

Machine Learning

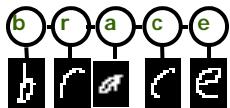
10-701/15-781, Fall 2011

Max-Margin Learning of Graphical Models

Eric Xing



Lecture 21, November 28, 2011



© Eric Xing @ CMU, 2006-2011

1

Advanced topics in Max-Margin Learning (cont.)



$$\max_{\alpha} \mathcal{J}(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\underbrace{\mathbf{x}_i^T \mathbf{x}_j}_{K(x_i, x_j)})$$

$$\underbrace{\mathbf{w}^T \mathbf{x}_{\text{new}} + b \leq 0}_{y \in \{-1, 1\}} \quad \underbrace{(w^T x + b)y \geq 0}_{f(y, x; w) \geq 0}$$

- Implicit nonlinear data transformation
 - The Kernel trick
- Point rule or average rule?
 - Maximum entropy discrimination
- Can we predict $\text{vec}(y)$?
 - Structured SVM, aka, Maximum Margin Markov Networks
- Putting everything all together !

$$w^* \rightarrow p(w) \quad f(\cdot, w) \rightarrow \int f(\cdot, w) p(w) d\omega$$

$y \in \{-1, 1\}$
 $\rightarrow y \in \{0, 1\}^K$
 $\downarrow \quad \downarrow \quad \downarrow$
 $x \quad y_1 \quad y_2 \quad \dots$

© Eric Xing @ CMU, 2006-2011

2

(2) Model averaging

- Inputs x , class $y = +1, -1$
- data $D = \{(x_1, y_1), \dots, (x_m, y_m)\}$

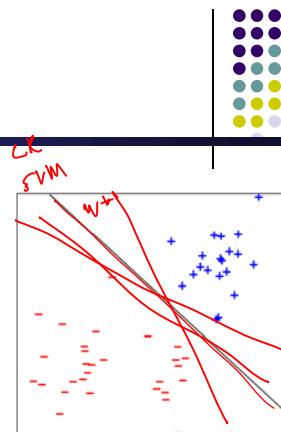
- Point Rule:

- learn $f^{\text{opt}}(x)$ discriminant function from $F = \{f\}$ family of discriminants
- classify $y = \text{sign } f^{\text{opt}}(x)$

- E.g., SVM

$$y = \underset{\text{argmax}}{f^{\text{opt}}(x)} = w^T x_{\text{new}} + b \rightarrow y = \sum_{w \in F} f^{\text{opt}}(w) g(w) \underset{y_{\text{new}}}{\downarrow}$$

© Eric Xing @ CMU, 2006-2011



3

Model averaging

- There exist many f with near optimal performance

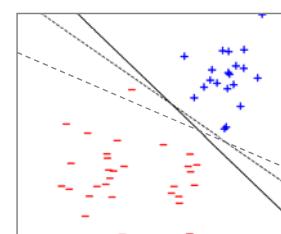
- Instead of choosing f^{opt} , average over all f in F

$Q(f) = \text{weight of } f$

$$\begin{aligned} y(x) &= \text{sign} \int_F Q(f) f(x) df \\ &= \text{sign} \langle f(x) \rangle_Q \end{aligned}$$

Expectation
rule

- How to specify:
 $F = \{f\}$ family of discriminant functions?
- How to learn $Q(f)$ distribution over F ?



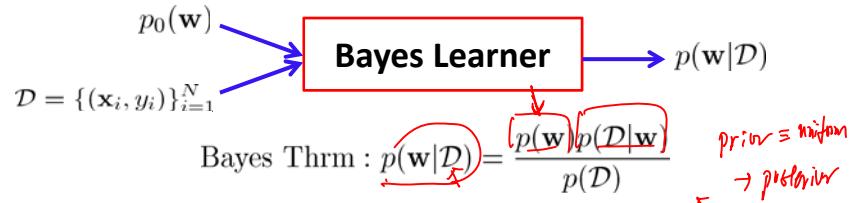
© Eric Xing @ CMU, 2006-2011

4

Recall Bayesian Inference

ML ↗

- Bayesian learning:



- Bayes Predictor (model averaging):

$$h_1(\mathbf{x}; p(\mathbf{w})) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \int p(\mathbf{w}) f(\mathbf{x}, \mathbf{y}; \mathbf{w}) d\mathbf{w}$$

Recall in SVM: $h_0(\mathbf{x}; \mathbf{w}) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} F(\mathbf{x}, \mathbf{y}; \mathbf{w})$

- What p_0 ?

© Eric Xing @ CMU, 2006-2011

5

How to score distributions?

