

## Lecture 9 — February 5th, Université de Montréal

Lecturer: Simon Lacoste-Julien

Scribe: Sarah Benamara, Elyes Lamouchi

## 9.1 Convex Optimization:

**Some definitions for context:** [S.Boyd and L.Vandenberghe, 2004], [Zhou, 2018]**Definition.** a set  $C$  is convex if  $\forall x, y \in C, \exists \theta \in [0, 1]$  such that  $\theta x + (1 - \theta)y \in C$ .**Definition.** a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if  $\begin{cases} i) \text{ the domain of } f \text{ is a convex set} \\ ii) \forall x, y \in \text{dom}(f), \theta \in [0, 1], f(\theta x + (1 - \theta)y) \\ \leq \theta f(x) + (1 - \theta)f(y) \end{cases}$ **Definition.**  $f$  is said to be  $\mu$ -strongly convex if  $\exists \mu > 0$  such that  $f(x) + \mu\|x\|^2$  is convex.Let  $f$  be a  $\mu$ -strongly convex, and  $g(x) = f(x) + \mu\|x\|^2$ .**Definition.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a convex function and let  $x \in \text{dom}(f)$ . An element  $u \in \mathbb{R}^n$  is called a **subgradient** of  $f$  at  $x$  if:

$$\langle u, y - x \rangle \leq f(y) - f(x), \forall y \in \text{dom}(f)$$

We call the collection of all subgradients of  $f(x)$  at  $x$  the **subdifferential** of  $f$  at  $x$ , and we denote it by  $\partial f(x)$ .**Definition.** The directional derivative of  $g$  is  $g'(x, v) \triangleq \lim_{\alpha \rightarrow 0} \frac{g(x + \alpha v) - g(x)}{\alpha}$ By convexity of  $g$  we have,  $\forall \alpha \in [0, 1]$ :

$$\begin{aligned} g(\alpha y + (1 - \alpha)x) &\leq \alpha g(y) + (1 - \alpha)g(x) \iff \frac{g(\alpha y + (1 - \alpha)x) - g(x)}{\alpha} \leq g(y) - g(x) \\ &\iff \lim_{\alpha \rightarrow 0} \frac{g(\alpha y + (1 - \alpha)x) - g(x)}{\alpha} \leq g(y) - g(x) \iff g'(x, y - x) \leq g(y) - g(x) \end{aligned}$$

*Note.* If  $g$  is convex in  $y - x$  and  $g$  is finite in  $x$  then the limit in  $*$  exists.Now the trick is to link the product of the subgradient and  $y - x$  with the directional derivative.**Proposition.**  $g(x + t(y - x)) = \sup_{v \in \partial g(x)} \langle v, y - x \rangle$ *Proof.*  $\forall v \in \partial g(x), z \in \text{dom}(g) \xrightarrow{\text{by definition}} \langle v, z - x \rangle \leq g(z) - g(x).$

$$\begin{aligned}
& \forall t \in [0, 1], \forall x, y \in \text{dom}(g) \text{ such that } z = ty + (1 - t)x \xrightarrow{\text{by convexity of } g} z \in \text{dom}(g). \\
& \iff \langle v, z - x \rangle = \langle v, (x + t(y - x)) - x \rangle \leq g(x + t(y - x)) - g(x) \\
& \iff t \langle v, y - x \rangle \leq g(x + t(y - x)) - g(x) \iff \langle v, y - x \rangle \leq \frac{g(x + t(y - x)) - g(x)}{t} = g'(x, y - x).
\end{aligned}$$

Thus we have the inequality,

$$\langle v, y - x \rangle \leq g(y) - g(x)$$

**Proposition.** If  $f, h$  are convex functions on  $\mathbb{R}^n$  and  $\lambda, \mu > 0$  then the subgradient of the function  $g(x) = \lambda f(x) + \mu h(x)$  satisfies  $\partial g(x) = \lambda \partial f(x) + \mu \partial h(x)$ .

Moreover  $f$  and  $h(x) = \|x\|^2$  are convex functions, thus the subgradient of the function  $g(x) = f(x) + \mu h(x)$  satisfies  $\partial g(x) = \partial f(x) + \mu \partial h(x)$ .

Let  $u, w$  be such that  $u \in \partial f(x)$  and  $w \in \partial h(x)$ ,

$$\langle v, y - x \rangle = \langle u + \mu w, y - x \rangle$$

Therefore, by the first order convexity condition on  $g$  i.e  $\langle v, y - x \rangle \leq g(y) - g(x)$ , we have

$$f(y) \geq f(x) + \langle u, y - x \rangle + \frac{\mu}{2} \|y - x\|^2, \forall u \in \partial f(x)$$

$\implies$  We have thus proven that saying that  $f(x)$  is  $\mu$ -strongly convex is equivalent to saying that,  $\forall x, y \in \text{dom}(f), f(y) \geq f(x) + \langle u, y - x \rangle + \frac{\mu}{2} \|y - x\|^2, \forall u \in \partial f(x)$ .

