

Annotated
Version

Optimization for Machine Learning

CS-439

Lecture 2: Gradient Descent

Martin Jaggi

EPFL – github.com/epfml/OptML_course

March 2, 2018

Recap

Convexity

recap,

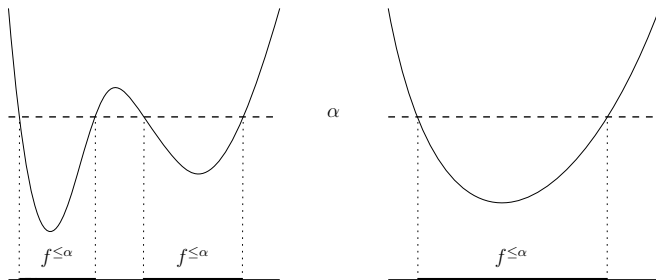
and short addition before we get to gradient descent...

Existence of a minimizer

Sublevel sets: Let $f : \text{dom}(f) \rightarrow \mathbb{R}$, $\alpha \in \mathbb{R}$. The set

$$f^{\leq \alpha} := \{\mathbf{x} \in \text{dom}(f) : f(\mathbf{x}) \leq \alpha\}$$


is the α -sublevel set of f ;



Weierstrass Theorem

Theorem

Let $f : \text{dom}(f) \rightarrow \mathbb{R}$ be a convex function, $\text{dom}(f)$ open, and suppose there is a nonempty and bounded sublevel set $f^{\leq \alpha}$. Then f has a global minimum.



Proof.



[github.com/epfml/...](https://github.com/epfml)

Chapter 2

Gradient Descent

The Algorithm

How to get near to a minimum \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

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

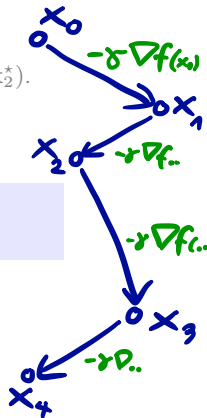
$\varepsilon > 0$ accuracy

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:

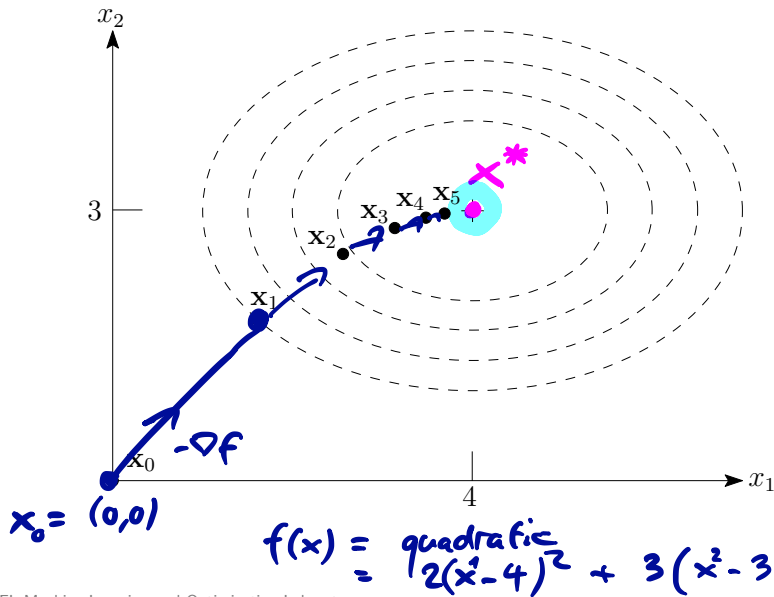
$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t),$$

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



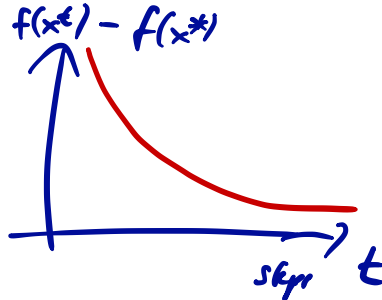
Example

$d=2$



Vanilla analysis

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

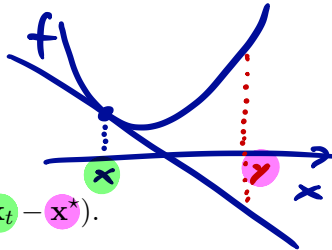


Vanilla analysis

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

- Convexity of f , for $\mathbf{x} = \mathbf{x}_t, \mathbf{y} = \mathbf{x}^*$, gives

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*).$$



$$f(y) \geq f(x) + \nabla f(x)^\top (y - x)$$

$$\begin{aligned} \mathbf{x}_{t+1} &:= \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t) \\ -(\mathbf{x}_{t+1} - \mathbf{x}_t) &= \gamma \nabla f(\mathbf{x}_t) \quad \text{GD} \end{aligned}$$

Vanilla analysis

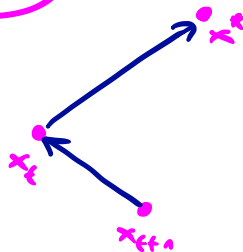
How to bound $f(\mathbf{x}_t) - f(\mathbf{x}^*)$?

