# Spectral Learning of Latent-Variable PCFGs

Karl Stratos[1]

Joint work with Shay Cohen[1], Michael Collins[1], Dean Foster[2], and Lyle Ungar[2]

[1]Columbia University

[2]University of Pennsylvania

# Latent-Variable Models for NLP and Speech

- Latent-variable models are of huge importance.
  - Speech recognition with HMMs
  - Gaussian mixture models
  - Machine translation with alignments as hidden variables
  - Latent-variable PCFGs (Matsuzaki et al., Petrov et al.)
  - Many many others
- The EM algorithm is remarkably successful. **But**:
  - No guarantee of reaching the global maximum of the likelihood function
  - Theoretical problem: parameter estimates not consistent
  - Practical problems: local optima difficult to deal with

# There is Hope

- Dasgupta (1999): Under separation conditions, it is possible to learn GMMs.
- Moitra and Valiant (2010): Arbitrary GMMs can be learned in polynomial time and sample complexity.
- Hsu, Kakade, and Zhang (2009): Under rank conditions, it is possible to learn HMMs efficiently and consistently.
- Kakade and Foster (2007): Under a wide class of models, CCA projections yield an optimal space for predicting hidden variables.

# This Work

- A spectral algorithm for learning latent-variable PCFGs
   L-PCFGs: Strong parsing performance (Petrov et al., 2006)

- Guaranteed to give consistent parameter estimates under assumptions on singular values

- Simple and efficient (SVD and matrix operations)

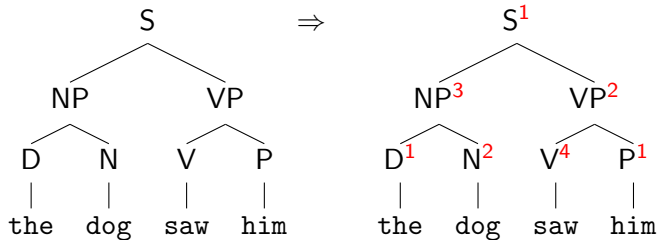# L-PCFGs (Matsuzaki et al., 2005, Petrov et al., 2006)

# The Probability of a Tree

$$p(\text{tree}, 1\ 3\ 1\ 2\ 2\ 4\ 1)$$
$$= \pi(\text{S}^1) \times$$
$$t(\text{S}^1 \rightarrow \text{NP}^3\ \text{VP}^2 | \text{S}^1) \times$$
$$t(\text{NP}^3 \rightarrow \text{D}^1\ \text{N}^2 | \text{NP}^3) \times$$
$$t(\text{VP}^2 \rightarrow \text{V}^4\ \text{P}^1 | \text{VP}^2) \times$$
$$q(\text{D}^1 \rightarrow \texttt{the} | \text{D}^1) \times$$
$$q(\text{N}^2 \rightarrow \texttt{dog} | \text{N}^2) \times$$
$$q(\text{V}^4 \rightarrow \texttt{saw} | \text{V}^4) \times$$
$$q(\text{P}^1 \rightarrow \texttt{him} | \text{P}^1)$$

```
            S¹
           /  \
       NP³      VP²
      /  \     /   \
    D¹    N²  V⁴    P¹
    |     |   |     |
   the   dog saw   him
```

$$p(\text{tree}) = \sum_{h_1 \ldots h_7} p(\text{tree}, h_1\ h_2\ h_3\ h_4\ h_5\ h_6\ h_7)$$

# The Probability of a Tree



$$p(\text{tree}, 1\ 3\ 1\ 2\ 2\ 4\ 1)$$
$$= \pi(\mathsf{S}^1) \times$$
$$\quad t(\mathsf{S}^1 \to \mathsf{NP}^3\ \mathsf{VP}^2 | \mathsf{S}^1) \times$$
$$\quad t(\mathsf{NP}^3 \to \mathsf{D}^1\ \mathsf{N}^2 | \mathsf{NP}^3) \times$$
$$\quad t(\mathsf{VP}^2 \to \mathsf{V}^4\ \mathsf{P}^1 | \mathsf{VP}^2) \times$$
$$\quad q(\mathsf{D}^1 \to \texttt{the} | \mathsf{D}^1) \times$$
$$\quad q(\mathsf{N}^2 \to \texttt{dog} | \mathsf{N}^2) \times$$
$$\quad q(\mathsf{V}^4 \to \texttt{saw} | \mathsf{V}^4) \times$$
$$\quad q(\mathsf{P}^1 \to \texttt{him} | \mathsf{P}^1)$$

$$p(\text{tree}) = \sum_{h_1 \ldots h_7} p(\textit{tree}, h_1\ h_2\ h_3\ h_4\ h_5\ h_6\ h_7)$$

