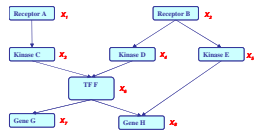


# Large Margin Structured Models

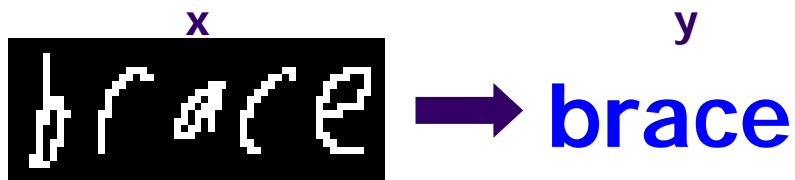
Probabilistic Graphical Models (10-708)

Lecture # 21  
11/28/2007  
Eric Xing

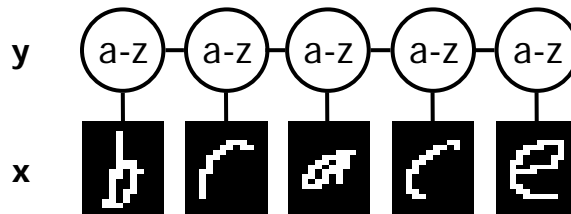


1

## OCR example



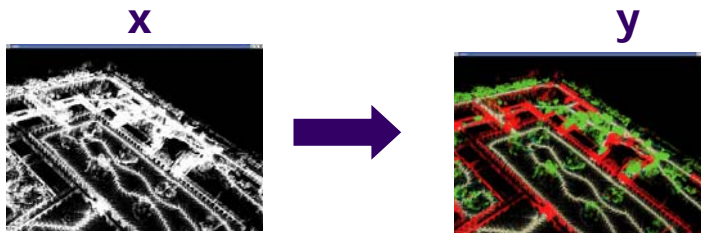
Sequential structure



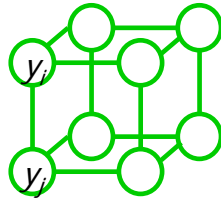
Eric Xing

2

# Image segmentation example



Spatial structure



Eric Xing

3

# Structured models

$$h(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} s(\mathbf{x}, \mathbf{y}) \quad \leftarrow \text{scoring function}$$

space of feasible outputs

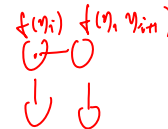
given  $\mathbf{x}$   
to get  $\mathbf{y}$ .

score =  $\frac{1}{2\sigma} \exp(w^T f)$   $\mathbf{y}^* = \arg \max_{\mathbf{y} \in \mathcal{Y}} \text{score}(\mathbf{x}, \mathbf{y})$

Assumptions:

$$\text{score}(\mathbf{x}, \mathbf{y}) = \mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}) = \sum_p \mathbf{w}^T \mathbf{f}(\mathbf{x}_p, \mathbf{y}_p)$$

linear combination of features



sum of part scores:

- index  $p$  represents a part in the structure

Eric Xing

4

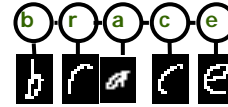
# Learning w



- Training examples  $(\mathbf{x}_i, \mathbf{y}_i)$

- Probabilistic approach:

$$P_{\mathbf{w}}(\mathbf{y} | \mathbf{x}) = \frac{1}{Z_{\mathbf{w}}(\mathbf{x})} \exp\{\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y})\}$$



- Computing  $Z_{\mathbf{w}}(\mathbf{x})$  can be NP-complete

- Tractable models but intractable estimation



- Large margin approach:

- Exact and efficient when prediction is tractable

# Learning Strategy



- Recall that in CRF

- We predict based on:

$$y^* | x = \arg \max_y p_{\theta}(y | x) = \frac{1}{Z(\theta, x)} \exp\left\{\sum_c \theta_c f_c(x, y_c)\right\}$$

- And we learn based on:

$$\theta_c^* | \{y_n, x_n\} = \arg \max_{\theta_c} \prod_n p_{\theta}(y_n | x_n) = \prod_n \frac{1}{Z(\theta, x_n)} \exp\left\{\sum_c \theta_c f_c(x_n, y_{n,c})\right\}$$

- MaxMargin:

- We predict based on:

$$y^* | x = \arg \max_y \sum_c \theta_c f_c(x, y_c) = \arg \max_y w^T F(x, y)$$

- And we learn based on:

$$w^* | \{y_n, x_n\} = \arg \max_w \left( \max_{y_n \neq y'_n, \forall n} w^T (F(y_n, x_n) - F(y'_n, x_n)) \right)$$

