
Homework 4 – Due by 11:59 PM on Sunday April 6

Instruction: For questions 1-3, you need to use the definition of subgradient. For question 4, use the first order optimality condition of prox. For question 5 (b,c,d), you need to submit your MATLAB (.m) files or Python (.py or .txt) files. When we run your code, the desired output should be displayed; otherwise, you lose points.

(1) Show that $\partial(\lambda f(x)) = \lambda \partial f(x)$ for any $\lambda > 0$.

(2) For a convex function f and subgradients $g_x \in \partial f(x)$, $g_y \in \partial f(y)$, prove that

$$(g_x - g_y)^T(x - y) \geq 0.$$

(3) x^* is the global minimum of f , iff $0 \in \partial f(x^*)$.

(4) Show that x^* minimizes a non-differentiable function f if and only if $x^* = \text{prox}_f(x^*)$.

(5) Let $g(x) = \sum_{i=1}^n \log(1 + \exp(-a_i^T x))$ and consider the following logistic regression problem:

$$\min_{x \in \mathbb{R}^m} f(x) \triangleq g(x) + \lambda \|x\|_1,$$

Note that $\nabla g(x) = \sum_{i=1}^n \frac{-a_i(\exp(-a_i^T x))}{1 + \exp(-a_i^T x)}$.

(a) Show that Lipschitz constant of $\nabla g(x)$ is $L = \frac{\|\sum_{i=1}^n a_i a_i^T\|}{4} = \frac{\|A\|^2}{4}$, where $A = [a_i^T]_{i=1}^n \in \mathbb{R}^{n \times m}$.

Problem Setup for parts (b)-(d). Fix the seed to 123. Generate $A = [a_i^T]_{i=1}^n \in \mathbb{R}^{n \times m}$ randomly with normal distribution, i.e., $A = \text{randn}(n, m)$; . Let $n = 500$, $m = 100$, $\lambda = 10^{-3}$, $x_0 = \mathbf{0} \in \mathbb{R}^m$ and set the total number of iterations as $\text{maxiter} = 200$.

(b) Solve the problem by subgradient method with diminishing stepsize $\alpha_k = 1/\sqrt{k}$. The output of the algorithm should be the best objective value.

(c) Solve the problem by proximal gradient method with stepsize $\alpha = 1/L$. The output of the algorithm should be the objective value of the last iterate.

(d) Solve the problem by FISTA with stepsize $\alpha = 1/L$. The output of the algorithm should be the objective value of the last iterate.

(e) Compare the output of the three algorithms. What can you conclude? i.e., Which one has the fastest and which one has the slowest convergence?

Students enrolled in SIE 549 must solve the following problem.
Students in SIE 449 will get extra credit by solving it.

(6) Let f be a continuously differentiable convex function over a closed and convex set $C \subseteq \mathbb{R}^n$. Show that if $x^* \in C$ is an optimal solution of

$$\min \{f(x) \mid x \in C\}$$

then $\nabla_x f(x)^T(x^* - x) \leq 0$ for all $x \in C$.

Hint: Use that fact that if f is convex then $(\nabla f(y) - \nabla f(x))^T(x - y) \leq 0$ for any $x, y \in C$.

Extra Credit Question. Consider the problem

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & g(x) \leq 0 \\ & x \in X, \end{aligned} \tag{P}$$

where f and g are convex functions over \mathbb{R}^n and $X \subseteq \mathbb{R}^n$ is a convex set. Suppose x^* is an optimal solution of (P) that satisfies $g(x^*) < 0$. Show that x^* is an optimal solution of

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & x \in X. \end{aligned} \tag{Q}$$

Hint: Prove by contradiction.

