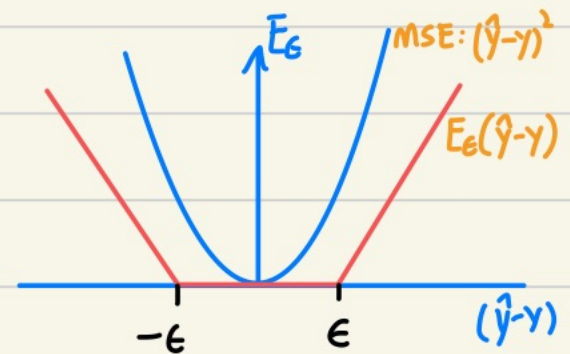


Support Vector Regression

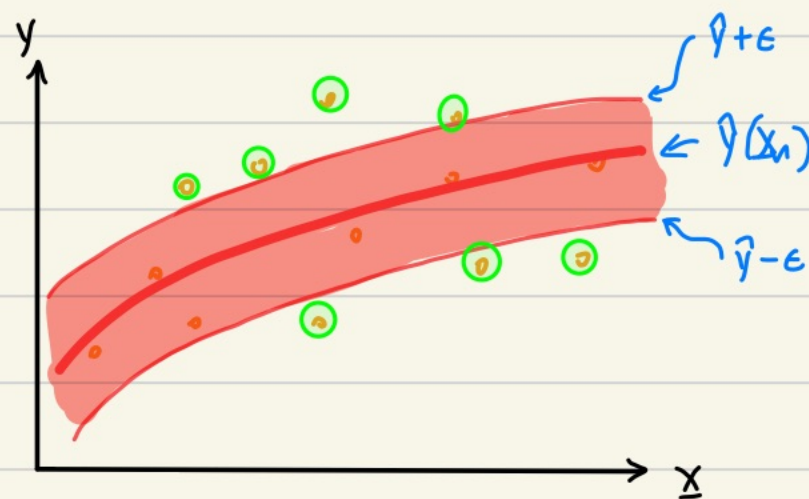
In Ridge Regression $J(\underline{w}) = \frac{1}{N} \|\underline{X}\underline{w} - \underline{y}\|_2^2 + \lambda \|\underline{w}\|_2^2$
 or $\frac{1}{N} \sum_{n=1}^N [\hat{y}(x_n) - y_n]^2 + \lambda \|\underline{w}\|_2^2 \quad (\lambda > 0)$

In SVR $J(\underline{w}) = C \cdot \sum_{n=1}^N E_\epsilon[\hat{y}(x_n) - y_n] + \frac{1}{2} \|\underline{w}\|_2^2$
 by convention

where $E_\epsilon(\hat{y} - y) = \begin{cases} 0 & \text{if } |\hat{y} - y| < \epsilon \\ |\hat{y} - y| - \epsilon & \text{if } |\hat{y} - y| \geq \epsilon \end{cases}$



Visual SVR:



Support vectors are those outside of tube

Similarly: add slack and constraints.