- Entropy

- Entropy $H(X)$ of a random variable X

$$H(X) = - \sum_{i=1}^N P(x = i) \log_2 P(x = i)$$

- $H(X)$ is the expected number of bits needed to encode a randomly drawn value of X (under most efficient code)

- Why?

Information theory:

Most efficient code assigns $-\log_2 P(X=i)$ bits to encode the message $X=i$,
So, expected number of bits to code one random X is:

$$-\sum_{i=1}^N P(x = i) \log_2 P(x = i)$$

© Eric Xing @ CMU, 2006-2011

6

Maximum Entropy Discrimination

- Given data set $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ find

$$Q_{\text{ME}} = \arg \max Q(Q) - C \sum_i \xi_i$$

s.t.

$$y^i \langle f(\mathbf{x}^i) \rangle_{Q_{\text{ME}}} \geq \xi_i, \quad \forall i$$

$$\xi_i \geq 0 \quad \forall i$$

$y^i f(\mathbf{x}^i, \mathbf{w}^*) \geq \xi_i$

- solution Q_{ME} correctly classifies \mathcal{D}
- among all admissible Q , Q_{ME} has max entropy
- max entropy \rightarrow "minimum assumption" about f

© Eric Xing @ CMU, 2006-2011

7

Introducing Priors

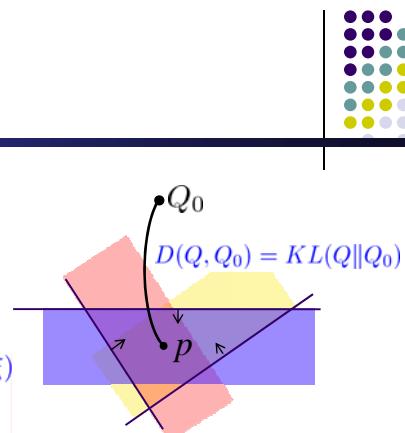
- Prior $Q_0(f)$
- Minimum Relative Entropy Discrimination

$$Q_{\text{MRE}} = \arg \min Q(Q) - C \sum_i \xi_i$$

s.t.

$$y^i \langle f(\mathbf{x}^i) \rangle_{Q_{\text{MRE}}} \geq \xi_i, \quad \forall i$$

$$\xi_i \geq 0 \quad \forall i$$



- Convex problem: Q_{MRE} unique solution
- MER \rightarrow "minimum additional assumption" over Q_0 about f

© Eric Xing @ CMU, 2006-2011

8

Solution: Q_{ME} as a projection

- Convex problem: Q_{ME} unique

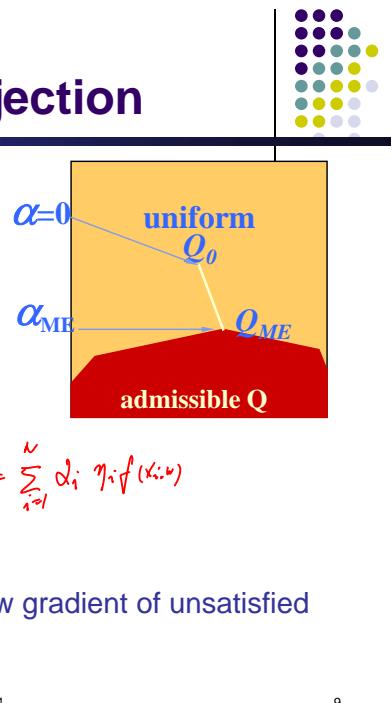
- Theorem:

Ans: $p(w|D) = p(D(w)p(w))$

$$Q_{MRE} \propto \exp\left\{\sum_{i=1}^N \alpha_i y_i f(x_i; w)\right\} Q_0(w)$$

$$= \sum_{i \in SV} \alpha_i y_i f(x_i; w)$$

$\alpha_i \geq 0$ Lagrange multipliers



- finding Q_M : start with $\alpha_i = 0$ and follow gradient of unsatisfied constraints

© Eric Xing @ CMU, 2006-2011

9

Solution to MED

- Theorem (Solution to MED):

- Posterior Distribution:

$$Q(\mathbf{w}) = \frac{1}{Z(\alpha)} Q_0(\mathbf{w}) \exp\left\{\sum_i \alpha_i y_i [f(\mathbf{x}_i; \mathbf{w})]\right\}$$

- Dual Optimization Problem:

$$\begin{aligned} D1 : \quad & \max_{\alpha} -\log Z(\alpha) - U^*(\alpha) \\ \text{s.t. } & \alpha_i(y) \geq 0, \forall i, \end{aligned}$$

$U^*(\cdot)$ is the conjugate of the $U(\cdot)$, i.e., $U^*(\alpha) = \sup_{\xi} (\sum_{i,y} \alpha_i(y) \xi_i - U(\xi))$

