

서포트 벡터 학습을 이용한 함수근사
(Function Approximation via Support Vector Regression)

2004년 10월 30일
한국 퍼지 및 지능시스템 학회 2004년 추계학술대회 발표자료

고려대학교 과학기술대학 제어계측공학과
박주영

Korea University

Jooyoung Park

Function Approximation via SVR

October 2004

1. Introduction to support vector learning

- ▶ We may use the support vector learning algorithm to construct the following:
 - RBF networks
 - Two-layer perceptrons (i.e. with a single hidden layer)

- ▶ The support vectors consist of a small subset of the training data extracted by the algorithm.

2. Introduction to support vector regression(SVR)

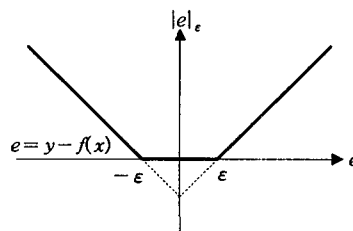
- ▶ Training data set:

$\{(x_i, y_i)\}_{i=1}^m$, where x_i : the input pattern, y_i : target output

- ▶ Epsilon-insensitive loss function:

$$|e|_\epsilon = \max(0, |e| - \epsilon)$$

$$|y - f(x)|_\epsilon = \max(0, |y - f(x)| - \epsilon)$$

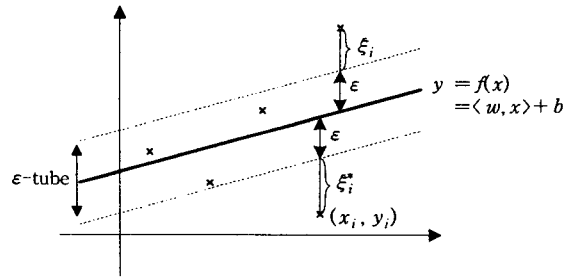


- ▶ The SVR utilizes a linear approximator $f(x) = w^T x + b$, which deviates least from the training data according to the ϵ -insensitive loss function.

- ▶ More precisely, the objective of ϵ -SVR is to find the "flat" function $f(x) = \langle w, x \rangle + b$ while keeping the approximation error $|y_i - f(x_i)|_\epsilon$ small.

- ▶ **Note:** This approximation problem can be formulated as an optimization problem, and its global min. can be found because

- ① its objective function is convex, and
- ② its constraints are linear in variables.



- Mathematical formulation for the SVR: Given the training data set $\{(x_i, y_i)\}_{i=1}^m$, solve the following quadratic problem:

$$\begin{aligned} \min \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \\ \text{s. t.} \quad & y_i - (\langle w, x_i \rangle + b) \leq \xi_i + \epsilon \\ & (\langle w, x_i \rangle + b) - y_i \leq \epsilon + \xi_i^* \\ & \xi_i, \xi_i^* \geq 0, \quad i = 1, \dots, m \end{aligned}$$

3. A brief overview of optimization theory

- **Theorem:** $f \in C^1$ has a min. at $x^* \Rightarrow \frac{\partial f}{\partial x}(x^*) = 0$.

This condition, together with convexity of f , is also a suff. cond.

- Example 1: $\min. f(x) = \frac{1}{2}(x_1^2 + x_2^2)$

Solution:

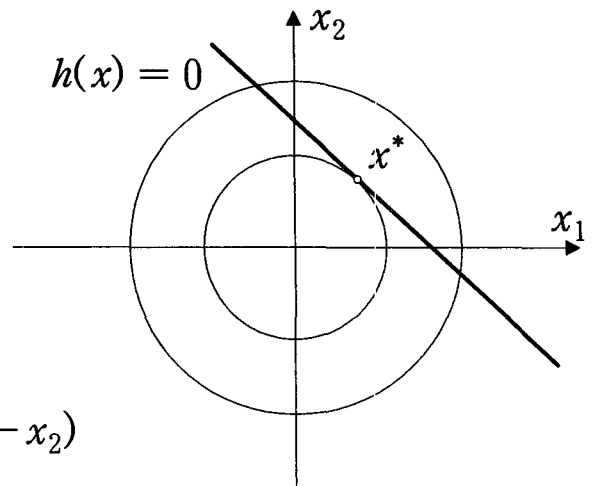
$$\begin{aligned} \frac{\partial f}{\partial x} = 0 & \Rightarrow \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \right] = [x_1 \quad x_2] = 0 \\ \therefore x^* &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} \end{aligned}$$

