

Lecture Note 1

1 When have we arrived at an optimal solution?

Consider the optimization problem

$$\min_{\mathbf{x} \in \Omega} f(\mathbf{x}). \quad (1)$$

Definition 1 (optimal solution)

- We say \mathbf{x}^* is a solution of (1) if $\mathbf{x}^* \in \Omega$ and for all $\mathbf{x} \in \Omega$ we have $f(\mathbf{x}) \geq f(\mathbf{x}^*)$.
- We say \mathbf{x}^* is a local solution of (1) if $\mathbf{x}^* \in \Omega$ and there is a neighborhood \mathcal{N} around \mathbf{x}^* such that for all $\mathbf{x} \in \Omega \cap \mathcal{N}$ we have $f(\mathbf{x}) \geq f(\mathbf{x}^*)$.

Interestingly there exists infinitely differentiable function f (in fact of the form $f(\mathbf{x}) = \sum_{i,j} Q_{ij}x_i^2x_j^2$) such that even checking local optimality is NP hard. This begs the question for which functions can we check optimality? One class of functions where we can check optimality is convex functions.

2 Convexity

Convex functions are functions for which the line between any two points lies above the graph. Formally

$$f(\theta\mathbf{x} + (1-\theta)\mathbf{y}) \leq \theta f(\mathbf{x}) + (1-\theta)f(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y} \in \Omega, \quad \theta \in [0, 1].$$

One of the reasons we like convex functions in optimization is the following theorem.

Theorem 1 For convex functions, local minima/maxima are global minima/maxima.

Proof Assume that local minima are not global minima. If $\mathbf{x}_1, \mathbf{x}_2$ are two local minima such that $f(\mathbf{x}_1) < f(\mathbf{x}_2)$, and f is convex, then

$$f(\theta\mathbf{x}_1 + (1-\theta)\mathbf{x}_2) \leq \theta f(\mathbf{x}_1) + (1-\theta)f(\mathbf{x}_2) < f(\mathbf{x}_2) \quad \forall \theta \in [0, 1],$$

which is in contradiction with local optimality of \mathbf{x}_1 . ■

3 Gradients, Hessians, and Taylor's theorem

For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, we define

$$\text{Gradient} \quad \nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_i} \right]_{i=1,2,\dots,n}$$

and

$$\text{Hessian} \quad \nabla^2 f(\mathbf{x}) = \left[\frac{\partial^2 f}{\partial x_i \partial x_j} \right]_{i,j=1,2,\dots,n}$$

The most fundamental theorem in all of optimization is one you learned in calculus.

Theorem 2 (Taylor's Theorem)

- If f is continuously differentiable, then

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(t\mathbf{x} + (1-t)\mathbf{y})^T(\mathbf{y} - \mathbf{x}) \quad \text{for some } t \in [0, 1].$$

- If f is twice continuously differentiable, then

$$\nabla f(\mathbf{y}) = \nabla f(\mathbf{x}) + \int_0^1 \nabla^2 f(t\mathbf{y} + (1-t)\mathbf{x})(\mathbf{y} - \mathbf{x}) dt$$

and

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \nabla^2 f(t\mathbf{y} + (1-t)\mathbf{x})(\mathbf{y} - \mathbf{x}) \quad \text{for some } t \in [0, 1].$$

4 Equivalent conditions for convexity

We now state two equivalent conditions for convexity.

Theorem 3 (first order condition for convexity) If f is continuously differentiable, then f is convex if and only if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ we have

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}).$$

Stated differently, the first order Taylor approximation is a global lower bound for f .

Theorem 4 (second order condition for convexity) If f is twice continuously differentiable on its domain ($\text{dom}(f)$), then f is convex if and only if for all $\mathbf{x} \in \text{dom}(f)$ we have $\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$.

4.1 Convexity through derivatives

A useful way to characterize convex function is through monotonicity of their gradients. This is the subject of the next definition.

Definition 2 A mapping $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is called monotone if for all $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$ we have

$$\langle g(\mathbf{x}) - g(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq 0.$$

It turns out that for differentiable functions convexity and monotonicity of gradients are equivalent.

Theorem 5 (monotonicity of gradients) A differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if and only if the gradient mapping ∇f is monotone.

Proof If f is convex by Theorem 3 for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ we have

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$$

and

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$$

Combining these two identities we conclude that

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq 0,$$

which implies monotonicity of the gradients.