- Algorithm: to computer α_t , $t = 1, \dots, T$

- start with $\alpha_t = 0$ (uniform distribution)
- iterative ascent on $J(\alpha)$ until convergence

© Eric Xing @ CMU, 2006-2011

10



Examples: SVMs

- Theorem

For $f(x) = w^T x + b$, $Q_0(w) = \text{Normal}(0, I)$, $Q_0(b) = \text{non-informative prior}$, the Lagrange multipliers α are obtained by maximizing $J(\alpha)$ subject to $0 \leq \alpha_t \leq C$ and $\sum_t \alpha_t y_t = 0$, where

$$J(\alpha) = \sum_t [\alpha_t + \log(1 - \alpha_t/C)] - \frac{1}{2} \sum_{s,t} \alpha_s \alpha_t y_s y_t x_s^T x_t$$

- Separable $D \rightarrow$ SVM recovered exactly
- Inseparable $D \rightarrow$ SVM recovered with different misclassification penalty

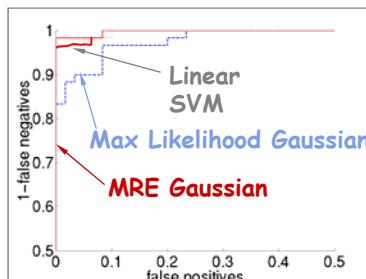
© Eric Xing @ CMU, 2006-2011

11



SVM extensions

- Example: Leptograpsus Crabs (5 inputs, $T_{\text{train}}=80$, $T_{\text{test}}=120$)



SVM
 → dist SVMs
 by MLE

© Eric Xing @ CMU, 2006-2011

12

(3) Structured Prediction

- Unstructured prediction



$$\mathbf{x} = \begin{pmatrix} x_{11} & x_{12} & \dots \\ x_{21} & x_{22} & \dots \\ \vdots & \vdots & \dots \end{pmatrix}$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \end{pmatrix}$$

$p(\mathbf{y} | \mathbf{x}) \rightarrow$
LR SVM

- Structured prediction

- Part of speech tagging

$\mathbf{x} = \text{"Do you want sugar in it?"} \Rightarrow \mathbf{y} = \langle \text{verb pron verb noun prep pron} \rangle$

- Image segmentation



$$\mathbf{x} = \begin{pmatrix} x_{11} & x_{12} & \dots \\ x_{21} & x_{22} & \dots \\ \vdots & \vdots & \dots \end{pmatrix}$$

$$\mathbf{y} = \begin{pmatrix} y_{11} & y_{12} & \dots \\ y_{21} & y_{22} & \dots \\ \vdots & \vdots & \dots \end{pmatrix}$$

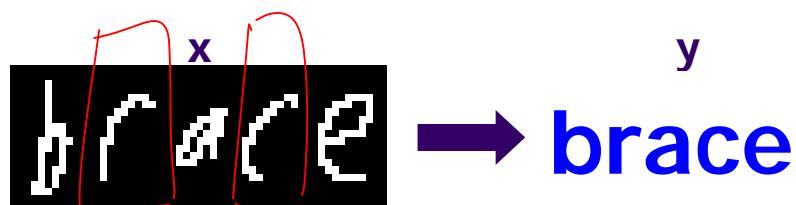
$p(\vec{\mathbf{y}} | \vec{\mathbf{x}}) \rightarrow ?$

CRF

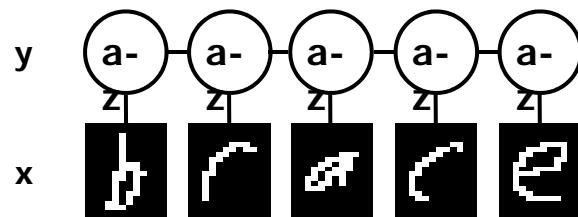
© Eric Xing @ CMU, 2006-2011

13

OCR example



Sequential structure



© Eric Xing @ CMU, 2006-2011

14

Classical Predictive Models

- Input and output space: $\mathcal{X} \triangleq \mathbb{R}^{M_x}$ $\mathcal{Y} \triangleq \{-1, +1\}$
- Predictive function $h(\mathbf{x}) : y^* = h(\mathbf{x}) \triangleq \arg \max_{y \in \mathcal{Y}} F(\mathbf{x}, y; \mathbf{w})$
- Examples: $F(\mathbf{x}, y; \mathbf{w}) = g(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, y))$
- Learning: $\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathcal{W}} \ell(\mathbf{x}, y; \mathbf{w}) + \lambda R(\mathbf{w})$
where $\ell(\cdot)$ represents a convex loss, and $R(\mathbf{w})$ is a regularizer preventing overfitting