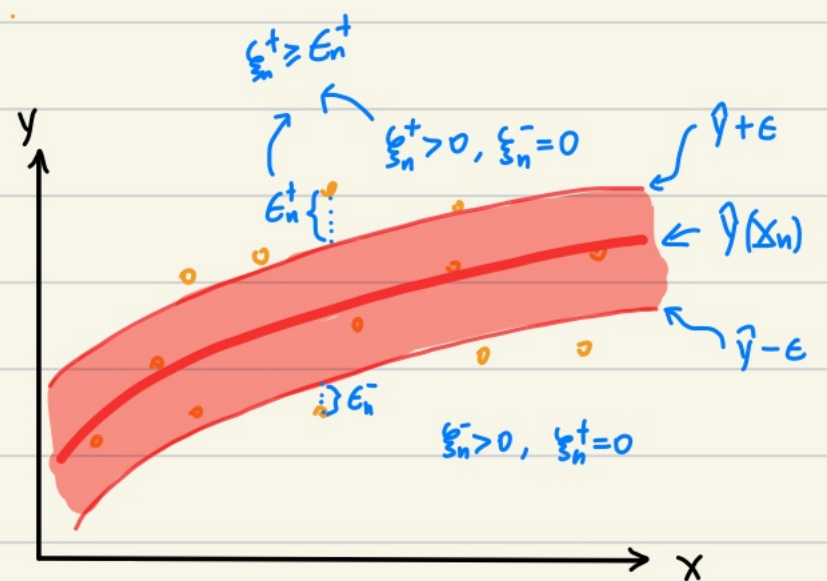
Inside tube: $\xi_n^+ = \xi_n^- = 0$

$\xi_n^+ > 0, \xi_n^- = 0$ iff $y_n > \hat{y} + \epsilon, \therefore y_n \leq \hat{y} + \epsilon + \xi_n^+$

$\xi_n^- > 0, \xi_n^+ = 0$ iff $y_n < \hat{y} - \epsilon, \therefore y_n \geq \hat{y} - \epsilon - \xi_n^-$

Now, $J(\underline{w}) = C \sum_{n=1}^N (\xi_n^+ + \xi_n^-) + \frac{1}{2} \|\underline{w}\|_2^2$

with constraints: $\begin{cases} \xi_n^+ \geq 0, \xi_n^- \geq 0 & \forall n \\ y_n \leq \hat{y}(x_n) + \epsilon + \xi_n^+ & \forall n \\ y_n \geq \hat{y}(x_n) - \epsilon - \xi_n^- & \forall n \end{cases}$



Derive Lagrange Optimization Equation (Primal form)

$$L(\underline{w}, w_0, \underline{\xi}^+, \underline{\xi}^-, \underline{\mu}^+, \underline{\mu}^-, \underline{\lambda}^+, \underline{\lambda}^-) = C \cdot \sum_{n=1}^N (\xi_n^+ + \xi_n^-) + \frac{1}{2} \|\underline{w}\|^2 - \sum_{n=1}^N (\underbrace{\mu_n^+}_{\text{multiplier}} \xi_n^+ + \underbrace{\mu_n^-}_{\text{multiplier}} \xi_n^-) \\ - \sum_{n=1}^N \underbrace{\lambda_n^+}_{\text{multiplier}} (\hat{y}_n + \epsilon + \xi_n^+ - y_n) - \sum_{n=1}^N \underbrace{\lambda_n^-}_{\text{multiplier}} (\epsilon + \xi_n^- - \hat{y}_n + y_n)$$

in which $\hat{y}_n = \hat{y}(x_n) = \underline{w}^T \phi(x_n) + w_0$,

$$\text{KKT conditions} \left\{ \begin{array}{ll} \xi_n^+ \geq 0, \xi_n^- \geq 0 & \forall n \\ y_n \leq \hat{y}(x_n) + \epsilon + \xi_n^+ & \forall n \\ y_n \geq \hat{y}(x_n) - \epsilon - \xi_n^- & \forall n \\ \mu_n^+, \mu_n^-, \lambda_n^+, \lambda_n^- \geq 0 & \forall n \end{array} \right. \quad \begin{array}{ll} \mu_n^+ \xi_n^+ = 0 & \forall n \\ \mu_n^- \xi_n^- = 0 & \forall n \\ \lambda_n^+ (\hat{y}_n + \epsilon + \xi_n^+ - y_n) = 0 & \forall n \\ \lambda_n^- (\epsilon + \xi_n^- - \hat{y}_n + y_n) = 0 & \forall n \end{array}$$

