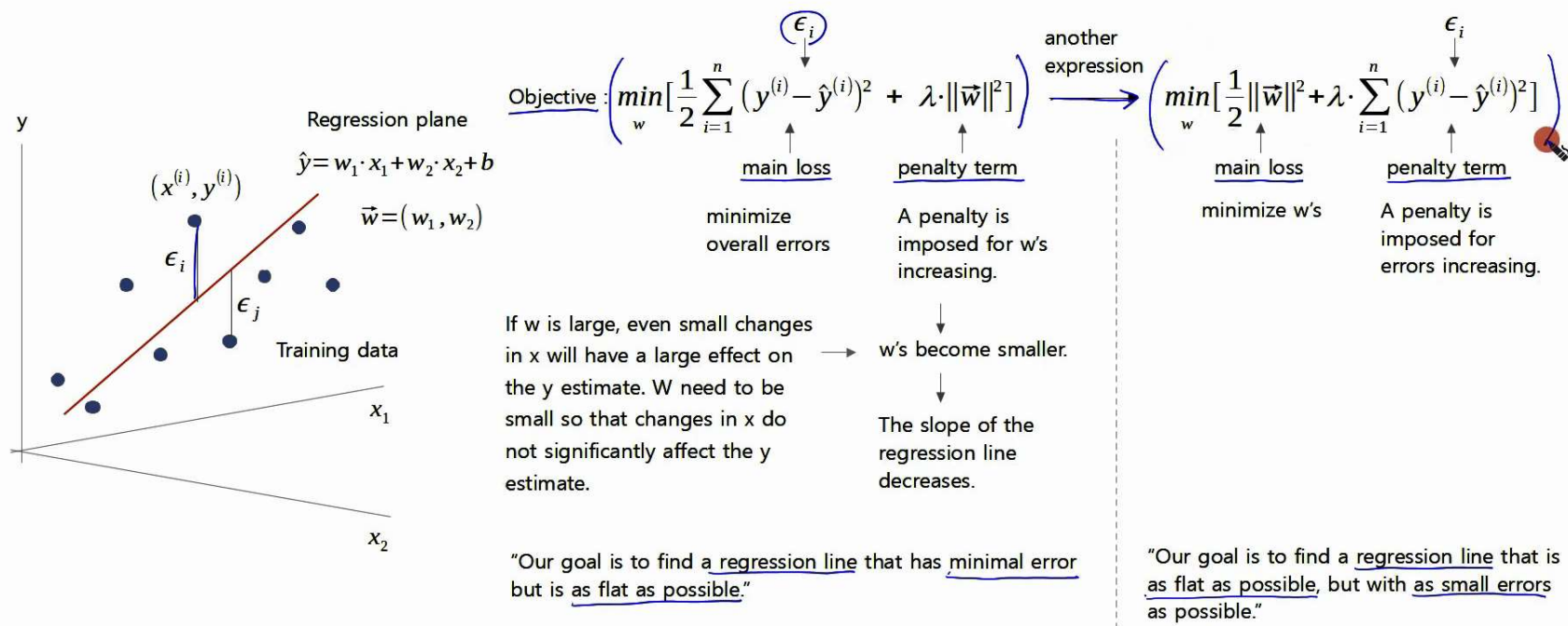


Support Vector Regression (SVR)

