

Learning From Data Locally and Globally

Kaizhu Huang

Supervisors:
Prof. Irwin King,
Prof. Michael R. Lyu



Outline

- Background
 - Linear Binary classifier
 - Global Learning
 - Bayes optimal Classifier
 - Local Learning
 - Support Vector Machine
- Contributions
- Minimum Error Minimax Probability Machine (MEMPM)
 - Biased Minimax Probability Machine (BMPPM)
- Maxi-Min Margin Machine (M^4)
 - Local Support Vector Regression (LSVR)
- Future work
- Conclusion



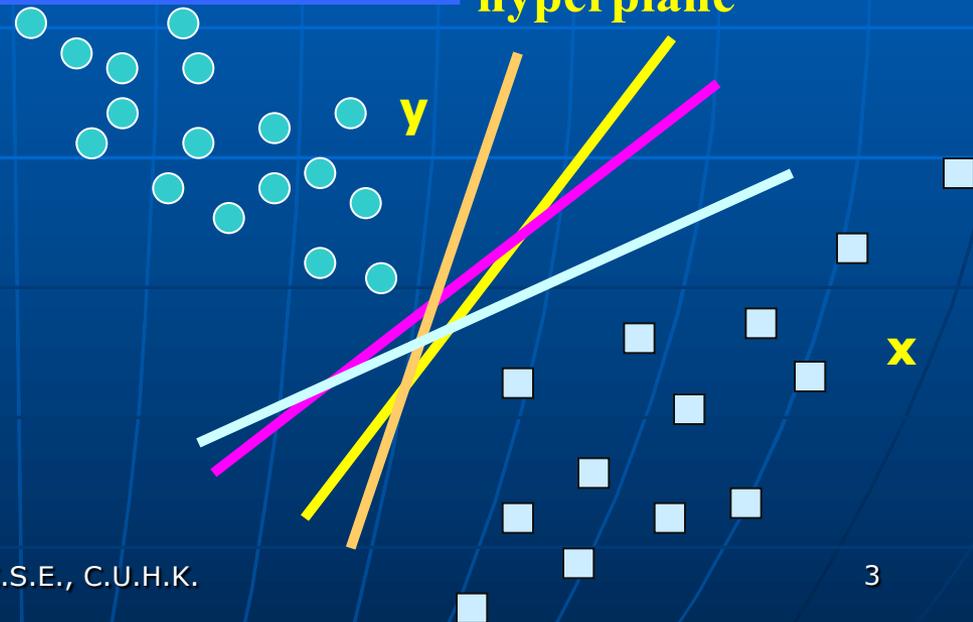
Background - Linear Binary Classifier

Given two classes of data sampled from \mathbf{x} and \mathbf{y} , we are trying to find a linear decision plane $\mathbf{w}^T \mathbf{z} + \mathbf{b} = 0$, which can correctly discriminate \mathbf{x} from \mathbf{y} .

$\mathbf{w}^T \mathbf{z} + \mathbf{b} < 0$, \mathbf{z} is classified as \mathbf{y} ;

$\mathbf{w}^T \mathbf{z} + \mathbf{b} > 0$, \mathbf{z} is classified as \mathbf{x} .

$\mathbf{w}^T \mathbf{z} + \mathbf{b} = 0$: decision hyperplane



Background - Global Learning (I)

■ Global learning

- **Basic idea**: Focusing on summarizing data usually by estimating a distribution
- **Example**
 - 1) Assume Gaussinity for the data
 - 2) Learn the parameters via MLE or other criteria
 - 3) Exploit Bayes theory to find the optimal thresholding for classification

Traditional Bayes Optimal Classifier



Background - Global Learning (II)

■ Problems

- I Usually have to assume specific models on data, which may NOT always coincide with data

“all models are wrong but some are useful...”—by George Box

- II Estimating distributions may be wasteful and imprecise

*Finding the ideal generator of the data, i.e., the distribution, is only an intermediate goal in many settings, e.g., in classification or regression. **Optimizing an intermediate objective may be inefficient or wasteful.***



Background- Local Learning (I)

■ Local learning

- **Basic idea**: Focus on exploiting part of information, which is directly related to the objective, e.g., the classification accuracy instead of describing data in a holistic way

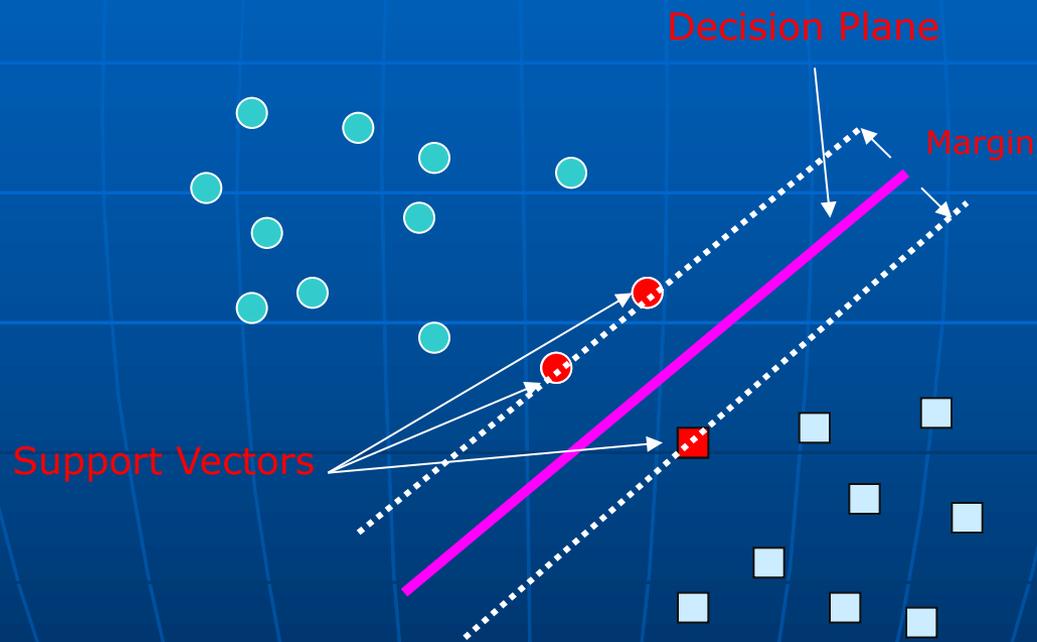
- **Example**

In classification, we need to accurately model the data around the (possible) separating plane, while inaccurate modeling of other parts is certainly acceptable (as is done in SVM).



Background - Local Learning (II)

- Support Vector Machine (SVM)
 - The current state-of-the-art classifier



Background - Local Learning (III)

■ Problems

- ③ The fact that the objective is exclusively determined by local information may lose the overall view of data



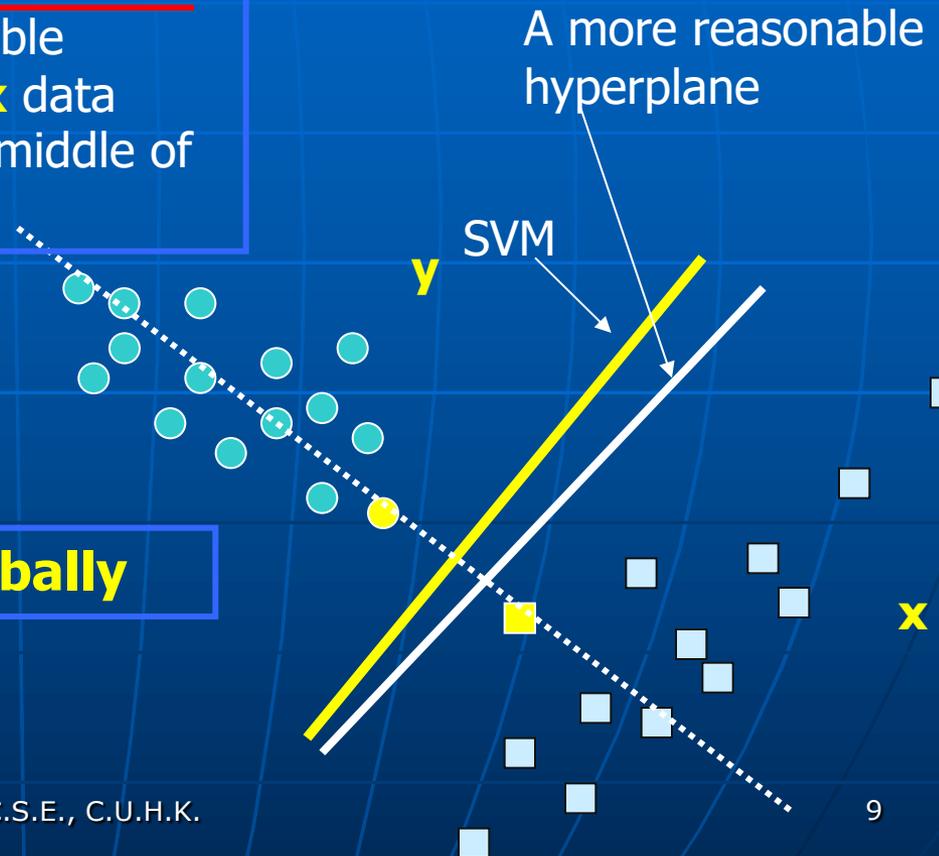
Background- Local Learning (IV)

An illustrative example

Along the dashed axis, the y data is obviously more likely to scatter than the x data. Therefore, a more reasonable hyperplane may lie closer to the x data rather than locating itself in the middle of two classes as in SVM.



Learning Locally and Globally



Learning Locally and Globally

- **Basic idea:** Focus on using **both local information** and **certain robust global information**
 - Do not try to estimate the distribution as in global learning, which may be inaccurate and indirect
 - Consider robust global information for providing a roadmap for local learning



Summary of Background

Global Learning

Problem  *Optimizing an intermediate objective*

Can we directly optimize the objective??

Local Learning, e.g., SVM 

Problem  *Assume specific models*

Without specific model assumption?

*Distribution-free Bayes optimal classifier ---
Minimum Error Minimax Probability Machine (MEMPM)*

Local Learning

Problem  *Focusing on local info may lose the roadmap of data*

Can we learn both globally and locally??

Maxi-Min Margin Machine (M⁴)

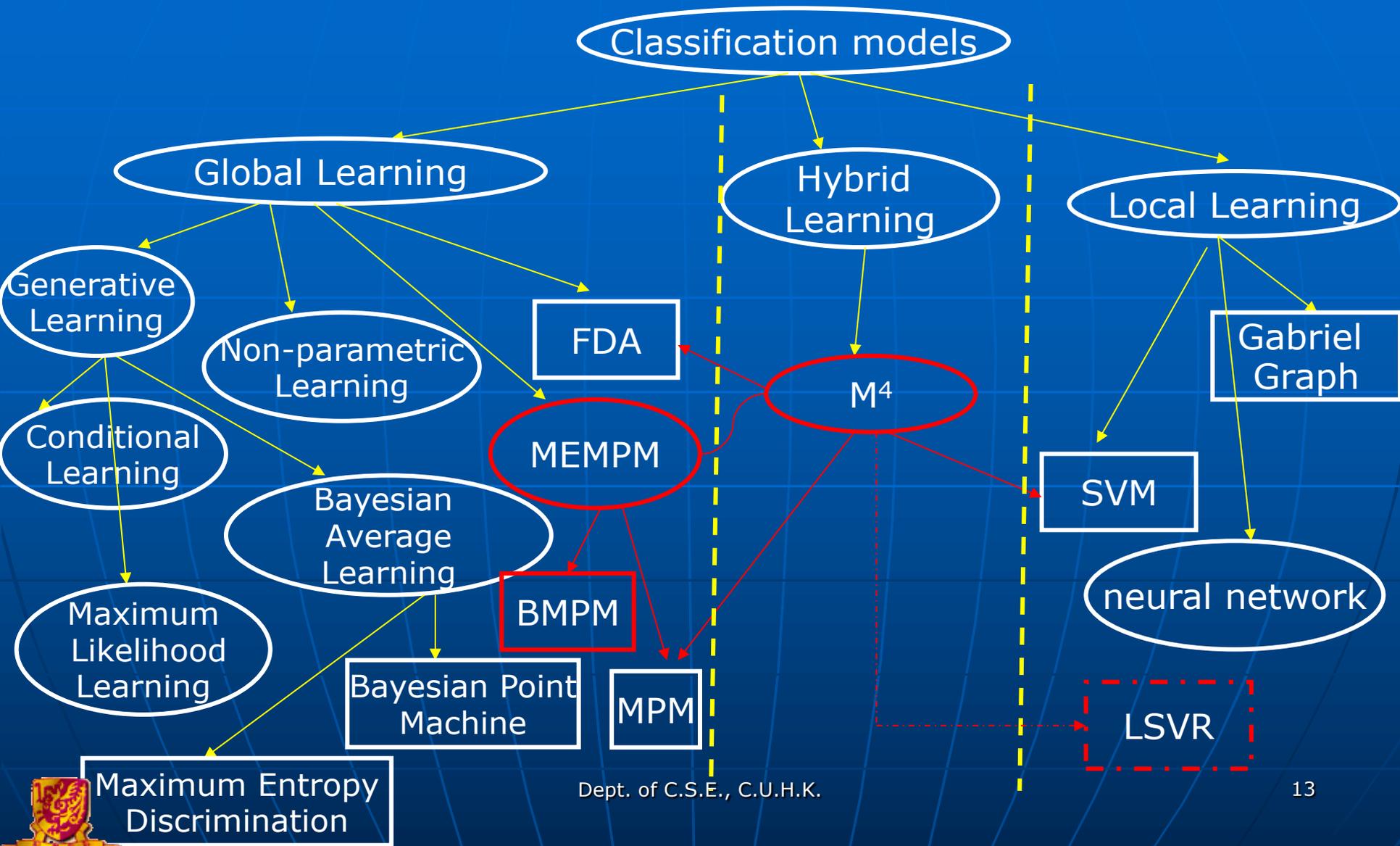


Contributions

- **Minimum Error Minimax Probability Machine**
(Accepted by JMLR 04)
 - A **worst-case distribution-free Bayes Optimal Classifier**
 - Containing Minimax Probability Machine (MPM) and **Biased Minimax Probability Machine (BMPPM)**(AMAI04,CVPR04) as special cases
- **Maxi-Min Margin Machine (M^4)** (ICML 04+Submitted)
 - A **unified framework** that learns locally and globally
 - Support Vector Machine (SVM)
 - Minimax Probability Machine (MPM)
 - Fisher Discriminant Analysis (FDA)
 - Can be linked with MEMPM
 - Can be extended into regression: **Local Support Vector Regression (LSVR)** (submitted)