- In a constrained min. problem,

$$f \in C^1 \text{ has a min. at } x^* \Rightarrow \frac{\partial f}{\partial x}(x^*) = 0$$

► Example 2:

$$\begin{aligned} \min . f(x) &= \frac{1}{2}(x_1^2 + x_2^2) \\ \text{s.t. } h(x) &= 1 - x_1 - x_2 = 0 \end{aligned}$$



Solution: Define the Lagrange function

$$\begin{aligned} L(x, \lambda) &= f(x) + \lambda h(x) \\ &= \frac{1}{2}(x_1^2 + x_2^2) + \lambda(1 - x_1 - x_2) \end{aligned}$$

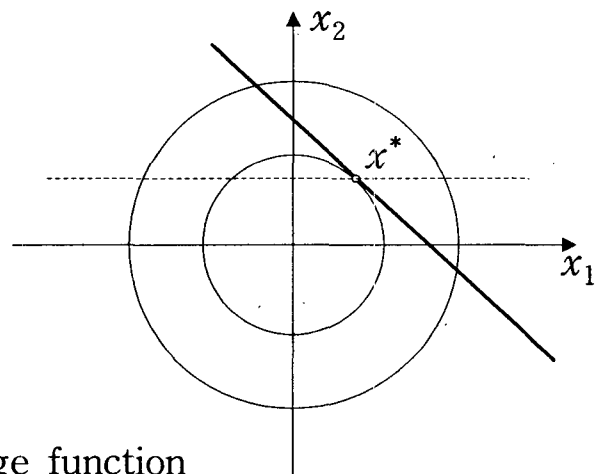
$$\frac{\partial L}{\partial x} = 0 \Rightarrow [x_1 - \lambda \quad x_2 - \lambda] = 0. \therefore x_1 = x_2 = \lambda.$$

$$\frac{\partial L}{\partial \lambda} = 0 \Rightarrow 1 - x_1 - x_2 = 0. \therefore 1 - 2\lambda = 0. \therefore \lambda = \frac{1}{2}.$$

$$\therefore x^* = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$

► Example 3:

$$\begin{aligned} \min . f(x) &= \frac{1}{2}(x_1^2 + x_2^2) \\ \text{s.t. } h(x) &= 1 - x_1 - x_2 = 0, \\ g(x) &= \frac{3}{4} - x_2 \leq 0 \end{aligned}$$



Solution: Define the generalized Lagrange function

$$\begin{aligned} L(x, \lambda, \alpha) &\triangleq f + \lambda h + \alpha g \\ &= \frac{1}{2}(x_1^2 + x_2^2) + \lambda(1 - x_1 - x_2) + \alpha\left(\frac{3}{4} - x_2\right), \quad \alpha \geq 0 \end{aligned}$$

$$\frac{\partial L}{\partial x} = 0 \Rightarrow [x_1 - \lambda, x_2 - \lambda - \alpha] = 0. \therefore x_1 = \lambda, x_2 = \lambda + \alpha$$

$$\frac{\partial L}{\partial \lambda} = 0 \Rightarrow 1 - x_1 - x_2 = 0. \therefore 2\lambda + \alpha = 1$$

Also, $\alpha \geq 0$ and $\frac{3}{4} - x_2 \leq 0$.

One more condition is needed to solve the problem.

→ The Kuhn-Tucker complementarity condition

$$\alpha \left(\frac{3}{4} - x_2 \right) = 0 \quad \text{i.e., } \alpha = 0 \quad \text{or} \quad x_2 = \frac{3}{4}$$

$$\textcircled{1} \text{ If } \alpha = 0, \text{ then } \lambda = \frac{1}{2}; \text{ thus } x_1 = x_2 = \frac{1}{2} \quad \otimes \quad (\because x_2 \geq \frac{3}{4})$$

$$\textcircled{2} \text{ If } x_2 = \frac{3}{4}, \text{ then } \begin{cases} \lambda + \alpha = \frac{3}{4} \\ 2\lambda + \alpha = 1 \end{cases} \therefore \begin{cases} \lambda = \frac{1}{4}, \alpha = \frac{1}{2} \\ x_1 = \frac{1}{4}, x_2 = \frac{3}{4} \end{cases}$$

$$\therefore x^* = \begin{bmatrix} 1/4 \\ 3/4 \end{bmatrix}$$

► **Theorem** (Kuhn-Tucker Theorem)

