quad (c_2)$$

---

## Bayesian model merging: PCFGs

### Nonterminal chunking

Merging alone cannot create CFG productions with the usual embedding structure, so we add a chunking operator as well

Given an ordered sequence of nonterminals $X_1 X_2 \ldots X_k$, create a new nonterminal Y that expands to $X_1 X_2 \ldots X_k$, and replace occurrences of $X_1 X_2 \ldots X_k$ on the RHS with Y

$$Z \rightarrow \lambda X_1 X_2 \ldots X_k \mu \quad (c)$$
$$\Downarrow \quad \text{chunk}(X_1 X_2 \ldots X_k) = Y$$
$$Z \rightarrow \lambda Y \mu \quad (c)$$
$$Y \rightarrow X_1 X_2 \ldots X_k \quad (c')$$

---

## Bayesian model merging: PCFGs

### More model details

**Likelihood**: probability of a string $x$ is the sum of the probabilities of all paths that generate $x$

$$P(x|M) = \sum_{x_1 \ldots x_l \in \mathcal{L}^*} p(q_0 \rightarrow q_1) p(q_1 \uparrow x_1) p(q_1 \rightarrow q_2) \ldots p(q_l \uparrow x_l) p(q_l \rightarrow q_f)$$

**Prior**:

Parameters on productions: Dirichlet distribution over the multinomial describing each nonterminal

Structural: Description length (e.g., length of each production assumed to have been drawn from a Poisson distribution)

---

## Bayesian model merging: PCFGs

### Performance: Induction is *very* difficult

Can try to improve the greedy search in a number of ways:

It still isn't very good – the size of the space is enormous, and has many local minima

Horning (1969) – suggested an algorithm, but it was an enumeration algorithm

Structure induction for PCFGs is a very difficult problem

---

## One possible solution: induce constituency more directly

## Constituent-Context Model (CCM)

Explicitly models constituent yields and contexts

System exploits the fact that long constituents often have short, common equivalents (proforms) that appear in similar context and whose constituency is similar (or guaranteed)
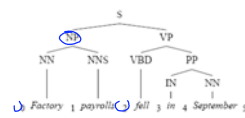
Pronoun aux adj

Det n aux adj

Det n prep det n aux adj

Model transfers the constituency of a sequence directly to its containing context, which is intended to then pressure new sequences that occur in that context into being parsed as constituents in the next round

## CCM: Representation

Model captures all *spans* (contiguous subsequences) of a sentence, and each span's *context* (the terminals preceding and following it
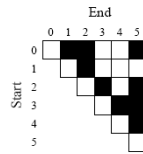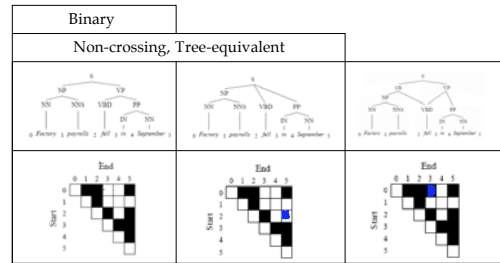
## CCM: Representation

The *bracketing* B of a sentence indicates which spans are constituents

## CCM: Representation

Different kinds of bracketings:

## CCM: Generative Model
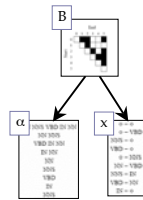
Choose a bracketing B according to some distribution P(B) and then generate the sentence S given B

The context and yield (content) of each span are independent of each other, and generated conditionally on the constituency of that span

$$P(S|B) = \prod_{\langle i,j \rangle \in spans(S)} P(\alpha_{ij}, x_{ij} | B_{ij})$$

$$= \prod_{\langle i,j \rangle} P(\alpha_{ij} | B_{ij}) P(x_{ij} | B_{ij})$$

## CCM: Generative Model

Choose a bracketing B according to some distribution P(B) and then generate the sentence S given B

P(B) = uniform

Soft clustering with two equal-prior classes (constituents & distituents)

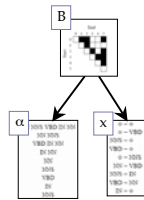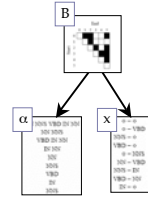## CCM: Generative Model

Choose a bracketing B according to some distribution P(B) and then generate the sentence S given B

