

Gradients and Hessians in Hilbert spaces

Jordan Bell

`jordan.bell@gmail.com`

Department of Mathematics, University of Toronto

July 26, 2015

1 Gradients

Let $(X, \langle \cdot, \cdot \rangle)$ be a real Hilbert space. The Riesz representation theorem says that the mapping

$$\Phi(x)(y) = \langle y, x \rangle, \quad \Phi : X \rightarrow X^*,$$

is an isometric isomorphism. Let U be a nonempty open subset of X and let $f : U \rightarrow \mathbb{R}$ be differentiable, with derivative $f' : U \rightarrow \mathcal{L}(X; \mathbb{R}) = X^*$.¹ The **gradient of f** is the function $\text{grad } f : U \rightarrow X$ defined by

$$\text{grad } f = \Phi^{-1} \circ f'.$$

Thus, for $x \in U$, $\text{grad } f(x)$ is the unique element of X satisfying

$$\langle \text{grad } f(x), y \rangle = f'(x)(y), \quad y \in X. \quad (1)$$

Because $\Phi^{-1} : X^* \rightarrow X$ is continuous, if $f \in C^1(U; \mathbb{R})$ then $\text{grad } f \in C(U; X)$, being a composition of two continuous functions.

For example, let T be a bounded self-adjoint operator on X and define $f : X \rightarrow \mathbb{R}$ by

$$f(x) = \frac{1}{2} \langle Tx, x \rangle, \quad x \in X.$$

For $x, h \in X$,

$$f(x+h) - f(x) = \frac{1}{2} \langle Tx, h \rangle + \frac{1}{2} \langle Th, x \rangle + \frac{1}{2} \langle Th, h \rangle = \langle Tx, h \rangle + \frac{1}{2} \langle Th, h \rangle.$$

Thus

$$|f(x+h) - f(x) - \langle Tx, h \rangle| = \frac{1}{2} |\langle Th, h \rangle| \leq \frac{1}{2} \|T\| \|h\|^2 = o(\|h\|),$$

¹See <http://individual.utoronto.ca/jordanbell/notes/weaksymplectic.pdf>, §3; <http://individual.utoronto.ca/jordanbell/notes/frechetderivatives.pdf>

which shows that f is differentiable at h , with $f'(x)(y) = \langle Tx, y \rangle$. Thus by (1), $\text{grad } f(x) = Tx$.

For example, let $T \in \mathcal{L}(X; X)$, let $h \in X$, and define $f : X \rightarrow \mathbb{R}$ by

$$f(x) = \frac{1}{2} \|Tx - h\|^2, \quad x \in X.$$

We calculate that

$$\text{grad } f(x) = T^*Tx - T^*h, \quad x \in X.$$

For $x_0 \in X$, define

$$\phi(t) = \exp(-tT^*T)x_0 + \int_0^t \exp(-(t-s)T^*T)T^*h ds, \quad t \geq 0.$$

It is proved² that ϕ satisfies

$$\phi'(t) = -(\text{grad } f)(\phi(t)), \quad \phi(0) = x_0.$$

For a function $F : X \rightarrow X$, we say that F is L **Lipschitz** if

$$\|F(x) - F(y)\| \leq L \|x - y\|, \quad x, y \in X.$$

The following is a useful inequality for functions whose gradients are Lipschitz.³

Lemma 1. *If $f : X \rightarrow \mathbb{R}$ is differentiable and $\text{grad } f : X \rightarrow X$ is L Lipschitz, then*

$$f(y) \leq f(x) + \langle \text{grad } f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2, \quad x, y \in X.$$

Proof. Let $h = y - x$ and define $g : [0, 1] \rightarrow \mathbb{R}$ by $g(t) = f(x + th)$. By the chain rule, for $0 < t < 1$,

$$g'(t) = f'(x + th)(h) = \langle \text{grad } f(x + th), h \rangle.$$

Thus by the fundamental theorem of calculus,

$$\int_0^1 \langle \text{grad } f(x + th), h \rangle dt = \int_0^1 g'(t) dt = g(1) - g(0) = f(x + h) - f(x) = f(y) - f(x),$$

²cf. J.W. Neuberger, *A Sequence of Problems on Semigroups*, p. 51, Problem 195.

