

SUPPORT VECTOR MACHINES & NEURAL NETWORKS

LECTURE 6 – SUPPORT VECTOR MACHINES PART # III

A. Bi-classification

History, LSVM, Approximate LSVM, Soft LSVM,
Kernel-based linear SVM, nonlinear SVM

B. Multi-classification

OVO, OVA, Twin SVM

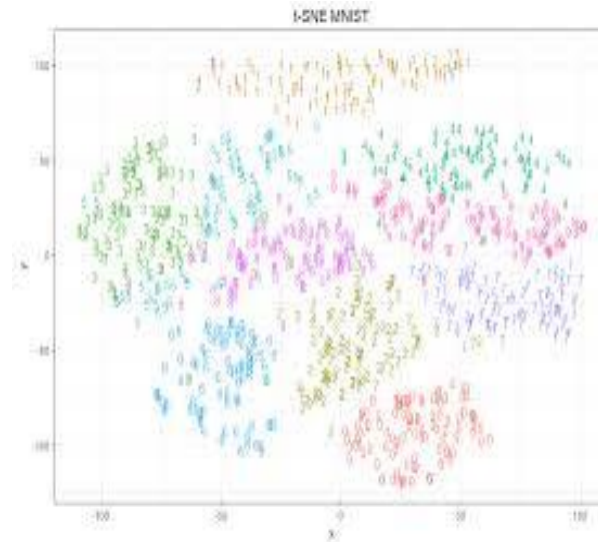
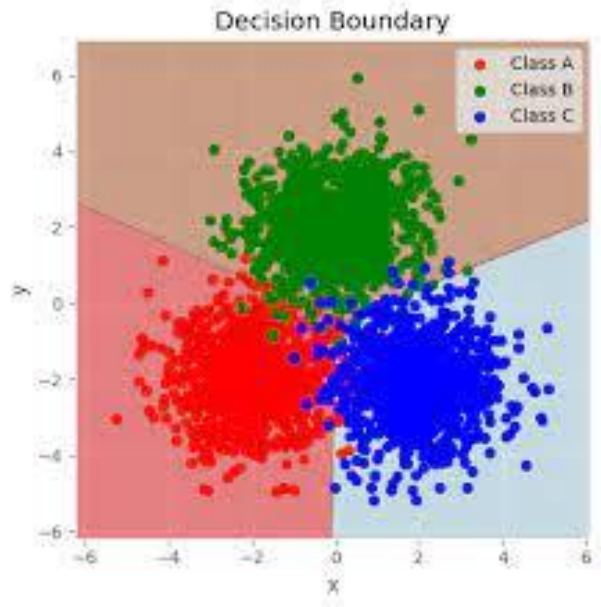
C. Prediction

Support Vector Regression (SVR)