Friday 24 May 2024 1:19 AM

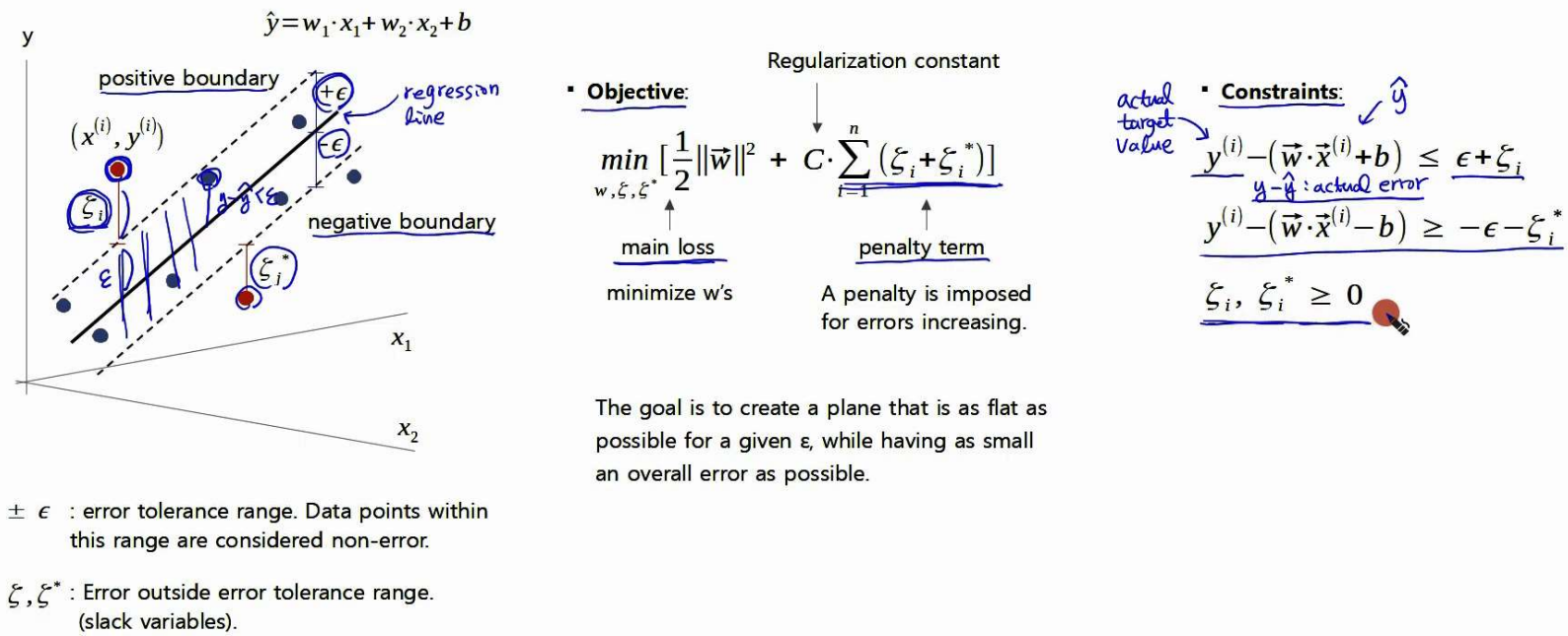
Support Vector Regression (SVR) : Regularized Linear Regression

- Regularized linear regression requires two goals to be considered. One is to minimize the overall errors, and the other is to make w's small. These two goals are a trade-off and must be balanced properly. To achieve both goals, you can use two expressions as follows.



Objective function for Support Vector regression

- To determine the optimal regression line (plane) for a given ϵ (error tolerance range), set the objective function and constraints as follows.



■ Lagrange primal and dual function.

■ Constrained optimization problem ■ Lagrange primal function

$$\begin{aligned} \min & \left[\frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) \right] \\ \text{s.t.} & \begin{cases} y^{(i)} - \vec{w} \cdot \vec{x}^{(i)} - b - \epsilon - \xi_i \leq 0 \\ -y^{(i)} + \vec{w} \cdot \vec{x}^{(i)} + b - \epsilon - \xi_i^* \leq 0 \\ -\xi_i \leq 0 \\ -\xi_i^* \leq 0 \end{cases} \end{aligned}$$

$$L_p = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) + \sum_{i=1}^n (-\eta_i \xi_i - \eta_i^* \xi_i^*) + \sum_{i=1}^n \lambda_i (y^{(i)} - \vec{w} \cdot \vec{x}^{(i)} - b - \epsilon - \xi_i) + \sum_{i=1}^n \lambda_i^* (-y^{(i)} + \vec{w} \cdot \vec{x}^{(i)} + b - \epsilon - \xi_i^*)$$

$$\frac{\partial L_p}{\partial \vec{w}} = \vec{w} - \sum_{i=1}^n (\lambda_i - \lambda_i^*) \vec{x}^{(i)} = 0$$

$$\frac{\partial L_p}{\partial \xi_i} = (C - \eta_i - \lambda_i) = 0 \rightarrow \lambda_i \leq C$$

