

1.1.

training data: $\{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$

where: $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}^k$ for $i = 1, 2, \dots, N$

consider: $\min_{w \in \mathbb{R}^{d \times k}, b \in \mathbb{R}^k} \sum_{i=1}^N \|y_i - w x_i - b\|_2^2$ (least-square problem)

Q1: set $\tilde{x}_i = \begin{bmatrix} x_i \\ 1 \end{bmatrix} \in \mathbb{R}^{d+1}$

$\tilde{w} = [w | b] \in \mathbb{R}^{k \times (d+1)}$

then origin question $\Leftrightarrow \min_{\tilde{w}} \sum_{i=1}^N \|y_i - \tilde{w} \tilde{x}_i\|_2^2$

set $X = \begin{bmatrix} \tilde{x}_1^T \\ \tilde{x}_2^T \\ \vdots \\ \tilde{x}_N^T \end{bmatrix} = \begin{bmatrix} x_1^T & 1 \\ x_2^T & 1 \\ \vdots & \vdots \\ x_N^T & 1 \end{bmatrix} \in \mathbb{R}^{N \times (d+1)}$; $Y = \begin{bmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_N^T \end{bmatrix} \in \mathbb{R}^{N \times k}$

$$\text{then } L(\tilde{w}) = \|Y - X \tilde{w}^T\|_F^2 = \text{tr}((Y - X \tilde{w}^T)^T (Y - X \tilde{w}^T)) \\ = \text{tr}(Y^T Y) - 2 \text{tr}(Y^T X \tilde{w}^T) + \text{tr}(X^T X \tilde{w}^T \tilde{w}^T)$$

$$\therefore \frac{\partial L}{\partial \tilde{w}} = -2(Y^T X)^T + 2(X^T X \tilde{w}^T)^T = -2X^T Y + 2X^T X \tilde{w}^T$$

$$\text{let } \frac{\partial L}{\partial \tilde{w}} = 0 \Rightarrow -2X^T Y + 2X^T X \tilde{w}^T = 0 \\ \Rightarrow X^T X \tilde{w}^T = X^T Y$$

assume $X^T X$ is invertible then

$$\tilde{w}^T = (X^T X)^{-1} X^T Y$$

$$\therefore \tilde{w} = [w | b]$$

$$\therefore w^* = \tilde{w}^T[:, :d], b^* = \tilde{w}^T[:, d]$$

where $\tilde{w}^T = (X^T X)^{-1} X^T Y \in \mathbb{R}^{(d+1) \times k}$

∴ final closed form $\begin{bmatrix} w^* \\ b^* \end{bmatrix} = (X^T X)^{-1} X^T Y \in \mathbb{R}^{(d+1) \times k}$ where $X = \begin{bmatrix} \tilde{x}_1^T \\ \tilde{x}_2^T \\ \vdots \\ \tilde{x}_N^T \end{bmatrix} = \begin{bmatrix} x_1^T & 1 \\ x_2^T & 1 \\ \vdots & \vdots \\ x_N^T & 1 \end{bmatrix} \in \mathbb{R}^{N \times (d+1)}$; $Y = \begin{bmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_N^T \end{bmatrix} \in \mathbb{R}^{N \times k}$

Q2:

for a single sample (x_i, y_i) , the loss is

$$l_i(w, b) = \|y_i - w x_i - b\|_2^2 = (y_i - w x_i - b)^T (y_i - w x_i - b)$$

we have

$$\frac{\partial l_i}{\partial w} = -2(y_i - w x_i - b)x_i^T; \quad \frac{\partial l_i}{\partial b} = -2(y_i - w x_i - b)$$

$$\therefore \frac{\partial L}{\partial w} = -2 \sum_{i=1}^N (y_i - w x_i - b)x_i^T \quad \text{and} \quad \frac{\partial L}{\partial b} = -2 \sum_{i=1}^N (y_i - w x_i - b)$$