*Copyright: Professor Shu-Cherng Fang of NCSU-ISE

Multi-class classification

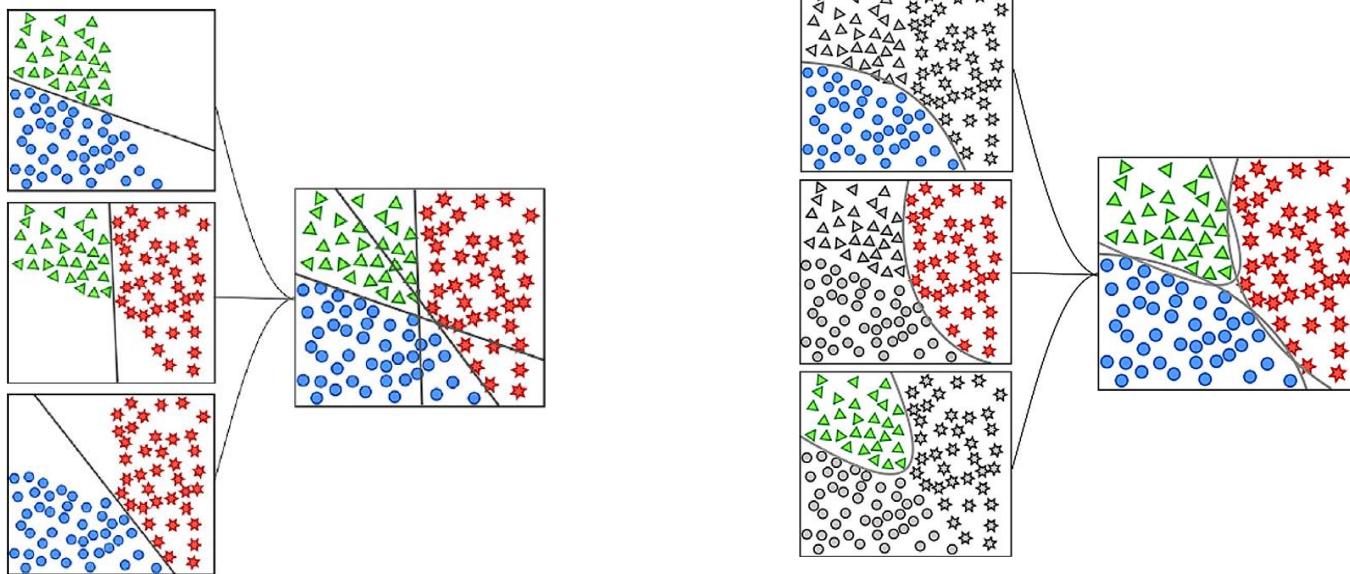
- **Multi-class classification** (multi-classification) is a problem of classifying instances into one of three or more classes.



- Pictures taken from Wikipedia

Basic ideas of multiclass classification

- If we are given a dataset in C classes and a new data-point $x \in \mathbb{R}^n$, we may consider two commonly adopted approaches to determine which class x belongs to:
 1. OVO (one vs. one) approach
 2. OVA (one vs. all) approach



Basic ideas of multiclass classification

- 1. OVO (one vs. one) approach

- Take each pair of different classes (C_i, C_j) .
- Label all data-points in C_i with a label +1;
Label all data-points in C_j with a label -1 .
- Apply a bi-class SVM classifier to find (\mathbf{w}, b)
and determine $\mathbf{x} \in C_i$ or $\mathbf{x} \in C_j$.
- Assign a voting score to each class of the pair

$$\text{Score}(C_i) = \begin{cases} 1, & \text{if } \mathbf{x} \in C_i \\ 0, & \text{if } \mathbf{x} \in C_j \end{cases}$$

- Sum up scores/votes over all pairs involving C_i .
- **Decision:** \mathbf{x} belongs to the class with the highest total score.

OVO related issues

- How many SVM involved?

of all possible pairs = $c(c-1)/2$

- Tie breaker?

secondary score?

- Better scoring measure?

Basic ideas of multiclass classification

- 2. OVA (one vs. all of the rest) approach
 - Take each class C_i .
 - Label all data-points **in** C_i with a label +1;
Label all data-points **not in** C_i with a label - 1.
 - Apply a bi-class SVM classifier to find (\mathbf{w}, b)
and determine $\mathbf{x} \in C_i$ or $\mathbf{x} \notin C_i$.
 - Assign a score to each class
 $Score(C_i) = \mathbf{w}^T \mathbf{x} + b$
 - **Decision:** \mathbf{x} belongs to the class with highest score.

OVA related issues

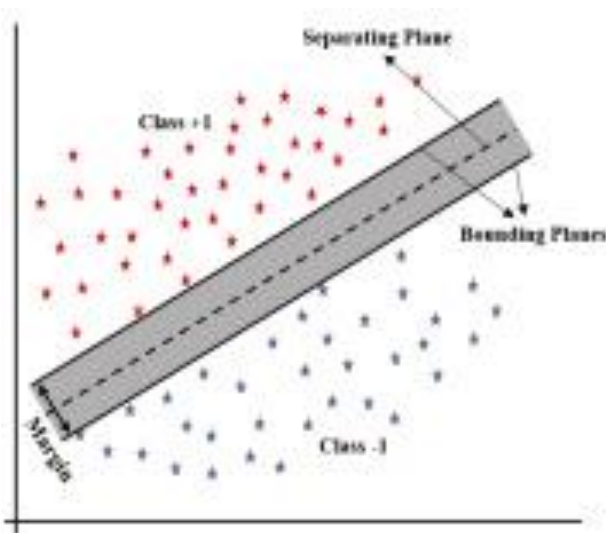
- How many SVM involved?
of all possible pairs = c ($\ll c(c-1)/2$)
- Imbalanced datasets
quality of results?
- Better scoring ?