P(B) = uniform over binary tree-equivalent bracketings

This turns distributional clustering into tree induction

---

## CCM: Inducing structure

EM: Sentences $S$ are observed, bracketings $B$ are unobserved. Parameters of the model are the constituency-conditional yield and context distributions ($P(\alpha|b)$ and $P(x|b)$, respectively)

**E-Step:** Find the conditional completion likelihoods $P(B|S, \Theta)$ according to the current $\Theta$.

**M-Step:** Fix $P(B|S, \Theta)$ and find the $\Theta'$ which maximizes $\sum_B P(B|S, \Theta) \log P(S, B|\Theta')$.

---

## CCM: Results

Corpus: 7422 sentences, Penn treebank Wall Street Journal that contains no more than 10 words

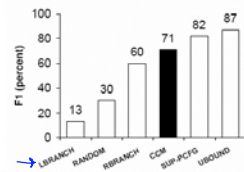Evaluation: F1 (harmonic mean of precision and recall)



Figure 4: $F_1$ for various models on WSJ-10.

---

## CCM: Extension

What if the model isn't given the POS tags?

Used the baseline method of word-type clustering (similar to Finch et. Al. (1993))

Performance is worse, but still better than next best (right-branching)

| Tags | Precision | Recall | $F_1$ | NP Recall | PP Recall | VP Recall | S Recall |
|---|---|---|---|---|---|---|---|
| Treebank | 63.8 | 80.2 | 71.1 | 83.4 | 78.5 | 78.6 | 40.7 |
| Induced | 56.8 | 71.1 | 63.2 | 52.8 | 56.2 | 90.0 | 60.5 |

---

## CCM: Generative Model with multiple classes

Choose a bracketing B according to some distribution P(B) and then generate the sentence S given B

Labeling is generated from that bracketing

The context and yield (content) of each span are independent of each other, and generated conditionally on the constituency of that span

---

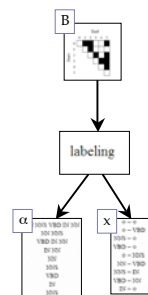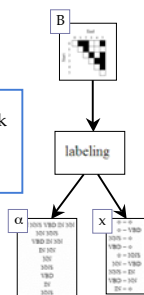## CCM: Generative Model with multiple classes

Choose a bracketing B according to some distribution P(B) and then generate the sentence S given B

Labeling is generated from that bracketing

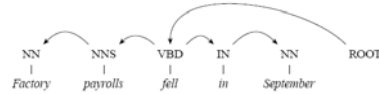This doesn't actually work that well: F1 = 70.9 (compared to 71.1)

The context and yield (content) of each span are independent of each other, and generated conditionally on the constituency of that span

## Another solution: use a different representation

## Dependency grammars



A *dependency* d is an arc: pair <h,a> of a head and argument, each of which is a word in a sentence s
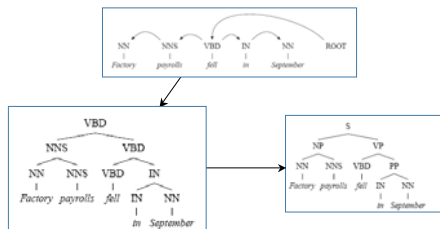
A *dependency structure* D over a sentence is a set of dependencies that form a planar, acyclic graph rooted at ROOT

Every D has an associated graph G:



## Dependency grammars

Isomorphic (in terms of strong generative capacity) to a restricted form of phrase structure grammar (Miller, 1999)



## Unsupervised Dependency Parsing

### Why use dependency grammars?

1) Most state-of-the-art supervised parsers make use of specific lexical information in addition to word-class level information: perhaps lexical information could be a useful source of information for unsupervised models
2) A central motivation for using tree structures is to enable the extraction of dependencies, and it might be more advantageous to do so directly
3) For languages like Chinese, which have few function words, and for which the definition of lexical categories is much less clear, dependency structures may be easier to detect

## Inducing dependency grammars (DEP-PCFG)

### Algorithm

0) Divide the corpus into two parts, the rule corpus and the training corpus
1) For all sentences in the rule corpus, generate all rules which might be used to generate (and/or parse) the sentence
2) Estimate the probabilities for the rules
3) Using the training corpus, improve our estimate of the probabilities
4) Delete all rules with small enough probability. What remains is the grammar

## Inducing dependency grammars (DEP-PCFG)

### Difficulties

1) There are way too many possible (CFG) rules that could lead to a sentence: a potentially infinite set of nonterminals
   - That's why we use a dependency grammar: it limits it to $n(2^{n-1}+1)$, where $n$ is length of sentence, if all terminals are distinct