$$(1) \partial f(\lambda x) = \lambda \partial f(x) \text{ for } \lambda > 0$$

$$\text{let } g \in \partial(\lambda f(x))$$

$$\lambda f(y) \geq \lambda f(x) + g^T(y-x) \quad \forall y$$

$$f(y) \geq f(x) + \left(\frac{g}{\lambda}\right)^T(y-x)$$

$$\frac{g}{\lambda} \in \partial f(x) \Rightarrow g \in \lambda \partial f(x)$$

$$\text{let } h \in \lambda \partial f(x), \text{ then } \frac{h}{\lambda} \in \partial f(x)$$

$$f(y) \geq f(x) + \left(\frac{h}{\lambda}\right)^T(y-x)$$

$$\lambda f(y) \geq \lambda f(x) + h^T(y-x)$$

$$h \in \partial(\lambda f(x))$$

$$\partial(\lambda f(x)) = \lambda(\partial f(x))$$

$$(2) \text{ Since } g_x \in \partial f(x), f(y) \geq f(x) + g_x^T(y-x)$$

$$\text{Since } g_y \in \partial f(y), f(x) \geq f(y) + g_y^T(x-y)$$

$$f(y) \geq f(x) + g_x^T(y-x)$$

$$+ f(x) \geq f(y) + g_y^T(x-y)$$

$$\underline{f(x) + f(y) \geq f(x) + f(y) + g_y^T(x-y) + g_x^T(y-x)}$$

$$0 \geq g_y^T(x-y) + g_x^T(y-x)$$

$$= -g_y^T(y-x) + g_x^T(y-x)$$

$$= (g_x - g_y)^T(y-x)$$

$$0 \leq (g_x - g_y)^T(x-y)$$

$$(3) \Rightarrow \text{if } x^* \text{ is optimal condition, then } f(y) \geq f(x^*) \text{ for } \forall y$$

$$\text{let } 0 \text{ be a vector in } \partial f(x^*),$$

$$\text{then } f(y) \geq f(x^*) + 0^T(y-x^*), \forall y$$

$$\text{and } x^* \text{ is global min.}$$

$$\Leftarrow \text{suppose } 0 \in \partial f(x^*)$$

$$\text{then for all } y, f(y) \geq f(x^*) + 0^T(y-x^*)$$

$$\text{and } x^* \text{ is global min.}$$

$$(4)$$

$$x^* \text{ minimizes } f(x) \text{ if and only if } 0 \in \partial f(x^*)$$

$$\text{prox}_f(v) = \arg\min_x \{f(x) + \frac{1}{2}\|x-v\|^2\}$$

$$\text{Setting } v \text{ to } x^*,$$

$$\text{prox}_f(x^*) = \arg\min_x \{f(x) + \frac{1}{2}\|x-x^*\|^2\}$$

$$\text{optimality condition:}$$

$$0 \in f(x^*) + (x^* - x^*)$$

$$0 \in \partial f(x^*)$$

$$x^* = \text{prox}_f(x^*) \text{ iff } 0 \in \partial f(x^*) \text{ iff } x^* \text{ minimizes } f(x)$$

$$(5) \nabla g(x) = \sum_{i=1}^n \frac{-a_i \exp(-a_i^T x)}{1 + \exp(-a_i^T x)} = \sum_{i=1}^n -a_i \sigma(-a_i^T x), \text{ where } \sigma(z) = \frac{e^z}{1+e^z} \text{ is sigmoid}$$

$$\nabla^2 g(x) = \sum_{i=1}^n a_i a_i^T \sigma(-a_i^T x) (1 - \sigma(-a_i^T x))$$

$$\text{Since } 0 \leq \sigma(z) \leq 1:$$

$$0 \leq \sigma(z) + (1 - \sigma(z)) \leq \frac{1}{4}$$

$$\text{Therefore:}$$

$$\|\nabla^2 g(x)\| \leq \sum_{i=1}^n \|a_i a_i^T\| \frac{1}{4} = \frac{1}{4} \sum_{i=1}^n a_i a_i^T$$

$$\text{Lipschitz Functions}$$

$$L = \frac{1}{4} \left\| \sum_{i=1}^n a_i a_i^T \right\| = \frac{\|A^T A\|}{4}$$