

# HW5

2020年11月25日 4:46



HW5Qu...

# MATH 3332 Data Analytic Tools

## Homework 5

Due date: 30 November, 6pm, Monday

1. Consider the following optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{Bx}\|_2^2,$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ , and  $\mathbf{B} \in \mathbb{R}^{p \times n}$  are given. Find the steepest descent (i.e. gradient descent with the exact line search) algorithm for solving this optimization.

In the following two problems, we consider variants of the gradient descent algorithm for solving

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad (1)$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a smooth function. Given  $\mathbf{x}^{(k)}$ , the gradient descent algorithm finds  $\mathbf{x}^{(k+1)}$  by

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}^{(k)}) + \langle \nabla f(\mathbf{x}^{(k)}), \mathbf{x} - \mathbf{x}^{(k)} \rangle + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2, \quad (2)$$

where  $\alpha_k > 0$  is a step size.

2. In the gradient descent (2), for distance between  $\mathbf{x}$  and  $\mathbf{x}^{(k)}$ , we used the 2-norm, which can be replaced by other norms to obtain variants of the gradient descent algorithm.

- (a) Prove that (2) is equivalent to

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \langle \nabla f(\mathbf{x}^{(k)}), \mathbf{x} - \mathbf{x}^{(k)} \rangle, \quad \text{subject to } \|\mathbf{x} - \mathbf{x}^{(k)}\|_2 \leq \alpha_k \|\nabla f(\mathbf{x}^{(k)})\|_2. \quad (3)$$

*(Hint: Use Cauchy-Schwartz inequality.)*

- (b) We may replace the 2-norm in (3) by 1-norm, i.e.,

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \langle \nabla f(\mathbf{x}^{(k)}), \mathbf{x} - \mathbf{x}^{(k)} \rangle, \quad \text{subject to } \|\mathbf{x} - \mathbf{x}^{(k)}\|_1 \leq \alpha_k \|\nabla f(\mathbf{x}^{(k)})\|_\infty. \quad (4)$$

Give an explicit expression of  $\mathbf{x}^{(k+1)}$  in (4). *(Hint: Use the inequality  $|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\|_\infty \|\mathbf{y}\|_1$  with equality attained if  $\mathbf{x} = c \cdot \text{sign}(\mathbf{y})$  or  $\mathbf{y} = c \cdot \mathbf{e}_{i_{\max}(\mathbf{x})}$  for some  $c \in \mathbb{R}$ . Here  $\text{sign}(\cdot)$  takes the sign of each entry of a vector, and  $i_{\max}(\mathbf{x}) = \arg \max_i |x_i|$ .)*

- (c) We may replace the 2-norm in (3) by  $\infty$ -norm, i.e.,

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \langle \nabla f(\mathbf{x}^{(k)}), \mathbf{x} - \mathbf{x}^{(k)} \rangle, \quad \text{subject to } \|\mathbf{x} - \mathbf{x}^{(k)}\|_\infty \leq \alpha_k \|\nabla f(\mathbf{x}^{(k)})\|_1. \quad (5)$$

Give an explicit expression of  $\mathbf{x}^{(k+1)}$  in (5).

3. Gradient depends on inner product. Therefore, one variant of gradient descent is obtained by choosing a weighted inner product of  $\mathbb{R}^n$  rather than the standard one. In particular, given a symmetric positive definite matrix  $\mathbf{W} \in \mathbb{R}^{n \times n}$ , we know that  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{W}} := \mathbf{x}^T \mathbf{W} \mathbf{y}$  defines an inner product of  $\mathbb{R}^n$ .
- Express  $\nabla_{\mathbf{W}} f$  (the gradient of  $f$  in  $\mathbb{R}^n$  with the weighted inner product  $\langle \cdot, \cdot \rangle_{\mathbf{W}}$ ) in terms of  $\mathbf{W}$  and  $\nabla f$  (the gradient of  $f$  in  $\mathbb{R}^n$  with the standard inner product).
  - Give the gradient descent algorithm for solving (1) under the weighted inner product  $\langle \cdot, \cdot \rangle_{\mathbf{W}}$  (i.e., find an explicit formula of  $\mathbf{x}^{(k+1)}$  when we replace  $\langle \cdot, \cdot \rangle$  and  $\|\cdot\|_2^2$  in (2) by  $\langle \cdot, \cdot \rangle_{\mathbf{W}}$  and  $\|\cdot\|_{\mathbf{W}}^2$  respectively). (*By choosing a suitable  $\mathbf{W}$ , we may obtain faster algorithms than the standard gradient descent. This technique is known as preconditioning.*)

