

Sparsity in Dependency Grammar Induction

Jennifer Gillenwater¹ Kuzman Ganchev¹ João Graça²
Ben Taskar¹ Fernando Pereira³

¹Computer & Information Science
University of Pennsylvania

²L²F INESC-ID, Lisboa, Portugal

³Google, Inc.

July 12, 2010

- A generative dependency parsing model

- A generative dependency parsing model
- The **ambiguity** problem this model faces

- A generative dependency parsing model
- The **ambiguity** problem this model faces
- Previous attempts to reduce ambiguity

- A generative dependency parsing model
- The **ambiguity** problem this model faces
- Previous attempts to reduce ambiguity
- How posteriors provide a good measure of ambiguity

- A generative dependency parsing model
- The **ambiguity** problem this model faces
- Previous attempts to reduce ambiguity
- How posteriors provide a good measure of ambiguity
- Applying **posterior regularization** to the likelihood objective

- A generative dependency parsing model
- The **ambiguity** problem this model faces
- Previous attempts to reduce ambiguity
- How posteriors provide a good measure of ambiguity
- Applying **posterior regularization** to the likelihood objective
- Success with respect to EM and parameter prior baselines

Dependency model with valence

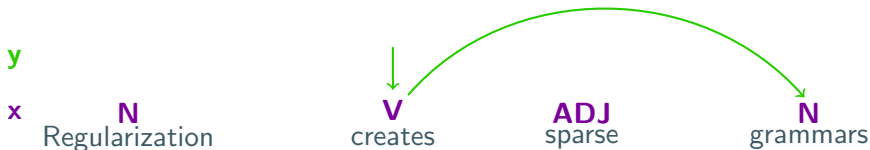
(Klein and Manning, ACL 2004)



$$p_{\theta}(\mathbf{x}, \mathbf{y}) = \theta_{\text{root}(V)}$$

Dependency model with valence

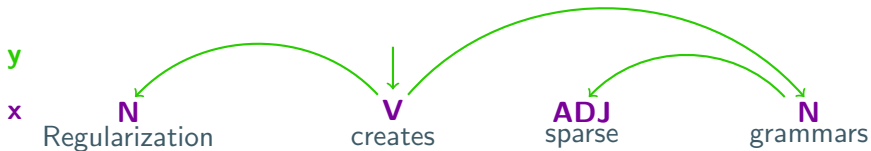
(Klein and Manning, ACL 2004)



$$p_{\theta}(\mathbf{x}, \mathbf{y}) = \theta_{root(V)} \cdot \theta_{stop(nostop|V, right, false)} \cdot \theta_{child(N|V, right)}$$

Dependency model with valence

(Klein and Manning, ACL 2004)



$$\begin{aligned} p_{\theta}(\mathbf{x}, \mathbf{y}) = & \theta_{root(\mathbf{V})} \\ & \cdot \theta_{stop(nostop|\mathbf{V},right,false)} \cdot \theta_{child(\mathbf{N}|\mathbf{V},right)} \\ & \cdot \theta_{stop(stop|\mathbf{V},right,true)} \cdot \theta_{stop(nostop|\mathbf{V},left,false)} \cdot \theta_{child(\mathbf{N}|\mathbf{V},left)} \\ & \dots \end{aligned}$$

- **Traditional objective:** marginal log likelihood

$$\max_{\theta} \mathcal{L}(\theta) = E_{\mathbf{X}}[\log p_{\theta}(\mathbf{x})] = E_{\mathbf{X}}[\log \sum_{\mathbf{y}} p_{\theta}(\mathbf{x}, \mathbf{y})]$$

- **Traditional objective:** marginal log likelihood

$$\max_{\theta} \mathcal{L}(\theta) = E_{\mathcal{X}}[\log p_{\theta}(\mathbf{x})] = E_{\mathcal{X}}[\log \sum_{\mathbf{y}} p_{\theta}(\mathbf{x}, \mathbf{y})]$$

- **Optimization method:** expectation maximization (EM)

- **Traditional objective:** marginal log likelihood

$$\max_{\theta} \mathcal{L}(\theta) = E_{\mathbf{X}}[\log p_{\theta}(\mathbf{x})] = E_{\mathbf{X}}[\log \sum_{\mathbf{y}} p_{\theta}(\mathbf{x}, \mathbf{y})]$$

- **Optimization method:** expectation maximization (EM)
- **Problem:** EM may learn a very ambiguous grammar
 - Too many non-zero probabilities
 - Ex: $V \rightarrow N$ should have non-zero probability,
but $V \rightarrow DET$, $V \rightarrow JJ$, $V \rightarrow PRP\$$, etc. should be 0