# The Probability of a Tree

$$S^1$$

$$NP^3 \quad VP^2$$

$$D^1 \quad N^2 \quad V^4 \quad P^1$$

the   dog   saw   him

$p(\text{tree}, 1\ 3\ 1\ 2\ 2\ 4\ 1)$

$= \pi(S^1) \times$

$\quad t(S^1 \rightarrow NP^3\ VP^2 | S^1) \times$

$\quad t(NP^3 \rightarrow D^1\ N^2 | NP^3) \times$
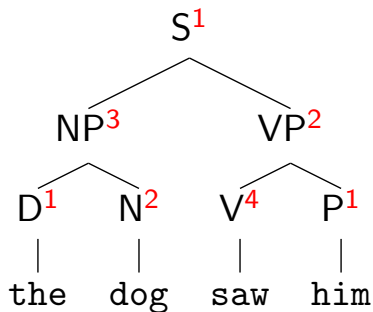
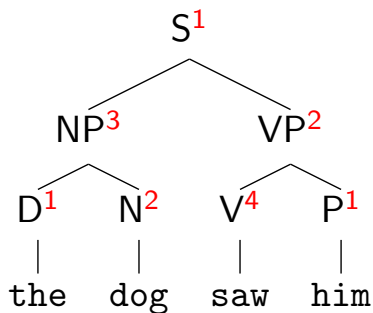$\quad t(VP^2 \rightarrow V^4\ P^1 | VP^2) \times$

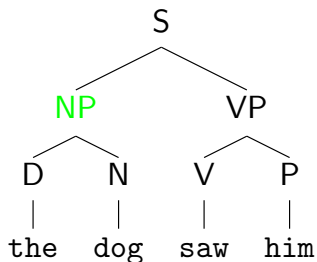$\quad q(D^1 \rightarrow \text{the} | D^1) \times$

$\quad q(N^2 \rightarrow \text{dog} | N^2) \times$

$\quad q(V^4 \rightarrow \text{saw} | V^4) \times$

$\quad q(P^1 \rightarrow \text{him} | P^1)$

$$p(\text{tree}) = \sum_{h_1 \ldots h_7} p(\text{tree}, h_1\ h_2\ h_3\ h_4\ h_5\ h_6\ h_7)$$

# The Probability of a Tree

$$S^1$$

$$NP^3 \qquad VP^2$$

$$D^1 \quad N^2 \quad V^4 \quad P^1$$

the    dog    saw    him

$p(\text{tree}, 1\ 3\ 1\ 2\ 2\ 4\ 1)$

$= \pi(S^1) \times$

$\quad t(S^1 \to NP^3\ \ VP^2 | S^1) \times$

$\quad t(NP^3 \to D^1\ \ N^2 | NP^3) \times$

$\quad t(VP^2 \to V^4\ \ P^1 | VP^2) \times$

$\quad q(D^1 \to \texttt{the} | D^1) \times$

$\quad q(N^2 \to \texttt{dog} | N^2) \times$

$\quad q(V^4 \to \texttt{saw} | V^4) \times$

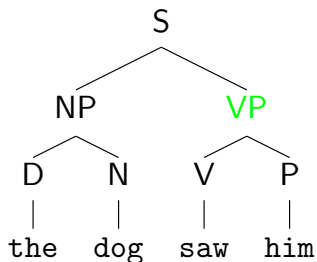$\quad q(P^1 \to \texttt{him} | P^1)$

$$p(\text{tree}) = \sum_{h_1 \ldots h_7} p(\text{tree}, h_1\ h_2\ h_3\ h_4\ h_5\ h_6\ h_7)$$