1. Consider the following optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{Bx}\|_2^2,$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ , and  $\mathbf{B} \in \mathbb{R}^{p \times n}$  are given. Find the steepest descent (i.e. gradient descent with the exact line search) algorithm for solving this optimization.

In the following two problems, we consider variants of the gradient descent algorithm for solving

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad (1)$$

where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a smooth function. Given  $\mathbf{x}^{(k)}$ , the gradient descent algorithm finds  $\mathbf{x}^{(k+1)}$  by

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}^{(k)}) + \langle \nabla f(\mathbf{x}^{(k)}), \mathbf{x} - \mathbf{x}^{(k)} \rangle + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2, \quad (2)$$

where  $\alpha_k > 0$  is a step size.

$$\text{let } f(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{Bx}\|_2^2.$$

We need calling 1 & 2 here to solve  $\nabla f(\mathbf{x})$ .

Calling 1: let  $\varphi(\mathbf{x}) = \mathbf{c}^\top \mathbf{x}$ ,  $\nabla \varphi(\mathbf{x}) = \mathbf{c}$

Calling 2: let  $\phi(\mathbf{x}) = \mathbf{x}^\top \mathbf{Bx}$ ,  $\nabla \phi(\mathbf{x}) = (\mathbf{B} + \mathbf{B}^\top) \mathbf{x}$

Proving calling 1:

$$\varphi(\mathbf{x}) = \sum_{j=1}^n c_j x_j$$

$$\frac{\partial \varphi}{\partial x_k} = \sum_{j=1}^n c_j \frac{\partial x_j}{\partial x_k} = \sum_{j=1}^n c_j \delta_{jk} = c_k$$

$$\Rightarrow \nabla \varphi(\mathbf{x}) = \mathbf{c}.$$

Proving calling 2:

$$\phi(\mathbf{x}) = \sum_{i=1}^n \sum_{j=1}^n b_{ij} x_i x_j$$

$$\frac{\partial \phi}{\partial x_k} = \sum_{i=1}^n \sum_{j=1}^n b_{ij} \frac{\partial (x_i x_j)}{\partial x_k}$$

$$= \sum_{i=1}^n \sum_{j=1}^n b_{ij} (s_{ik} x_j + x_i s_{jk})$$

$$= \sum_{j=1}^n b_{kj} x_j + \sum_{i=1}^n b_{ik} x_i$$

$$= (\mathbf{Bx})_k + \sum_{i=1}^n (\mathbf{B}^\top)_{ki} x_i$$

$$= (\mathbf{Bx})_k + (\mathbf{B}^\top \mathbf{x})_k$$

$$\therefore \nabla \phi(\mathbf{x}) = (\mathbf{B} + \mathbf{B}^\top) \mathbf{x},$$

$$\text{and also } \|\mathbf{Ax} - \mathbf{b}\|_2^2 = [(\mathbf{Ax} - \mathbf{b})^\top (\mathbf{Ax} - \mathbf{b})]$$