2) Even with a dependency grammar, this is a LOT of sentences. For instance, a sentence with 41 terminal symbols would have
$$(41(2^{40}+1) \approx 41((2^{10})^4) \approx 40((10^3)^4) \approx 4 \cdot 10^{13}$$

3) Deal with this by ordering sentences by length (since children see simpler language before more complex language). Once a rule has been eliminated, don't consider it again

## DEP-PCFG: Experiment

### Corpus generated by artificial grammar

| | | | | | | |
|---|---|---|---|---|---|---|
| 1.0 | S | → | $\cdot$ | 1.0 | $\cdot$ | → | $\overline{verb}$ . |
| 1.0 | $\overline{pron}$ | → | pron | .1 | $\overline{noun}$ | → | noun |
| .3 | $\overline{noun}$ | → | det noun | .1 | $\overline{noun}$ | → | $\overline{det\ adj}$ noun |
| .2 | $\overline{noun}$ | → | $\overline{det\ noun\ prep}$ | .2 | $\overline{noun}$ | → | $\overline{det\ noun}$ wh |
| .05 | $\overline{noun}$ | → | noun $\overline{prep}$ | .05 | $\overline{noun}$ | → | noun wh |
| 1.0 | $\overline{adj}$ | → | adj | 1.0 | $\overline{det}$ | → | det |
| 1.0 | $\overline{wh}$ | → | wh $\overline{verb}$ | | | | |
| .3 | $\overline{prep}$ | → | prep $\overline{pron}$ | .7 | $\overline{prep}$ | → | prep $\overline{noun}$ |
| .05 | $\overline{verb}$ | → | $\overline{noun}$ verb | .1 | $\overline{verb}$ | → | verb |
| .05 | $\overline{verb}$ | → | verb $\overline{noun}$ | .05 | $\overline{verb}$ | → | $\overline{pron}$ verb |
| .1 | $\overline{verb}$ | → | $\overline{noun}$ verb $\overline{noun}$ | .05 | $\overline{verb}$ | → | verb $\overline{pron}$ |
| .05 | $\overline{verb}$ | → | $\overline{noun}$ verb $\overline{pron}$ | .05 | $\overline{verb}$ | → | $\overline{pron}$ verb $\overline{noun\ noun}$ |
| .05 | $\overline{verb}$ | → | $\overline{noun}$ verb $\overline{noun\ noun}$ | .1 | $\overline{verb}$ | → | $\overline{noun}$ verb $\overline{noun\ prep}$ |
| .1 | $\overline{verb}$ | → | $\overline{pron}$ verb $\overline{noun\ prep}$ | .1 | $\overline{verb}$ | → | $\overline{pron}$ verb $\overline{pron\ prep}$ |
| .05 | $\overline{verb}$ | → | $\overline{noun}$ verb prep | .05 | $\overline{verb}$ | → | $\overline{pron}$ verb prep |

---

## DEP-PCFG: Results

### "Uniformly awful"

With 300 different random starting points, ended up with 300 different local minima

One segment of one grammar:

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| .220 | $\overline{pron}$ | → | pron $\overline{verb}$ | .117 | $\overline{pron}$ | → | $\overline{det\ verb}$ pron |
| .214 | $\overline{pron}$ | → | $\overline{prep}$ pron | .038 | $\overline{pron}$ | → | pron $\overline{verb\ noun}$ |
| .139 | $\overline{pron}$ | → | pron $\overline{verb\ det}$ | .023 | $\overline{pron}$ | → | $\overline{noun\ verb}$ pron |
| .118 | $\overline{pron}$ | → | $\overline{verb}$ pron | .013 | $\overline{pron}$ | → | pron $\overline{verb\ det\ det}$ |

---

## DEP-PCFG: Results

### "Uniformly awful"

Tried to make it better by giving the grammar a chart limiting what possible nonterminals could appear on the right-hand side of certain rules

| lhs | noun | verb | pron | det | prep | adj | wh | . |
|---|---|---|---|---|---|---|---|---|
| noun | o | o | o | O | O | O | O | o |
| verb | O | o | o | o | o | o | o | o |
| pron | o | x | o | o | o | o | o | o |
| det | x | o | o | o | o | x | o | o |
| prep | O | o | O | o | o | o | o | o |
| adj | x | o | o | x | o | o | o | o |
| wh | x | O | o | o | o | o | o | o |
| . | o | O | o | o | o | o | o | o |

