

Convex Optimization

R. Kyng & M. Probst

Problem Set 2 — Wednesday, March 1st

These exercises will not count toward your grade, but you are encouraged to solve them all. This exercise sheet contains exercises relating to lectures in Week 2.

To get feedback, you must hand in your solutions by 23.59 pm on March 9th. Both hand-written and L^AT_EX solutions are acceptable, but we will only attempt to read legible text.

Exercise 1

Prove that if a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric, then $\|\mathbf{A}\| = \max(|\lambda_{\max}(\mathbf{A})|, |\lambda_{\min}(\mathbf{A})|)$ and give an example of a non-symmetric matrix for which this is not true.

Solution.

By definition, we have

$$\|\mathbf{A}\| = \max_{\mathbf{x} \in \mathbb{R}^n, \mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|.$$

Based on the lecture notes we have that $\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top$, where the columns of \mathbf{V} form an orthogonal basis and $\Lambda_{ii} = \lambda_i(\mathbf{A})$. Thus, we will have

$$\begin{aligned} \|\mathbf{A}\|^2 &= \max_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|^2 \\ &= \max_{\|\mathbf{x}\|=1} \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} \\ &= \max_{\|\mathbf{x}\|=1} \mathbf{x}^\top (\mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top \mathbf{V}\mathbf{\Lambda}\mathbf{V}^\top) \mathbf{x} \\ &= \max_{\|\mathbf{x}\|=1} (\mathbf{V}^\top \mathbf{x})^\top \mathbf{\Lambda}^2 (\mathbf{V}^\top \mathbf{x}). \end{aligned}$$

Furthermore, since \mathbf{V} is orthogonal

$$\|\mathbf{V}^\top \mathbf{x}\|^2 = (\mathbf{V}^\top \mathbf{x})^\top (\mathbf{V}^\top \mathbf{x}) = \mathbf{x}^\top \mathbf{V} \mathbf{V}^\top \mathbf{x} = \mathbf{x}^\top \mathbf{x} = \|\mathbf{x}\|^2$$

which implies that $\|\mathbf{x}\| = \|\mathbf{V}^\top \mathbf{x}\|$.

Overall, we have

$$\|\mathbf{A}\|^2 = \max_{\|\mathbf{z}\|=1} \mathbf{z}^\top \mathbf{\Lambda}^2 \mathbf{z} = \max_{\|\mathbf{z}\|=1} \sum_{i=1}^n \lambda_i^2 z(i)^2$$

This implies $\|\mathbf{A}\|^2 = \max(\lambda_{\max}^2, \lambda_{\min}^2)$. Thus, we have $\|\mathbf{A}\| = \max(\lambda_{\max}, \lambda_{\min})$.

On the other hand, consider the matrix

$$\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$$

which is not symmetric. We observe that the eigenvalues of \mathbf{A} are zero, but $\|\mathbf{A}\| = 1$. Therefore, we get $\|\mathbf{A}\| \neq \max(\lambda_{\max}, \lambda_{\min})$.

Exercise 2

Consider a twice continuously differentiable function $f : S \rightarrow \mathbb{R}$, where $S \subset \mathbb{R}^n$ is a convex open set. Prove that f is β -gradient Lipschitz if and only if for all $\mathbf{x} \in S$ we have $\|\mathbf{H}_f(\mathbf{x})\| \leq \beta$.

Solution.

Assume that $S \subset \mathbb{R}^n$ is open and convex. Let us first prove that if f is β -gradient Lipschitz, then $\|\lambda_{\max}(\mathbf{H}_f(\mathbf{x}))\| \leq \beta$. For an arbitrary $\mathbf{x} \in S$ and any $\boldsymbol{\delta} \neq \mathbf{0} \in \mathbb{R}^n$ such that $\mathbf{x} + \boldsymbol{\delta} \in S$ we have

$$\nabla f(\mathbf{x} + \boldsymbol{\delta}) = \nabla f(\mathbf{x}) + \mathbf{H}_f(\mathbf{x})\boldsymbol{\delta} + q(\boldsymbol{\delta})$$

where $\lim_{\boldsymbol{\delta} \rightarrow \mathbf{0}} \frac{\|q(\boldsymbol{\delta})\|}{\|\boldsymbol{\delta}\|} = 0$.