Now to prove the other side assume that the mapping ∇f is monotone. Define $g(t) = f(\mathbf{x} + t(\mathbf{y} - \mathbf{x}))$. Thus, $g'(t) = \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle$. Note that monotonicity of ∇f implies that

$$t(g'(t) - g'(0)) = \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), t(\mathbf{y} - \mathbf{x}) \rangle \geq 0.$$

The latter implies that for all $0 \leq t \leq 1$

$$g'(t) \geq g'(0)$$

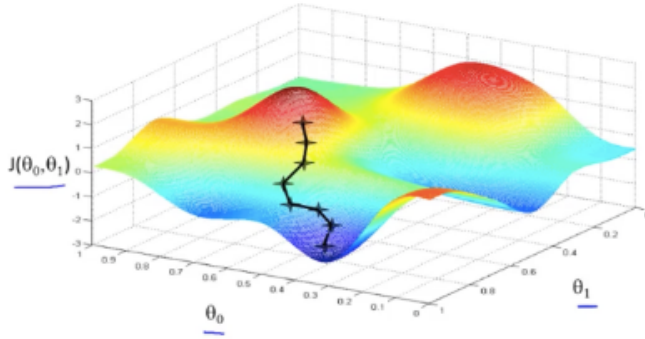
Thus,

$$f(\mathbf{y}) = g(1) = g(0) + \int_0^1 g'(t) dt \geq g(0) + g'(0) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$$

Thus we have shown that for all \mathbf{x}, \mathbf{y} monotonicity implies

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle.$$

The latter inequality together with Theorem 3 immediately implies that the function f is convex. ■



5 Gradient Descent

Gradient descent is perhaps the most widely used optimization algorithm. The strategy is simple: Go down hill, a.k.a. the descent direction.

Definition 3 (Descent direction) \mathbf{v} is a descent direction for f at \mathbf{x} if

$$f(\mathbf{x} + \mu \mathbf{v}) < f(\mathbf{x})$$

for some $\mu > 0$.

Lemma 6 For a continuously differentiable function f in a neighborhood of \mathbf{x} , if $\langle \mathbf{v}, \nabla f(\mathbf{x}) \rangle < 0$, then \mathbf{v} is a direction of descent.

Proof First note that by continuity of gradient there exists a μ_0 such that

$$\langle \nabla f(\mathbf{x} + \mu \mathbf{v}), \mathbf{v} \rangle < 0, \quad \text{for all } \mu \leq \mu_0.$$

By Taylor's theorem, we conclude that there exists $\tilde{\mu} \in [0, \mu]$ such that

$$f(\mathbf{x} + \mu \mathbf{v}) = f(\mathbf{x}) + \mu \langle \nabla f(\mathbf{x} + \tilde{\mu} \mathbf{v}), \mathbf{v} \rangle.$$

Thus using $\mu \langle \nabla f(\mathbf{x} + \tilde{\mu} \mathbf{v}), \mathbf{v} \rangle < 0$

$$f(\mathbf{x} + \mu \mathbf{v}) = f(\mathbf{x}) + \mu \langle \nabla f(\mathbf{x} + \tilde{\mu} \mathbf{v}), \mathbf{v} \rangle < f(\mathbf{x}).$$

Hence, \mathbf{v} is a descent direction. ■

5.1 Heuristic development of gradient descent

Let's focus on the approximation

$$f(\mathbf{x} + \mu \mathbf{v}) \approx f(\mathbf{x}) + \mu \mathbf{v}^\top \nabla f(\mathbf{x}).$$

We want to make $\mu \mathbf{v}^\top \nabla f(\mathbf{x})$ as negative as possible, so we pick \mathbf{v} such that

$$\mathbf{v} = - \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|_{\ell_2}}$$

Gradient Descent $\mathbf{x}_{t+1} = \mathbf{x}_t - \mu_t \nabla f(\mathbf{x}_t)$

- Choose $\mathbf{x}_0 \in \mathbb{R}^n$
- $\mathbf{x}_{t+1} = \mathbf{x}_t - \mu_t \nabla f(\mathbf{x}_t)$
- Repeat until convergence

5.2 Optimality conditions for gradient descent

Theorem 7 (Optimality Conditions) *The following are optimality conditions for the gradient descent algorithm:*

1. Assume f is continuously differentiable and \mathbf{x}^* is a local minimum $\Rightarrow \nabla f(\mathbf{x}^*) = \mathbf{0}$.
2. Assume $\nabla^2 f$ is continuous and \mathbf{x}^* is a local minimum $\Rightarrow \nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0}$.
3. If f is twice continuously differentiable, $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0} \Rightarrow \mathbf{x}^*$ is a local minimum.