Figure 9: For each left-hand side, non-terminals allowed on right

It helped, but they didn't really quantify how much

---

## DEP-PCFG: Results

### Why didn't this work well?

1) Random initialization is not good because the parameter space is riddled with local minima; try it with a uniform initialization
2) Dependency grammars are structurally unable to distinguish order information (e.g., whether subject or object should be attached to the verb first) – and therefore they are not that great at constituency
3) It did not encode valence: for instance, can learn that nouns to the left of a verb attach to it… but if given a NOUN NOUN VERB sequence, will think that both nouns attach to the verb. Cannot encode that verbs have exactly one subject.

---

## DMV: The generative model

### Includes a model of valence

Recursive generation process, beginning at the head $h$, defined according to the following rules:

For each direction L and R from $h$
- With P(STOP|$h,dir,adj$) stop generating additional arguments
- If don't stop, generate a new argument a with P($a$|$h,dir$)P(D($a$))
→ These probability factors are the model's parameters

---

## DMV: The generative model

Recursive generation process, beginning at the head $h$, defined according to the following rules:

For each direction L and R from $h$
- With P(STOP|$h,dir,adj$) stop generating additional arguments
- If don't stop, generate a new argument a with P($a$|$h,dir$)P(D($a$))
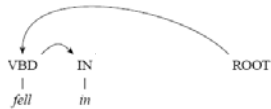→ These probability factors are the model's parameters

## DMV: The generative model

Recursive generation process, beginning at the head *h*, defined according to the following rules:

For each direction L and R from *h*
- With P(STOP | *h,dir,adj*) stop generating additional arguments
- If don't stop, generate a new argument a with P(*a* | *h,dir*)P(D(*a*))
→ These probability factors are the model's parameters

## DMV: The generative model

Recursive generation process, beginning at the head *h*, defined according to the following rules:

For each direction L and R from *h*
- With P(STOP | *h,dir,adj*) stop generating additional arguments
- If don't stop, generate a new argument a with P(*a* | *h,dir*)P(D(*a*))
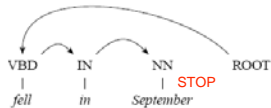→ These probability factors are the model's parameters

## DMV: The generative model

Recursive generation process, beginning at the head *h*, defined according to the following rules:

For each direction L and R from *h*
- With P(STOP | *h,dir,adj*) stop generating additional arguments
- If don't stop, generate a new argument a with P(*a* | *h,dir*)P(D(*a*))
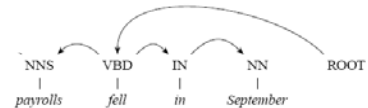→ These probability factors are the model's parameters

## DMV: The generative model

Recursive generation process, beginning at the head *h*, defined according to the following rules:

For each direction L and R from *h*
- With P(STOP | *h,dir,adj*) stop generating additional arguments
- If don't stop, generate a new argument a with P(*a* | *h,dir*)P(D(*a*))
→ These probability factors are the model's parameters

## DMV: The generative model

Recursive generation process, beginning at the head *h*, defined according to the following rules:

For each direction L and R from *h*
- With P(STOP | *h,dir,adj*) stop generating additional arguments
- If don't stop, generate a new argument a with P(*a* | *h,dir*)P(D(*a*))
→ These probability factors are the model's parameters

## DMV: The generative model

Recursive generation process, beginning at the head *h*, defined according to the following rules:

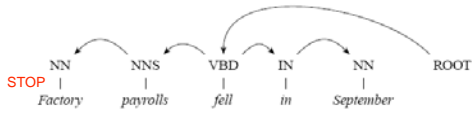For each direction L and R from *h*
- With P(STOP | *h,dir,adj*) stop generating additional arguments
- If don't stop, generate a new argument *a* with P(*a* | *h,dir*)P(D(*a*))
→ These probability factors are the model's parameters

## DMV: The generative model

Use inside-outside algorithm for reestimation of the parameters (the probability factors for STOP and CHOOSE)



STOP
| NN | NNS | VBD | IN | NN | ROOT |
| Factory | payrolls | fell | in | September | |