Now, we have

$$\|\nabla f(\mathbf{x} + \boldsymbol{\delta}) - \nabla f(\mathbf{x})\|_2 \leq \beta \|\boldsymbol{\delta}\|_2$$

By combing the aforementioned two inequalities, we get

$$\|\mathbf{H}_f(\mathbf{x})\boldsymbol{\delta} + q(\boldsymbol{\delta})\|_2 \leq \beta \|\boldsymbol{\delta}\|_2 \Rightarrow \frac{\|\mathbf{H}_f(\mathbf{x})\boldsymbol{\delta}\|_2}{\|\boldsymbol{\delta}\|_2} \leq \beta + \frac{\|q(\boldsymbol{\delta})\|_2}{\|\boldsymbol{\delta}\|_2}.$$

Since S is open, we can choose $\boldsymbol{\delta}$ such that

$$\frac{\|\mathbf{H}_f(\mathbf{x})\boldsymbol{\delta}\|_2}{\|\boldsymbol{\delta}\|_2} = \|\mathbf{H}_f(\mathbf{x})\|_2.$$

Furthermore, since $\mathbf{H}_f(\mathbf{x})$ is symmetric, we have

$$\|\mathbf{H}_f(\mathbf{x})\| = \max(|\lambda_{\max}(\mathbf{H}_f(\mathbf{x}))|, |\lambda_{\min}(\mathbf{H}_f(\mathbf{x}))|).$$

Then, as $\boldsymbol{\delta}$ tends to zero, we have bounded spectral norm, $\|\mathbf{H}_f(\mathbf{x})\| \leq \beta$.

Now, we prove the other direction. Assume that $\|\mathbf{H}_f(\mathbf{x})\| \leq \beta$ for all $\mathbf{x} \in S$, except a measure zero set. Consider $\mathbf{x} \neq \mathbf{y} \in S$ and let $\mathbf{x}_\theta = \mathbf{x} + \theta(\mathbf{y} - \mathbf{x})$. Since S is convex, $\mathbf{x}_\theta \in S$ for any $\theta \in [0, 1]$. By the fundamental theorem of calculus,

$$\nabla f(\mathbf{y}) = \nabla f(\mathbf{x}) + \int_0^1 \mathbf{H}_f(\mathbf{x}_\theta)(\mathbf{y} - \mathbf{x})d\theta.$$

Thus, we have

$$\begin{aligned}
\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|_2 &= \left\| \int_0^1 \mathbf{H}_f(\mathbf{x}_\theta)(\mathbf{y} - \mathbf{x}) d\theta \right\|_2 \\
&\leq \int_0^1 \|\mathbf{H}_f(\mathbf{x}_\theta)(\mathbf{y} - \mathbf{x})\|_2 d\theta \\
&\leq \int_0^1 \|\mathbf{H}_f(\mathbf{x}_\theta)\| \|\mathbf{y} - \mathbf{x}\|_2 d\theta \\
&\leq \int_0^1 \beta \|\mathbf{y} - \mathbf{x}\|_2 d\theta \\
&= \beta \|\mathbf{y} - \mathbf{x}\|_2.
\end{aligned}$$

Exercise 3

Prove that when running Gradient Descent, $\|\mathbf{x}_i - \mathbf{x}^*\|_2 \leq \|\mathbf{x}_0 - \mathbf{x}^*\|_2$ for all i .

Solution.

We assume f is β -Lipschitz continuous and convex, our prerequisites for gradient descent. Since f is convex, we have

$$f(\mathbf{x}_i) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_i)^\top (\mathbf{x}_i - \mathbf{x}^*).$$

By using $\mathbf{x}_{i+1} = \mathbf{x}_i - \frac{1}{\beta} \nabla f(\mathbf{x}_i)$, we will get

$$f(\mathbf{x}_i) - f(\mathbf{x}^*) \leq \beta(\mathbf{x}_i - \mathbf{x}_{i+1})^\top (\mathbf{x}_i - \mathbf{x}^*).$$