Derive Dual Representation by $\min_{\underline{w}, w_0, \underline{\xi}^+, \underline{\xi}^-} L(\underline{w}, w_0, \underline{\xi}^+, \underline{\xi}^-, \underline{\mu}^+, \underline{\mu}^-, \underline{\lambda}^+, \underline{\lambda}^-)$

$$\nabla_{\underline{w}} L = \underline{w} - \sum_{n=1}^N \lambda_n^+ \phi(x_n) + \sum_{n=1}^N \lambda_n^- \phi(x_n) = 0 \Rightarrow \underline{w} = \sum_{n=1}^N (\lambda_n^+ - \lambda_n^-) \phi(x_n)$$

$$\nabla_{w_0} L = -\sum_{n=1}^N \lambda_n^+ + \sum_{n=1}^N \lambda_n^- = 0 \Rightarrow \sum_{n=1}^N (\lambda_n^+ - \lambda_n^-) = 0$$

$$\nabla_{\xi_n^+} L = C - \mu_n^+ - \lambda_n^+ = 0 \Rightarrow \mu_n^+ + \lambda_n^+ = C$$

$$\nabla_{\xi_n^-} L = C - \mu_n^- - \lambda_n^- = 0 \Rightarrow \mu_n^- + \lambda_n^- = C$$

$$L = C \cdot \sum_{n=1}^N (\xi_n^+ + \xi_n^-) + \frac{1}{2} \|\underline{w}\|^2 - \sum_{n=1}^N (\mu_n^+ \xi_n^+ + \mu_n^- \xi_n^-) - \sum_{n=1}^N \lambda_n^+ (\hat{y}_n + \epsilon + \xi_n^+ - y_n) - \sum_{n=1}^N \lambda_n^- (\epsilon + \xi_n^- - \hat{y}_n + y_n)$$

$$= C \cdot \sum_{n=1}^N (\xi_n^+ + \xi_n^-) + \frac{1}{2} \sum_{n=1}^N (\lambda_n^+ - \lambda_n^-) \phi^T(x_n) \cdot \sum_{m=1}^N (\lambda_m^+ - \lambda_m^-) \phi(x_m) - \sum_{n=1}^N (\underbrace{\mu_n^+ + \lambda_n^+}_C) \xi_n^+ + (\underbrace{\mu_n^- + \lambda_n^-}_C) \xi_n^- \\ - \sum_{n=1}^N \lambda_n^+ \hat{y}_n - \lambda_n^- \hat{y}_n + \epsilon (\lambda_n^+ + \lambda_n^-) - y_n (\lambda_n^+ - \lambda_n^-)$$

$$= -\epsilon \sum_{n=1}^N (\lambda_n^+ + \lambda_n^-) + \sum_{n=1}^N y_n (\lambda_n^+ - \lambda_n^-) + \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (\lambda_n^+ \lambda_m^-) (\lambda_n^+ - \lambda_m^-) \phi^T(x_n) \phi(x_m) \\ - \sum_{n=1}^N (\lambda_n^+ - \lambda_n^-) \left\{ \left[\sum_{m=1}^N (\lambda_m^+ - \lambda_m^-) \phi(x_m) \right]^T \phi(x_n) + w_0 \right\}$$

$$= -\sum_{n=1}^N \sum_{m=1}^N (\lambda_n^+ - \lambda_n^-) (\lambda_m^+ - \lambda_m^-) \phi(x_n)^T \phi(x_m) - w_0 \sum_{n=1}^N (\lambda_n^+ - \lambda_n^-) = 0$$

$$\therefore L_D(\underline{\lambda}^+, \underline{\lambda}^-) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (\lambda_n^+ - \lambda_n^-) (\lambda_m^+ - \lambda_m^-) \underbrace{\phi(x_n)^T \phi(x_m)}_{K(x_n, x_m)} - \epsilon \sum_{n=1}^N (\lambda_n^+ + \lambda_n^-) + \sum_{n=1}^N y_n (\lambda_n^+ - \lambda_n^-)$$

↑

$K(x_n, x_m) \leftarrow \text{kernel}$

$\underline{\mu}^+$ & $\underline{\mu}^-$ disappear, but constraints exist

because $\mu_n^+ + \lambda_n^+ = C$ & $\mu_n^- + \lambda_n^- = C$

$\Rightarrow 0 \leq \lambda_n^+ \leq C$ & $0 \leq \lambda_n^- \leq C \leftarrow \text{similar to SVC case.}$

