

CS181 Practice Questions: Max-Margin Classification

When maximizing the margin, we seek to learn linear functions of the form

$$f(\mathbf{x}, \mathbf{w}, b) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}) + b$$

where \mathbf{w} is an M -dimensional column vector of weights and $\boldsymbol{\phi}(\mathbf{x})$ is a collection of feature maps (like the regression and neural network case). The training data set comprise of N input vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ with the corresponding target labels t_1, t_2, \dots, t_N , where $t_n \in \{-1, +1\}$.

1. Computing the Margin

What is the perpendicular distance from a data point \mathbf{x}_n to the decision boundary $y_w(\mathbf{x})$?

$$\frac{t_n(\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

2. Basic Maximization Problem

What is the optimization problem we write for maximizing the margin?

$$\boldsymbol{w}^*, b^* = \arg \max_{\boldsymbol{w}, b} \left\{ \frac{1}{\|\boldsymbol{w}\|} \min_n \left[t_n(\boldsymbol{w}^\top \boldsymbol{\phi}(x_n) + b) \right] \right\}$$

3. Constrained Minimization

What is the corresponding constrained quadratic minimization problem for maximizing the margin?

$$\boldsymbol{w}^*, b^* = \arg \min_{\boldsymbol{w}, b} \left\{ \frac{1}{2} \|\boldsymbol{w}\|_2^2 \right\} \text{ s.t. } t_n(\boldsymbol{w}^\top \boldsymbol{\phi}(x_n) + b) \geq 1$$

4. Equivalence

Explain (at a high level) why the constrained quadratic minimization of question (3) is equivalent to the unconstrained maximization in (2)

From above, we have that the max-margin problem is given by

$$\mathbf{w}^*, b^* = \arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_n [t_n(\mathbf{w}^\top \boldsymbol{\phi}(x_n) + b)] \right\}$$

and that the corresponding minimization problem is

$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} \left\{ \frac{1}{2} \|\mathbf{w}\|_2^2 \right\} \text{ s.t. } t_n(\mathbf{w}^\top \boldsymbol{\phi}(x_n) + b) \geq 1$$

Minimizing $\frac{1}{2} \|\mathbf{w}\|^2$ is equivalent to maximizing $\frac{1}{\|\mathbf{w}\|}$ because $\|\mathbf{w}\| \geq 0$. Furthermore, we note that the normalized orthogonal distance from a point to the decision boundary is invariant under scalar multiplication. To see this, we have

$$\frac{t_n(\mathbf{w}^\top \boldsymbol{\phi}(x_n) + b)}{\|\mathbf{w}\|^2} = \frac{\beta}{\beta} \cdot \frac{t_n(\mathbf{w}^\top \boldsymbol{\phi}(x_n) + b)}{\|\mathbf{w}\|^2} = \frac{t_n(\beta \mathbf{w}^\top \boldsymbol{\phi}(x_n) + (\beta b))}{\|\beta \mathbf{w}\|^2}$$

Thus, since the data is linearly separable, so there exists a decision boundary with non-zero, positive margin for each example, we do not lose any generality in imposing the restraint $t_n(\mathbf{w}^\top \boldsymbol{\phi}(x_n) + b) \geq 1$ (because we can just scale \mathbf{w} until the minimal value is ≥ 1).

5. Tightness

What happens to the inequalities $t_n(\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n) + b) \geq 1$ for the optimal solution?

At the optimal solution, some of the inequalities must be tight, i.e., $t_n(\mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n) + b) = 1$. If it wasn't, we could always choose a \mathbf{w} with a smaller norm that still satisfied the constraint.

6. Kernels

What is kernel function?

A kernel function is a scalar product on two vectors mapped by basis functions into a feature space. In general, we use kernels to map into higher dimensional feature spaces, using them to circumvent costly computations in high dimension spaces.