We know that $2\mathbf{v}^\top \mathbf{u} = \|\mathbf{v}\|_2^2 + \|\mathbf{u}\|_2^2 - \|\mathbf{v} - \mathbf{u}\|_2^2$ for two vectors \mathbf{v}, \mathbf{u} . Thus, we have

$$\begin{aligned}
f(\mathbf{x}_i) - f(\mathbf{x}^*) &\leq \frac{\beta}{2} (\|\mathbf{x}_i - \mathbf{x}_{i+1}\|_2^2 + \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2) \\
&= \frac{1}{2\beta} \|\nabla f(\mathbf{x}_i)\|_2^2 + \frac{\beta}{2} (\|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2).
\end{aligned}$$

Therefore, the above inequality yields

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

Since

$$f(\mathbf{x}_i) - f(\mathbf{x}^*) \geq f(\mathbf{x}_i) - f(\mathbf{x}_{i+1}) \geq \frac{\|\nabla f(\mathbf{x}_i)\|_2^2}{2\beta},$$

then

$$\|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2 \leq \|\mathbf{x}_i - \mathbf{x}^*\|_2^2.$$

Therefore, we can conclude that for any $i \geq 1$,

$$\|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2 \leq \|\mathbf{x}_0 - \mathbf{x}^*\|_2.$$

Exercise 4

Prove the following theorem.

Theorem. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be an β -gradient Lipschitz, convex function. Let \mathbf{x}_0 be a given starting point, and let $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ be a minimizer of f . The Gradient Descent algorithm given by $\mathbf{x}_{i+1} = \mathbf{x}_i - \frac{1}{\beta} \nabla f(\mathbf{x}_i)$ ensures that the k th iterate satisfies

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{2\beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{k+1}.$$

Hint: do an induction on $1/\text{gap}_i$.

Solution.

Let $C = \beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$. We want to prove that for any integer $i \geq 0$

$$\text{gap}_i \leq \frac{2C}{i+1}. \tag{1}$$

In the following, we assume $\text{gap}_i > 0$ for all i : if this is ever violated, the algorithm has reached the optimum and will stay there as the gradient is then zero – and in that case Equation (1) holds.

We will prove Equation (1) by proving by induction that $\frac{1}{\text{gap}_i} \geq \frac{i+1}{2C}$. For the base case, using that f is β -gradient Lipschitz, we have

$$\text{gap}_0 = f(\mathbf{x}_0) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}^*)^\top (\mathbf{x}_0 - \mathbf{x}^*) + \frac{\beta}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 = \frac{\beta}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 = C/2 \leq 2C.$$

This rearranges to $\frac{1}{\text{gap}_0} \geq \frac{0+1}{2C}$. From the lecture we know that

$$\text{gap}_{i+1} - \text{gap}_i \leq -\frac{\text{gap}_i^2}{2C}.$$

Since we can assume $\text{gap}_{i+1} > 0$, and we have $\text{gap}_i \geq \text{gap}_{i+1}$, dividing through by $\text{gap}_i \cdot \text{gap}_{i+1}$ gives

$$\frac{1}{\text{gap}_i} - \frac{1}{\text{gap}_{i+1}} \leq -\frac{\text{gap}_i^2}{2C \text{gap}_i \cdot \text{gap}_{i+1}} \leq -\frac{1}{2C}.$$

Finally, we have that $\frac{1}{\text{gap}_{i+1}} \geq \frac{1}{2C} + \frac{1}{\text{gap}_i} \geq \frac{(i+1)+1}{2C}$ using the induction hypothesis.

Exercise 5

1. For each of the following functions answer these questions:

- Is the function convex?
- Is the function β -gradient Lipschitz for some β ?
- If the function is β -gradient Lipschitz give an upper bound on β – the bound should be within a factor 4 of the true value.

- (a) $f(x) = |x|^{1.5}$ on $x \in \mathbb{R}$
- (b) $f(x) = \exp(x)$ on $x \in \mathbb{R}$
- (c) $f(x) = \exp(x)$ on $x \in (-1, 1)$
- (d) $f(x, y) = \sqrt{x+y}$ on $(x, y) \in (0, 1) \times (0, 1)$.
- (e) $f(x, y) = \sqrt{x+y}$ on $(x, y) \in (1/2, 1) \times (1/2, 1)$.
- (f) $f(x, y) = \sqrt{x^2 + y^2}$ on $(x, y) \in \mathbb{R}^2$.