³Juan Peypouquet, *Convex Optimization in Normed Spaces: Theory, Methods and Examples*, p. 15, Lemma 1.30.

and so, using the Cauchy-Schwarz inequality and the fact that $\text{grad } f$ is L Lipschitz,

$$\begin{aligned}
f(y) - f(x) &= \int_0^1 \langle \text{grad } f(x + th) - \text{grad } f(x) + \text{grad } f(x), h \rangle dt \\
&= \langle \text{grad } f(x), h \rangle + \int_0^1 \langle \text{grad } f(x + th) - \text{grad } f(x), h \rangle dt \\
&\leq \langle \text{grad } f(x), h \rangle + \int_0^1 \|\text{grad } f(x + th) - \text{grad } f(x)\| \|h\| dt \\
&\leq \langle \text{grad } f(x), h \rangle + \int_0^1 L \|th\| \|h\| dt \\
&= \langle \text{grad } f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2,
\end{aligned}$$

proving the claim. \square

2 Hessians

Let U be a nonempty open subset of X . We prove that if a function is C^2 then its gradient is C^1 .⁴

Theorem 2. *Let U be an open subset of X . If $f \in C^2(U; \mathbb{R})$, then $\text{grad } f \in C^1(U; X)$, and*

$$f''(x)(u)(v) = \langle v, (\text{grad } f)'(x)(u) \rangle, \quad x \in U, \quad u, v \in X. \quad (2)$$

Proof. That f is C^2 means that $f' : U \rightarrow X^*$ is C^1 . That is, for all $x \in U$, the map $f' : U \rightarrow X^*$ is continuous at x , there is $f''(x) \in \mathcal{L}(X; X^*)$ such that

$$\|f'(x + h) - f'(x) - f''(x)(h)\| = o(\|h\|), \quad (3)$$

as $h \rightarrow 0$, and the map $x \mapsto f''(x)$ is continuous $U \rightarrow \mathcal{L}(X; X^*)$.

Let $x \in U$ and let $h \in X$. Define $\phi_h \in X^*$ by

$$\phi_h(v) = f''(x)(h)(v), \quad v \in X.$$

Define $\nu_x(h) = \Phi^{-1}(\phi_h) \in X$, thus

$$f''(x)(h)(v) = \langle v, \nu_x(h) \rangle, \quad v \in X.$$

It is straightforward that ν_x is linear. Because Φ is an isometric isomorphism,

$$\|\nu_x(h)\| = \|\phi_h\| = \sup_{\|v\| \leq 1} |\phi_h(v)| = \sup_{\|v\| \leq 1} |f''(x)(h)(v)| \leq \|f''(x)\| \|h\|,$$

⁴Rodney Coleman, *Calculus on Normed Vector Spaces*, p. 139, Theorem 6.5.

where $(u, v) \mapsto f''(x)(u)(v)$ is a bilinear form, with

$$\|f''(x)\| = \sup_{\|u\| \leq 1, \|v\| \leq 1} |f''(x)(u)(v)|,$$

showing that $\nu_x : X \rightarrow X$ is a bounded linear operator with $\|\nu_x\| \leq \|f''(x)\|$. For h such that $x + h \in U$ and for $v \in X$,

$$(f'(x + h) - f'(x) - f''(x)(h))(v) = \langle v, \text{grad } f(x + h) - \text{grad } f(x) - \nu_x(h) \rangle,$$

so

$$\begin{aligned} \|f'(x + h) - f'(x) - f''(x)(h)\| &= \sup_{\|v\| \leq 1} |\langle v, \text{grad } f(x + h) - \text{grad } f(x) - \nu_x(h) \rangle| \\ &= \|\text{grad } f(x + h) - \text{grad } f(x) - \nu_x(h)\|. \end{aligned}$$