$$\frac{\partial L_p}{\partial \xi_i^*} = (C - \eta_i^* - \lambda_i^*) = 0 \rightarrow \lambda_i^* \leq C$$

$$\frac{\partial L_p}{\partial b} = \left(\sum_{i=1}^n (\lambda_i - \lambda_i^*) = 0 \right)$$

$$\frac{\partial L_p}{\partial \xi_i} = (C - \eta_i - \lambda_i) = 0 \rightarrow \lambda_i \leq C$$

$$\frac{\partial L_p}{\partial \xi_i^*} = (C - \eta_i^* - \lambda_i^*) = 0 \rightarrow \lambda_i^* \leq C$$

$$L_p = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) + \sum_{i=1}^n (-\eta_i \xi_i - \eta_i^* \xi_i^*) + \sum_{i=1}^n \lambda_i (y^{(i)} - \vec{w} \cdot \vec{x}^{(i)} - b - \epsilon - \xi_i) + \sum_{i=1}^n \lambda_i^* (-y^{(i)} + \vec{w} \cdot \vec{x}^{(i)} + b - \epsilon - \xi_i^*)$$

■ Lagrange dual function

$$L_D = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\lambda_i - \lambda_i^*) (\lambda_j - \lambda_j^*) \vec{x}^{(i)} \cdot \vec{x}^{(j)} - \epsilon \sum_{i=1}^n (\lambda_i + \lambda_i^*) + \sum_{i=1}^n y^{(i)} (\lambda_i - \lambda_i^*)$$

$$\text{s.t.} \begin{cases} \sum_{i=1}^n (\lambda_i - \lambda_i^*) = 0 \\ 0 \leq \lambda_i \leq C \\ 0 \leq \lambda_i^* \leq C \end{cases} \rightarrow \begin{cases} \lambda_i \geq 0, \lambda_i^* \leq C \\ \lambda_i^* \geq 0, \lambda_i \leq C \end{cases}$$

■ Decision function

$$\vec{w} = \sum_{i=1}^n (\lambda_i - \lambda_i^*) \vec{x}^{(i)}$$

$$\hat{y} = \vec{w} \cdot \vec{x} + b \leftarrow \text{Decision function (b will be calculated later)}$$

* Source : Alex J. Smola, Bernhard Scholkopf, 1998/2004, A tutorial on support vector regression

■ Dual function and Quadratic Programming (QP)

■ Dual function

$$L_D = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\lambda_i - \lambda_i^*) (\lambda_j - \lambda_j^*) \vec{x}^{(i)} \cdot \vec{x}^{(j)} - \epsilon \sum_{i=1}^n (\lambda_i + \lambda_i^*) + \sum_{i=1}^n y^{(i)} (\lambda_i - \lambda_i^*)$$

$$\text{subject to} \begin{cases} \sum_{i=1}^n (\lambda_i - \lambda_i^*) = 0 \\ -\lambda_i \leq 0 & \lambda_i \leq C \\ -\lambda_i^* \leq 0 & \lambda_i^* \leq C \end{cases}$$

$$\underset{\lambda, \lambda^*}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\lambda_i \lambda_j - \lambda_i \lambda_j^* - \lambda_i^* \lambda_j + \lambda_i^* \lambda_j^*) \vec{x}^{(i)} \cdot \vec{x}^{(j)} + \sum_{i=1}^n (\epsilon - y^{(i)}) \lambda_i + \sum_{i=1}^n (\epsilon + y^{(i)}) \lambda_i^*$$

■ Matrices : if (n = 2), i = (1, 2), j = (1, 2)

$$K = \begin{bmatrix} \vec{x}^{(1)} \cdot \vec{x}^{(1)} & \vec{x}^{(1)} \cdot \vec{x}^{(2)} \\ \vec{x}^{(2)} \cdot \vec{x}^{(1)} & \vec{x}^{(2)} \cdot \vec{x}^{(2)} \end{bmatrix}$$

$$P = \begin{bmatrix} K & -K \\ -K & K \end{bmatrix}$$

$$q = \begin{bmatrix} \epsilon - y^{(1)} \\ \epsilon - y^{(2)} \\ \epsilon + y^{(1)} \\ \epsilon + y^{(2)} \end{bmatrix}$$