# Calculating Tree Probability with Dynamic Programming

$$b_h^1 = \sum_{h_2, h_3} t(\mathsf{NP}^h \to \mathsf{D}^{h_2}\ \mathsf{N}^{h_3}|\mathsf{NP}^h) \times q(\mathsf{D}^{h_2} \to \mathtt{the}|\mathsf{D}^{h_2}) \times q(\mathsf{N}^{h_3} \to \mathtt{dog}|\mathsf{N}^{h_3})$$
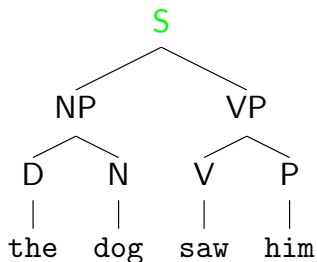
# Calculating Tree Probability with Dynamic Programming



$$b_h^1 = \sum_{h_2, h_3} t(\text{NP}^h \to \text{D}^{h_2} \ \text{N}^{h_3} | \text{NP}^h) \times q(\text{D}^{h_2} \to \texttt{the} | \text{D}^{h_2}) \times q(\text{N}^{h_3} \to \texttt{dog} | \text{N}^{h_3})$$

$$b_h^2 = \sum_{h_2, h_3} t(\text{VP}^h \to \text{V}^{h_2} \ \text{P}^{h_3} | \text{VP}^h) \times q(\text{V}^{h_2} \to \texttt{saw} | \text{V}^{h_2}) \times q(\text{P}^{h_3} \to \texttt{him} | \text{P}^{h_3})$$
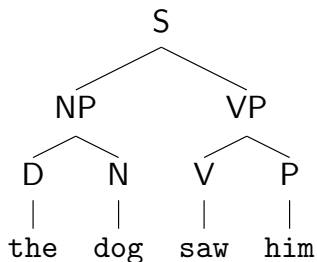
# Calculating Tree Probability with Dynamic Programming



$$b_h^1 = \sum_{h_2,h_3} t(\text{NP}^h \to \text{D}^{h_2} \ \text{N}^{h_3}|\text{NP}^h) \times q(\text{D}^{h_2} \to \texttt{the}|\text{D}^{h_2}) \times q(\text{N}^{h_3} \to \texttt{dog}|\text{N}^{h_3})$$

$$b_h^2 = \sum_{h_2,h_3} t(\text{VP}^h \to \text{V}^{h_2} \ \text{P}^{h_3}|\text{VP}^h) \times q(\text{V}^{h_2} \to \texttt{saw}|\text{V}^{h_2}) \times q(\text{P}^{h_3} \to \texttt{him}|\text{P}^{h_3})$$

$$b_h^3 = \sum_{h_2,h_3} t(\text{S}^h \to \text{NP}^{h_2} \ \text{VP}^{h_3}|\text{S}^h) \times b_{h_2}^1 \times b_{h_3}^2$$

# Calculating Tree Probability with Dynamic Programming



$$b_h^1 = \sum_{h_2,h_3} t(\text{NP}^h \to \text{D}^{h_2}\ \text{N}^{h_3}|\text{NP}^h) \times q(\text{D}^{h_2} \to \texttt{the}|\text{D}^{h_2}) \times q(\text{N}^{h_3} \to \texttt{dog}|\text{N}^{h_3})$$

$$b_h^2 = \sum_{h_2,h_3} t(\text{VP}^h \to \text{V}^{h_2}\ \text{P}^{h_3}|\text{VP}^h) \times q(\text{V}^{h_2} \to \texttt{saw}|\text{V}^{h_2}) \times q(\text{P}^{h_3} \to \texttt{him}|\text{P}^{h_3})$$

$$b_h^3 = \sum_{h_2,h_3} t(\text{S}^h \to \text{NP}^{h_2}\ \text{VP}^{h_3}|\text{S}^h) \times b_{h_2}^1 \times b_{h_3}^2$$

$$p(\text{tree}) = \sum_{h} \pi(\text{S}^h) \times b_h^3$$

# Marginals of a Sentence

- Given a sentence $x$, a marginal is defined as

$$\mu(a, i, j) = \sum_{t \in \tau(x):(a,i,j) \in t} p(t)$$

  for all $(a, i, j)$ tuples.