## DMV: Experiment

| Model | UP | UR | UF$_1$ | Dir | Undir |
|---|---|---|---|---|---|
| English (WSJ10 – 7422 Sentences) | | | | | |
| LBRANCH/RHEAD | 25.6 | 32.6 | 28.7 | 33.6 | 56.7 |
| RANDOM | 31.0 | 39.4 | 34.7 | 30.1 | 45.6 |
| RBRANCH/LHEAD | 55.1 | 70.0 | 61.7 | 24.0 | 55.9 |
| DMV | 46.6 | 59.2 | 52.1 | 43.2 | 62.7 |
| CCM | 64.2 | 81.6 | 71.9 | 23.8 | 43.3 |
| UBOUND | 78.8 | 100.0 | 88.1 | 100.0 | 100.0 |
| German (NEGRA10 – 2175 Sentences) | | | | | |
| LBRANCH/RHEAD | 27.4 | 48.8 | 35.1 | 32.6 | 51.2 |
| RANDOM | 27.9 | 49.6 | 35.7 | 21.8 | 41.5 |
| RBRANCH/LHEAD | 33.8 | 60.1 | 43.3 | 21.0 | 49.9 |
| DMV | 38.4 | 69.5 | 49.5 | 40.0 | 57.8 |
| CCM | 48.1 | 85.5 | 61.6 | 25.5 | 44.9 |
| UBOUND | 56.3 | 100.0 | 72.1 | 100.0 | 100.0 |
| Chinese (CTB10 – 2437 Sentences) | | | | | |
| LBRANCH/RHEAD | 26.3 | 48.8 | 34.2 | 30.2 | 43.9 |
| RANDOM | 27.3 | 50.7 | 35.5 | 35.9 | 47.3 |
| RBRANCH/LHEAD | 29.0 | 53.9 | 37.8 | 14.2 | 41.5 |
| DMV | 35.9 | 66.7 | 46.7 | 42.5 | 54.2 |
| CCM | 34.6 | 64.3 | 45.0 | 23.8 | 40.5 |
| UBOUND | 53.9 | 100.0 | 70.1 | 100.0 | 100.0 |

## Why not combine models?

Their strengths are complementary

Both CCM and DMV can be seen as models over lexicalized trees

Combine them by scoring each tree with the product of all the probabilities from the individual models
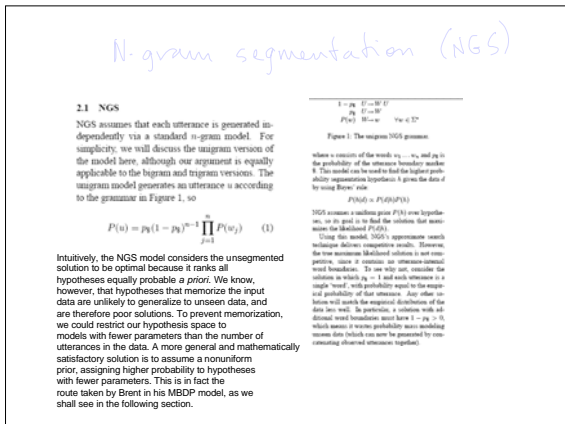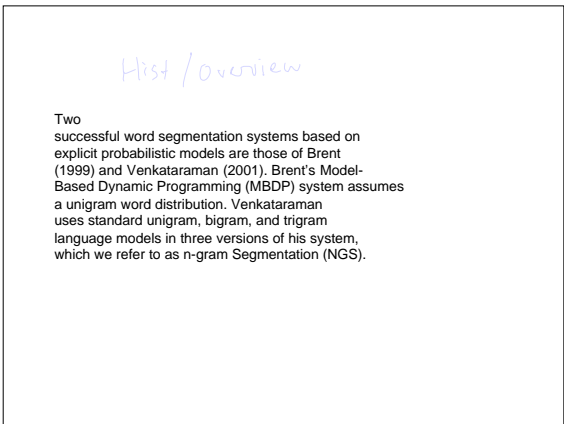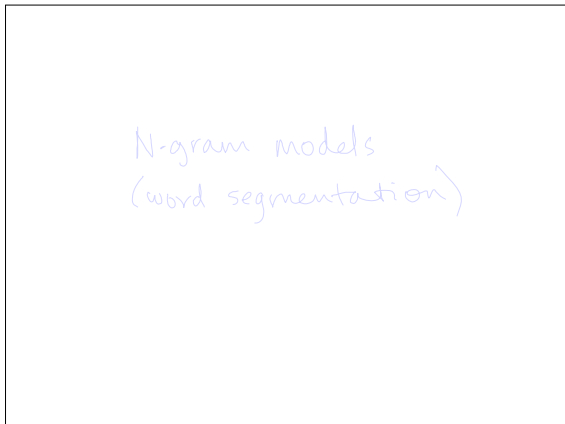
## DMV-CCM: Results

