

Introduction to Optimization

Lecture 05: Gradient descent. Convex functions.



university of
groningen

Recall: descent methods

Most algorithms are based on the idea of (sufficient) descent: given x_k , find x_{k+1} such that

$$(1) \quad f(x_{k+1}) \leq f(x_k) - \delta_k^2.$$

One way is to find $d_k \in \mathbb{R}^N$ and $\alpha_k > 0$, such that (1) holds with

$$x_{k+1} = x_k - \alpha_k d_k.$$

We say $-d_k$ is a **descent direction**, and α_k is the **step size**, **step length** or **learning rate** (in ML).

L -smoothness

A differentiable function $f : A \subset \mathbb{R}^N \rightarrow \mathbb{R}$ is **L -smooth**, with $L > 0$, if

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$$

for all $x, y \in A$.

L -smoothness

A differentiable function $f : A \subset \mathbb{R}^N \rightarrow \mathbb{R}$ is **L -smooth**, with $L > 0$, if

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$$

for all $x, y \in A$.

Proposition (Descent Lemma)

If f is L -smooth and A is convex, then

$$|f(y) - f(x) - \nabla f(x) \cdot (y - x)| \leq \frac{L}{2}\|x - y\|^2$$

for all $x, y \in A$.



Identifying descent directions

The L -smooth case

By the Descent Lemma, we have

$$f(x_n - \alpha_n d_n) \leq f(x_n) - \alpha_n \nabla f(x_n) \cdot d_n + \frac{\alpha_n^2 L}{2} \|d_n\|^2.$$

Identifying descent directions

The L -smooth case

By the Descent Lemma, we have

$$f(x_n - \alpha_n d_n) \leq f(x_n) - \alpha_n \nabla f(x_n) \cdot d_n + \frac{\alpha_n^2 L}{2} \|d_n\|^2.$$

It is therefore sufficient for α_n and d_n to satisfy

$$0 < \alpha_n L \|d_n\|^2 < 2 \nabla f(x_n) \cdot d_n.$$

Identifying descent directions

The L -smooth case

By the Descent Lemma, we have

$$f(x_n - \alpha_n d_n) \leq f(x_n) - \alpha_n \nabla f(x_n) \cdot d_n + \frac{\alpha_n^2 L}{2} \|d_n\|^2.$$

It is therefore sufficient for α_n and d_n to satisfy

$$0 < \alpha_n L \|d_n\|^2 < 2 \nabla f(x_n) \cdot d_n.$$

Gradient-consistency: $\tau \|d_n\|^2 \leq \|\nabla f(x_n)\|^2 \leq \sigma \nabla f(x_n) \cdot d_n.$

Identifying descent directions

The L -smooth case

By the Descent Lemma, we have

$$f(x_n - \alpha_n d_n) \leq f(x_n) - \alpha_n \nabla f(x_n) \cdot d_n + \frac{\alpha_n^2 L}{2} \|d_n\|^2.$$

It is therefore sufficient for α_n and d_n to satisfy

$$0 < \alpha_n L \|d_n\|^2 < 2 \nabla f(x_n) \cdot d_n.$$

Gradient-consistency: $\tau \|d_n\|^2 \leq \|\nabla f(x_n)\|^2 \leq \sigma \nabla f(x_n) \cdot d_n$.

Sufficient condition for descent: $0 < \inf \alpha_n \leq \sup \alpha_n < \frac{2\tau}{\sigma L}$.

Gradient-consistent methods

Constant stepsize $\alpha_n \equiv \alpha$, for simplicity

Proposition

Let f be L -smooth and bounded from below. Iterate $x_{n+1} = x_n - \alpha d_n$, where d_n is gradient-consistent and $0 < \alpha < \frac{2\tau}{\sigma L}$.

