

# Sparsity in Dependency Grammar Induction

Jennifer Gillenwater<sup>1</sup>   Kuzman Ganchev<sup>1</sup>   João Graça<sup>2</sup>  
Ben Taskar<sup>1</sup>   Fernando Pereira<sup>3</sup>

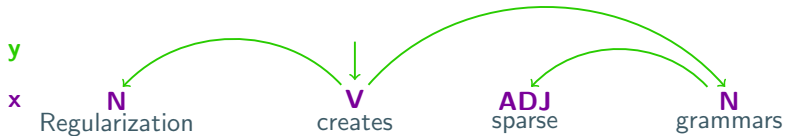
<sup>1</sup>Computer & Information Science  
University of Pennsylvania

<sup>2</sup>L<sup>2</sup>F INESC-ID, Lisboa, Portugal

<sup>3</sup>Google, Inc.

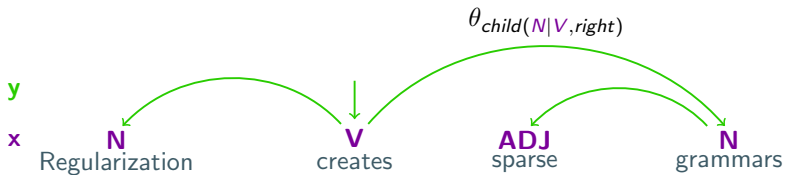
# Dependency model with valence (Klein and Manning, ACL 2004)

$$p_{\theta}(\mathbf{x}, \mathbf{y})$$



# Dependency model with valence (Klein and Manning, ACL 2004)

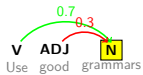
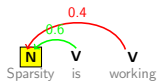
$$p_{\theta}(\mathbf{x}, \mathbf{y})$$






- **Traditional optimization:** expectation maximization (EM)

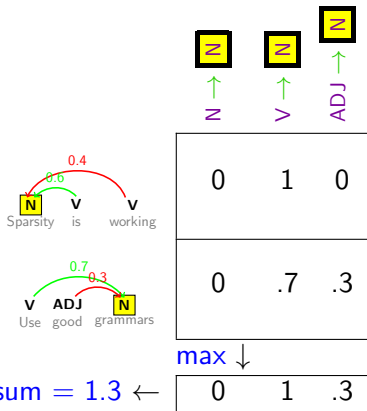
- **Traditional optimization:** expectation maximization (EM)
- **Problem:** EM may learn a very ambiguous grammar
  - $V \rightarrow N$  should have non-zero probability, but ...
  - $V \rightarrow DET, V \rightarrow JJ, V \rightarrow PRP\$,$  etc. should be 0

# Measuring ambiguity on distributions over trees



			
	↑	↑	↑
	N	>	ADJ
	0	1	0
	0	.7	.3

# Measuring ambiguity on distributions over trees



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



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

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

Learning Method	Accuracy		
	$\leq 10$	$\leq 20$	all
EM	45.8	40.2	35.9
Sparsifying Dirichlet Prior	46.4	40.9	36.5
PR ( $\sigma = 140$ )	<b>62.1</b>	<b>53.8</b>	<b>49.1</b>

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

Learning Method	Accuracy		
	$\leq 10$	$\leq 20$	all
EM	45.8	40.2	35.9
Sparsifying Dirichlet Prior	46.4	40.9	36.5
PR ( $\sigma = 140$ )	<b>62.1</b>	<b>53.8</b>	<b>49.1</b>

- 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		
	$\leq 10$	$\leq 20$	all
EM	45.8	40.2	35.9
Sparsifying Dirichlet Prior	46.4	40.9	36.5
PR ( $\sigma = 140$ )	<b>62.1</b>	<b>53.8</b>	<b>49.1</b>

- 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**