$$A = [1, 1, -1, -1] \quad b = 0$$

$$G\alpha = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_1^* \\ \lambda_2^* \\ -\lambda_1 \\ -\lambda_2 \\ -\lambda_1^* \\ -\lambda_2^* \end{bmatrix}$$

$$h = \begin{bmatrix} C \\ C \\ C \\ C \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{(Standard form of QP)}$$

$$\underset{\alpha}{\operatorname{argmin}} \frac{1}{2} \alpha^T \cdot P \cdot \alpha + q^T \cdot \alpha$$

$$\text{subject to} \begin{cases} A\alpha = b \\ G\alpha \leq h \end{cases}$$

■ Computing b

$$\text{■ Primal function : } L_p = \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) + \sum_{i=1}^n (-\eta_i \xi_i - \eta_i^* \xi_i^*) + \sum_{i=1}^n \lambda_i (y^{(i)} - \vec{w} \cdot \vec{x}^{(i)} - b - \epsilon - \xi_i) + \sum_{i=1}^n \lambda_i^* (-y^{(i)} + \vec{w} \cdot \vec{x}^{(i)} + b - \epsilon - \xi_i^*)$$

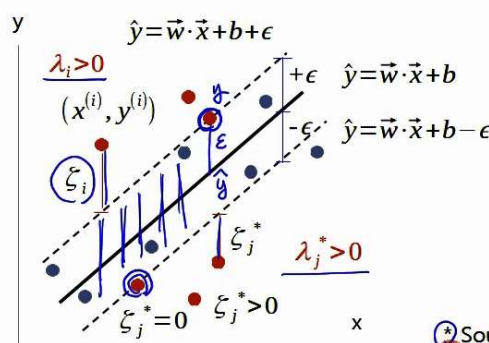
■ KKT condition : complementary slackness

$$\begin{aligned} (1) \quad \lambda_i (y^{(i)} - \vec{w} \cdot \vec{x}^{(i)} - b - \epsilon - \xi_i) &= 0 & \leftarrow \text{If } \lambda_i > 0 \text{ and } \xi_i = 0, \text{ the data point is on the positive boundary.} \\ (2) \quad \lambda_i^* (-y^{(i)} + \vec{w} \cdot \vec{x}^{(i)} + b - \epsilon - \xi_i^*) &= 0 & \leftarrow \text{If } \lambda_i^* > 0 \text{ and } \xi_i^* = 0, \text{ the data point is on the negative boundary.} \\ (3) \quad \eta_i \xi_i = 0 \rightarrow (C - \lambda_i) \xi_i = 0 & & \leftarrow \text{If } \lambda_i = C, \text{ then } \xi_i > 0 \text{ and the data point is outside the positive boundary.} \\ (4) \quad \eta_i^* \xi_i^* = 0 \rightarrow (C - \lambda_i^*) \xi_i^* = 0 & & \leftarrow \text{If } \lambda_i^* = C, \text{ then } \xi_i^* > 0 \text{ and the data point is outside the negative boundary.} \end{aligned}$$

① computing b

- $0 \leq \lambda_i \leq C$ and $0 \leq \lambda_i^* \leq C$
- In conditions 3) and 4), if λ_i and λ_i^* are less than C , then ξ_i and ξ_i^* must be 0. All of these data points lie between the positive and negative boundaries.
- The data points that satisfy the conditions $0 < \lambda_i < C$ and $0 < \lambda_i^* < C$ have ξ_i and ξ_i^* equal to 0.
- The data points that satisfy conditions 1) and 2) are on the positive and negative boundaries.

$$\begin{aligned} (1) \quad 0 < \lambda_i < C & \rightarrow \xi_i = 0 \rightarrow y^{(i)} - \vec{w} \cdot \vec{x}^{(i)} - b - \epsilon = 0 \\ (2) \quad 0 < \lambda_i^* < C & \rightarrow \xi_i^* = 0 \rightarrow -y^{(i)} + \vec{w} \cdot \vec{x}^{(i)} + b - \epsilon = 0 \end{aligned}$$



$$\begin{aligned} (1) \quad b &= y^{(i)} - \vec{w} \cdot \vec{x}^{(i)} - \epsilon : \text{For the data points (x, y) that satisfy the condition } 0 < \lambda_i < C. \\ (2) \quad b &= y^{(i)} - \vec{w} \cdot \vec{x}^{(i)} + \epsilon : \text{For the data points (x, y) that satisfy the condition } 0 < \lambda_i^* < C. \end{aligned}$$

* Take the average of b.