Gradient-consistent methods

Constant stepsize $\alpha_n \equiv \alpha$, for simplicity

Proposition

Let f be L -smooth and bounded from below. Iterate $x_{n+1} = x_n - \alpha d_n$, where d_n is gradient-consistent and $0 < \alpha < \frac{2\tau}{\sigma L}$. Then,

$$\textcircled{1} \quad \exists \lim_{n \rightarrow \infty} f(x_n) \in \mathbb{R}, \text{ and } \lim_{n \rightarrow \infty} \|\nabla f(x_n)\| = 0.$$

Gradient-consistent methods

Constant stepsize $\alpha_n \equiv \alpha$, for simplicity

Proposition

Let f be L -smooth and bounded from below. Iterate $x_{n+1} = x_n - \alpha d_n$, where d_n is gradient-consistent and $0 < \alpha < \frac{2\tau}{\sigma L}$. Then,

- ① $\exists \lim_{n \rightarrow \infty} f(x_n) \in \mathbb{R}$, and $\lim_{n \rightarrow \infty} \|\nabla f(x_n)\| = 0$.
- ② Cluster points are critical: if $x_{k_n} \rightarrow \hat{x}$, then $\nabla f(\hat{x}) = 0$.

Gradient-consistent methods

Constant stepsize $\alpha_n \equiv \alpha$, for simplicity

Proposition

Let f be L -smooth and bounded from below. Iterate $x_{n+1} = x_n - \alpha d_n$, where d_n is gradient-consistent and $0 < \alpha < \frac{2\tau}{\sigma L}$. Then,

- ① $\exists \lim_{n \rightarrow \infty} f(x_n) \in \mathbb{R}$, and $\lim_{n \rightarrow \infty} \|\nabla f(x_n)\| = 0$.
- ② Cluster points are critical: if $x_{k_n} \rightarrow \hat{x}$, then $\nabla f(\hat{x}) = 0$.
- ③ If f has no critical points, then $\lim_{n \rightarrow \infty} \|x_n\| = +\infty$.

Gradient-consistent methods

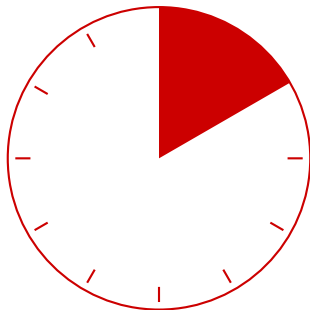
Constant stepsize $\alpha_n \equiv \alpha$, for simplicity

Proposition

Let f be L -smooth and bounded from below. Iterate $\mathbf{x}_{n+1} = \mathbf{x}_n - \alpha \mathbf{d}_n$, where \mathbf{d}_n is gradient-consistent and $0 < \alpha < \frac{2\tau}{\sigma L}$. Then,

- ① $\exists \lim_{n \rightarrow \infty} f(\mathbf{x}_n) \in \mathbb{R}$, and $\lim_{n \rightarrow \infty} \|\nabla f(\mathbf{x}_n)\| = 0$.
- ② Cluster points are critical: if $\mathbf{x}_{k_n} \rightarrow \hat{\mathbf{x}}$, then $\nabla f(\hat{\mathbf{x}}) = 0$.
- ③ If f has no critical points, then $\lim_{n \rightarrow \infty} \|\mathbf{x}_n\| = +\infty$.
- ④ There is $C > 0$ such that $\min \{\|\nabla f(\mathbf{x}_i)\| : 1 \leq i \leq n\} \leq \frac{C}{\sqrt{n}}$.

Break



Convex functions and the gradient method

Theorem

Let f be convex and L -smooth. Iterate $x_{n+1} = x_n - \alpha \nabla f(x_n)$ with $0 < \alpha < \frac{2}{L}$.

Convex functions and the gradient method

Theorem

Let f be convex and L -smooth. Iterate $x_{n+1} = x_n - \alpha \nabla f(x_n)$ with $0 < \alpha < \frac{2}{L}$.

① $\lim_{n \rightarrow \infty} f(x_n) = \inf(f).$



Convex functions and the gradient method

Theorem

Let f be convex and L -smooth. Iterate $x_{n+1} = x_n - \alpha \nabla f(x_n)$ with $0 < \alpha < \frac{2}{L}$.