$$\text{Consider } \|Ax-b\|_2^2 = [(Ax-b)^T(Ax-b)]$$

$$\|Ax-b\|_2^2 = [x^T A^T A x - (Ax)^T b - b^T Ax + b^T b]$$

$$\|Ax-b\|_2^2 = [x^T A^T A x - 2b^T Ax + b^T b]$$

$$\|Ax-b\|_2^2 = [x^T A^T A x - 2(A^T b)^T x + b^T b]$$

$$\text{Consider } \|Bx\|^2,$$

$$\|Bx\|_2^2 = (Bx)^T (Bx)$$

$$= x^T B^T B x = x^T (B^T B)^T x$$

$\Rightarrow$  putting  $A^T x$  in column ②, putting  $A^T b$  in column ①,

$$\text{we can get } \nabla(\|Ax-b\|_2^2) = (A^T A + (EPA)^T) x - 2A^T b$$

$$\text{and } (A^T A)^T = A^T (A^T)^T = A^T A.$$

$$\Rightarrow \nabla(\|Ax-b\|_2^2) = 2A^T A x - 2A^T b$$

$\Rightarrow$  putting  $B^T B$  in column ②, we can get  $\nabla(\|Bx\|_2^2)$

$$\nabla(\|Bx\|_2^2) = (B^T B + (B^T B)^T) x$$

$$= (B^T B + B^T B)x$$

$$= 2(B^T B)x$$

$$\therefore \nabla(f(x)) = \frac{1}{2}(2A^T A x - 2A^T b) + 2\alpha(B^T B)x$$

$$= A^T A x - A^T b + 2\alpha(B^T B)x$$

$$= (A^T A + 2\alpha(B^T B))x - A^T b$$

Also, it is obvious to check that  $f(x)$  is strictly convex and smooth function. since  $\alpha$ -norm function is continuous, strictly convex, thus  $f(x)$  is also continuous, strictly convex and derivative.

$$\therefore x^{(k+1)} = \arg \min_{x \in \mathbb{R}^n} f(x) \Leftrightarrow \nabla f(x^{(k+1)}) = 0$$

$$\nabla f(x^{(k+1)}) = 0 \Leftrightarrow \nabla f(x^{(k)}) + \frac{1}{\alpha k} (x^{(k+1)} - x^{(k)}) = 0$$

$$\nabla f(x^{(k+1)}) = 0 \iff \nabla f(x^{(k)}) + \frac{1}{\alpha_k} (x^{(k+1)} - x^{(k)}) = 0$$

$$x^{(k+1)} = x^{(k)} - \alpha_k \nabla f(x^{(k)})$$

also,  $\min_{x \in \mathbb{R}^n} f(x)$  has a unique solution

Next, we want to perform exact line search,

$$\alpha_k = \underset{\alpha \in \mathbb{R}}{\operatorname{argmin}} f(x^{(k)} - \alpha(A^T A + 2\alpha B^T B)x^{(k)} - A^T b)$$

$$\text{let } g(\alpha) = f(x^{(k)} - \alpha(A^T A + 2\alpha B^T B)x^{(k)} - A^T b)$$

It is obvious that  $g(\alpha)$  is convex since 2-norm function is convex.  
therefore we want to solve  $g'(\alpha_k) = 0$ ,

$$\text{recall } f(x) = \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|Bx\|_2^2,$$

$$f(x) = \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|Bx\|_2^2.$$

$$\alpha_k = \underset{\alpha \in \mathbb{R}}{\operatorname{argmin}} f(x^{(k)} - \alpha(A^T A + 2\alpha B^T B)x^{(k)} - A^T b)$$

$$g(\alpha) = \frac{1}{2} \|Ax^{(k)} - \alpha(AA^T A)x^{(k)} - 2\alpha A^T Bx^{(k)} + \alpha A^T b - b\|_2^2 + \lambda \|Bx^{(k)} - \alpha B^T Bx^{(k)} - 2\alpha B^T Bx^{(k)} + \alpha B^T b\|_2^2$$

$$= \frac{1}{2} \| -(\underbrace{AA^T A}_{\alpha(AA^T A + 2\alpha B^T B)x^{(k)} - A^T b)}x^{(k)} + A^T b - b\|_2^2 + \lambda \| -(\underbrace{B^T B}_{\alpha B^T Bx^{(k)} + B^T b})x^{(k)} + B^T b\|_2^2$$