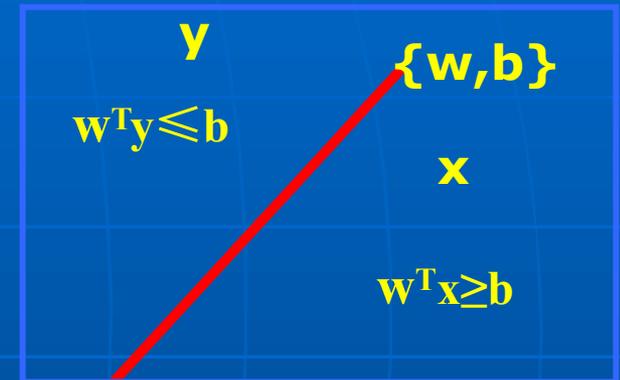
Hierarchy Graph of Related Models



Minimum Error Minimax Probability Machine (MEMPM)

Model Definition:

$$\begin{aligned} & \max_{\alpha, \beta, \mathbf{w} \neq \mathbf{0}, b} \quad \theta \alpha + (1 - \theta) \beta, \quad \text{s.t.} \\ & \inf_{\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma_{\mathbf{x}})} \quad \Pr \{ \mathbf{w}^T \mathbf{x} \geq b \} \geq \alpha, \\ & \inf_{\mathbf{y} \sim (\bar{\mathbf{y}}, \Sigma_{\mathbf{y}})} \quad \Pr \{ \mathbf{w}^T \mathbf{y} \leq b \} \geq \beta. \end{aligned}$$



- θ : prior probability of class x ; $a(\beta)$: represents the worst-case accuracy for class x (y)

$\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma_{\mathbf{x}})$: The class of distributions that have prescribed mean $\bar{\mathbf{x}}$ and covariance $\Sigma_{\mathbf{x}}$

$\mathbf{y} \sim (\bar{\mathbf{y}}, \Sigma_{\mathbf{y}})$: likewise



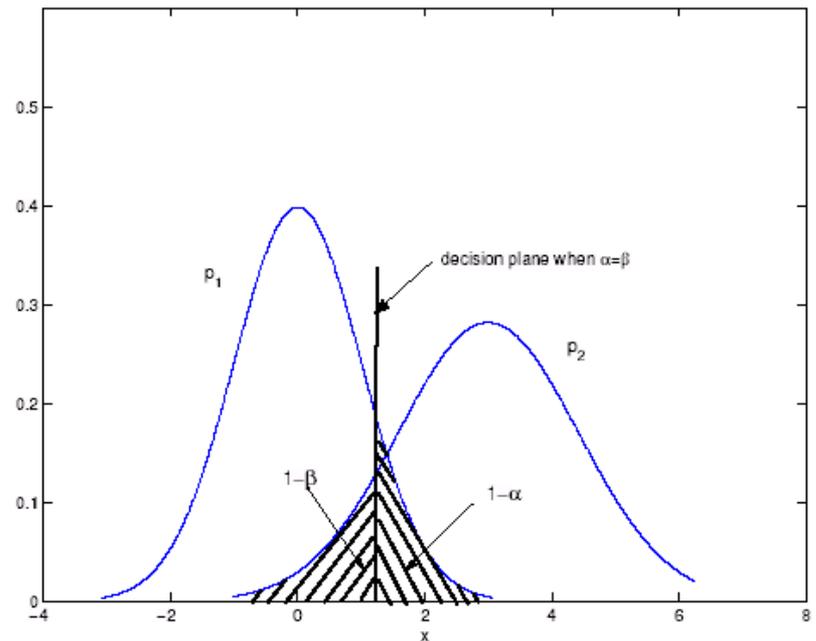
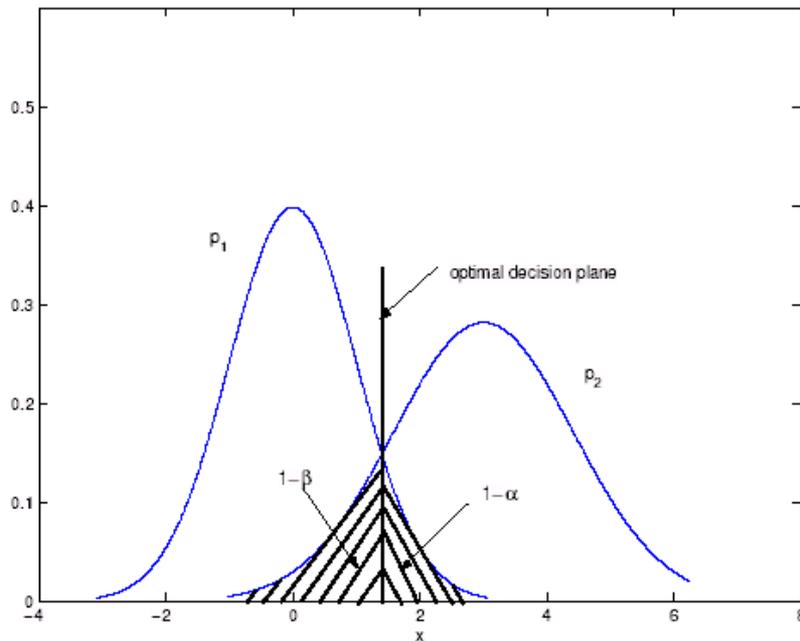
MEMPM: Model Comparison

MEMPM (JMLR04)

$$\begin{aligned} & \max_{\alpha, \beta, w \neq 0, b} \theta \alpha + (1 - \theta) \beta, \quad \text{s.t.} \\ & \inf_{x \sim (\bar{x}, \Sigma_x)} \Pr\{w^T x \geq b\} \geq \alpha, \\ & \inf_{y \sim (\bar{y}, \Sigma_y)} \Pr\{w^T y \leq b\} \geq \beta. \end{aligned}$$

MPM (Lanckriet et al. JMLR 2002)

$$\begin{aligned} & \max_{\alpha, w \neq 0, b} \alpha \quad \text{s.t.} \\ & \inf_{x \sim (\bar{x}, \Sigma_x)} \Pr\{w^T x \geq b\} \geq \alpha \\ & \inf_{y \sim (\bar{y}, \Sigma_y)} \Pr\{w^T y \leq b\} \geq \alpha \end{aligned}$$



MEMPM: Advantages

- A **distribution-free** Bayes optimal Classifier in the worst-case scenario
- Containing an **explicit accuracy bound**, namely, $\theta\alpha + (1-\theta)\beta$
- Subsuming a special case *Biased Minimax Probability Machine* **for biased classification**



MEMPPM: Biased MPPM

Biased Classification:

Diagnosis of epidemical disease: Classifying a patient who is infected with a disease into an opposite class results in more serious consequence than the other way around.

The classification accuracy should **be biased towards the class with disease.**

BMPPM

$$\begin{aligned} & \max_{\alpha, \beta, \mathbf{w} \neq \mathbf{0}, b} \alpha, \quad \text{s.t.} \\ & \inf_{\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma_x)} \Pr\{\mathbf{w}^T \mathbf{x} \geq b\} \geq \alpha, \\ & \inf_{\mathbf{y} \sim (\bar{\mathbf{y}}, \Sigma_y)} \Pr\{\mathbf{w}^T \mathbf{y} \leq b\} \geq \beta, \\ & \beta \geq \gamma \end{aligned}$$

An ideal model for biased classification.

A typical setting: We should **maximize the accuracy** for the important class as long as the accuracy for **the less important class is acceptable** (greater than an acceptable level γ).



MEMPM: Biased MPM (I)

- Objective

$$\begin{aligned} \max_{\alpha, \beta, \mathbf{a} \neq \mathbf{0}, b} \quad & \alpha \quad \text{s.t.} \quad \inf_{\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma_x)} \Pr\{\mathbf{a}^T \mathbf{x} \geq b\} \geq \alpha \\ & \inf_{\mathbf{y} \sim (\bar{\mathbf{y}}, \Sigma_y)} \Pr\{\mathbf{a}^T \mathbf{y} \leq b\} \geq \beta \\ & \beta \geq \gamma \end{aligned}$$

- Equivalently

$$\begin{aligned} \max_{\alpha, \beta, \mathbf{a} \neq \mathbf{0}, b} \quad & \kappa(\alpha) \quad \text{s.t.} \quad 1 = \kappa(\alpha) \sqrt{\mathbf{a}^T \Sigma_x \mathbf{a}} + \kappa(\beta) \sqrt{\mathbf{a}^T \Sigma_y \mathbf{a}} \\ & \mathbf{a}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1 \\ & \kappa(\beta) \geq \kappa(\gamma) \\ & \kappa(\alpha) = \sqrt{\frac{\alpha}{1-\alpha}}, \quad \kappa(\beta) = \sqrt{\frac{\beta}{1-\beta}} \end{aligned}$$



MEMPM: Biased MPM (II)

Objective

$$\max_{\alpha, \beta, \mathbf{w} \neq 0, b} \kappa(\alpha) \quad \text{s.t.} \quad 1 = \kappa(\alpha) \sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}} + \kappa(\beta) \sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}}$$

$$\mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1$$

$$\kappa(\beta) \geq \kappa(\gamma)$$

Equivalently,

$$\max_{\kappa(\beta), \mathbf{w} \neq 0} \frac{1 - \kappa(\beta) \sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}}}{\sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}}} \quad \text{s.t.} \quad \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1, \quad \kappa(\beta) \geq \kappa(\gamma)$$

Equivalently,

$$\max_{\mathbf{w} \neq 0} \frac{1 - \kappa(\gamma) \sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}}}{\sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}}} \quad \text{s.t.} \quad \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1$$

Conave-Convex Fractional Programming problem

1. Each local optimum is the **global optimum**
2. Can be solved in **$O(n^3 + Nn^2)$**

N : number of data points **n** : Dimension



MEMPM: Optimization (I)

- Objective

$$\begin{aligned} \max_{\alpha, \beta, \mathbf{w} \neq \mathbf{0}} \quad & \theta\alpha + (1-\theta)\beta, \quad \text{s.t.} \\ & 1 = \kappa(\alpha)\sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}} + \kappa(\beta)\sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}} \\ & \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1 \end{aligned}$$

- $$\begin{aligned} \min_{\kappa(\alpha), \kappa(\beta), \mathbf{w} \neq \mathbf{0}} \quad & \frac{\theta}{\kappa(\alpha)^2 + 1} + \frac{1-\theta}{\kappa(\beta)^2 + 1}, \quad \text{s.t.} \\ & 1 = \kappa(\alpha)\sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}} + \kappa(\beta)\sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}} \\ & \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1 \end{aligned}$$



MEMPM: Optimization (II)

- Objective

$$\max_{\alpha, \beta, \mathbf{w} \neq \mathbf{0}} \frac{\theta \kappa(\alpha)^2}{\kappa(\alpha)^2 + 1} + (1 - \theta) \beta, \quad \text{s.t.}$$

$$\mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1,$$

$$\text{where } \kappa(\alpha) = \frac{1 - \kappa(\beta) \sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}}}{\sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}}}$$

- Line search + BMPM method



MEMPM: Problems

- As a global learning approach, the decision plane is **exclusively** dependent on **global information**, i.e., up to second-order moments.
- These moments may **NOT be accurately** estimated! –We may need local information to neutralize the negative effect caused.



Learning Locally and Globally

Dept. of C.S.E., C.U.H.K.



