

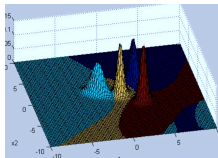
# Log-Linear Models, Logistic Regression and Conditional Random Fields

February 21, 2013

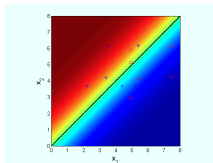
# Generative, Conditional and Discriminative

- Given  $\mathcal{D} = (x_t, y_t)_{t=1}^T$  sampled *iid* from unknown  $P(x, y)$
- Generative Learning (maximum likelihood Gaussians)
  - Choose family of functions  $p_\theta(x, y)$  parametrized by  $\theta$
  - Find  $\theta$  by **maximizing** likelihood:  $\prod_{t=1}^T p_\theta(x_i, y_i)$
  - Given  $x$ , output  $\hat{y} = \arg \max_y \frac{p_\theta(x, y)}{\sum_y p_\theta(x, y)}$
- Conditional Learning (logistic regression)
  - Choose family of functions  $p_\theta(y|x)$  parametrized by  $\theta$
  - Find  $\theta$  by **maximizing** conditional likelihood:  $\prod_{t=1}^T p_\theta(y_i|x_i)$
  - Given  $x$ , output  $\hat{y} = \arg \max_y p_\theta(y|x)$
- Discriminative Learning (support vector machines)
  - Choose family of functions  $y = f_\theta(x)$  parametrized by  $\theta$
  - Find  $\theta$  by **minimizing** classification error  $\sum_{t=1}^T \ell(y_i, f_\theta(x_i))$
  - Given  $x$ , output  $\hat{y} = f_\theta(x)$

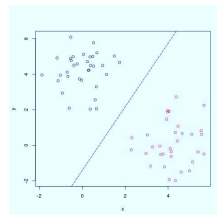
# Generative, Conditional and Discriminative



Generative



Conditional



Discriminative

# Generative: Maximum Entropy

Maximum entropy (or generally) minimum relative entropy  
 $\mathcal{RE}(p\|h) = \sum_y p(y) \ln \frac{p(y)}{h(y)}$  subject to linear constraints

$$\min_p \mathcal{RE}(p\|h) \text{ s.t. } \sum_y p(y) \mathbf{f}(y) = \mathbf{0}, \sum_y p(y) \mathbf{g}(y) \geq \mathbf{0}$$

Solution distribution looks like an exponential family model

$$p(y) = h(y) \exp \left( \boldsymbol{\theta}^\top \mathbf{f}(y) + \boldsymbol{\vartheta}^\top \mathbf{g}(y) \right) / Z(\boldsymbol{\theta}, \boldsymbol{\vartheta})$$

Maximize the dual (the negative log-partition) to get  $\boldsymbol{\theta}, \boldsymbol{\vartheta}$ .

$$\max_{\boldsymbol{\theta}, \boldsymbol{\vartheta} \geq \mathbf{0}} -\ln Z(\boldsymbol{\theta}, \boldsymbol{\vartheta}) = \max_{\boldsymbol{\theta}, \boldsymbol{\vartheta} \geq \mathbf{0}} -\ln \sum_y h(y) \exp \left( \boldsymbol{\theta}^\top \mathbf{f}(y) + \boldsymbol{\vartheta}^\top \mathbf{g}(y) \right)$$

# Generative: Exponential Family and Maximum Likelihood

All maximum entropy models give an exponential family form:

$$p(y) = h(y) \exp(\boldsymbol{\theta}^\top \mathbf{f}(y) - a(\boldsymbol{\theta}))$$

This is also a *log-linear model* over discrete  $y \in \Omega$  where  $|\Omega| = n$

$$p(y|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(y) \exp(\boldsymbol{\theta}^\top \mathbf{f}(y))$$

- Parameters are vector  $\boldsymbol{\theta} \in \mathbb{R}^d$
- Features are  $\mathbf{f} : \Omega \mapsto \mathbb{R}^d$  mapping each  $y$  to some vector
- Prior is  $h : \Omega \mapsto \mathbb{R}^+$  a fixed non-negative measure
- Partition function ensures that  $p(y|\boldsymbol{\theta})$  normalizes

$$Z(\boldsymbol{\theta}) = \sum_y h(y) \exp(\boldsymbol{\theta}^\top \mathbf{f}(y))$$

