

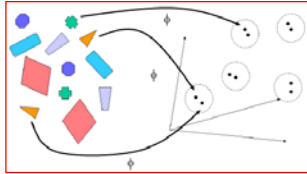
Machine Learning

10-701/15-781, Fall 2011

Advanced topics in Max-Margin Learning

Eric Xing

Lecture 20, November 21, 2011



© Eric Xing @ CMU, 2006-2010

1

Recap: the SVM problem

- We solve the following constrained opt problem:

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

$$\text{s.t.} \quad \alpha_i \geq 0, \quad i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y_i = 0.$$

- This is a **quadratic programming** problem.

- A global maximum of α_i can always be found.

- The solution:

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$$

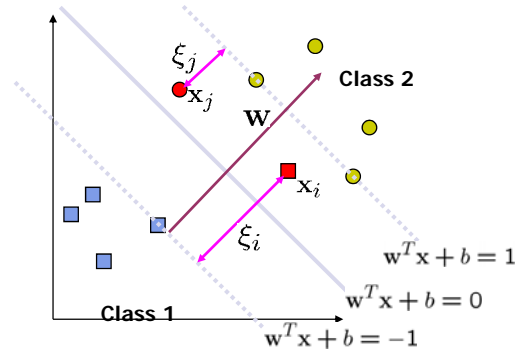
- How to predict:

$$\mathbf{w}^T \mathbf{x}_{\text{new}} + b \leq 0$$

© Eric Xing @ CMU, 2006-2010

2

Non-linearly Separable Problems



- We allow “error” ξ_i in classification; it is based on the output of the discriminant function $w^T x + b$
- ξ_i approximates the number of misclassified samples

© Eric Xing @ CMU, 2006-2010

3

Soft Margin Hyperplane

- Now we have a slightly different opt problem:

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

- ξ_i are “slack variables” in optimization
- Note that $\xi_i=0$ if there is no error for x_i
- ξ_i is an upper bound of the number of errors
- C : tradeoff parameter between error and margin

© Eric Xing @ CMU, 2006-2010

4

Lagrangian Duality, cont.



- Recall the Primal Problem:

$$\min_w \max_{\alpha, \beta, \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta)$$

- The Dual Problem:

$$\max_{\alpha, \beta, \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta)$$

- Theorem (weak duality):**

$$d^* = \max_{\alpha, \beta, \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta) \leq \min_w \max_{\alpha, \beta, \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta) = p^*$$

- Theorem (strong duality):**

Iff there exist a saddle point of $\mathcal{L}(w, \alpha, \beta)$, we have

$$d^* = p^*$$

© Eric Xing @ CMU, 2006-2011

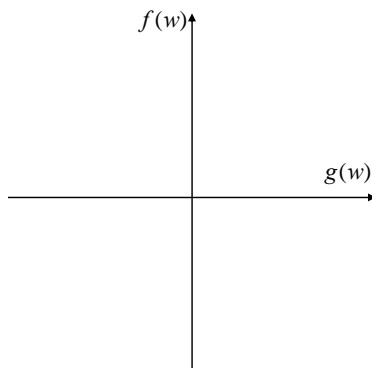
5

A sketch of strong and weak duality



- Now, ignoring $h(x)$ for simplicity, let's look at what's happening graphically in the duality theorems.

$$d^* = \max_{\alpha_i \geq 0} \min_w f(w) + \alpha^T g(w) \leq \min_w \max_{\alpha_i \geq 0} f(w) + \alpha^T g(w) = p^*$$



© Eric Xing @ CMU, 2006-2011

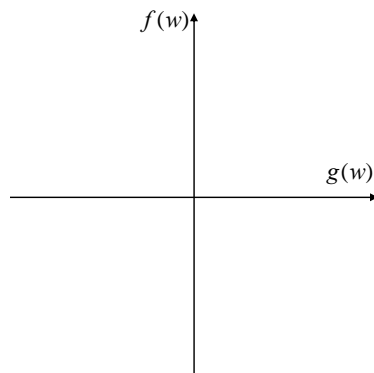
6

A sketch of strong and weak duality



- Now, ignoring $h(x)$ for simplicity, let's look at what's happening graphically in the duality theorems.

