

Solutions to Assignment 5 of CPSC 368/516 (Spring'23)

February 28, 2023

1 Problem 1

Problem 1.1. Gradient descent for strongly convex functions. *In this problem we analyze a gradient descent algorithm for minimizing a twice-differentiable convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, which satisfies for every $x \in \mathbb{R}^n$, one has $mI \preceq \nabla^2 f(x) \preceq MI$ for some $0 < m \leq M$.*

The algorithm starts with some $x_0 \in \mathbb{R}^n$ and at every step $t = 0, 1, 2, \dots$ it chooses the next point

$$x_{t+1} := x_t - \alpha_t \nabla f(x_t),$$

where α_t is chosen to minimize the value $f(x_t - \alpha \nabla f(x_t))$ over all $\alpha \in \mathbb{R}$ while fixing x_t . Let $y^ := \min\{f(x) : x \in \mathbb{R}^n\}$.*

1. *Prove that*

$$\forall x, y \in \mathbb{R}^n, \quad \frac{m}{2} \|y - x\|^2 \leq f(y) - f(x) + \langle \nabla f(x), x - y \rangle \leq \frac{M}{2} \|y - x\|^2.$$

2. *Prove that*

$$\forall x \in \mathbb{R}^n, \quad f(x) - \frac{1}{2m} \|\nabla f(x)\|^2 \leq y^* \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|^2.$$

3. *Prove that for every $t = 0, 1, 2, \dots$*

$$f(x_{t+1}) \leq f(x_t) - \frac{1}{2M} \|\nabla f(x_t)\|^2.$$

4. *Prove that for every $t = 0, 1, 2, \dots$*

$$f(x_t) - y^* \leq \left(1 - \frac{m}{M}\right)^t (f(x_0) - y^*).$$

What is the number of iterations t required to reach $f(x_t) - y^ \leq \varepsilon$?*

5. *Consider a linear system $Ax = b$, where $b \in \mathbb{R}^n$ is a vector and $A \in \mathbb{R}^{n \times n}$ is a symmetric positive definite matrix such that $\frac{L_n(A)}{L_1(A)} \leq \kappa$ (where $L_1(A)$ and $L_n(A)$ are the smallest and the largest eigenvalues of A respectively). Use the above framework to design an algorithm for approximately solving the system $Ax = b$ with logarithmic dependency on the error $\varepsilon > 0$ and polynomial dependency on κ . What is the running time?*

1.1 Part 1

Lemma 2.6 from [1], implies that for any $x, y \in \mathbb{R}^n$

$$\begin{aligned}
f(y) &= f(x) + \int_0^1 \langle \nabla f(x + t(y-x)), y-x \rangle dt \\
&= f(x) + \int_0^1 \left\langle \nabla f(x) + \left(\int_0^t t(y-x)^\top \nabla^2 f(x + rt(y-x)) dr \right), y-x \right\rangle dt \\
&\quad \text{(Using Lemma 2.6 to expand } \nabla f(x + t(y-x))) \\
&= f(x) + \int_0^1 \langle \nabla f(x), y-x \rangle dt + \int_0^1 \int_0^1 t(y-x)^\top \nabla^2 f(x + rt(y-x))(y-x) dr dt \\
&= f(x) + \langle \nabla f(x), y-x \rangle + \int_0^1 \int_0^1 t(y-x)^\top \nabla^2 f(x + rt(y-x))(y-x) dr dt. \tag{1}
\end{aligned}$$

Moreover, since $mI \preceq \nabla^2 f(x) \preceq MI$, for any vector $z \in \mathbb{R}^n$,

$$m \|z\|_2^2 \leq z^\top \nabla^2 f(x) z \leq M \|z\|_2^2.$$

Combining these inequalities with Equation (1), implies that

$$\begin{aligned}
f(y) &\leq f(x) + \langle \nabla f(x), y-x \rangle + \int_0^1 \int_0^1 tM \|y-x\|_2^2 dr dt \\
&\leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{M}{2} \|y-x\|_2^2, \tag{2}
\end{aligned}$$

$$\begin{aligned}
f(y) &\geq f(x) + \langle \nabla f(x), y-x \rangle + \int_0^1 \int_0^1 tm \|y-x\|_2^2 dr dt \\
&\geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{m}{2} \|y-x\|_2^2. \tag{3}
\end{aligned}$$

The required results follow by rearranging the above inequalities.

