# Subgradient method

Axel Böhm

September 23, 2021

1 Subgradient theory

2 Convergence subgradient

3 Smooth case

## Smooth vs. nonsmooth

$$\min_x f(x)$$

### f is *smooth* and convex

GD: $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$

$$f(x_k) - f^* = \mathcal{O}(\frac{1}{k})$$

if the stepsize fulfills
$\alpha_k \leq 1/L$.

### nonsmooth but convex: **subgradient method**

$$\left| \begin{array}{l} \text{pick } g_k \in \partial f(x_k) \\ x_{k+1} = x_k - \alpha_k g_k \end{array} \right.$$

$$f(x_k) - f^* = \mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$$

if stepsize $\alpha_k \approx 1/\sqrt{k}$.

**Subgradient theory**
○●○○○○○

Convergence subgradient
○○○○○○○○○○

Smooth case
○○○○○○○

# Subgradients
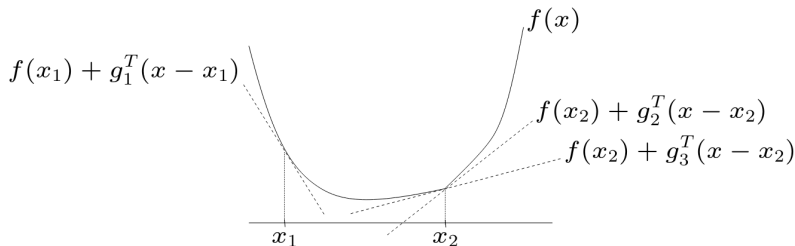
What if $f$ is not differentiable?

### Definition

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

$$f(y) \geq f(x) + g^T(y - x)$$

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

**Subgradient theory**
0●00000

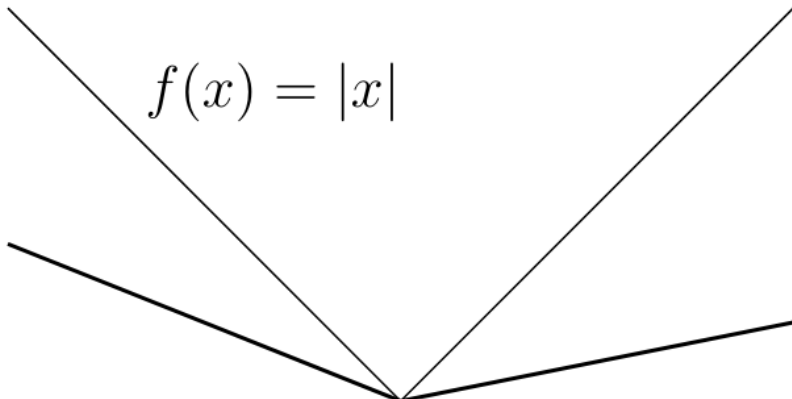Convergence subgradient
0000000000

Smooth case
0000000

# Subgradients II

### Definition

The subdifferential $\partial f(x)$ is the *set of all subgradients* of $f$ at $x$.
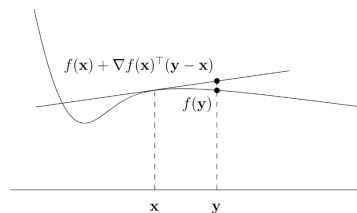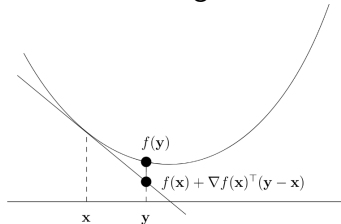
Example



$$f(x) = |x|$$

## Subgradients III

### Lemma

*If $f$ is differentiable at $x$ then $\partial f(x) \subset \{\nabla f(x)\}$*

So either one subgradient or none.