- Any approaches other than OVO and OVA ?
- tournament ?
- Which SVM model to use?

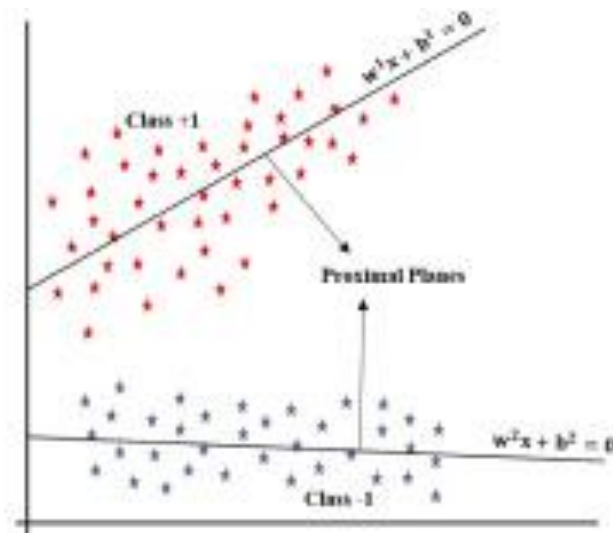
Beyond LSSVM and KSSVM

- **Motivation:** Who says we should use **only one** separation hyperplane (surface) for bi-classification?

Picture from ScienceDirect.com



(a)



(b)