$$d^* = \max_{\alpha_i \geq 0} \min_w f(w) + \alpha^T g(w) \leq \min_w \max_{\alpha_i \geq 0} f(w) + \alpha^T g(w) = p^*$$



© Eric Xing @ CMU, 2006-2011

7

The KKT conditions



- If there exists some saddle point of \mathcal{L} , then the saddle point satisfies the following "Karush-Kuhn-Tucker" (KKT) conditions:

$$\frac{\partial}{\partial w_i} \mathcal{L}(w, \alpha, \beta) = 0, \quad i = 1, \dots, k$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w, \alpha, \beta) = 0, \quad i = 1, \dots, l$$

$$\alpha_i g_i(w) = 0, \quad i = 1, \dots, m$$

$$g_i(w) \leq 0, \quad i = 1, \dots, m$$

$$\alpha_i \geq 0, \quad i = 1, \dots, m$$

- Theorem:** If w^* , α^* and β^* satisfy the KKT condition, then it is also a solution to the primal and the dual problems.

© Eric Xing @ CMU, 2006-2011

8

The Optimization Problem



- The dual of this new constrained optimization problem is

$$\begin{aligned} \max_{\alpha} \quad & \mathcal{J}(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y_i = 0. \end{aligned}$$

- This is very similar to the optimization problem in the linear separable case, except that there is an upper bound C on α_i now
- Once again, a QP solver can be used to find α_i

The SMO algorithm

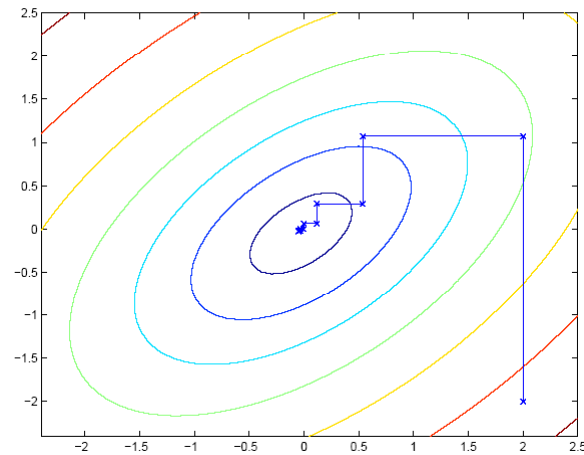


- Consider solving the **unconstrained** opt problem:

$$\max_{\alpha} W(\alpha_1, \alpha_2, \dots, \alpha_m)$$

- We've already see three opt algorithms!
 - Coordinate ascent
 - Gradient ascent
 - Newton-Raphson
- Coordinate ascend:

Coordinate ascend



© Eric Xing @ CMU, 2006-2010

11

Sequential minimal optimization

- Constrained optimization:

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

$$\text{s.t.} \quad 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m$$

$$\sum_{i=1}^m \alpha_i y_i = 0.$$

- Question: can we do coordinate along one direction at a time (i.e., hold all $\alpha_{[-i]}$ fixed, and update α_i ?)

© Eric Xing @ CMU, 2006-2010

12

The SMO algorithm



Repeat till convergence

1. Select some pair α_i and α_j to update next (using a heuristic that tries to pick the two that will allow us to make the biggest progress towards the global maximum).
2. Re-optimize $J(\alpha)$ with respect to α_i and α_j , while holding all the other α_k 's ($k \neq i, j$) fixed.

Will this procedure converge?