⇒ At iteration t , we have

$$w^{(t+1)} = w^{(t)} + 2\eta \sum_{i=1}^N (y_i - w^{(t)} x_i - b^{(t)}) x_i^T$$

$$b^{(t+1)} = b^{(t)} + 2\eta \sum_{i=1}^N (y_i - w^{(t)} x_i - b^{(t)})$$

where $\eta > 0$ is the learning rate.

in matrix form.

$$\text{Gradient: } \frac{\partial L}{\partial w} = -2X^T(Y - X\tilde{w}^T)$$

$$\text{Update rule: } \tilde{w}^{(t+1)} = \tilde{w}^{(t)} + 2\eta X^T(Y - X\tilde{w}^{(t)})$$

Stopping Criteria

1. Stop when the gradient norm is small enough:

$$\left\| \frac{\partial L}{\partial w} \right\|_F < \epsilon_g$$

2. Stop when the change in loss between consecutive iteration is small enough

$$|L(w^{(t+1)}, b^{(t+1)}) - L(w^{(t)}, b^{(t)})| < \epsilon_L \quad \text{or} \quad \frac{|L^{(t)} - L^{(t+1)}|}{|L^{(t)} + \epsilon_L|} < \epsilon_L$$

Q3.
consider loss for single sample

$$\ell_i(w) = -\sum_{k=1}^K y_{i,k} \log(\hat{y}_{i,k})$$

$$\text{then } \frac{\partial \ell_i}{\partial w} = \frac{\partial \ell_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial w}$$

$$\text{where } \frac{\partial \ell_i}{\partial z_{i,k}} = -\sum_{j=1}^K y_{i,j} \frac{\partial \log(\hat{y}_{i,j})}{\partial z_{i,k}}$$

$$\text{for } \frac{\partial \hat{y}_{i,j}}{\partial z_{i,k}} \text{ when } j=k, \frac{\partial \hat{y}_{i,j}}{\partial z_{i,k}} = \hat{y}_{i,k}(1-\hat{y}_{i,k}) \text{ for softmax}$$

$$\text{when } j \neq k, \frac{\partial \hat{y}_{i,j}}{\partial z_{i,k}} = -\hat{y}_{i,j} \hat{y}_{i,k}$$

combined, we have

$$\frac{\partial \hat{y}_{i,j}}{\partial z_{i,k}} = \hat{y}_{i,j} (\delta_{jk} - \hat{y}_{i,k}) \text{ where } \delta_{jk} \text{ is the Kronecker delta.}$$

$$\begin{aligned} \therefore \frac{\partial \ell_i}{\partial z_{i,k}} &= -\sum_{j=1}^K y_{i,j} \frac{1}{\hat{y}_{i,j}} \frac{\partial \hat{y}_{i,j}}{\partial z_{i,k}} \\ &= -\sum_{j=1}^K y_{i,j} \frac{1}{\hat{y}_{i,j}} (\delta_{jk} - \hat{y}_{i,k}) \cdot \hat{y}_{i,j} \\ &= -\sum_{j=1}^K y_{i,j} (\delta_{jk} - \hat{y}_{i,k}) \end{aligned}$$

$\because y_i$ is one-hot encoded ($\sum_{j=1}^K y_{i,j} = 1$)

$$\therefore \frac{\partial \ell_i}{\partial z_{i,k}} = \hat{y}_{i,k} - y_{i,k}$$

in vector form

$$\frac{\partial \ell_i}{\partial z_i} = \hat{y}_i - y_i$$

$$\because z_i = w^T x_i$$

$$\therefore z_{i,k} = \sum_{m=1}^M w_{m,k} x_{i,m}$$

$$\therefore \frac{\partial z_{i,k}}{\partial w_{m,n}} = x_{i,m} \delta_{kn}$$

in matrix form

$$\frac{\partial z_i}{\partial w} = x_i \otimes I_K \text{ where } \otimes \text{ is the outer product.}$$