Thus by (3),

$$\|\text{grad } f(x + h) - \text{grad } f(x) - \nu_x(h)\| = o(\|h\|)$$

as $h \rightarrow 0$, and because $\nu_x \in \mathcal{L}(X; X)$, this means that $\text{grad } f : U \rightarrow X$ is differentiable at x , with $(\text{grad } f)'(x) = \nu_x$. It remains to prove that $x \mapsto \nu_x$ is continuous $U \rightarrow \mathcal{L}(X; X)$, namely that $(\text{grad } f)'$ is continuous. For $x \in U$ and for h with $x + h \in U$,

$$\begin{aligned} \|\nu_{x+h} - \nu_x\| &= \sup_{\|u\| \leq 1} \|\nu_{x+h}(u) - \nu_x(u)\| \\ &= \sup_{\|u\| \leq 1} \sup_{\|v\| \leq 1} |\langle v, \nu_{x+h}(u) - \nu_x(u) \rangle| \\ &= \sup_{\|u\| \leq 1} \sup_{\|v\| \leq 1} |f''(x + h)(u)(v) - f''(x)(u)(v)| \\ &= \|f''(x + h) - f''(x)\|, \end{aligned}$$

and because f'' is continuous on U we get that $x \mapsto \nu_x$ is continuous on U , completing the proof. \square

If $f \in C^2(U; \mathbb{R})$, we proved in the above theorem that $\text{grad } f \in C^1(U; X)$. We call the derivative of $\text{grad } f$ the **Hessian of f** ,⁵

$$\text{Hess } f = (\text{grad } f)', \quad U \rightarrow \mathcal{L}(X; X),$$

and (2) then reads

$$f''(x)(u)(v) = \langle v, \text{Hess } f(x)(u) \rangle, \quad x \in U, \quad u, v \in X.$$

Furthermore, it is a fact that if $f \in C^2(U; \mathbb{R})$, then for each $x \in U$, the bilinear form

$$(u, v) \mapsto f''(x)(u)(v)$$

⁵cf. R. A. Tapia, *The differentiation and integration of nonlinear operators*, pp. 45–101, in *Nonlinear Functional Analysis and Applications* (Louis B. Rall, ed.)

is symmetric.⁶ Thus, for $x \in U$ and $u, v \in X$,

$$\langle v, \text{Hess } f(x)(u) \rangle = \langle u, \text{Hess } f(x)(v) \rangle.$$

Now, using that $\langle \cdot, \cdot \rangle$ is symmetric as X is a real Hilbert space, $(\text{Hess } f(x))^* \in \mathcal{L}(X; X)$ satisfies

$$\langle u, \text{Hess } f(x)(v) \rangle = \langle (\text{Hess } f(x))^* u, v \rangle = \langle v, (\text{Hess } f(x))^* u \rangle.$$

so

$$\langle v, \text{Hess } f(x)(u) \rangle = \langle v, (\text{Hess } f(x))^* u \rangle.$$

Because this is true for all v we have $\text{Hess } f(x)(u) = (\text{Hess } f(x))^* u$, and because this is true for all u we have $\text{Hess } f(x) = (\text{Hess } f(x))^*$, i.e. $\text{Hess } f(x)$ is self-adjoint.

