

Optimization for Machine Learning

CS-439

Lecture 4: Projected, Proximal and Subgradient Descent

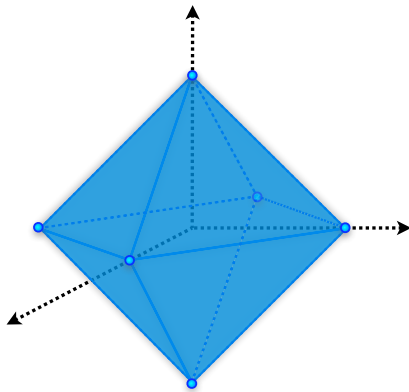
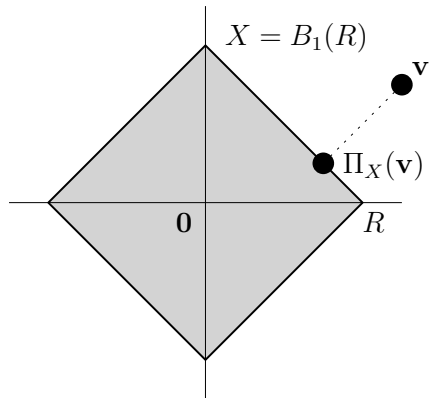
Martin Jaggi

EPFL – github.com/epfml/0ptML_course

March 15, 2019

Projecting onto ℓ_1 -balls

$$X = B_1(R) := \left\{ \mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i| \leq R \right\}$$



2^d facets!

Projecting onto ℓ_1 -balls

Theorem

Let $\mathbf{v} \in \mathbb{R}^d$, $R \in \mathbb{R}_+$, $X = B_1(R)$ the ℓ_1 -ball around $\mathbf{0}$ of radius R . The projection

$$\Pi_X(\mathbf{v}) = \operatorname{argmin}_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{v}\|^2$$

of \mathbf{v} onto $B_1(R)$ can be computed in time $\mathcal{O}(d \log d)$.

This can be improved to time $\mathcal{O}(d)$ by avoiding sorting.

Section 3.6

Proximal Gradient Descent

Composite optimization problems

Consider objective functions composed as

$$f(\mathbf{x}) := g(\mathbf{x}) + h(\mathbf{x})$$

where g is a “nice” function, where as h is a “simple” additional term, which however doesn't satisfy the assumptions of niceness which we used in the convergence analysis so far.

In particular, an important case is when h is not differentiable.

Idea

The classical gradient step for minimizing g :

$$\mathbf{x}_{t+1} = \operatorname{argmin}_{\mathbf{y}} g(\mathbf{x}_t) + \nabla g(\mathbf{x}_t)^\top (\mathbf{y} - \mathbf{x}_t) + \frac{1}{2\gamma} \|\mathbf{y} - \mathbf{x}_t\|^2 .$$

For the stepsize $\gamma := \frac{1}{L}$ it exactly minimizes the local quadratic model of g at our current iterate \mathbf{x}_t , formed by the smoothness property with parameter L .

Now for $f = g + h$, keep the same for g , and add h unmodified.

$$\begin{aligned} \mathbf{x}_{t+1} &:= \operatorname{argmin}_{\mathbf{y}} g(\mathbf{x}_t) + \nabla g(\mathbf{x}_t)^\top (\mathbf{y} - \mathbf{x}_t) + \frac{1}{2\gamma} \|\mathbf{y} - \mathbf{x}_t\|^2 + h(\mathbf{y}) \\ &= \operatorname{argmin}_{\mathbf{y}} \frac{1}{2\gamma} \|\mathbf{y} - (\mathbf{x}_t - \gamma \nabla g(\mathbf{x}_t))\|^2 + h(\mathbf{y}) , \end{aligned}$$

the **proximal gradient descent** update.

