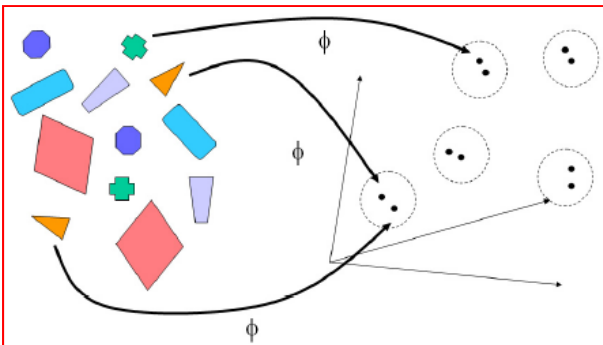# Machine Learning

**10-701, Fall 2016**

## Advanced topics in Max-Margin Learning
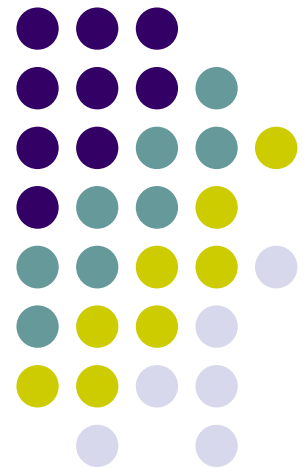
**Eric Xing**

**Lecture 7, September 28, 2016**

**Reading: class handouts**

1

# Recap: the SVM problem

- We solve the following constrained opt problem:

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

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

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

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

  - A global maximum of $\alpha_i$ can always be found.

  - The solution: $\quad w = \sum_{i=1}^{m} \alpha_i y_i \mathbf{x}_i$

  - How to predict: $\quad \mathbf{w}^T \mathbf{x}_{\text{new}} + b \lessgtr 0$

# The SMO algorithm

- Consider solving the unconstrained opt problem:

$$\vec{\alpha}^* = arg \quad \max_{\alpha} W(\alpha_1, \alpha_2, \ldots, \alpha_m)$$

$$a^{t+1} = a^t + f(\partial \alpha)$$
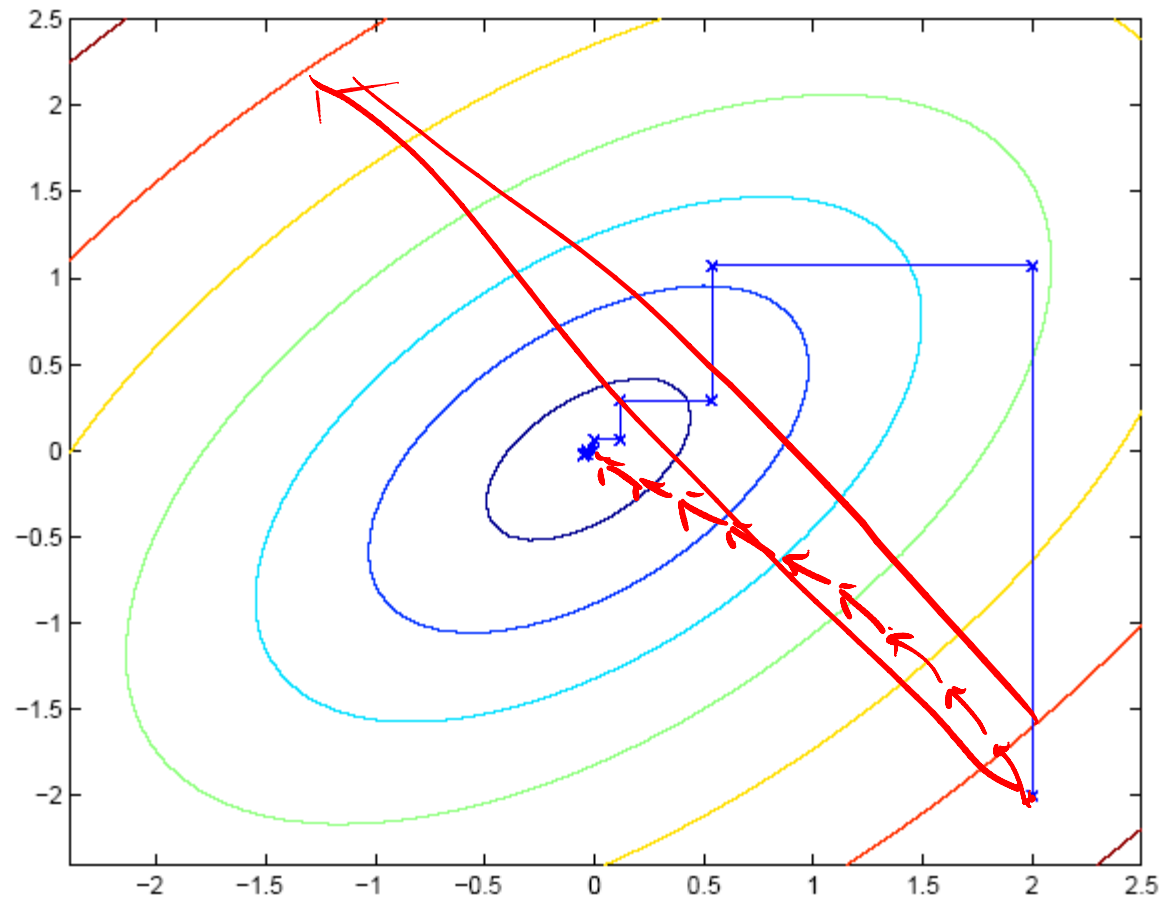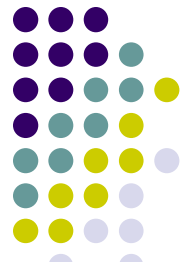
- We've already see three opt algorithms!
  - Coordinate ascent
  - Gradient ascent
  - Newton-Raphson

$$\Delta \alpha_i = \frac{\partial W}{\partial \alpha_i}$$

$$\Delta \vec{\alpha} = \frac{\partial W}{\partial \alpha} = \left( \frac{\partial W}{\partial \alpha_i}, \frac{\partial W}{\partial \alpha_2} \cdots \right)$$

$$\partial \vec{\alpha} = \frac{1}{\Delta W} \frac{\partial W}{\partial \alpha}$$