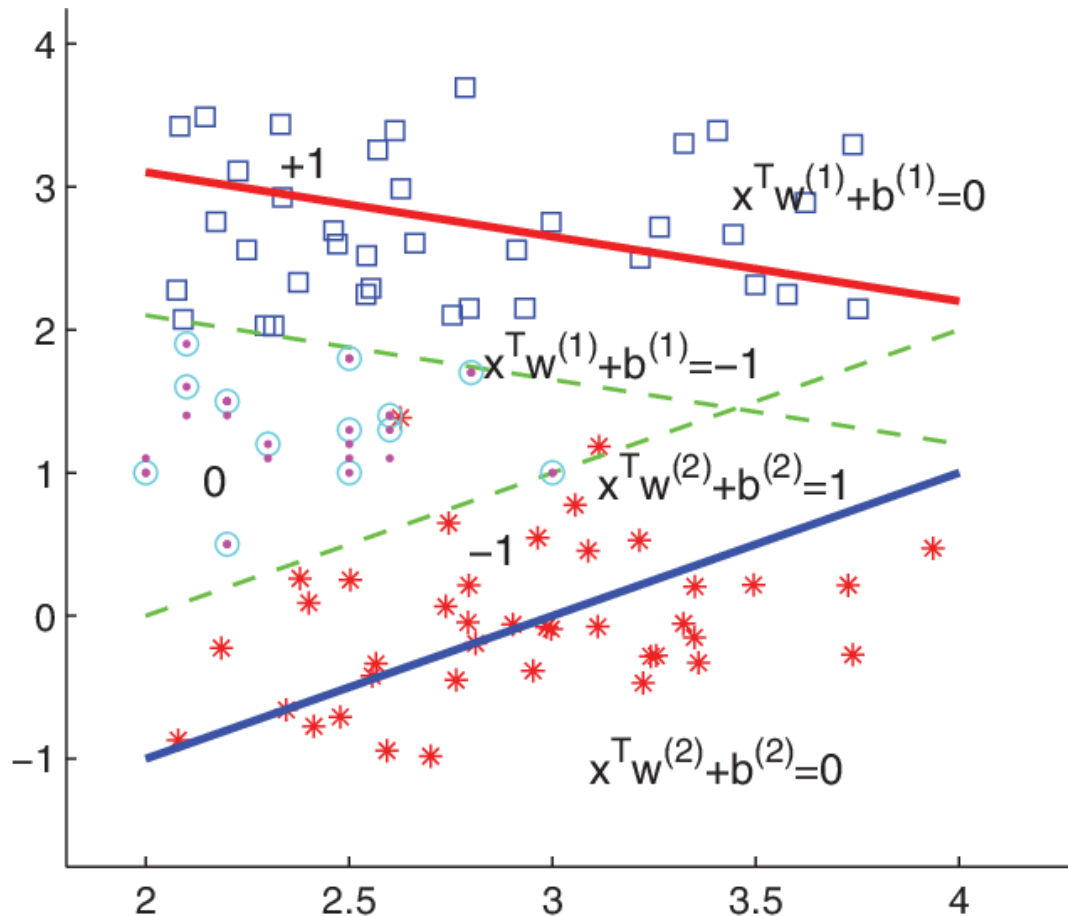
Twin support vector machines - TWSVM

- Basic ideas:

- (i) Use **two non-parallel SVM** to separate two distinct classes of data;
- (ii) All points in **one class center around** a corresponding SVM separation hyperplane while **all points in the other class are kept away** from this hyperplane for a safe distance;
- (iii) When a new point comes into the picture, it is classified based on its **“distance” to each hyperplane.**

Realization of TWSVM

- Picture from researchgate.net



Twin support vector machines

- **References:**

1. O.L. Mangasarian and E.W. Wild, "Multisurface Proximal Support Vector Classification via Generalized Eigenvalues," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 28, no. 1, pp. 69-74, [2006](#).
2. Jayadeva, R. Khemchandani and S. Chandra, "Twin Support Vector Machines for Pattern Classification," IEEE Trans. Pattern Analysis and Machine Intelligence, vol. 29, no. 5, pp. 905-910, [2007](#).

Twin support vector machines

Problem Setting:

- Dataset $S = \{\mathbf{x}^i \in \mathbb{R}^n \mid i = 1, 2, \dots, N\}$ of N data points of n attributes.
- Two classes: $S = S_A \cup S_B, S_A \cap S_B = \emptyset$ with $|S_A| = N_A$ and $|S_B| = N_B$.
- Denote $S_A = \{\mathbf{x}_A^i \in \mathbb{R}^n \mid i = 1, 2, \dots, N_A\}$,
 $S_B = \{\mathbf{x}_B^j \in \mathbb{R}^n \mid j = 1, 2, \dots, N_B\}$.
- Labels
 $y_i = +1$, for $\mathbf{x}_A^i \in S_A, i = 1, 2, \dots, N_A$
 $y_j = -1$, for $\mathbf{x}_B^j \in S_B, j = 1, 2, \dots, N_B$

Twin support vector machines

Problem Setting:

- Separation hyperplanes:

$$\text{For } S_A \text{ Class: } \mathbf{w}^{(1)T} \mathbf{x} + b^{(1)} = 0$$

$$\text{For } S_B \text{ Class: } \mathbf{w}^{(2)T} \mathbf{x} + b^{(2)} = 0$$

$$\mathbf{w}^{(i)} \in \mathbb{R}^n, b^{(i)} \in \mathbb{R}, i = 1, 2.$$

$$\text{Denote } \mathbf{u}^{(i)} = \begin{pmatrix} \mathbf{w}^{(i)} \\ b^{(i)} \end{pmatrix} \in \mathbb{R}^{n+1}, i = 1, 2.$$

- Data preparation

$$\text{data records of } S_A : X_A^T = [x_A^1, \dots, x_A^{N_A}] \in \mathbf{M}^{n \times N_A}$$

$$\text{data records of } S_B : X_B^T = [x_B^1, \dots, x_B^{N_B}] \in \mathbf{M}^{n \times N_B}$$

unit vector $\mathbf{e}_A \in \mathbb{R}^{N_A}$: a column vector with all elements being 1.