The proximal gradient descent algorithm

An iteration of proximal gradient descent is defined as

$$\mathbf{x}_{t+1} := \text{prox}_{h,\gamma}(\mathbf{x}_t - \gamma \nabla g(\mathbf{x}_t)) .$$

where the proximal mapping for a given function h , and parameter $\gamma > 0$ is defined as

$$\text{prox}_{h,\gamma}(\mathbf{z}) := \underset{\mathbf{y}}{\text{argmin}} \left\{ \frac{1}{2\gamma} \|\mathbf{y} - \mathbf{z}\|^2 + h(\mathbf{y}) \right\} .$$

The update step can be equivalently written as

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma G_\gamma(\mathbf{x}_t)$$

for $G_{h,\gamma}(\mathbf{x}) := \frac{1}{\gamma} \left(\mathbf{x} - \text{prox}_{h,\gamma}(\mathbf{x} - \gamma \nabla g(\mathbf{x})) \right)$ being the so called generalized gradient of f .

A generalization of gradient descent?

- ▶ $h \equiv 0$: recover gradient descent
- ▶ $h \equiv \iota_X$: recover projected gradient descent!

Given a closed convex set X , the indicator function of the set X is given as the convex function

$$\begin{aligned} \iota_X : \mathbb{R}^d &\rightarrow \mathbb{R} \cup +\infty \\ \mathbf{x} &\mapsto \iota_X(\mathbf{x}) := \begin{cases} 0 & \text{if } \mathbf{x} \in X, \\ +\infty & \text{otherwise.} \end{cases} \end{aligned}$$

Proximal mapping becomes

$$\text{prox}_{h,\gamma}(\mathbf{z}) := \underset{\mathbf{y}}{\text{argmin}} \left\{ \frac{1}{2\gamma} \|\mathbf{y} - \mathbf{z}\|^2 + \iota_X(\mathbf{y}) \right\} = \underset{\mathbf{y} \in X}{\text{argmin}} \|\mathbf{y} - \mathbf{z}\|^2$$

Convergence in $\mathcal{O}(1/\varepsilon)$ steps, and applications

Same convergence as vanilla case for smooth functions, but now for any h .

Cost: gradient step, plus computing the proximal mapping

Examples:

- ▶ ℓ_1 -norm, $g = \|\cdot\|_1$
 $\text{prox}_{h,\gamma}(\mathbf{z})$ is **soft thresholding operator**, cost $\mathcal{O}(d \log d)$
- ▶ Matrix nuclear norm, $g = \|\cdot\|_*$
 $\text{prox}_{h,\gamma}(\mathbf{Z})$ is **singular value thresholding operator**, costs same as SVD

Chapter 4

Subgradient Descent

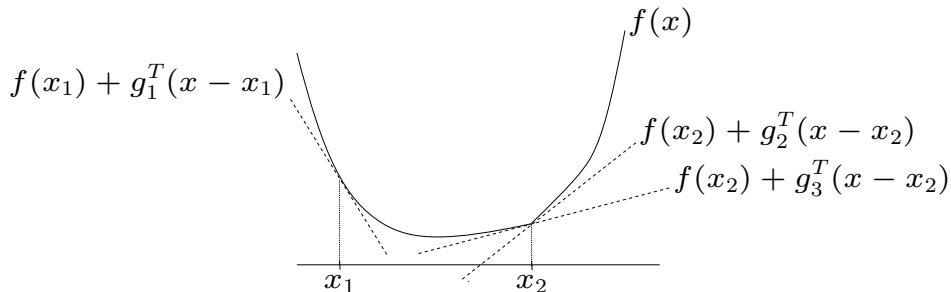
Subgradients

What if f is not differentiable?