Given an opt. prob. with convex domain $\Omega \subseteq R^n$

$$\left. \begin{array}{l} \min f(x), x \in \Omega \text{ (} x \text{ is primal variable)} \\ \text{s.t. } g_i(x) \leq 0, i = 1, \dots, k \\ h_j(x) = 0, j = 1, \dots, m \end{array} \right\} \text{primal opt. prob } \star$$

with $f \in C^1$ convex, and g_i, h_j affine, the following are necessary and sufficient condition for a point $x^* \in \Omega$ to be an opt.:

$$\text{For } L(x, \alpha, \lambda) \triangleq f(x) + \sum_{i=1}^k \alpha_i g_i(x) + \sum_{j=1}^m \lambda_j h_j(x) = f + \alpha^T g + \lambda^T h,$$

$$\exists \alpha^* \text{ and } \lambda^* \text{ s.t. } \frac{\partial L}{\partial x}(x^*, \alpha^*, \lambda^*) = 0, \quad \frac{\partial L}{\partial \lambda}(x^*, \alpha^*, \lambda^*) = 0$$

$$g_i(x^*) \leq 0 \text{ and } \alpha_i^* \geq 0 \text{ for } i = 1, \dots, k,$$

$$\text{and } \alpha_i^* g_i(x^*) = 0, \quad i = 1, \dots, k$$

↳ The Kuhn-Tucker complementarity condition

► **Definition:** The (Lagrangian) dual prob. of ★ is the following:

$$\begin{cases} \max & \Theta(\alpha, \lambda) \quad (\text{where } \alpha \text{ and } \lambda \text{ are dual variables}) \\ \text{s.t.} & \alpha \geq 0 \end{cases}$$

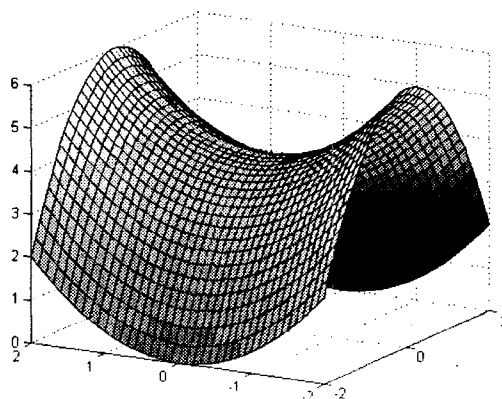
$$\text{where } \Theta(\alpha, \lambda) \triangleq \inf_{x \in \Omega} L(x, \alpha, \lambda)$$

Dual problem of ★ :

$$\begin{cases} \max & \Theta(\alpha, \lambda) \\ \text{s.t.} & \alpha \geq 0 \end{cases}$$

► **Remarks:**

- ① In the course, our attention is focused on convex quadratic programs. (목적함수-Convex 2차식, 제약조건-1차식)
- ② It can be shown that for convex quadratic programs, L has a saddle point w.r.t. the primal and dual variables at the opt. soln.



4. Back to the SVR

- ▶ To find the best linear approximator for the training data set $\{(x_i, y_i)\}_{i=1}^m$, we need to solve the following problem:

$$\begin{aligned} \min \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \\ \text{s.t.} \quad & y_i - (\langle w, x_i \rangle + b) \leq \xi_i + \varepsilon, \\ & (\langle w, x_i \rangle + b) - y_i \leq \varepsilon + \xi_i^*, \quad \xi_i, \xi_i^* \geq 0, \quad i = 1, \dots, m \end{aligned}$$

- ▶ Lagrange Function

$$\begin{aligned} L = \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) + \sum_{i=1}^m \alpha_i [y_i - (\langle w, x_i \rangle + b) - \varepsilon - \xi_i] \\ & + \sum_{i=1}^m \alpha_i^* [(\langle w, x_i \rangle + b) - y_i - \varepsilon - \xi_i^*] - \sum_{i=1}^m (\eta_i \xi_i + \eta_i^* \xi_i^*) \\ & \alpha_i^{(*)}, \eta_i^{(*)} \geq 0 \quad / \quad \text{Primal variable: } w, b, \xi, \xi^* \quad / \quad \text{Dual variable: } \alpha, \alpha^*, \eta, \eta^* \end{aligned}$$