- Coordinate ascend:

# Coordinate ascend

# Sequential minimal optimization

- Constrained optimization:

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

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

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

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

$$\Delta \alpha_i = \frac{\partial \bar{\mathcal{J}}}{\partial \alpha_i}$$

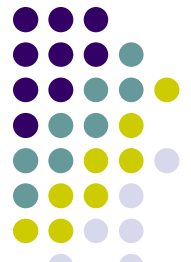$$\sum_{i \ne j} \alpha_i y_i + \alpha_j y_i = 0$$

# The SMO algorithm

$$\sum_{k \neq i,j} \alpha_k y_k = -\zeta$$

$$\alpha_i y_i + \alpha_j y_j = \zeta$$

Repeat till convergence

1. Select some pair $\alpha_i$ and $\alpha_j$ 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 $\alpha_i$ and $\alpha_j$, while holding all the other $\alpha_k$'s ($k \neq i; j$) fixed.

Will this procedure converge?

# Convergence of SMO

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

KKT:
$$\text{s.t.} \quad 0 \le \alpha_i \le C, \quad i = 1, \ldots, k$$
$$\sum_{i=1}^m \alpha_i y_i = 0.$$

- Let's hold $\alpha_3, \ldots, \alpha_m$ fixed and reopt J w.r.t. $\alpha_1$ and $\alpha_2$

# Convergence of SMO

- The constraints:

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

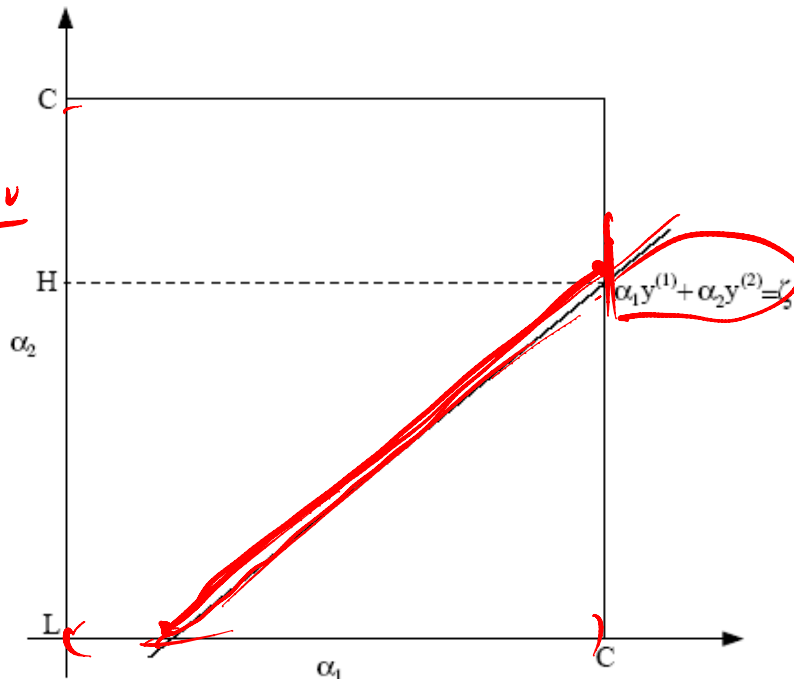$$0 \leq \alpha_1 \leq C$$

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

$$\alpha_1 = \frac{\zeta - \alpha_2 y_2}{y_1}$$



- The objective:

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

$$\frac{\partial \mathcal{J}}{\partial \alpha_2}$$

- Constrained opt:
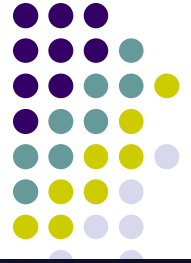
# Advanced topics in Max-Margin Learning

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

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

- Kernel

- Point rule or average rule
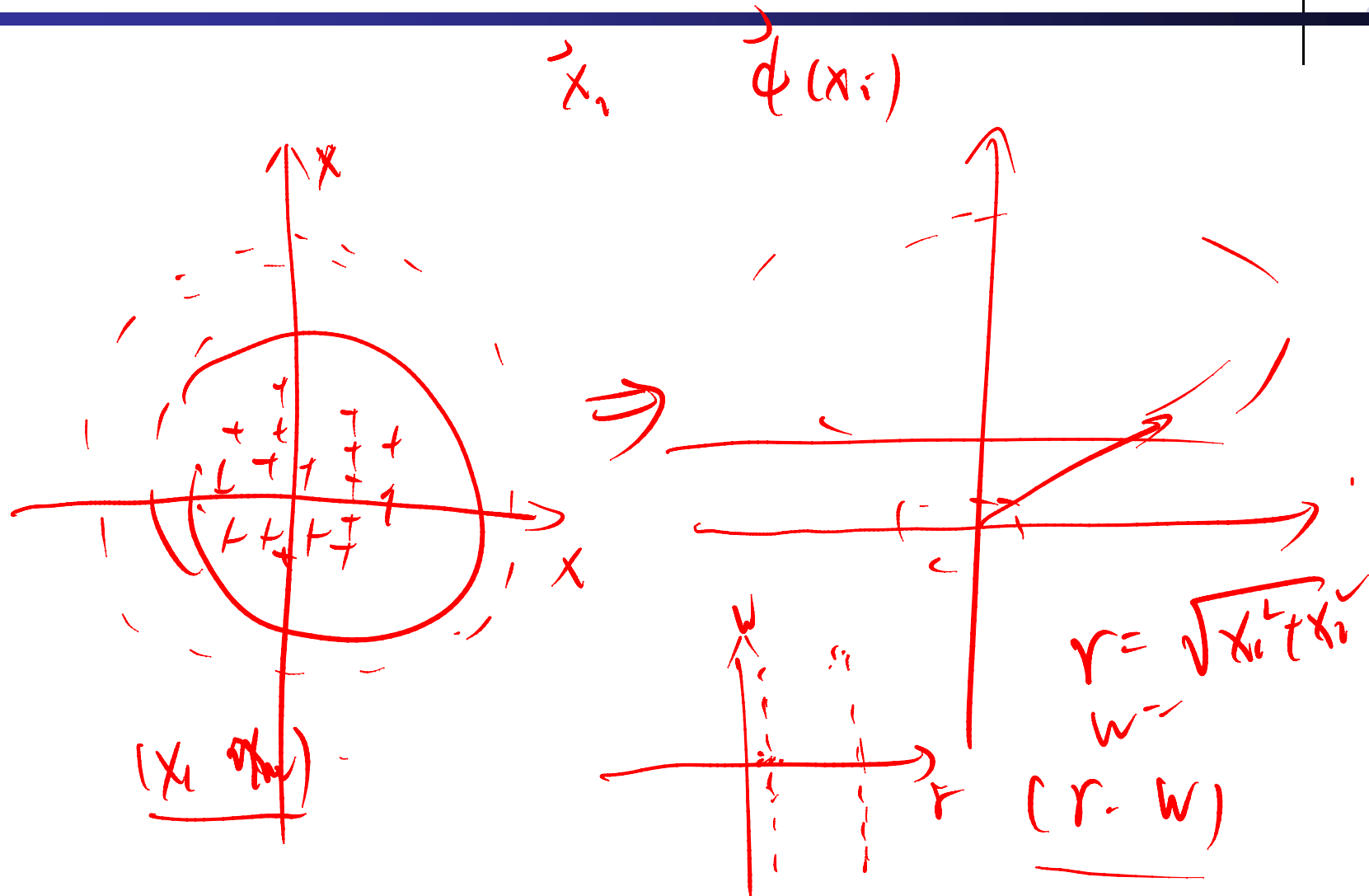
- Can we predict vec(y)?