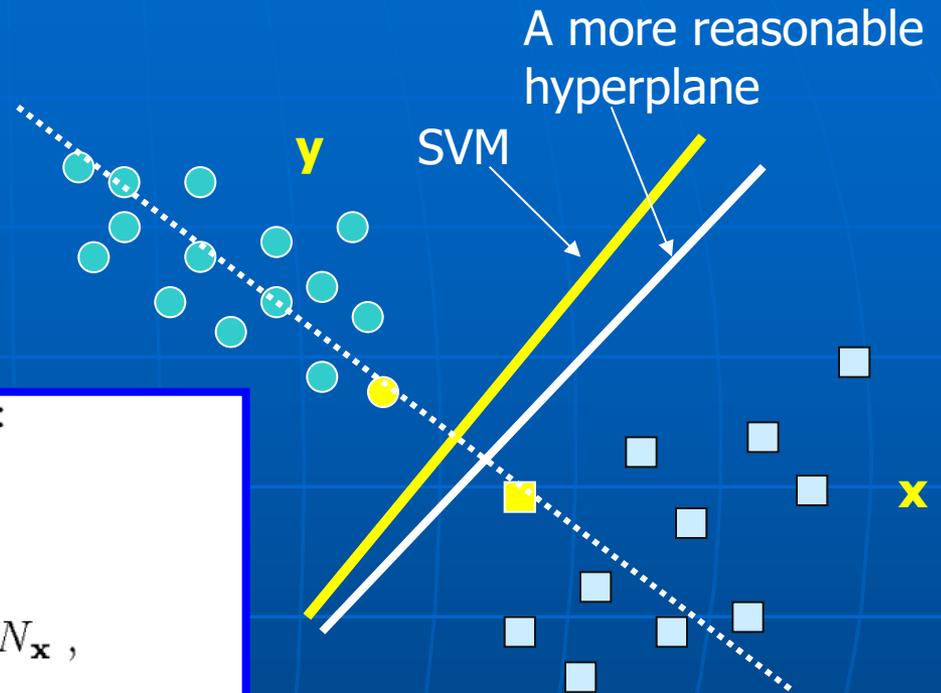
Learning Locally and Globally: Maxi-Min Margin Machine (M⁴)

Model Definition

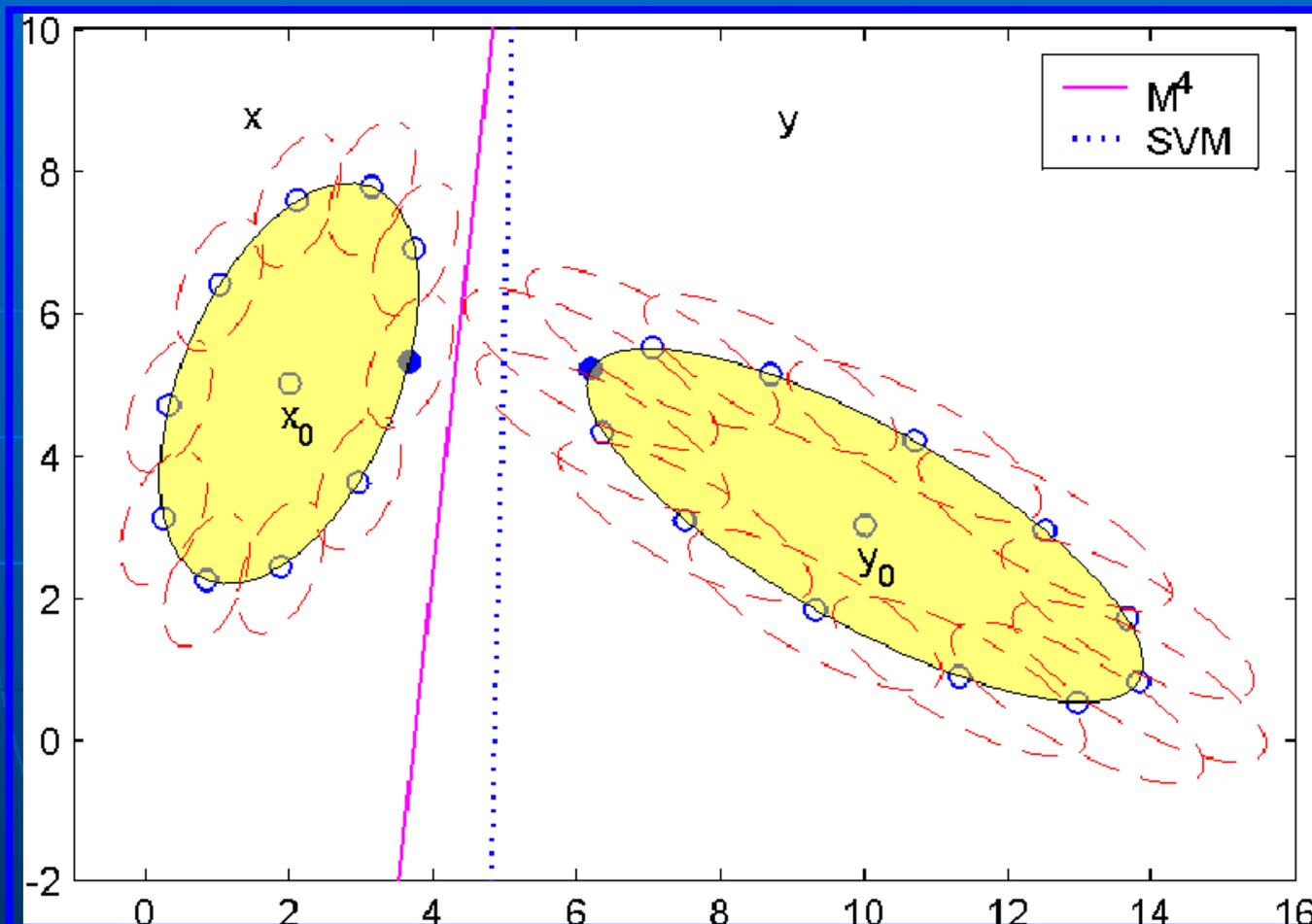
The formulation for M⁴ can be written as:

$$\begin{aligned} \max_{\rho, \mathbf{w} \neq \mathbf{0}, b} \quad & \rho \quad s.t. \\ \frac{(\mathbf{w}^T \mathbf{x}_i + b)}{\sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}}} & \geq \rho, \quad i = 1, 2, \dots, N_{\mathbf{x}}, \\ \frac{-(\mathbf{w}^T \mathbf{y}_j + b)}{\sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}} & \geq \rho, \quad j = 1, 2, \dots, N_{\mathbf{y}}, \end{aligned}$$

where $\Sigma_{\mathbf{x}}$ and $\Sigma_{\mathbf{y}}$ refer to the covariance matrices of the \mathbf{x} and the \mathbf{y} data, respectively.



M⁴: Geometric Interpretation



M⁴: Solving Method (I)

Divide and Conquer:

If we fix ρ to a specific ρ_n , the problem changes to check whether this ρ_n satisfies the following constraints:

$$\begin{aligned}(\mathbf{w}^T \mathbf{x}_i + b) &\geq \rho_n \sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}}, \quad i = 1, \dots, N_{\mathbf{x}}, \\ -(\mathbf{w}^T \mathbf{y}_j + b) &\geq \rho_n \sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}, \quad j = 1, \dots, N_{\mathbf{y}}.\end{aligned}$$

If yes, we increase ρ_n ; otherwise, we decrease it.



Second Order Cone Programming Problem!!!



M⁴: Solving Method (II)

Iterate the following two **Divide and Conquer** steps:

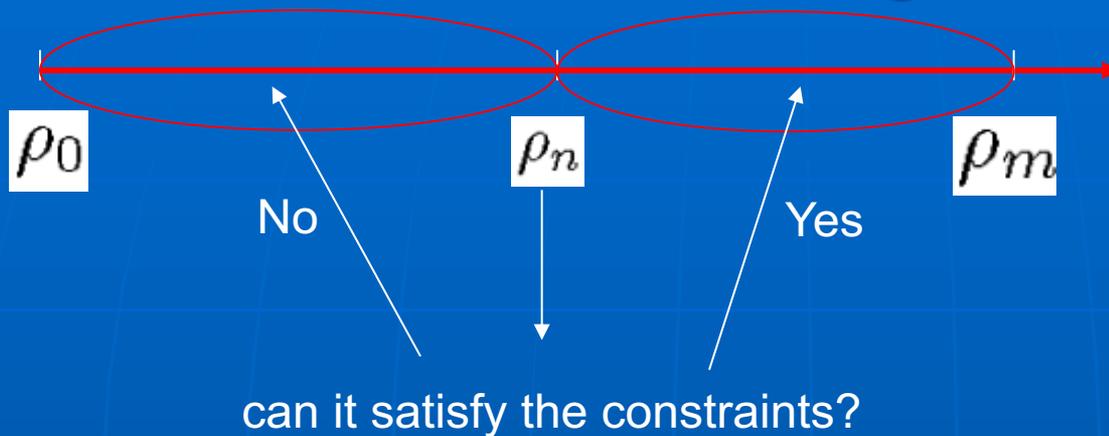
1. **Divide:** Set $\rho_n = (\rho_0 + \rho_m)/2$, where ρ_0 is a feasible ρ , ρ_m is an infeasible ρ , and $\rho_0 \leq \rho_m$.
2. **Conquer:** Call the Modified Second Order Cone Programming (MSOCP) procedure elaborated in the following to check whether ρ_n is a feasible ρ . If yes, set $\rho_0 = \rho_n$; otherwise, set $\rho_m = \rho_n$;



Sequential Second Order Cone Programming Problem!!!



M⁴: Solving Method (III)



- The worst-case iteration number is $\log(L/\varepsilon)$
 L : $\rho_{max} - \rho_{min}$ (search range)
 ε : The required precision
- Each iteration is a Second Order Cone Programming problem yielding $O(n^3)$
- Cost of forming the constraint matrix $O(N n^3)$
Total time complexity = $O(\log(L/\varepsilon) n^3 + N n^3) \approx O(N n^3)$
 N : number of data points n : Dimension



M⁴: Links with MPM (I)

$$\frac{(\mathbf{w}^T \mathbf{x}_i + b)}{\sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}}} \geq \rho, \quad i = 1, 2, \dots, N_{\mathbf{x}},$$

$$\frac{-(\mathbf{w}^T \mathbf{y}_j + b)}{\sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}} \geq \rho, \quad j = 1, 2, \dots, N_{\mathbf{y}},$$

Span all the data points and add them together

$$\begin{aligned} \mathbf{w}^T \sum_{i=1}^{N_{\mathbf{x}}} \mathbf{x}_i + N_{\mathbf{x}} b &\geq N_{\mathbf{x}} \rho \sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}}, \\ \Leftrightarrow \mathbf{w}^T \bar{\mathbf{x}} + b &\geq \rho \sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}}, \end{aligned}$$

$$\begin{aligned} -(\mathbf{w}^T \sum_{j=1}^{N_{\mathbf{y}}} \mathbf{y}_j + N_{\mathbf{y}} b) &\geq N_{\mathbf{y}} \rho \sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}, \\ \Leftrightarrow -(\mathbf{w}^T \bar{\mathbf{y}} + b) &\geq \rho \sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}, \end{aligned}$$

$$+$$

$$\begin{aligned} \max_{\rho, \mathbf{w} \neq \mathbf{0}, b} \quad &\rho \quad s.t. \\ \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) &\geq \rho (\sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}} + \sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}). \quad (15) \end{aligned}$$

Exactly MPM Optimization Problem!!!



M⁴: Links with MPM (II)



Remarks:

- The procedure is not reversible: MPM is a **special case of M⁴**
- MPM focuses on building decision boundary **GLOBALLY**, i.e., it exclusively depends on the means and covariances.
- However, means and covariances may not be accurately estimated.

The formulation for M⁴ can be written as:

$$\begin{aligned} \max_{\rho, \mathbf{w} \neq \mathbf{0}, b} \quad & \rho \quad s.t. \\ \frac{(\mathbf{w}^T \mathbf{x}_i + b)}{\sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}}} & \geq \rho, \quad i = 1, 2, \dots, N_{\mathbf{x}}, \\ \frac{-(\mathbf{w}^T \mathbf{y}_j + b)}{\sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}} & \geq \rho, \quad j = 1, 2, \dots, N_{\mathbf{y}}, \end{aligned}$$

where $\Sigma_{\mathbf{x}}$ and $\Sigma_{\mathbf{y}}$ refer to the covariance matrices of the \mathbf{x} and the \mathbf{y} data, respectively.



$$\begin{aligned} \max_{\rho, \mathbf{w} \neq \mathbf{0}, b} \quad & \rho \quad s.t. \\ \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) & \geq \rho (\sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}} + \sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}) \end{aligned}$$



M⁴: Links with SVM (I)

1

If one assumes $\Sigma_{\mathbf{x}} = \Sigma_{\mathbf{y}} = \Sigma$,

4

If one assumes $\Sigma = \mathbf{I}$

2

$$\begin{aligned} \max_{\rho, \mathbf{w} \neq \mathbf{0}, b} \quad & \rho \quad s.t. \\ & (\mathbf{w}^T \mathbf{x}_i + b) \geq \rho \sqrt{\mathbf{w}^T \Sigma \mathbf{w}}, \\ & -(\mathbf{w}^T \mathbf{y}_j + b) \geq \rho \sqrt{\mathbf{w}^T \Sigma \mathbf{w}}, \end{aligned}$$

where $i = 1, \dots, N_{\mathbf{x}}$ and $j = 1, \dots, N_{\mathbf{y}}$.