- These marginals can be computed using a variant of the inside-outside algorithm.

- A dynamic programming algorithm (Goodman, 1996) can be used to find the optimal parse defined as

$$t^* = \arg \max_{t \in \tau(x)} \sum_{(a,i,j) \in t} \mu(a, i, j)$$

# Parameter Estimation

- So this is a **parameter estimation** problem.
  - Given only skeletal trees, can we estimate $\pi$, $t$ and $q$?
- Past work used EM (Matsuzaki et al., 2005, Petrov et al. 2006).
  - No guarantee of converging to the correct distribution
  - Prone to local optima
- We present a spectral estimation method.
  - Under assumptions on singular values, gives consistent parameter estimates
  - Relatively simple, efficient

# Overview

# Inside and Outside Trees

At node VP:



Conditionally independent given the label and the hidden state

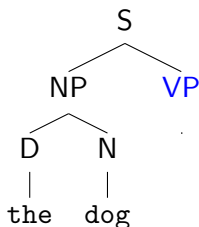$$p(o, t | \text{VP}, h) = p(o | \text{VP}, h) \times p(t | \text{VP}, h)$$

# Vector Representation of Inside and Outside Trees

Assume functions $Z$ and $Y$:

$Z$ maps any outside tree to a vector of length $m$.
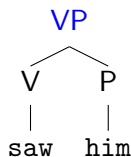
$Y$ maps any inside tree to a vector of length $m$.

Convention: $m$ is the number of hidden states under the L-PCFG.



Outside tree $o \Rightarrow$
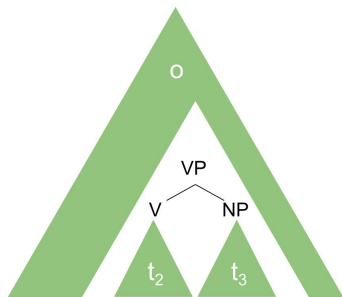$Z(o) = [1, 0.4, -5.3, \ldots, 72] \in \mathbb{R}^m$

Inside tree $t \Rightarrow$
$Y(t) = [-3, 17, 2, \ldots, 3.5] \in \mathbb{R}^m$

# Parameter Estimation for Binary Rules

Take $M$ samples of nodes with rule $VP \rightarrow V\ NP$.

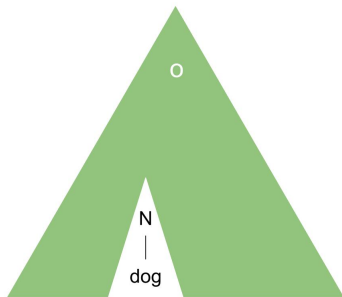

At sample $i$

- $o^{(i)} =$ outside tree at VP
- $t_2^{(i)} =$ inside tree at V
- $t_3^{(i)} =$ inside tree at NP

$$\hat{t}(VP^{h_1} \rightarrow V^{h_2}\ NP^{h_3}|VP^{h_1})$$

$$= \frac{\textbf{count}(VP \rightarrow V\ NP)}{\textbf{count}(VP)} \times \frac{1}{M} \sum_{i=1}^{M} \left( Z_{h_1}(o^{(i)}) \times Y_{h_2}(t_2^{(i)}) \times Y_{h_3}(t_3^{(i)}) \right)$$

# Parameter Estimation for Unary Rules

Take $M$ samples of nodes with rule $N \rightarrow \texttt{dog}$.



At sample $i$

- $o^{(i)} =$ outside tree at $N$

$$\hat{q}(N^h \rightarrow \texttt{dog}|N^h) = \frac{\textbf{count}(N \rightarrow \texttt{dog})}{\textbf{count}(N)} \times \frac{1}{M}\sum_{i=1}^{M} Z_h(o^{(i)})$$

# Parameter Estimation for the Root

Take $M$ samples of the root S.

