



# Data Application Lab

**Tuesday Machine Learning Basis Quiz Answer**

## Part I: Deeper understanding on SVM

### Problem 1: Non-separable SVM

(i) The original constraint optimization problem:

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & y_i (w^T x_i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0. \end{aligned}$$

The generalized Lagrangian:

$$\mathcal{L}(w, b, \alpha, \beta, \gamma) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i [1 - \xi_i - y_i (w^T x_i + b)] - \sum_{i=1}^m \gamma_i \xi_i \quad \textcircled{1}$$

$\min_w \max_{\alpha, \beta, \gamma} \mathcal{L}(w, b, \alpha, \beta, \gamma)$  has the same solution as the original problem.

$\Rightarrow$  Dual problem  $\max_{\alpha, \beta, \gamma} \min_w \mathcal{L}(w, b, \alpha, \beta, \gamma)$

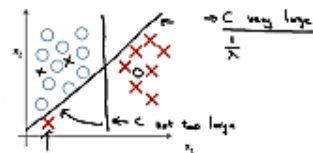
$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^m \alpha_i y_i x_i = 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i \quad \textcircled{2} \\ \frac{\partial \mathcal{L}}{\partial b} = - \sum_{i=1}^m \alpha_i y_i = 0 \Rightarrow \sum_{i=1}^m \alpha_i y_i = 0 \quad \textcircled{3} \\ \frac{\partial \mathcal{L}}{\partial \xi_i} = C - \alpha_i - \gamma_i = 0 \quad \textcircled{4} \end{cases}$$

$\textcircled{2} \textcircled{3} \textcircled{4} \Rightarrow \textcircled{1}$ :

$$\begin{aligned} \mathcal{L}(w^*, b^*, \alpha^*, \beta^*, \gamma^*) &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i [1 - \xi_i - y_i (w^T x_i + b)] - \sum_{i=1}^m \gamma_i \xi_i \\ &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i [1 - y_i (w^T x_i + b)] - \sum_{i=1}^m \alpha_i \xi_i - \sum_{i=1}^m \gamma_i \xi_i \\ &= \frac{1}{2} \|w\|^2 + \sum_{i=1}^m \alpha_i - w^T \left[ \sum_{i=1}^m \alpha_i y_i x_i \right] \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \|w\|^2 = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \end{aligned}$$

$\therefore$  the final result:

$$\begin{aligned} \max_{\alpha, \gamma} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \\ \text{s.t.} \quad & \alpha_i \geq 0, \alpha_i \leq C, \sum_{i=1}^m \alpha_i y_i = 0. \\ & \gamma_i \geq 0. \end{aligned}$$



(ii)  $C \rightarrow \infty$  the abnormal points will affect the classifier more, it will have a smaller margin hyperplane when it can separate the points correctly

$C \rightarrow 0$  the classifier will have a larger margin hyperplane even it will misclassify points.

## Problem 2: Kernels

Please try to understand the kernel tricks basing on this question.

$$\begin{aligned}
 \text{kernel } k(x, z) &= (x^T z + C)^2 \\
 &= (x_1 z_1 + x_2 z_2 + C)^2 \\
 &= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 + 2Cx_1 z_1 + 2Cx_2 z_2 + C^2 \\
 \phi(x) &= [x_1^2, \sqrt{2}x_1 x_2, x_2^2, \sqrt{2}Cx_1, \sqrt{2}Cx_2, C]^T \\
 \phi(z) &= [z_1^2, \sqrt{2}z_1 z_2, z_2^2, \sqrt{2}Cz_1, \sqrt{2}Cz_2, C]^T \\
 \langle \phi(x), \phi(z) \rangle &= x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2 + 2Cx_1 z_1 + 2Cx_2 z_2 + C^2 \\
 \therefore k(x, z) &= \langle \phi(x), \phi(z) \rangle.
 \end{aligned}$$

## Part II: The basic knowledge on ANN

### Problem 3: Backpropagation

$$\begin{aligned}
 \frac{\partial J(w, x, y)}{\partial w_{ji}^{(1)}} &= \frac{\partial J}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial h_j} \cdot \frac{\partial h_j}{\partial w_{ji}^{(1)}} && \text{sigmoid function:} \\
 &&& \sigma'(z) = \sigma(z)[1 - \sigma(z)] \\
 \frac{\partial J}{\partial \hat{y}_k} &= \hat{y}_k - y_k \\
 \frac{\partial \hat{y}_k}{\partial h_j} &= \frac{\partial (\sum_i h_j w_{kj}^{(2)})}{\partial h_j} = \sigma'(\sum_i h_j w_{kj}^{(2)}) \cdot w_{kj}^{(2)} = \hat{y}_k (1 - \hat{y}_k) w_{kj}^{(2)} \\
 \frac{\partial h_j}{\partial w_{ji}^{(1)}} &= \frac{\partial (\sum_i x_i w_{ji}^{(1)})}{\partial w_{ji}^{(1)}} = h_j (1 - h_j) x_i \\
 \therefore \frac{\partial J(w, x, y)}{\partial w_{ji}^{(1)}} &= \sum_k (\hat{y}_k - y_k) \hat{y}_k (1 - \hat{y}_k) w_{kj}^{(2)} h_j (1 - h_j) x_i
 \end{aligned}$$

### Problem 4: What are the advantages of ReLU over sigmoid function?

Sigmoid function has the problem of vanish gradient because the gradient of sigmoid becomes increasingly small as the absolute value of  $x$  increases. But ReLU can **reduce the likelihood of the gradient to vanish** and the constant gradient of ReLU when  $x > 0$  will **result in faster learning**.

Another advantage of ReLU is **sparsity**, which arises when  $x \leq 0$ . The more such units that exist in a layer the more sparse the resulting representation.