# Outline

- The Kernel trick

- Maximum entropy discrimination

- Structured SVM, aka, Maximum Margin Markov Networks
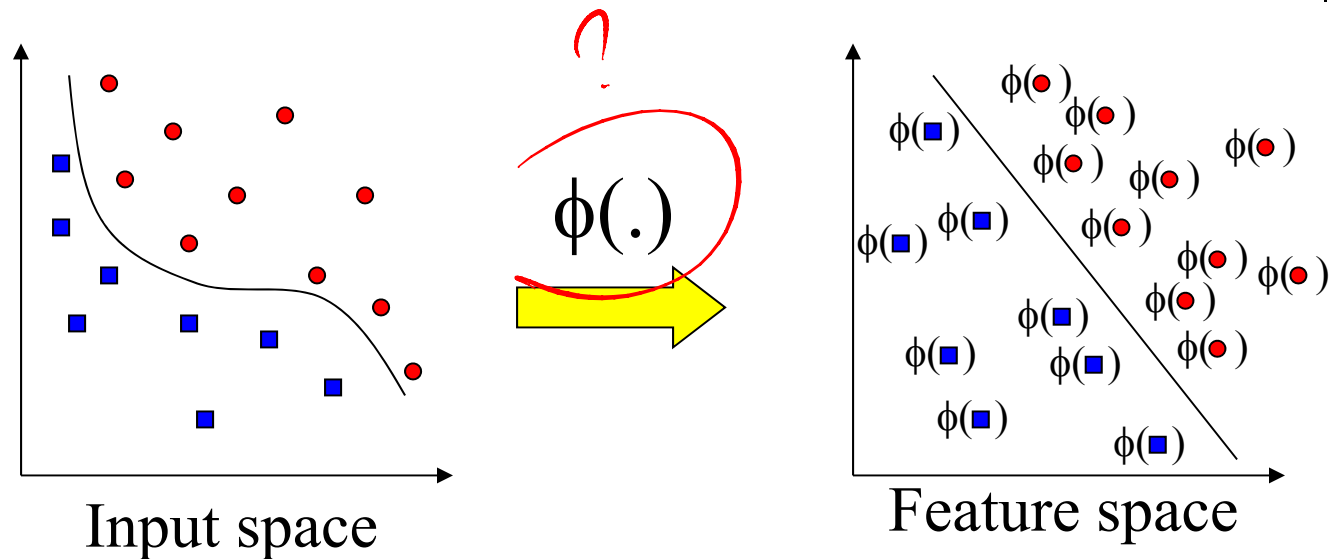
# (1) Non-linear Decision Boundary

- So far, we have only considered large-margin classifier with a linear decision boundary

- How to generalize it to become nonlinear?

- Key idea: transform $x_i$ to a higher dimensional space to "make life easier"
  - Input space: the space the point $x_i$ are located
  - Feature space: the space of $\phi(x_i)$ 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)

# Non-linear Decision Boundary

# Transforming the Data



Input space

$\phi(.)$

Feature space

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

# The Kernel Trick

$F(\vec{x}, \vec{x}_i)$

$K(x_i, x_j) = (1 + x_i x_j)^p$

$= \phi(x_i)\phi(x_j)$

- Recall the SVM optimization problem

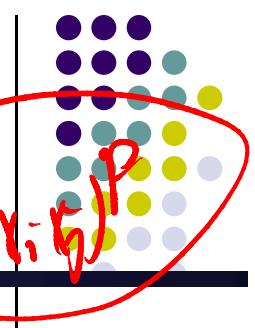$$\max_\alpha \quad \mathcal{J}(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j)$$

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

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

$\vec{W} \cdot \vec{X}_{new}$

$= \sum_{i \in SV} \alpha_i \vec{X}_i^T \cdot \vec{X}_{new}$

- 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 $\quad K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$