- Convexity of f , for $\mathbf{x} = \mathbf{x}_t, \mathbf{y} = \mathbf{x}^*$, gives

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*).$$

- Apply the definition of the iteration, $\nabla f(\mathbf{x}_t) = (\mathbf{x}_t - \mathbf{x}_{t+1})/\gamma$:

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



Vanilla analysis

$$f(\mathbf{x}^*) \geq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}^* - \mathbf{x}_t) + \mu \|\mathbf{x}^* - \mathbf{x}_t\|^2$$

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

- Convexity of f , for $\mathbf{x} = \mathbf{x}_t, \mathbf{y} = \mathbf{x}^*$, gives

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*) - \mu \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

- Apply the definition of the iteration, $\nabla f(\mathbf{x}_t) = (\mathbf{x}_t - \mathbf{x}_{t+1})/\gamma$:

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{1}{\gamma} (\mathbf{x}_t - \mathbf{x}_{t+1})^\top (\mathbf{x}_t - \mathbf{x}^*)$$

- Now we apply $2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$

$$\begin{aligned} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \\ &\stackrel{GD}{=} \frac{1}{2\gamma} (\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \end{aligned}$$

again by the definition gradient descent

Vanilla analysis, cont.

sum this over steps $t = 0, \dots, T - 1$:

$$\begin{aligned} & \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \\ & \leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \left(\|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2) \right) \\ & \leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \end{aligned}$$

an upper bound for the **average error** $f(\mathbf{x}_t) - f(\mathbf{x}^*)$, $t = 0 \dots T - 1$

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

Bounded gradients: $\mathcal{O}(1/\varepsilon^2)$ steps

Assume that all gradients of f are bounded in norm.

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 $\|\nabla f(\mathbf{x})\| \leq L$ for all \mathbf{x} . Choosing the stepsize

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

\Downarrow gradient descent yields

f being
Lipschitz continuous
with constant L

$$\frac{1}{T} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{RL}{\sqrt{T}} \in \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

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

Bounded gradients: $\mathcal{O}(1/\varepsilon^2)$ steps, II

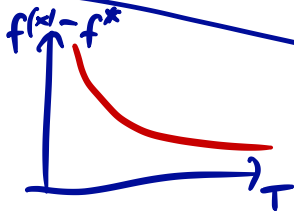
Proof. use (*)!

$$\sum_{t=0}^{T-1} f(x_t) - f^* \leq \underbrace{\frac{\sigma}{2} T L^2}_{\leq \|D\|} + \underbrace{\frac{1}{2\sigma} R^2}_{\|x_0 - x^*\|^2}$$

plug in $\sigma = \frac{R}{L\sqrt{T}}$:

$$= T \frac{LR}{\sqrt{T}}$$

□



divide by T

Bounded gradients: $\mathcal{O}(1/\varepsilon^2)$ steps, II

Advantages:

- ▶ dimension-independent! $d \gg 0$
- ▶ holds for both average, or best iterate

In Practice:

What if we don't know R and L ?

→ Exercise 13

\Leftrightarrow " $1/\varepsilon^2$ steps "

$$f(x_\varepsilon) - f^* \leq \mathcal{O}\left(\frac{1}{\sqrt{T}}\right) \leq \varepsilon$$

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

Convex, but not too convex?

Definition

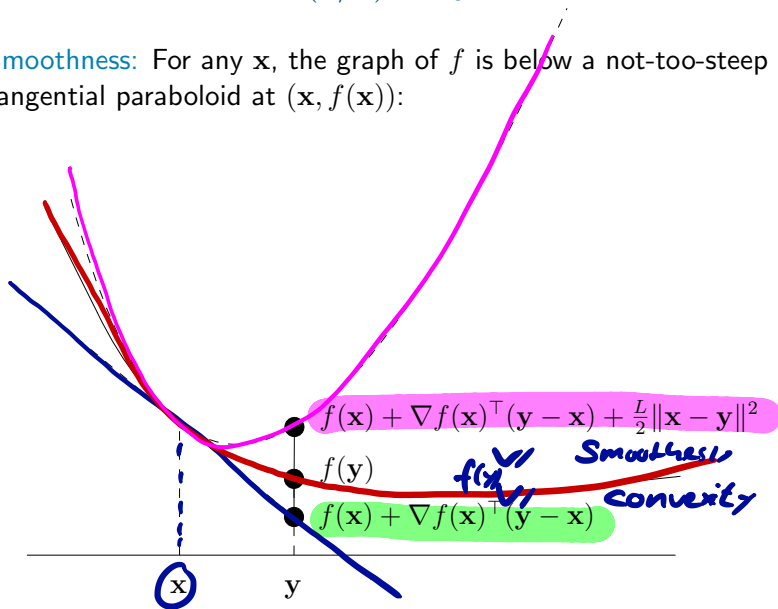
Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be convex and differentiable, $L \in \mathbb{R}_+$. f is called **smooth** (with parameter L) if

Smoothness

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

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
- ▶ Operations that preserve smoothness:

$$f(x) := \|x\|^2$$

Lemma (Exercise 15)

- (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 $g(x) = Ax + b$, for $A \in \mathbb{R}^{d \times m}$ and $b \in \mathbb{R}^d$. Then the convex function $f \circ g$ is smooth with parameter $L\|A\|^2$, where

$$\|A\| = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \max_{\|x\|=1} \|Ax\|$$

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

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

Convergence proof: See next lecture