Theorem 3. *If U is an open subset of X and $f \in C^2(U; \mathbb{R})$, then for each $x \in U$ it is the case that $\text{Hess } f(x) \in \mathcal{L}(X; X)$ is self-adjoint.*

3 Critical points

For an open set U in X for $k \geq 1$, and for $f \in C^{k+2}(U; \mathbb{R})$, we say that $x_0 \in U$ is a **critical point of f** if $f'(x_0) = 0$. If x_0 is a critical point of f , let us say that x_0 is a **nondegenerate critical point of f** if $\text{Hess } f(x_0) \in \mathcal{L}(X; X)$ is invertible. The **Morse-Palais lemma**⁷ states that if $f \in C^{k+2}(U; \mathbb{R})$ with $k \geq 1$, $f(0) = 0$, and 0 is a nondegenerate critical point of f , then there is some open subset V of U with $0 \in V$ and a C^k diffeomorphism $\phi : V \rightarrow V$, $\phi(0) = 0$, such that

$$f(x) = \frac{1}{2} \langle \text{Hess } f(0)(\phi(x)), \phi(x) \rangle, \quad x \in V.$$

If x is a critical point of a differentiable function $f : U \rightarrow \mathbb{R}$, we call $f(x)$ a **critical value of f** . If $k \geq n$ and $f \in C^k(\mathbb{R}^n; \mathbb{R})$, **Sard's theorem** tells us that the set of critical values of f has Lebesgue measure 0 and is meager.

For Banach spaces Y and Z , a **Fredholm operator**⁸ is a bounded linear operator $T : Y \rightarrow Z$ such that (i) $\alpha(T) = \dim \ker T < \infty$, (ii) $T(Y)$ is a closed subset of Z , and (iii) $\beta(T) = \dim \ker T^* < \infty$. The **index** of a Fredholm operator T is

$$\text{ind } T = \alpha(T) - \beta(T).$$

⁶Serge Lang, *Real and Functional Analysis*, third ed., p. 344, Theorem 5.3.

⁷Serge Lang, *Differential and Riemannian Manifolds*, p. 182, chapter VII, Theorem 5.1; Kung-ching Chang, *Infinite Dimensional Morse Theory and Multiple Solution Problems*, p. 33, Theorem 4.1; André Avez, *Calcul différentiel*, p. 87, §3; N. A. Bobylev, S. V. Emel'yanov, and S. K. Korovin, *Geometrical Methods in Variational Problems*, p. 360, Theorem 5.5.2; Hajime Urakawa, *Calculus of Variations and Harmonic Maps*, p. 87, chapter 3, §1, Theorem 1.10; Jean-Pierre Aubin and Ivar Ekeland, *Applied Nonlinear Analysis*, p. 52, Theorem 8; Melvyn S. Berger, *Nonlinearity and Functional Analysis: Lectures on Nonlinear Problems in Mathematical Analysis*, p. 355, Theorem 6.5.4.

⁸Martin Schechter, *Principles of Functional Analysis*, second ed., chapter 5.

For a differentiable function $f : U \rightarrow \mathbb{R}$, U an open subset of X , and for $x \in U$, $f'(x) \in \mathcal{L}(X; \mathbb{R}) = X^*$. $f'(x)$ is a Fredholm operator if and only if $\dim \ker f'(x) < \infty$. For U a connected open subset of X and for $f \in C^1(U; \mathbb{R})$, we call f a **Fredholm map** if $f'(x)$ is a Fredholm operator for each $x \in U$. It is a fact that $\text{ind } f'(x) = \text{ind } f'(y)$ for all $x, y \in U$, using that U is connected. We denote this common value by $\text{ind } f$. A generalization of Sard's theorem by Smale here tells us that if X is separable, U is a connected open subset of X , $f \in C^k(U; \mathbb{R})$ is a Fredholm map, and

$$k > \max\{\text{ind } f, 0\},$$

then the set of critical values of f is meager.⁹

A function $f \in C^1(X; \mathbb{R})$ is said to **satisfy the Palais-Smale condition** if (u_k) is a sequence in X such that (i) $\{f(u_k)\}$ is a bounded subset of \mathbb{R} and (ii) $\text{grad } f(u_k) \rightarrow 0$, then $\{u_k\}$ is a precompact subset of X : every subsequence of (u_k) itself has a Cauchy subsequence.