## 9.2 Fundamental Descent Lemma:

**Definition.** A function  $f$  has Lipschitz continuous gradient if:  $\forall x, y \in \text{dom}(f), \exists L > 0$  such that  $\|\nabla f(x) - \nabla f(y)\|^2 \leq L \|x - y\|^2$ .

We say that a function is  $L$ -smooth if its gradient is Lipschitz continuous for some  $L > 0$ .

**Theorem.**  $f$  is  $L$ -smooth  $\implies f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2$ .

*Proof.* [P.Bertsekas, ]: Let  $g$  be such that:

$$\begin{aligned}
g(\alpha) & \triangleq f(x + \alpha(y - x)) \longrightarrow f(y) - f(x) = g(1) - g(0) = \int_0^1 \frac{dg}{d\alpha}(\alpha) d\alpha = \int_0^1 \frac{df}{d\alpha}(x + \alpha(y - x)) d\alpha = \\
& \int_0^1 (y - x) \nabla f(x + \alpha(y - x)) d\alpha = \int_0^1 (y - x) [\nabla f(x + \alpha(y - x)) - \nabla f(x)] d\alpha + \int_0^1 (y - x) \nabla f(x) d\alpha \\
& \leq \left\| \int_0^1 (y - x) \nabla f(x + \alpha(y - x)) d\alpha \right\| + \int_0^1 (y - x) \nabla f(x) d\alpha
\end{aligned}$$

$$\begin{aligned}
& \stackrel{\text{by the triangle inequality}}{\leq} \int_0^1 \left\| (y-x) \nabla f(x + \alpha(y-x)) \right\| d\alpha + \int_0^1 (y-x) \nabla f(x) d\alpha \\
& \stackrel{\text{by Cauchy-Schwartz}}{\leq} \int_0^1 \|y-x\| \cdot \|\nabla f(x + \alpha(y-x))\| d\alpha + (y-x) \Delta f(x) \\
& \stackrel{\text{by Lipschitz continuity of } f}{\leq} \|y-x\| \int_0^1 L \alpha d\alpha + (y-x) \Delta f(x) \\
& = \langle \nabla f(x), y-x \rangle + L \|y-x\|^2
\end{aligned}$$

### Usefulness of this lemma:

**Proposition.** Let  $f$  be an  $L$ -Lipschitz function,  $\forall x \in \text{dom}(f), \exists \gamma$  such that  $x - \gamma \nabla f(x) \in \text{dom}(f) \implies f(x - \gamma \nabla f(x)) \leq f(x) - [\gamma(1 - \frac{\gamma L}{2})] \cdot \|\nabla f(x)\|^2$ .

*Proof.*  $f(x - \gamma \nabla f(x)) \leq f(x) - \gamma \langle \nabla f(x), \nabla f(x) \rangle + \frac{L}{2} \|x - \gamma \nabla f(x) - x\|^2 = f(x) - [\gamma(1 - \frac{\gamma L}{2})] \cdot \|\nabla f(x)\|^2$

**Gradient Descent** is an optimization method, in which we minimize an objective function  $f$  by following the direction of steepest descent (i.e the gradient) in each iteration :  $x_{t+1} \leftarrow x_t - \gamma \nabla f(x_t)$ . This is where the above proposition comes in handy:

$$\begin{aligned}
f(x - \gamma \nabla f(x)) - f(x) & \leq [\gamma(\frac{\gamma L}{2} - 1)] \cdot \|\nabla f(x)\|^2 \\
\iff \min_{\gamma} [f(x - \gamma \nabla f(x)) - f(x)] & \leq \min_{\gamma} [\gamma(\frac{\gamma L}{2} - 1)] \cdot \|\nabla f(x)\|^2 \\
\iff f(y_{\gamma*}) \leq f(x) - \frac{1}{2L} \|\nabla f(x)\|^2, & \text{ where } y_{\gamma*} = y_{\frac{1}{L}} = x - \frac{1}{L} \nabla f(x).
\end{aligned}$$

## 9.3 Convergence rates of gradient descent:

**The  $L$ -smooth,  $\mu$ -convex case:** [Liu, 2015]

As shown above, we take the step to be  $\gamma = \frac{1}{L}$  so that, at the  $k^{\text{th}}$  iteration we have:  $x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k) \iff -\frac{1}{L} \nabla f(x_k) = x_{k+1} - x_k$ .