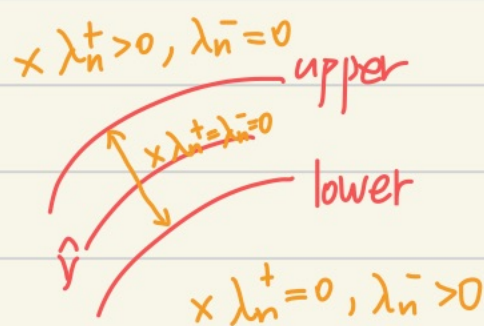
Now KKT conditions

| | | |
|---|--|--|
| $\xi_n^+ \geq 0, \xi_n^- \geq 0 \quad \forall n$ | $(C - \lambda_n^+) \xi_n^+ = 0 \quad \forall n$ | $\sum_{n=1}^N (\lambda_n^+ - \lambda_n^-) = 0$ |
| $y_n \leq \hat{y}(x_n) + \epsilon + \xi_n^+ \quad \forall n$ | $(C - \lambda_n^-) \xi_n^- = 0 \quad \forall n$ | |
| $y_n \geq \hat{y}(x_n) - \epsilon - \xi_n^- \quad \forall n$ | $\lambda_n^+ (\hat{y}_n + \epsilon + \xi_n^+ - y_n) = 0 \quad \forall n$ | |
| $(C - \lambda_n^+, C - \lambda_n^-), \lambda_n^+, \lambda_n^- \geq 0 \quad \forall n$ | $\lambda_n^- (\epsilon + \xi_n^- - \hat{y}_n + y_n) = 0 \quad \forall n$ | |

if x is on/above upper bound, $\lambda_n^+ > 0$, else = 0

if x is on/below lower bound, $\lambda_n^- > 0$, else = 0

$\therefore 3$ cases



SMD: $\max_{\underline{\lambda}^+, \underline{\lambda}^-} L_D(\underline{\lambda}^+, \underline{\lambda}^-)$

$$L_D(\underline{\lambda}^+, \underline{\lambda}^-) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (\lambda_n^+ - \lambda_n^-) (\lambda_m^+ - \lambda_m^-) K(x_n, x_m) - \epsilon \sum_{n=1}^N (\lambda_n^+ + \lambda_n^-) + \sum_{n=1}^N y_n (\lambda_n^+ - \lambda_n^-)$$

Similar to SVC, choose $\lambda_m^{+/-}, \lambda_n^{+/-}$ corresponding to two points.

then $L_D(\lambda_m^{+/-}, \lambda_n^{+/-}) = -\frac{1}{2} (\lambda_m^+ - \lambda_m^-) (\lambda_m^+ - \lambda_m^-) K_{mm} - \frac{1}{2} (\lambda_m^+ - \lambda_m^-) (\lambda_n^+ - \lambda_n^-) K_{mn}$

m, m m, n

$$\begin{aligned}
& -\frac{1}{2}(\lambda_m^+ - \lambda_m^-)(\lambda_m^+ - \lambda_m^-)K_{mm} - \frac{1}{2}(\lambda_n^+ - \lambda_n^-)(\lambda_n^+ - \lambda_n^-)K_{nn} - \frac{1}{2}(\lambda_m^+ - \lambda_m^-) \sum_{j=1}^N (\lambda_j^+ - \lambda_j^-) K_{mj} \\
& - \frac{1}{2}(\lambda_n^+ - \lambda_n^-) \sum_{j=1}^N (\lambda_j^+ - \lambda_j^-) K_{nj} - \frac{1}{2}(\lambda_m^+ - \lambda_m^-) \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) K_{im} - \frac{1}{2}(\lambda_n^+ - \lambda_n^-) \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) K_{in} \\
& - \epsilon(\lambda_m^+ + \lambda_m^- + \lambda_n^+ + \lambda_n^-) + \gamma_m(\lambda_m^+ - \lambda_m^-) + \gamma_n(\lambda_n^+ - \lambda_n^-) + K
\end{aligned}$$