Logistic Regression

- Max-likelihood (or MAP) estimation

$$\begin{aligned} & \ln p(y_i | \mathbf{x}_i) \\ &= \ln \left(\frac{1}{1+e^{-(\mathbf{w} \cdot \mathbf{x}_i + b)}} \right)^{y_i} \left(1 - \frac{1}{1+e^{-(\mathbf{w} \cdot \mathbf{x}_i + b)}} \right)^{1-y_i} \\ &= \ln(1+e^{-(\mathbf{w} \cdot \mathbf{x}_i + b)})^{-y_i} \end{aligned}$$

Support Vector Machines (SVM)

- Max-margin learning

$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{i=1}^m \xi_i \\ \text{s.t. } & \xi_i \geq 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b), \quad \forall i \\ & \xi_i \geq 0, \quad \forall i \end{aligned}$$

© Eric Xing @ CMU, 2006-2011

15

Classical Predictive Models

- Input and output space: $\mathcal{X} \triangleq \mathbb{R}^{M_x}$ $\mathcal{Y} \triangleq \{-1, +1\}$
- Predictive function $h(\mathbf{x}) : y^* = h(\mathbf{x}) \triangleq \arg \max_{y \in \mathcal{Y}} F(\mathbf{x}, y; \mathbf{w})$
- Examples: $F(\mathbf{x}, y; \mathbf{w}) = g(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, y))$
- Learning: $\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathcal{W}} \ell(\mathbf{x}, y; \mathbf{w}) + \lambda R(\mathbf{w})$

where $\ell(\cdot)$ represents a convex loss, and $R(\mathbf{w})$ is a regularizer preventing overfitting

Logistic Regression

- Max-likelihood (or MAP) estimation

$$\max_{\mathbf{w}} \mathcal{L}(\mathcal{D}; \mathbf{w}) \triangleq \sum_{i=1}^N \log p(y^i | \mathbf{x}^i; \mathbf{w}) + \mathcal{N}(\mathbf{w})$$

- Corresponds to a Log loss with L2 R

$$\ell_{LL}(\mathbf{x}, y; \mathbf{w}) \triangleq \ln \sum_{y' \in \mathcal{Y}} \exp\{\mathbf{w}^\top \mathbf{f}(\mathbf{x}, y')\} - \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y)$$

Support Vector Machines (SVM)

- Max-margin learning

$$\begin{aligned} & \min_{\mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{i=1}^N \xi_i \\ \text{s.t. } & \forall i, \forall y' \neq y^i : \mathbf{w}^\top \Delta \mathbf{f}_i(y') \geq 1 - \xi_i, \quad \xi_i \geq 0. \\ & \quad \text{Corresponds to a hinge loss with L2 R} \\ & \ell_{MM}(\mathbf{x}, y; \mathbf{w}) \triangleq \max_{y' \in \mathcal{Y}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y') - \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y) + \ell'(y', y) \end{aligned}$$

© Eric Xing @ CMU, 2006-2011

16



Classical Predictive Models

- Inputs:
 - a set of training samples $\mathcal{D} = \{(x^i, y^i)\}_{i=1}^N$
where $x^i = [x_1^i, x_2^i, \dots, x_d^i]^\top$ and $y^i \in C \triangleq \{c_1, c_2, \dots, c_L\}$
- Outputs:
 - a predictive function $h(x) \quad y^* = h(x) \triangleq \arg \max_y F(x, y; w)$

Logistic Regression

- Max-likelihood (or MAP) estimation

$$\max_w \mathcal{L}(\mathcal{D}; w) \triangleq \sum_{i=1}^N \log p(y^i | x^i; w) + \mathcal{N}(w)$$

- Corresponds to a Log loss with L2 R

$$\ell_{LL}(x, y; w) \triangleq \ln \sum_{y' \in \mathcal{Y}} \exp\{w^\top f(x, y')\} - w^\top f(x, y)$$

Support Vector Machines (SVM)

- Max-margin learning

$$\min_{w, \xi} \frac{1}{2} w^\top w + C \sum_{i=1}^N \xi_i;$$

s.t. $\forall i, \forall y' \neq y^i : w^\top \Delta f_i(y') \geq 1 - \xi_i, \xi_i \geq 0.$

- Corresponds to a hinge loss with L2 R

$$\ell_{MM}(x, y; w) \triangleq \max_{y' \in \mathcal{Y}} w^\top f(x, y') - w^\top f(x, y) + \ell'(y', y)$$

Advantages:

- Full probabilistic semantics
- Straightforward Bayesian or direct regularization
- Hidden structures or generative hierarchy

Advantages:

- Dual sparsity: few support vectors
- Kernel tricks
- Strong empirical results

© Eric Xing @ CMU, 2006-2011

17



Structured Prediction Graphical Models

- Input and output space $\mathcal{X} \triangleq \mathbb{R}_{X_1} \times \dots, \mathbb{R}_{X_K}$ $\mathcal{Y} \triangleq \mathbb{R}_{Y_1} \times \dots, \mathbb{R}_{Y_{K'}}$

• Conditional Random Fields (CRFs) (Lafferty et al 2001)

- Based on a Logistic Loss (LR)
- Max-likelihood estimation (point-estimate)

$$\mathcal{L}(\mathcal{D}; w) \triangleq \log \sum_{y'} \exp(w^\top f(x, y'))$$

$w^\top f(x, y')$

$-w^\top f(x, y)$

• Max-margin Markov Networks (M³Ns) (Taskar et al 2003)

- Based on a Hinge Loss (SVM)
- Max-margin learning (point-estimate)

$$\mathcal{L}(\mathcal{D}; w) \triangleq \log \max_y (w^\top f(x, y))$$

$w^\top f(x, y)$

$+R(w)$

$\mathcal{L}(\mathcal{M}) \triangleq \log \max_y (w^\top f(x, y))$

$w^\top f(x, y)$

© Eric Xing @ CMU, 2006-2011

18

Structured Models

$$h(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}(\mathbf{x})} F(\mathbf{x}, y)$$

↑
space of feasible outputs ↑
discriminant function

- Assumptions:

$$F(\mathbf{x}, \mathbf{y}) = \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) = \sum_p w_p f(x_p, y_p)$$

$\mathbf{f} = \begin{bmatrix} f(y_1 | x_1) \\ f(y_2 | x_2) \\ \vdots \\ f(y_n | x_n) \end{bmatrix}$

- Linear combination of features
- Sum of partial scores: index p represents a part in the structure
- Random fields or Markov network features:

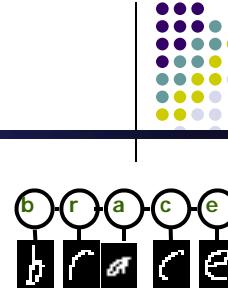
- Markov properties are encoded in the feature functions $f(\mathbf{x}, \mathbf{y})$

© Eric Xing @ CMU, 2006-2011

19

Learning w

- Training examples $(\mathbf{x}_i, \mathbf{y}_i)$
 - Probabilistic approach:
- $$P_w(\mathbf{y} | \mathbf{x}) = \frac{1}{Z_w(\mathbf{x})} \exp\{\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})\}$$
- Computing $Z_w(\mathbf{x})$ can be NP-complete
 - Tractable models but intractable estimation
 - Large margin approach:
 - Exact and efficient when prediction is tractable



© Eric Xing @ CMU, 2006-2011

20

E.g. Max-Margin Markov Networks

- Convex Optimization Problem:

$$P_0 (M^3N) : \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

s.t. $\forall i, \forall \mathbf{y} \neq \mathbf{y}_i :$

$$\mathbf{w}^\top \Delta \mathbf{f}_i(\mathbf{x}, \mathbf{y}) \geq \Delta \ell_i(\mathbf{y}) - \xi_i, \quad \xi_i \geq 0,$$

$$= f_i(\mathbf{x}, \mathbf{y}) - f_i(\mathbf{x}, \mathbf{y}_i)$$

- Feasible subspace of weights:

$$\mathcal{F}_0 = \{\mathbf{w} : \mathbf{w}^\top \Delta \mathbf{f}_i(\mathbf{x}, \mathbf{y}) \geq \Delta \ell_i(\mathbf{y}) - \xi_i; \forall i, \forall \mathbf{y} \neq \mathbf{y}_i\}$$

- Predictive Function:

$$h_0(\mathbf{x}; \mathbf{w}) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} F(\mathbf{x}, \mathbf{y}; \mathbf{w})$$

© Eric Xing @ CMU, 2006-2011

21

OCR Example

- We want:

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

- Equivalently:

$$\begin{aligned} \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"brace"}) &> \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"aaaaa"}) \\ \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"brace"}) &> \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"aaaab"}) \\ \dots \\ \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"brace"}) &> \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"zzzzz"}) \end{aligned}$$

*W must satisfy each of
the below constraints*

a lot!