unit vector $\mathbf{e}_B \in \mathbb{R}^{N_B}$: a column vector with all elements being 1.

$$\text{Denote } \hat{X}_A^T = \begin{pmatrix} X_A^T \\ \mathbf{e}_A^T \end{pmatrix} \in \mathbf{M}^{(n+1) \times N_A} \text{ and } \hat{X}_B^T = \begin{pmatrix} X_B^T \\ \mathbf{e}_B^T \end{pmatrix} \in \mathbf{M}^{(n+1) \times N_B}$$

Twin soft support vector machine model - TWSSVM

(TWSSVM-1)

$$\begin{aligned} \text{Min} \quad & \frac{1}{2} \sum_{i=1}^{N_A} \left(\mathbf{w}^{(1)T} \mathbf{x}_A^i + b^{(1)} \right)^2 + \mathcal{C}_1 \sum_{j=1}^{N_B} \xi_j^{(1)} \\ \text{s.t.} \quad & y_j \left(\mathbf{w}^{(1)T} \mathbf{x}_B^j + b^{(1)} \right) \geq 1 - \xi_j^{(1)}, j = 1, \dots, N_B \\ & \mathbf{w}^{(1)} \in \mathbb{R}^n, b^{(1)} \in \mathbb{R}, \xi^{(1)} \in \mathbb{R}_+^{N_B} \end{aligned}$$

(TWSSVM-2)

$$\begin{aligned} \text{Min} \quad & \frac{1}{2} \sum_{i=1}^{N_B} \left(\mathbf{w}^{(2)T} \mathbf{x}_B^i + b^{(2)} \right)^2 + \mathcal{C}_2 \sum_{i=1}^{N_A} \xi_i^{(2)} \\ \text{s.t.} \quad & y_i \left(\mathbf{w}^{(2)T} \mathbf{x}_A^i + b^{(2)} \right) \geq 1 - \xi_i^{(2)}, i = 1, \dots, N_A \\ & \mathbf{w}^{(2)} \in \mathbb{R}^n, b^{(2)} \in \mathbb{R}, \xi^{(2)} \in \mathbb{R}_+^{N_A} \end{aligned}$$

Observations

- Both (TWSSVM-1) and (TWSSVM-2) are **convex quadratic programming problems**.
- Compared to LSSVM, (TWSSVM-1) and (TWSSVM-2) have the **same number of variables but fewer constraints**.
- The **total complexity** of solving two smaller QPs is about $\frac{1}{4}$ of that solving LSSVM.
- The accuracy of TWSSVM is **no inferior** to that of LSSVM.

Example

- From reference [2]

Training Times (in Seconds)

Data Set	TWSVM (EXE file)	TWSVM (DLL file)	SVM (DLL file)
Hepatitis (155×19)	4.37	4.85	12.7
Sonar (208×60)	4.62	6.64	24.9
Heart-statlog (270×14)	4.72	11.3	50.9
Heart-c (303×14)	8.37	14.92	68.2
Ionosphere (351×34)	9.93	25.9	102.2
Votes (435×16)	12.8	45.8	189.4
Australian (690×14)	37.4	142.1	799.2
Pima-Indian (768×8)	56.9	231.5	1078.6
CMC (1473×9)	63.4	1737.9	6827.8

Data Set	TWSVM	GEPSVM	SVM
Heart-statlog (270×14)	84.44±4.32	84.81±3.87	84.07±4.40
Heart-c (303×14)	83.80±5.53	84.44±5.27	82.82±5.15
Hepatitis (155×19)	80.79±12.24	58.29±19.07	80.00±8.30
Ionosphere (351×34)	88.03±2.81	75.19±5.50	86.04 ±2.37
Sonar (208×60)	77.26±10.10	66.76±10.75	79.79±5.31
Votes (435×16)	96.08±3.29	91.93±3.18	94.50±2.71
Pima-Indian(768×8)	73.70±3.97	74.60±5.07	76.68±2.90
Australian (690×14)	85.80±5.05	85.65±4.60	85.51±4.58
CMC (1473×9)	67.28±2.21	65.99±2.30	67.82±2.63

Accuracies have been indicated as percentages.