* Source : Alex J. Smola, Bernhard Scholkopf, 1998, A tutorial on support vector regression. 1.4 Computing b, 식 (13)

■ Non-linear regression : Kernel trick [MXML-6-05]

■ Primal Lagrange

$$L_p = \frac{1}{2} \|\tilde{w}\|^2 + C \sum_{i=1}^n (\zeta_i + \zeta_i^*) + \sum_{i=1}^n (-\eta_i \zeta_i - \eta_i^* \zeta_i^*) + \sum_{i=1}^n \lambda_i (y^{(i)} - \tilde{w} \cdot \phi(\tilde{x}^{(i)}) - b - \epsilon - \zeta_i) + \sum_{i=1}^n \lambda_i^* (-y^{(i)} + \tilde{w} \cdot \phi(\tilde{x}^{(i)}) + b - \epsilon - \zeta_i^*) \quad (\eta_i, \eta_i^*, \lambda_i, \lambda_i^* \geq 0)$$

just replace $\tilde{x}^{(i)}$ with $\phi(\tilde{x}^{(i)})$

$$\frac{\partial L_p}{\partial \tilde{w}} = \tilde{w} - \sum_{i=1}^n (\lambda_i - \lambda_i^*) \phi(\tilde{x}^{(i)}) = 0 \quad \frac{\partial L_p}{\partial b} = \left(\sum_{i=1}^n (\lambda_i^* - \lambda_i) = 0 \right) \quad \frac{\partial L_p}{\partial \zeta_i} = (C - \eta_i - \lambda_i = 0) \rightarrow \lambda_i \leq C \quad \frac{\partial L_p}{\partial \zeta_i^*} = (C - \eta_i^* - \lambda_i^* = 0) \rightarrow \lambda_i^* \leq C$$

$$L_p = \frac{1}{2} \|\tilde{w}\|^2 + C \sum_{i=1}^n (\cancel{\zeta_i} + \cancel{\zeta_i^*}) + \sum_{i=1}^n (-\cancel{\eta_i \zeta_i} - \cancel{\eta_i^* \zeta_i^*}) + \sum_{i=1}^n \lambda_i (y^{(i)} - \tilde{w} \cdot \phi(\tilde{x}^{(i)}) - \cancel{b} - \epsilon - \cancel{\zeta_i}) + \sum_{i=1}^n \lambda_i^* (-y^{(i)} + \tilde{w} \cdot \phi(\tilde{x}^{(i)}) + \cancel{b} - \epsilon - \cancel{\zeta_i^*})$$

■ Linear SVR

$$\tilde{w} = \sum_{i=1}^n (\lambda_i - \lambda_i^*) \tilde{x}^{(i)} \\ b = y^{(i)} - \tilde{w} \cdot \tilde{x} \pm \epsilon$$

■ Dual Lagrange

$\tilde{x}^{(i)} \cdot \tilde{x}^{(j)}$ was replaced by $\phi(\tilde{x}^{(i)}) \cdot \phi(\tilde{x}^{(j)})$

$$L_D = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\lambda_i - \lambda_i^*) (\lambda_j - \lambda_j^*) \phi(\tilde{x}^{(i)}) \cdot \phi(\tilde{x}^{(j)}) - \epsilon \sum_{i=1}^n (\lambda_i + \lambda_i^*) + \sum_{i=1}^n y^{(i)} (\lambda_i - \lambda_i^*)$$

$$\text{s.t.} \quad \begin{cases} \sum_{i=1}^n (\lambda_i - \lambda_i^*) = 0 \\ 0 \leq \lambda_i \leq C \\ 0 \leq \lambda_i^* \leq C \end{cases} \rightarrow \begin{cases} \lambda_i \geq 0, \lambda_i \leq C \\ \lambda_i^* \geq 0, \lambda_i^* \leq C \end{cases}$$

We can find \hat{y} because we can find $\phi(\tilde{x}^{(i)}) \cdot \phi(\tilde{x}^{(j)})$.