Convergence of SMO



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

$$\text{KKT:} \quad \begin{aligned} \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, k \\ & \sum_{i=1}^m \alpha_i y_i = 0. \end{aligned}$$

- Let's hold $\alpha_3, \dots, \alpha_m$ fixed and reopt J w.r.t. α_1 and α_2

Convergence of SMO

- The constraints:

$$\alpha_1 y_1 + \alpha_2 y_2 = \xi$$

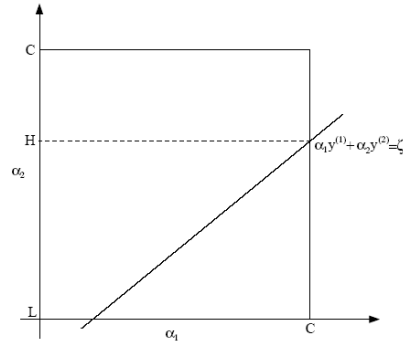
$$0 \leq \alpha_1 \leq C$$

$$0 \leq \alpha_2 \leq C$$

- The objective:

$$\mathcal{J}(\alpha_1, \alpha_2, \dots, \alpha_m) = \mathcal{J}((\xi - \alpha_2 y_2) y_1, \alpha_2, \dots, \alpha_m)$$

- Constrained opt:



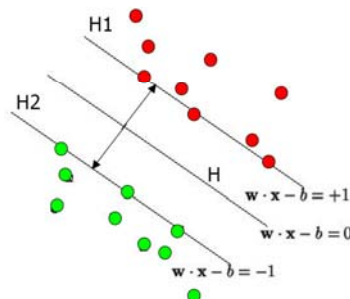
© Eric Xing @ CMU, 2006-2010

15

Cross-validation error of SVM

- The leave-one-out cross-validation error does not depend on the dimensionality of the feature space but only on the # of support vectors!

$$\text{Leave-one-out CV error} = \frac{\# \text{ support vectors}}{\# \text{ of training examples}}$$



© Eric Xing @ CMU, 2006-2010

16

Advanced topics in Max-Margin Learning



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

$$\mathbf{w}^T \mathbf{x}_{\text{new}} + b \leq 0$$

- Kernel
- Point rule or average rule
- Can we predict $\text{vec}(y)$?

© Eric Xing @ CMU, 2006-2010

17

Outline



- The Kernel trick
- Maximum entropy discrimination
- Structured SVM, aka, Maximum Margin Markov Networks

© Eric Xing @ CMU, 2006-2010

18

(1) Non-linear Decision Boundary



- So far, we have only considered large-margin classifier with a linear decision boundary
- How to generalize it to become nonlinear?
- Key idea: transform \mathbf{x}_i to a higher dimensional space to “make life easier”
 - Input space: the space the point \mathbf{x}_i are located
 - Feature space: the space of $\phi(\mathbf{x}_i)$ after transformation
- Why transform?
 - Linear operation in the feature space is equivalent to non-linear operation in input space
 - Classification can become easier with a proper transformation. In the XOR problem, for example, adding a new feature of x_1x_2 make the problem linearly separable (homework)

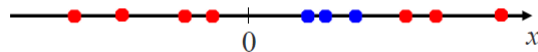
© Eric Xing @ CMU, 2006-2010

19

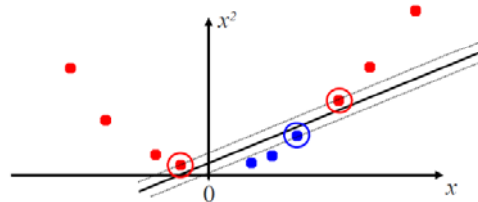
The Kernel Trick



- Is this data linearly-separable?



- How about a quadratic mapping $\phi(x_i)$?