TWSSVM

- TWSSVM classifier

$$\mathit{class}_{TWSSVM}(\mathbf{x}) = \mathit{argmin} \{ |f_A(\mathbf{x})|, |f_B(\mathbf{x})| \}$$

- Primal version TWSSVM

$$f_A(\mathbf{x}) = \mathbf{w}^{(1)T} \mathbf{x} + b^{(1)}$$

$$f_B(\mathbf{x}) = \mathbf{w}^{(2)T} \mathbf{x} + b^{(2)}$$

Dual TWSSVM and kernel-based TWSSVM

- Basic approach:

1. Following the same procedure of finding dual LSSVM, we can derive dual TWSSVM.
2. Following the same procedure of finding kernel-based LSSVM, we can derive kernel-based TWSSVM.

TWSSVM model – vector form

(TWSSVM-1)

$$\begin{aligned} \text{Min} \quad & \frac{1}{2} \|X_A \mathbf{w}^{(1)} + \mathbf{e}_A b^{(1)}\|_2^2 + C_1 \mathbf{e}_B^T \boldsymbol{\xi}^{(1)} \\ \text{s.t.} \quad & (-1)(X_B \mathbf{w}^{(1)} + \mathbf{e}_B b^{(1)}) \geq \mathbf{e}_B - \boldsymbol{\xi}^{(1)} \\ & \mathbf{w}^{(1)} \in \mathbb{R}^n, b^{(1)} \in \mathbb{R}, \boldsymbol{\xi}^{(1)} \in \mathbb{R}_+^{N_B} \end{aligned}$$

(TWSSVM-2)

$$\begin{aligned} \text{Min} \quad & \frac{1}{2} \|X_B \mathbf{w}^{(2)} + \mathbf{e}_B b^{(2)}\|_2^2 + C_2 \mathbf{e}_A^T \boldsymbol{\xi}^{(2)} \\ \text{s.t.} \quad & (+1)(X_A \mathbf{w}^{(2)} + \mathbf{e}_A b^{(2)}) \geq \mathbf{e}_A - \boldsymbol{\xi}^{(2)} \\ & \mathbf{w}^{(2)} \in \mathbb{R}^n, b^{(2)} \in \mathbb{R}, \boldsymbol{\xi}^{(2)} \in \mathbb{R}_+^{N_A} \end{aligned}$$

Dual TWSSVM-1 model

- Lagrangian

$$L(\mathbf{w}^{(1)}, b^{(1)}, \boldsymbol{\xi}^{(1)}, \boldsymbol{\alpha}, \boldsymbol{\theta}) = \frac{1}{2} \|X_A \mathbf{w}^{(1)} + \mathbf{e}_A b^{(1)}\|_2^2 + \mathcal{C}_1 \mathbf{e}_B^T \boldsymbol{\xi}^{(1)} + \boldsymbol{\alpha}^T (\mathbf{e}_B - \boldsymbol{\xi}^{(1)} + X_B \mathbf{w}^{(1)} + \mathbf{e}_B b^{(1)}) - \boldsymbol{\theta}^T \boldsymbol{\xi}^{(1)}$$

- K-K-T conditions

(i) $X_A^T (X_A \mathbf{w}^{(1)} + \mathbf{e}_A b^{(1)}) + X_B^T \boldsymbol{\alpha} = 0;$

(ii) $\mathbf{e}_A^T (X_A \mathbf{w}^{(1)} + \mathbf{e}_A b^{(1)}) + \mathbf{e}_B^T \boldsymbol{\alpha} = 0;$

(iii) $\mathcal{C}_1 \mathbf{e}_B - \boldsymbol{\alpha} - \boldsymbol{\theta} = 0;$

(iv) $-(X_B \mathbf{w}^{(1)} + \mathbf{e}_B b^{(1)}) + \boldsymbol{\xi}^{(1)} \geq \mathbf{e}_B, \boldsymbol{\xi}^{(1)} \geq 0;$