Solution.

- (a) Function $f(x) = |x|^{1.5}$ on $x \in \mathbb{R}$ is convex but it is not β -gradient Lipschitz. We have

$$\nabla f(x) = \frac{3}{2}|x|^{-\frac{1}{2}}x \quad \text{and} \quad \mathbf{H}_f(x) = \frac{3}{4}|x|^{-\frac{1}{2}}.$$

Since $\mathbf{H}_f(x)$ is positive semidefinite, function f is convex. However, it is not β -Lipschitz gradient for any β because $\|\mathbf{H}_f(x)\|_2$ tends to infinity when we let x go to zero, which implies that $\|\mathbf{H}_f(x)\|_2$ cannot be upper-bounded.

- (b) Function $f(x) = \exp(x)$ on $x \in \mathbb{R}$ is convex but it is not β -gradient Lipschitz. We have $\mathbf{H}_f(x) = \exp(x)$. Function f is convex because $\exp(x) \geq 0$ for any $x \in \mathbb{R}$, but it is not β -gradient Lipschitz because $\exp(x)$ is not bounded on \mathbb{R} .

- (c) Function $f(x) = \exp(x)$ on $x \in (-1, 1)$ is convex and β -gradient Lipschitz for $\beta = \exp(1)$. We know that $\mathbf{H}_f(x) = \exp(x)$. Since $\exp(x) \geq 0$ for any $x \in (-1, 1)$, then it is convex. Furthermore, since $\exp(-1) < \exp(x) < \exp(1)$ for $x \in (-1, 1)$, then f is β -gradient Lipschitz for $\beta = \exp(1)$. If we set $x = 0$, then the eigenvalue of $\mathbf{H}_f(x)$ is equal to 1 which implies that $\beta \geq 1$ by using $\|\lambda_{\max}(\mathbf{H}_f(\mathbf{x}))\|_2 \leq \beta$ from Exercise 1. Therefore, the upper bound of $\exp(1)$ is within a factor $\exp(1) < 4$ from the true value (actually, $1 + \epsilon$ for any $\epsilon > 0$).

- (d) Function $f(x, y) = \sqrt{x+y}$ on $(x, y) \in (0, 1) \times (0, 1)$ is neither convex nor β -gradient Lipschitz for any β . We have

$$\mathbf{H}_f(x, y) = \begin{bmatrix} -\frac{1}{4(x+y)^{3/2}} & -\frac{1}{4(x+y)^{3/2}} \\ -\frac{1}{4(x+y)^{3/2}} & -\frac{1}{4(x+y)^{3/2}} \end{bmatrix}.$$

For a non-zero vector $\mathbf{z} = [z_1, z_2]^\top$, we have $\mathbf{z}^\top \mathbf{H}_f(x, y) \mathbf{z} = -(z_1 + z_2)^2 / 4(x+y)^{3/2}$ which is negative for any $(x, y) \in (0, 1) \times (0, 1)$. Thus, f is not convex. Furthermore, it is not β -gradient Lipschitz because $\|\mathbf{H}_f(x, y)\|_2 = \frac{1}{2(x+y)^{3/2}}$ which is not bounded on $(0, 1) \times (0, 1)$.

- (e) Function $f(x, y) = \sqrt{x+y}$ on $(x, y) \in (1/2, 1) \times (1/2, 1)$ is not convex but it is $\frac{1}{2}$ -gradient Lipschitz. We know that

$$\mathbf{H}_f(x, y) = \begin{bmatrix} -\frac{1}{4(x+y)^{3/2}} & -\frac{1}{4(x+y)^{3/2}} \\ -\frac{1}{4(x+y)^{3/2}} & -\frac{1}{4(x+y)^{3/2}} \end{bmatrix}.$$

With the same argument as above, f is not convex. However, it is β -gradient Lipschitz for $\beta = 1/2$ because $\frac{1}{2^{5/2}} \leq \|\mathbf{H}_f(x, y)\|_2 = \frac{1}{2(x+y)^{3/2}} \leq \frac{1}{2}$ on $(x, y) \in (1/2, 1) \times (1/2, 1)$.

