### GRAPH-BASED POSTERIOR REGULARIZATION FOR SEMI-SUPERVISED STRUCTURED PREDICTION

Luheng HeJennifer Gillenwater<sup>†</sup>Ben Taskar†University of PennsylvaniaUniversity of Washington

### OVERVIEW

#### Structured Prediction (CRF)

Graph Propagation

A Joint Objective



 $\mathcal{J}(q, \mathbf{p}_{\theta})$ 



#### This is recognized as a face











#### This is recognized as a face . VERB

# The painting is considered as a work of genius . ?













### GRAPH LAPLACIAN REGULARIZER

**NOUN** ... a run along ...

Prob (NOUN | ninth run for) = 0.6 Prob (VERB | ninth run for) = 0.4 Prob (ADV | ninth run for) = 0 Prob (DET | ninth run for) = 0



### GRAPH LAPLACIAN REGULARIZER

**NOUN** ... a run along ...

 $\begin{aligned} \operatorname{Prob}(\operatorname{tag} \mid \operatorname{ninth} \operatorname{run} \operatorname{for}) &= \\ \operatorname{arg\,min}_{m} & 0.4 \times \| \, m - \operatorname{Prob}(\operatorname{tag} \mid \operatorname{a} \operatorname{run} \operatorname{along}) \|_{2}^{2} \\ &+ 0.8 \times \| \, m - \operatorname{Prob}(\operatorname{tag} \mid \operatorname{a} \operatorname{run} \operatorname{for}) \|_{2}^{2} \\ &+ 0.8 \times \| \, m - \operatorname{Prob}(\operatorname{tag} \mid \operatorname{luck} \operatorname{run} \operatorname{out}) \|_{2}^{2} \end{aligned}$ 



### GRAPH LAPLACIAN REGULARIZER

**NOUN** ... a run along ...

 $\min_{\mathbf{m}} \operatorname{Lap}(\mathbf{m}) = \sum_{a \in \text{Unlab}} \sum_{b \in \text{Neighbors}(a)} \sum_{k \in \text{Tags}} w_{ab} (m_{a,k} - m_{b,k})^2$ 

 $m_{a,k}$  : The proportion of time trigram **a** has tag **k** 



## COMBININGTHETWO







**Our work**: retains <u>efficiency</u> while optimizing an <u>extendible, joint</u> objective.

## HOWTO COMBINE?

introduce auxiliary variables q

## HOW TO COMBINE?

introduce auxiliary variables q

$$p_{\theta}(\mathbf{y} \mid \mathbf{x}^{i}) = \frac{1}{Z_{\theta}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} \theta^{\top} \mathbf{f}(y_{t}, y_{t-1}, \mathbf{x}^{i})\right]$$
$$q_{y}^{i} = \frac{1}{Z_{q}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}(y_{t}, y_{t-1})\right]$$

## HOWTO COMBINE?

introduce auxiliary variables q

I. Normalized

$$p_{\theta}(\mathbf{y} \mid \mathbf{x}^{i}) = \frac{1}{Z_{\theta}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} \theta^{\top} \mathbf{f}(y_{t}, y_{t-1}, \mathbf{x}^{i})\right]$$
$$q_{y}^{i} = \frac{1}{Z_{q}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}(y_{t}, y_{t-1})\right]$$

## HOWTO COMBINE?

introduce auxiliary variables q

I. Normalized 2. Decomposed into local factors

$$p_{\theta}(\mathbf{y} \mid \mathbf{x}^{i}) = \frac{1}{Z_{\theta}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} \theta^{\top} \mathbf{f}(y_{t}, y_{t-1}, \mathbf{x}^{i})\right]$$
$$q_{y}^{i} = \frac{1}{Z_{q}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}(y_{t}, y_{t-1})\right]$$





 $\mathcal{J}(q, p_{\theta}) = \operatorname{Lap}(q) + \operatorname{NLik}(p_{\theta}) + \operatorname{KL}(q || p_{\theta})$ 







$$\mathrm{KL}(q \parallel p_{\theta}) = \sum_{i=1}^{n} \sum_{\mathbf{y}} q_{\mathbf{y}}^{i} \log \frac{q_{\mathbf{y}}^{i}}{p_{\theta}(\mathbf{y} \mid \mathbf{x}^{i})}$$

 $\min_{q,\boldsymbol{\theta}} \mathcal{J}(q,\boldsymbol{p_{\boldsymbol{\theta}}})$ 

 $\min_{\substack{q,\theta \\ \Delta}} \mathcal{J}(q, p_{\theta})$ 

*p* update:

$$\theta' = \theta - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial \theta}$$

$$\min_{\substack{q,\theta \\ \Delta}} \mathcal{J}(q, p_{\theta})$$

*p* update:

$$\theta' = \theta - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial \theta}$$

#### ${\it q}$ update:

projection is hard  $\sum_{\mathbf{y}} q_{\mathbf{y}}^{i} = 1$ no compact form  $(\# \text{ tags})^{(i's \text{ length})}$  values