| Model | UP | UR | UF$_1$ | Dir | Undir |
|---|---|---|---|---|---|
| English (WSJ10 – 7422 Sentences) | | | | | |
| LBRANCH/RHEAD | 25.6 | 32.6 | 28.7 | 33.6 | 56.7 |
| RANDOM | 31.0 | 39.4 | 34.7 | 30.1 | 45.6 |
| RBRANCH/LHEAD | 55.1 | 70.0 | 61.7 | 24.0 | 55.9 |
| DMV | 46.6 | 59.2 | 52.1 | 43.2 | 62.7 |
| CCM | 64.2 | 81.6 | 71.9 | 23.8 | 43.3 |
| DMV+CCM (POS) | 69.3 | 88.0 | 77.6 | 47.5 | 64.5 |
| DMV+CCM (DISTR.) | 65.2 | 82.8 | 72.9 | 42.3 | 60.4 |
| UBOUND | 78.8 | 100.0 | 88.1 | 100.0 | 100.0 |
| German (NEGRA10 – 2175 Sentences) | | | | | |
| LBRANCH/RHEAD | 27.4 | 48.8 | 35.1 | 32.6 | 51.2 |
| RANDOM | 27.9 | 49.6 | 35.7 | 21.8 | 41.5 |
| RBRANCH/LHEAD | 33.8 | 60.1 | 43.3 | 21.0 | 49.9 |
| DMV | 38.4 | 69.5 | 49.5 | 40.0 | 57.8 |
| CCM | 48.1 | 85.5 | 61.6 | 25.5 | 44.9 |
| DMV+CCM | 49.6 | 89.7 | 63.9 | 50.6 | 64.7 |
| UBOUND | 56.3 | 100.0 | 72.1 | 100.0 | 100.0 |
| Chinese (CTB10 – 2437 Sentences) | | | | | |
| LBRANCH/RHEAD | 26.3 | 48.8 | 34.2 | 30.2 | 43.9 |
| RANDOM | 27.3 | 50.7 | 35.5 | 35.9 | 47.3 |
| RBRANCH/LHEAD | 29.0 | 53.9 | 37.8 | 14.2 | 41.5 |
| DMV | 35.9 | 66.7 | 46.7 | 42.5 | 54.2 |
| CCM | 34.6 | 64.3 | 45.0 | 23.8 | 40.5 |
| DMV+CCM | 33.3 | 62.0 | 43.3 | 55.2 | 60.3 |
| UBOUND | 53.9 | 100.0 | 70.1 | 100.0 | 100.0 |

## Conclusion

### Grammar induction (of structure) is difficult

- Expressivity/learnability tradeoff
- Finite-state grammars are easier, and there are some useful linguistic domains that they are reasonable models for
  - Word segmentation
  - Phonetics
  - Syntactic categories
- More structured grammars are difficult
  - As we saw in Mark Johnson's talk, PCFGs aren't great models for language, but even they are quite hard
  - Constituent-context model
  - Dependency grammars

N-gram models
(word segmentation)

Hist / Overview

Two
successful word segmentation systems based on
explicit probabilistic models are those of Brent
(1999) and Venkataraman (2001). Brent's Model-
Based Dynamic Programming (MBDP) system assumes
a unigram word distribution. Venkataraman
uses standard unigram, bigram, and trigram
language models in three versions of his system,
which we refer to as n-gram Segmentation (NGS).

N-gram segmentation (NGS)

### 2.1 NGS

NGS assumes that each utterance is generated in-
dependently via a standard $n$-gram model. For
simplicity, we will discuss the unigram version of
the model here, although our argument is equally
applicable to the bigram and trigram versions. The
unigram model generates an utterance $u$ according
to the grammar in Figure 1, so

$$P(u) = p_\$ (1 - p_\$)^{n-1} \prod_{j=1}^{n} P(w_j) \quad (1)$$

Intuitively, the NGS model considers the unsegmented
solution to be optimal because it ranks all
hypotheses equally probable *a priori*. We know,
however, that hypotheses that memorize the input
data are unlikely to generalize to unseen data, and
are therefore poor solutions. To prevent memorization,
we could restrict our hypothesis space to
models with fewer parameters than the number of
utterances in the data. A more general and mathematically
satisfactory solution is to assume a nonuniform
prior, assigning higher probability to hypotheses
with fewer parameters. This is in fact the
route taken by Brent in his MBDP model, as we
shall see in the following section.

