

Coursework (4) for *Introductory Lectures on Optimization*

Zhao Hongjue
3190104515

Nov. 17, 2022

Exercise 1. Prove the following results. For proximal point method, if f is closed and convex and optimal value f^* is finite and attained at \mathbf{x}_* . We have

$$f(\mathbf{x}_k) - f^* \leq \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2}{2 \sum_{i=0}^{k-1} t_i}, \quad \text{for } k \geq 1.$$

Proof of Exercise 1:

Lemma 1. $\mathbf{u} = \text{prox}_h(\mathbf{x})$ is equivalent to the following

1. $\mathbf{x} - \mathbf{u} \in \partial h(\mathbf{u})$,
2. $h(\mathbf{z}) \geq h(\mathbf{u}) + (\mathbf{x} - \mathbf{u})^\top (\mathbf{z} - \mathbf{u})$ for all \mathbf{z} .

The updating rule of proximal point method can be concluded as

$$\mathbf{x}' = \text{prox}_{tf}(\mathbf{x}).$$

According to Lemma. 1, we can obtain

$$f(\mathbf{z}) \geq f(\mathbf{x}') + \frac{1}{t}(\mathbf{x} - \mathbf{x}')^\top (\mathbf{z} - \mathbf{x}')$$

Let $\mathbf{z} = \mathbf{x}$, we can obtain

$$f(\mathbf{x}') \leq f(\mathbf{x}) - \frac{1}{t}\|\mathbf{x} - \mathbf{x}'\|_2^2.$$

Therefore, this algorithm is a descent method. Let $\mathbf{z} = \mathbf{x}_*$, we can get

$$\begin{aligned} f(\mathbf{x}') - f(\mathbf{x}_*) &\leq \frac{1}{t}(\mathbf{x} - \mathbf{x}')^\top (\mathbf{x}' - \mathbf{x}_*) \\ &= \frac{1}{t}[(\mathbf{x} - \mathbf{x}_*) - (\mathbf{x}' - \mathbf{x}_*)]^\top (\mathbf{x}' - \mathbf{x}_*) \\ &= \frac{1}{t}[(\mathbf{x} - \mathbf{x}_*)^\top (\mathbf{x}' - \mathbf{x}_*) - \|\mathbf{x}' - \mathbf{x}_*\|_2^2] \\ &\leq \frac{1}{t}[\|\mathbf{x} - \mathbf{x}_*\|_2 \|\mathbf{x}' - \mathbf{x}_*\|_2 - \|\mathbf{x}' - \mathbf{x}_*\|_2^2] \\ &\leq \frac{1}{t} \left[\frac{\|\mathbf{x} - \mathbf{x}_*\|_2^2 + \|\mathbf{x}' - \mathbf{x}_*\|_2^2}{2} - \|\mathbf{x}' - \mathbf{x}_*\|_2^2 \right] \\ &= \frac{1}{2t} (\|\mathbf{x} - \mathbf{x}_*\|_2^2 - \|\mathbf{x}' - \mathbf{x}_*\|_2^2) \end{aligned}$$

Let $t = t_i$, $\mathbf{x} = \mathbf{x}_i$, $\mathbf{x}' = \mathbf{x}_{i+1}$. For $i = 0, \dots, k-1$

$$t_i(f(\mathbf{x}_{i+1}) - f(\mathbf{x}_*)) \leq \frac{1}{2} \left(\|\mathbf{x}_i - \mathbf{x}_*\|_2^2 + \|\mathbf{x}_{i+1} - \mathbf{x}_*\|_2^2 \right).$$

Adding inequalities from $i = 0$ to $k-1$ gives

$$\left(\sum_{i=0}^{k-1} t_i \right) (f(\mathbf{x}_k) - f(\mathbf{x}_*)) \leq \sum_{i=0}^{k-1} t_i (f(\mathbf{x}_{i+1}) - f(\mathbf{x}_*)) \leq \frac{1}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2.$$

Thus we have $f(\mathbf{x}_k) - f^* \leq (\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2) / (2 \sum_{i=0}^{k-1} t_i)$ for $k \geq 1$. □

Excercise 2. Derive the the dual problem of hard margin SVM.

Solution of Excercise 2: First we derive the primal optimization problem of hard margin SVM. For hard margin SVM, it assumes the existence of a hyperplane that perfectly separates the training sample into two populations of positively and negatively labeled points.

