

1.1 (i) Derivative with respect to w_0 Denote the error term as $e = y - Xw - w_0 \mathbf{1}$ \therefore The cost function is: $J(w, w_0) = e^T e + \lambda w^T w$

$$\frac{\partial J}{\partial w_0} = \frac{\partial}{\partial w_0} (e^T e) + \frac{\partial}{\partial w_0} (\lambda w^T w) = -2e^T \mathbf{1} = -2(y - Xw - w_0 \mathbf{1})^T \mathbf{1}$$

$$\text{Let } \frac{\partial J}{\partial w_0} = 0: (y - Xw - w_0 \mathbf{1})^T \mathbf{1} = 0$$

 \therefore The transpose of a scalar is the scalar itself: $\mathbf{1}^T (y - Xw - w_0 \mathbf{1}) = 0$

$$\text{Expand: } \mathbf{1}^T y - \mathbf{1}^T X w - w_0 \mathbf{1}^T \mathbf{1} = 0$$

 $\mathbf{0}^T (x \text{ is centered, each column sums to zero})$ n (the number of samples)

$$\therefore w_0 = \frac{1}{n} \mathbf{1}^T y = \bar{y} \quad \therefore \hat{w}_0 = \bar{y}$$

(ii) Derivative with respect to w

$$\frac{\partial J}{\partial w} = \frac{\partial}{\partial w} (e^T e) + \frac{\partial}{\partial w} (\lambda w^T w) = 2e^T \frac{\partial e}{\partial w} + 2\lambda w$$

$$\frac{\partial}{\partial w} (e^T e) = 2e^T (-X) = -2(y - Xw - w_0 \mathbf{1})^T X$$

$$\therefore \nabla_w [e^T e] = -2X^T (y - Xw - w_0 \mathbf{1})$$

$$\frac{\partial}{\partial w} (\lambda w^T w) = 2\lambda w$$

$$\nabla_w J = -2X^T (y - Xw - w_0 \mathbf{1}) + 2\lambda w$$

$$\text{Set to 0: } X^T (y - Xw - w_0 \mathbf{1}) = \lambda w$$

Substitute $w_0 = \bar{y}$:

$$X^T (y - Xw - \bar{y} \mathbf{1}) = \lambda w$$

 $\therefore X^T \mathbf{1}$ is a vector where the j -th element is $\sum_{i=1}^n X_{ij} = 0 \quad \therefore X^T \mathbf{1} = 0$

$$\therefore X^T (y - \bar{y} \mathbf{1}) = X^T y - \bar{y} X^T \mathbf{1} = X^T y = X^T X w + \lambda w = (X^T X + \lambda I) w \quad \therefore w = (X^T X + \lambda I)^{-1} X^T y$$

1.2 (i) Compute the feature vectors:

$$\text{For } x_1 = 0: \phi(x_1) = [1, 0, 0]^T$$

$$\text{For } x_2 = 1: \phi(x_2) = [1, \sqrt{2}, 1]^T$$

$$\phi(x_2) - \phi(x_1) = [0, \sqrt{2}, 1]^T$$

In SVM theory, $w = \sum \alpha_i y_i \phi(x_i)$ For two support vectors with equal $\alpha_i = \alpha$:

$$w = \alpha (y_1 \phi(x_1) + y_2 \phi(x_2)) = \alpha (-[1, 0, 0]^T + [1, \sqrt{2}, 1]^T) = \alpha [0, \sqrt{2}, 1]^T$$

 $\therefore w$ is parallel to $[0, \sqrt{2}, 1]^T$ 

$$\text{For } y_1 = -1: y_1 (w^T \phi(x_1) + w_0) = 1 \quad \therefore w^T \phi(x_1) + w_0 = -1$$

$$\text{For } y_2 = 1: y_2 (w^T \phi(x_2) + w_0) = 1 \quad \therefore w^T \phi(x_2) + w_0 = 1$$