# An Example for feature mapping and kernels

- Consider an input $\mathbf{x} = [x_1, x_2]$

- Suppose $\phi(.)$ is given as follows

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

- An inner product in the feature space is

$$\left\langle \phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right), \phi\left(\begin{bmatrix} x_1' \\ x_2' \end{bmatrix}\right) \right\rangle = 1 + 2x_1x_1' + 2x_2x_2' + x_1^2x_1'^2 + (x_2^2x_2'^2) + 2x_1x_1'x_2x_2'$$

$$= (1 + x'^T x)^2$$

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

$$K(\mathbf{x}, \mathbf{x}') = \left(1 + \mathbf{x}^T \mathbf{x}'\right)^2 \qquad \left(1 + x^{\dagger}x'\right)^r$$

# More examples of kernel functions

$J \left( x_i^T, x_j \right)$

$K($  $)$

$K\overline{1}$

- Linear kernel (we've seen it)

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

$\phi(\cdot): \begin{pmatrix} n \\ p \end{pmatrix} \sim O(n^p)$

- Polynomial kernel (we just saw an example)

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

$O(n+p)$

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

$\phi(\cdot)$

- Radial basis kernel

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

$\infty$

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

# 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?
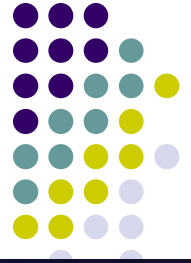  - How many operations it takes to compute K()?

  $\phi(\ )\phi(\ )$

- 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?

# Kernel matrix

- Suppose for now that $K$ is indeed a valid kernel corresponding to some feature mapping $\phi$, then for $x_1, \ldots, x_m$, 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:
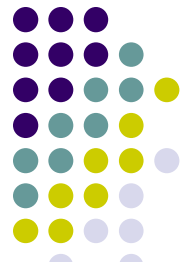
  - Symmetry $\qquad\qquad\qquad K = K^T$

    proof $\qquad K_{i,j} = \phi(x_i)^T \phi(x_j) = \phi(x_j)^T \phi(x_i) = K_{j,i}$
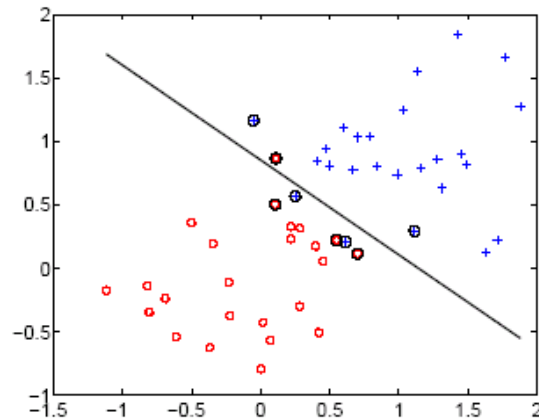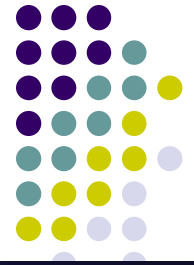
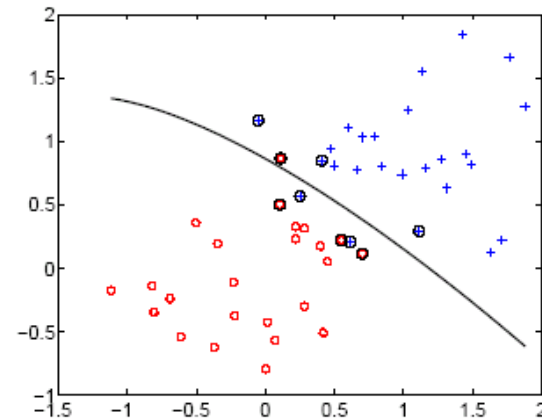  - Positive –semidefinite $\qquad y^T K y \geq 0 \quad \forall y$

    proof?

# 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 $\{x_i, \ldots, x_m\}$, $(m < \infty)$, the corresponding kernel matrix is symmetric positive semi-denite.
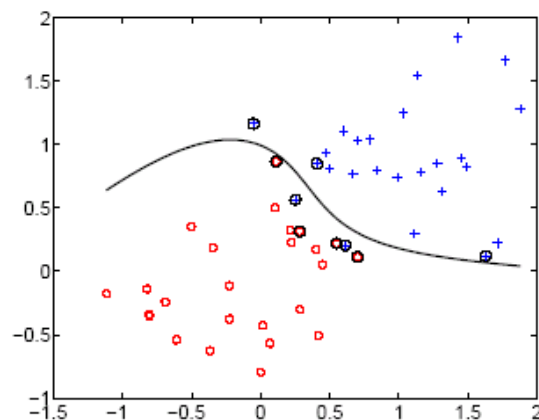
# SVM examples