Support Vector Machines

3

$$\begin{aligned} \min_{\mathbf{w} \neq \mathbf{0}, b} \quad & \mathbf{w}^T \Sigma \mathbf{w} \quad s.t. \\ & (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \\ & -(\mathbf{w}^T \mathbf{y}_j + b) \geq 1, \end{aligned}$$

where $i = 1, \dots, N_{\mathbf{x}}$ and $j = 1, \dots, N_{\mathbf{y}}$.

The magnitude of w can scale up without influencing the optimization. Assume

$$\rho(\mathbf{w}^T \Sigma \mathbf{w})^{0.5} = 1$$

SVM is the special case of M⁴

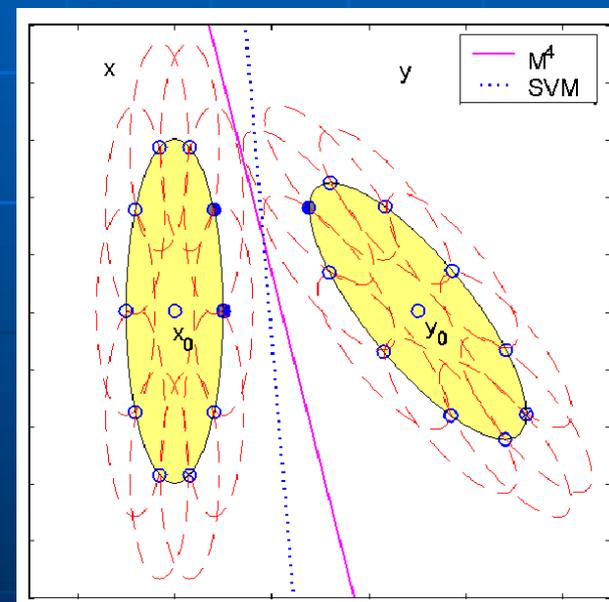
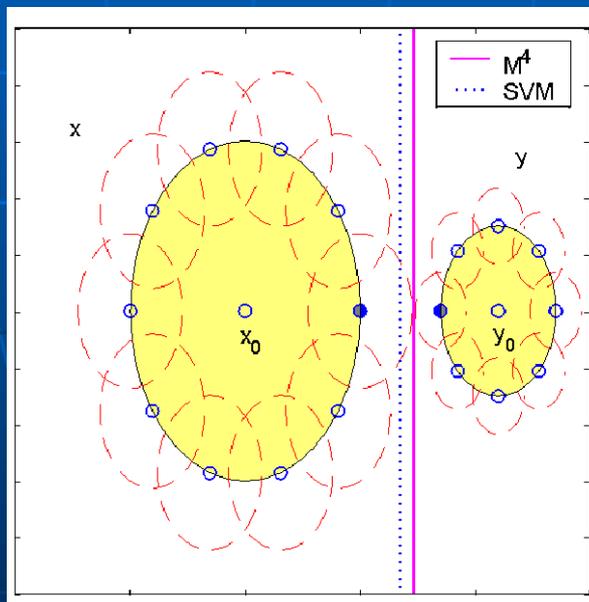


M⁴: Links with SVM (II)

Assumption 1 → If one assumes $\Sigma_x = \Sigma_y = \Sigma$,

Assumption 2 → If one assumes $\Sigma=I$

These two assumptions of SVM are **inappropriate**



M⁴: Links with FDA (I)

If one assumes
 $\Sigma_x = \Sigma_y = (\Sigma_y^* + \Sigma_x^*)/2$

$$\begin{aligned} \max_{\rho, \mathbf{w} \neq \mathbf{0}, b} \quad & \rho \quad s.t. \\ \frac{(\mathbf{w}^T \mathbf{x}_i + b)}{\sqrt{\mathbf{w}^T \Sigma_x \mathbf{w} + \mathbf{w}^T \Sigma_y \mathbf{w}}} & \geq \rho, \\ \frac{-(\mathbf{w}^T \mathbf{y}_j + b)}{\sqrt{\mathbf{w}^T \Sigma_x \mathbf{w} + \mathbf{w}^T \Sigma_y \mathbf{w}}} & \geq \rho, \end{aligned}$$

Perform a procedure similar to
MPM...

$$\begin{aligned} \max_{\rho, \mathbf{w} \neq \mathbf{0}, b} \quad & \rho \quad s.t. \\ \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) & \geq \rho \sqrt{\mathbf{w}^T \Sigma_x \mathbf{w} + \mathbf{w}^T \Sigma_y \mathbf{w}}. \end{aligned}$$

FDA



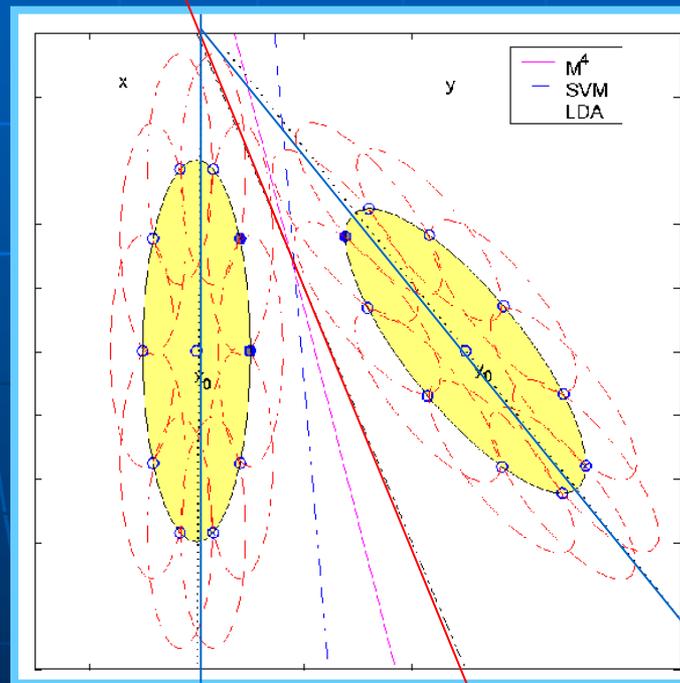
M⁴: Links with FDA (II)

Assumption



$$\Sigma_x = \Sigma_y = (\Sigma^* y + \Sigma^* x) / 2$$

Still inappropriate



M⁴: Links with MEMPM

M⁴ (a globalized version)

$$\begin{aligned} \max_{\mathbf{w} \neq \mathbf{0}, b} \quad & \theta t + (1 - \theta)s \quad \text{s.t.} \\ & \frac{\mathbf{w}^T \bar{\mathbf{x}} + b}{\sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}}} \geq t, \\ & -\frac{\mathbf{w}^T \bar{\mathbf{y}} + b}{\sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}}} \geq s, \end{aligned}$$

MEMPM

$$\begin{aligned} \max_{\alpha, \beta, \mathbf{a} \neq \mathbf{0}, b} \quad & \theta \alpha + (1 - \theta) \beta \quad \text{s.t.} \\ & \inf_{\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma_x)} \Pr\{\mathbf{a}^T \mathbf{x} \geq b\} \geq \alpha, \\ & \inf_{\mathbf{y} \sim (\bar{\mathbf{y}}, \Sigma_y)} \Pr\{\mathbf{a}^T \mathbf{y} \leq b\} \geq \beta. \end{aligned}$$

$$\begin{aligned} \max_{\alpha, \beta, \mathbf{a} \neq \mathbf{0}, b} \quad & \theta \alpha + (1 - \theta) \beta \quad \text{s.t.} \\ & -b + \mathbf{a}^T \bar{\mathbf{x}} \geq \kappa(\alpha) \sqrt{\mathbf{a}^T \Sigma_x \mathbf{a}}, \\ & b - \mathbf{a}^T \bar{\mathbf{y}} \geq \kappa(\beta) \sqrt{\mathbf{a}^T \Sigma_y \mathbf{a}}. \end{aligned}$$

T and ***s*** \longleftrightarrow ***K*(α)** and ***K*(β)** :

The margin from the mean to the decision plane

The globalized M⁴ **maximizes the weighted margin**, while MEMPM **Maximizes the weighted worst-case accuracy**.



M⁴ : Nonseparable Case

Introducing slack variables

$$\begin{aligned} \max_{\rho, \mathbf{w} \neq \mathbf{0}, b, \boldsymbol{\xi}} \quad & \rho - C \sum_{k=1}^{N_x + N_y} \xi_k \quad s.t. \\ (\mathbf{w}^T \mathbf{x}_i + b) & \geq \rho \sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}} - \xi_i, \\ -(\mathbf{w}^T \mathbf{y}_j + b) & \geq \rho \sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}} - \xi_{j+N_x}, \\ \xi_k & \geq 0, \end{aligned}$$

How to solve?? **Line Search+Second Order Cone Programming**

Step 1. Generate a new ρ_n from three previous ρ_1, ρ_2, ρ_3 by using the Quadratic Interpolation method.

Step 2. Fix $\rho = \rho_n$, perform the optimization based on SOCP algorithms.
Update ρ_1, ρ_2, ρ_3 .



M⁴ : Extended into Regression---

Local Support Vector Regression (LSVR)

Regression: Find a function $f(x) = \mathbf{w}^T \mathbf{x} + b$, $\mathbf{w}, \mathbf{x} \in \mathbb{R}^d, b \in \mathbb{R}$. to approximate the data $\{(x_i, y_i) \mid x_i \in \mathbb{R}^d, y_i \in \mathbb{R}, i = 1, \dots, N\}$

LSVR Model Definition

$$\begin{aligned} \min_{\mathbf{w}, b, \xi_i, \xi_i^*} \quad & \frac{1}{N} \sum_{i=1}^N \sqrt{\mathbf{w}^T \Sigma_i \mathbf{w}} + C \sum_{i=1}^N (\xi_i + \xi_i^*), \\ \text{s.t.} \quad & y_i - (\mathbf{w}^T \mathbf{x}_i + b) \leq \epsilon \sqrt{\mathbf{w}^T \Sigma_i \mathbf{w}} + \xi_i, \\ & (\mathbf{w}^T \mathbf{x}_i + b) - y_i \leq \epsilon \sqrt{\mathbf{w}^T \Sigma_i \mathbf{w}} + \xi_i^*, \\ & \xi_i \geq 0, \quad \xi_i^* \geq 0, \quad i = 1, \dots, N, \end{aligned}$$

SVR Model Definition

$$\begin{aligned} \min_{\mathbf{w}, b, \xi_i, \xi_i^*} \quad & \|\mathbf{w}\| + C \sum_{i=1}^N (\xi_i + \xi_i^*), \\ \text{s.t.} \quad & y_i - (\mathbf{w}^T \mathbf{x}_i + b) \leq \epsilon + \xi_i, \\ & (\mathbf{w}^T \mathbf{x}_i + b) - y_i \leq \epsilon + \xi_i^*, \\ & \xi_i \geq 0, \quad \xi_i^* \geq 0, \quad i = 1, \dots, N. \end{aligned}$$

$$\because y_i = \mathbf{w}^T \mathbf{x}_i + b, \quad \bar{y}_i = \mathbf{w}^T \bar{\mathbf{x}}_i + b$$

$$\because \Delta_i = \frac{1}{2k+1} \sum_{j=-k}^k (y_{i+j} - \bar{y}_i)^2 = \frac{1}{2k+1} \sum_{j=-k}^k [\mathbf{w}^T (\mathbf{x}_{i+j} - \bar{\mathbf{x}}_i)]^2 = \mathbf{w}^T \Sigma_i \mathbf{w}$$



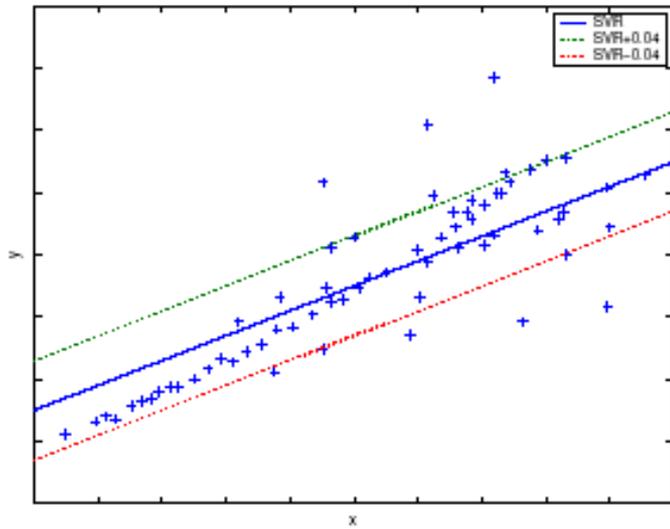
Local Support Vector Regression (LSVR)

- *When supposing $\Sigma_i = \mathbf{I}$ for each observation, LSVR is equivalent with l_1 -SVR under a mild assumption.*

