
Exercise 08 of Machine Learning [IN 2064]

Name: Yiman Li
Matr-Nr: 03724352
cooperate with Kejia Chen(03729686)

Problem 1

According to [1], both perceptron algorithm and support machine vector are solutions that separate the classes exactly.

However, the solution that perceptron algorithm finds will be dependent on the arbitrary initial values chosen for w and b as well as on the order in which the data points are presented. If there are multiple solutions all of which classify the training data set exactly, then we should try to find the one that will give the smallest generalization error.

On the other hand, the support vector machine approaches this problem through the concept of the margin, which is defined to be the smallest distance between the decision boundary and any of the samples.

Problem 2

a) Recall that $X \in \mathcal{R}^{N \times D}$, $y \in \mathcal{R}^N$, so we have

$$\begin{aligned} g(\alpha) &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j x_i^T x_j \\ &= \frac{1}{2} \alpha^T Q \alpha + \alpha^T \mathbf{1}_N \\ &= \frac{1}{2} \alpha^T (-X X^T \odot y y^T) \alpha + \alpha^T \mathbf{1}_N \end{aligned} \quad (1)$$

so we get

$$Q = -X X^T \odot y y^T \quad (2)$$

b) Firstly we have

$$Q = - \begin{bmatrix} x_1^T x_1 y_1 y_1 & x_1^T x_2 y_1 y_2 & \cdots & x_1^T x_N y_1 y_N \\ x_2^T x_1 y_2 y_1 & x_2^T x_2 y_2 y_2 & \cdots & x_2^T x_N y_2 y_N \\ \vdots & \vdots & \ddots & \vdots \\ x_N^T x_1 y_N y_1 & x_N^T x_2 y_N y_2 & \cdots & x_N^T x_N y_N y_N \end{bmatrix} \quad (3)$$

note that

$$Q_{(i,j)} = -(x_i^T x_j y_i y_j) = -(x_j^T x_i y_j y_i) = Q_{(j,i)} \quad (4)$$

which means that Q is symmetric. Since $X^T X$ and $y y^T$ are all *Gram Matrix*, so $-Q$ are always negative (semi-)definite matrix. Considering the linear constraints on α , so the function $g(\alpha)$ is concave, and the local maximizer of g is just its global maximizer. So we can search for a local maximizer of g to find its maximum.

Problem 3

Intuitively, the misclassification rate of LOOCV can be understood as below:

$$\epsilon = \frac{\text{Number of misclassification samples}}{\text{Number of all samples}} \quad (5)$$

Since we use support vector machine here with the number of training samples being s , so the number of misclassification samples must be smaller than s , so we have

$$\epsilon \leq \frac{s}{N} \quad (6)$$

Problem 4

To be seen in the end.

Problem 5

Firstly, we have

$$\begin{aligned} k(\mathbf{x}_m^T, \mathbf{x}_n^T) &= \sum_{i=1}^N a_i (\mathbf{x}_m^T \mathbf{x}_n)^i + a_0 \\ &= \sum_{i=1}^N a_i (\mathbf{x}_n^T \mathbf{x}_m)^i + a_0 \\ &= k(\mathbf{x}_n^T, \mathbf{x}_m^T) \end{aligned} \quad (7)$$

Which means the matrix is symmetric. Then the kernel matrix

$$\begin{aligned} \mathbf{K} &= \begin{bmatrix} \sum_{i=1}^N a_i (\mathbf{x}_1^T \mathbf{x}_1)^i + a_0 & \sum_{i=1}^N a_i (\mathbf{x}_1^T \mathbf{x}_2)^i + a_0 & \cdots & \sum_{i=1}^N a_i (\mathbf{x}_1^T \mathbf{x}_N)^i + a_0 \\ \sum_{i=1}^N a_i (\mathbf{x}_2^T \mathbf{x}_1)^i + a_0 & \sum_{i=1}^N a_i (\mathbf{x}_2^T \mathbf{x}_2)^i + a_0 & \cdots & \sum_{i=1}^N a_i (\mathbf{x}_2^T \mathbf{x}_N)^i + a_0 \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^N a_i (\mathbf{x}_N^T \mathbf{x}_1)^i + a_0 & \sum_{i=1}^N a_i (\mathbf{x}_N^T \mathbf{x}_2)^i + a_0 & \cdots & \sum_{i=1}^N a_i (\mathbf{x}_N^T \mathbf{x}_N)^i + a_0 \end{bmatrix} \\ &= a_1 \mathbf{X} \mathbf{X}^T + a_2 (\mathbf{X} \mathbf{X}^T)^2 + \cdots + a_N (\mathbf{X} \mathbf{X}^T)^N + a_0 \end{aligned} \quad (8)$$

which is the positive sum of a series of **Gram Matrix** plus a constant a_0 . Since **Gram Matrix** is always positive semi-definite, so we have that the kernel matrix is also positive semi-definite for any input data \mathbf{X} .

Since the function here gives rise to a symmetric and positive semi-definite kernel matrix \mathbf{K} , so we can say that the function k here is a valid kernel.

Problem 6

For $x_1, x_2 \in (0, 1)$, consider a kernel function in an infinite-dimensional feature space, we have

$$\lim_{n \rightarrow \infty} x_1^n = 0, \quad \lim_{n \rightarrow \infty} x_2^n = 0, \quad \lim_{n \rightarrow \infty} (x_1 x_2)^n = 0 \quad (9)$$

so use the feature transformation function

$$\begin{aligned} \phi(x) &= 1 + x^1 + x^2 + \cdots + x^n \\ &= \frac{1 - x^{n+1}}{1 - x} \\ &\approx \frac{1}{1 - x} \end{aligned} \quad (10)$$

then we have

$$\begin{aligned}k(x_1, x_2) &= \phi(x_1)^T \phi(x_2) \\&= 1 + x_1 x_2 + (x_1 x_2)^2 + \cdots + (x_1 x_2)^n \\&= \frac{1 - (x_1 x_2)^{n+1}}{1 - x_1 x_2} \\&\approx \frac{1}{1 - x_1 x_2}\end{aligned}\tag{11}$$

References

- [1] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, 1 edition, 2007.