

# Convex Optimization

(EE227A: UC Berkeley)

## Lecture 12

(Subgradient methods)

28 Feb, 2013



Stefanie J. (for Suvrit Sra)

# Announcements

---

- ▶ Midterm: final date is 19th March

# ALGORITHMS

# Unconstrained convex problem

---

$$\min_x f(x)$$

How to optimize?

# Unconstrained convex problem

---

$$\min_x f(x)$$

How to optimize?

- 1 Start with some guess  $x^0$ ; set  $k = 0$

# Unconstrained convex problem

---

$$\min_x f(x)$$

How to optimize?

- 1 Start with some guess  $x^0$ ; set  $k = 0$
- 2 If  $0 \in \partial f(x^k)$ , **stop**; output  $x^k$

# Unconstrained convex problem

---

$$\min_x f(x)$$

How to optimize?

- 1 Start with some guess  $x^0$ ; set  $k = 0$
- 2 If  $0 \in \partial f(x^k)$ , **stop**; output  $x^k$
- 3 Otherwise, generate next guess  $x^{k+1}$

# Unconstrained convex problem

---

$$\min_x f(x)$$

How to optimize?

- 1 Start with some guess  $x^0$ ; set  $k = 0$
- 2 If  $0 \in \partial f(x^k)$ , **stop**; output  $x^k$
- 3 Otherwise, generate next guess  $x^{k+1}$
- 4 Repeat above procedure



# Unconstrained convex problem

---

$$\min_x f(x)$$

How to optimize?

- 1 Start with some guess  $x^0$ ; set  $k = 0$
- 2 If  $0 \in \partial f(x^k)$ , **stop**; output  $x^k$
- 3 Otherwise, generate next guess  $x^{k+1}$
- 4 Repeat above procedure
  - ▶ In reality: we stop in finite time
  - ▶ Only solve problem approximately
  - ▶  $f(x^k) \leq f(x^*) + \varepsilon$
  - ▶ **shorthand**  $f^k \leq f^* + \varepsilon$

# Subgradient method

---

$$x^{k+1} = x^k - \alpha_k g^k$$

where  $g^k \in \partial f(x^k)$  is **any** subgradient

# Subgradient method

---

$$x^{k+1} = x^k - \alpha_k g^k$$

where  $g^k \in \partial f(x^k)$  is **any** subgradient

**Stepsize**  $\alpha_k > 0$  **must be chosen**

# Subgradient method

---

$$x^{k+1} = x^k - \alpha_k g^k$$

where  $g^k \in \partial f(x^k)$  is **any** subgradient

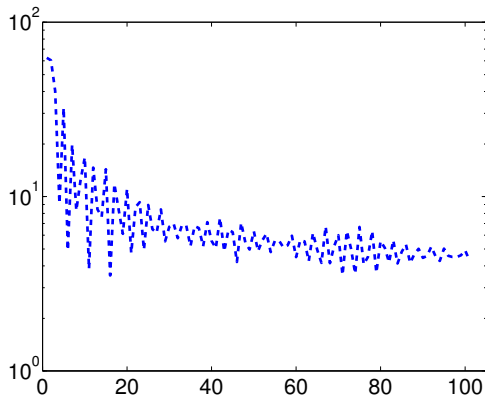
**Stepsize**  $\alpha_k > 0$  **must be chosen**

- ▶ Method generates sequence  $\{x^k\}_{k \geq 0}$
- ▶ Does this sequence converge to an optimal solution  $x^*$ ?
- ▶ If yes, then how fast?
- ▶ What if have constraints:  $x \in \mathcal{C}$ ?

# Example

---

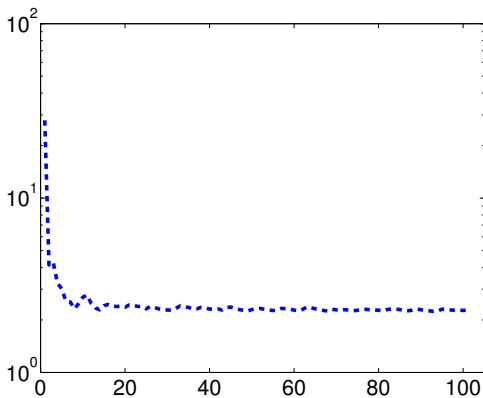
$$\min \quad \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1$$
$$x^{k+1} = x^k - \alpha_k (A^T (Ax^k - b) + \lambda \operatorname{sgn}(x^k))$$



## Example

---

$$\min \quad \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1$$
$$x^{k+1} = x^k - \alpha_k (A^T (Ax^k - b) + \lambda \operatorname{sgn}(x^k))$$



(More careful implementation)