Often when speaking about ordinary differential equations in \mathbb{R}^d , we deal with differentiable functions whose derivatives are locally Lipschitz. \mathbb{R}^d has the Heine-Borel property: a subset K of \mathbb{R}^d is compact if and only if K is closed and bounded. In fact no infinite dimensional Banach space has the Heine-Borel property.¹⁰ Thus a locally Lipschitz function need not be Lipschitz on a bounded subset of X . (On a compact set, the set is covered by balls on which the function is Lipschitz, and then the function is Lipschitz on the compact set with Lipschitz constant equal to the maximum of finitely many Lipschitz constants on the balls.) We denote by \mathcal{C} the set of function $f : X \rightarrow \mathbb{R}$ that are differentiable and such that for each bounded subset A of X , the restriction of $\text{grad } f$ to A is Lipschitz.

The **mountain pass theorem**¹¹ states that if (i) $I \in \mathcal{C}$, (ii) I satisfies the Palais-Smale condition, (iii) $I(0) = 0$, (iv) there are $r, a > 0$ such that $I(u) \geq a$ when $\|u\| = r$, and (v) there is some $v \in X$ satisfying $\|v\| > r$ and $I(v) \leq 0$, then

$$\inf_{g \in \Gamma_v} \sup_{0 \leq t \leq 1} (I \circ g)(t)$$

is a critical value of I , where

$$\Gamma_v = \{g \in C([0, 1]; X) : g(0) = 0, g(1) = v\}.$$

⁹Eberhard Zeidler, *Nonlinear Functional Analysis and its Applications, IV: Applications to Mathematical Physics*, p. 829, Theorem 78.A; Melvyn S. Berger, *Nonlinearity and Functional Analysis: Lectures on Nonlinear Problems in Mathematical Analysis*, p. 125, Theorem 3.1.45.

¹⁰Some Fréchet spaces have the Heine-Borel property, like the space of holomorphic functions on the open unit disc, which is what Montel's theorem says: <http://individual.utoronto.ca/jordanbell/notes/holomorphic.pdf>, Theorem 13.

¹¹Lawrence C. Evans, *Partial Differential Equations*, p. 480, Theorem 2; Antonio Ambrosetti and David Arcoya Ivarez, *An Introduction to Nonlinear Functional Analysis and Elliptic Problems*, p. 48, §5.3.

4 Convexity

We prove that a critical point of a differentiable convex function on an open convex set is a minimum.¹²

Theorem 4. *If A is an open convex set, $f : A \rightarrow \mathbb{R}$ is differentiable and convex, and $x_0 \in A$ is a critical point of f , then $f(x_0) \leq f(x)$ for all $x \in A$.*

Proof. Because f is convex, for $0 < t < 1$,

$$f(tx + (1-t)x_0) \leq tf(x) + (1-t)f(x_0),$$

i.e.

$$\frac{f(x_0 + t(x - x_0)) - f(x_0)}{t} \leq f(x) - f(x_0).$$

Taking $t \rightarrow 0$,

$$f'(x_0)(x - x_0) \leq f(x) - f(x_0),$$

and because x_0 is a critical point,

$$0 \leq f(x) - f(x_0),$$

i.e. $f(x_0) \leq f(x)$. □

We establish equivalent conditions for a differentiable function to be convex.¹³

Theorem 5. *If A is an open convex subset of X and $f : A \rightarrow \mathbb{R}$ is differentiable, then the following are equivalent:*

1. f is convex.
2. $f(y) \geq f(x) + \langle \text{grad } f(x), y - x \rangle$, $x, y \in A$.
3. $\langle \text{grad } f(x) - \text{grad } f(y), x - y \rangle \geq 0$, $x, y \in A$.

Proof. Suppose (1). For $x, y \in A$ and $0 < t < 1$, that f is convex means $f(ty + (1-t)x) \leq tf(y) + (1-t)f(x)$, i.e.

$$\frac{f(x + t(y - x)) - f(x)}{t} \leq f(y) - f(x),$$

and taking $t \rightarrow 0$ yields

$$f'(x)(y - x) \leq f(y) - f(x),$$

¹²N. A. Bobylev, S. V. Emel'yanov, and S. K. Korovin, *Geometrical Methods in Variational Problems*, p. 39, Theorem 2.1.4.