Proof

1. $-\nabla f(\mathbf{x}^*)$ is a descent direction $\Rightarrow \nabla f(\mathbf{x}^*) = \mathbf{0}$
2. \mathbf{x}^* is a local minimum $\Rightarrow f(\mathbf{x}^* + t\mathbf{u}) \geq f(\mathbf{x}^*)$ for all \mathbf{u} and small enough t . By Taylor's theorem, $\exists \tilde{t} \in [0, t]$ such that

$$f(\mathbf{x}^* + t\mathbf{u}) = f(\mathbf{x}^*) + t\langle \nabla f(\mathbf{x}^*), \mathbf{u} \rangle + \frac{1}{2}t^2 \mathbf{u}^\top \nabla^2 f(\mathbf{x}^* + \tilde{t}\mathbf{u}) \mathbf{u}$$

By 1., we have

$$\begin{aligned} f(\mathbf{x}^* + t\mathbf{u}) - f(\mathbf{x}^*) &= \frac{1}{2}t^2 \mathbf{u}^\top \nabla^2 f(\mathbf{x}^* + \tilde{t}\mathbf{u}) \mathbf{u} \\ \Rightarrow \forall \mathbf{u}, \quad \mathbf{u}^\top \nabla^2 f(\mathbf{x}^* + \tilde{t}\mathbf{u}) \mathbf{u} &\geq 0 \\ \Rightarrow \forall \mathbf{u}, \quad \nabla^2 f(\mathbf{x}^* + \tilde{t}\mathbf{u}) &\succeq \mathbf{0} \\ \Rightarrow \nabla^2 f(\mathbf{x}^*) &\succeq \mathbf{0} \end{aligned}$$

3. $\exists r$ such that $\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ for all $\mathbf{x} \in B(\mathbf{x}^*, r)$ (ball with radius r). Pick \mathbf{u} with $\|\mathbf{u}\|_{\ell_2} < r$, then we have $t \in [0, 1]$ such that

$$f(\mathbf{x}^* + \mathbf{u}) = f(\mathbf{x}^*) + \mathbf{u}^\top \nabla f(\mathbf{x}^*) + \frac{1}{2} \mathbf{u}^\top \nabla^2 f(\mathbf{x}^* + t\mathbf{u}) \mathbf{u} \geq f(\mathbf{x}^*)$$

■

For the following, we will assume f to be convex.

Lemma 8 *For $f : \mathbb{R}^n \rightarrow \mathbb{R}$ continuously differentiable, \mathbf{x} is a global minimizer $\iff \nabla f(\mathbf{x}^*) = \mathbf{0}$.*

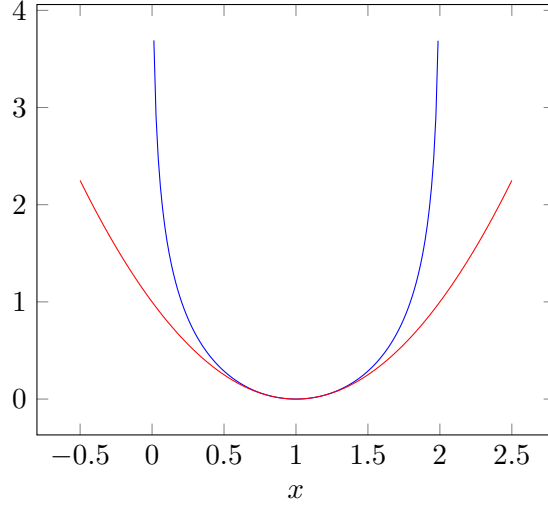


Figure 1: $-(\ln(x) + \ln(2 - x))$ descends fast (blue), while $(x - 1)^2$ descends slowly in comparison (red).

Proof We will assume $\nabla f(\mathbf{x}^*) = 0$. Then,

$$f \text{ continuously differentiable} \iff f(\mathbf{x}) \geq f(\mathbf{x}^*) + \langle \nabla f(\mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle \quad \forall \mathbf{x},$$

which implies that, if $\nabla f(\mathbf{x}^*) = 0$, then

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) \quad \forall \mathbf{x}.$$

If \mathbf{x} is a global minimum, then $\nabla f(\mathbf{x}) = 0$. ■

6 Convergence of gradient descent

Gradient descent doesn't work well when functions change too rapidly. A way to quantify the rate of change of a function is its "Lipschitzness". Figure 1 shows a examples of functions that change and do not change too rapidly.