However, there are then infinitely many such separating hyperplanes. Consider the training dataset $\{\mathbf{x}_i, y_i\}_{i=1}^m$ for $\mathbf{x} \in \mathcal{X} \in \mathbb{R}^N$ and $y_i \in \mathcal{Y} = \{+1, -1\}$, we define the geometric margin of linear classifier $h : \mathbf{x} \mapsto \mathbf{w}^\top \mathbf{x} + b$ at \mathbf{x} as follows:

$$\rho_h(\mathbf{x}) = \frac{|\mathbf{w}^\top \mathbf{x} + b|}{\|\mathbf{w}\|_2}.$$

Then we select hyperplane which can maximize the margin ρ :

$$\rho = \max_{\mathbf{w}, b} \min_{i \in [m]} \frac{y_i(\mathbf{w}^\top \mathbf{x}_i + b)}{\|\mathbf{w}\|_2}.$$

Observe that the last expression is invariant to multiplication (\mathbf{w}, b) by a positive scalar. Thus we restrict ourselves to pairs (\mathbf{w}, b) scaled such that $\min_{i \in [m]} y_i(\mathbf{w}^\top \mathbf{x}_i + b) = 1$:

$$\rho = \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|_2}.$$

Since maximizing $1/\|\mathbf{w}\|_2$ is equivalent to minimizing $\frac{1}{2}\|\mathbf{w}\|_2^2$, the primal optimization problem can be concluded as

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|_2^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1, \quad i = 1, \dots, m \end{aligned} \tag{1}$$

The Lagrangian of problem. 1 can then be defined for all $\mathbf{w} \in \mathbb{R}^N$, $b \in \mathbb{R}$, and $\boldsymbol{\alpha} \in \mathbb{R}_+^m$, by

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|_2^2 - \sum_{i=1}^m \alpha_i [y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1]. \tag{2}$$

The KKT conditions are obtained by setting the gradient of the Lagrangian with respect to the primal variables \mathbf{w} and b to zero and by writing the complementarity conditions:

$$\nabla_{\mathbf{w}} L = \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i = 0 \implies \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \tag{3}$$

$$\nabla_b L = - \sum_{i=1}^m \alpha_i y_i = 0 \implies \sum_{i=1}^m \alpha_i y_i = 0 \tag{4}$$

$$\forall i, \alpha_i[y_i(\mathbf{w}^\top \mathbf{x} + b) - 1] = 0 \implies \alpha_i = 0 \vee y_i(\mathbf{w}^\top \mathbf{x} + b) = 1 \quad (5)$$

To derive the dual form of the constrained optimization problem 1, we plug into the Lagrangian the definition of \mathbf{w} in terms of the dual variables as expressed in (3) and apply the constraint (4). This yields

$$\begin{aligned} L &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 - \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^\top \mathbf{x}_j) - \sum_{i=1}^m \alpha_i y_i b + \sum_{i=1}^m \alpha_i \\ &= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^\top \mathbf{x}_j). \end{aligned}$$

This leads to the following dual optimization problem for hard margin SVM:

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^\top \mathbf{x}_j) \\ \text{s.t.} \quad & \alpha_i \geq 0 \wedge \sum_{i=1}^m \alpha_i y_i = 0, \quad i = 1, \dots, m. \end{aligned}$$

Obviously, the objective function is infinitely differentiable. Since the constraints are affine and convex, this dual problem is a convex optimization problem. According to the KKT conditions, we can obtain the solution of the primal optimization problem. \square

Exercise 3. For KL divergence defined on the probability simplex, prove that the upper bound of $\Delta_\psi(x^*, x_1)$ is $\log n$, for $x_1 = [\frac{1}{n}, \dots, \frac{1}{n}]$.

Proof of Exercise 3:

$$\begin{aligned} \Delta_\psi(\mathbf{x}^*, \mathbf{x}_1) &= \sum_{i=1}^n x_i^* \log \frac{x_i^*}{x_{1i}} \\ &= \sum_{i=1}^n x_i^* \log(n x_i^*) \\ &= \sum_{i=1}^n x_i^* (\log n + \log x_i^*) \\ &= \log n \underbrace{\sum_{i=1}^n x_i^*}_{=1} + \underbrace{\sum_{i=1}^n x_i^* \log x_i^*}_{\leq 0} \\ &\leq \log n. \end{aligned}$$

\square