By  $L$ -smoothness of  $f$ , we have:

$$\begin{aligned}
f(x_{k+1}) & \leq f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2 = f(x_k) - \frac{1}{L} \nabla f(x_k) + \frac{L}{2} \left\| -\frac{1}{L} \nabla f(x_k) \right\|^2 \\
& = f(x_k) - \frac{1}{2L} \|\nabla f(x_k)\|^2
\end{aligned}$$

And let  $\Delta_k = f(x_k) - f^*$ , where  $f^* = \min_x f(x)$ . By  $\mu$ -strong convexity we have:

$$\begin{aligned} f(x) - f(x_k) &\geq \langle \nabla f(x_k), x - x_k \rangle + \frac{\mu}{2} \|x - x_k\|^2 \iff \min_x f(x) - f(x_k) \geq \langle \nabla f(x_k), x^* - x_k \rangle + \frac{\mu}{2} \|x^* - x_k\|^2 \\ &\iff \Delta_{k+1} \leq \frac{1}{2L} \|\nabla f(x_k)\|^2 \iff \Delta_k - \frac{1}{2L} \|\nabla f(x_k)\|^2 \leq \Delta_{k+1} \leq \left[1 - \frac{\mu}{L}\right] \Delta_k \end{aligned}$$

Thus we have,

$$\lim_{k \rightarrow \infty} \frac{\Delta_{k+1}}{\Delta_k} = 1 - \frac{\mu}{L}$$

Therefore, taking  $\mu \leq L$  we have  $1 - \frac{\mu}{L} \in (0, 1)$ . A convergence rate of this sort is called a **linear rate**.

### The weakly-convex case:

From convexity it follows,

$$f(x^*) - f(x_k) \geq \langle \nabla f(x_k), x^* - x_k \rangle \iff \Delta_k \leq \langle \nabla f(x_k), x_k - x^* \rangle \leq \|x_k - x^*\| \|\nabla f(x_k)\|$$

By Cauchy-Schwartz

Moreover,

$$\begin{aligned} \|x_k - x^*\|^2 &= \|x_k - \frac{1}{L} \nabla f(x_k) - x^*\|^2 = \|x_k - x^*\|^2 - \frac{2}{L} \langle \nabla f(x_k), x_k - x^* \rangle + \frac{1}{L^2} \|\nabla f(x_k)\|^2 \\ &= \|x_k - x^*\|^2 + \frac{1}{L^2} \|\nabla f(x_k)\|^2 - \frac{2}{L} \Delta_k \end{aligned}$$

By recursion, we have:

$$\Delta_{k+1} \leq \Delta_k - \frac{1}{2L \underbrace{\|x_0 - x^*\|^2}_{r_0^2}} \Delta_k^2 \iff \frac{1}{\Delta_{k+1}} \geq \frac{1}{\Delta_k} + \frac{1}{2L r_0^2} \geq \dots \geq \frac{1}{\Delta_0} + \frac{1}{2L r_0^2} (k+1)$$

Thus we have  $\lim_{k \rightarrow \infty} \frac{\Delta_{k+1}}{\Delta_k} = 1$ . And a rate of convergence of  $O(\frac{k}{L r_0^2})$ . Such a convergence is said to have a **linear rate**.

## 9.4 Second order methods:

**Newton's method:** [S.Boyd and L.Vandenberghe, 2004]

By gradient descent, we can also minimize an objective function  $f(x)$  by using its second derivative i.e its Hessian matrix  $\nabla^2 f(x)$ . **The Newton step:** For  $x \in \text{dom}(f)$ , the vector  $\Delta x = -\nabla^2 f(x)^{-1} \nabla f(x)$  is called **the Newton step**.

**Proposition.** The Hessian matrix is semi definite.

The implication of the above proposition is that,

$$\nabla f(x)^T \Delta x = -\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \underset{\text{positive definiteness}}{<} 0$$

Thus going from one iteration to the next, we are always in descent direction.

# Bibliography

- [Liu, 2015] Liu, J. (2015). A simple proof of chernoff's bound, 1998. In *Class notes*, pages 23:1,23:3.
- [P.Bertsekas, ] P.Bertsekas, D. Convergence analysis of gradient methods. In *Class notes*, page 8.
- [S.Boyd and L.Vandenberghe, 2004] S.Boyd and L.Vandenberghe (2004). Newton's method. In *Convex Optimization*, pages 67,484. Cambridge University Press.
- [Zhou, 2018] Zhou, X. (2018). On the fenchel duality between strong convexity and lipschitz continuous gradient. pages 2,3.