$$h(\alpha) = \alpha \|A^T Ax^{(k)} + 2\alpha A^T Bx^{(k)} - A^T b\|_2^2 - (A^T Ax^{(k)} + 2\alpha A^T Bx^{(k)} - A^T b)^T (A^T Ax^{(k)} - b)$$

$$k(\alpha) = 2\alpha \|B^T Ax^{(k)} + \alpha B^T Bx^{(k)} - B^T b\|_2^2 - 2\alpha (B^T Ax^{(k)} + \alpha B^T Bx^{(k)} - B^T b)^T (B^T Ax^{(k)})$$

$$j(\alpha) = 0 \iff \alpha \left( \|AA^T Ax^{(k)} + 2\alpha A^T Bx^{(k)} - A^T b\|_2^2 + 2\alpha B^T Ax^{(k)} + 4\alpha^2 B^T Bx^{(k)} - 2\alpha B^T b \right) = \langle AA^T Ax^{(k)} + 2\alpha A^T Bx^{(k)} - A^T b, Ax^{(k)} - b \rangle + 2\alpha \langle B^T Ax^{(k)} + 2\alpha B^T Bx^{(k)} - B^T b, B^T Ax^{(k)} \rangle$$

$$\iff \alpha \left( \langle AA^T A + 2\alpha A^T B + 4\alpha^2 B^T B, Ax^{(k)} - b \rangle + 2\alpha \langle B^T A + 2\alpha B^T B, B^T Ax^{(k)} \rangle \right) = \langle AA^T Ax^{(k)} + 2\alpha A^T Bx^{(k)} - A^T b, Ax^{(k)} - b \rangle + 2\alpha \langle B^T Ax^{(k)} + 2\alpha B^T Bx^{(k)} - B^T b, B^T Ax^{(k)} \rangle$$

$$\iff \alpha \left( \langle AA^T A + 2\alpha A^T B + 4\alpha^2 B^T B, Ax^{(k)} - b \rangle + 2\alpha \langle B^T A + 2\alpha B^T B, B^T Ax^{(k)} \rangle \right) = \langle AA^T Ax^{(k)} + 2\alpha A^T Bx^{(k)} - A^T b, Ax^{(k)} - b \rangle + 2\alpha \langle B^T Ax^{(k)} + 2\alpha B^T Bx^{(k)} - B^T b, B^T Ax^{(k)} \rangle$$

$$\alpha = \frac{\langle AA^T Ax^{(k)} + 2\alpha A^T Bx^{(k)} - A^T b, Ax^{(k)} - b \rangle - \langle AA^T Ax^{(k)} + 2\alpha A^T Bx^{(k)} - A^T b, Ax^{(k)} - b \rangle}{\langle AA^T A + 2\alpha A^T B + 4\alpha^2 B^T B, Ax^{(k)} - b \rangle - \langle AA^T A + 2\alpha A^T B + 4\alpha^2 B^T B, B^T Ax^{(k)} \rangle}$$

$$= \frac{\langle (AA^T A + 2\alpha A^T B + 4\alpha^2 B^T B) x^{(k)} - (A^T A + 2\alpha A^T B) A^T b, Ax^{(k)} - b \rangle}{\langle (AA^T A + 2\alpha A^T B + 4\alpha^2 B^T B) x^{(k)} - (A^T A + 2\alpha A^T B) A^T b, B^T Ax^{(k)} \rangle}$$

$$\lambda = \frac{\langle (A^T A + 2\alpha B^T B)^2 x^{(k)}, x^{(k)} \rangle - \langle (A^T A + 2\alpha B^T B) x^{(k)}, b \rangle + b^T A^T b}{\| (A^T A + 2\alpha (B^T A + A B^T) + 4\alpha^2 B^T B) x^{(k)} - (A + 2\alpha B) A^T b \|_2^2}$$