© Eric Xing @ CMU, 2006-2010

20

The Kernel Trick



- Recall the SVM optimization problem

$$\begin{aligned} \max_{\alpha} \quad & \mathcal{J}(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y_i = 0. \end{aligned}$$

- The data points only appear as **inner product**
- As long as we can calculate the inner product in the feature space, we do not need the mapping explicitly
- Many common geometric operations (angles, distances) can be expressed by inner products
- Define the kernel function K by $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$

© Eric Xing @ CMU, 2006-2010

21

II. The Kernel Trick



- Computation depends on feature space
 - Bad if its dimension is much larger than input space

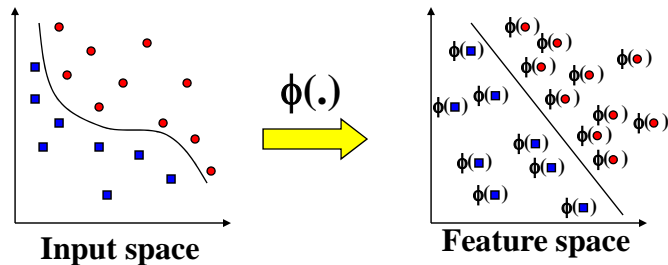
$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad i = 1, \dots, k \\ & \sum_{i=1}^m \alpha_i y_i = 0. \end{aligned}$$

Where $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ $y^*(z) = \text{sign} \left(\sum_{i \in SV} \alpha_i y_i K(\mathbf{x}_i, z) + b \right)$

© Eric Xing @ CMU, 2006-2010

22

Transforming the Data



Note: feature space is of higher dimension than the input space in practice

- Computation in the feature space can be costly because it is high dimensional
 - The feature space is typically infinite-dimensional!
- The kernel trick comes to rescue

© Eric Xing @ CMU, 2006-2010

23

An Example for feature mapping and kernels

- Consider an input $\mathbf{x}=[x_1, x_2]$
- Suppose $\phi(\cdot)$ is given as follows

$$\phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = 1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2$$

- An inner product in the feature space is

$$\left\langle \phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right), \phi\left(\begin{bmatrix} x_1' \\ x_2' \end{bmatrix}\right) \right\rangle =$$

- So, if we define the **kernel function** as follows, there is no need to carry out $\phi(\cdot)$ explicitly

$$K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^2$$

© Eric Xing @ CMU, 2006-2010

24

More examples of kernel functions



- Linear kernel (we've seen it)

$$K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

- Polynomial kernel (we just saw an example)

$$K(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^p$$

where $p = 2, 3, \dots$ To get the feature vectors we concatenate all p th order polynomial terms of the components of \mathbf{x} (weighted appropriately)

- Radial basis kernel

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|^2\right)$$

In this case the feature space consists of functions and results in a non-parametric classifier.

© Eric Xing @ CMU, 2006-2010

25

The essence of kernel



- Feature mapping, but “without paying a cost”

- E.g., polynomial kernel

$$K(x, z) = (x^T z + c)^d$$

- How many dimensions we've got in the new space?
- How many operations it takes to compute $K()$?

- Kernel design, any principle?

- $K(x, z)$ can be thought of as a similarity function between x and z
- This intuition can be well reflected in the following “Gaussian” function (Similarly one can easily come up with other $K()$ in the same spirit)

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$$

- Is this necessarily lead to a “legal” kernel?
(in the above particular case, $K()$ is a legal one, do you know how many dimension $\phi(x)$ is?)

© Eric Xing @ CMU, 2006-2010

26

Kernel matrix



- Suppose for now that K is indeed a valid kernel corresponding to some feature mapping ϕ , then for x_1, \dots, x_m , we can compute an $m \times m$ matrix $K = \{K_{i,j}\}$, where $K_{i,j} = \phi(x_i)^T \phi(x_j)$
- This is called a **kernel matrix**!
- Now, if a kernel function is indeed a valid kernel, and its elements are dot-product in the transformed feature space, it must satisfy:
 - Symmetry $K = K^T$
proof $K_{i,j} = \phi(x_i)^T \phi(x_j) = \phi(x_j)^T \phi(x_i) = K_{j,i}$
 - Positive –semidefinite $y^T K y \geq 0 \quad \forall y$
proof?