## Subgradient method – stepsizes

---

- ▶ **Constant** Set  $\alpha_k = \alpha > 0$ , for  $k \geq 0$
- ▶ **Scaled constant**  $\alpha_k = \alpha / \|g^k\|_2$  ( $\|x^{k+1} - x^k\|_2 = \alpha$ )

## Subgradient method – stepsizes

---

- ▶ **Constant** Set  $\alpha_k = \alpha > 0$ , for  $k \geq 0$
- ▶ **Scaled constant**  $\alpha_k = \alpha / \|g^k\|_2$  ( $\|x^{k+1} - x^k\|_2 = \alpha$ )
- ▶ **Square summable but not summable**

$$\sum_k \alpha_k^2 < \infty, \quad \sum_k \alpha_k = \infty$$

- ▶ **Diminishing scalar**

$$\lim_k \alpha_k = 0, \quad \sum_k \alpha_k = \infty$$

- ▶ **Adaptive stepsizes** (not covered)

Not a descent method!

Work with best  $f^k$  so far:  $f_{\min}^k := \min_{0 \leq i \leq k} f^i$



# Convergence analysis

---

## Assumptions

- ▶ Min is attained:  $f^* := \inf_x f(x) > -\infty$ , with  $f(x^*) = f^*$

# Convergence analysis

---

## Assumptions

- ▶ Min is attained:  $f^* := \inf_x f(x) > -\infty$ , with  $f(x^*) = f^*$
- ▶ Bounded subgradients:  $\|g\|_2 \leq G$  for all  $g \in \partial f$   
( $f(x) - f(y) = \langle g_\xi, x - y \rangle$ ; use Cauchy-Schwarz or Hölder)

# Convergence analysis

---

## Assumptions

- ▶ Min is attained:  $f^* := \inf_x f(x) > -\infty$ , with  $f(x^*) = f^*$
- ▶ Bounded subgradients:  $\|g\|_2 \leq G$  for all  $g \in \partial f$   
( $f(x) - f(y) = \langle g_\xi, x - y \rangle$ ; use Cauchy-Schwarz or Hölder)
- ▶ Bounded domain:  $\|x^0 - x^*\|_2 \leq R$

# Convergence analysis

---

## Assumptions

- ▶ Min is attained:  $f^* := \inf_x f(x) > -\infty$ , with  $f(x^*) = f^*$
- ▶ Bounded subgradients:  $\|g\|_2 \leq G$  for all  $g \in \partial f$   
( $f(x) - f(y) = \langle g_\xi, x - y \rangle$ ; use Cauchy-Schwarz or Hölder)
- ▶ Bounded domain:  $\|x^0 - x^*\|_2 \leq R$

Convergence results for:  $f_{\min}^k := \min_{0 \leq i \leq k} f^i$

# Subgradient method – convergence

---

Lyapunov function: Distance to  $x^*$ , not function values

# Subgradient method – convergence

---

**Lyapunov function:** Distance to  $x^*$ , not function values

$$\|x^{k+1} - x^*\|_2^2 = \|x^k - \alpha_k g^k - x^*\|_2^2$$

# Subgradient method – convergence

---

**Lyapunov function:** Distance to  $x^*$ , not function values

$$\begin{aligned}\|x^{k+1} - x^*\|_2^2 &= \|x^k - \alpha_k g^k - x^*\|_2^2 \\ &= \|x^k - x^*\|_2^2 + \alpha_k^2 \|g^k\|_2^2 - 2\langle \alpha_k g^k, x^k - x^* \rangle\end{aligned}$$

# Subgradient method – convergence

---

**Lyapunov function:** Distance to  $x^*$ , not function values

$$\begin{aligned}\|x^{k+1} - x^*\|_2^2 &= \|x^k - \alpha_k g^k - x^*\|_2^2 \\ &= \|x^k - x^*\|_2^2 + \alpha_k^2 \|g^k\|_2^2 - 2\langle \alpha_k g^k, x^k - x^* \rangle \\ &\leq \|x^k - x^*\|_2^2 + \alpha_k^2 \|g^k\|_2^2 - 2\alpha_k (f(x^k) - f^*),\end{aligned}$$

since  $f^* = f(x^*) \geq f(x^k) + \langle g^k, x^* - x^k \rangle$



# Subgradient method – convergence

---

**Lyapunov function:** Distance to  $x^*$ , not function values

$$\begin{aligned}\|x^{k+1} - x^*\|_2^2 &= \|x^k - \alpha_k g^k - x^*\|_2^2 \\ &= \|x^k - x^*\|_2^2 + \alpha_k^2 \|g^k\|_2^2 - 2\langle \alpha_k g^k, x^k - x^* \rangle \\ &\leq \|x^k - x^*\|_2^2 + \alpha_k^2 \|g^k\|_2^2 - 2\alpha_k (f(x^k) - f^*),\end{aligned}$$