¹³Juan Peypouquet, *Convex Optimization in Normed Spaces: Theory, Methods and Examples*, p. 38, Proposition 3.10.

i.e.

$$\langle \text{grad } f(x), y - x \rangle \leq f(y) - f(x).$$

Suppose (2) and let $x, y \in A$, for which

$$\langle \text{grad } f(x), y - x \rangle \leq f(y) - f(x), \quad \langle \text{grad } f(y), x - y \rangle \leq f(x) - f(y).$$

Adding these inequalities,

$$\langle \text{grad } f(x), y - x \rangle - \langle \text{grad } f(y), y - x \rangle \leq 0.$$

Suppose (3), let $x, y \in A$, and define $\phi : [0, 1] \rightarrow \mathbb{R}$ by

$$\phi(t) = f(tx + (1 - t)y) - tf(x) - (1 - t)f(y).$$

$\phi(0) = 0$ and $\phi(1) = 0$, and for $0 < t < 1$, using the chain rule gives

$$\begin{aligned} \phi'(t) &= f'(tx + (1 - t)y)(x - y) - f(x) + f(y) \\ &= \langle \text{grad } f(tx + (1 - t)y), x - y \rangle - f(x) + f(y). \end{aligned}$$

Let $0 < s < t < 1$, let $u = sx + (1 - s)y$ and $v = tx + (1 - t)y$, which both belong to A because A is convex, and so the above reads

$$\phi'(s) = \langle \text{grad } f(u), x - y \rangle - f(x) + f(y), \quad \phi'(t) = \langle \text{grad } f(v), x - y \rangle - f(x) + f(y),$$

so

$$\phi'(s) - \phi'(t) = \langle \text{grad } f(u) - \text{grad } f(v), x - y \rangle.$$

And

$$(s - t)(x - y) = u - y - (v - y) = u - v,$$

so

$$\phi'(s) - \phi'(t) = \frac{1}{s - t} \langle \text{grad } f(u) - \text{grad } f(v), u - v \rangle.$$

But (3) tells us

$$\langle \text{grad } f(u) - \text{grad } f(v), u - v \rangle \geq 0,$$

so, as $s - t < 0$,

$$\phi'(s) - \phi'(t) \leq 0,$$

showing that ϕ' is nondecreasing. On the other hand, because $\phi(0) = 0$ and $\phi(1) = 0$, by the mean value theorem there is some $0 < t_0 < 1$ for which $\phi'(t_0) = 0$. Therefore, because ϕ' is nondecreasing it holds that

$$\phi'(t) \leq 0, \quad 0 \leq t \leq t_0,$$

and

$$\phi'(t) \geq 0, \quad t_0 \leq t \leq 1.$$

That is, ϕ is nonincreasing on $[0, t_0]$, and with $\phi(0) = 0$ this yields $\phi(t) \leq 0$ for $t \in [0, t_0]$, and ϕ is nondecreasing on $[t_0, 1]$, and with $\phi(1) = 0$ this yields $\phi(t) \leq 0$ for $t \in [t_0, 1]$. Therefore $\phi(t) \leq 0$ for $t \in [0, 1]$, which means that

$$f(tx + (1 - t)y) - tf(x) - (1 - t)f(y) \leq 0, \quad 0 \leq t \leq 1,$$

showing that f is convex. □

Theorem 6. *If A is an open convex subset of X and $f : A \rightarrow \mathbb{R}$ is twice differentiable, then the following are equivalent:*

1. f is convex.
2. $\langle \text{Hess } f(x)(v), v \rangle \geq 0$, $x \in A$, $v \in X$.