S

t
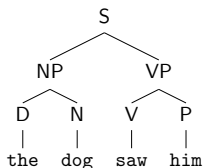
At sample $i$

- $t^{(i)} =$ inside tree at S

$$\hat{\pi}(S^h) = \frac{\textbf{count}(\text{root=S})}{\textbf{count}(\text{root})} \times \frac{1}{M} \sum_{i=1}^{M} Y_h(t^{(i)})$$

# Calculating Tree Probability with Dynamic Programming: Revisited



$$\hat{b}_h^1 = \sum_{h_2, h_3} \hat{t}(\mathsf{NP}^h \to \mathsf{D}^{h_2}\ \mathsf{N}^{h_3}|\mathsf{NP}^h) \times \hat{q}(\mathsf{D}^{h_2} \to \mathtt{the}|\mathsf{D}^{h_2}) \times \hat{q}(\mathsf{N}^{h_3} \to \mathtt{dog}|\mathsf{N}^{h_3})$$

$$\hat{b}_h^2 = \sum_{h_2, h_3} \hat{t}(\mathsf{VP}^h \to \mathsf{V}^{h_2}\ \mathsf{P}^{h_3}|\mathsf{VP}^h) \times \hat{q}(\mathsf{V}^{h_2} \to \mathtt{saw}|\mathsf{V}^{h_2}) \times \hat{q}(\mathsf{P}^{h_3} \to \mathtt{him}|\mathsf{P}^{h_3})$$

$$\hat{b}_h^3 = \sum_{h_2, h_3} \hat{t}(\mathsf{S}^h \to \mathsf{NP}^{h_2}\ \mathsf{VP}^{h_3}|\mathsf{S}^h) \times \hat{b}_{h_2}^1 \times \hat{b}_{h_3}^2$$

$$p(\text{tree}) = \sum_h \hat{\pi}(S^h) \times \hat{b}_h^3$$

# Overview

# Deriving $Z$ and $Y$

Design functions $\psi$ and $\phi$:

$\psi$ maps any outside tree to a vector of length $d'$

$\phi$ maps any inside tree to a vector of length $d$



Outside tree $o \Rightarrow$
$\psi(o) = [0, 1, 0, 0, \ldots, 0, 1] \in \mathbb{R}^{d'}$

Inside tree $t \Rightarrow$
$\phi(t) = [1, 0, 0, 0, \ldots, 1, 0] \in \mathbb{R}^{d}$

$Z$ and $Y$ will be reduced dimensional representations of $\psi$ and $\phi$.

# Reducing Dimensions via a Singular Value Decomposition

Have $M$ samples of a node with non-terminal $a$. At sample $i$, $o^{(i)}$ is the outside tree rooted at $a$ and $t^{(i)}$ is the inside tree rooted at $a$.

- Compute a matrix $\hat{\Omega}^a \in \mathbb{R}^{d \times d'}$ with entries

$$[\hat{\Omega}^a]_{j,k} = \frac{1}{M} \sum_{i=1}^{M} \phi_j(t^{(i)}) \psi_k(o^{(i)})$$

# Reducing Dimensions via a Singular Value Decomposition

Have $M$ samples of a node with non-terminal $a$. At sample $i$, $o^{(i)}$ is the outside tree rooted at $a$ and $t^{(i)}$ is the inside tree rooted at $a$.

▶ Compute a matrix $\hat{\Omega}^a \in \mathbb{R}^{d \times d'}$ with entries

$$[\hat{\Omega}^a]_{j,k} = \frac{1}{M} \sum_{i=1}^{M} \phi_j(t^{(i)}) \psi_k(o^{(i)})$$

▶ An SVD:

$$\underbrace{\hat{\Omega}^a}_{d \times d'} \approx \underbrace{U^a}_{d \times m} \underbrace{\Sigma^a}_{m \times m} \underbrace{(V^a)^T}_{m \times d'}$$

# Reducing Dimensions via a Singular Value Decomposition

Have $M$ samples of a node with non-terminal $a$. At sample $i$, $o^{(i)}$ is the outside tree rooted at $a$ and $t^{(i)}$ is the inside tree rooted at $a$.