$$\lambda = \frac{\langle (A^T A + 2\alpha B^T B)^2 x^{(k)}, x^{(k)} \rangle - \langle (A^T A + 2\alpha B^T B) x^{(k)}, b \rangle + b^T A^T b}{\| (A^T A + 2\alpha (B^T A + A B^T) + 4\alpha^2 B^T B) x^{(k)} - (A + 2\alpha B) A^T b \|_2^2}$$

$$\lambda = \frac{\| (A^T A + 2\alpha B^T B) x^{(k)} \|_2^2 - 2 \langle A (A^T A + 2\alpha B^T B) x^{(k)}, b \rangle + b^T A^T b}{\| (A^T A + 2\alpha (B^T A + A B^T) + 4\alpha^2 B^T B) x^{(k)} - (A + 2\alpha B) A^T b \|_2^2}$$

$$\lambda = \frac{\| (A^T A + 2\alpha B^T B) x^{(k)} \|_2^2 - 2 \langle A (A^T A + 2\alpha B^T B) x^{(k)}, b \rangle + b^T A^T b}{\| (A^T A + 2\alpha (B^T A + A B^T) + 4\alpha^2 B^T B) x^{(k)} - (A + 2\alpha B) A^T b \|_2^2}$$

$$\lambda = \frac{\| (A^T A + 2\alpha B^T B) x^{(k)} \|_2^2 - 2 \langle A (A^T A + 2\alpha B^T B) x^{(k)}, A^T b \rangle + \| A^T b \|_2^2}{\| (A + 2\alpha B) (A^T A + 2\alpha B^T B) x^{(k)} - (A + 2\alpha B) A^T b \|_2^2}$$

$$\lambda = \frac{\| (A^T A + 2\alpha B^T B) x^{(k)} \|_2^2 + \| (A^T A + 2\alpha B^T B) x^{(k)} - A^T b \|_2^2 - \| (A^T A + 2\alpha B^T B) x^{(k)} \|_2^2 - \| A^T b \|_2^2 + \| A^T b \|_2^2}{\| (A + 2\alpha B) ((A^T A + 2\alpha B^T B) x^{(k)} - A^T b) \|_2^2}$$

$$\lambda = \frac{\| (A^T A + 2\alpha B^T B) x^{(k)} - A^T b \|_2^2}{\| (A + 2\alpha B) ((A^T A + 2\alpha B^T B) x^{(k)} - A^T b) \|_2^2}$$

∴ steepest gradient descent algorithm:

1. Initialize  $x^{(0)}$

2. for  $k=0, 1, 2, \dots$

$$g_k = (A^T A + 2\alpha B^T B) x^{(k)} - A^T b$$

$$\alpha_k = \frac{\| g_k \|_2^2}{\| (A + 2\alpha B) g_k \|_2^2}$$

$$x^{(k+1)} = x^{(k)} - \alpha_k g_k$$

end.

2. In the gradient descent (2), for distance between  $\mathbf{x}$  and  $\mathbf{x}^{(k)}$ , we used the 2-norm, which can be replaced by other norms to obtain variants of the gradient descent algorithm.

- (a) Prove that (2) is equivalent to

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \langle \nabla f(\mathbf{x}^{(k)}), \mathbf{x} - \mathbf{x}^{(k)} \rangle, \quad \text{subject to} \quad \|\mathbf{x} - \mathbf{x}^{(k)}\|_2 \leq \alpha_k \|\nabla f(\mathbf{x}^{(k)})\|_2. \quad (3)$$

(Hint: Use Cauchy-Schwartz inequality.)

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}^{(k)}) + \langle \nabla f(\mathbf{x}^{(k)}), \mathbf{x} - \mathbf{x}^{(k)} \rangle + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2, \quad (2)$$

(Q2a). We need to show that  $x = x^{(k)} - \alpha_k \nabla f(x^{(k)})$  is the optimal solution to

$$\min_{x \in \mathbb{R}^n} \langle \nabla f(x^{(k)}), x - x^{(k)} \rangle \text{ subject to } \|x - x^{(k)}\|_2 \leq \alpha_k \|\nabla f(x^{(k)})\|_2 - (\star)$$