(f) Function $f(x, y) = \sqrt{x^2 + y^2}$ on $(x, y) \in \mathbb{R}^2$ is convex but it is not β -gradient Lipschitz for any β . We have

$$\mathbf{H}_f(x, y) = (x^2 + y^2)^{-\frac{3}{2}} \begin{bmatrix} y^2 & -xy \\ -xy & x^2 \end{bmatrix}$$

for $x, y \neq 0$ and

$$\lim_{(x,y) \rightarrow (0,0)} \mathbf{H}_f(x, y) = [\infty, \infty]^\top.$$

Furthermore, for any $\mathbf{z} = [z_1, z_2]^\top$, we have

$$\mathbf{z}^\top \mathbf{H}_f(x, y) \mathbf{z} = \frac{(z_1 y - z_2 x)^2}{(x^2 + y^2)^{\frac{3}{2}}} \geq 0$$

for all $(x, y) \in \mathbb{R}^2$. Therefore, $\mathbf{H}_f(x, y)$ is positive semidefinite for any $(x, y) \in \mathbb{R}^2$, which implies that f is convex.

Bonus Exercise 6: Strongly Convex Functions

This longer exercise will teach you about *strongly convex* functions. In it, you will show that, with a small tweak, gradient descent quickly converges to a highly accurate solutions on these functions.

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. Assume f is twice continuously (Fréchet) differentiable and that its first and second (Fréchet) derivatives are integrable (basically, don't worry that weird stuff is happening with the derivatives). Assume that for all \mathbf{x} , we have for some constant $\mu > 0$, that $\lambda_{\min}(H_f(\mathbf{x})) \geq \mu$. When this holds, we say that f is μ -strongly convex.

Part A. Prove that for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

Part B. Prove that there is value $L \in \mathbb{R}$ such that for all $\mathbf{x} \in \mathbb{R}^n$, we have $f(\mathbf{x}) \geq L$. In other words, the function is not unbounded below.

Part C. Prove that f is *strictly convex* as per Definition 3.2.8 in Chapter 3. Prove also that the minimizer $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ of f is unique.

Part D. Let \mathbf{x}_0 be a given starting point and \mathbf{x}^* be the minimizer of f . Suppose we have an algorithm DECENTDESCENT which takes a starting point \mathbf{x}_0 , and a step count $t \in \mathbb{N}$. DECENTDESCENT(\mathbf{x}_0, t) runs for t steps and returns $\tilde{\mathbf{x}} \in \mathbb{R}^n$ such that

$$f(\tilde{\mathbf{x}}) - f(\mathbf{x}^*) \leq \frac{\gamma \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{t + 1}$$

where $\gamma > 0$ is a positive number.

Assume that the *cost* of running DECENTDESCENT for t steps is t . Explain how, with a total cost of at most $\frac{8\gamma}{\mu} \log(\|\mathbf{x}_0 - \mathbf{x}^*\|_2 / \delta)$, we can produce a point $\hat{\mathbf{x}} \in \mathbb{R}^n$ such that $\|\hat{\mathbf{x}} - \mathbf{x}^*\|_2 \leq \delta$ for $\delta > 0$.

Part E. Consider a function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ which is both μ -strongly convex and β -gradient Lipschitz. Give an algorithm that returns \mathbf{x}' with

$$h(\mathbf{x}') - h(\mathbf{x}^*) \leq \epsilon$$

by computing the gradient of h at at most $\frac{32\beta}{\mu} \log(2\beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 / \epsilon)$ points.

Solution.

Part A. Let \mathbf{z} be an arbitrary vector in \mathbb{R}^n . Since $\mathbf{H}_f(\mathbf{z})$ is symmetric, by applying the Courant-Fischer theorem, we have

$$\lambda_{\min}(\mathbf{H}_f(\mathbf{z})) = \min_{\mathbf{u} \in \mathbb{R}^n, \mathbf{u} \neq \mathbf{0}} \frac{\mathbf{u}^\top \mathbf{H}_f(\mathbf{z}) \mathbf{u}}{\mathbf{u}^\top \mathbf{u}}.$$

We know that $\lambda_{\min}(\mathbf{H}_f(\mathbf{z})) \geq \mu$. Therefore, for any $\mathbf{u} \in \mathbb{R}^n$

$$\mathbf{u}^\top \mathbf{H}_f(\mathbf{z}) \mathbf{u} \geq \mu \|\mathbf{u}\|_2^2. \quad (2)$$

(Note Equation (2) trivially holds for $\mathbf{u} = \mathbf{0}$.)