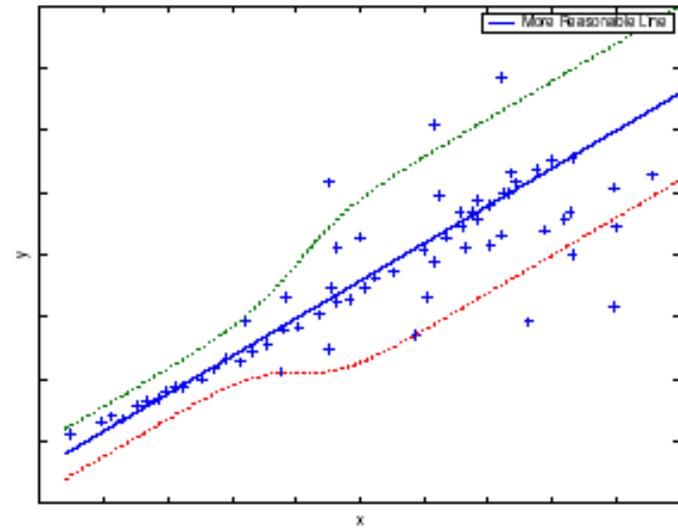
Lemma The LSVR model with setting $\Sigma_i = \mathbf{I}$ is equivalent to the l_1 -norm SVR in the sense that: (1) Assuming a unique ϵ_1^* exists for making l_1 -norm SVR optimal (i.e. setting ϵ to ϵ_1^* will make the objective function minimal), if for ϵ_1^* the l_1 -norm SVR achieves a solution $\{\mathbf{w}^*, b^*\} = \text{SVR}(\epsilon_1^*)$, then the LSVR can produce the same solution by setting the parameter $\epsilon = \frac{\epsilon_1^*}{\|\mathbf{w}_1^*\|}$, i.e., $\text{LSVR}(\frac{\epsilon_1^*}{\|\mathbf{w}_1^*\|}) = \text{SVR}(\epsilon_1^*)$; (2) Assuming a unique ϵ_2^* exists for making the special case of LSVR optimal (i.e. setting ϵ to ϵ_2^* will make the objective function minimal), if for ϵ_2^* the special case of LSVR achieves a solution $\{\mathbf{w}_2^*, b_2^*\} = \text{LSVR}(\epsilon_2^*)$, then the l_1 -norm SVR can produce the same solution by setting the parameter $\epsilon = \epsilon_2^* \|\mathbf{w}_2^*\|$, i.e., $\text{SVR}(\epsilon_2^* \|\mathbf{w}_2^*\|) = \text{LSVR}(\epsilon_2^*)$.



SVR vs. LSVR



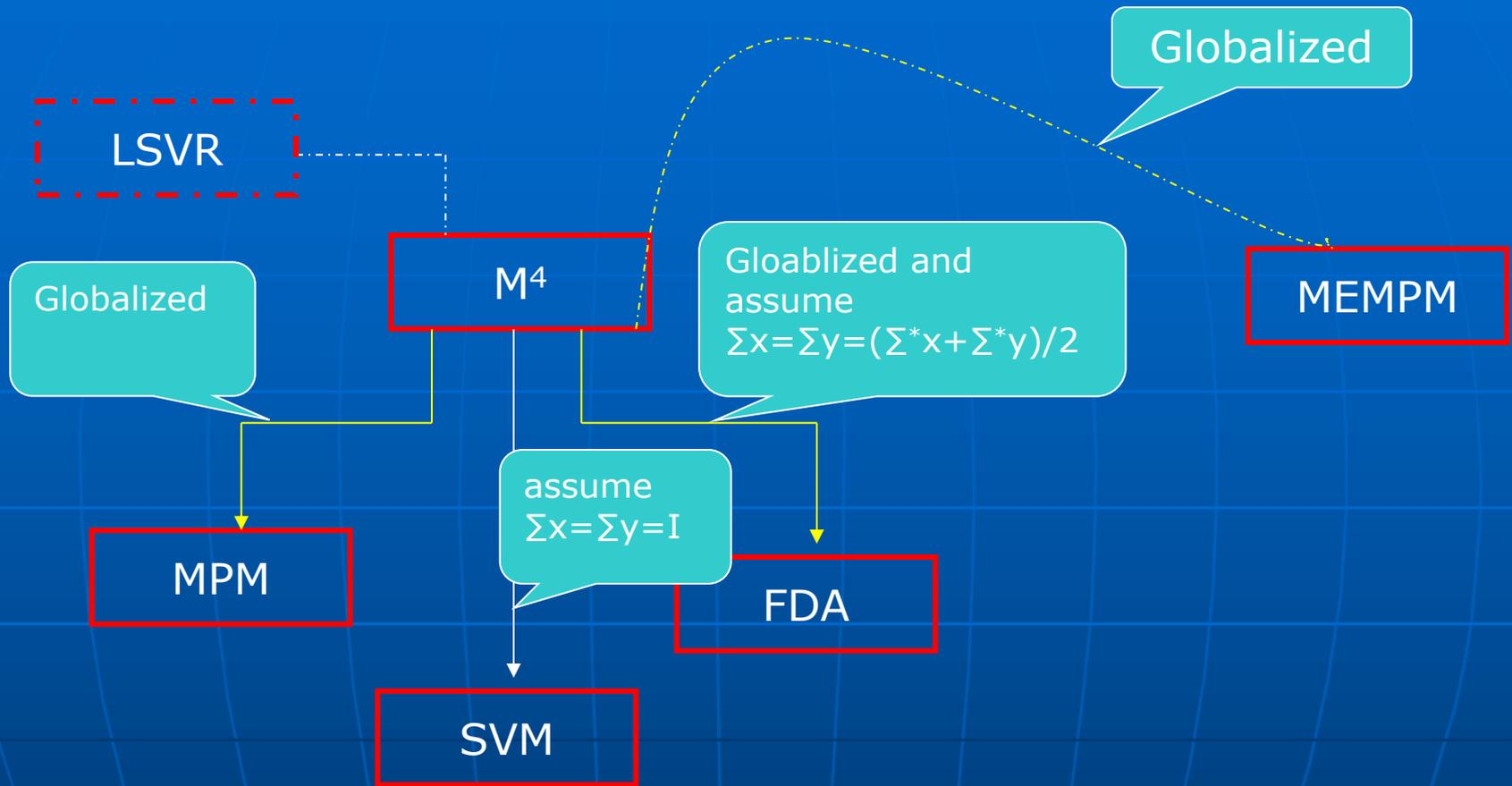
(a)



(b)



Short Summary



Non-linear Classifier : Kernelization (I)

- Previous discussions of MEMPM, BMPM, M^4 , and LSVR are conducted in the scope of linear classification.
- How about non-linear classification problems?



Using Kernelization techniques



Non-linear Classifier : Kernelization (II)

- In the next slides, we mainly discuss the kernelization on M^4 , while the proposed kernelization method is also applicable for MEMPM, BMPM, and LSVR.



Nonlinear Classifier: Kernelization (III)

- Map data to higher dimensional feature space \mathbf{R}^f

$$\mathbf{x}_i \rightarrow \varphi(\mathbf{x}_i)$$

$$\mathbf{y}_i \rightarrow \varphi(\mathbf{y}_i)$$

- Construct the linear decision plane $\mathbf{f}(\mathbf{y}, \mathbf{b}) = \mathbf{y}^T \mathbf{z} + \mathbf{b}$ in the feature space \mathbf{R}^f , with $\mathbf{y} \in \mathbf{R}^f$, $\mathbf{b} \in \mathbf{R}$

- In \mathbf{R}^f , we need to solve

$$\begin{aligned} & \max_{\rho, \gamma \neq 0, b} \quad \rho \quad s.t. \\ & \frac{(\gamma^T \varphi(\mathbf{x}_i) + b)}{\sqrt{\gamma^T \Sigma_{\varphi(\mathbf{x})} \gamma}} \geq \rho, \quad i = 1, 2, \dots, N_x, \\ & \frac{-(\gamma^T \varphi(\mathbf{y}_j) + b)}{\sqrt{\gamma^T \Sigma_{\varphi(\mathbf{y})} \gamma}} \geq \rho, \quad j = 1, 2, \dots, N_y. \end{aligned}$$

- However, we do not want to solve this in an explicit form of φ . Instead, we want to solve it in a kernelization form

$$\mathbf{K}(\mathbf{z}_1, \mathbf{z}_2) = \varphi(\mathbf{z}_1)^T \varphi(\mathbf{z}_2)$$



Nonlinear Classifier: Kernelization (IV)

Corollary If the estimates of means and covariance matrices are given in M^4 as the following plug-in estimates:

$$\begin{aligned}\overline{\varphi(\mathbf{x})} &= \frac{1}{N_x} \sum_{i=1}^{N_x} \varphi(\mathbf{x}_i), & \overline{\varphi(\mathbf{y})} &= \frac{1}{N_y} \sum_{j=1}^{N_y} \varphi(\mathbf{y}_j), \\ \Sigma_{\varphi(\mathbf{x})} &= \frac{1}{N_x} \sum_{i=1}^{N_x} (\varphi(\mathbf{x}_i) - \overline{\varphi(\mathbf{x})})(\varphi(\mathbf{x}_i) - \overline{\varphi(\mathbf{x})})^T, \\ \Sigma_{\varphi(\mathbf{y})} &= \frac{1}{N_y} \sum_{j=1}^{N_y} (\varphi(\mathbf{y}_j) - \overline{\varphi(\mathbf{y})})(\varphi(\mathbf{y}_j) - \overline{\varphi(\mathbf{y})})^T,\end{aligned}$$

then the optimal γ in (4.37-4.39) lies in the space spanned by the training points.

$$\gamma = \sum_{i=1}^{N_x} \mu_i \varphi(\mathbf{x}_i) + \sum_{j=1}^{N_y} \nu_j \varphi(\mathbf{y}_j),$$



Nonlinear Classifier: Kernelization (V)

Kernelization Theorem of M^4 *The optimal decision hyperplane for M^4 involves solving the following optimization problem:*

$$\begin{aligned} \max_{\rho, \boldsymbol{\eta} \neq \mathbf{0}, b} \quad & \rho \quad \text{s.t.} \\ \frac{(\boldsymbol{\eta}^T \mathbf{K}_i + b)}{\sqrt{\frac{1}{N_x} \boldsymbol{\eta}^T \tilde{\mathbf{K}}_x^T \tilde{\mathbf{K}}_x \boldsymbol{\eta}}} & \geq \rho, \quad i = 1, 2, \dots, N_x, \\ \frac{-(\boldsymbol{\eta}^T \mathbf{K}_{j+N_x} + b)}{\sqrt{\frac{1}{N_y} \boldsymbol{\eta}^T \tilde{\mathbf{K}}_y^T \tilde{\mathbf{K}}_y \boldsymbol{\eta}}} & \geq \rho, \quad j = 1, 2, \dots, N_y. \end{aligned}$$

Notation

$$\begin{aligned} \boldsymbol{\eta} &:= [\mu_1, \dots, \mu_{N_x}, v_1, \dots, v_{N_y}]^T \\ \tilde{\mathbf{k}}_x, \tilde{\mathbf{k}}_y &\in \mathbb{R}^{N_x + N_y} & [\tilde{\mathbf{k}}_x]_i &:= \frac{1}{N_x} \sum_{j=1}^{N_x} \mathbf{K}(\mathbf{x}_j, \mathbf{z}_i) \cdot \\ & & [\tilde{\mathbf{k}}_y]_i &:= \frac{1}{N_y} \sum_{j=1}^{N_y} \mathbf{K}(\mathbf{y}_j, \mathbf{z}_i) \cdot \\ \mathbf{1}_{N_x} &\in \mathbb{R}^{N_x} & \mathbf{1}_i &:= 1 \quad i = 1, 2, \dots, N_x \cdot \\ \mathbf{1}_{N_y} &\in \mathbb{R}^{N_y} & \mathbf{1}_i &:= 1 \quad i = 1, 2, \dots, N_y \cdot \\ \tilde{\mathbf{K}} &:= \begin{pmatrix} \tilde{\mathbf{K}}_x \\ \tilde{\mathbf{K}}_y \end{pmatrix} & &:= \begin{pmatrix} \mathbf{K}_x - \mathbf{1}_{N_x} \tilde{\mathbf{k}}_x^T \\ \mathbf{K}_y - \mathbf{1}_{N_y} \tilde{\mathbf{k}}_y^T \end{pmatrix} \cdot \end{aligned}$$



Experimental Results ---MEMPM (I)

Six benchmark data sets From UCI Repository

	Dataset	Attributes #	Instances#
1	Twonorm	20	7400
2	Breast	9	699
3	Ionosphere	34	351
4	Pima	8	768
5	Heart-disease	13	270
6	Vote	16	435



Platform: Windows 2000
Developing tool: Matlab 6.5

Evaluate both the linear and the Gaussian kernel with the wide parameter for Gaussian chosen by cross validations.