# Subgradient characterization of convexity

### Lemma

*A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if and only if $\partial f(x)$ is not empty for all $x$.*

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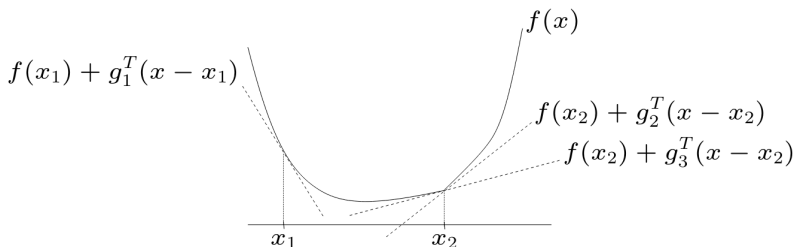


Figure: Subgradients at every point.

## Lipschitz = bounded subgradients

### Definition

We call $f$ $L$-Lipschitz (continuous) if

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

### Lemma

Let $f$ be convex. Then the following two are equivalent.

1. All **subgradients are uniformly bounded**.

$$\|g\| \leq L \quad \forall x, \forall g \in \partial f(x)$$

2. $f$ is $L$-Lipschitz

# Subgradient optimality condition

### Lemma

*Let $0 \in \partial f(\bar{x})$, then $\bar{x}$ is a global minimum.*

### Proof.

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

$$f(y) \geq f(\bar{x}) + g^T(y - \bar{x}) = f(\bar{x}).$$

$\square$

## Convergence statement

### Theorem

*f is convex, subgradients are bounded $\|g(x)\| \leq G$ for all $g(x) \in \partial f(x)$. Then,*

$$f(\bar{x}_k) - f^* \leq \frac{\|x_1 - x^*\|^2 G}{\sqrt{k}}$$

*for the averaged iterates $\bar{x}_k = \frac{\sum_{i=1}^k \alpha_i x_i}{\sum_{i=1}^k \alpha_i}$*

- Also holds for the "best" iterate.
- **Dimension independent!** (no $d$)

## Proof

$$\|x_{k+1} - x^*\|^2 \leq \|x_k - \alpha_k g_k - x^*\|^2$$
$$= \|x_k - x^*\|^2 + 2\alpha_k \langle g_k, x^* - x_k \rangle + \alpha^2 \|g_k\|^2.$$

Using the subgradient ineq. $\langle g_k, x^* - x_k \rangle \leq f(x^*) - f(x_k)$ we deduce

$$\|x_{k+1} - x^*\|^2 \leq \|x_k - x^*\|^2 + 2\alpha_k(f(x^*) - f(x_k)) + \alpha_k^2 \|g_k\|^2.$$

Summing up (telescoping) yields

$$2\sum_{i=1}^{k} \alpha_i(f(x_i) - f(x^*)) + \|x_{k+1} - x^*\|^2 \leq \|x_1 - x^*\|^2 + \sum_{i=1}^{k} \alpha_i^2 \|g_k\|^2.$$

Via the *bounded subgradient* assumption

$$2\sum_{i=1}^{k} \alpha_i(f(x_i) - f(x^*)) + \|x_{k+1} - x^*\|^2 \leq \|x_1 - x^*\|^2 + \sum_{i=1}^{k} \alpha_i^2 G^2.$$

# Proof [contd]

Using Jensens inequality

$$\sum_i \lambda_i f(x_i) \geq \sum_i f\left(\frac{\sum_i \lambda_i x_i}{\sum_i \lambda_i}\right)$$

we obtain

$$2\sum_{i=1}^{k}(f(\bar{x}_k) - f(x^*)) + \|x_{k+1} - x^*\|^2 \leq \|x_1 - x^*\|^2 + \sum_{i=1}^{k}\alpha_i^2 G^2.$$

## How to choose the stepsize?

$$f(\bar{x}_k) - f^* \leq \frac{\|x_1 - x^*\|^2 + G^2 \sum_{i=1}^k \alpha_i^2}{2 \sum_{i=1}^k \alpha_i}$$