# Generative: Exponential Family and Maximum Likelihood

We are given some *iid* data  $y_1, \dots, y_T$  where  $y \in \{0, 1\}$ . If we wanted to find the best parameters of an exponential family distribution known as the Bernoulli distribution:

$$\begin{aligned} p(y|\theta) &= h(y) \exp(\boldsymbol{\theta}^\top \mathbf{f}(y) - a(\boldsymbol{\theta})) \\ &= \theta^y (1 - \theta)^{1-y} \end{aligned}$$

This is unsupervised generative learning

We simply find the  $\boldsymbol{\theta}$  that maximizes the likelihood

$$L(\boldsymbol{\theta}) = \prod_{t=1}^T p(y_t|\boldsymbol{\theta}) = \theta^{\sum_t y_t} (1 - \theta)^{T - \sum_t y_t}$$

Taking log then derivatives and setting to zero gives  $\theta = \frac{1}{T} \sum_t y_t$ .

## Conditional: Logistic Regression

Given input-output *iid* data  $(x_1, y_1), \dots, (x_T, y_T)$  where  $y \in \{0, 1\}$ .  
Binary logistic regression computes a probability for  $y = 1$  by

$$p(y = 1|x, \vartheta) = \frac{1}{1 + \exp(-\vartheta^\top \phi(x))}.$$

And the probability for  $p(y = 0|x, \theta) = 1 - p(y = 1|x, \theta)$ .  
This is supervised conditional learning.

We find the  $\theta$  that maximizes the *conditional* likelihood

$$L(\vartheta) = \prod_{t=1}^T p(y_t|x_t, \vartheta)$$

We can maximize this by doing gradient ascent.

Logistic regression is an example of a *log-linear model*.

## Conditional: Log-linear Models

Like an exponential family, but allow  $Z$ ,  $h$  and  $\mathbf{f}$  also depend on  $x$

$$p(y|x, \boldsymbol{\theta}) = \frac{1}{Z(x, \boldsymbol{\theta})} h(x, y) \exp\left(\boldsymbol{\theta}^\top \mathbf{f}(x, y)\right)$$

- Parameters are just one long vector  $\boldsymbol{\theta} \in \mathbb{R}^d$
- Functions  $\mathbf{f} : \Omega_x \times \Omega_y \mapsto \mathbb{R}^d$  map  $x, y$  to a vector
- Prior is  $h : \Omega_x \times \Omega_y \mapsto \mathbb{R}^+$  a fixed non-negative measure
- Partition function ensures that  $p(y|x, \boldsymbol{\theta})$  normalizes

To make a prediction, we simply output

$$\hat{y} = \arg \max_y p(y|x, \boldsymbol{\theta}).$$

Let's mimic (multi-class) logistic regression with this form.



# Conditional: Log-linear Models

In multi-class logistic regression, we have  $y \in \{1, \dots, n\}$ .

$$p(y|x, \theta) = \frac{1}{Z(x, \theta)} h(x, y) \exp(\theta^\top \mathbf{f}(x, y))$$

If  $\phi(x) \in \mathbb{R}^k$ , then  $\mathbf{f}(x, y) \in \mathbb{R}^{kn}$ .

Choose the following for the feature function

$$\mathbf{f}(x, y) = \left[ \delta[y = 1] \phi(x)^\top \quad \delta[y = 2] \phi(x)^\top \quad \dots \quad \delta[y = n] \phi(x)^\top \right]^\top.$$

If  $n = 2$  and  $h(x, y) = 1$ , get traditional binary logistic regression!

## Conditional: Log-linear Models

Rewrite binary logistic regression  $p(y = 1|x, \vartheta) = \frac{1}{1 + \exp(-\vartheta^\top \phi(x))}$  as a log-linear model with  $n = 2$ ,  $h(x, y) = 1$  and  $\mathbf{f}(x, y)$  as before

$$\begin{aligned}
 p(y|x, \boldsymbol{\theta}) &= \frac{h(x, y) \exp(\boldsymbol{\theta}^\top \mathbf{f}(x, y))}{Z(x, \boldsymbol{\theta})} \\
 &= \frac{\exp(\mathbf{f}(x, y)^\top \boldsymbol{\theta})}{\sum_{y=0}^1 \exp(\mathbf{f}(x, y)^\top \boldsymbol{\theta})} \\
 p(y = 1|x, \boldsymbol{\theta}) &= \frac{\exp([\mathbf{0} \ \phi(x)^\top] \boldsymbol{\theta})}{\exp([\phi(x)^\top \ \mathbf{0}] \boldsymbol{\theta}) + \exp([\mathbf{0} \ \phi(x)^\top] \boldsymbol{\theta})} \\
 &= \frac{1}{1 + \exp([\phi(x)^\top \ \mathbf{0}] \boldsymbol{\theta} - [\mathbf{0} \ \phi(x)^\top] \boldsymbol{\theta})}
 \end{aligned}$$

Can you see how to write  $\vartheta$  in terms of  $\boldsymbol{\theta}$ ?