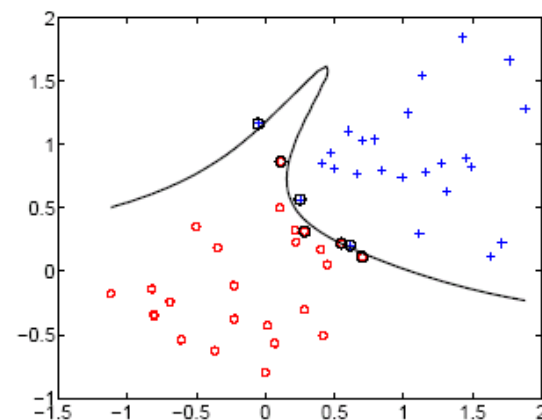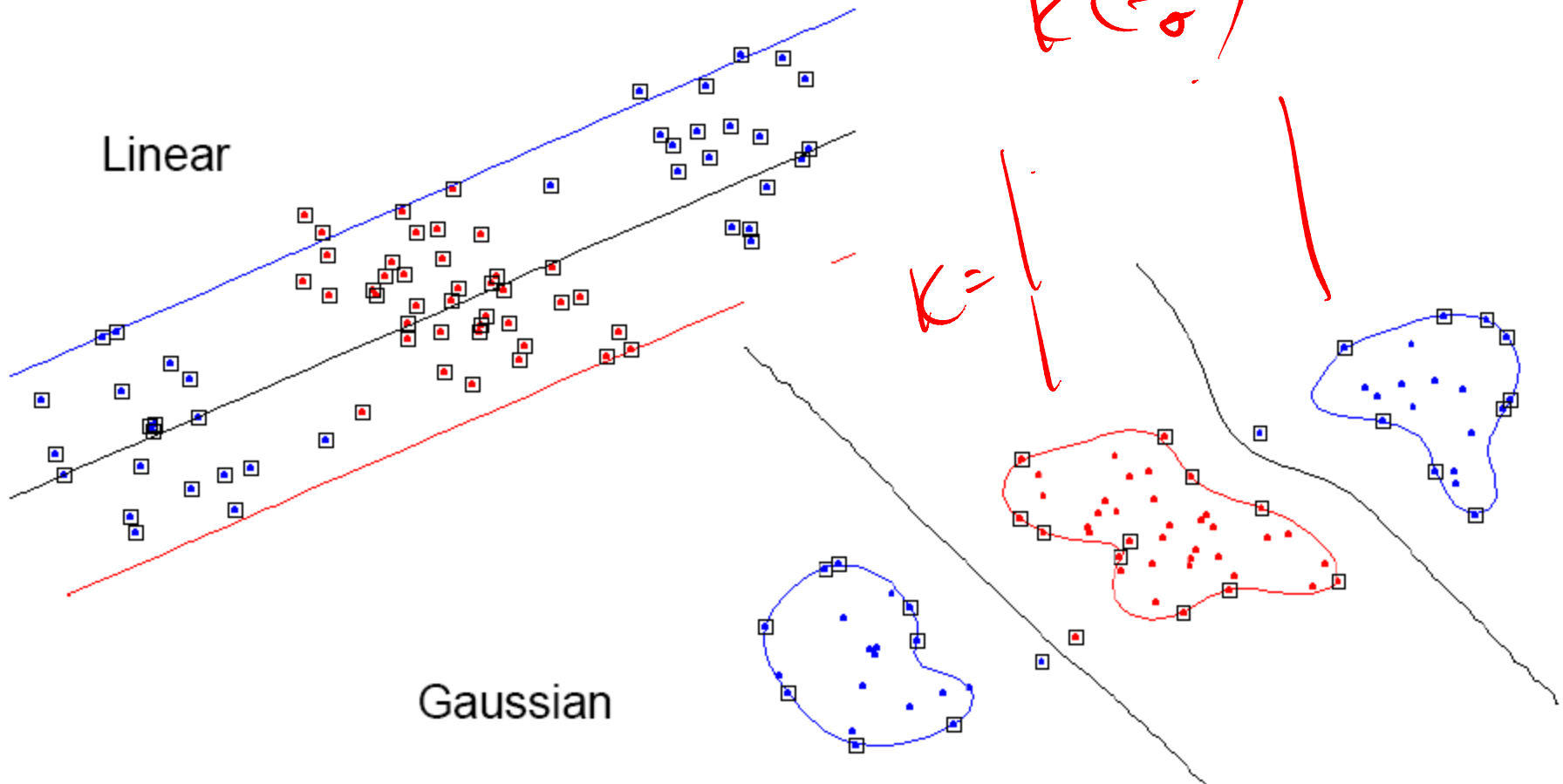linear

$2^{nd}$ order polynomial

$4^{th}$ order polynomial

$8^{th}$ order polynomial

# Examples for Non Linear SVMs – Gaussian Kernel



Linear

Gaussian

$$K\left(\frac{x_i - x_j}{\sigma}\right)$$

$$K = \begin{vmatrix} \\ \\ \end{vmatrix}$$

# (2) Model averaging

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

- Point Rule: $\quad w^*$

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

  - classify $y = sign\ f^{opt}(x)$

- E.g., SVM

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

# Model averaging

- There exist many **f** with near optimal performance

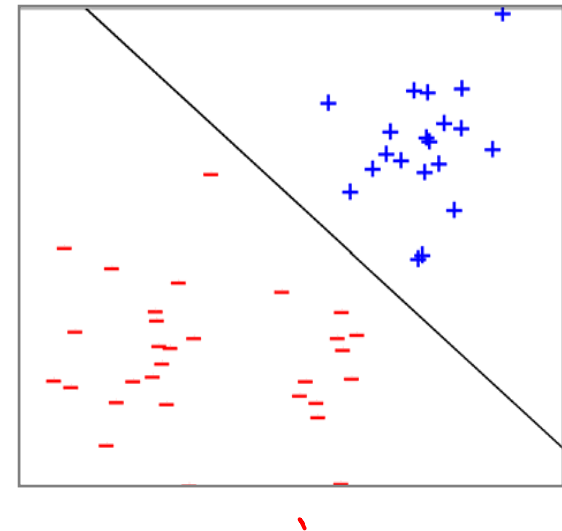- Instead of <u>choosing</u> **f**$^{opt}$,
  <u>average</u> over all **f** in **F**

  **Q(f)** = weight of **f**

  $$y(x) = \text{sign} \int_F Q(f)f(x)\,df$$
  $$= \text{sign}\langle f(x)\rangle_Q$$

- How to specify:
  **F** = { **f** } family of discriminant functions?

- How to learn **Q(f)** distribution over **F**?

# Recall Bayesian Inference

- Bayesian learning:

$$p_0(\mathbf{w})$$

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

Bayes Learner $\longrightarrow p(\mathbf{w}|\mathcal{D})$

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

$f($ $|w)$

- 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 $p_0$?