▶ Compute a matrix $\hat{\Omega}^a \in \mathbb{R}^{d \times d'}$ with entries

$$[\hat{\Omega}^a]_{j,k} = \frac{1}{M} \sum_{i=1}^{M} \phi_j(t^{(i)}) \psi_k(o^{(i)})$$

▶ An SVD:

$$\underbrace{\hat{\Omega}^a}_{d \times d'} \approx \underbrace{U^a}_{d \times m} \underbrace{\Sigma^a}_{m \times m} \underbrace{(V^a)^T}_{m \times d'}$$

▶ Projection:

$$Y(t^{(i)}) = \underbrace{(U^a)^T}_{m \times d} \underbrace{\phi(t^{(i)})}_{d \times 1} \in \mathbb{R}^m$$

$$Z(o^{(i)}) = \underbrace{(\Sigma^a)^{-1}}_{m \times m} \underbrace{(V^a)^T}_{m \times d'} \underbrace{\psi(o^{(i)})}_{d' \times 1} \in \mathbb{R}^m$$

# Consistency and Sample Complexity

If the $d \times d'$ matrix

$$\Omega^a = \mathbf{E}[\phi(T)\psi(O)^T | \text{label} = a]$$

has rank $m$, these projections yield consistent parameter estimates with high probability. The required number of samples grows polynomially in

- $m$: the number of hidden states
- $\log R$: where $R$ is the number of rules
- Spectral properties of the grammar (e.g., $\max \frac{1}{\sigma^a}$ where $\sigma^a$ is the $m^{th}$ largest singular value of $\Omega^a$)

# A Summary of the Algorithm

1. Design feature functions $\phi$ and $\psi$ for inside and outside trees.
2. Use SVD to compute vectors
   $Y(t) \in \mathbb{R}^m$ for inside trees
   $Z(o) \in \mathbb{R}^m$ for outside trees
3. Estimate the parameters $\hat{t}$, $\hat{q}$, and $\hat{\pi}$ from the training data.
4. Parse a new sentence by computing its marginals with these parameters.

# Overview

L-PCFGs

The Spectral Algorithm for Parameter Estimation
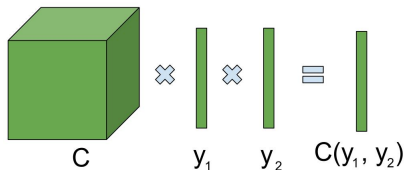### Calculating Parameter Estimates
### SVD and Projection

Justification

# Tensor Definition

A third-order tensor $C \in \mathbb{R}^{m \times m \times m}$ is a set of $m^3$ values $[C]_{j,k,l}$. It can be viewed as a function $C : \mathbb{R}^m \times \mathbb{R}^m \to \mathbb{R}^m$ that takes two vectors $y_1, y_2 \in \mathbb{R}^m$ as input and returns a vector $C(y_1, y_2) \in \mathbb{R}^m$ as output. The output vector has entries

$$[C(y_1, y_2)]_h = \sum_{h_2, h_3} \left( [C]_{h, h_2, h_3} \times [y_1]_{h_2} \times [y_2]_{h_3} \right)$$



C       $y_1$       $y_2$       C($y_1$, $y_2$)

# Tensor Form of the Parameters

For each non-terminal $a$, define a vector $\pi^a \in \mathbb{R}^m$ with entries
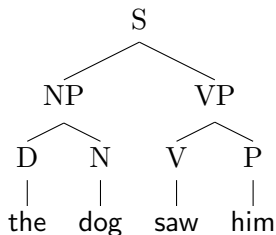
$$[\pi^a]_h = \pi(a^h)$$

For each rule $a \to x$, define a vector $q_{a \to x} \in \mathbb{R}^m$ with entries

$$[q_{a \to x}]_h = q_{a \to x}(a^h \to x | a^h)$$

For each rule $a \to b\ c$, define a tensor $T^{a \to b\ c} \in \mathbb{R}^{m \times m \times m}$ with entries

$$[T^{a \to b\ c}]_{h_1, h_2, h_3} = t(a^{h_1} \to b^{h_2}\ c^{h_3} | a^{h_1})$$

# Dynamic Programming in Tensor Form



$$T^{\text{S}\to\text{NP VP}}(T^{\text{NP}\to\text{D N}}(q_{\text{D}\to\text{the}}, q_{\text{N}\to\text{dog}}), T^{\text{VP}\to\text{V P}}(q_{\text{V}\to\text{saw}}, q_{\text{P}\to\text{him}}))\, \pi^{\text{S}}$$

$$|||$$

$$p(\text{tree}) = \sum_{h_1 \ldots h_7} p(\textit{tree}, h_1\ h_2\ h_3\ h_4\ h_5\ h_6\ h_7)$$

# Thought Experiment

- We want the parameters (in tensor form)

$$\pi^a \in \mathbb{R}^m$$

$$q_{a \to x} \in \mathbb{R}^m$$
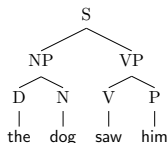
$$T^{a \to b\ c}(y_2, y_3) \in \mathbb{R}^m$$

- What if we had an invertible matrix $G^a \in \mathbb{R}^{m \times m}$ for every non-terminal $a$?