© Eric Xing @ CMU, 2006-2011

22



Large Margin Estimation

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

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

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

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

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

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

*Taskar et al. 03

© Eric Xing @ CMU, 2006-2011

22



Large Margin Estimation

- Recall from SVMs:

- Maximizing margin γ 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. $\mathbf{w}^\top \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \geq \mathbf{w}^\top (\mathbf{f}(\mathbf{x}_i, \mathbf{y}'_i) + \ell(\mathbf{y}_i, \mathbf{y}'_i)), \quad \forall i, \mathbf{y}'_i \in \mathcal{Y}_i$

© Eric Xing @ CMU, 2006-2011

24



Min-max Formulation

- Brute force enumeration of constraints:

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

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

- The constraints are exponential in the size of the structure

- Alternative: min-max formulation

- add only the most violated constraint

$$\mathbf{y}' = \arg \max_{\mathbf{y} \neq \mathbf{y}^*} [\mathbf{w}^\top \mathbf{f}(\mathbf{x}_i, \mathbf{y}) + \ell(\mathbf{y}_i, \mathbf{y})]$$

$$\text{add to QP : } \mathbf{w}^\top \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \geq \mathbf{w}^\top \mathbf{f}(\mathbf{x}_i, \mathbf{y}') + \ell(\mathbf{y}_i, \mathbf{y}')$$

- Handles more general loss functions
- Only polynomial # of constraints needed
- Several algorithms exist ...

© Eric Xing @ CMU, 2006-2011

25



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

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

discrete optim.

continuous optim.

- How to do this conversion?

- Linear chain example in the next slides →

© Eric Xing @ CMU, 2006-2011

26

$y \Rightarrow z$ map for linear chain structures



OCR example: $y = 'ABABB'$;

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 1 . 0	0 0 . 0	0 1 . 0	0 0 . 0
B	0 0 . 0	1 0 . 0	0 0 . 0	0 1 . 0
:	. . . 0	. . . 0	. . . 0	. . . 0
B	0 0 0 0	0 0 0 0	0 0 0 0	0 0 0 0

A B . B	A B . B	A B . B	A B . B
---------	---------	---------	---------

© Eric Xing @ CMU, 2006-2011

27

$y \Rightarrow z$ map for linear chain structures



Rewriting the maximization function in terms of indicator variables:

$$\max_z \sum_{j,m} z_j(m) [w^\top f_{\text{node}}(x_j, m) + \ell_j(m)] + \sum_{jk,m,n} z_{jk}(m, n) [w^\top f_{\text{edge}}(x_{jk}, m, n) + \ell_{jk}(m, n)]$$

$\left. \begin{matrix} \\ \end{matrix} \right\} (\mathbf{F}^\top \mathbf{w} + \ell)^\top \mathbf{z}$

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

$\left. \begin{matrix} z_j(m) \\ z_{jk}(m, n) \end{matrix} \right\} \text{normalization } \sum_m z_j(m) = 1$

$\left. \begin{matrix} z_k(n) \\ z_{jk}(m, n) \end{matrix} \right\} \text{agreement } \sum_n z_{jk}(m, n) = z_j(m)$

$\left. \begin{matrix} \max_{Az=b} (\mathbf{F}^\top \mathbf{w} + \ell)^\top \mathbf{z} \\ Az=b \end{matrix} \right\}$

© Eric Xing @ CMU, 2006-2011

28

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:

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

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

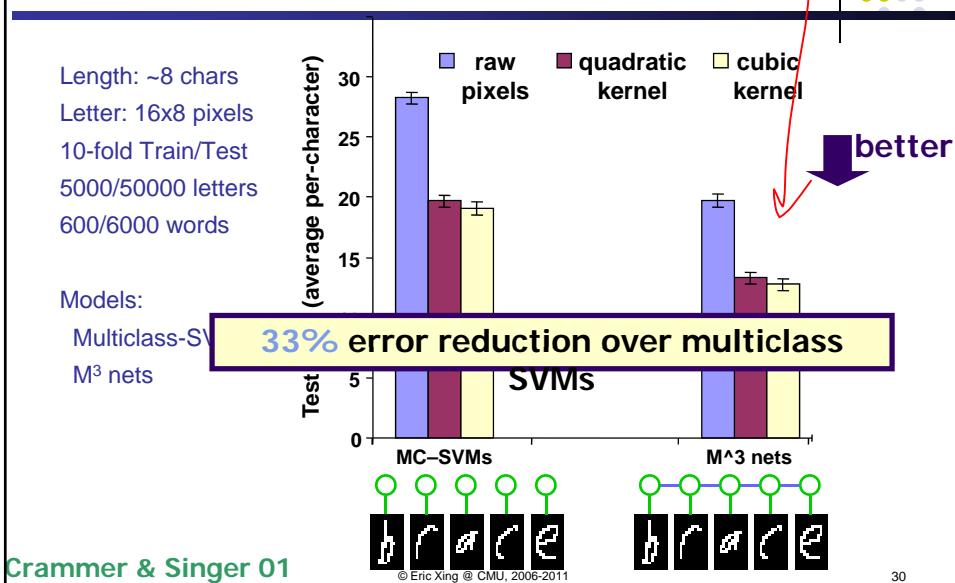
$$\begin{aligned} \mathbf{z} &\geq 0 \\ \sum z_i &= 1 \\ \sum z_i &= 1 \end{aligned}$$