# 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:

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

# Maximum Entropy Discrimination

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

$$Q_{\mathrm{ME}} = \arg\max \quad \mathrm{H}(Q)$$

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

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

$$f( \cdot )$$

$$Q(w)$$

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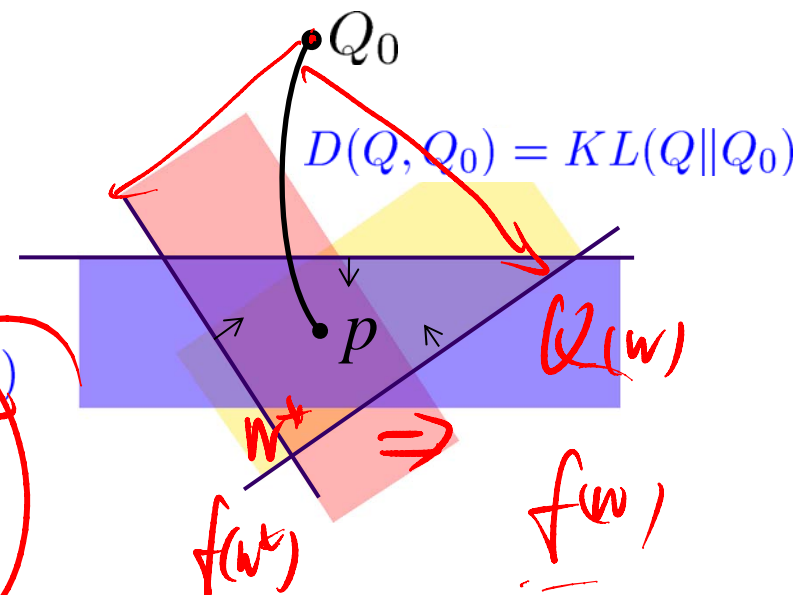- max entropy $\Longrightarrow$ "minimum assumption" about $f$

# Introducing Priors

- Prior $Q_0(f)$

- <u>Minimum Relative Entropy Discrimination</u>

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

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

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

$$Q_0$$

$$D(Q, Q_0) = KL(Q\|Q_0)$$

$$p$$

- Convex problem: $Q_{\mathrm{MRE}}$ unique solution
- MER $\Rightarrow$ "minimum *additional* assumption" over $Q_0$ about $f$
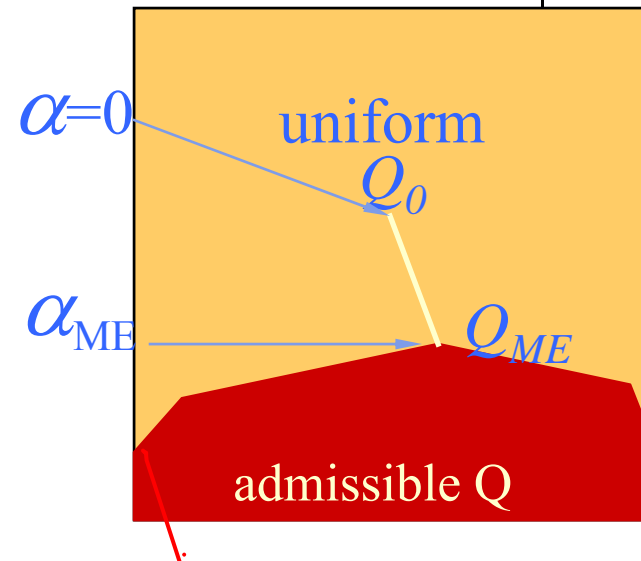
# Solution: Q ME as a projection

- Convex problem: $Q_{ME}$ unique

- Theorem: $f(w) \propto P(w) \, P_D(w)$

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

$\alpha_i \geq 0$ Lagrange multipliers



$\alpha=0$  uniform $Q_0$

$\alpha_{ME}$  $Q_{ME}$

admissible Q

- 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:

    $$D1: \quad \max_\alpha \ -\log Z(\alpha) - U^\star(\alpha)$$

    $$\text{s.t. } \alpha_i(\mathbf{y}) \geq 0, \ \forall i,$$

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

- Algorithm: to computer $\alpha_t$, $t = 1,...\mathrm{T}$

    - 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 $\alpha$ are obtained by maximizing $J(\alpha)$ subject to $0 \leq \alpha_t \leq C$ and $\sum_t \alpha_t y_t = 0$, where
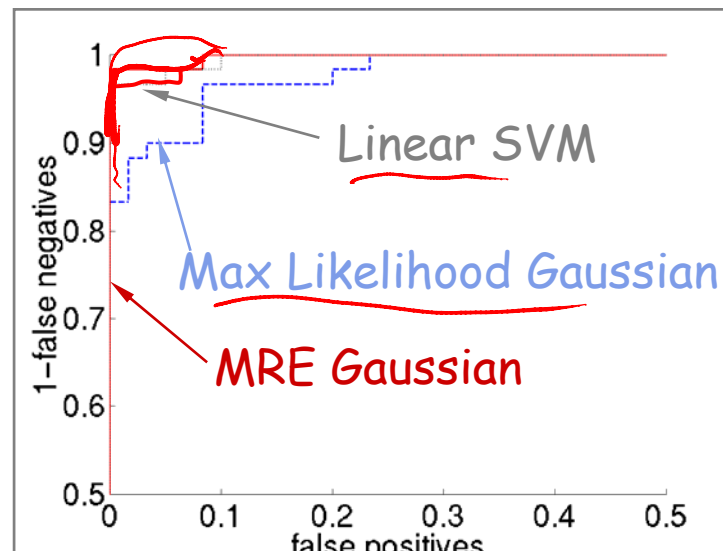
$$J(\alpha) = \sum_t \left[ \alpha_t + \log(1 - \alpha_t/C) \right] - \frac{1}{2} \sum_{s,t} \alpha_s \alpha_t y_s y_t x_s^T x_t$$

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

# SVM extensions
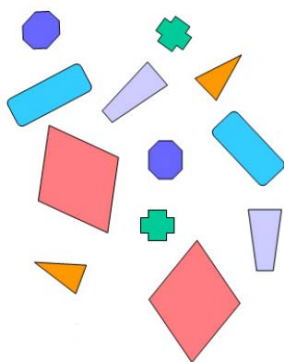
- Example: Leptograpsus Crabs (5 inputs, $T_{train}$=80, $T_{test}$=120)

# (3) Structured Prediction

- **Unstructured prediction**

$$\mathbf{x} = \begin{pmatrix} x_{11} & x_{12} & \cdots \\ x_{21} & x_{22} & \cdots \\ \vdots & \vdots & \cdots \end{pmatrix} \qquad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \end{pmatrix}$$