Experimental Results ---MEMPM(II)

Dataset	MEMPM				MPM	
	α	β	$\theta\alpha + (1 - \theta)\beta$	Accuracy	α	Accuracy
Twonorm(%)	80.3 ± 0.2%	79.9 ± 0.1%	80.1 ± 0.1%	97.9 ± 0.1%	80.1 ± 0.1%	97.9 ± 0.1%
Breast(%)	77.8 ± 0.8%	91.4 ± 0.5%	86.7 ± 0.5%	96.9 ± 0.3%	84.4 ± 0.5%	97.0 ± 0.2%
Ionosphere(%)	95.9 ± 1.2%	36.5 ± 2.6%	74.5 ± 0.8%	88.5 ± 1.0%	✓ 63.4 ± 1.1%	84.8 ± 0.8%
Pima(%)	0.9 ± 0.0%	62.9 ± 1.1%	41.3 ± 0.8%	76.8 ± 0.6%	✓ 32.0 ± 0.8%	76.1 ± 0.6%
Heart-disease(%)	43.6 ± 2.5%	66.5 ± 1.5%	56.3 ± 1.4%	84.2 ± 0.7%	✓ 54.9 ± 1.4%	83.2 ± 0.8%
Vote(%)	82.6 ± 1.3%	84.6 ± 0.7%	83.9 ± 0.9%	94.9 ± 0.4%	83.8 ± 0.9%	94.8 ± 0.4%

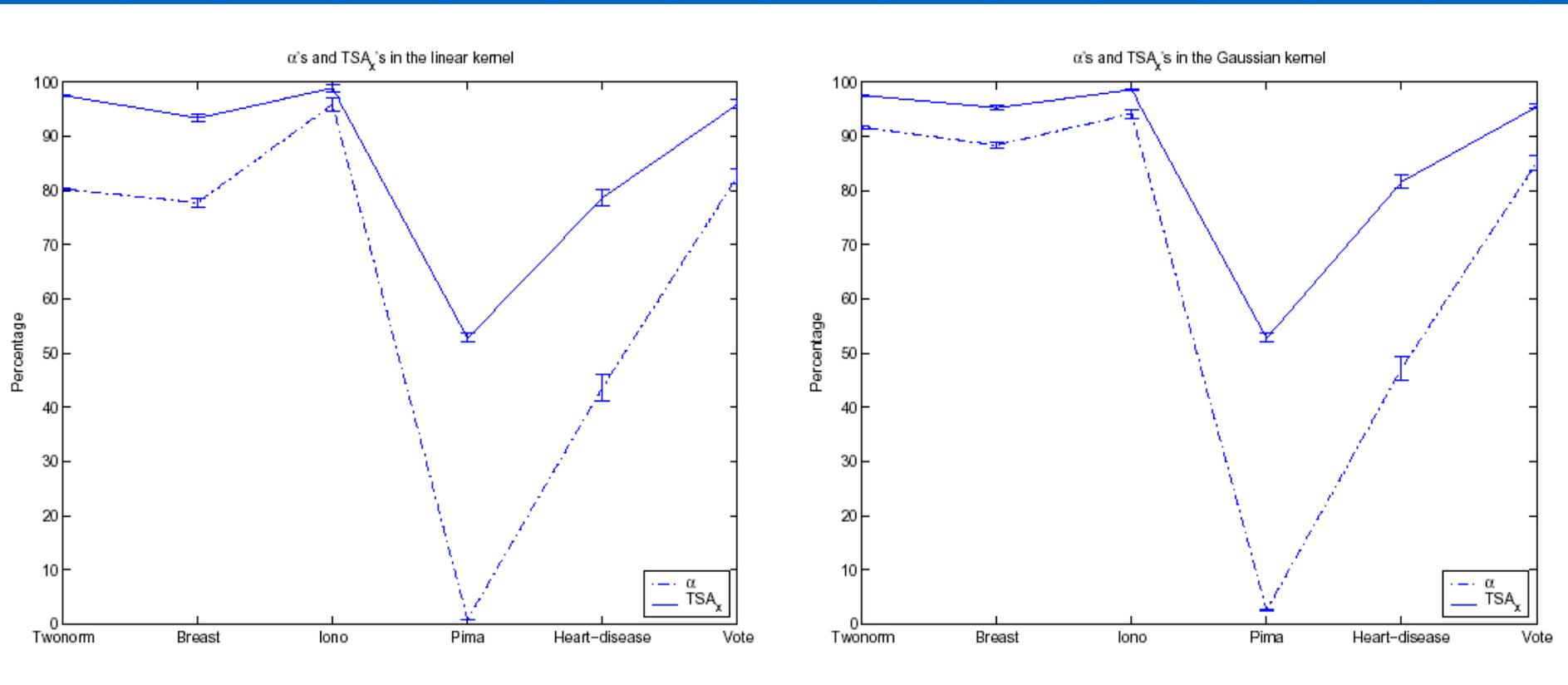
Table 2: Lower bound α , β , and test accuracy compared to MPM in the linear setting.

Dataset	MEMPM				MPM	
	α	β	$\theta\alpha + (1 - \theta)\beta$	Accuracy	α	Accuracy
Twonorm(%)	91.7 ± 0.2%	91.7 ± 0.2%	91.7 ± 0.2%	97.9 ± 0.1%	91.7 ± 0.2%	97.9 ± 0.1%
Breast(%)	88.4 ± 0.6%	90.7 ± 0.4%	89.9 ± 0.4%	96.9 ± 0.2%	89.9 ± 0.4%	96.9 ± 0.3%
Ionosphere(%)	94.2 ± 0.8%	80.9 ± 3.0%	89.4 ± 0.8%	93.8 ± 0.4%	✓ 89.0 ± 0.8%	92.2 ± 0.4%
Pima(%)	2.6 ± 0.1%	62.3 ± 1.6%	41.4 ± 1.1%	77.0 ± 0.7%	✓ 32.1 ± 1.0%	76.2 ± 0.6%
Heart-disease(%)	47.1 ± 2.2%	66.6 ± 1.4%	58.0 ± 1.5%	83.9 ± 0.9%	✓ 57.4 ± 1.6%	83.1 ± 1.0%
Vote(%)	85.1 ± 1.3%	84.3 ± 0.7%	84.7 ± 0.8%	94.7 ± 0.5%	84.4 ± 0.8%	94.6 ± 0.4%

Table 3: Lower bound α , β , and test accuracy compared to MPM with the Gaussian kernel.

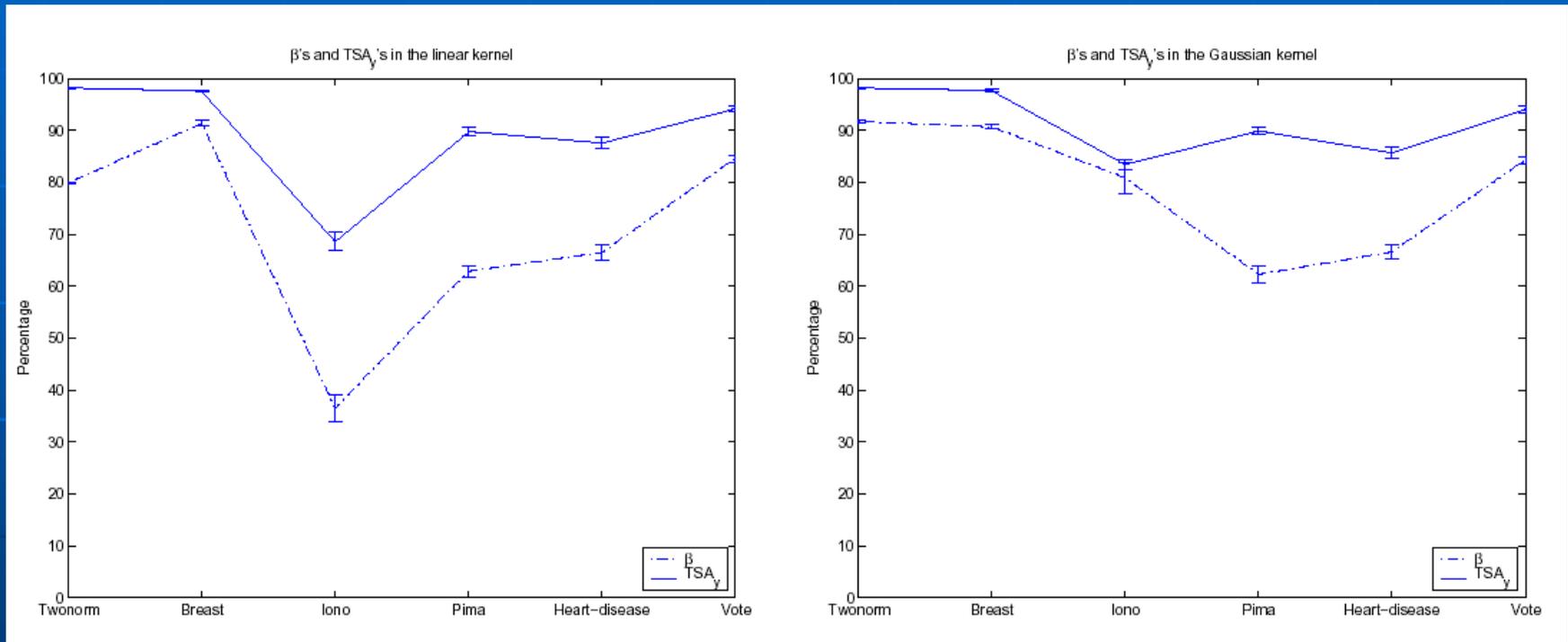
Experimental Results ---MEMPM (III)

α vs. The test-set accuracy for x (TSA_x)



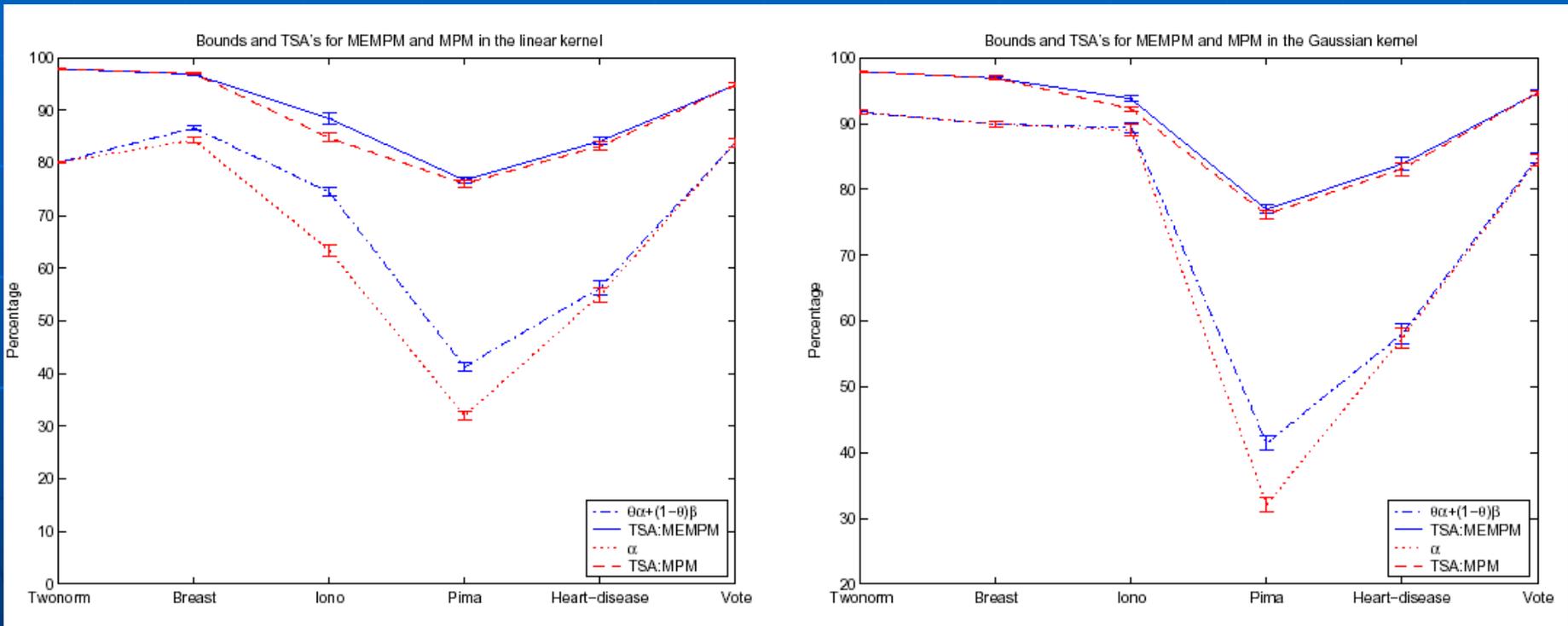
Experimental Results ---MEMPM (IV)

β vs. The test-set accuracy for y (TSA_y)



Experimental Results ---MEMPM (V)