# Conditional Random Fields (CRFs)

- Conditional random fields generalize maximum entropy
- Trained on *iid* data  $\{(x_1, y_1), \dots, (x_t, y_t)\}$
- A CRF is just a log-linear model with big  $n$

$$p(y|x_j, \theta) = \frac{1}{Z(x_j, \theta)} h(x_j, y) \exp(\theta^\top \mathbf{f}(x_j, y))$$

- Maximum conditional log-likelihood objective function is

$$J(\theta) = \sum_{j=1}^t \ln \frac{h(x_j, y_j)}{Z(x_j, \theta)} + \theta^\top \mathbf{f}(x_j, y_j) \quad (1)$$

- Regularized conditional maximum likelihood is

$$J(\theta) = \sum_{j=1}^t \ln \frac{h(x_j, y_j)}{Z(x_j, \theta)} + \theta^\top \mathbf{f}(x_j, y_j) - \frac{t\lambda}{2} \|\theta\|^2 \quad (2)$$

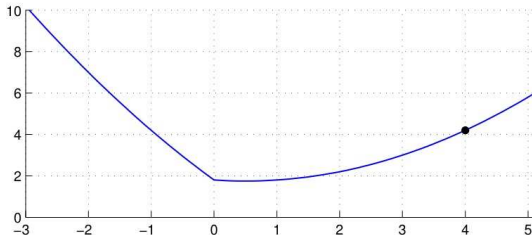
# Conditional Random Fields (CRFs)

- To train a CRF, we maximize (regularized) conditional likelihood
- Traditionally, maximum entropy, log-linear models and CRFs were trained using *majorization* (the EM algorithm is a majorization method)
- The algorithms were called *improved iterative scaling (IIS)* or *generalized iterative scaling (GIS)*
  - Maximum entropy [Jaynes '57]
  - Conditional random fields [Lafferty, et al. '01]
  - Log-linear models [Darroch & Ratcliff '72]

# Majorization

If cost function  $\theta^* = \arg \min_{\theta} C(\theta)$  has no closed form solution  
Majorization uses with a surrogate  $Q$  with closed form update  
to monotonically minimize the cost from an initial  $\theta_0$

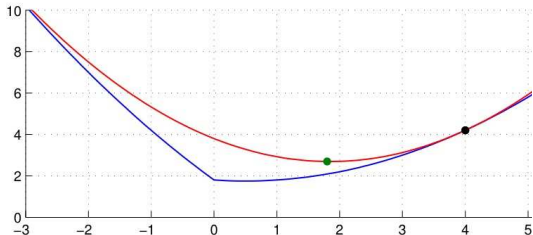
- Find bound  $Q(\theta, \theta_i) \geq C(\theta)$  where  $Q(\theta_i, \theta_i) = C(\theta_i)$
- Update  $\theta_{i+1} = \arg \min_{\theta} Q(\theta, \theta_i)$
- Repeat until converged



# Majorization

If cost function  $\theta^* = \arg \min_{\theta} C(\theta)$  has no closed form solution  
Majorization uses with a surrogate  $Q$  with closed form update  
to monotonically minimize the cost from an initial  $\theta_0$

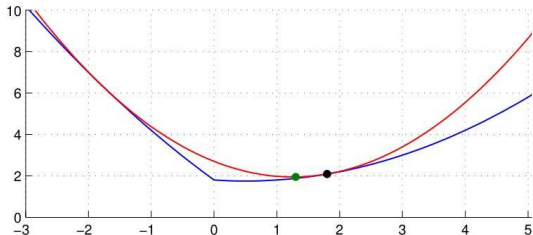
- Find bound  $Q(\theta, \theta_i) \geq C(\theta)$  where  $Q(\theta_i, \theta_i) = C(\theta_i)$
- Update  $\theta_{i+1} = \arg \min_{\theta} Q(\theta, \theta_i)$
- Repeat until converged