6.1 Lipschitzness

Definition 4 (L-Lipschitz/Bounded Gradient) *The function f is L -Lipschitz if $\forall x, y \in \text{dom}(f)$, $|f(x) - f(y)| \leq L\|x - y\|_{\ell_2}$. Furthermore, if f is also differentiable, then*

$$\|\nabla f(\mathbf{x})\| \leq L.$$

Lemma 9 *We have*

- (1) *If $|f(\mathbf{x}) - f(\mathbf{y})| \leq L\|\mathbf{x} - \mathbf{y}\|_{\ell_2}$ and f is differentiable then $\|\nabla f(\mathbf{x})\|_{\ell_2} \leq L$.*
- (2) *If $\|\nabla f(\mathbf{x})\|_{\ell_2} \leq L$ then $|f(\mathbf{x}) - f(\mathbf{y})| \leq L\|\mathbf{x} - \mathbf{y}\|_{\ell_2}$.*

Proof The first part follows from the mean value theorem. For the second part, by the Fundamental Theorem of Calculus, we have

$$f(\mathbf{x}) - f(\mathbf{y}) = \int_0^1 \langle \nabla f(t\mathbf{x} + (1-t)\mathbf{y}), (\mathbf{x} - \mathbf{y}) \rangle dt$$

Therefore,

$$\begin{aligned} \|f(\mathbf{x}) - f(\mathbf{y})\| &= \left| \int_0^1 \langle \nabla f(t\mathbf{x} + (1-t)\mathbf{y}), (\mathbf{x} - \mathbf{y}) \rangle dt \right| \\ &\leq \int_0^1 |\langle \nabla f(t\mathbf{x} + (1-t)\mathbf{y}), (\mathbf{x} - \mathbf{y}) \rangle| dt \\ &\leq \int_0^1 \|\nabla f(t\mathbf{x} + (1-t)\mathbf{y})\| \cdot \|\mathbf{x} - \mathbf{y}\| dt \\ &= \left(\int_0^1 \|\nabla f(t\mathbf{x} + (1-t)\mathbf{y})\| dt \right) \cdot \|\mathbf{x} - \mathbf{y}\| \\ &\leq L\|\mathbf{x} - \mathbf{y}\| \end{aligned}$$

■

Theorem 10 *Let f be a convex and L -Lipschitz function. Furthermore, assume the function is L -Lipschitz and the initial estimate obeys $\|\mathbf{x}_1 - \mathbf{x}^*\| \leq R$. Then, for a step-size obeying $\mu = \frac{R}{L\sqrt{t}}$ we have*

$$f\left(\frac{1}{t} \sum_{s=1}^t \mathbf{x}_s\right) - f(\mathbf{x}^*) \leq \frac{RL}{\sqrt{t}}$$

Proof

$$\begin{aligned} f(\mathbf{x}_s) - f(\mathbf{x}^*) &\leq \langle \nabla f(\mathbf{x}_s), \mathbf{x}_s - \mathbf{x}^* \rangle \\ &= \frac{1}{\mu} \langle \mathbf{x}_s - \mathbf{x}_{s+1}, \mathbf{x}_s - \mathbf{x}^* \rangle \\ &\stackrel{(a)}{=} \frac{1}{2\mu} (\|\mathbf{x}_s - \mathbf{x}^*\|^2 + \|\mathbf{x}_s - \mathbf{x}_{s+1}\|^2 - \|\mathbf{x}_{s+1} - \mathbf{x}^*\|^2) \\ &\stackrel{(b)}{=} \frac{1}{2\mu} (\|\mathbf{x}_s - \mathbf{x}^*\|^2 - \|\mathbf{x}_{s+1} - \mathbf{x}^*\|^2) + \frac{\mu}{2} \|\nabla f(\mathbf{x}_s)\|^2 \end{aligned}$$

Where (a) holds from the fact that

$$2\langle \mathbf{a}, \mathbf{b} \rangle = \|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 - \|\mathbf{a} - \mathbf{b}\|^2,$$

and by taking $\mathbf{a} = \mathbf{x}_s - \mathbf{x}^*$ and $\mathbf{b} = \mathbf{x}_s - \mathbf{x}_{s+1}$, and (b) holds because $\|\mathbf{x}_s - \mathbf{x}_{s+1}\|^2 = \mu^2 \|\nabla f(\mathbf{x}_s)\|^2$. Therefore, summing over all s , and canceling the telescopic sum values, we have

$$\begin{aligned} \sum_{s=1}^t (f(\mathbf{x}_s) - f(\mathbf{x}^*)) &\leq \frac{1}{2\mu} (\|\mathbf{x}_1 - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) + \frac{\mu}{2} \sum_{s=1}^t \|\nabla f(\mathbf{x}_s)\|^2 \\ &\leq \frac{1}{2\mu} \|\mathbf{x}_1 - \mathbf{x}^*\|^2 + \frac{\mu}{2} L^2 t \\ &\leq \frac{R^2}{2\mu} + \frac{\mu}{2} L^2 t. \end{aligned}$$