- **Structured prediction**

  - Part of speech tagging

    $\mathbf{x} =$ "Do you want sugar in it?"  $\Rightarrow$  $\mathbf{y} =$ verb pron verb noun prep pron>
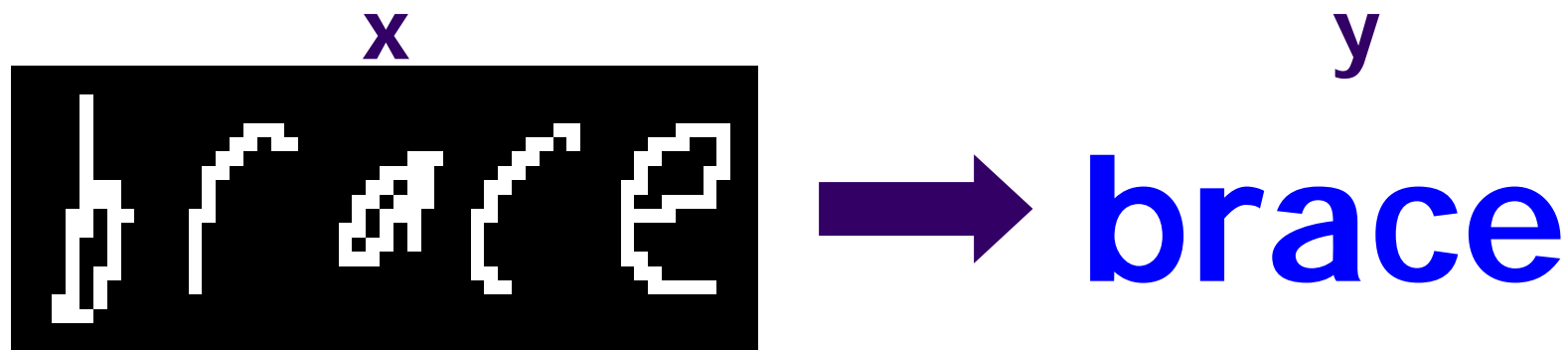
  - Image segmentation

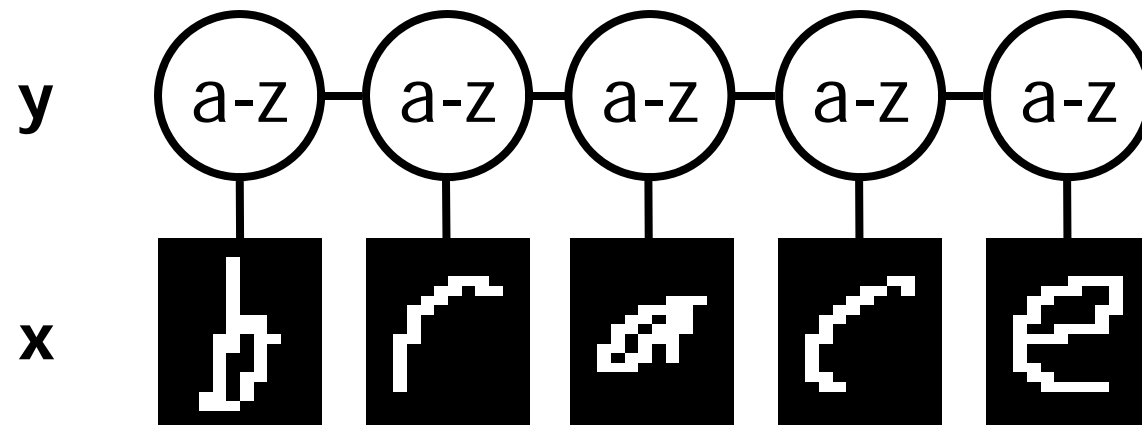$$\mathbf{x} = \begin{pmatrix} x_{11} & x_{12} & \cdots \\ x_{21} & x_{22} & \cdots \\ \vdots & \vdots & \cdots \end{pmatrix} \qquad \mathbf{y} = \begin{pmatrix} y_{11} & y_{12} & \cdots \\ y_{21} & y_{22} & \cdots \\ \vdots & \vdots & \cdots \end{pmatrix}$$

# OCR example



x → **brace**

Sequential structure

# Classical Classification Models

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

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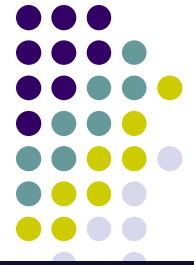  $$F(\mathbf{x}, y) = \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y)$$

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

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

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

# Structured Models

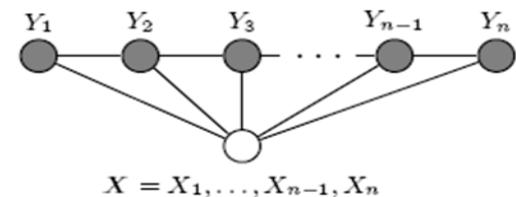$$h(\mathbf{x}) = \arg\max_{y \in \mathcal{Y}(\mathbf{x})} \ F(\mathbf{x}, \mathbf{y})$$

<span style="color:blue">discriminant function</span>

<span style="color:blue">space of feasible outputs</span>

- 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:



$$X = X_1, \ldots, X_{n-1}, X_n$$

# Discriminative Learning Strategies

- ## Max Conditional Likelihood

  - ### We predict based on:

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

  - ### And we learn based on:

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

- ## Max Margin:

  - ### We predict based on:

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

  - ### And we learn based on:

$$\mathbf{w}^* \mid \{\mathbf{y}_i, \mathbf{x}_i\} = \arg\max_{\mathbf{w}} \left(\min_{\mathbf{y} \neq \mathbf{y}^i, \forall i} \mathbf{w}^T \left(f(\mathbf{y}_i, \mathbf{x}_i) - f(\mathbf{y}, \mathbf{x}_i)\right)\right)$$

# E.g. Max-Margin Markov Networks

- Convex Optimization Problem:

$$\text{P0 (M}^3\text{N)} : \quad \min_{\mathbf{w},\xi} \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{N}\xi_i$$

$$\text{s.t. } \forall i, \forall \mathbf{y} \neq \mathbf{y}_i : \quad \mathbf{w}^\top \Delta \mathbf{f}_i(\mathbf{x},\mathbf{y}) \geq \Delta \ell_i(\mathbf{y}) - \xi_i, \ \xi_i \geq 0 \,,$$

- 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:

argmax$_{word}$ $\mathbf{w}^T \mathbf{f}(\text{[brace]}$ , **word**$)$ = "brace"

- Equivalently:

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

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

...