Clearly $\alpha_i = \ell_2 \; \ell_1$ leads convergence, for example $1/i$. However, $\alpha_i = \mathcal{O}(1/\sqrt{i})$ gives

$$\sum \alpha_i = (\frac{1}{\sqrt{1}} + \frac{1}{\sqrt{2}} + \cdots + \frac{1}{\sqrt{k}}) > \sqrt{k} \sum \alpha_i^2 = (\frac{1}{1} + \frac{1}{2} + \cdots + \frac{1}{k}) \approx \log(k)$$

$$f(\bar{x}_k) - f^* \leq \frac{\|x_1 - x^*\|^2 + G^2 \log(k)}{2\sqrt{k}}$$

gives the rate

$$\mathcal{O}\left(\frac{\log(k)}{k}\right) =: \tilde{\mathcal{O}}\left(\frac{1}{k}\right)$$

## Complexity

For convex Lipschitz functions we require $\mathcal{O}(\epsilon^{-2})$ iterations. For $D := \|x_1 - x^*\|$

$$f(\bar{x}_k) - f^* \leq \frac{DG}{\sqrt{k}}$$

**Q**: How many iterations to get

$$f(\bar{x}_k) - f^* \leq \epsilon?$$

**A**: We get this if

$$\frac{DG}{\sqrt{k}} \leq \epsilon$$

Equivalently

$$k \geq \frac{D^2 G^2}{\epsilon^2}.$$

## Projected subgradient method

$$\text{(constrained setting)} \quad \min_x f(x)$$

---

**Algorithm** Projected subgradient method

---

1: **for** $k = 1, 2, \ldots$ **do**
2:     Pick $g_k \in \partial f(x_k)$
3:     $x_{k+1} = P_C(x_k - \alpha_k g_k)$

---

By using the fact that the projection is a contraction

$$\|P_C(x) - P_C(y)\| \leq \|x - y\|$$

## Projected subgradient method II

### Proof.

We can deduce the exact same inequality as before

$$
\begin{aligned}
\|x_{k+1} - x^*\|^2 &= \|P_C(x_k - \alpha_k g_k) - x^*\|^2 \\
&\leq \|x_k - \alpha_k g_k - x^*\|^2 \\
&= \|x_k - x^*\|^2 + 2\alpha_k \langle g_k, x^* - x_k \rangle + \alpha^2 \|g_k\|^2 \\
&\leq \|x_k - x^*\|^2 + 2\alpha_k(f^* - f(x_k)) + \alpha^2 \|g_k\|^2.
\end{aligned}
$$

□

## Polyak stepsize

Let's revisit the convergence proof of the subgradient method

$$
\begin{aligned}
\|x_{k+1} - x^*\|^2 &\leq \|x_k - \alpha_k g_k - x^*\|^2 \\
&= \|x_k - x^*\|^2 + 2\alpha_k \langle g_k, x^* - x_k \rangle + \alpha^2 \|g_k\|^2 \\
&\leq \|x_k - x^*\|^2 + 2\alpha_k (f^* - f(x_k)) + \alpha^2 \|g_k\|^2.
\end{aligned}
$$

Can we pick $\alpha_k$ such that the RHS is minimized?

$$
\min_\alpha \alpha^2 \|g_k\|^2 + 2\alpha_k (f^* - f(x_k))
$$

gives

$$
\alpha^* = \frac{f(x_k) - f^*}{\|g_k\|^2}
$$

$$
\|x_{k+1} - x^*\|^2 = \|x_k - x^*\|^2 - \left( \frac{f(x_k) - f^*}{\|g_k\|} \right)^2
$$

## Polyak stepsize [contd]

- Requires us to know the optimal objective function value
- can be the case in certain setting: separable data, feasibility problems
- modern deep learning interpolation setting

## Polyak stepsize [contd]

- Requires us to know the optimal objective function value
- can be the case in certain setting: separable data, feasibility problems
- modern deep learning interpolation setting

Figure: Interpolation / overparametrization regime

## Can we do better?

If $f$ is in addition *strongly convex* the rate improves to

$$f(\bar{x}_k) - f(x^*) \leq \frac{L\|x_1 - x^*\|^2}{\mu T}$$

by choosing the stepsize $\alpha_k \approx \frac{1}{T}$.