# Majorization

If cost function  $\theta^* = \arg \min_{\theta} C(\theta)$  has no closed form solution  
Majorization uses with a surrogate  $Q$  with closed form update  
to monotonically minimize the cost from an initial  $\theta_0$

- Find bound  $Q(\theta, \theta_i) \geq C(\theta)$  where  $Q(\theta_i, \theta_i) = C(\theta_i)$
- Update  $\theta_{i+1} = \arg \min_{\theta} Q(\theta, \theta_i)$
- Repeat until converged



# Majorization

IIS and GIS were preferred until [Wallach '03, Andrew & Gao '07]

Method	Iterations	LL Evaluations	Time (s)
IIS	$\geq 150$	$\geq 150$	$\geq 188.65$
Conjugate gradient (FR)	19	99	124.67
Conjugate gradient (PRP)	27	140	176.55
L-BFGS	22	22	29.72

Gradient descent appears to be faster

But newer majorization methods are faster still



# Gradient Ascent for CRFs

We have the following model

$$p(y|x, \boldsymbol{\theta}) = \frac{1}{Z(x, \boldsymbol{\theta})} h(x, y) \exp(\boldsymbol{\theta}^\top \mathbf{f}(x, y))$$

We want to maximize the conditional (log) likelihood:

$$\begin{aligned} \log L(\boldsymbol{\theta}) &= \sum_{t=1}^T \log p(y_t|x_t, \boldsymbol{\theta}) \\ &= \sum_{t=1}^T -\log Z(x_t, \boldsymbol{\theta}) + \log(h(x_t, y_t)) + \boldsymbol{\theta}^\top \mathbf{f}(x_t, y_t) \\ &= \text{const} - \sum_{t=1}^T \log Z(x_t, \boldsymbol{\theta}) + \boldsymbol{\theta}^\top \sum_{t=1}^T \mathbf{f}(x_t, y_t) \end{aligned}$$

Same as minimizing the sum of log partition functions plus linear!

## Gradient Ascent for CRFs

$$\begin{aligned}
\frac{\partial \log L}{\partial \boldsymbol{\theta}} &= \frac{\partial}{\partial \boldsymbol{\theta}} \left( \boldsymbol{\theta}^\top \sum_{t=1}^T \mathbf{f}(x_t, y_t) - \sum_{t=1}^T \log Z(x_t, \boldsymbol{\theta}) \right) \\
&= \sum_{t=1}^T \mathbf{f}(x_t, y_t) - \sum_{t=1}^T \frac{1}{Z(x_t, \boldsymbol{\theta})} \sum_y h(x_t, y) \frac{\partial}{\partial \boldsymbol{\theta}} \exp \left( \boldsymbol{\theta}^\top \mathbf{f}(x_t, y) \right) \\
&= \sum_{t=1}^T \mathbf{f}(x_t, y_t) - \sum_{t=1}^T \sum_y \frac{h(x_t, y)}{Z(x_t, \boldsymbol{\theta})} \exp \left( \boldsymbol{\theta}^\top \mathbf{f}(x_t, y) \right) \mathbf{f}(x_t, y) \\
&= \sum_{t=1}^T \mathbf{f}(x_t, y_t) - \sum_{t=1}^T \sum_y \mathbf{f}(x_t, y) p(y|x_t, \boldsymbol{\theta})
\end{aligned}$$

The gradient is the difference between the feature vectors at the true labels minus the expected feature vectors under the current distribution. To update,  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \frac{\partial \log L}{\partial \boldsymbol{\theta}}$ .

# Stochastic Gradient Ascent for CRFs

Given current  $\theta$ , update by taking a small step along the gradient

$$\theta \leftarrow \theta + \eta \frac{\partial \log L}{\partial \theta}.$$

We can use the full derivative:

$$\frac{\partial \log L}{\partial \theta} = \sum_{t=1}^T \mathbf{f}(x_t, y_t) - \sum_{t=1}^T \sum_y \mathbf{f}(x_t, y) p(y|x_t, \theta)$$

Or do stochastic gradient with only a single random datapoint  $t$ :

$$\frac{\partial \log L}{\partial \theta} = \mathbf{f}(x_t, y_t) - \sum_y \mathbf{f}(x_t, y) p(y|x_t, \theta)$$