(v) $\boldsymbol{\alpha}^T (-(X_B \mathbf{w}^{(1)} + \mathbf{e}_B b^{(1)}) + \boldsymbol{\xi}^{(1)} - \mathbf{e}_B) = 0, \boldsymbol{\theta}^T \boldsymbol{\xi}^{(1)} = 0;$

(vi) $\boldsymbol{\alpha} \geq 0, \boldsymbol{\theta} \geq 0.$

Dual TWSSVM-1 model

(iii) says $C_1 \mathbf{e}_B \geq \boldsymbol{\alpha} \geq 0$,

$$\Leftrightarrow 0 \leq \alpha_j \leq C_1 \text{ for } j = 1, \dots, N_B$$

(i)+(ii) says

$$\begin{pmatrix} X_A^T \\ \mathbf{e}_A^T \end{pmatrix} (X_A, \mathbf{e}_A) \begin{pmatrix} \mathbf{w}^{(1)} \\ b^{(1)} \end{pmatrix} + \begin{pmatrix} X_B^T \\ \mathbf{e}_B^T \end{pmatrix} \boldsymbol{\alpha} = \mathbf{0}$$

$$\Leftrightarrow \hat{X}_A^T \hat{X}_A \mathbf{u}^{(1)} + \hat{X}_B^T \boldsymbol{\alpha} = \mathbf{0}$$

$$\Leftrightarrow \mathbf{u}^{(1)} = - (\hat{X}_A^T \hat{X}_A)^{-1} \hat{X}_B^T \boldsymbol{\alpha} \text{ [practically, use } (\hat{X}_A^T \hat{X}_A + \varepsilon I)^{-1}]$$

• dual objective function

$$h(\boldsymbol{\alpha}, \boldsymbol{\theta}) = \text{Min}_{\mathbf{w}^{(1)} \in \mathbb{R}^n, b^{(1)} \in \mathbb{R}, \boldsymbol{\xi}^{(1)} \in \mathbb{R}_+^{N_B}} L(\mathbf{w}^{(1)}, b^{(1)}, \boldsymbol{\zeta}^{(1)}, \boldsymbol{\alpha}, \boldsymbol{\theta})$$

$$= -\frac{1}{2} \boldsymbol{\alpha}^T \hat{X}_B (\hat{X}_A^T \hat{X}_A)^{-1} \hat{X}_B^T \boldsymbol{\alpha} + \mathbf{e}_B^T \boldsymbol{\alpha}$$

Dual TWSSVM model

(DTWSSVM-1)

$$\text{Max} \quad -\frac{1}{2} \boldsymbol{\alpha}^T \hat{X}_B (\hat{X}_A^T \hat{X}_A)^{-1} \hat{X}_B^T \boldsymbol{\alpha} + \mathbf{e}_B^T \boldsymbol{\alpha}$$

$$\text{s.t.} \quad 0 \leq \alpha_j \leq C_1 \text{ for } j = 1, \dots, N_B$$

$$\text{dual-primal conversion: } \begin{pmatrix} \mathbf{w}^{(1)} \\ b^{(1)} \end{pmatrix} = \mathbf{u}^{(1)} = - (\hat{X}_A^T \hat{X}_A)^{-1} \hat{X}_B^T \boldsymbol{\alpha}$$

*A simple convex quadratic programming problem.

Similarly, we have

(DTWSSVM-2)

$$\text{Max} \quad -\frac{1}{2} \boldsymbol{\gamma}^T \hat{X}_A (\hat{X}_B^T \hat{X}_B)^{-1} \hat{X}_A^T \boldsymbol{\gamma} + \mathbf{e}_A^T \boldsymbol{\gamma}$$

$$\text{s.t.} \quad 0 \leq \gamma_i \leq C_2 \text{ for } i = 1, \dots, N_A$$

$$\text{dual-primal conversion: } \begin{pmatrix} \mathbf{w}^{(2)} \\ b^{(2)} \end{pmatrix} = \mathbf{u}^{(2)} = - (\hat{X}_B^T \hat{X}_B)^{-1} \hat{X}_A^T \boldsymbol{\gamma}$$