Figure 1: The unigram NGS grammar.

## Model-based dynamic programming (Brent) MBDP

MBDP assumes a corpus of utterances is generated as a single probabilistic event with four steps:
1. Generate L, the number of lexical types.
2. Generate a phonemic representation for each type (except the utterance boundary type, $).
3. Generate a token frequency for each type.
4. Generate an ordering for the set of tokens.
In a final deterministic step, the ordered tokens are concatenated to create an unsegmented corpus. This means that certain segmented corpora will produce the observed data with probability 1, and all others will produce it with probability 0. The posterior probability of a segmentation given the data is thus proportional to its prior probability under the generative model, and the best segmentation is that with the highest prior probability.
There are two important points to note about the MBDP model. First, the distribution over L assigns higher probability to models with fewer lexical items. We have argued that this is necessary to avoid memorization, and indeed the unsegmented corpus is not the optimal solution under this model, as we will show in Section 3. Second, the factorization into four separate steps makes it theoretically possible to modify each step independently in order to investigate the effects of the various modeling assumptions. However, the mathematical statement of the model and the approximations necessary for the search procedure make it unclear how to modify the model in any interesting way. In particular, the fourth step uses a uniform distribution, which creates a unigram constraint that cannot easily be changed.

*Prob don't want to go over this— unne confusing— suffice to know Brent assumes prior (as they do) thus solves prob from doing MLE*

## Dirichlet model

[dense academic paper text, unreadable at this resolution]

Second, the parameter alpha0 can be used to control how sparse the solutions found by the model are. This parameter determines the total probability of generating *any* novel word, a probability that decreases as more data is observed, but never disappears. Finally, the parameter P0 can be used to encode expectations about the nature of the lexicon, since it defines a probability distribution across different novel words. The fact that this distribution is defined separately from the distribution on word frequencies gives the model additional flexibility, since either distribution can be modified independently of the other.

## Accuracy of unigram model

In Table 1(a), we compare the results of our system to those of MBDP and NGS.4 Although our system has higher lexicon accuracy than the others, its token accuracy is much worse. This result occurs because our system often mis-analyzes frequently occurring words. In particular, many of these words occur in common collocations such as *what's that* and *do you*, which the system interprets as a single words. It turns out that a full 31% of the proposed lexicon and nearly 30% of tokens consist of these kinds of errors.
Upon reflection, it is not surprising that a unigram language model would segment words in this way. Collocations violate the unigram assumption in the model, since they exhibit strong word-toword dependencies. The only way the model can capture these dependencies is by assuming that these collocations are in fact words themselves. Why don't the MBDP and NGS unigram models exhibit these problems? We have already shown that NGS's results are due to its search procedure rather than its model. The same turns out to be true for MBDP. Table 2 shows the probabilities under each model of various segmentations of the corpus. From these figures, we can see that the MBDP model assigns higher probability to the solution found by our Gibbs sampler than to the solution found by Brent's own incremental search algorithm. In other words, Brent's model *does* prefer the lower-accuracy collocation solution, but his search algorithm instead finds a higher-accuracy but lower-probability solution.

## Bigram model

[dense academic paper text with equations, unreadable at this resolution]

---

*Then he looks at grammar induction*
*- well, sorta: prior on θ makes it sort of like also choosing prods*

*His slides also talk about unigram/bigram word segmentation, though I don't remember it*

*He also has some slides about non-local dependencies (e.g. adding features to PCFG)*

- *Introduced HMMs, PCFGs, formalizn (a bit on Chomsky hierarchy but not much)*
- *Talked about how to parse sentences w/ PCFGs or use to predict sentences*
  - *and how to estimate params from the data (EM, inside-outside)*
- *Unsupervised inference for θ (and thus parse trees)*
  - *Bayesian prior on θ*
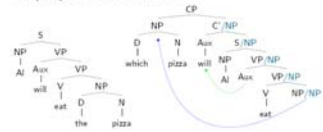  - *in a way is learning grammar (as in sesotho example) but still very constrained*

- *Did BP an... -Talked abo... of hidden struct...*

Slide 1:

# Problems with PCFGs

Slide 2:

*from Johnson*

## Nonlocal "movement" constructions

"Movement" constructions involve a phrase appearing far from its normal location.