# Better Majorization for CRFs

Recall log-linear model over discrete  $y \in \Omega$  where  $|\Omega| = n$

$$p(y|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(y) \exp(\boldsymbol{\theta}^\top \mathbf{f}(y))$$

- Parameters are vector  $\boldsymbol{\theta} \in \mathbb{R}^d$
- Features are  $\mathbf{f} : \Omega \mapsto \mathbb{R}^d$  mapping each  $y$  to some vector
- Prior is  $h : \Omega \mapsto \mathbb{R}^+$  a fixed non-negative measure
- Partition function ensures that  $p(y|\boldsymbol{\theta})$  normalizes

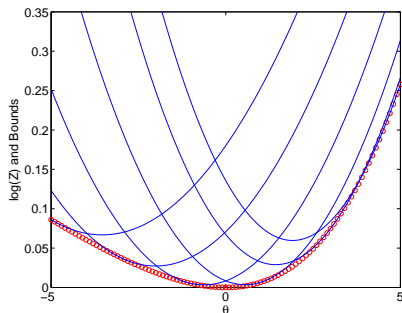
$$Z(\boldsymbol{\theta}) = \sum_y h(y) \exp(\boldsymbol{\theta}^\top \mathbf{f}(y))$$

Problem: it's ugly to minimize (unlike a quadratic function)

# Better Majorization for CRFs

The bound  $\ln Z(\theta) \leq \ln z + \frac{1}{2}(\theta - \tilde{\theta})^\top \Sigma (\theta - \tilde{\theta}) + (\theta - \tilde{\theta})^\top \mu$  is tight at  $\tilde{\theta}$  and holds for parameters given by

Input $\tilde{\theta}, \mathbf{f}(y), h(y) \forall y \in \Omega$
Init $z \rightarrow 0^+, \mu = \mathbf{0}, \Sigma = z\mathbf{I}$
For each $y \in \Omega$ {
$\alpha = h(y) \exp(\tilde{\theta}^\top \mathbf{f}(y))$
$\mathbf{l} = \mathbf{f}(y) - \mu$
$\Sigma \vdash = \frac{\tanh(\frac{1}{2} \ln(\alpha/z))}{2 \ln(\alpha/z)} \mathbf{ll}^\top$
$\mu \vdash = \frac{\alpha}{z+\alpha} \mathbf{l}$
$z \vdash = \alpha$ }
Output $z, \mu, \Sigma$



# Better Majorization for CRFs

## Bound Proof.

- 1) Start with bound  $\log(e^\theta + e^{-\theta}) \leq c\theta^2$  [Jaakkola & Jordan '99]
- 2) Prove scalar bound via Fenchel dual using  $\theta = \sqrt{\vartheta}$
- 3) Make bound multivariate  $\log(e^{\theta^\top \mathbf{1}} + e^{-\theta^\top \mathbf{1}})$
- 4) Handle scaling of exponentials  $\log(h_1 e^{\theta^\top \mathbf{f}_1} + h_2 e^{-\theta^\top \mathbf{f}_2})$
- 5) Add one term  $\log(h_1 e^{\theta^\top \mathbf{f}_1} + h_2 e^{-\theta^\top \mathbf{f}_2} + h_3 e^{-\theta^\top \mathbf{f}_3})$
- 6) Repeat extension for  $n$  terms



## Better Majorization for CRFs (Bound also Finds Gradient)

Init  $z \rightarrow 0^+, \boldsymbol{\mu} = \mathbf{0}, \boldsymbol{\Sigma} = z\mathbf{I}$

For each  $y \in \Omega$  {

$$\alpha = h(y) \exp(\tilde{\boldsymbol{\theta}}^\top \mathbf{f}(y))$$

$$\mathbf{l} = \mathbf{f}(y) - \boldsymbol{\mu}$$

$$\boldsymbol{\Sigma} + = \frac{\tanh(\frac{1}{2} \ln(\alpha/z))}{2 \ln(\alpha/z)} \mathbf{l} \mathbf{l}^\top$$

$$\boldsymbol{\mu} + = \frac{\alpha}{z+\alpha} \mathbf{l}$$

$$z + = \alpha \quad \}$$

Output  $z, \boldsymbol{\mu}, \boldsymbol{\Sigma}$

$$\text{Recall gradient } \frac{\partial \log L}{\partial \boldsymbol{\theta}} = \sum_{t=1}^T \mathbf{f}(x_t, y_t) - \sum_{t=1}^T \sum_y \mathbf{f}(x_t, y) p(y|x_t, \boldsymbol{\theta})$$