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

© Eric Xing @ CMU, 2006-2011

29

Results: Handwriting Recognition



MLE versus max-margin learning



- **Likelihood-based estimation**

- Probabilistic (joint/conditional likelihood model)
- Easy to perform Bayesian learning, and incorporate prior knowledge, latent structures, missing data
- Bayesian or direct regularization
- Hidden structures or generative hierarchy

- **Max-margin learning**

- Non-probabilistic (concentrate on input-output mapping)
- Not obvious how to perform Bayesian learning or consider prior, and missing data
- Support vector property, sound theoretical guarantee with limited samples
- Kernel tricks

- **Maximum Entropy Discrimination (MED) (Jaakkola, et al., 1999)**

- Model averaging $\hat{y} = \text{sign} \int p(\mathbf{w}) F(x; \mathbf{w}) d\mathbf{w} \quad (y \in \{+1, -1\})$
- The optimization problem (binary classification)

$$\min_{p(\Theta)} KL(p(\Theta) || p_0(\Theta))$$

$$\text{s.t. } \int p(\Theta) [y_i F(x; \mathbf{w}) - \xi_i] d\Theta \geq 0, \forall i,$$

where Θ is the parameter \mathbf{w} when ξ are kept fixed or the pair (\mathbf{w}, ξ) when we want to optimize over ξ

© Eric Xing @ CMU, 2006-2011

31

Maximum Entropy Discrimination Markov Networks



- Structured MaxEnt Discrimination (SMED):

$$P1 : \min_{p(\mathbf{w}), \xi} KL(p(\mathbf{w}) || p_0(\mathbf{w})) + U(\xi)$$

$$\text{s.t. } p(\mathbf{w}) \in \mathcal{F}_1, \xi_i \geq 0, \forall i.$$

generalized maximum entropy or *regularized* KL-divergence

- Feasible subspace of weight distribution:

$$\mathcal{F}_1 = \left\{ p(\mathbf{w}) : \int p(\mathbf{w}) [\Delta F_i(\mathbf{y}; \mathbf{w}) - \Delta \ell_i(\mathbf{y})] d\mathbf{w} \geq -\xi_i, \forall i, \forall \mathbf{y} \neq \mathbf{y}^i \right\},$$

expected margin constraints.

- Average from distribution of M³Ns

$$h_1(\mathbf{x}; p(\mathbf{w})) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \int p(\mathbf{w}) F(\mathbf{x}, \mathbf{y}; \mathbf{w}) d\mathbf{w}$$

© Eric Xing @ CMU, 2006-2011

32



Solution to MaxEnDNet

- Theorem:

- Posterior Distribution:

$$p(\mathbf{w}) = \frac{1}{Z(\alpha)} p_0(\mathbf{w}) \exp \left\{ \sum_{i,y} \alpha_i(y) [\Delta F_i(y; \mathbf{w}) - \Delta \ell_i(y)] \right\}$$

- Dual Optimization Problem:

$$\begin{aligned} D1 : \quad & \max_{\alpha} -\log Z(\alpha) - U^*(\alpha) \\ \text{s.t.} \quad & \alpha_i(y) \geq 0, \forall i, \forall y, \end{aligned}$$

$U^*(\cdot)$ is the conjugate of the $U(\cdot)$, i.e., $U^*(\alpha) = \sup_{\xi} (\sum_{i,y} \alpha_i(y) \xi_i - U(\xi))$

© Eric Xing @ CMU, 2006-2011

33



Gaussian MaxEnDNet (reduction to M³N)



- Theorem

- Assume

$$F(\mathbf{x}, \mathbf{y}; \mathbf{w}) = \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}), U(\xi) = C \sum_i \xi_i, \text{ and } p_0(\mathbf{w}) = \mathcal{N}(\mathbf{w} | 0, I)$$

- Posterior distribution: $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mu_{\mathbf{w}}, I)$, where $\mu_{\mathbf{w}} = \sum_{i,y} \alpha_i(y) \Delta \mathbf{f}_i(\mathbf{y})$
- Dual optimization:
$$\begin{aligned} \max_{\alpha} \quad & \sum_{i,y} \alpha_i(y) \Delta \ell_i(y) - \frac{1}{2} \left\| \sum_{i,y} \alpha_i(y) \Delta \mathbf{f}_i(\mathbf{y}) \right\|^2 \\ \text{s.t.} \quad & \sum_y \alpha_i(y) = C; \alpha_i(y) \geq 0, \forall i, \forall y, \end{aligned}$$
- Predictive rule:
$$h_1(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \int p(\mathbf{w}) F(\mathbf{x}, \mathbf{y}; \mathbf{w}) d\mathbf{w} = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \mu_{\mathbf{w}}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$$

- Thus, MaxEnDNet subsumes M³Ns and admits all the merits of max-margin learning
- Furthermore, MaxEnDNet has at least three advantages ...