since  $f^* = f(x^*) \geq f(x^k) + \langle g^k, x^* - x^k \rangle$

Apply same argument to  $\|x^k - x^*\|_2^2$  recursively

# Subgradient method – convergence

---

**Lyapunov function:** Distance to  $x^*$ , not function values

$$\begin{aligned}\|x^{k+1} - x^*\|_2^2 &= \|x^k - \alpha_k g^k - x^*\|_2^2 \\ &= \|x^k - x^*\|_2^2 + \alpha_k^2 \|g^k\|_2^2 - 2\langle \alpha_k g^k, x^k - x^* \rangle \\ &\leq \|x^k - x^*\|_2^2 + \alpha_k^2 \|g^k\|_2^2 - 2\alpha_k (f(x^k) - f^*),\end{aligned}$$

since  $f^* = f(x^*) \geq f(x^k) + \langle g^k, x^* - x^k \rangle$

Apply same argument to  $\|x^k - x^*\|_2^2$  recursively

$$\|x^{k+1} - x^*\|_2^2 \leq \|x^0 - x^*\|_2^2 + \sum_{t=1}^k \alpha_t^2 \|g^t\|_2^2 - 2 \sum_{t=1}^k \alpha_t (f^t - f^*).$$

# Subgradient method – convergence

---

**Lyapunov function:** Distance to  $x^*$ , not function values

$$\begin{aligned}\|x^{k+1} - x^*\|_2^2 &= \|x^k - \alpha_k g^k - x^*\|_2^2 \\ &= \|x^k - x^*\|_2^2 + \alpha_k^2 \|g^k\|_2^2 - 2\langle \alpha_k g^k, x^k - x^* \rangle \\ &\leq \|x^k - x^*\|_2^2 + \alpha_k^2 \|g^k\|_2^2 - 2\alpha_k (f(x^k) - f^*),\end{aligned}$$

since  $f^* = f(x^*) \geq f(x^k) + \langle g^k, x^* - x^k \rangle$

Apply same argument to  $\|x^k - x^*\|_2^2$  recursively

$$\|x^{k+1} - x^*\|_2^2 \leq \|x^0 - x^*\|_2^2 + \sum_{t=1}^k \alpha_t^2 \|g^t\|_2^2 - 2 \sum_{t=1}^k \alpha_t (f^t - f^*).$$

Now use our convenient assumptions!

## Subgradient method – convergence

---

$$\|x^{k+1} - x^*\|_2^2 \leq R^2 + G^2 \sum_{t=1}^k \alpha_t^2 - 2 \sum_{t=1}^k \alpha_t (f^t - f^*).$$

► To get a bound on the last term, simply notice (for  $t \leq k$ )

$$f^t \geq f_{\min}^t \geq f_{\min}^k \quad \text{since} \quad f_{\min}^t := \min_{0 \leq i \leq t} f(x^i)$$

## Subgradient method – convergence

---

$$\|x^{k+1} - x^*\|_2^2 \leq R^2 + G^2 \sum_{t=1}^k \alpha_t^2 - 2 \sum_{t=1}^k \alpha_t (f^t - f^*).$$

- To get a bound on the last term, simply notice (for  $t \leq k$ )

$$f^t \geq f_{\min}^t \geq f_{\min}^k \quad \text{since} \quad f_{\min}^t := \min_{0 \leq i \leq t} f(x^i)$$

- Plugging this in yields the bound

$$2 \sum_{t=1}^k \alpha_t (f^t - f^*) \geq 2(f_{\min}^k - f^*) \sum_{t=1}^k \alpha_t.$$

## Subgradient method – convergence

---

$$\|x^{k+1} - x^*\|_2^2 \leq R^2 + G^2 \sum_{t=1}^k \alpha_t^2 - 2 \sum_{t=1}^k \alpha_t (f^t - f^*).$$

► To get a bound on the last term, simply notice (for  $t \leq k$ )

$$f^t \geq f_{\min}^t \geq f_{\min}^k \quad \text{since} \quad f_{\min}^t := \min_{0 \leq i \leq t} f(x^i)$$

► Plugging this in yields the bound

$$2 \sum_{t=1}^k \alpha_t (f^t - f^*) \geq 2(f_{\min}^k - f^*) \sum_{t=1}^k \alpha_t.$$

► So that we finally have

$$0 \leq \|x^{k+1} - x^*\|_2^2 \leq R^2 + G^2 \sum_{t=1}^k \alpha_t^2 - 2(f_{\min}^k - f^*) \sum_{t=1}^k \alpha_t$$