This implies that

$$\frac{1}{t} \sum_{s=1}^t f(\mathbf{x}_s) - f(\mathbf{x}^*) \leq \frac{R^2}{2\mu t} + \frac{\mu}{2} L^2$$

Therefore, using Jensen's inequality, we conclude that

$$f\left(\frac{1}{t} \sum_{s=1}^t \mathbf{x}_s\right) - f(\mathbf{x}^*) \leq \frac{1}{t} \sum_{s=1}^t (f(\mathbf{x}_s) - f(\mathbf{x}^*)) \leq \frac{R^2}{2\mu t} + \frac{\mu}{2} L^2 \leq \frac{RL}{\sqrt{t}}$$

with $\mu = \frac{R}{L\sqrt{t}}$ and convexity. ■

The previous result implies that to achieve an error $\leq \epsilon$, we need $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$ iterations. A few remarks are in order

- This result cannot be improved,
- It is independent of n (the dimension of the data!)
- We can vary the step size in each iteration: $\mathbf{x}_{s+1} = \mathbf{x}_s - \mu_s \nabla f(\mathbf{x}_s)$, $\mu_s = \frac{R}{L\sqrt{s}}$. This convergence rate is optimal for L-Lipschitz functions.

6.2 Projection onto Sets

We have looked at optimization without constraints. But, what happens if we do have them?

Proposition 11 *Assume Ω is convex and closed, we have the following optimization problem with constraints*

$$\begin{aligned} \min_x & f(\mathbf{x}) \\ \text{s.t. } & \mathbf{x} \in \Omega \end{aligned}$$

Definition 5 (Projection onto Ω) *The projection onto Ω is defined as*

$$\Pi_{\Omega}(\mathbf{x}) = \arg \min_{\mathbf{y} \in \Omega} \|\mathbf{x} - \mathbf{y}\|$$

Examples of projections are:

- $\Omega = \{\mathbf{x} : \mathbf{x} \geq 0\} \Rightarrow$ All $\mathbf{x} < \mathbf{0}$ are set to $\mathbf{0}$.
- $\Omega = \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1\} \Rightarrow \Pi_{\Omega}(\mathbf{y}) = \begin{cases} \mathbf{y}, & \|\mathbf{y}\| \leq 1 \\ \frac{\mathbf{y}}{\|\mathbf{y}\|}, & \|\mathbf{y}\| > 1 \end{cases}$

Lemma 12 *Let $\mathbf{x} \in \Omega$ (convex) and $\mathbf{y} \in \mathbb{R}^n$. Then*

$$\langle \Pi_{\Omega}(\mathbf{y}) - \mathbf{x}, \Pi_{\Omega}(\mathbf{y}) - \mathbf{y} \rangle \leq 0.$$

Corollary 13 *Following from Lemma 12,*

$$\|\Pi_{\Omega}(\mathbf{y}) - \mathbf{x}\|^2 + \|\mathbf{y} - \Pi_{\Omega}(\mathbf{y})\|^2 \leq \|\mathbf{y} - \mathbf{x}\|^2$$

Corollary 14 *For $\mathbf{x} \in \Omega$, $\forall \mathbf{y}$,*

$$\|\Pi_{\Omega}(\mathbf{y}) - \mathbf{x}\| \leq \|\mathbf{y} - \mathbf{x}\|$$

For gradient descent, the iterations look like:

$$\mathbf{x}_{t+1} = \Pi_{\Omega}(\mathbf{x}_t - \mu \nabla f(\mathbf{x}_t))$$

Let $\mathbf{y}_{t+1} = \mathbf{x}_t - \mu \nabla f(\mathbf{x}_t)$, and $\mathbf{x}^* \in \Omega$. Then,

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\| \leq \|\mathbf{y}_{t+1} - \mathbf{x}^*\|$$

Our starting point for Proof of Theorem 10 was $f(\mathbf{x}_s) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}_s), \mathbf{x}_s - \mathbf{x}^* \rangle^2$. The exact same proof works with projections.

Lemma 15 *All of the following are equivalent (\iff):*

- f is convex and L -smooth
- $g(\mathbf{x}) = \frac{L}{2} \mathbf{x}^\top \mathbf{x} - f(\mathbf{x})$ is convex
- $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2$
- $\nabla^2 f(\mathbf{x}) \leq LI$, where I is the identity matrix