Proof. Suppose (1) and let $x \in A$. From Theorem 5, $v \in X$ and for $t > 0$ with which $x + tv \in A$,

$$\langle \text{grad } f(x + tv) - \text{grad } f(x), tv \rangle \geq 0,$$

i.e.

$$\frac{f'(x + tv)(v) - f'(x)(v)}{t} \geq 0.$$

Taking $t \rightarrow 0$,

$$f''(x)(v)(v) \geq 0,$$

i.e.

$$\langle \text{Hess } f(x)(v), v \rangle \geq 0.$$

Suppose (2), let $x, y \in A$ and define $\phi : [0, 1] \rightarrow \mathbb{R}$ by

$$\phi(t) = f(tx + (1 - t)y) - tf(x) - (1 - t)f(y).$$

Applying the chain rule, for $0 < t < 1$,

$$\phi''(t) = f''(tx + (1 - t)y)(x - y)(x - y),$$

i.e.

$$\phi''(t) = \langle \text{Hess } f(tx + (1 - t)y)(x - y), x - y \rangle \geq 0,$$

showing that ϕ' is nondecreasing. In the proof of Theorem 5 we deduced from ϕ' being nondecreasing and satisfying $\phi(0) = 0$, $\phi(1) = 0$, that f is convex, and the same reasoning yields here that f is convex. \square

We call a function $F : X \rightarrow X$ **β co-coercive** if

$$\langle F(x) - F(y), x - y \rangle \geq \beta \|F(x) - F(y)\|^2.$$

We prove conditions under which the gradient of a differentiable convex function is co-coercive.¹⁴

Theorem 7 (Baillon-Haddad theorem). *Let $f : X \rightarrow \mathbb{R}$ be differentiable and convex and let $L > 0$. Then $\text{grad } f$ is L Lipschitz if and only if $\text{grad } f$ is $\frac{1}{L}$ co-coercive.*

¹⁴Juan Peypouquet, *Convex Optimization in Normed Spaces: Theory, Methods and Examples*, p. 40, Theorem 3.13.

Proof. Suppose that $\text{grad } f$ is L Lipschitz and for $x \in X$, define $h_x : X \rightarrow \mathbb{R}$ by

$$h_x(y) = f(y) - f'(x)(y) = f(y) - \langle \text{grad } f(x), y \rangle.$$

For $y, z \in X$ and $0 < t < 1$, because f is convex,

$$\begin{aligned} h_x(tz + (1-t)y) &= f(tz + (1-t)y) - \langle \text{grad } f(x), tz + (1-t)y \rangle \\ &\leq tf(z) + (1-t)f(y) - \langle \text{grad } f(x), tz + (1-t)y \rangle \\ &= th_x(z) + (1-t)h_x(y), \end{aligned}$$

showing that h_x is convex. For $y, z \in X$,¹⁵

$$h'_x(y)(z) = f'(y)(z) - f'(x)(z),$$

and in particular $\text{grad } h_x(x) = 0$. Thus by Theorem 4,

$$h_x(x) \leq h_x(y), \quad y \in X. \quad (4)$$

For $x, y, z \in X$, by Lemma 1,

$$f(z) \leq f(x) + \langle \text{grad } f(x), z - x \rangle + \frac{L}{2} \|z - x\|^2,$$

so

$$h_y(z) \leq f(x) - \langle \text{grad } f(y), z \rangle + \langle \text{grad } f(x), z - x \rangle + \frac{L}{2} \|z - x\|^2,$$

i.e.

$$h_y(z) \leq h_x(x) + \langle \text{grad } f(x) - \text{grad } f(y), z \rangle + \frac{L}{2} \|z - x\|^2,$$

and applying (4),

$$h_y(y) \leq h_x(x) + \langle \text{grad } f(x) - \text{grad } f(y), z \rangle + \frac{L}{2} \|z - x\|^2. \quad (5)$$