Furthermore, based on Taylor's theorem for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, there exists $\mathbf{z} \in \mathbb{R}^n$ such that

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{x})^\top \mathbf{H}_f(\mathbf{z}) (\mathbf{y} - \mathbf{x}). \quad (3)$$

Combining Equations (2) and (3), we conclude that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2} \mu \|\mathbf{y} - \mathbf{x}\|_2^2.$$

Part B. Consider a fixed vector $\mathbf{x} \in \mathbb{R}^n$ with bounded $f(\mathbf{x})$ and $\nabla f(\mathbf{x})$. From Part A, we know that for any $\mathbf{y} \in \mathbb{R}^n$

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{1}{2} \mu \|\mathbf{y} - \mathbf{x}\|_2^2 \\ &= f(\mathbf{x}) + \frac{1}{2} \mu \left(\mathbf{y} - \mathbf{x} + \frac{1}{\mu} \nabla f(\mathbf{x}) \right)^\top \left(\mathbf{y} - \mathbf{x} + \frac{1}{\mu} \nabla f(\mathbf{x}) \right) - \frac{1}{2} \mu \left\| \frac{1}{\mu} \nabla f(\mathbf{x}) \right\|_2^2 \\ &\geq f(\mathbf{x}) - \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2. \end{aligned}$$

Since \mathbf{x} is fixed, $L = f(\mathbf{x}) - \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2$ is a constant lower bound on $f(\mathbf{y})$.

Part C. Consider two arbitrary vectors $\mathbf{x} \neq \mathbf{y} \in \mathbb{R}^n$ and $\theta \in (0, 1)$. Let $\mathbf{z} = \theta \mathbf{x} + (1 - \theta) \mathbf{y}$, which implies that $\mathbf{x} - \mathbf{z} = (1 - \theta)(\mathbf{x} - \mathbf{y})$ and $\mathbf{y} - \mathbf{z} = \theta(\mathbf{y} - \mathbf{x})$. Now, by applying Part A we have

$$f(\mathbf{x}) \geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top ((1 - \theta)(\mathbf{x} - \mathbf{y})) + \frac{\mu}{2} (1 - \theta) \|\mathbf{x} - \mathbf{y}\|_2^2 \quad (4)$$

and

$$f(\mathbf{y}) \geq f(\mathbf{z}) + \nabla f(\mathbf{z})^\top (\theta(\mathbf{y} - \mathbf{x})) + \frac{\mu}{2} \theta \|\mathbf{x} - \mathbf{y}\|_2^2 \quad (5)$$

By multiplying Equation (4) by θ and Equation (5) by $1 - \theta$ and summing them up we get

$$\theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y}) \geq f(\mathbf{z}) + \mu\theta(1 - \theta)\|\mathbf{x} - \mathbf{y}\|_2^2.$$

Since $\mu > 0$, $\theta \in (0, 1)$, and $\mathbf{y} \neq \mathbf{x}$, the term $\mu\theta(1 - \theta)\|\mathbf{y} - \mathbf{x}\|_2^2$ is strictly positive. Therefore, we have that

$$\theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y}) > f(\theta\mathbf{x} + (1 - \theta)\mathbf{y})$$

which implies that f is strictly convex.

Next, we show that if $\arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ is non-empty, then it has exactly one element. Assume that f attains a minimum f^* at both \mathbf{x}_1 and \mathbf{x}_2 , where $\mathbf{x}_1 \neq \mathbf{x}_2$. Since f is strictly convex, we have that

$$f\left(\frac{1}{2}\mathbf{x}_1 + \frac{1}{2}\mathbf{x}_2\right) < \frac{1}{2}f(\mathbf{x}_1) + \frac{1}{2}f(\mathbf{x}_2) = \frac{1}{2}f^* + \frac{1}{2}f^* = f^*.$$

This is a contradiction.