1.2 Part 2

By the definition of y^* , for all $y \in \mathbb{R}^n$, $y^* \leq f(y)$. Since $y^* \leq f(y)$ and Equation (2) holds, for all $x, y \in \mathbb{R}^n$

$$y^* \leq f(y) \leq f(x) + \langle \nabla f(x), x-y \rangle + \frac{M}{2} \|y-x\|_2^2.$$

Further, as the above inequality holds for all $y \in \mathbb{R}^n$, it implies that

$$\begin{aligned}
y^* &\leq \inf_{y \in \mathbb{R}^n} f(x) + \langle \nabla f(x), x-y \rangle + \frac{M}{2} \|y-x\|_2^2 \\
&= \inf_{z \in \mathbb{R}^n: \|z\|=1} \inf_{t \in \mathbb{R}} f(x) + \frac{M}{2} t^2 + t \nabla f(x)^\top z \quad \text{(Substituting } y = x + tz \text{ for some } t \in \mathbb{R} \text{ and } z \in \mathbb{R}^n) \\
&= \inf_{z \in \mathbb{R}^n: \|z\|=1} f(x) - \frac{1}{2M} \langle \nabla f(x), z \rangle^2 \\
&= f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2.
\end{aligned}$$

Minimizing the RHS in Equation (3), implies that for all $x, y \in \mathbb{R}^n$

$$\begin{aligned}
f(y) &\geq \inf_{w \in \mathbb{R}^n} f(x) + \langle \nabla f(x), w-x \rangle + \frac{m}{2} \|w-x\|_2^2 \\
&= \inf_{z \in \mathbb{R}^n: \|z\|_2=1} \inf_{t \in \mathbb{R}} f(x) + t \langle \nabla f(x), z \rangle + \frac{m}{2} t^2 \quad \text{(Substituting } y = x + tz \text{ for some } t \in \mathbb{R} \text{ and } z \in \mathbb{R}^n) \\
&= \inf_{z \in \mathbb{R}^n: \|z\|_2=1} f(x) - \frac{2}{m} \langle \nabla f(x), z \rangle^2 \\
&= f(x) - \frac{2}{m} \|\nabla f(x)\|_2^2.
\end{aligned}$$

1.3 Part 3

Let $y = x_t - \alpha \nabla f(x_t)$ for some $\alpha \in \mathbb{R}$. Substituting $x = x_t$ in Equation (2), implies that:

$$\begin{aligned} f(y) &\leq f(x_t) + \langle \nabla f(x_t), y - x_t \rangle + \frac{M}{2} \|y - x_t\|_2^2 \\ &\leq f(x_t) + \langle \nabla f(x_t), -\alpha \nabla f(x_t) \rangle + \frac{M}{2} \|\alpha \nabla f(x_t)\|_2^2 \\ &\leq f(x_t) - \alpha \|\nabla f(x_t)\|_2^2 + \frac{M\alpha^2}{2} \|\nabla f(x_t)\|_2^2. \end{aligned}$$

Since x_{t+1} is defined such that $f(x_{t+1}) = \inf_{\alpha \in \mathbb{R}} f(x_t - \alpha \nabla f(x_t))$, for any $\alpha \in \mathbb{R}$

$$f(x_{t+1}) \leq f(y) \leq f(x_t) - \alpha \|\nabla f(x_t)\|_2^2 + \frac{M\alpha^2}{2} \|\nabla f(x_t)\|_2^2.$$

In particular, minimizing the RHS over α implies that

$$\begin{aligned} f(x_{t+1}) &\leq \inf_{\alpha \in \mathbb{R}} f(x_t) - \alpha \|\nabla f(x_t)\|_2^2 + \frac{M\alpha^2}{2} \|\nabla f(x_t)\|_2^2 \\ &\leq f(x_t) - \frac{1}{2M} \|\nabla f(x_t)\|_2^2. \end{aligned} \tag{4}$$

1.4 Part 4

Subtracting y^* from both sides of Equation (4), implies that for all $t = 0, 1, \dots$

$$\begin{aligned} f(x_{t+1}) - y^* &\leq f(x_t) - y^* - \frac{1}{2M} \|\nabla f(x_t)\|_2^2 \\ &\leq f(x_t) - y^* - \frac{m}{M} (y^* - f(x_t)) \\ &= \left(1 - \frac{m}{M}\right) \cdot (f(x_t) - y^*). \end{aligned}$$

Chaining the above inequality for all $t = 0, 1, \dots, T-1$, implies that

$$f(x_T) - y^* \leq \left(1 - \frac{m}{M}\right)^T \cdot (f(x_0) - y^*).$$

Recall that we require that $f(x_T) - y^* \leq \varepsilon$, or equivalently that

$$f(x_T) - y^* \leq \left(1 - \frac{m}{M}\right)^{T-1} (f(x_0) - y^*) \leq \varepsilon \iff T \geq \log \left(\frac{\varepsilon}{f(x_0) - y^*} \right) \cdot \left(\log \left(1 - \frac{m}{M} \right) \right)^{-1} + 1.$$

Thus, $\Theta \left(\frac{M}{m} \cdot \log \left(\frac{f(x_0) - y^*}{\varepsilon} \right) \right)$ iterations suffice.