The bound's  $\boldsymbol{\mu}$  give part of gradient (can skip  $\boldsymbol{\Sigma}$  updates).

$$\boldsymbol{\mu} = \sum_y \mathbf{f}(x_t, y) p(y|x_t, \boldsymbol{\theta})$$

## Better Majorization for CRFs

Input $x_j, y_j$ and functions $h_{x_j}, \mathbf{f}_{x_j}$ for $j=1, \dots, t$ Input regularizer $\lambda \in \mathbb{R}^+$
Initialize $\theta_0$ anywhere and set $\tilde{\theta} = \theta_0$ While not converged <ul style="list-style-type: none"> <li>For <math>j = 1</math> to <math>t</math> compute bound for <math>\mu_j, \Sigma_j</math> from <math>h_{x_j}, \mathbf{f}_{x_j}, \tilde{\theta}</math></li> <li>Set <math>\tilde{\theta} = \arg \min_{\theta} \sum_j \frac{1}{2} (\theta - \tilde{\theta})^\top (\Sigma_j + \lambda \mathbf{I}) (\theta - \tilde{\theta}) + \sum_j \theta^\top (\mu_j - \mathbf{f}_{x_j}(y_j) + \lambda \tilde{\theta})</math></li> </ul>
Output $\hat{\theta} = \tilde{\theta}$

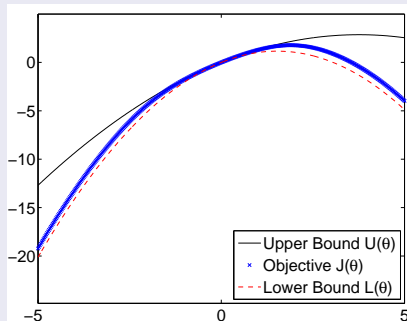
### Theorem

If  $\|\mathbf{f}(x_j, y)\| \leq r$  get  $J(\hat{\theta}) - J(\theta_0) \geq (1 - \epsilon) \max_{\theta} (J(\theta) - J(\theta_0))$   
 within  $\left\lceil \ln(1/\epsilon) / \ln \left( 1 + \frac{\lambda \log n}{2r^2 n} \right) \right\rceil$  steps



# Convergence Proof

Proof.



**Figure:** Quadratic bounding sandwich. Compare upper and lower bound curvatures to bound maximum # of iterations.



## Experiments - Multi-Class Classification &amp; Linear Chains

Data-set	SRBCT		Tumors		Text		SecStr		CoNLL		PennTree	
Size	$n = 4$ $t = 83$ $d = 9236$ $\lambda = 10^1$		$n = 26$ $t = 308$ $d = 390260$ $\lambda = 10^1$		$n = 2$ $t = 1500$ $d = 23922$ $\lambda = 10^2$		$n = 2$ $t = 83679$ $d = 632$ $\lambda = 10^1$		$m = 9$ $t = 1000$ $d = 33615$ $\lambda = 10^1$		$m = 45$ $t = 1000$ $d = 14175$ $\lambda = 10^1$	
Algorithm	time	iter	time	iter	time	iter	time	iter	time	iter	time	iter
LBFGS	6.10	42	3246.83	8	15.54	7	881.31	47	25661.54	17	62848.08	7
Grad	7.27	43	18749.15	53	153.10	69	1490.51	79	93821.72	12	156319.31	12
Congrad	40.61	100	14840.66	42	57.30	23	667.67	36	88973.93	23	76332.39	18
Bound	<b>3.67</b>	<b>8</b>	<b>1639.93</b>	<b>4</b>	<b>6.18</b>	<b>3</b>	<b>27.97</b>	<b>9</b>	<b>16445.93</b>	<b>4</b>	<b>27073.42</b>	<b>2</b>

**Table:** Time in seconds and iterations to match LBFGS solution for multi-class logistic regression (on SRBCT, Tumors, Text and SecStr data-sets where  $n$  is the number of classes) and Markov CRFs (on CoNLL and PennTree data-sets, where  $m$  is the number of classes). Here,  $t$  is the number of samples,  $d$  is the dimensionality of the feature vector and  $\lambda$  is the cross-validated regularization setting.