## Subgradient method – convergence

$$\|x^{k+1} - x^*\|_2^2 \leq R^2 + G^2 \sum_{t=1}^k \alpha_t^2 - 2 \sum_{t=1}^k \alpha_t (f^t - f^*).$$

► To get a bound on the last term, simply notice (for  $t \leq k$ )

$$f^t \geq f_{\min}^t \geq f_{\min}^k \quad \text{since} \quad f_{\min}^k := \min_{0 \leq i \leq t} f(x^i)$$

► Plugging this in yields the bound

$$2 \sum_{t=1}^k \alpha_t (f^t - f^*) \geq 2(f_{\min}^k - f^*) \sum_{t=1}^k \alpha_t.$$

► So that we finally have

$$0 \leq \|x^{k+1} - x^*\|_2^2 \leq R^2 + G^2 \sum_{t=1}^k \alpha_t^2 - 2(f_{\min}^k - f^*) \sum_{t=1}^k \alpha_t$$

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t}$$

## Subgradient method – convergence

---

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t}$$

**Exercise:** Analyze  $\lim_{k \rightarrow \infty} f_{\min}^k - f^*$  for the different choices of stepsize that we mentioned.



## Subgradient method – convergence

---

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t}$$

**Exercise:** Analyze  $\lim_{k \rightarrow \infty} f_{\min}^k - f^*$  for the different choices of stepsize that we mentioned.

**Constant step:**  $\alpha_k = \alpha$ ; We obtain

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 k \alpha^2}{2k\alpha}$$

## Subgradient method – convergence

---

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t}$$

**Exercise:** Analyze  $\lim_{k \rightarrow \infty} f_{\min}^k - f^*$  for the different choices of stepsize that we mentioned.

**Constant step:**  $\alpha_k = \alpha$ ; We obtain

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 k \alpha^2}{2k\alpha} \rightarrow \frac{G^2 \alpha}{2} \quad \text{as } k \rightarrow \infty.$$

# Subgradient method – convergence

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t}$$

**Exercise:** Analyze  $\lim_{k \rightarrow \infty} f_{\min}^k - f^*$  for the different choices of stepsize that we mentioned.

**Constant step:**  $\alpha_k = \alpha$ ; We obtain

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 k \alpha^2}{2k\alpha} \rightarrow \frac{G^2 \alpha}{2} \quad \text{as } k \rightarrow \infty.$$

**Square summable, not summable:**  $\sum_k \alpha_k^2 < \infty$ ,  $\sum_k \alpha_k = \infty$

# Subgradient method – convergence

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t}$$

**Exercise:** Analyze  $\lim_{k \rightarrow \infty} f_{\min}^k - f^*$  for the different choices of stepsize that we mentioned.

**Constant step:**  $\alpha_k = \alpha$ ; We obtain

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 k \alpha^2}{2k\alpha} \rightarrow \frac{G^2 \alpha}{2} \quad \text{as } k \rightarrow \infty.$$

**Square summable, not summable:**  $\sum_k \alpha_k^2 < \infty$ ,  $\sum_k \alpha_k = \infty$   
As  $k \rightarrow \infty$ , numerator  $< \infty$  but denominator  $\rightarrow \infty$ ; so  $f_{\min}^k \rightarrow f^*$

# Subgradient method – convergence

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t}$$

**Exercise:** Analyze  $\lim_{k \rightarrow \infty} f_{\min}^k - f^*$  for the different choices of stepsize that we mentioned.

**Constant step:**  $\alpha_k = \alpha$ ; We obtain

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 k \alpha^2}{2k\alpha} \rightarrow \frac{G^2 \alpha}{2} \quad \text{as } k \rightarrow \infty.$$

**Square summable, not summable:**  $\sum_k \alpha_k^2 < \infty$ ,  $\sum_k \alpha_k = \infty$   
As  $k \rightarrow \infty$ , numerator  $< \infty$  but denominator  $\rightarrow \infty$ ; so  $f_{\min}^k \rightarrow f^*$

In practice, fair bit of stepsize tuning needed, e.g.  $\alpha_t = a/(b+t)$

## Subgradient method – convergence

---

- ▶ Suppose we want  $f_{\min}^k - f^* \leq \varepsilon$ , how big should  $k$  be?