■ Decision function

$$\tilde{w} = \sum_{i=1}^n (\lambda_i - \lambda_i^*) \phi(\tilde{x}^{(i)}) \quad \leftarrow \text{We cannot find } w \text{ because we don't know } \phi(\tilde{x}^{(i)})$$

$$\tilde{w} \cdot \phi(\tilde{x}) = \sum_{i=1}^n (\lambda_i - \lambda_i^*) \phi(\tilde{x}^{(i)}) \cdot \phi(\tilde{x}) \quad b = y^{(i)} - \tilde{w} \cdot \phi(\tilde{x}) \pm \epsilon$$

$$\hat{y} = \sum_{i=1}^n (\lambda_i - \lambda_i^*) \phi(\tilde{x}^{(i)}) \cdot \phi(\tilde{x}_{\text{test}}) + b \quad \leftarrow \text{Decision function}$$

$$\hat{y} = \sum_{i=1}^n (\lambda_i - \lambda_i^*) k(\tilde{x}^{(i)}, \tilde{x}_{\text{test}}) + b$$

The method for calculating b is similar to linear SVM.

* Source : Alex J. Smola, Bernhard Scholkopf, 1998/2004, A tutorial on support vector regression

■ Dual function

$$L_D = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\lambda_i - \lambda_i^*) (\lambda_j - \lambda_j^*) \phi(\tilde{x}^{(i)}) \cdot \phi(\tilde{x}^{(j)}) - \epsilon \sum_{i=1}^n (\lambda_i + \lambda_i^*) + \sum_{i=1}^n y^{(i)} (\lambda_i - \lambda_i^*)$$

subject to $\sum_{i=1}^n (\lambda_i - \lambda_i^*) = 0$

$$\underset{\lambda, \lambda^*}{\text{argmin}} \quad \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (\lambda_i \lambda_j - \lambda_i \lambda_j^* - \lambda_i^* \lambda_j + \lambda_i^* \lambda_j^*) \phi(\tilde{x}^{(i)}) \cdot \phi(\tilde{x}^{(j)}) + \sum_{i=1}^n (\epsilon - y^{(i)}) \lambda_i + \sum_{i=1}^n (\epsilon + y^{(i)}) \lambda_i^*$$

$$\begin{aligned} -\lambda_i &\leq 0 & \lambda_i &\leq C \\ -\lambda_i^* &\leq 0 & \lambda_i^* &\leq C \end{aligned}$$

■ Matrices : if $n = 2$, $i = (1, 2)$, $j = (1, 2)$

$$k(\tilde{x}^{(i)}, \tilde{x}^{(j)}) = \phi(\tilde{x}^{(i)}) \cdot \phi(\tilde{x}^{(j)})$$

$$K = \begin{bmatrix} k(\tilde{x}^{(1)}, \tilde{x}^{(1)}) & k(\tilde{x}^{(1)}, \tilde{x}^{(2)}) \\ k(\tilde{x}^{(2)}, \tilde{x}^{(1)}) & k(\tilde{x}^{(2)}, \tilde{x}^{(2)}) \end{bmatrix}$$

$$P = \begin{bmatrix} K & -K \\ -K & K \end{bmatrix}$$

$$q = \begin{bmatrix} \epsilon - y^{(1)} \\ \epsilon - y^{(2)} \\ \epsilon + y^{(1)} \\ \epsilon + y^{(2)} \end{bmatrix}$$

$$A = [1, 1, -1, -1] \quad b = 0$$

$$G = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \end{bmatrix}$$

$$G\alpha = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_1^* \\ \lambda_2^* \\ -\lambda_1 \\ -\lambda_2 \\ -\lambda_1^* \\ -\lambda_2^* \end{bmatrix}$$

$$h = \begin{bmatrix} C \\ C \\ C \\ C \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

■ (Standard form of QP)

$$\underset{\alpha}{\text{argmin}} \quad \frac{1}{2} \alpha^T \cdot P \cdot \alpha + q^T \cdot \alpha$$

$$\text{subject to} \quad A\alpha = b$$

$$G\alpha \leq h$$