Using Cauchy-Schwarz inequality we have

$$|\langle \alpha_k \nabla f(x^{(k)}), x - x^{(k)} \rangle| \leq \alpha_k \|\nabla f(x^{(k)})\|_2 \|x - x^{(k)}\|_2 \leq \alpha_k^2 \|\nabla f(x^{(k)})\|_2^2$$

$$\Rightarrow \langle \nabla f(x^{(k)}), x - x^{(k)} \rangle \geq -\alpha_k \|\nabla f(x^{(k)})\|_2^2$$

$$\Rightarrow \langle \nabla f(x^{(k)}), x \rangle \geq \langle \nabla f(x^{(k)}), x^{(k)} \rangle - \alpha_k \|\nabla f(x^{(k)})\|_2^2$$

Since equality occurs only if  $\alpha_k \nabla f(x^{(k)}) = x^{(k)} - x$

i. The optimal solution to  $(\star)$  is  $x = x^{(k)} - \alpha_k \nabla f(x^{(k)})$ , the given gradient descent algorithm.

b).

$$x^{(k+1)} = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \langle \nabla f(x^{(k)}), x - x^{(k)} \rangle, \text{ subject to } \|x - x^{(k)}\|_1 \leq \alpha_k \|\nabla f(x^{(k)})\|_\infty$$

$$|\langle \alpha_k \nabla f(x^{(k)}), x - x^{(k)} \rangle| \leq \alpha_k \|\nabla f(x^{(k)})\|_\infty \|x - x^{(k)}\|_1 \leq \alpha_k^2 \|\nabla f(x^{(k)})\|_\infty^2$$

With equality attained if  $x - x^{(k)} = c \cdot e_{i_{\max}(\nabla f(x^{(k)}))}$  for some  $c \in \mathbb{R}$ .

$$\Rightarrow \langle \nabla f(x^{(k)}), x - x^{(k)} \rangle \geq -\alpha_k \|\nabla f(x^{(k)})\|_\infty^2$$

$$\Rightarrow \langle \nabla f(x^{(k)}), x \rangle \geq \langle \nabla f(x^{(k)}), x^{(k)} \rangle - \alpha_k \|\nabla f(x^{(k)})\|_\infty^2$$

Since equality occurs only if  $x - x^{(k)} = c \cdot e_{i_{\max}(\nabla f(x^{(k)}))}$  for some  $c \in \mathbb{R}$ ,

$$x^{(k+1)} = x^{(k)} + c \cdot e_{i_{\max}(\nabla f(x^{(k)}))} \text{ for some } c \in \mathbb{R}.$$

c). Similarly,  $|\langle \alpha_k \nabla f(x^{(k)}), x - x^{(k)} \rangle| \leq \alpha_k \|\nabla f(x^{(k)})\|_1 \|x - x^{(k)}\|_\infty \leq \alpha_k^2 \|\nabla f(x^{(k)})\|_1$

$$\Rightarrow \langle \nabla f(x^{(k)}), x - x^{(k)} \rangle \geq -\alpha_k \|\nabla f(x^{(k)})\|_1^2$$

$$\Rightarrow \langle \nabla f(x^{(k)}), x \rangle \geq \langle \nabla f(x^{(k)}), x^{(k)} \rangle - \alpha_k \|\nabla f(x^{(k)})\|_1^2$$

Since equality occurs only if  $x - x^{(k)} = c \cdot \operatorname{sign}(\nabla f(x^{(k)}))$  for some  $c \in \mathbb{R}$ ,

$$x^{(k+1)} = x^{(k)} + c \cdot \operatorname{sign}(\nabla f(x^{(k)})) \text{ for some } c \in \mathbb{R}.$$