# Subgradient method – convergence

---

- ▶ Suppose we want  $f_{\min}^k - f^* \leq \varepsilon$ , how big should  $k$  be?
- ▶ Optimize the bound for  $\alpha_t$ : want

$$f_{\min}^k - f^* \leq \varepsilon$$

## Subgradient method – convergence

---

- ▶ Suppose we want  $f_{\min}^k - f^* \leq \varepsilon$ , how big should  $k$  be?
- ▶ Optimize the bound for  $\alpha_t$ : want

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t} \leq \varepsilon$$



## Subgradient method – convergence

---

- ▶ Suppose we want  $f_{\min}^k - f^* \leq \epsilon$ , how big should  $k$  be?
- ▶ Optimize the bound for  $\alpha_t$ : want

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t} \leq \epsilon$$

- ▶ For fixed  $k$ : best possible stepsize is constant  $\alpha$

$$\frac{R^2 + G^2 k \alpha^2}{2k\alpha} \leq \epsilon \quad \Rightarrow \quad \alpha = \frac{R}{G\sqrt{k}}$$

## Subgradient method – convergence

---

- ▶ Suppose we want  $f_{\min}^k - f^* \leq \epsilon$ , how big should  $k$  be?
- ▶ Optimize the bound for  $\alpha_t$ : want

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t} \leq \epsilon$$

- ▶ For fixed  $k$ : best possible stepsize is constant  $\alpha$

$$\frac{R^2 + G^2 k \alpha^2}{2k\alpha} \leq \epsilon \quad \Rightarrow \quad \alpha = \frac{R}{G\sqrt{k}}$$

- ▶ Then, after  $k$  steps  $f_{\min}^k - f^* \leq RG/\sqrt{k}$ .
- ▶ For accuracy  $\epsilon$ , we need at least  $(RG/\epsilon)^2 = O(1/\epsilon^2)$  steps

## Subgradient method – convergence

---

- ▶ Suppose we want  $f_{\min}^k - f^* \leq \epsilon$ , how big should  $k$  be?
- ▶ Optimize the bound for  $\alpha_t$ : want

$$f_{\min}^k - f^* \leq \frac{R^2 + G^2 \sum_{t=1}^k \alpha_t^2}{2 \sum_{t=1}^k \alpha_t} \leq \epsilon$$

- ▶ For fixed  $k$ : best possible stepsize is constant  $\alpha$

$$\frac{R^2 + G^2 k \alpha^2}{2k\alpha} \leq \epsilon \quad \Rightarrow \quad \alpha = \frac{R}{G\sqrt{k}}$$

- ▶ Then, after  $k$  steps  $f_{\min}^k - f^* \leq RG/\sqrt{k}$ .
- ▶ For accuracy  $\epsilon$ , we need at least  $(RG/\epsilon)^2 = O(1/\epsilon^2)$  steps
- ▶ (quite slow)

# Exercise

---

## Support vector machines

- ▶ Let  $\mathcal{D} := \{(x_i, y_i) \mid x_i \in \mathbb{R}^n, y_i \in \{\pm 1\}\}$
- ▶ We wish to find  $w \in \mathbb{R}^n$  and  $b \in \mathbb{R}$  such that

$$\min_{w,b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \max[0, 1 - y_i(w^T x_i + b)]$$

- ▶ Derive and implement a subgradient method
- ▶ Plot evolution of objective function
- ▶ Experiment with different values of  $C > 0$
- ▶ Plot and keep track of  $f_{\min}^k := \min_{0 \leq t \leq k} f(x^t)$

## Polyak's stepsize

---

- ▶ Assume  $f^*$  is known (or can be estimated). Then use

$$\alpha_t = \frac{f^t - f^*}{\|g^t\|_2^2}$$

# Polyak's stepsize

---

- ▶ Assume  $f^*$  is known (or can be estimated). Then use

$$\alpha_t = \frac{f^t - f^*}{\|g^t\|_2^2}$$

- ▶ Motivation: recall bound

$$\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2 - 2\alpha_t(f^t - f^*) + \alpha_t^2\|g^t\|^2$$

and minimize RHS.

# Polyak's stepsize

---

- ▶ Assume  $f^*$  is known (or can be estimated). Then use

$$\alpha_t = \frac{f^t - f^*}{\|g^t\|_2^2}$$

- ▶ Motivation: recall bound

$$\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2 - 2\alpha_t(f^t - f^*) + \alpha_t^2\|g^t\|^2$$

and minimize RHS.

- ▶ Let's plug in  $\alpha_t$ :

$$\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2 - \frac{(f^t - f^*)^2}{\|g_t\|^2}$$

## Polyak's stepsize

---

$$\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2 - \frac{(f^t - f^*)^2}{\|g_t\|^2}$$



# Polyak's stepsize

---

$$\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2 - \frac{(f^t - f^*)^2}{\|g_t\|^2}$$

► **Observation 1**  $\|x^t - x^*\|$  decreases

# Polyak's stepsize

---

$$\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2 - \frac{(f^t - f^*)^2}{\|g_t\|^2}$$

