

# Maximum Margin Markov Networks and Constraint Generation continued

Optimization - 10725

Carlos Guestrin

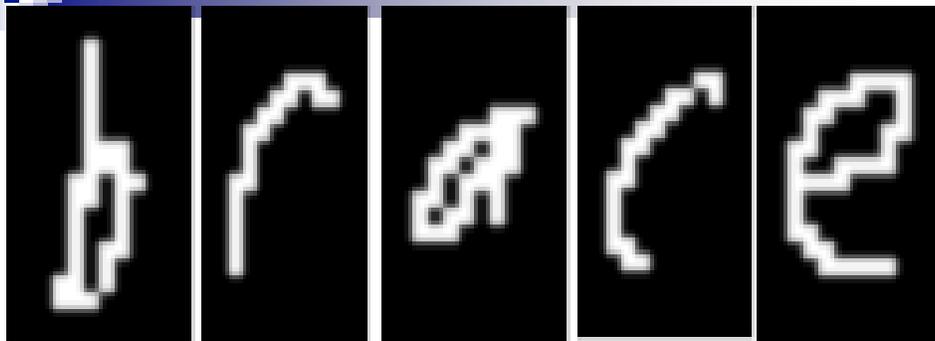
Carnegie Mellon University

February 18<sup>th</sup>, 2008

©2008 Carlos Guestrin

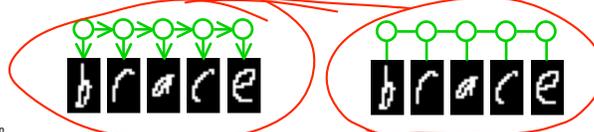
1

## Handwriting Recognition 2



Graphical models: HMMs, MNs

Linear in length



©2008 Carlos Guestrin

2

An example of a CRF  $\leftarrow$  Conditional random field

# Chain Markov Net (aka CRF\*)

potentials

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_i \phi(\mathbf{x}_i, y_i) \prod_i \phi(y_i, y_{i+1})$$

$\phi(\mathbf{x}_i, y_i) = \exp\{\sum_{\alpha} w_{\alpha} f_{\alpha}(\mathbf{x}_i, y_i)\}$  *Similar to logistic regression*

$\phi(y_i, y_{i+1}) = \exp\{\sum_{\beta} w_{\beta} f_{\beta}(y_i, y_{i+1})\}$

$f_{\beta}(y, y') = I(y='z', y'='a')$

$f_{\alpha}(\mathbf{x}, y) = I(x_p=1, y='z')$

©2008 Carlos Guestrin \*Lafferty et al. 01 3

## CRF - short notation

$e^{w_1 f_1 + w_2 f_2} = e^{w_1 f_1 + w_2 f_2}$

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_i \phi(\mathbf{x}_i, y_i) \prod_i \phi(y_i, y_{i+1}) = \frac{1}{Z(\mathbf{x})} \exp\{\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y})\}$$

$\phi(\mathbf{x}_i, y_i) = \exp\{\sum_{\alpha} w_{\alpha} f_{\alpha}(\mathbf{x}_i, y_i)\}$

$\phi(y_i, y_{i+1}) = \exp\{\sum_{\beta} w_{\beta} f_{\beta}(y_i, y_{i+1})\}$

$\mathbf{w} = \begin{bmatrix} w_{\alpha_1} \\ \vdots \\ w_{\alpha_n} \\ w_{\beta_1} \\ \vdots \\ w_{\beta_m} \end{bmatrix}$

$\mathbf{f}(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} f_{\alpha_1}(\mathbf{x}_1, y_1) \\ \vdots \\ f_{\alpha_n}(\mathbf{x}_n, y_n) \\ f_{\beta_1}(y_1, y_2) \\ \vdots \\ f_{\beta_m}(y_{i-1}, y_i) \end{bmatrix}$

©2008 Carlos Guestrin \*Lafferty et al. 01 4

$\operatorname{argmax}_y f(y) = \operatorname{argmax}_y \log f(y)$

# Max (Conditional) Likelihood

D

$\mathbf{x}^1, \mathbf{t}(\mathbf{x}^1)$   
image for forward ... label for first

$\mathbf{x}^m, \mathbf{t}(\mathbf{x}^m)$

**Estimation**

maximize  $\underline{w}$

$\sum_{\mathbf{x} \in D} \log P_w(\underline{t}(\mathbf{x}) | \underline{x})$

parameters  $w$

**Classification**

test case  $\underline{x}_\ell$

$\underline{\operatorname{argmax}}_y P_w(y | \underline{x}_\ell)$

**f(x, y)**

0-0-0-0-0

At classification time

$\operatorname{argmax}_y P_w(y | x) = \operatorname{argmax}_y \frac{1}{z(x)} e^{w'f(y, x)}$  doesn't depend y irrelevant for decision

$= \operatorname{argmax}_y \log \frac{1}{z(x)} e^{w'f(y, x)} = \operatorname{argmax}_y w'f(y, x) - \log z(x)$

$= \operatorname{argmax}_y w'f(y, x)$

why do I care about  $z(x)$

©2008 Carlos Guestrin 5

## OCR Example

- We want: given input

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

- Equivalently: true word wins over all other possibilities

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

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

...

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

...

}

number of constraints is exponential in length of word

strictly greater than?  
ho w???

©2008 Carlos Guestrin 6

# Max Margin Estimation

- **Goal:** find  $\mathbf{w}$  such that

$$\mathbf{w}^\top \mathbf{f}(\mathbf{x}, t(\mathbf{x})) > \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) \quad \forall \mathbf{x} \in \mathcal{D} \quad \forall \mathbf{y} \neq t(\mathbf{x})$$

$$\mathbf{w}^\top [\mathbf{f}(\mathbf{x}, t(\mathbf{x})) - \mathbf{f}(\mathbf{x}, \mathbf{y})] > 0$$

$$\mathbf{w}^\top \Delta \mathbf{f}_x(\mathbf{y}) > 0$$

©2008 Carlos Guestrin

7

# Not all margins are equal

- **Goal:** find  $\mathbf{w}$  such that

$$\mathbf{w}^\top \Delta \mathbf{f}_x(\mathbf{y}) \geq \gamma \quad \forall \mathbf{x} \in \mathcal{D} \quad \forall \mathbf{y} \neq t(\mathbf{x})$$

- Gain over  $\mathbf{y}$  grows with # of mistakes in  $\mathbf{y}$ :  $\Delta t_x(\mathbf{y})$

$$\Delta t_{\text{brace}}(\text{"craze"})$$

$$\Delta t_{\text{brace}}(\text{"zzzzz"})$$

$$\mathbf{w}^\top \Delta \mathbf{f}_{\text{brace}}(\text{"craze"})$$

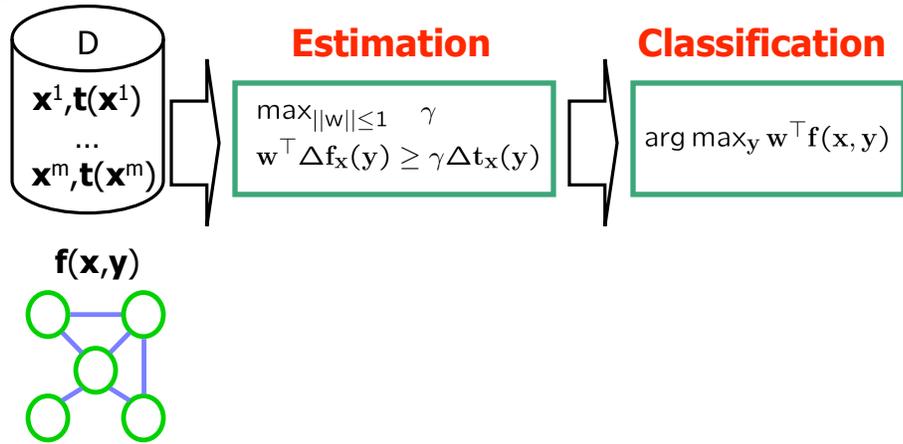
$$\mathbf{w}^\top \Delta \mathbf{f}_{\text{brace}}(\text{"zzzzz"})$$

©2008 Carlos Guestrin

8

# Maximum Margin Markov Nets

[Taskar, Guestrin, Koller '03]



■ BTW. Just like SVMs, there are “non-linearly separable” cases, must add slack variables...

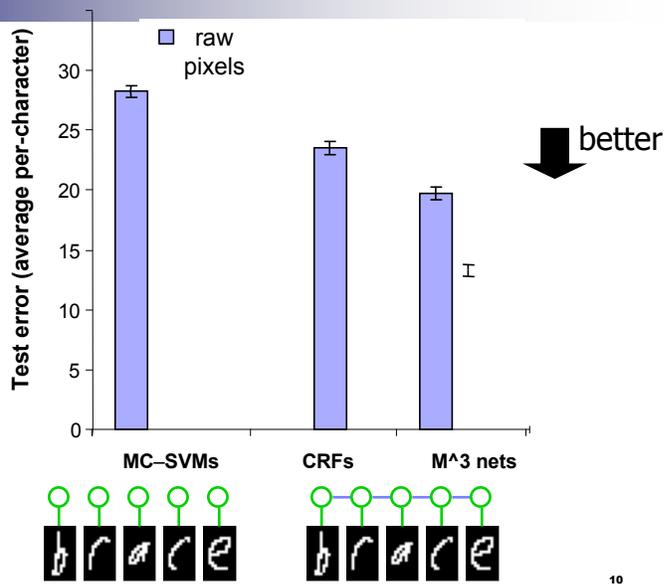
©2008 Carlos Guestrin ■ see store for details

9

# 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



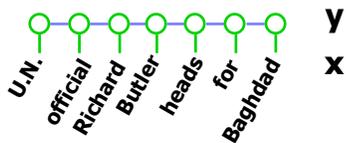
\*Crammer & Singer 01

©2008 Carlos Guestrin

10

# Named Entity Recognition

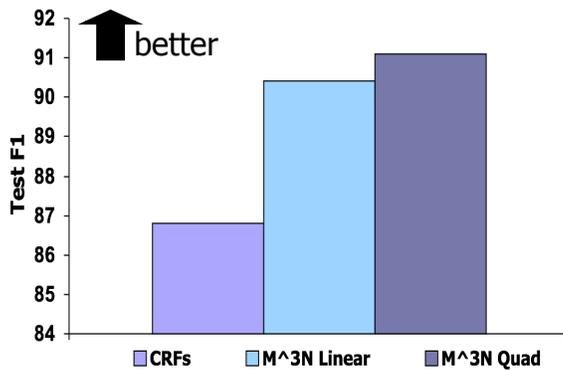
- Locate and classify named entities in sentences:
  - 4 categories: organization, person, location, misc.
  - e.g. "U.N. official Richard Butler heads for Baghdad".
- CoNLL 03 data set (200K words train, 50K words test)



$y_i = \text{org/per/loc/misc/none}$

$f(y_i, x) = [\dots,$   
 $I(y_i=\text{org}, x_i=\text{"U.N."}),$   
 $I(y_i=\text{per}, x_i=\text{capitalized}),$   
 $I(y_i=\text{loc}, x_i=\text{known city}),$

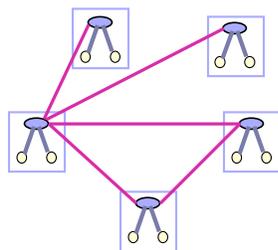
©2008 Carlos Guestrin



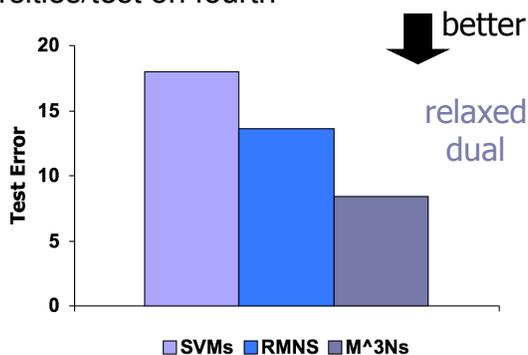
11

# 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



loopy belief propagation



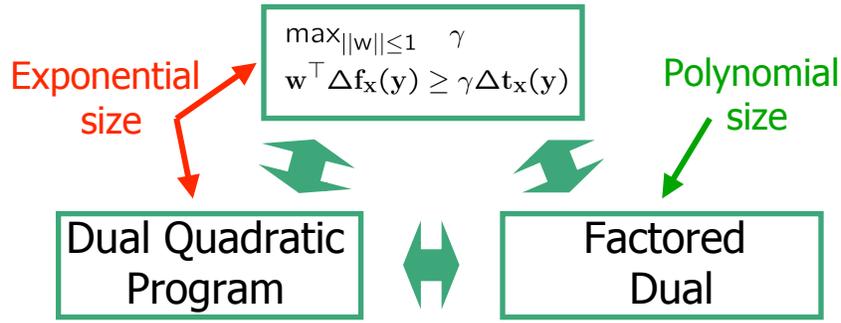
\*Taskar et al 02

©2009 Carlos Guestrin

12

# Solving M<sup>3</sup>Ns [Taskar, Guestrin, Koller '03]

## Estimation



©2008 Carlos Guestrin

13

# Other ways to solve M<sup>3</sup>Ns

- Sequential minimal optimization (SMO) [Taskar, Guestrin, Koller '03]
- Exponentiated gradient [Bartlett, Collins, Taskar, McAllester '04]
- Subgradient method [Ratliff, Bagnell, Zinkevich '07]
- ...
- Today
  - Simple constraint generation
    - (Other methods will perform better in many practical problems)
    - (Other methods are better suited to adding kernels)
    - (Other methods use similar principles to simpler constraint generation)

©2008 Carlos Guestrin

14

# Constraint generation overview

- Minimize  $w^2$ 
  - Subject to:  
 $w^T f(\text{brace}, \text{"brace"}) \geq w^T f(\text{brace}, \text{"aaaaa"}) + \Delta(\text{brace}, \text{aaaaa})$   
...
- General form:  
 $w^T f(\text{brace}, \text{"brace"}) \geq w^T f(\text{brace}, \text{"zzzzz"}) + \Delta(\text{brace}, \text{zzzzz})$
- Subset of constraints:
- Constraint generation:
  - 
  -

©2008 Carlos Guestrin

15

# Generating a constraint (simpler setting)

- Form of constraint  
 $w^T f(\text{brace}, \text{"brace"}) \geq w^T f(\text{brace}, \text{"aaaaa"}) + \Delta(\text{brace}, \text{aaaaa})$
- Another way of expressing:
- Given  $w$ , are any constraints violated?
- Separation oracle question:
- $\Delta(\text{brace}, \text{aaaaa})$  seems hard, simpler question:
- Exponentially many possibilities...

©2008 Carlos Guestrin

16

## Generating a constraint with hamming margin part ( $\Delta(\text{brace,aaaaa})$ )

$$P(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_i \phi(\mathbf{x}_i, y_i) \prod_i \phi(y_i, y_{i+1}) = \frac{1}{Z(\mathbf{x})} \exp\{\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y})\}$$

- Without  $\Delta(\text{brace,aaaaa}) \rightarrow$  standard (MAP or MPE) inference in graphical models
  - Solve with dynamic programming
  - For chains, it's called the Viterbi algorithm
- What do we do about  $\Delta(\text{brace,aaaaa})$  ?
- Reformulation:
- Same inference algorithm!!!
  - (slightly different potentials)

©2008 Carlos Guestrin

17

## Overview of constraint generation for $M^3Ns$

- Problem we want to solve:
- Maintain subset of “runner up labels” for each training example:
- Obtain some value for weights  $\mathbf{w}$
- Separation oracle:
  - Reformulate model to include hamming margin  $\Delta(\text{brace,aaaaa})$
  - Dynamic programming (inference in graphical models)
  - Apply to each data point

©2008 Carlos Guestrin

18

## Some reasons $M^3Ns$ are cool... :)

- Often perform better
- Can use kernels easily, and get sparsity
- Can be learned exactly in many problems where CRFs require approximate inference techniques
  - E.g., image segmentation (graph cuts)
- Can be generalized to other optimization problems
  - E.g., [Taskar, Chatalbashev, Koller, Guestrin '05]
  - Matching problems
  - Paths
  - Pretty much any optimization technique in the inner loop