## Experiments - Linear Chains

Model	Error	oov Error
Hidden Markov Model	5.69%	45.59%
Maximum Entropy Markov Model	6.37%	54.61%
Conditional Random Field	5.55%	48.05%

**Table:** Accuracy on Penn tree-bank data-set for parts-of-speech tagging with training on half of the 1.1 million word corpus. Note, the oov rate is the error rate on out-of-vocabulary words.

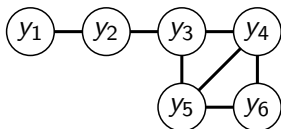
Parts of speech data-set where there are 45 labels per word, e.g.

PRP	VBD	DT	NN	IN	DT	NN
I	saw	the	man	with	the	telescope

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{Z} \psi(y_1, y_2) \psi(y_2, y_3) \psi(y_3, y_4) \psi(y_4, y_5) \psi(y_5, y_6) \psi(y_6, y_7)$$

How big is  $y$ ? Recall graphical models for large spaces...

# Bounding Graphical Models with Large $n$

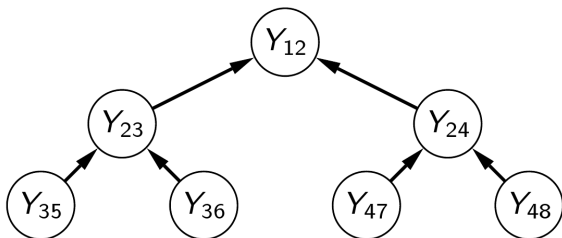


- Each iteration is  $\mathcal{O}(tn)$ , but what if  $n$  is large?
- Graphical model: an undirected graph  $G$  representing a distribution  $p(Y)$  where  $Y = \{y_1, \dots, y_n\}$  and  $y_i \in \mathbb{Z}$
- $p(Y)$  factorizes as product of  $\{\psi_1, \dots, \psi_C\}$  functions over  $\{Y_1, \dots, Y_C\}$  subsets of variables over the maximal cliques of  $G$

$$p(y_1, \dots, y_n) = \frac{1}{Z} \prod_{c \in C} \psi_c(Y_c)$$

- E.g.  $p(y_1, \dots, y_6) = \psi(y_1, y_2)\psi(y_2, y_3)\psi(y_3, y_4, y_5)\psi(y_4, y_5, y_6)$

# Bounding Graphical Models with Large $n$



- Instead of enumerating over all  $n$ , exploit graphical model
- Build junction tree and run a *Collect* algorithm
- Useful for computing  $Z(\theta)$ ,  $\frac{\partial \log Z(\theta)}{\partial \theta}$  and  $\Sigma$  efficiently
- Bound needs  $\mathcal{O}(t \sum_c |Y_c|)$  rather than  $\mathcal{O}(tn)$
- For an HMM, this is  $\mathcal{O}(TM^2)$  instead of  $\mathcal{O}(M^T)$

Bounding Graphical Models with Large  $n$ 

for  $c = 1, \dots, m$  {

$$Y_{both} = Y_c \cap Y_{pa(c)}; \quad Y_{solo} = Y_c \setminus Y_{pa(c)}$$

for each  $u \in Y_{both}$  {

$$\text{initialize } z_{c|x} \leftarrow 0^+, \quad \mu_{c|x} = \mathbf{0}, \quad \Sigma_{c|x} = z_{c|x} \mathbf{I}$$

for each  $v \in Y_{solo}$  {

$$w = u \otimes v; \quad \alpha_w = h_c(w) e^{\tilde{\theta}^\top f_c(w)} \prod_{b \in ch(c)} z_{b|w}$$

$$\mathbf{l}_w = \mathbf{f}_c(w) - \mu_{c|u} + \sum_{b \in ch(c)} \mu_{b|w}$$

$$\Sigma_{c|u} += \sum_{b \in ch(c)} \Sigma_{b|w} + \frac{\tanh\left(\frac{1}{2} \ln\left(\frac{\alpha_w}{z_{c|u}}\right)\right)}{2 \ln\left(\frac{\alpha_w}{z_{c|u}}\right)} \mathbf{l}_w \mathbf{l}_w^\top$$

$$\mu_{c|u} += \frac{\alpha_w}{z_{c|u} + \alpha_w} \mathbf{l}_w; \quad z_{c|u} += \alpha_w \quad \}}}$$