Linguists believed CFGs could not generate them, and posited "movement transformations" to produce them (Chomsky 1957)

But CFGs can generate them via "feature passing" using a conspiracy of rules (Gazdar 1984)



Slide 3:

## Learning PCFGs using Expectation Maximization

- ATIS treebank consists of 1,300 hand-constructed parse trees
- input consists of POS tags rather than words
- about 1,000 PCFG rules are needed to build these trees

## Probability of training strings

### Why didn't it learn the right grammar?

- higher likelihood ≠ parse accuracy
  - *probabilistic model and/or estimation procedure are wrong*
- Bayesian prior preferring smaller grammars doesn't h
- What could be wrong?
  - Wrong model of grammar (Klein and Manning)
  - Wrong estimation procedure (Smith and Eisner)
  - Wrong training data (Yang)
  - Predicting word strings is wrong objective
  - Grammar *ignores semantics* (Zettlemoyer and Colli
- de Marken (1995) "Lexical heads, phrase structure and the indu grammar"

## Accuracy of parses

## The PCF

- Parse a
  - higher likelihood ≠ better parses
  - *the statistical model is wrong*
- Initialized EM with correct parse trees
  - started with true rules and their probabilities
    ⇒ poor performance not due to search error
- Evaluated on training data
  - poor performance not due to over-learning

Slide 4:

*Johnson's slides on unigram/bigram seg*

## Unigram model of word segmentation

- Unigram model: each word is generated independently
- Input is *unsegmented broad phonemic transcription* (Brent)
  Example y u w a n t t u s i D 6 b u k
- CRP for Word non-terminal caches previously seen words

Words → Word
Words → Word Words
Word → Chars

## Morphology and word segmentation

Words → Word
Words → Word Words
Word → Stem Suffix
Word → Stem
Stem → Chars
Suffix → Chars

## Unigram model often finds collocations

- Unigram word segmentation model.
  generated independently
- But there are strong inter-word depe
- Unigram model can only capture such
  analyzing collocations as words

## Bigram segmentation model

- Implemented using Gibbs sampling
  - ith component is "word boundary at position i"
  - sampling amounts to possibly splitting a word at position i or joining the two words abutting at position i
- Performs significantly better than unigram model
  - bigram: 77% token f-score, 63% type f-score
  - unigram: 54% token f-score, 59% type f-score
- Number of CRPs is number of words, which is not known in advance
  ⇒ cannot be formulated as an adaptor grammar (which have a CRP per non-terminal)

## Hierarchical CRP bigram word segmentation

based on preceding word
to predict following word

$\sim$ BIGRAM
$\sim$ DP($a_1$, UNIGRAM)
$\sim$ DP($a_1$, $P_0$)

Slide 5:

*Johnson's slides on adding to PCFGs*

## Make dependencies local – GPSG-style

| rule | count | rel freq |
| --- | --- | --- |
| S → NP VP | 2 | 2/3 |
| S → NP VP | 1 | 1/3 |

## 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

## "Head to head" dependencies

Rules:
S → NP VP

## Generative language model (Charniak 2001)

The changes allow

- 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
- Final sense way of accurately estimating models in the presence of unmodeled dependencies
  ⇒ exponential models

Slide 6:

*More extensions of PCFGs*

## Exponential models

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

Let $f_1, \dots$ weights: $\lambda$

## PCFGs are exponential models

$\Psi$ = set of all trees generated by PCFG G
$f_j(\psi)$ = number of times the jth rule is used in $\psi$
$p(r_j)$ = probability of jth rule in G
Set weight $\lambda_j = \log p(r_j)$

## Advantages of exponential models

- Exponential models are very flexible
- Features $f$ can be *any function of parse*

## Modeling dependencies

## MLE of exponential models from visible da

Parses $\Psi = \psi_1, \dots, \psi_n$

- It's usually difficult
  a particular set of
  - probability of tl
    of *conditional pl*
  - non-local depe
    GPSG-style lead
- It's easy to make e
  dependencies
  - add a new feat
  - *figuring out wh*
    *incorporating th*

## 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 property they encode *need not to directly represented in the parse*
  - very different from PCFGs, where the tree label and shape determines the generalizations extracted

## Coarse to fine parsing

- Parsing with a gra
  nonterminals is slo
  programming algo
- Coarse to fine pars
  features of the coar
  classes of the finely
- The parses produce
  constrain the searc
- The Charniak gene
  PCFG to identify m
  the fine-grained PC