- ▶ The saddle point condition:

$$\frac{\partial L}{\partial w} = 0 \quad \Leftrightarrow \quad w - \sum_{i=1}^m \alpha_i x_i + \sum_{i=1}^m \alpha_i^* x_i = 0 \quad \Leftrightarrow \quad w = \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i$$

$$\frac{\partial L}{\partial b} = 0 \quad \Leftrightarrow \quad \sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0$$

$$\frac{\partial L}{\partial \xi_i^{(*)}} = 0 \quad \Leftrightarrow \quad C - \alpha_i^{(*)} - \eta_i^{(*)} = 0 \quad \Leftrightarrow \quad \alpha_i^{(*)} + \eta_i^{(*)} = C \quad \therefore \alpha_i^{(*)} \in [0, C]$$

- ▶ Substitute the above into L to remove all primal variable $(w, b, \xi_i^{(*)})$
- ▶ Dual Problem:

$$\begin{aligned} \max \quad D = \quad & -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) \langle x_i, x_j \rangle \\ & + \sum_{i=1}^m (\alpha_i - \alpha_i^*) y_i - \sum_{i=1}^m (\alpha_i + \alpha_i^*) \varepsilon \\ \text{s.t.} \quad & \sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0, \quad \alpha_i^{(*)} \in [0, C], \quad i = 1, \dots, m \end{aligned}$$

$$\triangleright w = \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i$$

$$\therefore f(x) = \langle w, x \rangle + b = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \langle x_i, x \rangle + b$$

► The KT condition

$$\textcircled{1} \alpha_i [y_i - (\langle w, x_i \rangle + b) - \varepsilon - \xi_i] = 0, \quad \forall i$$

$$\textcircled{2} \alpha_i^* [(\langle w, x_i \rangle + b) - y_i - \varepsilon - \xi_i^*] = 0, \quad \forall i$$

$$\textcircled{3} \eta_i \xi_i = 0 \text{ i.e., } (C - \alpha_i) \xi_i = 0, \quad \forall i$$

$$\textcircled{4} \eta_i^* \xi_i^* = 0 \text{ i.e., } (C - \alpha_i^*) \xi_i^* = 0, \quad \forall i$$

► Remark:

$$\textcircled{1} \alpha_i \alpha_i^* = 0 \quad \because \alpha_i \neq 0 \Rightarrow y_i - (\langle w, x_i \rangle + b) - \varepsilon - \xi_i = 0$$

$$\Rightarrow y_i - (\langle w, x_i \rangle + b) = \varepsilon + \xi_i$$

$$\Rightarrow \frac{(\langle w, x_i \rangle + b) - y_i - \varepsilon - \xi_i}{= -\varepsilon - \xi_i} < 0$$

$$\Rightarrow \alpha_i^* = 0$$

Similarly for the other way

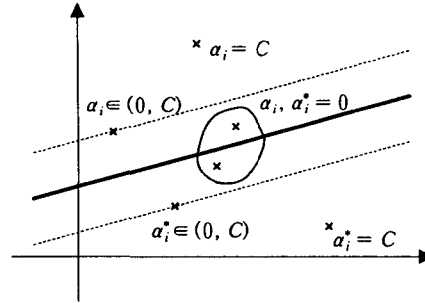
$$\textcircled{2} \alpha_i \in (0, C) \Rightarrow \xi_i = 0 \quad \text{and} \quad y_i - (\langle w, x_i \rangle + b) = \varepsilon$$

$$\textcircled{3} \alpha_i^* \in (0, C) \Rightarrow \xi_i^* = 0 \quad \text{and} \quad (\langle w, x_i \rangle + b) - y_i = \varepsilon$$

$$\textcircled{4} \xi_i > 0 \Rightarrow \alpha_i = C$$

$$\textcircled{5} \xi_i^* > 0 \Rightarrow \alpha_i^* = C$$

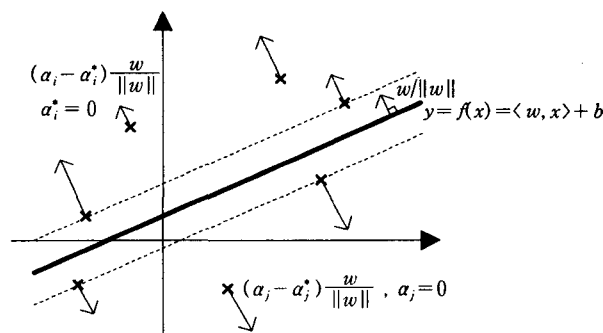
- ⑥ ϵ -tube 내부에 training data에 대응하는 $\alpha_i, \alpha_i^* = 0$.
 (\because the KT condition ①, ② 모두 두 번째항이 nonzero 이므로)



▶ Note:

$$\begin{aligned}
 \text{목적함수 } D = & -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \langle x_i, x_j \rangle \\
 & + \sum_{i=1}^m (\alpha_i - \alpha_i^*) y_i - \sum_{i=1}^m (\alpha_i + \alpha_i^*) \epsilon
 \end{aligned}$$

▶ Mechanical Interpretation



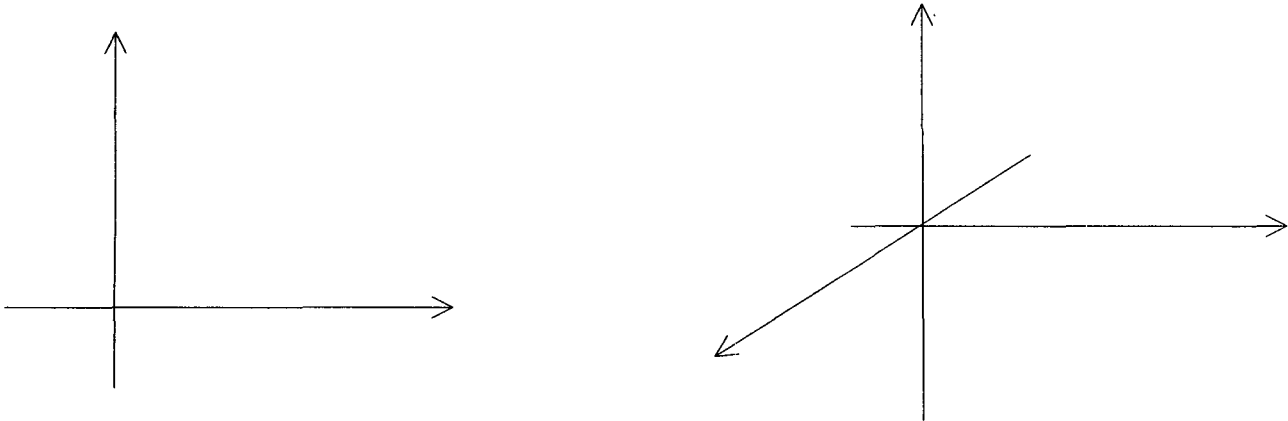
$$\sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0 : \text{Force balance (힘 균형)}$$

$$w = \sum_{i=1}^m (\alpha_i - \alpha_i^*) x_i ,$$

$$\therefore \sum_{i=1}^m x_i \times (\alpha_i - \alpha_i^*) \frac{w}{\|w\|} = w \times \frac{w}{\|w\|} = 0 : \text{Torque balance (회전 모멘트 균형)}$$

트 균형)

- An extension for nonlinear approximators:



- Feature space : Given training data $\{(x_i, y_i)\}_{i=1}^m$,

where $x_i \in R^l$, $y_i \in R$, preprocess the data with

$$\begin{aligned} \phi: R^l &\rightarrow F \\ x &\mapsto \phi(x), \end{aligned} \quad \text{where } l \ll \dim(F),$$

to get $(\phi(x_i), y_i) \in F \times R$, $i = 1, \dots, m$.

- Example:

$$\begin{aligned} \phi: R^2 &\rightarrow R^3 \\ (x_1, x_2) &\mapsto (x_1^2, \sqrt{2}x_1x_2, x_2^2) \end{aligned}$$

- The kernel trick

$$\begin{aligned} \langle \phi(x), \phi(y) \rangle &= (x_1^2, \sqrt{2}x_1x_2, x_2^2) \cdot (y_1^2, \sqrt{2}y_1y_2, y_2^2) \\ &= x_1^2y_1^2 + 2x_1x_2y_1y_2 + x_2^2y_2^2 \\ &= (x_1y_1 + x_2y_2)^2 \\ &= (\langle x, y \rangle)^2 \\ &= K(x, y) \end{aligned}$$

Note : $\langle \phi(x), \phi(y) \rangle = (\langle x, y \rangle)^2 = K(x, y)$ can be computed easily on the input space !!!

- Question: Which kernel function $K(x, y)$ has a corresponding inner product $\langle \phi(x), \phi(y) \rangle$ in some feature space?

Answer: Mercer kernels.

- **Mercer's Theorem** (1909): Roughly speaking,

$$\text{if } \int \int K(x, y) f(x) f(y) dx dy \geq 0 \text{ for } \forall f \in L_2,$$

$$\text{then } \exists \phi: R^l \rightarrow F \text{ s.t. } K(x, y) = \langle \phi(x), \phi(y) \rangle.$$

- Examples of Mercer kernels:

- RBF: $K(x, y) = \exp\left(-\frac{1}{2} \frac{\|x-y\|^2}{\sigma^2}\right)$

- MLP: $K(x, y) = \tanh(\gamma \langle x, y \rangle + \theta)$, where $\gamma, \theta > 0$

- Polynomial: $K(x, y) = (\langle x, y \rangle + c)^d$

- Note: Recall the note of p. 17.

- Kernel-based nonlinear approximators:

If we use the nonlinear approximator $f(x) = \langle w, \phi(x) \rangle + b$, we have

$$w = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \phi(x_i),$$

$$f(x) = \langle w, \phi(x) \rangle + b = \sum_{i=1}^m (\alpha_i - \alpha_i^*) k(x_i, x) + b$$

∴ Dual problem

$$\Leftrightarrow \max_{\alpha_i, \alpha_i^*} -\frac{1}{2} \sum_i^m \sum_j^m (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) k(x_i, x_j) \\ + \sum_i^m (\alpha_i - \alpha_i^*) y_i - \sum_i^m (\alpha_i + \alpha_i^*) \varepsilon$$

$$\text{s.t. } \sum_i^m (\alpha_i - \alpha_i^*) = 0, \alpha_i^{(*)} \in [0, C], i = 1, \dots, m$$

$$\Leftrightarrow \min_{\underline{\alpha}_i, \underline{\alpha}_i^*} \frac{1}{2} (\underline{\alpha} - \underline{\alpha}^*)^T \mathbf{K} (\underline{\alpha} - \underline{\alpha}^*) + \begin{bmatrix} \varepsilon \mathbf{1} \\ \varepsilon \mathbf{1} \end{bmatrix}^T \begin{bmatrix} \underline{\alpha} \\ \underline{\alpha}^* \end{bmatrix} - \begin{bmatrix} \underline{y} \\ -\underline{y} \end{bmatrix}^T \begin{bmatrix} \underline{\alpha} \\ \underline{\alpha}^* \end{bmatrix}$$

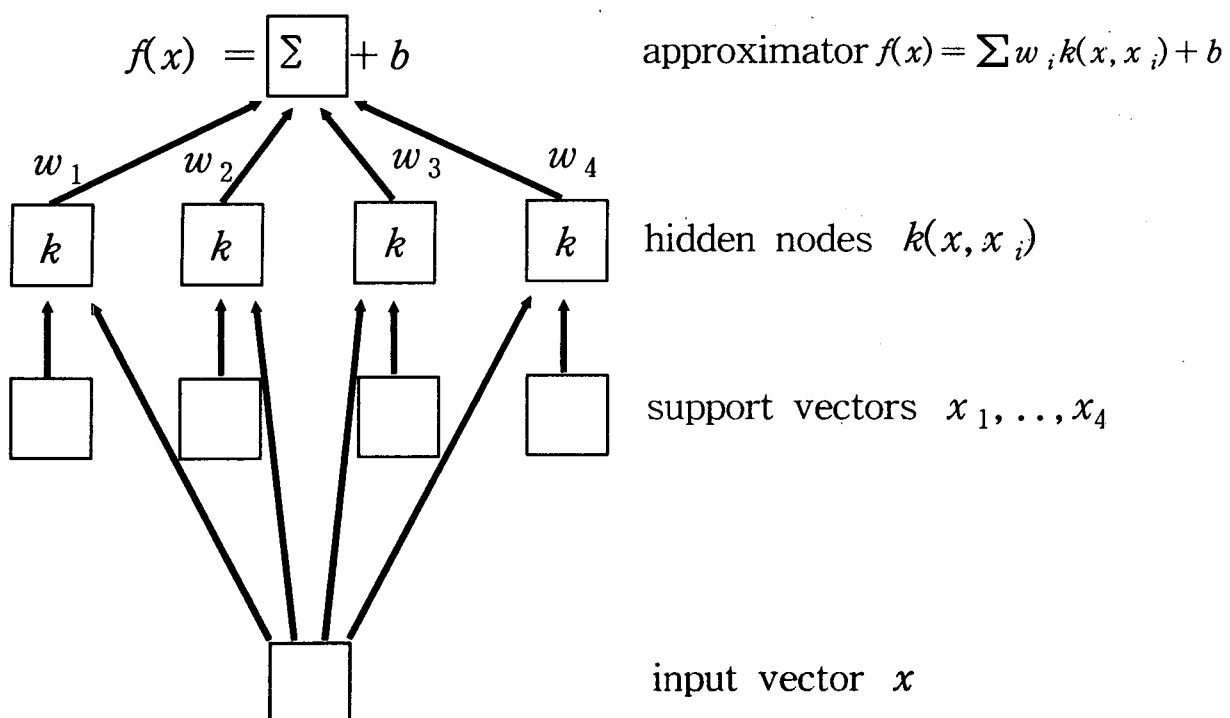
$$s.t. \begin{bmatrix} -\mathbf{1} \\ \mathbf{1} \end{bmatrix}^T \begin{bmatrix} \underline{\alpha} \\ \underline{\alpha}^* \end{bmatrix} = 0, \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \leq \begin{bmatrix} \underline{\alpha} \\ \underline{\alpha}^* \end{bmatrix} \leq \begin{bmatrix} C \\ \vdots \\ C \end{bmatrix}$$

$$\Leftrightarrow \min_{\underline{\alpha}_i, \underline{\alpha}_i^*} \frac{1}{2} \begin{bmatrix} \underline{\alpha} \\ \underline{\alpha}^* \end{bmatrix}^T \begin{bmatrix} K & -K \\ -K & K \end{bmatrix} \begin{bmatrix} \underline{\alpha} \\ \underline{\alpha}^* \end{bmatrix} + \begin{bmatrix} \varepsilon \mathbf{1} - \underline{y} \\ \varepsilon \mathbf{1} + \underline{y} \end{bmatrix}^T \begin{bmatrix} \underline{\alpha} \\ \underline{\alpha}^* \end{bmatrix}$$

$$s.t. \begin{bmatrix} -\mathbf{1} \\ \mathbf{1} \end{bmatrix}^T \begin{bmatrix} \underline{\alpha} \\ \underline{\alpha}^* \end{bmatrix} = 0, \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \leq \begin{bmatrix} \underline{\alpha} \\ \underline{\alpha}^* \end{bmatrix} \leq \begin{bmatrix} C \\ \vdots \\ C \end{bmatrix}$$

(MATLAB의 “QP solver quadprog”에 적합한 꼴)

5. The SVR architecture



6. Discussion

Given the data set $\{(x_i, y_i)\}_{i=1}^m$, we have the following kinds of training problems:

- ▶ Least square problem: w_i ?

$$f(x) = \sum_{i=1}^{m_1} w_i g(x; \mu_i, \sigma_i), \text{ where } m_1, \mu_i, \sigma_i \text{ are fixed}$$

- ▶ Learning via back-propagation: w_i, μ_i, σ_i ?

$$f(x) = \sum_{i=1}^{m_2} w_i g(x; \mu_i, \sigma_i), \text{ where } m_2 \text{ is fixed}$$

- ▶ Support vector learning: w_i ?

$$f(x) = \sum_{i=1}^m w_i g(x; x_i, \sigma_i), \text{ where } \sigma_i \text{ is fixed and } \#(\text{nonzero } w_i) \ll m$$

7. References

- [1] Bernhard Schölkopf, Alexander J. Smola, *Learning with Kernels*, MIT Press, 2002.
- [2] Nello Cristianini, John Shawe-Taylor, *An Introduction to Support Vector Machines*, Cambridge University Press, 2000.
- [3] <http://www.kernel-machines.org>
- [4] <http://www.isis.ecs.soton.ac.kr/resources/svminfo/>