- ▶ **Observation 1**  $\|x^t - x^*\|$  decreases
- ▶ Recursion:

$$\sum_{t=1}^k \frac{(f^t - f^*)^2}{\|g_t\|^2} \leq \|x^1 - x^*\|^2 \leq R^2$$

# Polyak's stepsize

---

$$\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2 - \frac{(f^t - f^*)^2}{\|g_t\|^2}$$

- ▶ **Observation 1**  $\|x^t - x^*\|$  decreases
- ▶ Recursion:

$$\sum_{t=1}^k \frac{(f^t - f^*)^2}{\|g_t\|^2} \leq \|x^1 - x^*\|^2 \leq R^2$$

- ▶ Now use  $\|g^t\| \leq G$

$$\sum_{t=1}^k (f^t - f^*)^2 \leq R^2 G^2$$

# Polyak's stepsize

---

$$\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2 - \frac{(f^t - f^*)^2}{\|g_t\|^2}$$

- ▶ **Observation 1**  $\|x^t - x^*\|$  decreases
- ▶ Recursion:

$$\sum_{t=1}^k \frac{(f^t - f^*)^2}{\|g_t\|^2} \leq \|x^1 - x^*\|^2 \leq R^2$$

- ▶ Now use  $\|g^t\| \leq G$

$$\sum_{t=1}^k (f^t - f^*)^2 \leq R^2 G^2$$

- ▶ **Observation 2**  $f^t \rightarrow f^*$

# Polyak's stepsize

---

$$\|x^{t+1} - x^*\|^2 \leq \|x^t - x^*\|^2 - \frac{(f^t - f^*)^2}{\|g_t\|^2}$$

- ▶ **Observation 1**  $\|x^t - x^*\|$  decreases
- ▶ Recursion:

$$\sum_{t=1}^k \frac{(f^t - f^*)^2}{\|g_t\|^2} \leq \|x^1 - x^*\|^2 \leq R^2$$

- ▶ Now use  $\|g^t\| \leq G$

$$\sum_{t=1}^k (f^t - f^*)^2 \leq R^2 G^2$$

- ▶ **Observation 2**  $f^t \rightarrow f^*$
- ▶ for accuracy  $\varepsilon$ , need  $k = (RG/\varepsilon)^2$

# Constrained optimization

---

$$\min f(x) \quad \text{s.t.} \quad x \in \mathcal{C}$$

# Constrained optimization

---

$$\min f(x) \quad \text{s.t.} \quad x \in \mathcal{C}$$

- Previously:

$$x^{t+1} = x^t - \alpha_t g^t$$

- This could be infeasible!  
Solution: projection

## Projected subgradient method

---

$$x^{k+1} = P_{\mathcal{C}}(x^k - \alpha_k g^k)$$

where  $g^k \in \partial f(x^k)$  is any subgradient



# Projected subgradient method

---

$$x^{k+1} = P_{\mathcal{C}}(x^k - \alpha_k g^k)$$

where  $g^k \in \partial f(x^k)$  is any subgradient

- **Projection** closest feasible point

$$P_{\mathcal{C}}(x) = \arg \min_{y \in \mathcal{C}} \|x - y\|^2$$

(Assume  $\mathcal{C}$  is closed and convex, then projection is unique)

# Projected subgradient method

---

$$x^{k+1} = P_{\mathcal{C}}(x^k - \alpha_k g^k)$$

where  $g^k \in \partial f(x^k)$  is any subgradient

- ▶ **Projection** closest feasible point

$$P_{\mathcal{C}}(x) = \arg \min_{y \in \mathcal{C}} \|x - y\|^2$$

(Assume  $\mathcal{C}$  is closed and convex, then projection is unique)

- ▶ Great as long as projection is “easy”
- ▶ Same questions as before:
  - Does it converge?
  - For which stepsizes?
  - How fast?

# Convergence

---

## Assumptions

- ▶ Min is attained:  $f^* := \inf_x f(x) > -\infty$ , with  $f(x^*) = f^*$
- ▶ Bounded subgradients:  $\|g\|_2 \leq G$  for all  $g \in \partial f$
- ▶ Bounded domain:  $\|x^0 - x^*\|_2 \leq R$

# Convergence

---

## Assumptions

- ▶ Min is attained:  $f^* := \inf_x f(x) > -\infty$ , with  $f(x^*) = f^*$
- ▶ Bounded subgradients:  $\|g\|_2 \leq G$  for all  $g \in \partial f$
- ▶ Bounded domain:  $\|x^0 - x^*\|_2 \leq R$

## Analysis

- ▶ Let  $z^{t+1} = x^t - \alpha_t g^t$ .
- ▶ Then  $x^{t+1} = P_{\mathcal{C}}(z^{t+1})$ .

