

# Optimization for Machine Learning

## CS-439

### Lecture 2: Gradient Descent

**Martin Jaggi**

EPFL – [github.com/epfml/OptML\\_course](https://github.com/epfml/OptML_course)

March 1, 2019

# Chapter 2

## Gradient Descent

# The Algorithm

Get near to a minimum  $\mathbf{x}^*$  / close to the optimal value  $f(\mathbf{x}^*)$ ?

(Assumptions:  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  convex, differentiable, has a global minimum  $\mathbf{x}^*$ )

**Goal:** Find  $\mathbf{x} \in \mathbb{R}^d$  such that

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \varepsilon.$$

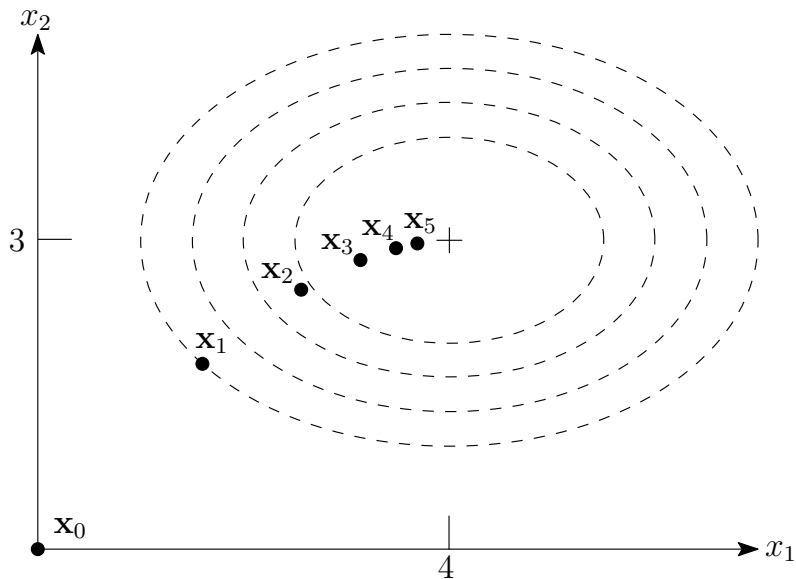
Note that there can be several minima  $\mathbf{x}_1^* \neq \mathbf{x}_2^*$  with  $f(\mathbf{x}_1^*) = f(\mathbf{x}_2^*)$ .

**Iterative Algorithm:**

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t),$$

for **timesteps**  $t = 0, 1, \dots$ , and **stepsize**  $\gamma \geq 0$ .

## Example



# Vanilla analysis

How to bound  $f(\mathbf{x}_t) - f(\mathbf{x}^*)$  ?

- ▶ Abbreviate  $\mathbf{g}_t := \nabla f(\mathbf{x}_t)$ , and consider (using the definition of gradient descent)

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{1}{\gamma} (\mathbf{x}_t - \mathbf{x}_{t+1})^\top (\mathbf{x}_t - \mathbf{x}^*).$$

- ▶ Apply  $2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$  to rewrite

$$\begin{aligned}\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) &= \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \\ &= \frac{\gamma}{2} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2)\end{aligned}$$

- ▶ Sum this up over the iterations  $t$ :

$$\sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2)$$

## Vanilla analysis, II

- ▶ Now we invoke convexity of  $f$  with  $\mathbf{x} = \mathbf{x}_t, \mathbf{y} = \mathbf{x}^\star$ :

$$f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^\star)$$

giving

$$\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^\star)) \leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2,$$

an upper bound for the **average error**  $f(\mathbf{x}_t) - f(\mathbf{x}^\star)$  over the steps

- ▶ last iterate is not necessarily the best one
- ▶ stepsize is crucial

## Lipschitz convex functions: $\mathcal{O}(1/\varepsilon^2)$ steps

Assume that all gradients of  $f$  are bounded in norm.

- Equivalent to  $f$  being Lipschitz (**Exercise 11**).

### Theorem

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable with a global minimum  $\mathbf{x}^*$ ; furthermore, suppose that  $\|\mathbf{x}_0 - \mathbf{x}^*\| \leq R$  and  $\|\nabla f(\mathbf{x})\| \leq B$  for all  $\mathbf{x}$ . Choosing the stepsize

$$\gamma := \frac{R}{B\sqrt{T}},$$

*gradient descent yields*

$$\frac{1}{T} \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{RB}{\sqrt{T}}.$$

## Lipschitz convex functions: $\mathcal{O}(1/\varepsilon^2)$ steps, II

Proof.