- Structural annealing¹

1 Smith and Eisner, ACL 2006

2 Headden et al., NAACL 2009

3 Liang et al., EMNLP 2007; Johnson et al., NIPS 2007; Cohen et al., NIPS 2008, NAACL 2009

Previous approaches to improving performance

- Structural annealing¹
- $\mathcal{L}(\theta')$: Model extension²

1 Smith and Eisner, ACL 2006

2 Headden et al., NAACL 2009

3 Liang et al., EMNLP 2007; Johnson et al., NIPS 2007; Cohen et al., NIPS 2008, NAACL 2009

Previous approaches to improving performance

- Structural annealing¹
- $\mathcal{L}(\theta')$: Model extension²
- $\mathcal{L}(\theta) + \log p(\theta)$: Parameter regularization³
 - Tend to **reduce unique # of children per parent**, rather than directly **reducing # of unique parent-child pairs**
 - $\theta_{child(Y|X,dir)} \neq \text{posterior}(X \rightarrow Y)$

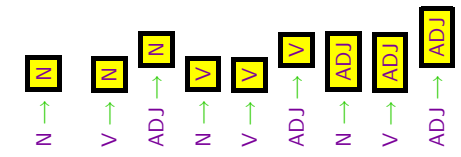
1 Smith and Eisner, ACL 2006

2 Headden et al., NAACL 2009

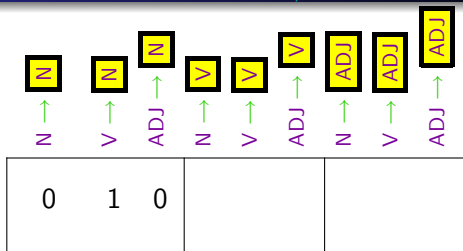
3 Liang et al., EMNLP 2007; Johnson et al., NIPS 2007; Cohen et al., NIPS 2008, NAACL 2009

Ambiguity measure using posteriors: $L_{1/\infty}$

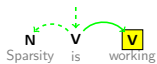
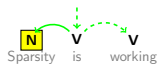
Intuition: True # of unique parent tags for a child tag is small



Ambiguity measure using posteriors: $L_{1/\infty}$

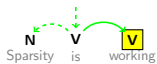
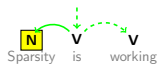


Ambiguity measure using posteriors: $L_{1/\infty}$



	N	N	N	V	V	V	ADJ	ADJ	ADJ
	↑	↑	↑	↑	↑	↑	↑	↑	↑
	N	V	ADJ	N	V	ADJ	N	V	ADJ
0	1	0							
			0	1	0				

Ambiguity measure using posteriors: $L_{1/\infty}$



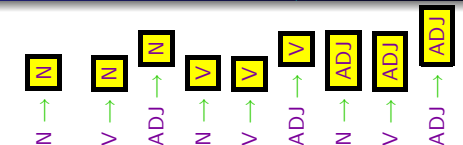
	N	N	N	V	V	V	ADJ	ADJ	ADJ
	↑	↑	↑	↑	↑	↑	↑	↑	↑
	N	V	ADJ	N	V	ADJ	N	V	ADJ
0	1	0							
			0	1	0				
0	1	0							

Ambiguity measure using posteriors: $L_{1/\infty}$



	N	V	ADJ	N	V	V	ADJ	ADJ	ADJ
	N	V	ADJ	N	V	ADJ	N	V	ADJ
	0	1	0						
				0	1	0			
	0	1	0						
							1	0	0

Ambiguity measure using posteriors: $L_{1/\infty}$



0	1	0					
			0	1	0		
0	1	0					
					1	0	0

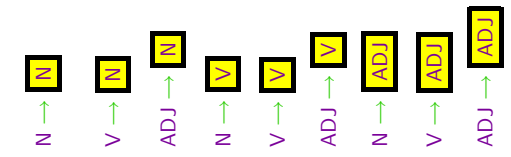
max ↓

sum = 3 ←

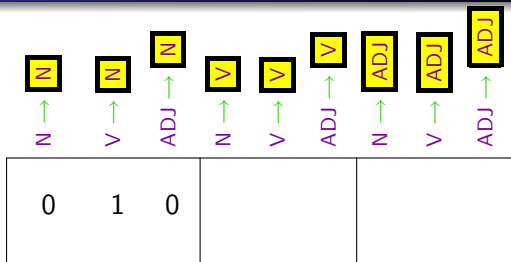
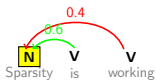
0	1	0	0	1	0	1	0	0
---	---	---	---	---	---	---	---	---

Measuring ambiguity on distributions over trees

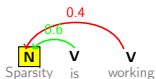
For a distribution $p_{\theta}(\mathbf{y} \mid \mathbf{x})$ instead of gold trees:



Measuring ambiguity on distributions over trees

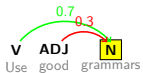
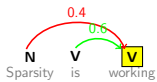
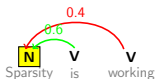


Measuring ambiguity on distributions over trees



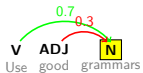
	N	N	N	V	V	V	ADJ	ADJ	ADJ
	↑	↑	↑	↑	↑	↑	↑	↑	↑
	N	V	ADJ	N	V	ADJ	N	V	ADJ
	0	1	0						
				.4	.6	0			