\uparrow constant

$$\begin{aligned}
L_D(\lambda_m^+, \lambda_n^+) &= -\frac{1}{2}(\lambda_m^+ - \lambda_m^-)^2 K_{mm} - \frac{1}{2}(\lambda_n^+ - \lambda_n^-)^2 K_{nn} - (\lambda_m^+ - \lambda_m^-)(\lambda_n^+ - \lambda_n^-) K_{mn} \\
& - (\lambda_m^+ - \lambda_m^-) \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) K_{im} - (\lambda_n^+ - \lambda_n^-) \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) K_{in} \\
& - \epsilon(\lambda_m^+ + \lambda_m^- + \lambda_n^+ + \lambda_n^-) + \gamma_m(\lambda_m^+ - \lambda_m^-) + \gamma_n(\lambda_n^+ - \lambda_n^-) + K
\end{aligned}$$

\uparrow constant

Form dual representation $KK^T: \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) = 0$

$$\therefore (\lambda_m^+ - \lambda_m^-) + (\lambda_n^+ - \lambda_n^-) = - \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) \leftarrow \text{let it be } \alpha$$

$$\therefore (\lambda_m^+ - \lambda_m^-) = \alpha - (\lambda_n^+ - \lambda_n^-)$$

$$\lambda_m^+ + \lambda_m^- = |\alpha - (\lambda_n^+ - \lambda_n^-)| \quad \text{because only } \lambda_m^+ \text{ or } \lambda_m^- > 0$$

$$\text{Also, let } V_m = \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) K_{im} \quad V_n = \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) K_{in}$$

$$\begin{aligned}
L_D(\lambda_n^+, \lambda_n^-) &= -\frac{1}{2}(\alpha - \lambda_n^+ + \lambda_n^-)^2 K_{nn} - \frac{1}{2}(\lambda_n^+ - \lambda_n^-)^2 K_{nn} - (\alpha - \lambda_n^+ + \lambda_n^-)(\lambda_n^+ - \lambda_n^-) K_{nn} \\
& - (\alpha - \lambda_n^+ + \lambda_n^-) V_m - (\lambda_n^+ - \lambda_n^-) V_n - \epsilon \cdot |\alpha - (\lambda_n^+ - \lambda_n^-)| - \epsilon(\lambda_n^+ + \lambda_n^-) \\
& + \gamma_m(\alpha - \lambda_n^+ + \lambda_n^-) + \gamma_n(\lambda_n^+ - \lambda_n^-) + K
\end{aligned}$$

Combine two variable λ_i^+ & λ_i^- into one β_i $\begin{cases} \beta_i = \lambda_i^+ - \lambda_i^- \\ |\beta_i| = \lambda_i^+ + \lambda_i^- \end{cases}$

$$\text{Then, } L_D(\beta_n) = -\frac{1}{2}(\alpha - \beta_n)^2 K_{mm} - \frac{1}{2}\beta_n^2 K_{nn} - (\alpha - \beta_n)\beta_n K_{mn} \\ - (\alpha - \beta_n)V_m - \beta_n V_n - \epsilon \cdot |\alpha - \beta_n| - \epsilon \cdot |\beta_n| + \gamma_m(\alpha - \beta_n) + \gamma_n \beta_n + K$$

$$\frac{\partial L_D(\beta_n)}{\partial \beta_n} = (\alpha - \beta_n)K_{mm} - \beta_n K_{nn} + \beta_n K_{mn} - (\alpha - \beta_n)K_{mn} + V_m - V_n \\ + \epsilon \cdot \text{sign}(\beta_m) - \epsilon \cdot \text{sign}(\beta_n) - \gamma_m + \gamma_n$$