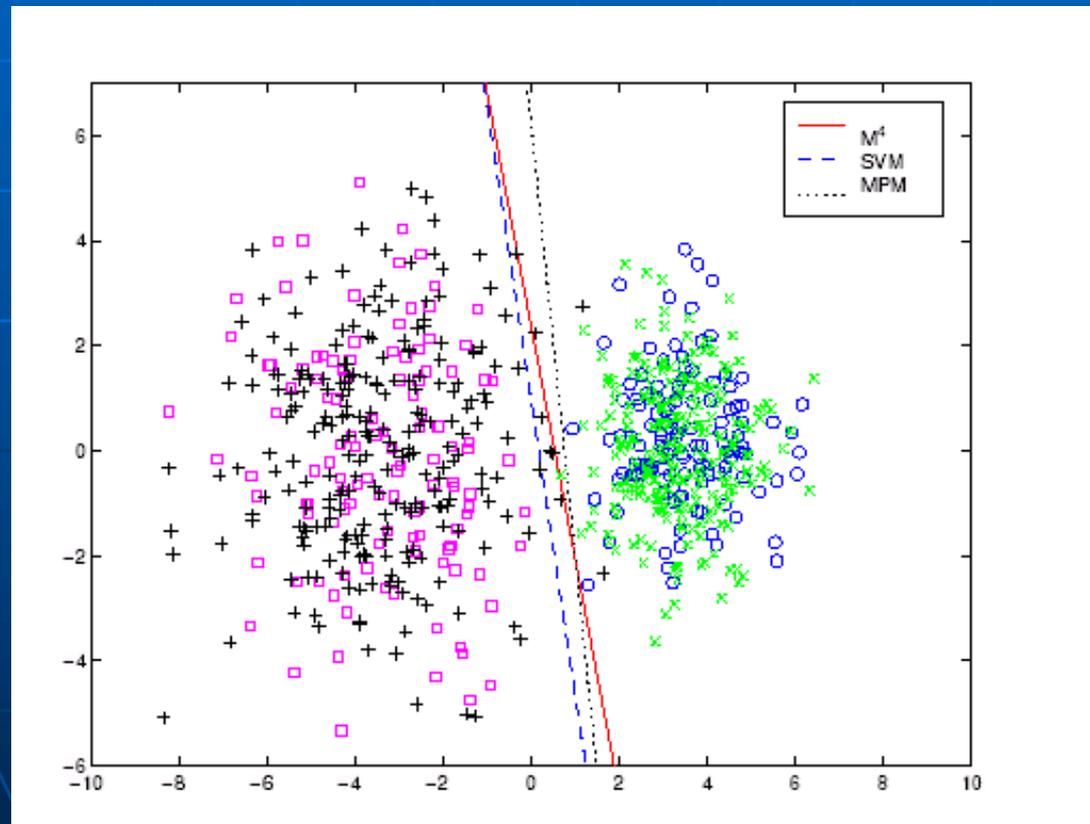
$\theta\alpha + (1-\theta)\beta$ vs. The overall test-set accuracy (TSA)



Experimental Results ---M⁴ (I)

- Synthetic Toy Data (1)

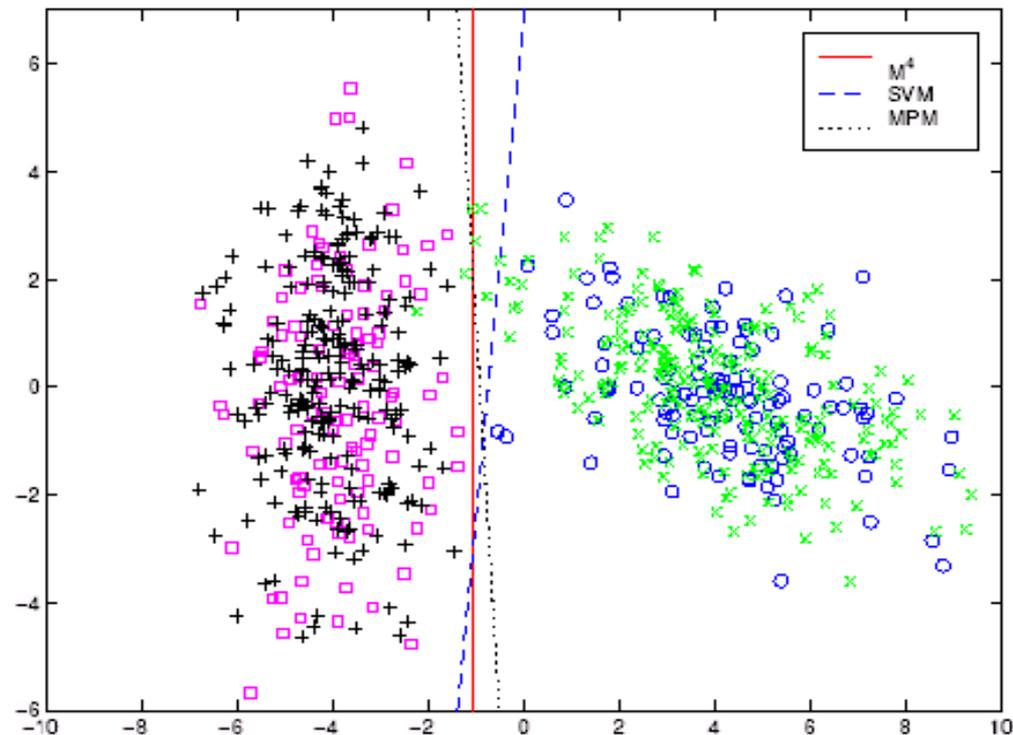
Two types of data with **the same data orientation** but **different data magnitude**



Experimental Results --- M^4 (II)

- Synthetic Toy Data (2)

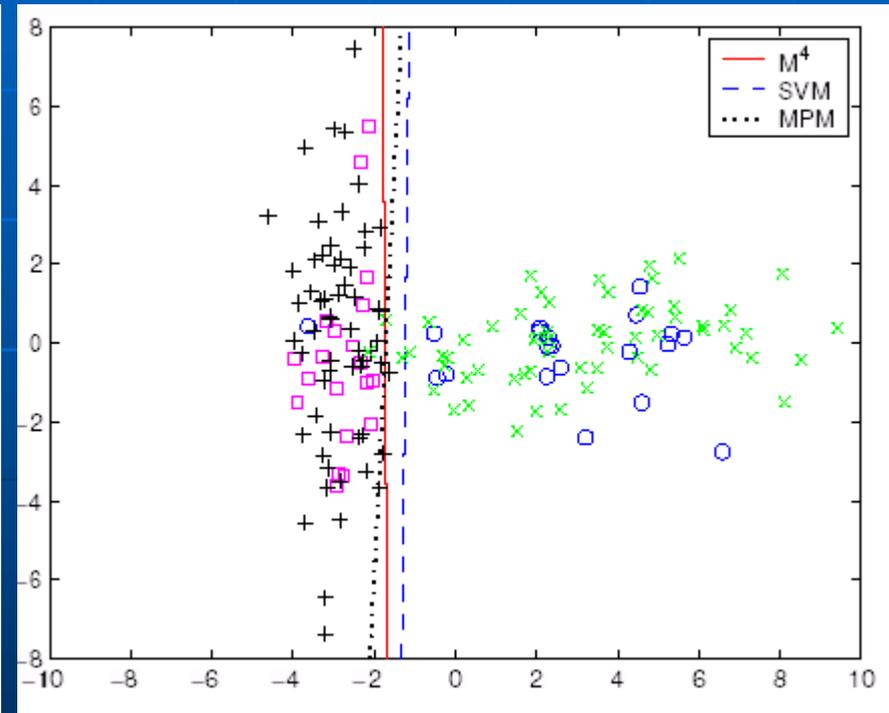
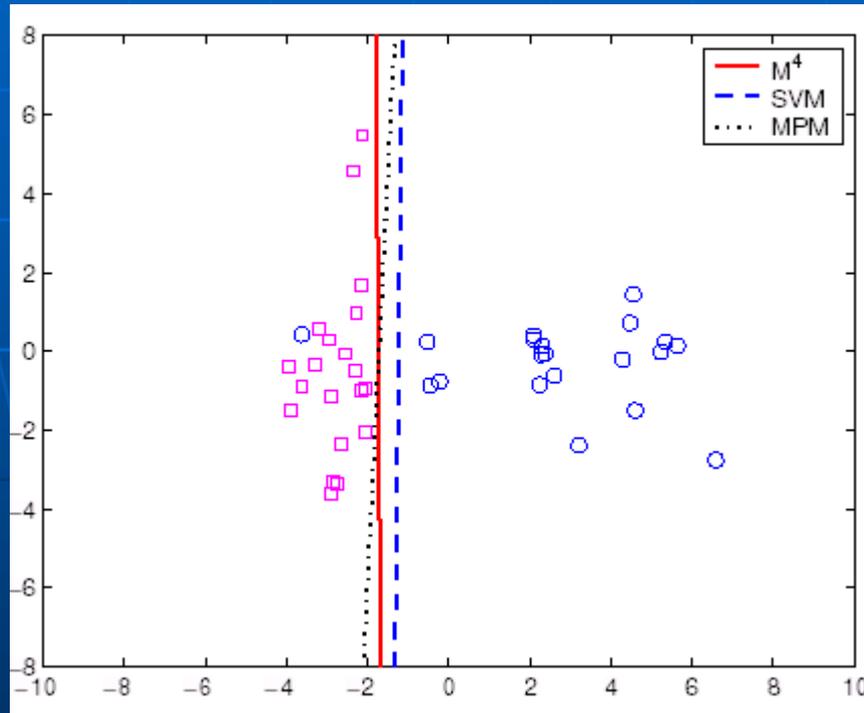
Two types of data with **the same data magnitude** but **different data orientation**



Experimental Results ---M⁴ (III)

- Synthetic Toy Data (3)

Two types of data with **the different data magnitude** and **different data orientation**



Experimental Results ---M⁴ (IV)

- Benchmark Data from UCI

Data set	Linear kernel			Gaussian kernel		
	M ⁴	SVM	MPM	M ⁴	SVM	MPM
Twonorm(%)	96.5 ± 0.6	95.1 ± 0.7	97.6 ± 0.5	96.5 ± 0.7	96.1 ± 0.4	97.6 ± 0.5
Breast(%)	97.5 ± 0.7	96.6 ± 0.5	96.9 ± 0.8	97.5 ± 0.6	96.7 ± 0.4	96.9 ± 0.8
Ionosphere(%)	87.7 ± 0.8	86.9 ± 0.6	84.8 ± 0.8	94.5 ± 0.4	94.2 ± 0.3	92.3 ± 0.6
Pima(%)	77.7 ± 0.9	77.9 ± 0.7	76.1 ± 1.2	77.6 ± 0.8	78.0 ± 0.5	76.2 ± 1.2
Sonar(%)	77.6 ± 1.2	76.2 ± 1.1	75.5 ± 1.1	84.9 ± 1.2	86.5 ± 1.1	87.3 ± 0.8
Vote(%)	96.1 ± 0.5	95.1 ± 0.4	94.8 ± 0.4	96.2 ± 0.5	95.9 ± 0.6	94.6 ± 0.4
Heart-disease(%)	86.6 ± 0.8	84.1 ± 0.7	83.2 ± 0.8	86.2 ± 0.8	83.8 ± 0.5	83.1 ± 1.0

Table 2: Comparisons of classification accuracies among M⁴, SVM, and MPM.



Future Work

- Speeding up M^4 and MEMPM
 - Contain support vectors—can we employ its sparsity as has been done in SVM?
 - Can we reduce redundant points??
- How to impose constraints on the kernelization for keeping the topology of data?
- Generalization error bound?
 - SVM and MPM have both error bounds.
- How to extend to multi-category classifications?
 - One vs. One or One vs. All?
 - Or seeking a principled way to construct multi-way boundary **in a step??**



Conclusion

- We propose a general global learning model MEMPM
 - A Worst-case distribution-free Bayes Optimal classifier
 - Containing an explicit error bound for future data
 - Subsuming BMPM which is idea for biased classification
- We propose a hybrid framework M^4 by learning from data locally and globally
 - This model subsumes three important models as special cases
 - SVM
 - MPM
 - FDA
 - Extended into regression tasks



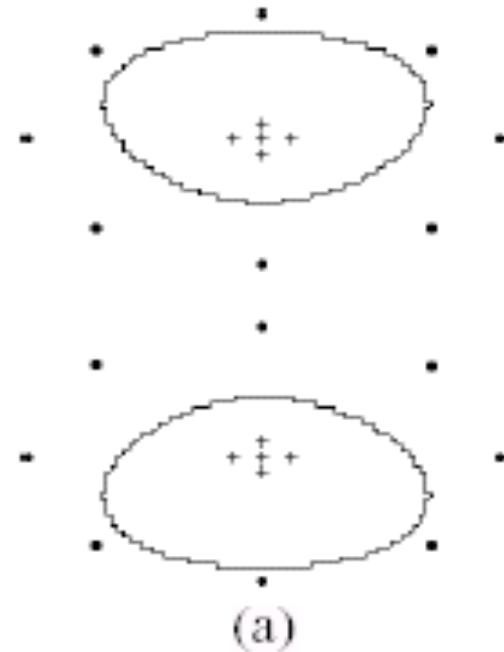
Discussion (I)

- In **linear** cases, M^4 **outperforms SVM and MPM**
- In **Gaussian** cases, M^4 is **slightly better or comparable than SVM**
 - (1) Sparsity in the feature space results in inaccurate estimation of covariance matrices
 - (2) Kernelization may not keep data topology of the original data.— Maximizing Margin in the feature space does not necessarily maximize margin in the original space



Discussion (II)

An example to illustrate that maximizing the margin in the feature space does not necessarily maximize the margin in the original space



(a) SVM using degree 4 polynomial kernel.
From Simon Tong et al. *Restricted Bayesian Optimal classifiers*, AAAI, 2000.