$$\therefore \frac{\partial \ell_i}{\partial w} = x_i (\hat{y}_i - y_i)^T$$

$$\therefore \frac{\partial J}{\partial w} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \ell_i}{\partial w}$$

$$= \frac{1}{N} \sum_{i=1}^N x_i (\hat{y}_i - y_i)^T$$

$$\text{set } X = \begin{bmatrix} \hat{x}_1^T \\ \hat{x}_2^T \\ \vdots \\ \hat{x}_N^T \end{bmatrix} \in R^{N \times d}; Y = \begin{bmatrix} y_1^T \\ y_2^T \\ \vdots \\ y_N^T \end{bmatrix} \in R^{N \times k}, \hat{Y} = \begin{bmatrix} \hat{y}_1^T \\ \hat{y}_2^T \\ \vdots \\ \hat{y}_N^T \end{bmatrix}$$

$$\text{then } \sum_{i=1}^N x_i (\hat{y}_i - y_i)^T = X^T (\hat{Y} - Y)$$

$$\therefore \nabla_w J = \frac{1}{N} X^T (\hat{Y} - Y)$$

1.2.

Q1.

Slack Variables

$\epsilon_{i,0}$: measures the violation when $y_i - f(x_i) > \epsilon$

$\epsilon_{i,0}^*$: measures the violation when $f(x_i) - y_i > \epsilon$

\therefore origin problem becomes: $\min_{w, \epsilon, \epsilon^*} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\epsilon_i + \epsilon_i^*)$

$$\begin{aligned} \text{Subject to} \\ y_i - w^T x_i &\leq \epsilon + \xi_i^* \\ w^T x_i - y_i &\leq C + \xi_i \end{aligned}$$

Q2.

introduce
 $\alpha_i \geq 0$ for constraints $w^T x_i - y_i \leq \epsilon + \xi_i^*$

$\alpha_i^* \geq 0$ for constraints $y_i - w^T x_i \leq \epsilon + \xi_i^*$

$\lambda_i \geq 0$ for constraints $\xi_i \geq 0$

$\lambda_i^* \geq 0$ for constraints $\xi_i^* \geq 0$

$$\begin{aligned} \therefore L(w, \xi, \xi^*, \alpha, \alpha^*, \lambda, \lambda^*) &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) \\ &+ \sum_{i=1}^n \alpha_i (w^T x_i - y_i - \epsilon - \xi_i) \\ &+ \sum_{i=1}^n \alpha_i^* (w^T x_i - y_i - \epsilon - \xi_i^*) \\ &- \sum_{i=1}^n \lambda_i \xi_i - \sum_{i=1}^n \lambda_i^* \xi_i^* \end{aligned}$$

Q3.

We have

$$\begin{aligned} L(w, \xi, \xi^*, \alpha, \alpha^*, \lambda, \lambda^*) &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n (\xi_i + \xi_i^*) \\ &+ \sum_{i=1}^n \alpha_i (w^T x_i - y_i - \epsilon - \xi_i) \\ &+ \sum_{i=1}^n \alpha_i^* (w^T x_i - y_i - \epsilon - \xi_i^*) \\ &- \sum_{i=1}^n \lambda_i \xi_i - \sum_{i=1}^n \lambda_i^* \xi_i^* \end{aligned}$$

\because KKT

$$\begin{aligned} \therefore \frac{\partial L}{\partial w} &= w + \sum_{i=1}^n \alpha_i x_i - \sum_{i=1}^n \alpha_i^* x_i = 0 \\ \Rightarrow w &= \sum_{j=1}^n (\alpha_j^* - \alpha_j) x_j \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial \xi_i} &= C - \alpha_i - \lambda_i = 0 \\ \Rightarrow C &= \alpha_i + \lambda_i \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial \xi_i^*} &= C - \alpha_i^* - \lambda_i^* = 0 \\ \Rightarrow C &= \alpha_i^* + \lambda_i^* \end{aligned}$$