$$\text{where } \text{sign}(x) = \begin{cases} 1, & x > 0 \\ 0, & x = 0 \\ -1, & x < 0 \end{cases}$$

$$\text{Let } \frac{\partial L_D(\beta_n)}{\partial \beta_n} = 0 \Rightarrow \beta_n K_{mm} + \beta_n K_{nn} - 2\beta_n K_{mn} = \alpha K_{mm} - \alpha K_{mn} - \gamma_m + \gamma_n \\ + \epsilon [\text{sign}(\beta_m) - \text{sign}(\beta_n)] + V_m - V_n$$

$$\text{because } V_m = \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) K_{im} \quad V_n = \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) K_{in} \quad \& \quad \underline{w} = \sum_{n=1}^N (\lambda_n^+ - \lambda_n^-) \phi(x_n) \\ \& \quad f(x) = \underline{w}^T x + w_0$$

$$\therefore V_i = f(x_i) - \beta_m K_{mi} - \beta_n K_{ni} - w_0 = \sum_{n=1}^N (\lambda_n^+ - \lambda_n^-) K_{ni} + w_0$$

$$\therefore V_m - V_n = f(x_m) - f(x_n) - \beta_m (K_{mm} - K_{mn}) - \beta_n (K_{nm} - K_{nn})$$

$$\therefore \underbrace{\beta_n}_{\text{new}} (K_{mm} + K_{nn} - 2K_{mn}) = \alpha (K_{mm} - K_{mn}) - \gamma_m + \gamma_n + \epsilon [\text{sign}(\beta_m) - \text{sign}(\beta_n)] \\ + f(x_m) - f(x_n) - \underbrace{\beta_m}_{\text{old}} (K_{mm} - K_{mn}) - \underbrace{\beta_n}_{\text{old}} (K_{nm} - K_{nn})$$

$$\text{because } \alpha = - \sum_{i=1}^N (\lambda_i^+ - \lambda_i^-) = \underbrace{\beta_m}_{\text{old}} + \underbrace{\beta_n}_{\text{old}}$$

$$\therefore \beta_n^{\text{new}} \eta = \gamma_n - \gamma_m + \epsilon [\text{sign}(\beta_m) - \text{sign}(\beta_n)] + f(x_m) - f(x_n) + \beta_n^{\text{old}} \eta$$

$$\text{in which } \eta = (K_{mm} + K_{nn} - 2K_{mn})$$

$$\Rightarrow \beta_n^{\text{new}} = \beta_n^{\text{old}} + \frac{1}{\eta} (y_n - y_m + \epsilon [\text{sign}(\beta_m) - \text{sign}(\beta_n)] + f(x_m) - f(x_n))$$

where

$$\begin{cases} f(x) = \sum_{n=1}^N (\lambda_n^+ - \lambda_n^-) K_{nk} + w_0 = \sum_{n=1}^N \beta_n^{\text{old}} K_{nk} + w_0 \\ \eta = (K_{mm} + K_{nn} - 2K_{mn}) \end{cases}$$

Note*: In above β_n update function, the β_m & β_n in $\text{sign}(\cdot)$ function should be **new** instead of old. But this makes it a recursive function.

Finding solutions ✧ [Efficient SVM Regression Training with SMO.
Gary & Steve. 2002]

① if $\beta_m < 0 < \beta_n$, $\text{sign}(\beta_m) - \text{sign}(\beta_n) = -2$

② if $\beta_n < 0 < \beta_m$, $\text{sign}(\beta_m) - \text{sign}(\beta_n) = 2$

③ if $\beta_m, \beta_n \overset{\text{both}}{\geq} 0$, $\text{sign}(\beta_m) - \text{sign}(\beta_n) = 0$