3. Gradient depends on inner product. Therefore, one variant of gradient descent is obtained by choosing a weighted inner product of  $\mathbb{R}^n$  rather than the standard one. In particular, given a symmetric positive definite matrix  $\mathbf{W} \in \mathbb{R}^{n \times n}$ , we know that  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{W}} := \mathbf{x}^T \mathbf{W} \mathbf{y}$  defines an inner product of  $\mathbb{R}^n$ .
- (a) Express  $\nabla_{\mathbf{W}} f$  (the gradient of  $f$  in  $\mathbb{R}^n$  with the weighted inner product  $\langle \cdot, \cdot \rangle_{\mathbf{W}}$ ) in terms of  $\mathbf{W}$  and  $\nabla f$  (the gradient of  $f$  in  $\mathbb{R}^n$  with the standard inner product).
  - (b) Give the gradient descent algorithm for solving (1) under the weighted inner product  $\langle \cdot, \cdot \rangle_{\mathbf{W}}$  (i.e., find an explicit formula of  $\mathbf{x}^{(k+1)}$  when we replace  $\langle \cdot, \cdot \rangle$  and  $\|\cdot\|_2^2$  in (2) by  $\langle \cdot, \cdot \rangle_{\mathbf{W}}$  and  $\|\cdot\|_{\mathbf{W}}^2$  respectively). (*By choosing a suitable  $\mathbf{W}$ , we may obtain faster algorithms than the standard gradient descent. This technique is known as preconditioning.*)

Q3. From the definition of gradient, we have:

$$\lim_{x \rightarrow x_0} \frac{|f(x) - f(x_0) - \langle v, x - x_0 \rangle_w|}{\|x - x_0\|_w} = 0, \text{ where } v \text{ is gradient } \nabla f(x) \text{ at } x = x_0.$$

Consider  $\langle v, x - x_0 \rangle$ ,

$$\langle v, x - x_0 \rangle_w = v^T w (x - x_0) = \langle v, w(x - x_0) \rangle = \langle w^T v, x - x_0 \rangle$$

$$\text{If } v = \nabla w f, \quad \nabla f = w^T v$$

$$\Leftrightarrow \nabla f = w^T \nabla w f \quad \Leftrightarrow \nabla w f = (w^{-1})^T \nabla f$$

b). As  $w$  is S.P.D.,

$$\therefore w = U \Lambda U^T \text{ for some matrix } U.$$

$$\|x\|_w^2 = x^T w x = (Ux)^T U x = \langle Ux, Ux \rangle = \|Ux\|_2^2$$

$$\langle x, y \rangle_w = \langle w^T x, y \rangle$$

For gradient descent algorithm, we have:

$$x^{(k+1)} = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} f(x^{(k)}) + \langle \nabla w f(x^{(k)}), x - x^{(k)} \rangle_w + \frac{1}{2\alpha_k} \|x - x^{(k)}\|_w^2$$

$$x^{(k+1)} = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \langle w^T \nabla w f(x^{(k)}), x - x^{(k)} \rangle + \frac{1}{2\alpha_k} \|U(x - x^{(k)})\|_2^2$$

$$\text{let } F(x) = \langle w^T \nabla w f(x^{(k)}), x - x^{(k)} \rangle + \frac{1}{2\alpha_k} \|U(x - x^{(k)})\|_2^2$$

Solving  $\nabla F(x) = 0$ , we have:

$$w^T \nabla w f(x^{(k)}) + \frac{1}{2\alpha_k} (2U \cdot (Ux - Ux^{(k)})) = 0$$

$$\underbrace{w^T}_{= w^T(w^{-1})^{-1}} \underbrace{\nabla f}_{= \nabla f} = \underbrace{\nabla f}_{= \nabla f} \quad \underbrace{U^T}_{= W} \underbrace{U}_{= W}$$

$$\Leftrightarrow \nabla f(x^{(k)}) + \frac{1}{\alpha_k} w(x - x^{(k)}) = 0$$

$$w x = w x^{(k)} - \alpha_k \nabla f(x^{(k)})$$

$$\Leftrightarrow x = x^{(k)} - \alpha_k w^{-1} \nabla f(x^{(k)})$$

$\therefore$  we have  $x^{(k+1)} = x^{(k)} - \alpha_k w^{-1} \nabla f(x^{(k)})$  for gradient descent.