Definition

$\mathbf{g} \in \mathbb{R}^d$ is a **subgradient** of f at \mathbf{x} if

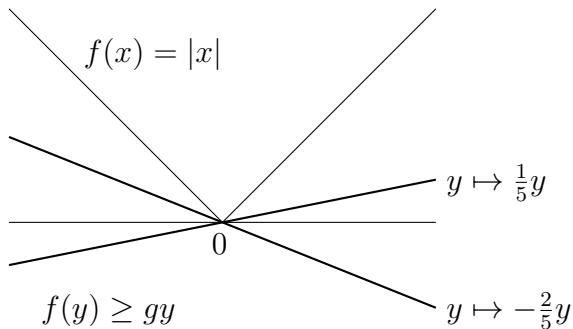
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{y} - \mathbf{x}) \quad \text{for all } \mathbf{y} \in \text{dom}(f)$$



$\partial f(\mathbf{x}) \subseteq \mathbb{R}^d$ is the **subdifferential**, the set of subgradients of f at \mathbf{x} .

Subgradients II

Example:



Subgradient condition at $x = 0$: $f(y) \geq f(0) + g(y - 0) = gy$.

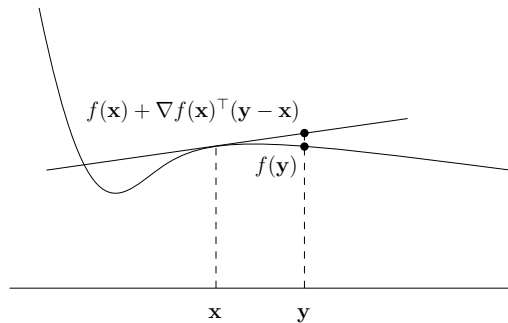
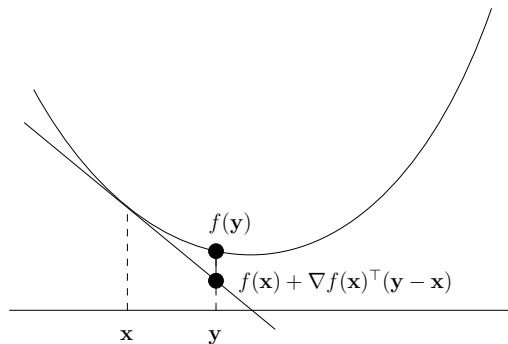
$$\partial f(0) = [-1, 1]$$

Subgradients III

Lemma (Exercise 23)

If $f : \text{dom}(f) \rightarrow \mathbb{R}$ is differentiable at $\mathbf{x} \in \text{dom}(f)$, then $\partial f(\mathbf{x}) \subseteq \{\nabla f(\mathbf{x})\}$.

Either exactly one subgradient $\nabla f(\mathbf{x})$or no subgradient at all.

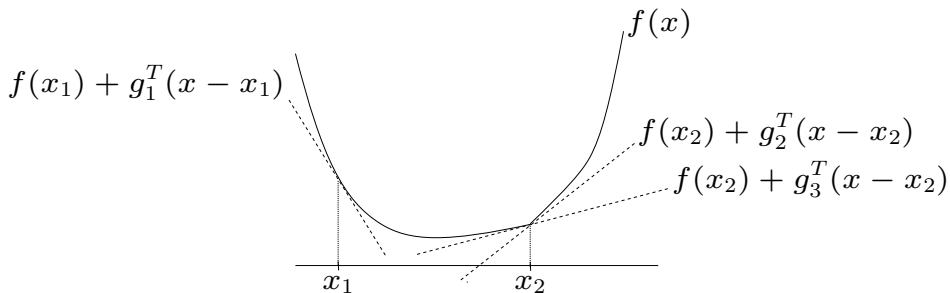


Subgradient characterization of convexity

“convex = subgradients everywhere”

Lemma (Exercise 24)

A function $f : \mathbf{dom}(f) \rightarrow \mathbb{R}$ is convex if and only if $\mathbf{dom}(f)$ is convex and $\partial f(\mathbf{x}) \neq \emptyset$ for all $\mathbf{x} \in \mathbf{dom}(f)$.