① $\lim_{n \rightarrow \infty} f(x_n) = \inf(f).$

② If f has minimizers, then x_n converges to one of them, and

$$f(x_n) - \min(f) \leq \frac{D^2}{\alpha(2 - \alpha L)n},$$

where D is the distance from x_0 to its closest minimizer.

Moreover, $\lim_{n \rightarrow \infty} n [f(x_n) - \min(f)] = 0.$

An important tool

Proposition (Baillon-Haddad Lemma)

If f is convex and L -smooth, then

$$\frac{1}{L} |\nabla f(y) - \nabla f(x)|^2 \leq (\nabla f(y) - \nabla f(x)) \cdot (y - x)$$

for all $x, y \in A$.

This will be proved in the tutorial.

Sketch of the proof

First, use

$$(2) \quad \|x_{n+1} - u\|^2 = \alpha^2 \|\nabla f(x_n)\|^2 + \|x_n - u\|^2 - 2\alpha \nabla f(x_n) \cdot (x_n - u)$$

and Baillon-Haddad Lemma to show that

$$\frac{\alpha}{L}(2 - \alpha L) \sum_{n=0}^k \|\nabla f(x_n)\|^2 \leq \|x_0 - u\|^2.$$

Then, use (2) and convexity to prove that

$$2\alpha(f(x_n) - \min(f)) \leq \|x_n - u\|^2 - \|x_{n+1}\|^2 + \alpha^2 \|\nabla f(x_n)\|^2.$$

Sum over n , and combine the two inequalities, to conclude that

$$2\alpha k(f(x_k) - \min(f)) \leq \|x_0 - u\|^2 \left(1 + \frac{\alpha L}{2 - \alpha L}\right).$$

Sketch of the proof, continued

For the convergence, use that cluster points are critical, and that $\|x_n - u\|$ is nonincreasing, to deduce that (x_n) cannot have more than one cluster point. This is sufficient for convergence because (x_n) is bounded.

For the last statement, use the following lemma with $e_n = f(x_n) - \min(f)$:

Lemma

Let (e_n) be a positive, nonincreasing sequence such that $\sum_{n=0}^{\infty} e_n < +\infty$. Then $\lim_{n \rightarrow \infty} ne_n = 0$.

Other choices for α_n

- **Constant:** $\alpha_n \equiv \alpha \in (0, 2/L)$.

Other choices for α_n

- **Constant:** $\alpha_n \equiv \alpha \in (0, 2/L)$.
- **Exact minimization:** $\min_{\alpha > 0} f(x_n - \alpha d_n)$.

Other choices for α_n

- **Constant:** $\alpha_n \equiv \alpha \in (0, 2/L)$.
- **Exact minimization:** $\min_{\alpha > 0} f(x_n - \alpha d_n)$.
- **Limited minimization:** $\min_{\alpha \in (0, A]} f(x_n - \alpha d_n)$.

Other choices for α_n

- **Constant:** $\alpha_n \equiv \alpha \in (0, 2/L)$.
- **Exact minimization:** $\min_{\alpha > 0} f(x_n - \alpha d_n)$.
- **Limited minimization:** $\min_{\alpha \in (0, A]} f(x_n - \alpha d_n)$.
- **Vanishing:** $\alpha_n \rightarrow 0, \sum \alpha_n = \infty$.

Other choices for α_n

- **Constant:** $\alpha_n \equiv \alpha \in (0, 2/L)$.
- **Exact minimization:** $\min_{\alpha > 0} f(x_n - \alpha d_n)$.
- **Limited minimization:** $\min_{\alpha \in (0, A]} f(x_n - \alpha d_n)$.
- **Vanishing:** $\alpha_n \rightarrow 0, \sum \alpha_n = \infty$.
- **Backtracking** (Armijo, Goldstein).



Some remarks

The simplest example, revisited

We had applied the gradient method to the function $f(x) = x^2$. Were the hypotheses on α and the rate of convergence consistent with the previous theorem?

Some remarks

The simplest example, revisited

We had applied the gradient method to the function $f(x) = x^2$. Were the hypotheses on α and the rate of convergence consistent with the previous theorem?

Strong convexity

If the objective function is strongly convex, can we expect the gradient method to converge faster?