Setup

- Three concerns:

- Binary classification data sets

- For easy comparison. MPM (*Lanckriet et al. JMLR 02 or nips02*) also uses these data sets.

- Medium or smaller size Data sets

Appendix A: MEMPM- BMPM (I)

1

Lemma Given $\mathbf{w} \neq \mathbf{0}$ and b , such that $\mathbf{w}^T \mathbf{y} \leq b$ and $\beta \in [0, 1)$, the condition

$$\inf_{\mathbf{y} \sim (\bar{\mathbf{y}}, \Sigma_{\mathbf{y}})} \Pr\{\mathbf{w}^T \mathbf{y} \leq b\} \geq \beta,$$

holds if and only if $b - \mathbf{w}^T \bar{\mathbf{y}} \geq \kappa(\beta) \sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}$ with $\kappa(\beta) = \sqrt{\frac{\beta}{1-\beta}}$.

2

$$\begin{aligned} \max_{\alpha, \mathbf{w} \neq \mathbf{0}, b} \quad & \alpha \quad \text{s.t.} \\ & -b + \mathbf{w}^T \bar{\mathbf{x}} \geq \kappa(\alpha) \sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}}, \\ & b - \mathbf{w}^T \bar{\mathbf{y}} \geq \kappa(\beta_0) \sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}, \end{aligned}$$

3

$$\mathbf{w}^T \bar{\mathbf{y}} + \kappa(\beta_0) \sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}} \leq b \leq \mathbf{w}^T \bar{\mathbf{x}} - \kappa(\alpha) \sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}}.$$

4

$$\begin{aligned} \max_{\alpha, \mathbf{w} \neq \mathbf{0}} \quad & \alpha \quad \text{s.t.} \\ & 1 \geq \kappa(\alpha) \sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}} + \kappa(\beta_0) \sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}, \\ & \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1. \end{aligned}$$

5

$$\max_{\mathbf{w} \neq \mathbf{0}} \frac{1 - \kappa(\beta_0) \sqrt{\mathbf{w}^T \Sigma_{\mathbf{y}} \mathbf{w}}}{\sqrt{\mathbf{w}^T \Sigma_{\mathbf{x}} \mathbf{w}}} \quad \text{s.t.} \quad \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1.$$



Fractional Programming



Appendix A: MEMPM- BMPM (II)

Solving Fractional Programming problem

- Parametric Method

Find \mathbf{w} by solving

$$\max_{\mathbf{w} \neq 0} 1 - \kappa(\gamma) \sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}} - \lambda \sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}} \quad \text{s.t.} \quad \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1$$

Update

$$\lambda \leftarrow \frac{1 - \kappa(\gamma) \sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}}}{\sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}}}$$

- Equivalently

$$\min_{\mathbf{w} \neq 0} \kappa(\gamma) \sqrt{\mathbf{w}^T \Sigma_y \mathbf{w}} + \lambda \sqrt{\mathbf{w}^T \Sigma_x \mathbf{w}} \quad \text{s.t.} \quad \mathbf{w}^T (\bar{\mathbf{x}} - \bar{\mathbf{y}}) = 1$$

■ Least-squares approach



Appendix B: Optimization of LSVR(I)

$$\begin{aligned} \min_{\mathbf{w}, b, t_i, \xi_i, \xi_i^*} \quad & \frac{1}{N} \sum_{i=1}^N t_i + C \sum_{i=1}^N (\xi_i + \xi_i^*), \\ \text{s.t.} \quad & y_i - (\mathbf{w}^T \mathbf{x}_i + b) \leq \epsilon \sqrt{\mathbf{w}^T \Sigma_i \mathbf{w}} + \xi_i, \\ & (\mathbf{w}^T \mathbf{x}_i + b) - y_i \leq \epsilon \sqrt{\mathbf{w}^T \Sigma_i \mathbf{w}} + \xi_i^*, \\ & \sqrt{\mathbf{w}^T \Sigma_i \mathbf{w}} \leq t_i, \end{aligned}$$



Hard to be solved...



Appendix B: Optimization of LSVR(II)

Can be relaxed as the following:

$$\begin{aligned} \min_{\mathbf{w}, b, t_i, \xi_i, \xi_i^*} \quad & \frac{1}{N} \sum_{i=1}^N t_i + C \sum_{i=1}^N (\xi_i + \xi_i^*), \\ \text{s.t.} \quad & y_i - (\mathbf{w}^T \mathbf{x}_i + b) \leq \epsilon t_i + \xi_i, \\ & (\mathbf{w}^T \mathbf{x}_i + b) - y_i \leq \epsilon t_i + \xi_i^*, \\ & \sqrt{\mathbf{w}^T \Sigma_i \mathbf{w}} \leq t_i, \\ & t_i \geq 0, \xi_i \geq 0, \xi_i^* \geq 0, i = 1, \dots, N. \end{aligned}$$



Second-Order Cone Programming



Appendix C: Convex Optimization

Linear Program:

$$\min_x c^T x \quad \text{s.t.} \quad Ax \preceq b \\ Fx = g$$

LP SVM
(Mangasarian, Bennett)

Quadratic Program:

$$\min_x x^T P x + 2q^T x + r \quad \text{s.t.} \quad Ax \preceq b \\ Fx = g$$

SVM
(Vapnik)

Quadratic Constrained Quadratic Program:

$$\min_x x^T P_0 x + 2q_0^T x + r_0 \quad \text{s.t.} \quad x^T P_i x + 2q_i^T x + r_i \leq 0 \\ i = 1, \dots, L \quad (\text{convex if } P_i \succeq 0)$$

Kernel Fisher Discriminant
(Mika et al.)

Second Order Cone Program:

$$\min_x c^T x \quad \text{s.t.} \quad \|A_i x + b_i\|_2 \leq e_i^T x + d_i \\ i = 1, \dots, L$$

Minimax Probability Machine
(Lanckriet et al.)

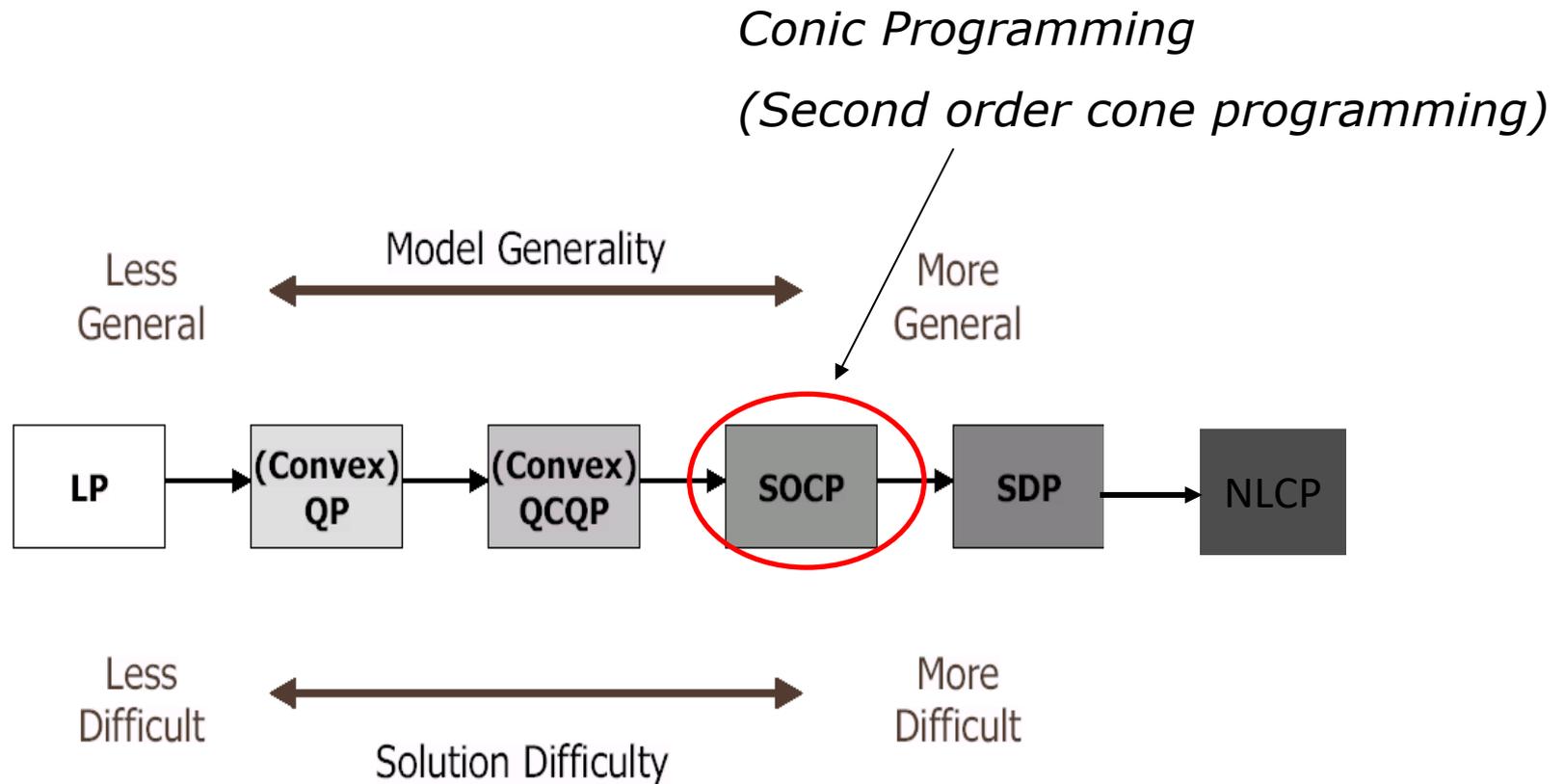
Semi-Definite Program:

$$\min_x c^T x \quad \text{s.t.} \quad A(x) = A_0 + \sum_{i=1}^n x_i A_i \succeq 0 \quad (A_i = A_i^T \in \mathbb{R}^{p \times p}) \\ Fx = g$$

Kernel matrix learning
(Lanckriet, Cristianini et al.)



Appendix C: Convex Optimization



Appendix C: Convex Optimization - SOCP

$$\begin{aligned} \min \quad & f^T x \\ \text{s.t.} \quad & \|C_i x + d_i\| \leq a_i^T x + b_i, \quad i = 1, \dots, N \end{aligned}$$

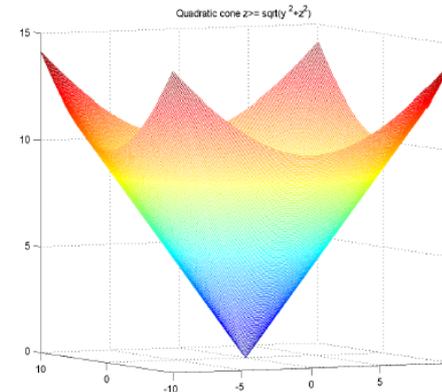
$$x \in \mathbb{R}^n \quad f, a_i \in \mathbb{R}^n \quad C_i \in \mathbb{R}^{(k_i-1) \times n} \quad d_i \in \mathbb{R}^{k_i-1} \quad b_i \in \mathbb{R}$$

Equivalent to conic program

- **Linear constraints:** cone dimension $k=1$
- **Cone constraints:** change of variables
(vector) $y = C_i x + d_i, \quad z = a_i^T x + b_i$

8

Quadratic cone C sometimes also called **Lorentz cone** (or ice cream cone)



Trivial Quadratic Cone:

$$z \geq \sqrt{x^2 + y^2}$$

7



Appendix C: SOCP-Solver

Sedumi (MATLAB)

Loqo (C, MATLAB)

MOSEK (C, MATLAB)

SDPT3 (MATLAB+C or FORTRAN)

The worst-case cost is $O(n^3)$



Time Complexity

Models	Time Complexity
MEMPM	$O(Ln^3 + Nn^2)$
BMPM	$O(n^3 + Nn^2)$
M ⁴	$O(Nn^3)$
LS-SVM	$O(n^3 + Nn^2)$
LSVR	$O(Nn^3)$
LS-SVR	$O(n^3 + Nn^2)$



Time Complexity

Thus we believe that for practical purposes the cost of solving an SOCP is roughly equal to the cost of solving a modest number (5–50) of systems of the form (40). If no special structure in the problem data is exploited, the cost of solving the system is $O(n^3)$, and the cost of forming the system matrix is $O(n^2 \sum_{i=1}^N n_i)$. In practice, special problem structure (*e.g.*, sparsity) often allows forming the equations faster, or solving the systems (39) or (40) more efficiently.

----“*Applications of Second Order Cone Programming*”,
Lobo, Boyd et al. in Linear Algebra and Applications.