# Convergence

---

## Assumptions

- ▶ Min is attained:  $f^* := \inf_x f(x) > -\infty$ , with  $f(x^*) = f^*$
- ▶ Bounded subgradients:  $\|g\|_2 \leq G$  for all  $g \in \partial f$
- ▶ Bounded domain:  $\|x^0 - x^*\|_2 \leq R$

## Analysis

- ▶ Let  $z^{t+1} = x^t - \alpha_t g^t$ .
- ▶ Then  $x^{t+1} = P_{\mathcal{C}}(z^{t+1})$ .
- ▶ Recall analysis of unconstrained method:

$$\begin{aligned}\|z^{t+1} - x^*\|_2^2 &= \|x^t - \alpha_t g^t - x^*\|_2^2 \\ &\leq \|x^t - x^*\|_2^2 + \alpha_t^2 \|g^t\|_2^2 - 2\alpha_t (f(x^t) - f^*) \\ &\dots\end{aligned}$$

- ▶ Need to relate to  $\|x^{t+1} - x^*\|_2^2$ , the rest of the proof is the same as above.

# Projection Theorem

---

Let  $\mathcal{C}$  be nonempty, closed and convex.

- Optimality conditions:  $y^* = P_{\mathcal{C}}(z)$  iff

$$\langle z - y^*, y - y^* \rangle \leq 0 \text{ for all } y \in \mathcal{C}$$

- The projection is nonexpansive:

$$\|P_{\mathcal{C}}(x) - P_{\mathcal{C}}(z)\| \leq \|x - z\| \quad \text{for all } x, z \in \mathbb{R}^n.$$

# Convergence

---

- Use nonexpansiveness of projection:

$$\begin{aligned} & \|x^t - \alpha_t g^t - x^*\|_2^2 \\ & \leq \|x^t - x^*\|_2^2 + \alpha_t^2 \|g^t\|_2^2 - 2\alpha_t (f(x^t) - f^*) \\ & \dots \end{aligned}$$

# Convergence

---

- ▶ Use nonexpansiveness of projection:

$$\begin{aligned}\|x^{t+1} - x^*\|_2^2 &= \|P_C(x^t - \alpha_t g^t) - x^*\|_2^2 \\ &\leq \|x^t - \alpha_t g^t - x^*\|_2^2 \\ &\leq \|x^t - x^*\|_2^2 + \alpha_t^2 \|g^t\|_2^2 - 2\alpha_t (f(x^t) - f^*) \\ &\dots\end{aligned}$$

Same convergence results as in unconstrained case:

- ▶ within neighborhood of optimal for constant step size
- ▶ converges for diminishing non-summable



# Examples

---

$$\begin{aligned} \min \quad & \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 \\ \text{s.t.} \quad & x \in \mathcal{C} \end{aligned}$$

► **Nonnegativity**  $x \geq 0$

$$P_{\mathcal{C}}(z) = [z]_+$$

$$\text{Update step: } x^{k+1} = [x^k - \alpha_k (A^T (Ax^k - b) + \lambda \operatorname{sgn}(x^k))]_+$$

# Examples

---

$$\begin{aligned} \min \quad & \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 \\ \text{s.t.} \quad & x \in \mathcal{C} \end{aligned}$$

- ▶ **Nonnegativity**  $x \geq 0$

$$P_{\mathcal{C}}(z) = [z]_+$$

$$\text{Update step: } x^{k+1} = [x^k - \alpha_k (A^T (Ax^k - b) + \lambda \text{sgn}(x^k))]_+$$

- ▶  **$l_\infty$ -ball**  $\|x\|_\infty \leq 1$

$$\text{Projection: } \min \|x - z\|^2 \text{ s.t. } x \leq 1 \text{ and } x \geq -1$$

# Examples

---

$$\begin{aligned} \min \quad & \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 \\ \text{s.t.} \quad & x \in \mathcal{C} \end{aligned}$$

- ▶ **Nonnegativity**  $x \geq 0$

$$P_{\mathcal{C}}(z) = [z]_+$$

$$\text{Update step: } x^{k+1} = [x^k - \alpha_k (A^T (Ax^k - b) + \lambda \text{sgn}(x^k))]_+$$

- ▶  **$\ell_\infty$ -ball**  $\|x\|_\infty \leq 1$

$$\text{Projection: } \min \|x - z\|^2 \text{ s.t. } x \leq 1 \text{ and } x \geq -1$$

this is separable, so do it coordinate-wise:

$$P_{\mathcal{C}}(z) = y \text{ where } y_i = \text{sgn}(z_i) \min\{|z_i|, 1\}$$

# Examples

---

$$\begin{aligned} \min \quad & \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 \\ \text{s.t.} \quad & x \in \mathcal{C} \end{aligned}$$