Part D. Set $t' = 8\gamma/\mu$. We observe that $\text{DECENTDESCENT}(\mathbf{x}_0, t')$ returns $\mathbf{x}_{t'}$ such that

$$f(\mathbf{x}_{t'}) - f(\mathbf{x}^*) \leq \frac{\gamma\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{t' + 1} \leq \frac{\mu\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{8}.$$

Furthermore, from Part A we know that

$$\frac{\mu}{2}\|\mathbf{x}_{t'} - \mathbf{x}^*\|_2^2 \leq f(\mathbf{x}_{t'}) - f(\mathbf{x}^*)$$

where we used $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

Combining the above two equations implies that

$$\frac{\mu}{2}\|\mathbf{x}_{t'} - \mathbf{x}^*\|_2^2 \leq \frac{\mu\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{8} \Rightarrow \|\mathbf{x}_{t'} - \mathbf{x}^*\|_2 \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2}{2}.$$

Now, by setting \mathbf{x}_0 to be $\mathbf{x}_{t'}$ and applying the above argument iteratively we can conclude that for any $k \geq 1$

$$\|\mathbf{x}_{kt'} - \mathbf{x}^*\|_2 \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2}{2^k}.$$

Thus, for $k' = \log(\frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2}{\delta})$, we have

$$\|\mathbf{x}_{k't'} - \mathbf{x}^*\|_2 \leq \delta.$$

Therefore, with a total cost of at most $t'k' = \frac{8\gamma}{\mu} \log(\|\mathbf{x}_0 - \mathbf{x}^*\|_2/\delta)$, we can produce a point $\hat{\mathbf{x}} = \mathbf{x}_{t'k'} \in \mathbb{R}^n$ such that $\|\hat{\mathbf{x}} - \mathbf{x}^*\|_2 \leq \delta$.

Part E. We know that after t steps of Gradient Descent, each of which requires one computation of the gradient of function h , we obtain a point \mathbf{x}_t such that

$$h(\mathbf{x}_t) - h(\mathbf{x}^*) \leq \frac{2\beta\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{t + 1}.$$

Notice that this is the same as DECENTDESCENT from Part D for $\gamma = 2\beta$. Therefore, if we set $\delta = \sqrt{4\epsilon/\gamma}$, after

$$\begin{aligned}
t' &= \frac{8\gamma}{\mu} \log\left(\frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2}{\delta}\right) \\
&= \frac{16\beta}{\mu} \log\left(\sqrt{\frac{\beta}{2\epsilon}} \|\mathbf{x}_0 - \mathbf{x}^*\|_2\right) \\
&= \frac{8\beta}{\mu} \log\left(\frac{\beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2\epsilon}\right) \\
&\leq \frac{32\beta}{\mu} \log\left(\frac{2\beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{\epsilon}\right)
\end{aligned}$$

steps, we obtain \mathbf{x}' such that

$$\|\mathbf{x}' - \mathbf{x}^*\|_2 \leq \sqrt{\frac{2\epsilon}{\beta}}. \quad (6)$$

Furthermore, since h is β -gradient Lipschitz, we have

$$h(\mathbf{x}') - h(\mathbf{x}^*) \leq \nabla h(\mathbf{x}^*)^\top (\mathbf{x}' - \mathbf{x}^*) + \frac{\beta}{2} \|\mathbf{x}' - \mathbf{x}^*\|_2^2 = \frac{\beta}{2} \|\mathbf{x}' - \mathbf{x}^*\|_2^2 \quad (7)$$

where we used that $\nabla h(\mathbf{x}^*) = \mathbf{0}$.

Combining Equation (6) and Equation (7) yields

$$h(\mathbf{x}') - h(\mathbf{x}^*) \leq \frac{\beta}{2} \left(\sqrt{\frac{2\epsilon}{\beta}}\right)^2 = \epsilon.$$

Thus, by computing the gradient of h for at most $\frac{32\beta}{\mu} \log(2\beta \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 / \epsilon)$ times, we can find \mathbf{x}' such that $h(\mathbf{x}') - h(\mathbf{x}^*) \leq \epsilon$.