- And what if we had instead

$$c^a = G^a \pi^a$$

$$c_{a \to x} = q_{a \to x}(G^a)^{-1}$$

$$C^{a \to b\ c}(y_2, y_3) = T^{a \to b\ c}(y_2 G^b, y_3 G^c)(G^a)^{-1}$$

# Cancellation of the Linear Operators



$$C^{\mathrm{S\to NP\,VP}}(C^{\mathrm{NP\to D\,N}}(c_{\mathrm{D\to the}}, c_{\mathrm{N\to dog}}), C^{\mathrm{VP\to V\,P}}(c_{\mathrm{V\to saw}}, c_{\mathrm{P\to him}}))\ c^{\mathrm{S}}$$

$$|||$$

$$T^{\mathrm{S\to NP\,VP}}(T^{\mathrm{NP\to D\,N}}(q_{\mathrm{D\to the}}(G^{\mathrm{D}})^{-1}G^{\mathrm{D}}, q_{\mathrm{N\to dog}}(G^{\mathrm{N}})^{-1}G^{\mathrm{N}})(G^{\mathrm{NP}})^{-1}G^{\mathrm{NP}},$$
$$T^{\mathrm{VP\to V\,P}}(q_{\mathrm{V\to saw}}(G^{\mathrm{V}})^{-1}G^{\mathrm{V}}, q_{\mathrm{P\to him}}(G^{\mathrm{P}})^{-1}G^{\mathrm{P}})(G^{\mathrm{VP}})^{-1}G^{\mathrm{VP}})(G^{\mathrm{S}})^{-1}G^{\mathrm{S}}\pi^{\mathrm{S}}$$

$$|||$$

$$T^{\mathrm{S\to NP\,VP}}(T^{\mathrm{NP\to D\,N}}(q_{\mathrm{D\to the}}, q_{\mathrm{N\to dog}}), T^{\mathrm{VP\to V\,P}}(q_{\mathrm{V\to saw}}, q_{\mathrm{P\to him}}))\ \pi^{\mathrm{S}}$$

$$|||$$

$$p(\mathrm{tree}) = \sum_{h_1\ldots h_7} p(\mathit{tree},\ h_1\ h_2\ h_3\ h_4\ h_5\ h_6\ h_7)$$

# Estimation Guarantees

- Basic argument: If $\Omega^a$ has rank $m$, parameters $\hat{C}^{a \to b\,c}$, $\hat{c}_{a \to x}$, and $\hat{c}^a$ converge to

$$C^{a \to b\,c}(y_2, y_3) = T^{a \to b\,c}(y_2 G^b, y_3 G^c)(G^a)^{-1}$$
$$c_{a \to x} = q_{a \to x}(G^a)^{-1}$$
$$c^a = G^a \pi^a$$

for some $G^a$ that is invertible.

- Because the parameters converge, the estimated distribution $\hat{p}(\text{tree})$ converges to the true distribution $p(\text{tree})$, and the estimated marginal $\hat{\mu}(a, i, j)$ converges to the true marginal $\mu(a, i, j)$.

# Preliminary Experiments

The algorithm is much faster than EM.

- ► SVD: modern algorithms are very efficient
- ► Parameter calculation: takes less time than a single iteration of EM

A straightforward implementation lags behind EM by about 1-2% in F1 measure.

Current work: experiments focused on understanding the method and improving performance

# Summary

We presented a spectral algorithm that yields a consistent estimator for L-PCFGs

- Simple and efficient: SVD and standard matrix operations

Future work includes

- Pushing the empirical side of the algorithm
- Deriving spectral algorithms for other latent-variable models in NLP