

Machine Learning Exercise 5

何舜成

2012011515

1. Show KKT condition is necessary condition and sufficient in some cases:

Proof:

A standard form of optimal problem can be described as follow

$$\begin{aligned} \min \quad & f_0(x) \\ \text{s.t.} \quad & f_i(x) \leq 0, i = 1, \dots, m \\ & h_i(x) = 0, i = 1, \dots, p \end{aligned}$$

where $x \in \mathcal{R}^n$ and the domain $\mathcal{D} = \bigcap_{i=1}^m \text{dom} f_i \cap \bigcap_{i=1}^p \text{dom} h_i$ not empty. And the optimal solution is p^* .

Define Lagrangian function

$$L(x, \lambda, \mu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \mu_i h_i(x)$$

where λ and μ are called Lagrangian multiplier.

Define Lagrangian dual function

$$g(\lambda, \mu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \mu) = \inf_{x \in \mathcal{D}} (f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \mu_i h_i(x))$$

Suppose \tilde{x} is feasible to primal problem, then for all $\lambda \geq 0$ and μ

$$\sum_{i=1}^m \lambda_i f_i(\tilde{x}) + \sum_{i=1}^p \mu_i h_i(\tilde{x}) \leq 0$$

$$L(\tilde{x}, \lambda, \mu) = f_0(\tilde{x}) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \mu_i h_i(x) \leq f_0(\tilde{x})$$

therefore

$$g(\lambda, \mu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \mu) \leq L(\tilde{x}, \lambda, \mu) \leq f_0(\tilde{x})$$

Consider the Lagrangian dual problem

$$\begin{aligned} \max \quad & g(\lambda, \mu) \\ \text{s.t.} \quad & \lambda_i \geq 0, i = 1, \dots, m \end{aligned}$$

the dual problem is convex optimal problem. If (λ^*, μ^*) is the optimal solution and the target function reaches d^* at (λ^*, μ^*) , we easily get this inequality

$$d^* \leq p^*$$

and we call this weak duality. Likewise, if

$$d^* = p^*$$

holds true, we call this strong duality.

In condition of strong duality and the existence of optimal solution x^* of primal problem, we have

$$f_0(x^*) = g(\lambda^*, \mu^*) \tag{1}$$

$$= \inf_x (f_0(x) + \sum_{i=1}^m \lambda_i^* f_i(x) + \sum_{i=1}^p \mu_i^* h_i(x)) \tag{2}$$

$$\leq f_0(x^*) + \sum_{i=1}^m \lambda_i^* f_i(x^*) + \sum_{i=1}^p \mu_i^* h_i(x^*) \tag{3}$$

$$\leq f_0(x^*) \tag{4}$$

Therefore the inequalities are forced to be equal. We can infer

$$\left. \frac{\partial L(x, \lambda^*, \mu^*)}{\partial x} \right|_{x=x^*} = 0$$

from the third equation (supposing f_i and h_i are differentiable), and infer

$$\sum_{i=1}^m \lambda_i^* f_i(x^*) = 0$$

or $\lambda_i^* f_i(x^*) = 0, i = 1, \dots, m$ from the fourth equation.

To conclude, if x^* and (λ^*, μ^*) are optimal solutions of the primal and dual problem respectively, and strong duality is sufficed, KKT condition

(1) x^* is primal feasible

$$f_i(x) \leq 0, i = 1, \dots, m$$

$$h_i(x) = 0, i = 1, \dots, p$$

(2) (λ^*, μ^*) are dual feasible

$$\lambda^* \geq 0, i = 1, \dots, m$$

(3) complementary slackness

$$\lambda^* f_i(x^*) = 0, i = 1, \dots, m$$

(4) stationary

$$\nabla f_0(x^*) + \sum_{i=1}^m \lambda_i^* \nabla f_i(x^*) + \sum_{i=1}^p \mu_i^* \nabla h_i(x^*) = 0$$

will hold true simultaneously. KKT condition is **necessary**.

KKT condition is also sufficient if

$$\left. \frac{\partial L(x, \lambda^*, \mu^*)}{\partial x} \right|_{x=x^*} = 0 \Rightarrow L(x^*, \lambda^*, \mu^*) = \inf_{x \in \mathcal{D}} (x, \lambda^*, \mu^*)$$

Proof:

Combining the condition above with the (4) equation in KKT condition, we get

$$g(\lambda^*, \mu^*) = L(x^*, \lambda^*, \mu^*)$$

According to other 3 conditions, we know

$$g(\lambda^*, \mu^*) = f_0(x^*)$$

Therefore we can infer x^* and (λ^*, μ^*) are optimal solutions of primal and dual problem respectively from weak duality. KKT condition is **sufficient**.

2. Give the dual problem of SVM when linear inseparable

Slackness variables are introduced when data are linear inseparable.

The primal problem

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

The Lagrangian function is

$$L(w, b, \xi, \alpha, r) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i(w^T x_i + b) - 1 + \xi_i) - \sum_{i=1}^n r_i \xi_i$$

for all $i \in [n], \alpha_i \geq 0, r_i \geq 0$.

Set the partial derivatives to zero, and we get

$$\begin{aligned} \frac{\partial L}{\partial w} = 0 \quad & \Rightarrow \quad w = \sum_{i=1}^n \alpha_i y_i x_i \\ \frac{\partial L}{\partial b} = 0 \quad & \Rightarrow \quad \sum_{i=1}^n \alpha_i y_i = 0 \\ \frac{\partial L}{\partial \xi_i} = 0 \quad & \Rightarrow \quad C - \alpha_i - r_i = 0, \forall i \in [n] \end{aligned}$$

According to the previous exercise,

$$L(w, b, \xi, \alpha, r) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$

Therefore the dual can be described as follow

$$\begin{aligned} \max \quad & \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C, \forall i \in [n] \\ & \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

3. Design an algorithm to calculate gradient of the loss function of neural network

Mathematical model of artificial neural network:

$$x^l = f(u^l), u^l = (W^{l-1})^T x^{l-1} + b^l$$

where l denotes the current layer with the output layer designated to be layer L and the input layer designated to be layer 1. Function $f(\cdot)$ is a nonlinear function (i.e. sigmoid or hyperbolic tangent). Define the loss function as

$$E(x^L, t)$$

where x^L is the network output and t is the target output. Usually we choose

$$E(x^L, t) = \frac{1}{2} \|t - x^L\|^2$$

Since

$$E(x^L, t) = E(f((W^{L-1})^T x^{L-1}), t)$$

we can write the derivatives w.r.t. W^{L-1}

$$\frac{\partial E}{\partial W^{L-1}} = x^{L-1} (f'(u^L) \star \frac{\partial E}{\partial x^L})^T$$

where \star denotes elementwise multiplying, and if we define

$$\delta^L = f'(u^L) \star \frac{\partial E}{\partial x^L}$$

we get

$$\frac{\partial E}{\partial W^{L-1}} = x^{L-1}(\delta^L)^T$$

If we calculate the δ term recursively

$$\delta^l = f'(u^l) \star ((W^l)^T \delta^{l+1}), l = L-1, \dots, 2$$

it is easy to write

$$\frac{\partial E}{\partial W^l} = x^l(\delta^{l+1})^T, l = L-2, \dots, 1$$