

Writing Assignment 2

Issued: Friday 22nd October, 2021

Due: Wednesday 3rd November, 2021

2.1. (Kernel Regression Least Square) Suppose we are given a dataset $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^m$ consisting of m independent samples, where $\mathbf{x}^{(i)} \in \mathbb{R}^n$ is a n -dimension vector, and $y^{(i)} \in \mathbb{R}$. Now, we aim to learn a linear model $f(\mathbf{x}) = \boldsymbol{\theta}^T \phi(\mathbf{x})$ in a given feature space, i.e. $\phi(\mathbf{x}) : \mathcal{X} \rightarrow \mathcal{F}$, with regularization term $\lambda \|\boldsymbol{\theta}\|_2^2$. The loss function of the linear regression problem can be given as,

$$\sum_{i=1}^m [y^{(i)} - (\boldsymbol{\theta}^T \phi(\mathbf{x}^{(i)}))]^2 + \lambda \|\boldsymbol{\theta}\|_2^2. \quad (1)$$

- (a) (2 points) Prove that the optimal parameter $\boldsymbol{\theta}^*$ is in the span of features $\{\phi(\mathbf{x}^{(i)})\}_{i=1}^m$, i.e. $\boldsymbol{\theta}^* = \sum_{i=1}^m c_i \phi(\mathbf{x}^{(i)})$, where c_i then becomes the term needed to be optimized. (Hint: Set the differentiation of the loss over $\boldsymbol{\theta}$ to 0)
- (b) (2 points) The mapping function $\phi(\cdot)$ can often result to a high-dimensional or infinite feature $\phi(\mathbf{x})$. Thus, we adopt a kernel $\mathbf{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \stackrel{\text{def}}{=} \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle$ to make the calculation easier. Based on (a), we know that

$$f(\mathbf{x}) = \boldsymbol{\theta}^T \phi(\mathbf{x}) = \sum_{i=1}^m c_i \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}) \rangle \quad (2)$$

$$\|\boldsymbol{\theta}\|_2^2 = \left\langle \sum_{i=1}^m c_i \phi(\mathbf{x}^{(i)}), \sum_{j=1}^m c_j \phi(\mathbf{x}^{(j)}) \right\rangle = \mathbf{c}^T \mathbf{K} \mathbf{c}, \quad (3)$$

where $\mathbf{c} \stackrel{\text{def}}{=} [c_1, \dots, c_m]^T$ and $\mathbf{K} \in \mathbb{R}^{m \times m}$ with the (i, j) -th entry defined as $\mathbf{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$. Now please rewrite the loss function (1) using \mathbf{c} and \mathbf{K} , and give the optimal parameter \mathbf{c}^* .

2.2. (Least-Squares SVM) Suppose we are given a training dataset $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^m$ consisting of m independent examples, where $\mathbf{x}^{(i)} \in \mathbb{R}^n$ is n -dimension vector, and $y^{(i)} \in \{-1, 1\}$. The Least-Squares Support Vector Machine (LS-SVM) aims to construct a linear model $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$ in a given feature space, i.e. $\phi(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{F}$, that is able to distinguish between examples drawn from different categories \mathcal{C}^- and \mathcal{C}^+ , such that

$$\mathbf{x} \in \begin{cases} \mathcal{C}^+, & f(\mathbf{x}) \geq 0 \\ \mathcal{C}^-, & o.w. \end{cases}.$$

The optimal model parameters (\mathbf{w}^*, b^*) are given by solving a constrained optimization problem,

$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{1}{2\mu} \sum_{i=1}^m \epsilon_i^2, \\ & \text{subject to} && y_i = \mathbf{w}^T \phi(\mathbf{x}_i) + b + \epsilon_i, \quad i = 1, \dots, m, \end{aligned} \quad (4)$$

where μ is a regularization hyper-parameter. The primal Lagrangian for this optimisation problem (4) gives the unconstrained minimisation problem,

$$\mathcal{L} = \frac{1}{2} \|\mathbf{w}\|_2^2 + \frac{1}{2\mu} \sum_{i=1}^m \epsilon_i^2 - \sum_{i=1}^m \alpha_i [\mathbf{w}^T \phi(\mathbf{x}_i) + b + \epsilon_i - y_i], \quad (5)$$

where $\boldsymbol{\alpha} \stackrel{\text{def}}{=} [\alpha_1, \dots, \alpha_m]^T$ is a vector of Lagrange multipliers.

(a) (1 point) Give the KKT optimality conditions for this problem.

(Hint: Set $\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial \epsilon_i} = \frac{\partial \mathcal{L}}{\partial \alpha_i} = 0$)

(b) (2 points) Denoting that $\mathbf{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \stackrel{\text{def}}{=} \langle \phi(\mathbf{x}^{(i)}), \phi(\mathbf{x}^{(j)}) \rangle$, prove that

$$\begin{bmatrix} \mathbf{K} + \mu \mathbf{I} & \mathbf{1} \\ \mathbf{1}^T & 0 \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha}^* \\ b^* \end{bmatrix} = \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix}.$$

(c) (1 point) Let $\mathbf{M} \stackrel{\text{def}}{=} \mathbf{K} + \mu \mathbf{I}$, prove that

$$\boldsymbol{\alpha}^* = \mathbf{M}^{-1}(\mathbf{y} - b^* \mathbf{1}), \quad b^* = \frac{\mathbf{1}^T \mathbf{M}^{-1} \mathbf{y}}{\mathbf{1}^T \mathbf{M}^{-1} \mathbf{1}},$$

where $\mathbf{y} \stackrel{\text{def}}{=} [y_1, \dots, y_m]^T$ and $\mathbf{1} \stackrel{\text{def}}{=} [1, \dots, 1]^T$.

2.3. (Kernel SVM) Suppose we are given a training dataset $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^m$ consisting of m independent examples, where $\mathbf{x}^{(i)} \in \mathbb{R}^n$ is n -dimension vector, and $y^{(i)} \in \{-1, +1\}$. When the data are not linearly separable, consider the Kernel-SVM given by

$$\begin{aligned} & \underset{\mathbf{w}, b}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|_2^2 \\ & \text{subject to} && y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1, \quad i = 1, \dots, m, \end{aligned} \quad (6)$$

where $\phi(\mathbf{x})$ is a mapping function $\phi(\mathbf{x}) : (x_1, x_2) \mapsto (x_1^2, \sqrt{2}x_1x_2, x_2^2)$.

(a) (1 point) Prove that $\mathbf{K}(\mathbf{x}_i, \mathbf{x}_j) \stackrel{\text{def}}{=} \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$ is positive semi-definite symmetric, i.e. for any vector $\mathbf{v} \in \mathbb{R}^m$, $\mathbf{v}^T \mathbf{K} \mathbf{v} \geq 0$.

(b) (2 points) Given data set $\{((1, \sqrt{2})^T, 1), ((\sqrt{2}, 1)^T, 1), ((2, \sqrt{2})^T, -1)\}$, derive the optimal value of \mathbf{w}^* and b^* in (6).

(c) (1 point) In (b), for new sample $(4\sqrt{2}, 1)^T$, make your decision of classification.