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

**a lot!**

# Min-max Formulation

- Brute force enumeration of constraints:

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

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

  - The constraints are exponential in the size of the structure

- Alternative: min-max formulation

  - add only the most violated constraint

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

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

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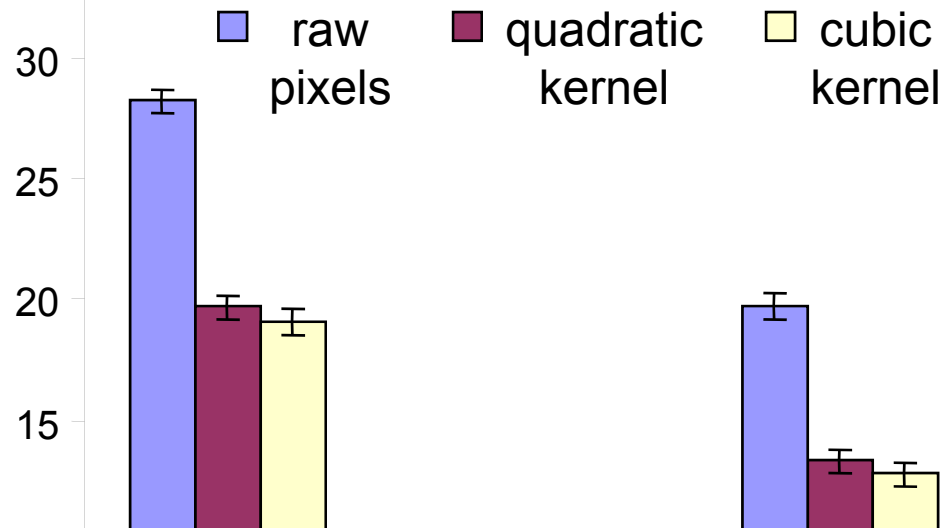  - Several algorithms exist …

# Results: Handwriting Recognition

Length: ~8 chars

Letter: 16x8 pixels

10-fold Train/Test

5000/50000 letters

600/6000 words
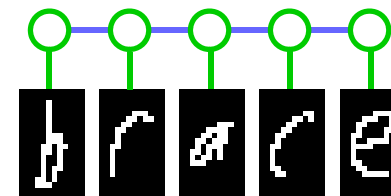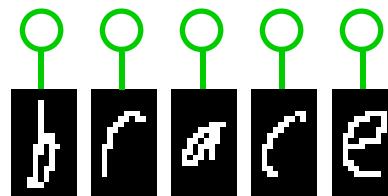
Models:

Multiclass-SV

$M^3$ nets

**33%** error reduction over multiclass SVMs

better



MC–SVMs

M^3 nets

*Crammer & Singer 01

# Discriminative Learning Paradigms

**SVM**



$$y = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$$

$$\min_{\mathbf{w},\xi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{m}\xi_i$$
$$y^i(\mathbf{w}^\top \mathbf{x}^i + b) \geq 1 - \xi_i, \quad \forall i$$

**M³N**



$$y = \arg\max_{\mathbf{y}\in\mathcal{Y}(\mathbf{x})} F(\mathbf{x},\mathbf{y};\mathbf{w})$$
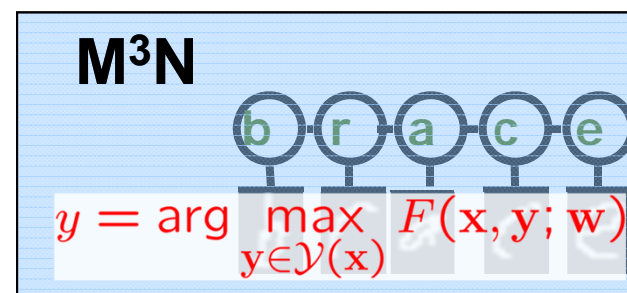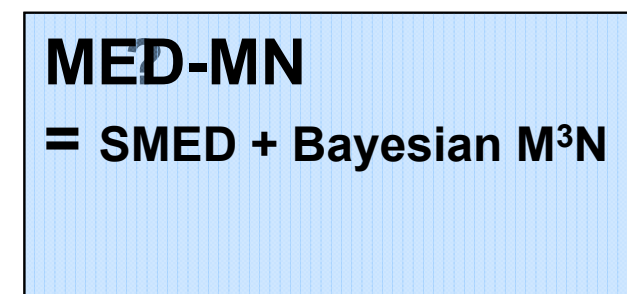
$$\min_{\mathbf{w},\xi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{m}\xi_i$$
$$\mathbf{w}^\top[\mathbf{f}(\mathbf{x}^i,\;) - \mathbf{f}(\mathbf{x}^i,\mathbf{y})] \geq \ell(\mathbf{y}^i,\mathbf{y}) - \xi_i, \quad \forall i, \forall \mathbf{y} \neq \mathbf{y}^i$$

**MED**



$$y = \text{sign}(\langle f(\mathbf{x},\mathbf{w})\rangle_{Q(\mathbf{w})})$$

$$\min_{Q} \quad \text{KL}(Q\|Q_0)$$
$$y^i\langle f(\mathbf{x}^i)\rangle_Q \geq \xi_i, \quad \forall i$$

**MED-MN**

**= SMED + Bayesian M³N**

See [Zhu and Xing, 2008]

# Summary

- ## Maximum margin nonlinear separator
  - Kernel trick
  - Project into linearly separatable 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 w who defines the separation hyperplane
  - P(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 Y
  - Variable in the outputs are not independent of each other (structured input/output)
  - Margin constraint over every possible configuration of Y (exponentially many!)