## MLE of Feature Based UGMs



- Scaled likelihood function

$$\begin{aligned}\tilde{\ell}(\theta; D) &= \ell(\theta; D) / N = \frac{1}{N} \sum_n \log p(x_n | \theta) \\ &= \sum_x \tilde{p}(x) \log p(x | \theta) \\ &= \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \log Z(\theta)\end{aligned}$$

$\propto \frac{1}{Z(\theta)} \exp(\theta^T f(x))$

- Instead of optimizing this objective directly, we attack its lower bound

- The logarithm has a linear upper bound ...  
 $\log Z(\theta) \leq \mu Z(\theta) - \log \mu - 1$
- This bound holds for all  $\mu$ , in particular, for  $\mu = Z^{-1}(\theta^{(t)})$

- Thus we have

$$\tilde{\ell}(\theta; D) \geq \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{Z(\theta)}{Z(\theta^{(t)})} - \log Z(\theta^{(t)}) + 1$$



Eric Xing

7

## Generalized Iterative Scaling (GIS)



- Lower bound of scaled loglikelihood

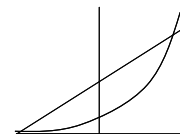
$$\tilde{\ell}(\theta; D) \geq \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{Z(\theta)}{Z(\theta^{(t)})} - \log Z(\theta^{(t)}) + 1$$

- Define  $\Delta\theta_i^{(t)} \stackrel{\text{def}}{=} \theta_i - \theta_i^{(t)}$

$$\begin{aligned}\tilde{\ell}(\theta; D) &\geq \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{1}{Z(\theta^{(t)})} \sum_x \exp\left\{\sum_i \theta_i f_i(x)\right\} - \log Z(\theta^{(t)}) + 1 \\ &= \sum_i \theta_i \sum_x \tilde{p}(x) f_i(x) - \frac{1}{Z(\theta^{(t)})} \sum_x \exp\left\{\sum_i \theta_i^{(t)} f_i(x)\right\} \exp\left\{\sum_i \Delta\theta_i^{(t)} f_i(x)\right\} - \log Z(\theta^{(t)}) + 1 \\ &= \sum_i \theta_i \sum_x \tilde{p}(x) f_i(x) - \sum_x p(x | \theta^{(t)}) \exp\left\{\sum_i \Delta\theta_i^{(t)} f_i(x)\right\} - \log Z(\theta^{(t)}) + 1\end{aligned}$$

- Relax again

- Assume  $f_i(x) \geq 0$ ,  $\sum_i f_i(x) = 1$
- Convexity of exponential:  $\exp\left(\sum_i \pi_i x_i\right) \leq \sum_i \pi_i \exp(x_i)$



- We have:

$$\tilde{\ell}(\theta; D) \geq \sum_i \theta_i \sum_x \tilde{p}(x) f_i(x) - \sum_x p(x | \theta^{(t)}) \sum_i f_i(x) \exp(\Delta\theta_i^{(t)}) - \log Z(\theta^{(t)}) + 1 \stackrel{\text{def}}{=} \Lambda(\theta)$$

Eric Xing

8

# GIS



- Lower bound of scaled loglikelihood

$$\tilde{\ell}(\theta; D) \geq \sum_i \theta_i \sum_x \tilde{p}(x) f_i(x) - \sum_x p(x | \theta^{(t)}) \sum_i f_i(x) \exp(\Delta \theta_i^{(t)}) - \log Z(\theta^{(t)}) + 1 \stackrel{\text{def}}{=} \Lambda(\theta)$$

- Take derivative:  $\frac{\partial \Lambda}{\partial \theta_i} = \sum_x \tilde{p}(x) f_i(x) - \exp(\Delta \theta_i^{(t)}) \sum_x p(x | \theta^{(t)}) f_i(x)$

- Set to zero

$$e^{\Delta \theta_i^{(t)}} = \frac{\sum_x \tilde{p}(x) f_i(x)}{\sum_x p(x | \theta^{(t)}) f_i(x)} = \frac{\sum_x \tilde{p}(x) f_i(x)}{\sum_x p^{(t)}(x) f_i(x)} Z(\theta^{(t)})$$