Convex and Lipschitz = bounded subgradients

Lemma (Exercise 25)

Let $f : \mathbf{dom}(f) \rightarrow \mathbb{R}$ be convex, $\mathbf{dom}(f)$ open, $B \in \mathbb{R}_+$. Then the following two statements are equivalent.

- (i) $\|\mathbf{g}\| \leq B$ for all $\mathbf{x} \in \mathbf{dom}(f)$ and all $\mathbf{g} \in \partial f(\mathbf{x})$.
- (ii) $|f(\mathbf{x}) - f(\mathbf{y})| \leq B\|\mathbf{x} - \mathbf{y}\|$ for all $\mathbf{x}, \mathbf{y} \in \mathbf{dom}(f)$.

Subgradient optimality condition

Lemma

Suppose that $f : \mathbf{dom}(f) \rightarrow \mathbb{R}$ and $\mathbf{x} \in \mathbf{dom}(f)$. If $\mathbf{0} \in \partial f(\mathbf{x})$, then \mathbf{x} is a global minimum.

Proof.

By definition of subgradients, $\mathbf{g} = \mathbf{0} \in \partial f(\mathbf{x})$ gives

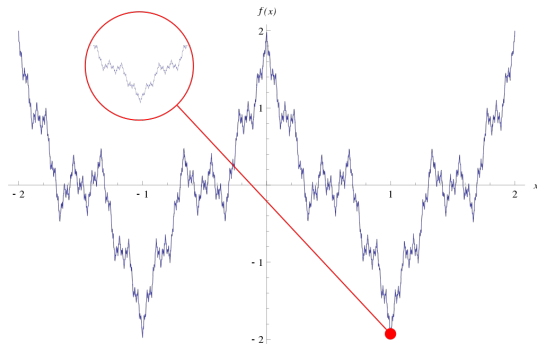
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^\top(\mathbf{y} - \mathbf{x}) = f(\mathbf{x})$$

for all $\mathbf{y} \in \mathbf{dom}(f)$, so \mathbf{x} is a global minimum. □

Differentiability of convex functions

How “wild” can a non-differentiable convex function be?

Weierstrass function: a function that is continuous **everywhere** but differentiable **nowhere**



<https://commons.wikimedia.org/wiki/File:WeierstrassFunction.svg>

Differentiability of convex functions

Theorem ([Roc97, Theorem 25.5])

A *convex* function $f : \mathbf{dom}(f) \rightarrow \mathbb{R}$ is differentiable *almost everywhere*.

In other words:

- ▶ Set of points where f is non-differentiable has measure 0 (no volume).
- ▶ For all $\mathbf{x} \in \mathbf{dom}(f)$ and all $\varepsilon > 0$, there is a point \mathbf{x}' such that $\|\mathbf{x} - \mathbf{x}'\| < \varepsilon$ and f is differentiable at \mathbf{x}' .

The subgradient descent algorithm

Subgradient descent: choose $\mathbf{x}_0 \in \mathbb{R}^d$.

$$\text{Let } \mathbf{g}_t \in \partial f(\mathbf{x}_t)$$

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma_t \mathbf{g}_t$$

for **times** $t = 0, 1, \dots$, and **stepsizes** $\gamma_t \geq 0$.

Stepsize can vary with time!

This is possible in (projected) gradient descent as well.

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

Theorem

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

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

subgradient descent yields

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

Proof is identical to the one of Theorem 2.1, except...

- ▶ In vanilla analysis, now use $\mathbf{g}_t \in \partial f(\mathbf{x}_t)$ instead of $\mathbf{g}_t = \nabla f(\mathbf{x}_t)$.
- ▶ Inequality $f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*)$ now follows from subgradient property instead of first-order characterization of convexity.

Bibliography



R. Tyrrell Rockafellar.

Convex Analysis.

Princeton Landmarks in Mathematics. Princeton University Press, 1997.