© Eric Xing @ CMU, 2006-2010

27

Mercer kernel

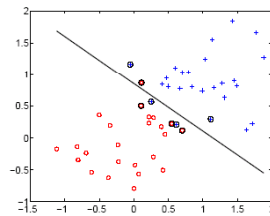


Theorem (Mercer): Let $K: \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}$ be given. Then for K to be a valid (Mercer) kernel, it is necessary and sufficient that for any $\{x_i, \dots, x_m\}$, ($m < \infty$), the corresponding kernel matrix is symmetric positive semi-definite.

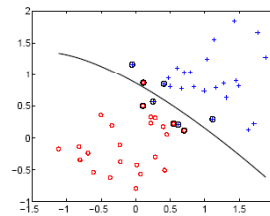
© Eric Xing @ CMU, 2006-2010

28

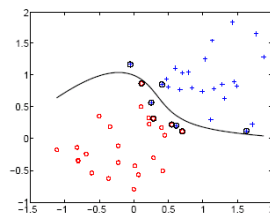
SVM examples



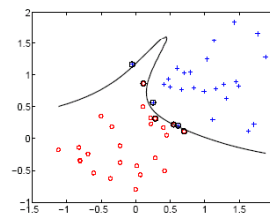
linear



2nd order polynomial



4th order polynomial



8th order polynomial

© Eric Xing @ CMU, 2006-2010

29

(2) Model averaging

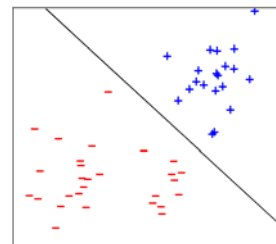
- Inputs \mathbf{x} , class $y = +1, -1$
- data $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)\}$

- Point Rule:

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

- E.g., SVM

$$f^{\text{opt}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}_{\text{new}} + b$$



© Eric Xing @ CMU, 2006-2010

30

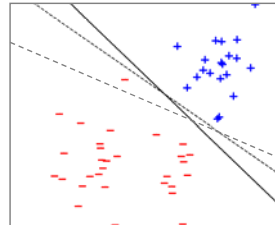
Model averaging

- There exist many f with near optimal performance

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

$Q(f)$ = weight of f

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



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

© Eric Xing @ CMU, 2006-2010

31

Recall Bayesian Inference

- Bayesian learning:



$$\text{Bayes Thrm : } p(w|\mathcal{D}) = \frac{p(w)p(\mathcal{D}|w)}{p(\mathcal{D})}$$

- Bayes Predictor (model averaging):

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

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

- What p_0 ?

© Eric Xing @ CMU, 2006-2010

32

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, 2006-2010

33

Maximum Entropy Discrimination



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

$$\begin{aligned} Q_{\text{ME}} &= \arg \max Q \quad H(Q) \\ \text{s.t.} \quad &y^i \langle f(\mathbf{x}^i) \rangle_{Q_{\text{ME}}} \geq \xi_i, \quad \forall i \\ &\xi_i \geq 0 \quad \forall i \end{aligned}$$

- 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, 2006-2010

34

Introducing Priors

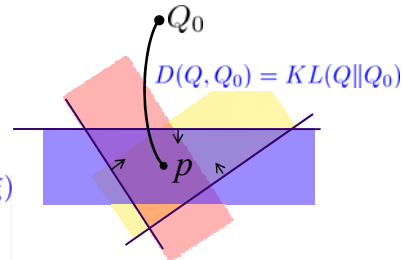
- Prior $Q_0(f)$

- Minimum Relative Entropy Discrimination

$$Q_{\text{MRE}} = \arg \min Q \text{ KL}(Q \| Q_0) + U(\xi)$$

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

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



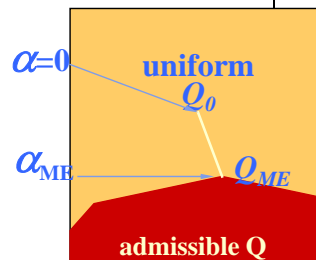
- Convex problem: Q_{MRE} unique solution
- MER \rightarrow "minimum *additional* assumption" over Q_0 about f