Now,

$$\|\text{grad } f(x) - \text{grad } f(y)\| = \sup_{\|v\| \leq 1} \langle \text{grad } f(x) - \text{grad } f(y), v \rangle$$

so for each $\epsilon > 0$ there is some $v_\epsilon \in X$ with $\|v_\epsilon\| \leq 1$ and

$$\langle \text{grad } f(x) - \text{grad } f(y), v_\epsilon \rangle \geq \|\text{grad } f(x) - \text{grad } f(y)\| - \epsilon.$$

Let $R = \frac{\|\text{grad } f(x) - \text{grad } f(y)\|}{L}$, and applying (5) with $z = x - Rv_\epsilon$ yields

$$\begin{aligned} h_y(y) &\leq h_x(x) + \langle \text{grad } f(x) - \text{grad } f(y), x - Rv_\epsilon \rangle + \frac{L}{2} \|Rv_\epsilon\|^2 \\ &= h_x(x) + \langle \text{grad } f(x) - \text{grad } f(y), x \rangle - R \langle \text{grad } f(x) - \text{grad } f(y), v_\epsilon \rangle \\ &\quad + \frac{1}{2L} \|\text{grad } f(x) - \text{grad } f(y)\|^2 \|v_\epsilon\|^2 \\ &\leq h_x(x) + \langle \text{grad } f(x) - \text{grad } f(y), x \rangle - R \|\text{grad } f(x) - \text{grad } f(y)\| + R\epsilon \\ &\quad + \frac{1}{2L} \|\text{grad } f(x) - \text{grad } f(y)\|^2 \\ &= h_x(x) + \langle \text{grad } f(x) - \text{grad } f(y), x \rangle - \frac{1}{2L} \|\text{grad } f(x) - \text{grad } f(y)\|^2 + R\epsilon. \end{aligned}$$

¹⁵Henri Cartan, *Differential Calculus*, p. 29, Proposition 2.4.2.

Likewise, because R does not change when x and y are switched,

$$h_x(x) \leq h_y(y) + \langle \text{grad } f(y) - \text{grad } f(x), y \rangle - \frac{1}{2L} \|\text{grad } f(y) - \text{grad } f(x)\|^2 + R\epsilon.$$

Adding these inequalities,

$$\begin{aligned} 0 &\leq \langle \text{grad } f(x) - \text{grad } f(y), x \rangle + \langle \text{grad } f(y) - \text{grad } f(x), y \rangle \\ &\quad - \frac{1}{L} \|\text{grad } f(x) - \text{grad } f(y)\|^2 + 2R\epsilon, \end{aligned}$$

i.e.

$$\frac{1}{L} \|\text{grad } f(x) - \text{grad } f(y)\|^2 \leq \langle \text{grad } f(x) - \text{grad } f(y), x - y \rangle + 2R\epsilon.$$

This is true for all $\epsilon > 0$, so

$$\frac{1}{L} \|\text{grad } f(x) - \text{grad } f(y)\|^2 \leq \langle \text{grad } f(x) - \text{grad } f(y), x - y \rangle,$$

showing that $\text{grad } f$ is $\frac{1}{L}$ co-coercive.

Suppose that $\text{grad } f$ is $\frac{1}{L}$ co-coercive and let $x, y \in X$. Then applying the Cauchy-Schwarz inequality,

$$\begin{aligned} \|\text{grad } f(x) - \text{grad } f(y)\|^2 &\leq L \langle \text{grad } f(x) - \text{grad } f(y), x - y \rangle \\ &\leq L \|\text{grad } f(x) - \text{grad } f(y)\| \|x - y\|. \end{aligned}$$

If $\|\text{grad } f(x) - \text{grad } f(y)\| = 0$ then certainly $\|\text{grad } f(x) - \text{grad } f(y)\| \leq L \|x - y\|$. Otherwise, dividing by $\|\text{grad } f(x) - \text{grad } f(y)\|$ gives

$$\|\text{grad } f(x) - \text{grad } f(y)\| \leq L \|x - y\|,$$

showing that $\text{grad } f$ is L Lipschitz. □