DTWSSVM

- TWSSVM classifier

$$\mathit{class}_{TWSSVM}(\mathbf{x}) = \mathit{argmin} \{ |f_A(\mathbf{x})|, |f_B(\mathbf{x})| \}$$

- Dual version TWSSVM

Denote $\hat{\mathbf{x}} = \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}$

$$f_A(\mathbf{x}) = -\left((\hat{X}_A^T \hat{X}_A)^{-1} \hat{X}_B^T \boldsymbol{\alpha} \right)^T \hat{\mathbf{x}}$$

$$f_B(\mathbf{x}) = -\left((\hat{X}_B^T \hat{X}_B)^{-1} \hat{X}_A^T \boldsymbol{\gamma} \right)^T \hat{\mathbf{x}}$$

Kernel-based twin soft SVM - KTWSSVM

- Using a *feature map* $\phi(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^l$ ($l \geq n$) to transform the problem to a higher dimensional space for linear separability.

(KTWSSVM-1)

$$\begin{aligned} \text{Min} \quad & \frac{1}{2} \sum_{i=1}^{N_A} \left(\mathbf{w}^{(1)T} \phi(\mathbf{x}_A^i) + b^{(1)} \right)^2 + \mathcal{C}_1 \sum_{j=1}^{N_B} \xi_j^{(1)} \\ \text{s.t.} \quad & y_j \left(\mathbf{w}^{(1)T} \phi(\mathbf{x}_B^j) + b^{(1)} \right) \geq 1 - \xi_j^{(1)}, j = 1, \dots, N_B \\ & \mathbf{w}^{(1)} \in \mathbb{R}^l, b^{(1)} \in \mathbb{R}, \boldsymbol{\xi}^{(1)} \in \mathbb{R}_+^{N_B} \end{aligned}$$

(KTWSSVM-2)

$$\begin{aligned} \text{Min} \quad & \frac{1}{2} \sum_{i=1}^{N_B} \left(\mathbf{w}^{(2)T} \phi(\mathbf{x}_B^i) + b^{(2)} \right)^2 + \mathcal{C}_2 \sum_{i=1}^{N_A} \xi_i^{(2)} \\ \text{s.t.} \quad & y_i \left(\mathbf{w}^{(2)T} \phi(\mathbf{x}_A^i) + b^{(2)} \right) \geq 1 - \xi_i^{(2)}, i = 1, \dots, N_A \\ & \mathbf{w}^{(2)} \in \mathbb{R}^l, b^{(2)} \in \mathbb{R}, \boldsymbol{\xi}^{(2)} \in \mathbb{R}_+^{N_A} \end{aligned}$$

Dual kernel-based TWSSVM model

(DKTWSSVM-1)

$$\text{Max} \quad -\frac{1}{2} \boldsymbol{\alpha}^T \phi(X_B) (\phi(X_A)^T \phi(X_A))^{-1} \phi(X_B)^T \boldsymbol{\alpha} + \mathbf{e}_B^T \boldsymbol{\alpha}$$

$$\text{s.t.} \quad 0 \leq \alpha_j \leq C_1 \text{ for } j = 1, \dots, N_B$$

** Kernel matrix $K_A \triangleq \phi(X_A)^T \phi(X_A)$

(DKTWSSVM-2)

$$\text{Max} \quad -\frac{1}{2} \boldsymbol{\gamma}^T \phi(X_A) (\phi(X_B)^T \phi(X_B))^{-1} \phi(X_A)^T \boldsymbol{\gamma} + \mathbf{e}_A^T \boldsymbol{\gamma}$$

$$\text{s.t.} \quad 0 \leq \gamma_i \leq C_2 \text{ for } i = 1, \dots, N_A$$

** Kernel matrix $K_B \triangleq \phi(X_B)^T \phi(X_B)$

* $\phi(X_A)^T \in \mathbf{M}^{(l+1) \times N_A}$ formed by $\{\phi(x_A^i)\}$ with the last row being all 1's.

$\phi(X_B)^T \in \mathbf{M}^{(l+1) \times N_B}$ form by $\{\phi(x_B^j)\}$ with the last row being all 1's.