Solution: Q_{ME} as a projection

- Convex problem: Q_{ME} unique

- Theorem:

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



$\alpha_i \geq 0$ Lagrange multipliers

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

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:

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

$U^*(\cdot)$ 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

Examples: SVMs



- Theorem

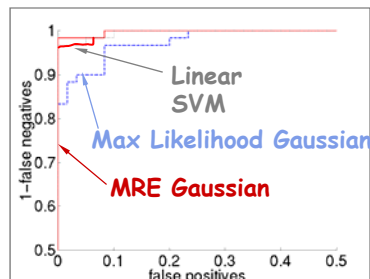
For $f(x) = w^T x + b$, $Q_0(w) = \text{Normal}(0, I)$, $Q_0(b) = \text{non-informative prior}$, the Lagrange multipliers α are obtained by maximizing $J(\alpha)$ subject to $0 \leq \alpha_t \leq C$ and $\sum_t \alpha_t y_t = 0$, where

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

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

SVM extensions

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



© Eric Xing @ CMU, 2006-2010

39

(3) Structured Prediction

- Unstructured prediction



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

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

- 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



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

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

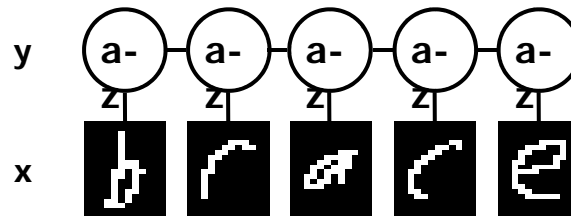
© Eric Xing @ CMU, 2006-2010

40

OCR example



Sequential structure



© Eric Xing @ CMU, 2006-2010

41

Classical Classification Models

- Inputs:

- a set of training samples $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, where $x_i = [x_i^1, x_i^2, \dots, x_i^d]^\top$ and $y_i \in C \triangleq \{c_1, c_2, \dots, c_L\}$

- Outputs:

- a predictive function $h(x)$: $y^* = h(x) \triangleq \arg \max_y F(x, y)$
 $F(x, y) = \mathbf{w}^\top \mathbf{f}(x, y)$

- Examples:

- SVM: $\max_{\mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{i=1}^N \xi_i; \text{ s.t. } \mathbf{w}^\top \Delta \mathbf{f}_i(y) \geq 1 - \xi_i, \forall i, \forall y.$

- Logistic Regression: $\max_{\mathbf{w}} \mathcal{L}(\mathcal{D}; \mathbf{w}) \triangleq \sum_{i=1}^N \log p(y_i | x_i)$

where

$$p(y|x) = \frac{\exp\{\mathbf{w}^\top \mathbf{f}(x, y)\}}{\sum_{y'} \exp\{\mathbf{w}^\top \mathbf{f}(x, y')\}}$$

© Eric Xing @ CMU, 2006-2010

42

Structured Models



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

↑
space of feasible outputs

↑
discriminant function

- Assumptions:

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

- Linear combination of features
- Sum of partial scores: index p represents a part in the structure

- Random fields or Markov network features:



© Eric Xing @ CMU, 2006-2010

43

Discriminative Learning Strategies



- Max Conditional Likelihood

- We predict based on:

$$\mathbf{y}^* | \mathbf{x} = \arg \max_{\mathbf{y}} p_{\mathbf{w}}(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{w}, \mathbf{x})} \exp \left\{ \sum_c w_c f_c(\mathbf{x}, \mathbf{y}_c) \right\}$$

- And we learn based on:

$$\mathbf{w}^* | \{\mathbf{y}_i, \mathbf{x}_i\} = \arg \max_{\mathbf{w}} \prod_i p_{\mathbf{w}}(\mathbf{y}_i | \mathbf{x}_i) = \prod_i \frac{1}{Z(\mathbf{w}, \mathbf{x}_i)} \exp \left\{ \sum_c w_c f_c(\mathbf{x}_i, \mathbf{y}_i) \right\}$$