$$\therefore \alpha_i (w^T x_i - y_i - \epsilon - \xi_i) = 0$$

$$\alpha_i^* (y_i - w^T x_i - \epsilon - \xi_i^*) = 0$$

$$\lambda_i \xi_i = 0$$

$$\lambda_i^* \xi_i^* = 0$$

$$\therefore \frac{1}{2} \|w\|^2 = \frac{1}{2} \left\| \sum_{i=1}^n (\alpha_i^* - \alpha_i) x_i \right\|^2 = \frac{1}{2} \sum_{i,j=1}^n (\alpha_i^* - \alpha_i)(\alpha_j^* - \alpha_j) x_i^T x_j$$

$$\therefore C = \alpha_i + \lambda_i \quad C = \alpha_i^* + \lambda_i^*$$

$$\therefore C \sum_{i=1}^n (\xi_i + \xi_i^*) - \sum_{i=1}^n \lambda_i \xi_i - \sum_{i=1}^n \lambda_i^* \xi_i^* = \sum_{i=1}^n \alpha_i \xi_i + \sum_{i=1}^n \alpha_i^* \xi_i^*$$

$$\therefore L = -\frac{1}{2} \sum_{i,j=1}^n (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) x_i^T x_j + \sum_{i=1}^n (\alpha_i^* - \alpha_i) y_i - C \sum_{i=1}^n (\alpha_i + \alpha_i^*)$$

\therefore we have Dual Form of SVR

Maximize

$$L_D(\alpha^*) = -\frac{1}{2} \sum_{i,j=1}^n (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) x_i^T x_j + \sum_{i=1}^n (\alpha_i^* - \alpha_i) y_i - C \sum_{i=1}^n (\alpha_i + \alpha_i^*)$$

Subject to:

$$0 \leq \alpha_i \leq C, i=1, \dots, n$$

$$0 \leq \alpha_i^* \leq C, i=1, \dots, n$$

$$\text{Using } w = \sum_{i=1}^n (\alpha_i^* - \alpha_i) x_i \text{ we have } f(x) = \sum_{i=1}^n (\alpha_i^* - \alpha_i) x_i^T x$$

Q4. yes:

Dual Problem re-written:

Minimize

$$\frac{1}{2} \sum_{i,j=1}^n (\alpha_i - \alpha_i^*) (\alpha_j - \alpha_j^*) x_i^T x_j - \sum_{i=1}^n (\alpha_i^* - \alpha_i) y_i + \epsilon \sum_{i=1}^n (\alpha_i + \alpha_i^*)$$

Subject to:

$$0 \leq \alpha_i \leq C, i=1, \dots, n$$

$$0 \leq \alpha_i^* \leq C, i=1, \dots, n$$

We define

$$u = [\alpha_1 \dots \alpha_n, \alpha_1^* \dots, \alpha_n^*]^T \in \mathbb{R}^{2n}$$

$$Q = \begin{bmatrix} K & -K \\ -K & K \end{bmatrix} \text{ where } K_{ij} = x_i^T x_j$$

$$p = [-y_1 + \epsilon, \dots, -y_n + \epsilon, y_1 + \epsilon, \dots, y_n + \epsilon]^T$$

then we have

$$\min \frac{1}{2} u^T Q u + p^T u$$

Subject to

$$0 \leq u \leq C$$

Q5.

define margin Support Vectors:

Points where $0 < \alpha_i < C$ or $0 < \alpha_i^* < C$

$$\text{For } \alpha_i > 0, w^T x_i - y_i = \epsilon$$

$$\text{For } \alpha_i^* > 0, y_i - w^T x_i = \epsilon$$

$$\text{and } |f(x_i) - y_i| = \epsilon$$

$\because KKT$

$$\therefore \alpha_i > 0, w^T x_i - y_i = \epsilon + \epsilon_i$$

$$\alpha_i^* > 0, y_i - w^T x_i = \epsilon + \epsilon_i^*$$

$$\text{for } 0 < \alpha_i < C, \epsilon_i = 0$$

$$\text{for } \alpha_i = C, \epsilon_i > 0$$