## Can we do better if the function is smooth?

### Definition

We call a function *L-smooth* if

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2.$$
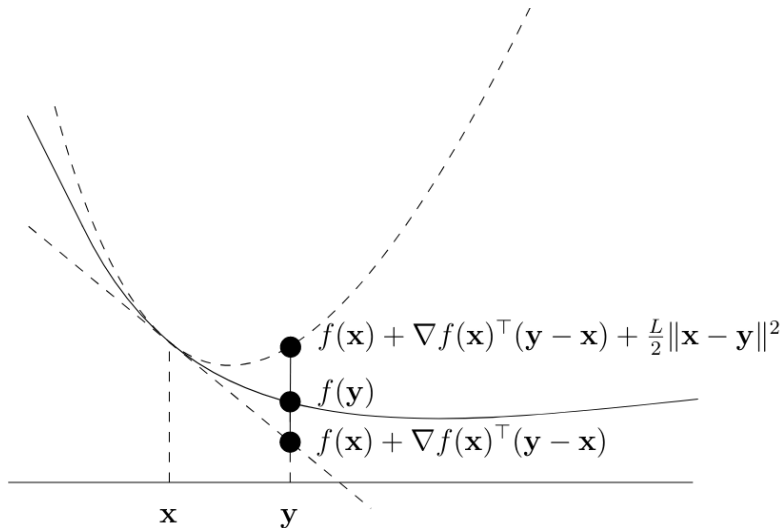
*Can be upper bounded by a quadratic.*

### Lemma

*Is implied by L-Lipschitz gradients:*

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

Note: Definition does not require convexity.

## Smoothness

If $f$ is convex we get upper and lower bound:



$$f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2}\|\mathbf{x} - \mathbf{y}\|^2$$

$$f(\mathbf{y})$$

$$f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x})$$

$\mathbf{x}$          $\mathbf{y}$

- Bounded (sub)gradients $\Leftrightarrow$ Lipschitz continuity of $f$
- Smoothness $\Leftrightarrow$ Lipschitz continuity of $\nabla f$ (if convex)

### Lemma

*Let $f$ be convex and differentiable, then the following are equivalent*

1. *$f$ is smooth with parameter $L$*
2. *$\nabla f$ is $L$-Lipschitz*

## Sufficient decrease

### Lemma

*If $f$ is L-smooth with stepsize $\alpha = 1/L$, then gradient descent satisfies*

$$f(x_{k+1}) \leq f(x_k) - \frac{1}{2L}\|\nabla f(x_k)\|^2$$

### Proof.

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

$\square$

Subgradient theory
ooooooo

Convergence subgradient
oooooooooo

Smooth case
oooo●oo

## Smooth convex functions

### Theorem

Let $f : \mathbb{R}^d \to \mathbb{R}$ be convex and L-smooth and the stepsize $\alpha = 1/L$, then gradient descent yields

$$f(x_k) - f^* \leq \frac{L}{2k}\|x_1 - x^*\|^2$$

- holds for last iterate
- independet of dimension $d$

## Complexity of gradient method

Denote $D^2 := \|x_1 - x^*\|^2$

$$\text{iteration} k \geq \frac{D^2 L}{2\epsilon} \Rightarrow \text{error} \leq \frac{LD^2}{2k} \leq \epsilon$$

Given error $\epsilon = 0.01$ results in

- $50 \cdot D^2 L$ iterations for *smooth* case
- $10000 \cdot D^2 G^2$ for nonsmooth but Lipschitz

What if we don't know $L$?

# Proof of $\mathcal{O}(\epsilon^{-1})$ for smooth functions

### Proof.

Subgradient analysis gave us

$$2\alpha \sum_{i=1}^{k} (f(x_i) - f(x^*)) + \|x_{k+1} - x^*\|^2 \leq \|x_1 - x^*\|^2 + \alpha^2 \sum_{i=1}^{k} \|g_k\|^2,$$

see (**??**) with $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ □