

Maximum Margin Markov Networks and Constraint Generation continued

Optimization - 10725

Carlos Guestrin

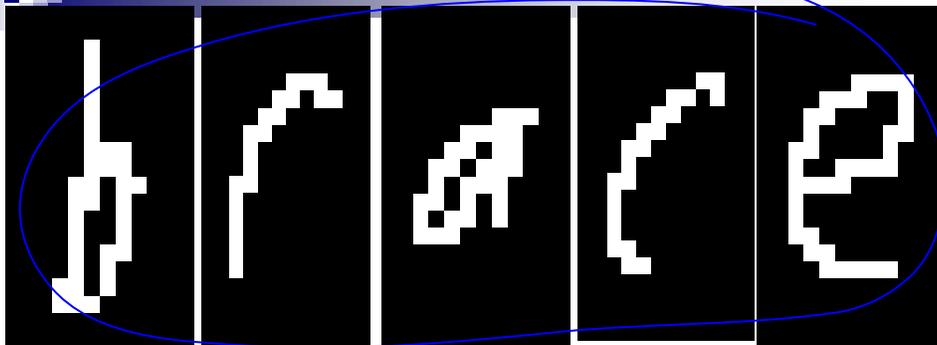
Carnegie Mellon University

February 18th, 2008

©2008 Carlos Guestrin

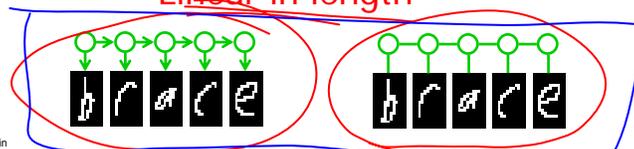
1

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

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

$Z(\mathbf{x})$: normal function (Partition function)
 $\phi(\mathbf{x}_i, y_i) = \exp\{\sum_{\alpha} w_{\alpha} f_{\alpha}(\mathbf{x}_i, y_i)\}$ (position, pixels)
 $\phi(y_i, y_{i+1}) = \exp\{\sum_{\beta} w_{\beta} f_{\beta}(y_i, y_{i+1})\}$ (neighboring letters)
 $f_{\beta}(y, y') = I(y='z', y'='a')$ (letter is z, following letter is a)
 $f_{\alpha}(\mathbf{x}, y) = I(x_p=1, y='z')$ (indexes pixel coordinate, pixel on which letter is z)

©2008 Carlos Guestrin *Lafferty et al. 01 3

CRF - short notation

$e^{w \cdot f} \cdot e^{w' \cdot f'} = e^{w \cdot f + w' \cdot f'}$

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

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

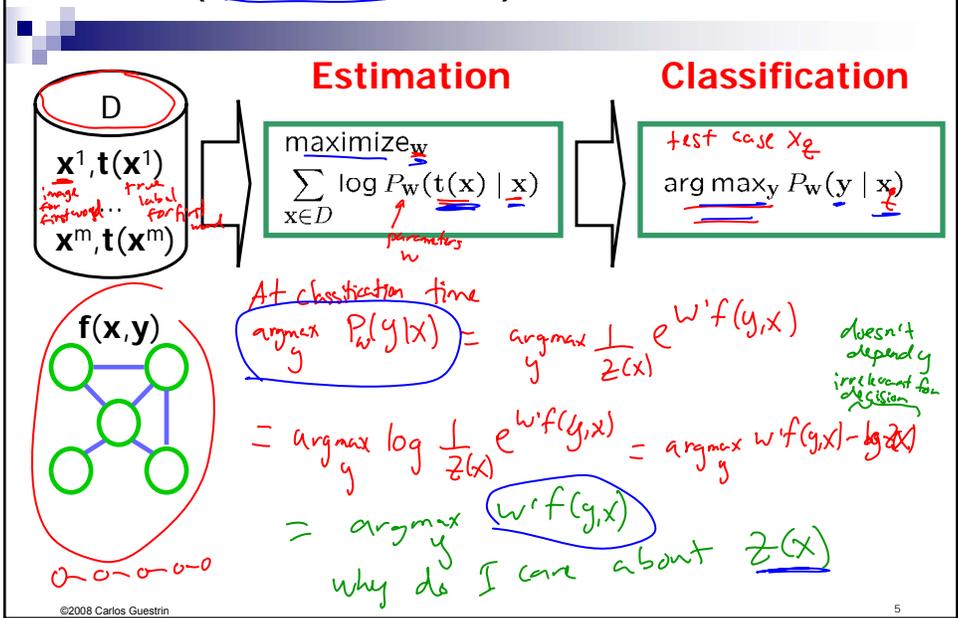
$\mathbf{w} = \begin{bmatrix} w_{\alpha_1} \\ \vdots \\ w_{\alpha_n} \\ w_{\beta_1} \\ \vdots \\ w_{\beta_m} \end{bmatrix}$

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

©2008 Carlos Guestrin *Lafferty et al. 01 4

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

Max (Conditional) Likelihood



OCR Example

- We want: given input
 $\operatorname{argmax}_{\text{word}} w^T f(\text{word}) = \text{"brace"}$
 - Equivalently: true but wins over all other possibilities
 $w^T f(\text{"brace"}, \text{"brace"}) > w^T f(\text{"brace"}, \text{"aaaaa"})$
 $w^T f(\text{"brace"}, \text{"brace"}) > w^T f(\text{"brace"}, \text{"aaaab"})$
 ...
 $w^T f(\text{"brace"}, \text{"brace"}) > w^T f(\text{"brace"}, \text{"zzzzz"})$
} number of constraints is exponential in length of word
- strictly greater than?
no w???
- ©2008 Carlos Guestrin 6

Max Margin Estimation

- Goal: find w such that

$$w^T f(x, t(x)) > w^T f(x, y) \quad \forall x \in D \quad \forall y \neq t(x)$$

(true label) (other labels)

$$w^T [f(x, t(x)) - f(x, y)] > 0$$

today:

"linearly separable"

can add slack vars just like SVMs

$$w^T \Delta f_x(y) > 0 \quad \forall x \in D, \forall y \neq t(x)$$

$$\max_{w, \gamma} \gamma$$

$$w^T \Delta f_x(y) \geq \gamma \leftarrow \text{"margin"}$$

"Solution 0"

$$\forall x \in D, \forall y \neq t(x)$$

$$\|w\|_2 \leq 1$$

©2008 Carlos Guestrin

7

Not all margins are equal

- Goal: find w such that

$$w^T \Delta f_x(y) \geq \gamma \Delta t_x(y) \quad \forall x \in D \quad \forall y \neq t(x)$$

→ if $y = t(x) \Rightarrow$ rhs: $\gamma \cdot \Delta t_x(y) = 0$

lhs: $w^T \Delta f_x(y) = w^T [f(x, t(x)) - f(x, y)] = 0$

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

$$\Delta t_{\text{brace}}(\text{"craze"}) = 2$$

$$\Delta t_{\text{brace}}(\text{"zzzzz"}) = 5$$

$$w^T \Delta f_{\text{brace}}(\text{"craze"}) \geq 2\gamma$$

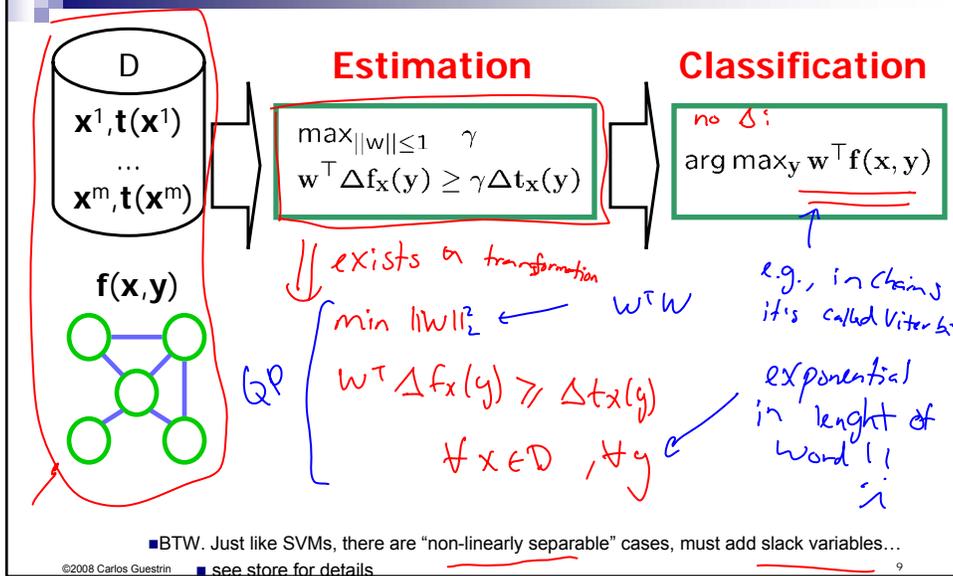
$$w^T \Delta f_{\text{brace}}(\text{"zzzzz"}) \geq 5\gamma$$

©2008 Carlos Guestrin

8

Maximum Margin Markov Nets

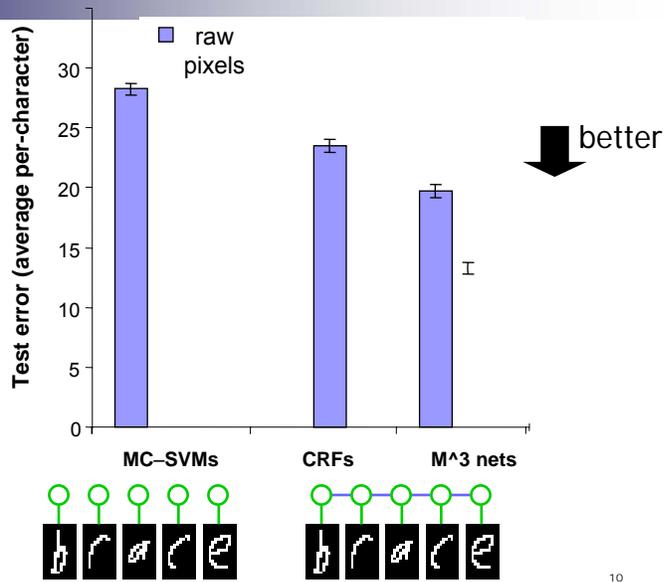
[Taskar, Guestrin, Koller '03]



Handwriting Recognition

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

Models:
 Multiclass-SVMs*
 CRFs
 M³ nets



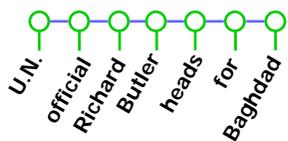
*Crammer & Singer 01

©2008 Carlos Guestrin

10

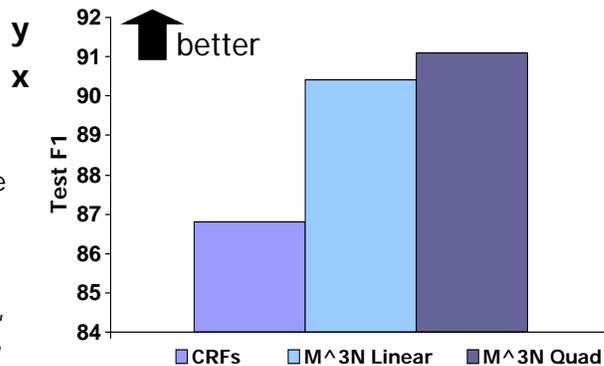
Named Entity Recognition

- 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)



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

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

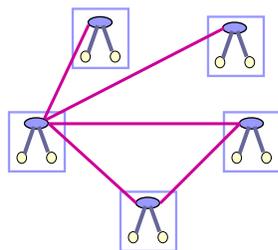


©2008 Carlos Guestrin

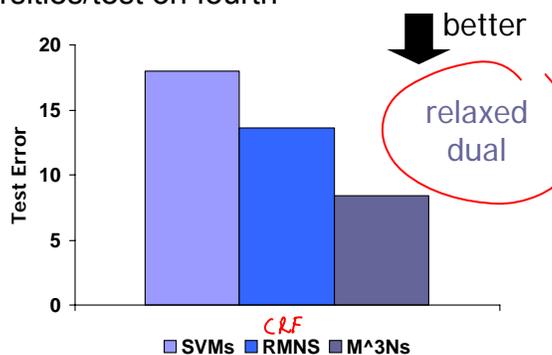
11

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

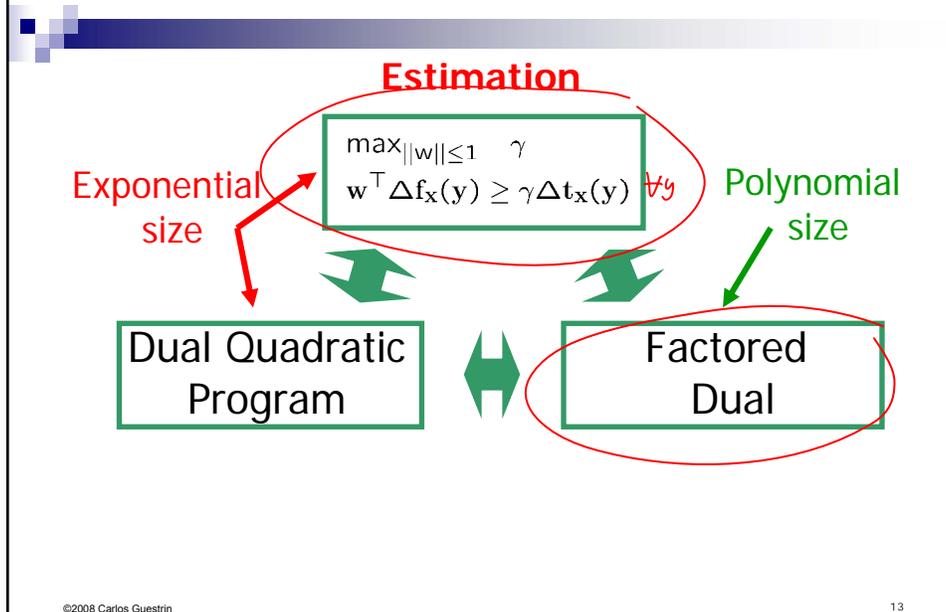


*Taskar et al 02

©2008 Carlos Guestrin

12

Solving M³Ns [Taskar, Guestrin, Koller '03]



Other ways to solve M³Ns

- 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)

Constraint generation overview

- Minimize $w^2 = w^T v$
 - Subject to:
 - $w^T f(\text{brace} \text{ "brace" }) \geq w^T f(\text{brace} \text{ "aaaaa" }) + \Delta(\text{brace, aaaaa})$ $\forall y$
 - ...
 - $w^T f(\text{brace} \text{ "brace" }) \geq w^T f(\text{brace} \text{ "zzzzz" }) + \Delta(\text{brace, zzzzz})$
- General form: $\min w^2$

$$w^T f(x, t(x)) \geq w^T f(x, y) + \Delta t_x(y) \quad \forall y \in \mathcal{Y}_x$$
- Subset of constraints: $\mathcal{R}_x \subset \mathcal{Y}_x$

$$\min w^2$$

$$w^T f(x, t(x)) \geq w^T f(x, y) + \Delta t_x(y) \quad \forall y \in \mathcal{R}_x$$
- Constraint generation:
 - Solve for w for a given set of \mathcal{R}_x
 - Find a violated constraint:
 - for each x :

$$w^T f(x, t(x)) \geq \max_y [w^T f(x, y) + \Delta t_x(y)]$$

separation oracle

©2008 Carlos Guestrin

15

Generating a constraint (simpler setting)

- Form of constraint

$$w^T f(\text{brace} \text{ "brace" }) \geq w^T f(\text{brace} \text{ "aaaaa" }) + \Delta(\text{brace, aaaaa}) \quad \forall y$$
- Another way of expressing:

$$w^T f(x, t(x)) \geq \max_y [w^T f(x, y) + \Delta t_x(y)]$$
- Given w , are any constraints violated?
- Separation oracle question:

$$\arg \max_y [w^T f(x, y) + \Delta t_x(y)]$$
- $\Delta(\text{brace, aaaaa})$ seems hard, simpler question:

$$\arg \max_y w^T f(x, y)$$
- Exponentially many possibilities...
 - same question as at classification (testing)
 - chains visits
 - time
 - more generally! inference for GM, e.g., variable elimination

©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 $\Delta(\text{brace,aaaaa})$? $\Delta(\text{brace,aaaaa}) = \neq$

$$\Delta t_x(y) = \sum_i \Delta t_{x_i}(y_i) \quad \Delta(\text{brace,aaaaa}) = \neq$$

$$= \Delta(b,a) + \Delta(v,a) + \Delta(a,a) + \Delta(c,a) + \Delta(a,c)$$

- Reformulation:

$$\prod_i \phi(x_i, y_i) \prod_i \phi(y_i, y_{i+1}) \cdot e^{\sum_i \Delta t_{x_i}(y_i)}$$

$$= \prod_i \underbrace{\phi(x_i, y_i)}_{\phi_i(x_i, y_i)} \cdot e^{\Delta t_{x_i}(y_i)} \cdot \prod_i \phi(y_i, y_{i+1}) = e$$

wanted to argmax $w^T f(x, y) + \Delta t_x(y)$

- Same inference algorithm!!!
 - (slightly different potentials)

©2008 Carlos Guestrin

17

Overview of constraint generation for M³Ns

- Problem we want to solve:

$$\min \|w\| = w^T w$$

$$w^T f(x, t(x)) \geq w^T f(x, y) + \Delta t_x(y)$$

$\forall x, \forall y \in \mathcal{Y}_x$
- Maintain subset of "runner up labels" for each training example:

$$\mathcal{L}_x \subset \mathcal{Y}_x$$
- Obtain some value for weights 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 M³Ns are cool... :)

- Often perform better
- Can use kernels easily, and get sparsity *like SVMs*
- 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