1.5 Part 5

Notice that the solution to $Ax = b$ is also the optimal solution to the optimization problem

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \langle A, xx^\top \rangle - b^\top x.$$

This holds because $\frac{1}{2} \langle A, xx^\top \rangle - b^\top x$ is convex (as A is PD) and its gradient is zero iff $Ax = b$. Define

$$f(x) := x^\top \frac{A}{2} x - b^\top x.$$

One can observe that $\nabla^2 f(x) = A \succ 0$. Hence, it follows that $\nabla^2 f(x) = A \preceq \lambda_n(A)$ and that $\nabla^2 f(x) = A \succeq \lambda_1(A)$. We initialize the algorithm developed in the previous parts as follows:

- $m = \lambda_1(A)$ and $M = \lambda_n(A)$
- The gradient of f is $\nabla f(x) = -b + Ax$
- $f(x_t - \alpha \nabla f(x_t))$ is a convex function in α . Minimizing it over $\alpha \in \mathbb{R}$ and using the definition of α_t , we get that α_t (as a function of x_t) is

$$\alpha_t := \frac{(b - Ax_t)^\top (b - Ax_t)}{(b - Ax_t)^\top A (b - Ax_t)}$$

- We can initialize the starting point as $x_0 = 0$. This implies that $f(x_0) - y^* \leq \frac{1}{2} b^\top A^{-1} b$.

Thus,

$$\begin{aligned} T &= \log \left(\frac{2\varepsilon}{b^\top A^{-1} b} \right) \cdot \left(\log \left(1 - \frac{1}{\kappa} \right) \right)^{-1} + 1 \\ &= O \left(\kappa \log \left(\frac{2\varepsilon}{b^\top A^{-1} b} \right) \right). \end{aligned} \quad \text{(Using that } \log(1 - \frac{1}{\kappa}) = -O(\frac{1}{\kappa}) \text{)}$$

In every step, it takes $O(n^2)$ arithmetic operations in total to compute $m, M, \nabla f(x_t), \alpha_t$ and x_{t+1} , and hence, the amount of arithmetic steps are

$$O \left(n^2 \kappa \log \frac{2\varepsilon}{b^\top A^{-1} b} \right).$$

2 Problem 2

Problem 2.1. Let $G = (V, E)$ be an undirected graph with n vertices and m edges. Let $B \in \mathbb{R}^{n \times m}$ be the vertex-edge incidence matrix of G . Assume that G is connected and let $\Pi := B^\top (BB^\top)^+ B$. Prove that, given a vector $g \in \mathbb{R}^m$, if we let x_g denote the projection of g on the subspace $K := \{x \in \mathbb{R}^m : Bx = 0\}$, then it holds that

$$x_g = g - \Pi g.$$

Recall that the projection of a point $h \in \mathbb{R}^n$ on a closed convex nonempty set $K \subseteq \mathbb{R}^n$ is defined as the unique point $x_h \in K$ that minimizes the Euclidean distance to h :

$$x_h := \operatorname{argmin}_{y \in K} \|y - h\|_2. \quad (5)$$

First, observe that for any $h \in \mathbb{R}^m$ x_h is feasible for the instance Equation (5): To see this observe that

$$\begin{aligned} Bx_h &= Bh - B\Pi h \\ &= Bh - BB^\top (BB^\top)^+ Bh && \text{(Using that } \Pi = B^\top (BB^\top)^+ B \text{)} \\ &= Bh - Bh \\ &= 0. \end{aligned}$$

Thus, $x_h \in K$ and, hence, feasible for Equation (5).

Next, observe K is a linear space, and hence, the projection of h on K is the unique point $x_h \in K$ such that $h - x_h$ is orthogonal to all points $z \in K$. Further, observe that K is the null space of B , and hence, K is also the null space of $B^\top B$ (exercise: prove this). Furthermore, since K is the null space of $B^\top B$, it is the span of the eigenvectors corresponding to the 0 eigenvalue of $B^\top B$. Since $B^\top B$ is the Laplacian of a connected graph, the all-ones vector is the unique eigenvector corresponding to the eigenvalue 0 (here, connectedness is required to guarantee uniqueness). Since $z \in K$ must be parallel to the all-ones vector,

to check if $h - x_h$ is orthogonal to z , it suffices to show that $h - x_h$ is orthogonal to $1 \in \mathbb{R}^m$. This follows because

$$\begin{aligned}
\langle 1, h - x_h \rangle &= \langle 1, \Pi h \rangle && \text{(Using that } x_h = h - \Pi h \text{)} \\
&= 1^\top B^\top (BB^\top)^+ Bh && \text{(Using that } \Pi = B^\top (BB^\top)^+ B \text{)} \\
&= 0^\top (BB^\top)^+ Bh && \text{(Using that } B1 = 0 \text{)} \\
&= 0.
\end{aligned}$$

References

- [1] Nisheeth K. Vishnoi. *Algorithms for Convex Optimization*. Cambridge University Press, 2021.