 ${\bf q}$  can be represented by local factors  ${\bf r}$ 

$$q_{\mathbf{y}}^{i} = \frac{1}{Z_{q}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}(y_{t}, y_{t-1})\right]$$

**q** can be represented by local factors **r** 

$$q_{\mathbf{y}}^{i} = \frac{1}{Z_{q}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}(y_{t}, y_{t-1})\right]$$

doing an additive gradient update

$$q_{\mathbf{y}}^{i \prime} = q_{\mathbf{y}}^{i} - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^{i}}$$

 ${\bf q}$  can be represented by local factors  ${\bf r}$ 

$$q_{\mathbf{y}}^{i} = \frac{1}{Z_{q}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}(y_{t}, y_{t-1})\right]$$

doing an additive gradient update

$$q_{\mathbf{y}}^{i \prime} = q_{\mathbf{y}}^{i} - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^{i}}$$

q' cannot be written as product of local factors!

multiplicative gradient update:

$$q_{\mathbf{y}}^{i\,\prime} = \frac{1}{Z_{q^{\prime}}(\mathbf{x}^{i})} q_{\mathbf{y}}^{i} \exp\left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^{i}}\right]$$

multiplicative gradient update:

$$q_{\mathbf{y}}^{i \prime} = \frac{1}{Z_{q'}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}(y_{t}, y_{t-1})\right] \exp\left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^{i}}\right]$$

multiplicative gradient update:

$$q_{\mathbf{y}}^{i\,\prime} = \frac{1}{Z_{q^{\prime}}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}(y_{t}, y_{t-1})\right] \exp\left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^{i}}\right]$$
  
decompose into local factors

multiplicative gradient update:

$$q_{\mathbf{y}'}^{i\,\prime} = \frac{1}{Z_{q'}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}(y_{t}, y_{t-1})\right] \exp\left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^{i}}\right]$$
  
decompose into local factors  
$$= \frac{1}{Z_{q'}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}'(y_{t}, y_{t-1})\right]$$

multiplicative gradient update:

$$q_{\mathbf{y}'}^{i\,\prime} = \frac{1}{Z_{q'}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}(y_{t}, y_{t-1})\right] \exp\left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^{i}}\right]$$
decompose into local factors
$$= \frac{1}{Z_{q'}(\mathbf{x}^{i})} \exp\left[\sum_{t=1}^{T} r_{i,t}'(y_{t}, y_{t-1})\right]$$
only updating  $(\# \text{tags})^{2} \times (i\text{'s length})$  variables!

### SUMMARY



 $\mathcal{J}(q, p_{\theta}) = \operatorname{Lap}(q) + \operatorname{NLik}(p_{\theta}) + \operatorname{KL}(q || p_{\theta})$ 

### SUMMARY



M-step:  $\theta' = \theta - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial \theta}$ E-step:  $q_{\mathbf{y}}^{i \prime} = \frac{1}{Z_q(\mathbf{x}^i)} q_{\mathbf{y}}^{i} \exp \left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^{i}}\right]$ (update each  $r_{i,t}$  in practice)

### SUMMARY



 $\mathcal{J}(q, p_{\theta}) = \operatorname{Lap}(q) + \operatorname{NLik}(p_{\theta}) + \operatorname{KL}(q || p_{\theta})$ 

M-step: 
$$\theta' = \theta - \eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial \theta}$$
  
E-step:  $q_{\mathbf{y}}^{i\,\prime} = \frac{1}{Z_q(\mathbf{x}^i)} q_{\mathbf{y}}^{i} \exp\left[-\eta \frac{\partial \mathcal{J}(q, p_{\theta})}{\partial q_{\mathbf{y}}^{i}}\right]$ 

Theorem: Converges to a local optimum of  $\mathcal{J}(q, p_{\theta})$ 

(update each  $r_{i,t}$  in practice)

## EXPERIMENT SETTING

10 Languages (CoNLL-X and CoNLL-2007) 100 Randomly sampled labeled sentences Averaged over 10 sampling runs Universal POS Tags (Petrov et al. 2011) Second Order CRF Model  $f(y_t, y_{t-1}, y_{t-2}, \mathbf{x})$ 







GP



GP



 $\mathsf{GP} \qquad \mathsf{GP} \rightarrow \mathsf{CRF}$ 

12 10 8 SL Language DE ES PT SV ΕN DA EL IT NL GP  $GP \rightarrow CRF$ 





ninth run for





 $\mathsf{GP} \qquad \qquad \mathsf{GP} \rightarrow \mathsf{CRF}$ 

CRF







## CONCLUSION





### $\mathcal{J}(q, p_{\theta}) = \operatorname{Lap}(q) + \operatorname{NLik}(p_{\theta}) + \operatorname{KL}(q || p_{\theta})$

## CONCLUSION



any convex, differentiable regularizer

## CONCLUSION



any convex, differentiable regularizer

Code: https://code.google.com/p/pr-graph/