## Lipschitz convex functions: $\mathcal{O}(1/\varepsilon^2)$ steps, III

$$T \geq \frac{R^2 B^2}{\varepsilon^2} \quad \Rightarrow \quad \text{average error} \leq \frac{RB}{\sqrt{T}} \leq \varepsilon.$$

### Advantages:

- ▶ dimension-independent!
- ▶ holds for both average, or best iterate

### In Practice:

What if we don't know  $R$  and  $B$ ?

→ **Exercise 13**

# Smooth functions

## “Not too curved”

### Definition

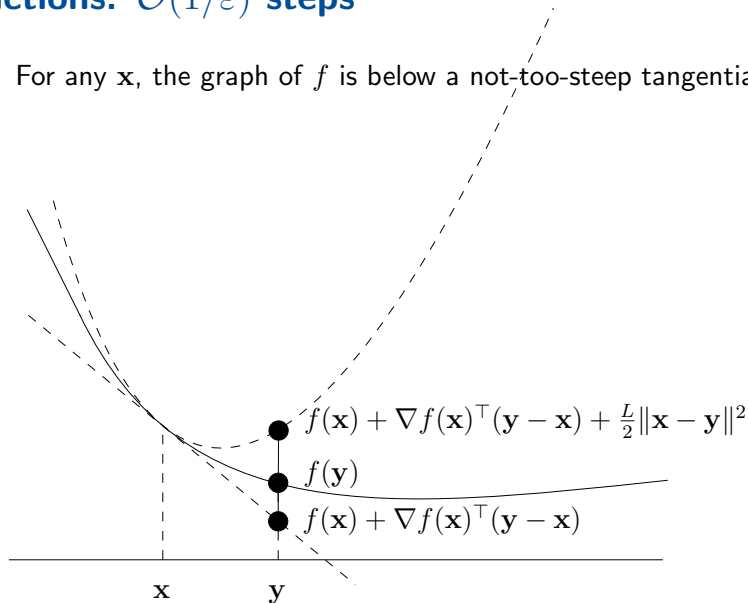
Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable.  $f$  is called **smooth** (with parameter  $L \geq 0$ ) if

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d.$$

Definition does not require convexity (useful later)

## Smooth functions: $\mathcal{O}(1/\varepsilon)$ steps

**Smoothness:** For any  $\mathbf{x}$ , the graph of  $f$  is below a not-too-steep tangential paraboloid at  $(\mathbf{x}, f(\mathbf{x}))$ :



## Smooth functions: $\mathcal{O}(1/\varepsilon)$ steps

- ▶ Quadratic functions are smooth (**Exercise 11**)
- ▶ Operations that preserve smoothness:

### Lemma (Exercise 14)

- (i) Let  $f_1, f_2, \dots, f_m$  be convex functions that are smooth with parameters  $L_1, L_2, \dots, L_m$ , and let  $\lambda_1, \lambda_2, \dots, \lambda_m \in \mathbb{R}_+$ . Then the convex function  $f := \sum_{i=1}^m \lambda_i f_i$  is smooth with parameter  $\sum_{i=1}^m \lambda_i L_i$ .
- (ii) Let  $f$  be convex and smooth with parameter  $L$ , and let  $g(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$ , for  $A \in \mathbb{R}^{d \times m}$  and  $\mathbf{b} \in \mathbb{R}^d$ . Then the convex function  $f \circ g$  is smooth with parameter  $L\|A\|^2$ , where

$$\|A\| = \max_{\mathbf{x} \neq 0} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|}$$

is the **2-norm** (or spectral norm) of  $A$ .

# Smooth vs Lipschitz

- ▶ Bounded gradients  $\Leftrightarrow$  Lipschitz continuity of  $f$ ,
- ▶ Now: smoothness  $\Leftrightarrow$  Lipschitz continuity of  $\nabla f$ .

## Lemma

*Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable. The following two statements are equivalent.*

- (i)  *$f$  is smooth with parameter  $L$ .*
- (ii)  *$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ .*

## Smooth functions: $\mathcal{O}(1/\varepsilon)$ steps

### Theorem

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable with a global minimum  $\mathbf{x}^*$ ; furthermore, suppose that  $f$  is smooth with parameter  $L$ . Choosing

$$\gamma := \frac{1}{L},$$

gradient descent with arbitrary  $\mathbf{x}_0$  satisfies

(i) Function values are monotone decreasing:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2, \quad t \geq 0.$$

(ii) Use the fact that

$2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$  to obtain

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

# Smooth functions: $\mathcal{O}(1/\varepsilon)$ steps. Proof

Proof.



## Smooth functions: $\mathcal{O}(1/\varepsilon)$ steps

In Practice:

What if we don't know the smoothness parameter  $L$ ?

→ **Exercise 15**