- where  $p^{(t)}(x)$  is the unnormalized version of  $p(x | \theta^{(t)})$

- Update  $\theta_i^{(t+1)} = \theta_i^{(t)} + \Delta \theta_i^{(t)} \Rightarrow p^{(t+1)}(x) = p^{(t)}(x) e^{\Delta \theta_i^{(t)} f_i(x)}$

$$\begin{aligned} p^{(t+1)}(x) &= \frac{p^{(t)}(x)}{Z(\theta^{(t)})} \prod_i \left( \frac{\sum_x \tilde{p}(x) f_i(x)}{\sum_x p^{(t)}(x) f_i(x)} Z(\theta^{(t)}) \right)^{f_i(x)} \\ \Rightarrow &= \frac{p^{(t)}(x)}{Z(\theta^{(t)})} \prod_i \left( \frac{\sum_x \tilde{p}(x) f_i(x)}{\sum_x p^{(t)}(x) f_i(x)} \right)^{f_i(x)} (Z(\theta^{(t)}))^{\sum_i f_i(x)} \\ &= p^{(t)}(x) \prod_i \left( \frac{\sum_x \tilde{p}(x) f_i(x)}{\sum_x p^{(t)}(x) f_i(x)} \right)^{f_i(x)} \end{aligned}$$

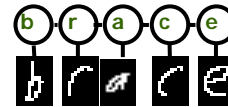
Eric Xing

9

# Learning w



- Training examples  $(\mathbf{x}_i, \mathbf{y}_i)$



- Probabilistic approach:

$$P_{\mathbf{w}}(\mathbf{y} | \mathbf{x}) = \frac{1}{Z_{\mathbf{w}}(\mathbf{x})} \exp\{\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y})\}$$

- Computing  $Z_{\mathbf{w}}(\mathbf{x})$  can be NP-complete
  - Tractable models but intractable estimation

- Large margin approach:
  - Exact and efficient when prediction is tractable

Eric Xing

10

## OCR Example



- We want:

$$\operatorname{argmax}_{\text{word}} \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{word}) = \text{"brace"}$$

- Equivalently:

$$\mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"brace"}) > \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"aaaaa"})$$

$$\mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"brace"}) > \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"aaaab"})$$

...

$$\mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"brace"}) > \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"zzzzz"})$$

} a lot!

$$|z^5 - 1|$$

Eric Xing

11

## Large Margin Estimation



- Given training example  $(\mathbf{x}, \mathbf{y}^*)$ , we want:

$$\operatorname{arg max}_{\mathbf{y}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) = \mathbf{y}^*$$

$$\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}^*) > \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) \quad \forall \mathbf{y} \neq \mathbf{y}^*$$

$$\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}^*) \geq \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) + \gamma \ell(\mathbf{y}^*, \mathbf{y}) \quad \forall \mathbf{y}$$

- Maximize margin  $\gamma$
- Mistake weighted margin:  $\gamma \ell(\mathbf{y}^*, \mathbf{y})$

$$\ell(\mathbf{y}^*, \mathbf{y}) = \sum_i I(y_i^* \neq y_i) \quad \# \text{ of mistakes in } \mathbf{y}$$

Eric Xing

\*Taskar et al, 03

# Large Margin Estimation



- Recall from SVMs:
  - Maximizing margin  $\gamma$  is equivalent to minimizing the square of the L2-norm of the weight vector  $\mathbf{w}$ :
- New objective function:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

$$s.t. \mathbf{w}^T \mathbf{f}(x_i, y_i) \geq \mathbf{w}^T \mathbf{f}(x_i, y'_i) + \ell(y_i, y'_i), \quad \forall i, y'_i \in \mathcal{Y}_i$$

# Min-max Formulation

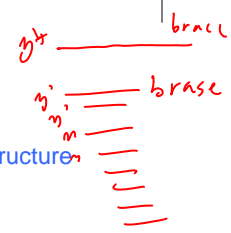


- Brute force enumeration of constraints:

$$\min \frac{1}{2} \|\mathbf{w}\|^2$$

$$\mathbf{w}^T \mathbf{f}(\mathbf{x}, y^*) \geq \mathbf{w}^T \mathbf{f}(\mathbf{x}, y) + \ell(y^*, y), \quad \forall y$$

- The constraints are exponential in the size of the structure
- Alternative: min-max formulation
  - add only the most violated constraint



$$y' = \arg \max_{y \neq y^*} [\mathbf{w}^T \mathbf{f}(x^i, y) + \ell(y^i, y)]$$

$$\text{add to QP : } \mathbf{w}^T \mathbf{f}(x^i, y^i) \geq \mathbf{w}^T \mathbf{f}(x^i, y') + \ell(y^i, y')$$

- Handles more general loss functions
- Only polynomial # of constraints needed

# Min-max formulation

$$\min \frac{1}{2} \|\mathbf{w}\|^2$$

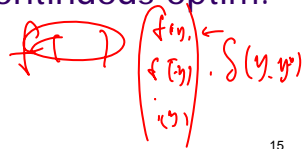
$$\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}^*) \geq \max_{\mathbf{y} \neq \mathbf{y}^*} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) + \ell(\mathbf{y}^*, \mathbf{y})$$

- Key step: convert the maximization in the constraint from discrete to continuous
  - This enables us to plug it into a QP

$$\max_{\mathbf{y} \neq \mathbf{y}^*} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) + \ell(\mathbf{y}^*, \mathbf{y}) \iff \max_{\mathbf{z} \in \mathcal{Z}} (\mathbf{F}^\top \mathbf{w} + \ell)^\top \mathbf{z}$$

discrete optim. continuous optim.

- How to do this conversion?
  - Linear chain example in the next slides →



Eric Xing

15

# y ⇒ z map for linear chain structures

OCR example:  $\mathbf{y} = \text{'ABABB'}$ ;

$\mathbf{z}$ 's are the indicator variables for the corresponding classes (alphabet)

	$z_1(m)$	$z_2(m)$	$z_3(m)$	$z_4(m)$	$z_5(m)$
A	1	0	1	0	0
B	0	1	0	1	1
⋮	⋮	⋮	⋮	⋮	⋮
B	0	0	0	0	0

	$z_{12}(m, n)$	$z_{23}(m, n)$	$z_{34}(m, n)$	$z_{45}(m, n)$
A	0	0	0	0
B	1	0	0	0
⋮	⋮	⋮	⋮	⋮
B	0	0	0	0

A	B	.	B
A	B	.	B
A	B	.	B
A	B	.	B

Eric Xing

16



## y ⇒ z map for linear chain structures

Rewriting the maximization function in terms of indicator variables:

$$\max_{\mathbf{z}} \sum_{j,m} z_j(m) [\mathbf{w}^\top \mathbf{f}_{\text{node}}(\mathbf{x}_j, m) + \ell_j(m)] + \sum_{jk,m,n} z_{jk}(m,n) [\mathbf{w}^\top \mathbf{f}_{\text{edge}}(\mathbf{x}_{jk}, m, n) + \ell_{jk}(m,n)]$$

$\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$   
 $(\mathbf{F}^\top \mathbf{w} + \ell)^\top \mathbf{z}$

$z_j(m) \geq 0; z_{jk}(m,n) \geq 0;$

normalization  $\sum_m z_j(m) = 1$

agreement  $\sum_n z_{jk}(m,n) = z_j(m)$

$\mathbf{z}_i(l) = \begin{pmatrix} u \\ v \\ y \end{pmatrix}$

$\max_{\mathbf{z}} (\mathbf{F}^\top \mathbf{w} + \ell)^\top \mathbf{z}$   
 $\mathbf{Az} = \mathbf{b}$

$z_k(n)$ 

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

  
 $z_j(m)$ 

0
0
1
0

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

  
 $z_{jk}(m,n)$ 

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

Eric Xing

17

## Min-max formulation

- Original problem:

$$\min \frac{1}{2} \|\mathbf{w}\|^2$$

$$\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}^*) \geq \max_{\mathbf{y}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) + \ell(\mathbf{y}^*, \mathbf{y})$$

- Transformed problem:

$$\min \frac{1}{2} \|\mathbf{w}\|^2$$

$$\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}^*) \geq \max_{\substack{\mathbf{z} \geq 0; \\ \mathbf{Az} = \mathbf{b};}} \mathbf{q}^\top \mathbf{z} \quad \text{where } \mathbf{q}^\top = \mathbf{w}^\top \mathbf{F} + \ell^\top$$

- Has integral solutions  $\mathbf{z}$  for chains, trees
- Can be fractional for untriangulated networks

Eric Xing

18

## Min-max formulation



- Using strong Lagrangian duality:  
(beyond the scope of this lecture)

$$\max_{\substack{z \geq 0; \\ \mathbf{A}z = \mathbf{b};}} \mathbf{q}^\top \mathbf{z} = \min_{\mathbf{A}^\top \boldsymbol{\mu} \geq \mathbf{q}} \mathbf{b}^\top \boldsymbol{\mu}$$

- Use the result above to minimize jointly over  $\mathbf{w}$  and  $\boldsymbol{\mu}$ :

$$\begin{aligned} \min_{\mathbf{w}, \boldsymbol{\mu}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}^*) \geq \mathbf{b}^\top \boldsymbol{\mu}; \\ & \mathbf{A}^\top \boldsymbol{\mu} \geq \mathbf{q}; \end{aligned}$$

Eric Xing

19

## Min-max formulation



$$\begin{aligned} \min_{\mathbf{w}, \boldsymbol{\mu}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}^*) \geq \mathbf{b}^\top \boldsymbol{\mu}; \\ & \mathbf{A}^\top \boldsymbol{\mu} \geq (\mathbf{w}^\top \mathbf{F} + \boldsymbol{\ell})^\top \end{aligned}$$

- Formulation produces compact QP for
  - Low-treewidth Markov networks
  - Associative Markov networks
  - Context free grammars
  - Bipartite matchings
  - Any problem with compact LP inference

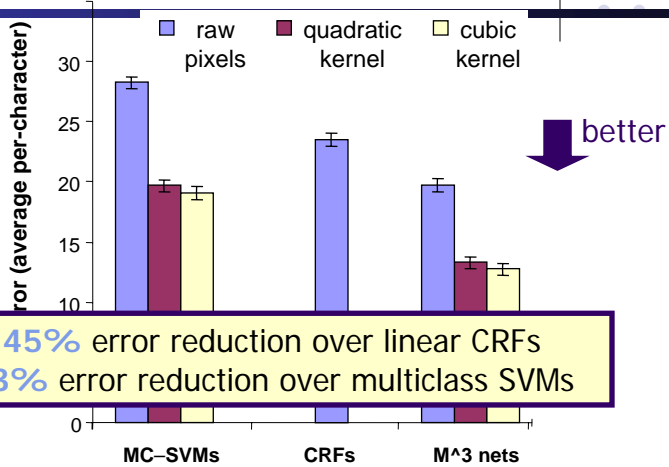
Eric Xing

20

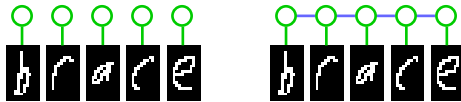
## Results: Handwriting Recognition

Length: ~8 chars  
 Letter: 16x8 pixels  
 10-fold Train/Test  
 5000/50000 letters  
 600/6000 words

Models:  
 Multiclass-SVMs  
 CRFs  
 M<sup>3</sup> nets



45% error reduction over linear CRFs  
 33% error reduction over multiclass SVMs

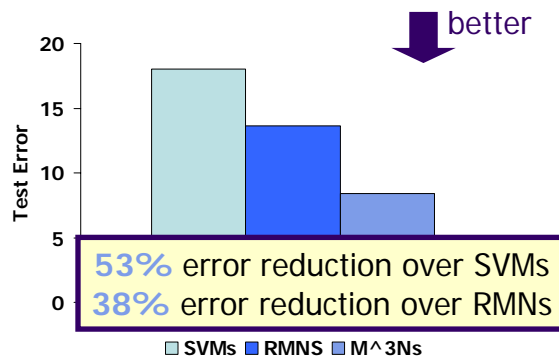
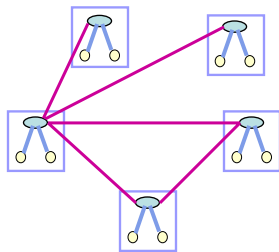


\*Crammer & Singer 01

21

## Results: Hypertext Classification

- WebKB dataset
  - Four CS department websites: 1300 pages/3500 links
  - Classify each page: faculty, course, student, project, other
  - Train on three universities/test on fourth



53% error reduction over SVMs  
 38% error reduction over RMNs

\*Taskar et al 02

22

## Summary



- Structured Maximum Margin Networks
  - Concise representation
  - Efficient QP algorithms
  - Applications:
    - Sequences
    - Trees
    - Matchings ....
  - Strong empirical results
- Acknowledgments:
  - Adopted from two different presentations given by Ben Taskar.

Eric Xing

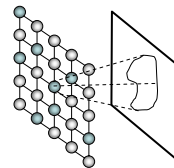
23

## Open Problems



- Unsupervised CRF learning and MaxMargin Learning

- We want to recognize a pattern that is maximally different from the rest!



- What does margin or conditional likelihood mean in these cases?  
Given only  $\{X_n\}$ , how can we define the cost function?

$$p_{\theta}(y|x) = \frac{1}{Z(\theta, x)} \exp\left\{\sum_c \theta_c f_c(x, y_c)\right\}$$

$$\text{margin} = w^T (F(y_n, x_n) - F(y'_n, x_n))$$

- Algorithmic challenge

Eric Xing

24