© Eric Xing @ CMU, 2006-2011

34

Three Advantages

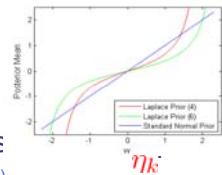
- An averaging Model: PAC-Bayesian prediction error guarantee

$$\Pr_Q(M(h, \mathbf{x}, \mathbf{y}) \leq 0) \leq \Pr_D(M(h, \mathbf{x}, \mathbf{y}) \leq \gamma) + O\left(\frac{\sqrt{\gamma^{-2}KL(p||p_0)\ln(N|\mathcal{Y}|)} + \ln N + \ln \delta^{-1}}{N}\right).$$

- Error

- Standard Normal prior => reduction to standard M³N (we've seen it)

- Laplace prior => Posterior shrinkage effects (sparse M³N)



- Integration, Concavity and Dis

- Incorporate latent variables and structures (PoMEN)
- Semisupervised learning (with partially labeled data)

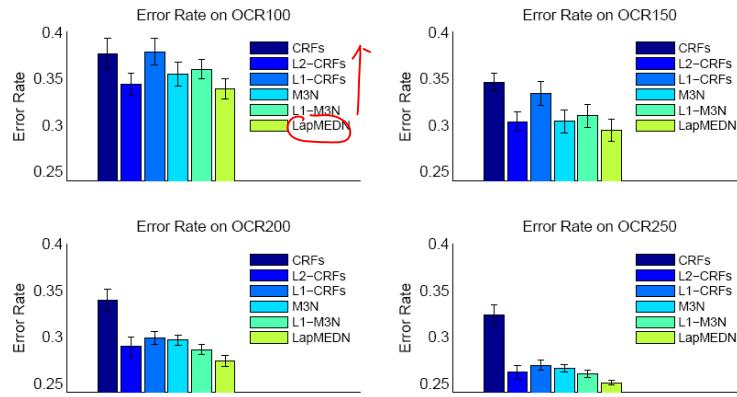
© Eric Xing @ CMU, 2006-2011

35

Experimental results on OCR datasets

(CRFs, L₁ – CRFs, L₂ – CRFs, M³Ns, L₁ – M³Ns, and LapMEDN)

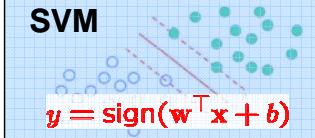
- We randomly construct OCR100, OCR150, OCR200, and OCR250 for 10 fold CV.



© Eric Xing @ CMU, 2006-2011

36

Discriminative Learning Paradigms



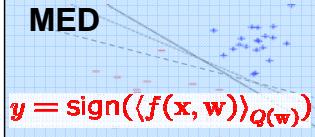
$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

$$y^i (\mathbf{w}^\top \mathbf{x}^i + b) \geq 1 - \xi_i, \quad \forall i$$



$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

$$\mathbf{w}^\top [\mathbf{f}(\mathbf{x}^i)] - \mathbf{f}(\mathbf{x}^i, \mathbf{y}) \geq \ell(\mathbf{y}^i, \mathbf{y}) - \xi_i, \quad \forall i, \forall \mathbf{y} \neq \mathbf{y}^i$$



$$\min_Q \text{KL}(Q || Q_0)$$

$$y^i \langle f(\mathbf{x}^i) \rangle_Q \geq \xi_i, \quad \forall i$$



See [Zhu and Xing, 2008]

© Eric Xing @ CMU, 2006-2011

37

Summary

- Maximum margin nonlinear separator
 - Kernel trick
 - Project into linearly separable space (possibly high or infinite dimensional)
 - No need to know the explicit projection function
- Max-entropy discrimination
 - Average rule for prediction,
 - Average taken over a posterior distribution of \mathbf{w} who defines the separation hyperplane
 - $P(\mathbf{w})$ is obtained by max-entropy or min-KL principle, subject to expected marginal constraints on the training examples
- Max-margin Markov network
 - Multi-variate, rather than uni-variate output \mathbf{Y}
 - Variable in the outputs are not independent of each other (structured input/output)
 - Margin constraint over every possible configuration of \mathbf{Y} (exponentially many!)

© Eric Xing @ CMU, 2006-2011

38