- **Nonnegativity**  $x \geq 0$

$$P_{\mathcal{C}}(z) = [z]_+$$

$$\text{Update step: } x^{k+1} = [x^k - \alpha_k (A^T (Ax^k - b) + \lambda \operatorname{sgn}(x^k))]_+$$

- **$\ell_\infty$ -ball**  $\|x\|_\infty \leq 1$

$$\text{Projection: } \min \|x - z\|^2 \text{ s.t. } x \leq 1 \text{ and } x \geq -1$$

this is separable, so do it coordinate-wise:

$$P_{\mathcal{C}}(z) = y \text{ where } y_i = \operatorname{sgn}(z_i) \min\{|z_i|, 1\}$$

Update step:

$$z^{k+1} = x^k - \alpha_k (A^T (Ax^k - b) + \lambda \operatorname{sgn}(x^k))$$

$$x_i^{k+1} = \operatorname{sgn}(z_i^{k+1}) \min\{|z_i^{k+1}|, 1\}$$

## Examples

---

- ▶ **Linear equality constraints**  $Ax = b$  ( $A \in \mathbb{R}^{n \times m}$  has rank  $n$ )

$$\begin{aligned} P_{\mathcal{C}}(x) &= z - A^{\top}(AA^{\top})^{-1}(Az - b) \\ &= (I - A^{\top}(A^{\top}A)^{-1}A)z + A^{\top}(AA^{\top})^{-1}b \end{aligned}$$

## Examples

---

- **Linear equality constraints**  $Ax = b$  ( $A \in \mathbb{R}^{n \times m}$  has rank  $n$ )

$$\begin{aligned}P_{\mathcal{C}}(x) &= z - A^{\top}(AA^{\top})^{-1}(Az - b) \\ &= (I - A^{\top}(A^{\top}A)^{-1}A)z + A^{\top}(AA^{\top})^{-1}b\end{aligned}$$

Update step, using  $Ax^t = b$ :

$$\begin{aligned}x^{t+1} &= P_{\mathcal{C}}(x^t - \alpha_t g^t) \\ &= x^t - \alpha_t(I - A^{\top}(AA^{\top})^{-1}A)g^t\end{aligned}$$

# Examples

---

- ▶ **Linear equality constraints**  $Ax = b$  ( $A \in \mathbb{R}^{n \times m}$  has rank  $n$ )

$$\begin{aligned}P_{\mathcal{C}}(x) &= z - A^{\top}(AA^{\top})^{-1}(Az - b) \\ &= (I - A^{\top}(A^{\top}A)^{-1}A)z + A^{\top}(AA^{\top})^{-1}b\end{aligned}$$

Update step, using  $Ax^t = b$ :

$$\begin{aligned}x^{t+1} &= P_{\mathcal{C}}(x^t - \alpha_t g^t) \\ &= x^t - \alpha_t (I - A^{\top}(AA^{\top})^{-1}A)g^t\end{aligned}$$

- ▶ **Simplex**  $x^{\top}1 = 1$  and  $x \geq 0$   
more complex but doable, similarly  $\ell_1$ -norm ball

## Some remarks

---

- ▶ Why care?
  - simple
  - low-memory
  - stochastic version possible



## Some remarks

---

- ▶ Why care?
  - simple
  - low-memory
  - stochastic version possible
- ▶ Another perspective

$$x^{k+1} = \min_{x \in \mathcal{C}} \langle x, g^k \rangle + \frac{1}{2\alpha_k} \|x - x_k\|^2$$

Mirror Descent

## Some remarks

---

- ▶ Why care?
  - simple
  - low-memory
  - stochastic version possible
- ▶ Another perspective

$$x^{k+1} = \min_{x \in \mathcal{C}} \langle x, g^k \rangle + \frac{1}{2\alpha_k} \|x - x_k\|^2$$

Mirror Descent

- ▶ Improvements using more information (heavy-ball, filtered subgradient, ...)

## Some remarks

---

- ▶ Why care?
  - simple
  - low-memory
  - stochastic version possible
- ▶ Another perspective

$$x^{k+1} = \min_{x \in \mathcal{C}} \langle x, g^k \rangle + \frac{1}{2\alpha_k} \|x - x_k\|^2$$

### Mirror Descent

- ▶ Improvements using more information (heavy-ball, filtered subgradient, ...)
- ▶ Don't forget the dual!
  - may be more amenable to optimization
  - duality gap

## What we did not cover

---

- ♠ Adaptive stepsize tricks
- ♠ Space dilation methods, quasi-Newton style subgrads
- ♠ Barrier subgradient method
- ♠ Sparse subgradient method
- ♠ Ellipsoid method, center of gravity, etc. as subgradient methods
- ♠ And many more

# References

---

- ♠ S. Boyd, EE364b Slides
- ♠ Bertsekas, Nonlinear Programming