In these two cases, we don't really do this

if $\beta_m = 0$, $\beta_n = \alpha - \beta_m = \alpha = \beta_m^{\text{old}} + \beta_n^{\text{old}}$
 $\text{sign}(\beta_m) - \text{sign}(\beta_n) = -\text{sign}(\beta_n) = -\text{sign}(\alpha)$

if $\beta_n = 0$, $\beta_m = \alpha - \beta_n = \alpha = \beta_m^{\text{old}} + \beta_n^{\text{old}}$
 $\text{sign}(\beta_m) - \text{sign}(\beta_n) = \text{sign}(\beta_m) = \text{sign}(\alpha)$

Because $\text{sign}(\cdot)$ function makes the update function a recursive function, there is no way to calculate; below is a reasonable efficient algorithm that is proved by exps.

For a single step of SMO (update β)

1 we have $\alpha = \beta_m^{\text{old}} + \beta_n^{\text{old}}$

$$\eta = K_{mm} + K_{nn} - 2K_{mn}$$

$$\Delta = \frac{2\epsilon}{\eta} \quad \leftarrow \text{correction term [1] [2]}$$

2 $\beta_n = \beta_n^{\text{old}} + \frac{1}{\eta} [y_n - y_m + \underbrace{f(x_m)}_{\text{with old } \beta} - \underbrace{f(x_n)}_{\text{with old } \beta}]$

$$\beta_m = \alpha - \beta_n \quad \leftarrow \text{if } \text{sign}(\beta_m) == \text{sign}(\beta_n), \text{ jump 4.}$$

3 if $\beta_m \cdot \beta_n < 0$, use correction term. [3]

$$\text{if } (|\beta_m| \geq \Delta \ \&\& \ |\beta_n| \geq \Delta) \quad \leftarrow \text{sign won't be affected by } \Delta \text{ term}$$

$$\beta_n = \beta_n + \text{sign}(\beta_m) \cdot \Delta$$

else \leftarrow signs of β_m & β_n will be affected, one of them is set to be 0, another be α .

$$\beta_n = \alpha \quad \text{if } |\beta_n| > |\beta_m| \text{ else } 0$$

4 Group step by $[0 \leq \lambda_i^+, \lambda_i^- \leq C] \quad [-C \leq \beta_i = \lambda_i^+ - \lambda_i^- \leq C]$

$$L = \max(\alpha - C, -C)$$

$$H = \min(C, \alpha + C)$$

$$\beta_n^{\text{new}} = \min(\max(\beta_n, L), H)$$

$$\beta_m^{\text{new}} = \alpha - \beta_n^{\text{new}}$$

5 update w_0 , to make it easy:

if $\beta_m^{\text{new}} = 0$, force it have $f(x_m) = y_m$, then $y_m = \sum_{i=1}^N \beta_i^{\text{new}} K_{im} + w_0^{\text{new}}$
(within tube)

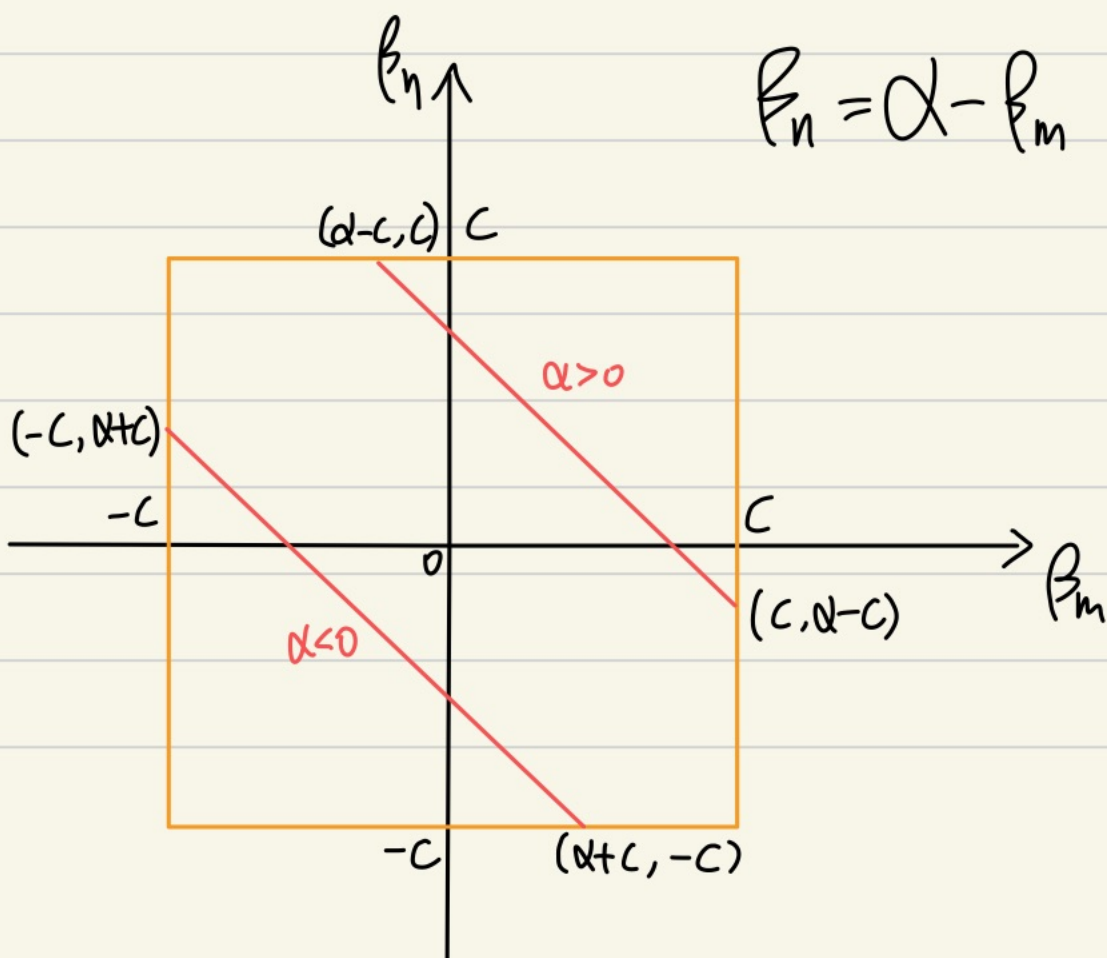
$$\begin{aligned} \text{i) } w_0^{\text{new}} &= y_m - \sum_{i=1, i \neq m, n}^N \beta_i^{\text{new}} K_{im} - \beta_m^{\text{new}} K_{mm} - \beta_n^{\text{new}} K_{nm} \\ &= y_m - [f^{\text{old}}(x_m) - w_0^{\text{old}} - \beta_m^{\text{old}} K_{mm} - \beta_n^{\text{old}} K_{nm}] - \beta_m^{\text{new}} K_{mm} - \beta_n^{\text{new}} K_{nm} \\ &= w_0^{\text{old}} + y_m - f^{\text{old}}(x_m) + (\beta_m^{\text{old}} - \beta_m^{\text{new}}) K_{mm} + (\beta_n^{\text{old}} - \beta_n^{\text{new}}) K_{nm} \end{aligned}$$

if $\beta_n^{\text{new}} = 0$, force it have $f(x_n) = y_n$
(within tube)

$$\text{ii) } w_0^{\text{new}} = w_0^{\text{old}} + y_n - f^{\text{old}}(x_n) + (\beta_m^{\text{old}} - \beta_m^{\text{new}}) K_{mn} + (\beta_n^{\text{old}} - \beta_n^{\text{new}}) K_{nn}$$

if neither, average over i) & ii).

Crop Visualization



$$\therefore \alpha < 0, -C \leq \beta_n \leq \alpha + C$$

$$\alpha > 0, \alpha - C \leq \beta_n \leq C$$

$$\Rightarrow L = \max(\alpha - C, -C)$$

$$H = \min(\alpha + C, C)$$

for β_n