- Max Margin:

- We predict based on:

$$\mathbf{y}^* | \mathbf{x} = \arg \max_{\mathbf{y}} \sum_c w_c f_c(\mathbf{x}, \mathbf{y}_c) = \arg \max_{\mathbf{y}} \mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y})$$

- And we learn based on:

$$\mathbf{w}^* | \{\mathbf{y}_i, \mathbf{x}_i\} = \arg \max_{\mathbf{w}} \left(\min_{\mathbf{y} \neq \mathbf{y}', \forall i} \mathbf{w}^T (f(\mathbf{y}_i, \mathbf{x}_i) - f(\mathbf{y}, \mathbf{x}_i)) \right)$$

© Eric Xing @ CMU, 2006-2010

44

E.g. Max-Margin Markov Networks



- Convex Optimization Problem:

$$\begin{aligned} P0 (M^3N) : \quad & \min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t. } \forall i, \forall \mathbf{y} \neq \mathbf{y}_i : \quad & \mathbf{w}^\top \Delta \mathbf{f}_i(\mathbf{x}, \mathbf{y}) \geq \Delta \ell_i(\mathbf{y}) - \xi_i, \quad \xi_i \geq 0, \end{aligned}$$

- 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:

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

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!

Min-max Formulation

- Brute force enumeration of constraints:

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

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

- The constraints are exponential in the size of the structure

- Alternative: min-max formulation

- add only the most violated constraint

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

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

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

© Eric Xing @ CMU, 2006-2010

47

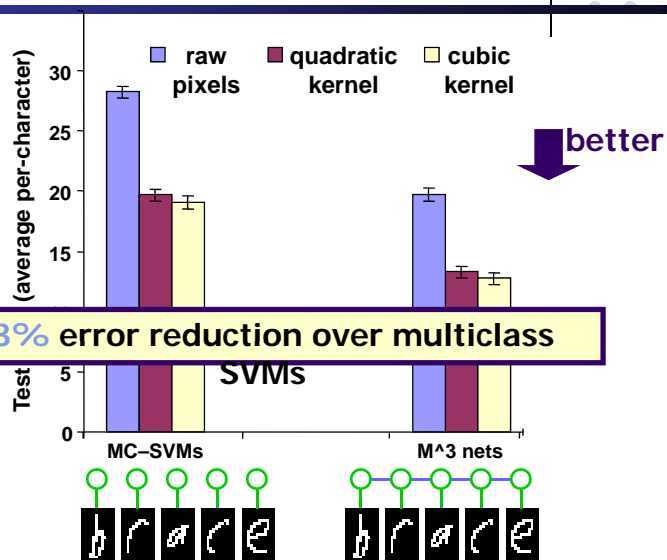
Results: Handwriting Recognition

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

Models:

Multiclass-SVMs

M³ nets

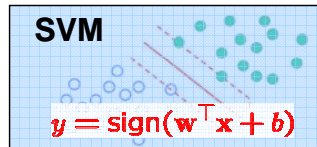


Crammer & Singer 01

© Eric Xing @ CMU, 2006-2010

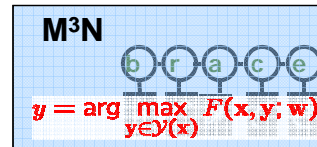
48

Discriminative Learning Paradigms



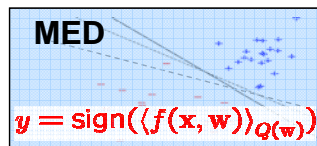
$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

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



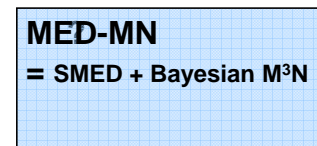
$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i$$

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



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

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



See [Zhu and Xing, 2008]

© Eric Xing @ CMU, 2006-2010

49

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, 2006-2010

50