Measuring ambiguity on distributions over trees



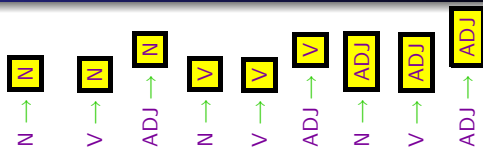
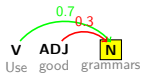
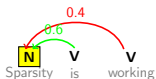
	N	N	N	V	V	V	ADJ	ADJ	ADJ
	↑	↑	↑	↑	↑	↑	↑	↑	↑
	N	V	ADJ	N	V	ADJ	N	V	ADJ
	0	1	0						
				.4	.6	0			
	0	.7	.3						

Measuring ambiguity on distributions over trees



	N	N	N	V	V	V	ADJ	ADJ	ADJ
	↑	↑	↑	↑	↑	↑	↑	↑	↑
	N	V	ADJ	N	V	ADJ	N	V	ADJ
0	1	0							
				.4	.6	0			
0	.7	.3							
							.4	.6	0

Measuring ambiguity on distributions over trees



0	1	0					
			.4	.6	0		
0	.7	.3					
					.4	.6	0

max ↓

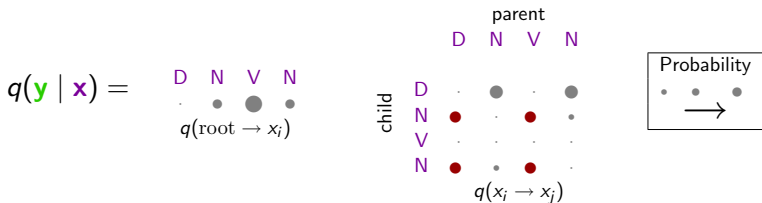
sum = 3.3 ←

0	1	.3	.4	.6	0	.4	.6	0
---	---	----	----	----	---	----	----	---

E-Step $q^t(\mathbf{y} \mid \mathbf{x}) = \arg \min_{q(\mathbf{y} \mid \mathbf{x})} KL(q \parallel p_{\theta^t})$

Minimizing ambiguity through posterior regularization

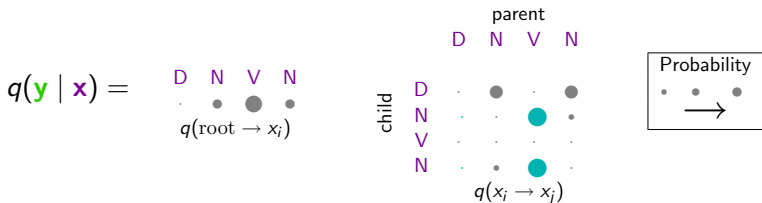
E-Step $q^t(\mathbf{y} \mid \mathbf{x}) = \arg \min_{q(\mathbf{y} \mid \mathbf{x})} KL(q \parallel p_{\theta^t})$



Minimizing ambiguity through posterior regularization

Apply E-step penalty $L_{1/\infty}$ on posteriors $q(\mathbf{y} | \mathbf{x})$ to induce sparsity
(Graca et al., NIPS 2007 & 2009)

$$\mathbf{E}\text{-Step} \quad q^t(\mathbf{y} | \mathbf{x}) = \arg \min_{q(\mathbf{y}|\mathbf{x})} KL(q \parallel p_{\theta^t}) + \sigma L_{1/\infty}(q(\mathbf{y} | \mathbf{x}))$$



- English from Penn Treebank: state-of-the-art accuracy

Learning Method	Accuracy		
	≤ 10	≤ 20	all
PR ($\sigma = 140$)	62.1	53.8	49.1
LN families	59.3	45.1	39.0
SLN TieV & N	61.3	47.4	41.4
PR ($\sigma = 140, \lambda = 1/3$)	64.4	55.2	50.5
DD ($\alpha = 1, \lambda$ learned)	65.0 (± 5.7)		

- English from Penn Treebank: state-of-the-art accuracy

Learning Method	Accuracy		
	≤ 10	≤ 20	all
PR ($\sigma = 140$)	62.1	53.8	49.1
LN families	59.3	45.1	39.0
SLN TieV & N	61.3	47.4	41.4
PR ($\sigma = 140, \lambda = 1/3$)	64.4	55.2	50.5
DD ($\alpha = 1, \lambda$ learned)	65.0 (± 5.7)		

- 11 other languages from CoNLL-X:
 - Dirichlet prior baseline: **1.5%** average gain over EM
 - Posterior regularization: **6.5%** average gain over EM

- English from Penn Treebank: state-of-the-art accuracy

Learning Method	Accuracy		
	≤ 10	≤ 20	all
PR ($\sigma = 140$)	62.1	53.8	49.1
LN families	59.3	45.1	39.0
SLN TieV & N	61.3	47.4	41.4
PR ($\sigma = 140, \lambda = 1/3$)	64.4	55.2	50.5
DD ($\alpha = 1, \lambda$ learned)	65.0 (± 5.7)		

- 11 other languages from CoNLL-X:
 - Dirichlet prior baseline: **1.5%** average gain over EM
 - Posterior regularization: **6.5%** average gain over EM
- Come see the poster for more details**