Substitute the feature vectors ($w = [w_1, w_2, w_3]^T$):

$$\begin{cases} w^T \phi(x_1) = [w_1, w_2, w_3] \cdot [1, 0, 0] = w_1 & \therefore w_1 + w_0 = -1 \quad (1) \\ w^T \phi(x_2) = [w_1, w_2, w_3] \cdot [1, \sqrt{2}, 1] = w_1 + \sqrt{2} w_2 + w_3 & \therefore w_1 + \sqrt{2} w_2 + w_3 + w_0 = 1 \quad (2) \end{cases}$$

$$\text{Substitute (1) into (2): } \sqrt{2} w_2 + w_3 = 2 \quad (3)$$

$$\therefore w \text{ is parallel to } [0, \sqrt{2}, 1]^T, \text{ let } w = k [0, \sqrt{2}, 1]^T = [0, \sqrt{2}k, k]^T \quad \therefore \begin{cases} w_1 = 0 \\ w_2 = \sqrt{2}k \\ w_3 = k \end{cases}$$

$$\begin{cases} \text{From (1)}: W_0 = -1 \\ \text{From (3)}: k = \frac{2}{3} \end{cases} \therefore W = \left[0, \frac{2\sqrt{2}}{3}, \frac{2}{3} \right]^T$$

(ii) $\text{Margin} = \frac{1}{\|W\|} = \frac{1}{\frac{2}{\sqrt{3}}} = \frac{\sqrt{3}}{2}$

(iv) $W_0 = -1$

(v) $f(x) = W^T \phi(x) + W_0 = \left[0, \frac{2\sqrt{2}}{3}, \frac{2}{3} \right]^T [1, \sqrt{2}x, x^2] - 1 = \frac{2}{3}x^2 + \frac{4}{3}x - 1$

2.1 Definitions:

Input X , shape (N, D) , where $N=150$ (samples), $D=3$ (features including bias)
 Weights W , shape (D, C) , where $C=3$ (classes)
 Logits $Z = XW$, shape (N, C)
 Predictions $\hat{Y} = \text{softmax}(Z)$, where $\hat{y}_{ic} = \frac{\exp(z_{ic})}{\sum_{j=1}^C \exp(z_{ij})}$
 True labels: Y_{onehot} , shape (N, C) , one-hot encoded.
 Loss: $\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log(\hat{y}_{ic})$

Now need $\frac{\partial \mathcal{L}}{\partial W_{dk}}$, where W_{dk} is the weight for feature d and class k .

$$\mathcal{L}_i = -\sum_{c=1}^C y_{ic} \log(\hat{y}_{ic}) \quad \text{Total loss: } \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_i$$

$$\frac{\partial \mathcal{L}}{\partial W_{dk}} = \sum_{i=1}^N \left(\frac{\partial \mathcal{L}}{\partial \hat{y}_{ic}} \right) \left(\frac{\partial \hat{y}_{ic}}{\partial z_{ik}} \right) \left(\frac{\partial z_{ik}}{\partial W_{dk}} \right) = \frac{1}{N} \sum_{i=1}^N (\hat{y}_{ik} - y_{ik}) x_{id}$$

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \hat{y}_{ic}} &= -\frac{1}{N} \frac{y_{ic}}{\hat{y}_{ic}} \\ \text{if } c=k: \frac{\partial \hat{y}_{ic}}{\partial z_{ik}} &= \hat{y}_{ic} (1 - \hat{y}_{ic}) \\ \text{if } c \neq k: \frac{\partial \hat{y}_{ic}}{\partial z_{ik}} &= -\hat{y}_{ic} y_{ik} \\ z_{ik} &= \sum_{d=1}^D x_{id} w_{dk} \\ \frac{\partial z_{ik}}{\partial w_{dk}} &= x_{id} \end{aligned}$$

Matrix form: $\frac{\partial \mathcal{L}}{\partial W} = \frac{1}{N} \underbrace{X^T}_{\text{shape}(D, N)} (\underbrace{\hat{Y} - Y_{\text{onehot}}}_{\text{shape}(N, C)})$

\Rightarrow result: shape (D, C) , matching W .