

*Annotated  
Version*

# Optimization for Machine Learning

## CS-439

### Lecture 4: Projected and Proximal Gradient Descent

**Martin Jaggi**

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

March 16, 2018

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

## Theorem

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable. Let  $X \subseteq \mathbb{R}^d$  be a closed convex set, and assume that there is a minimizer  $\mathbf{x}^*$  of  $f$  over  $X$ ; furthermore, suppose that  $f$  is  $L$ -smooth over  $X$ . When choosing the stepsize

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

projected gradient descent with  $\mathbf{x}_0 \in X$  satisfies:

(i) Function values are monotone decreasing:

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

Exercise 19

(ii)

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad T > 0.$$

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

Proof. Use smoothness:

Notation  $x^+ := x_{t+1}$   
 $x := x_t$

$$f(x^+) \leq f(x) + \underbrace{\nabla f(x)}_{//}^T (x^+ - x) + \frac{L}{2} \|x^+ - x\|^2$$

$\delta = \frac{1}{L}$

$$= f(x) - L(\underbrace{y - x}_{//})^T (x^+ - x) + \frac{L}{2} \|x^+ - x\|^2$$

$$2u^T v = \|u\|^2 + \|v\|^2 - \|u - v\|^2$$

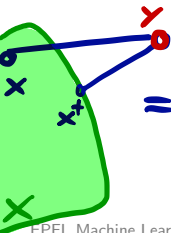
$$= f(x) - \frac{L}{2} \left( \underbrace{\|y - x\|^2}_{=} + \underbrace{\|x^+ - x\|^2}_{\times} - \underbrace{\|y - x^+\|^2}_{\times} \right) + \frac{L}{2} \|x^+ - x\|^2$$

$$= \frac{1}{L} \|\nabla f(x)\|^2$$

$$= f(x) - \frac{1}{2L} \|\nabla f(x)\|^2 + \frac{L}{2} \|y - x^+\|^2 \Rightarrow ;)$$

□

For ii): Mimic smooth unconstrained ...



# Strongly convex constrained minimization:

$\mathcal{O}(\log(1/\varepsilon))$  steps

## Theorem

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable. Let  $X \subseteq \mathbb{R}^d$  be a closed and convex set and suppose that  $f$  is smooth over  $X$  with parameter  $L$  and strongly convex over  $X$  with parameter  $\mu > 0$ .

Choosing

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

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

(i)

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

(ii)

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^t \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

same as unconstrained

# Strongly convex constrained minimization:

$\mathcal{O}(\log(1/\varepsilon))$  steps

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

Proof.

Strengthen the “constrained” vanilla bound

$$\underbrace{f(\mathbf{x}_t^+) - f(\mathbf{x}^*)}_{\leq} \leq \frac{1}{2\gamma} (\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}^+ - \mathbf{x}^*\|^2 - \|\mathbf{y}^+ - \mathbf{x}^+\|^2)$$

to

$$\frac{1}{2\gamma} (\gamma^2 \|\nabla f(\mathbf{x}_t)\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}^+ - \mathbf{x}^*\|^2 - \|\mathbf{y}^+ - \mathbf{x}^+\|^2) - \frac{\mu}{2} \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

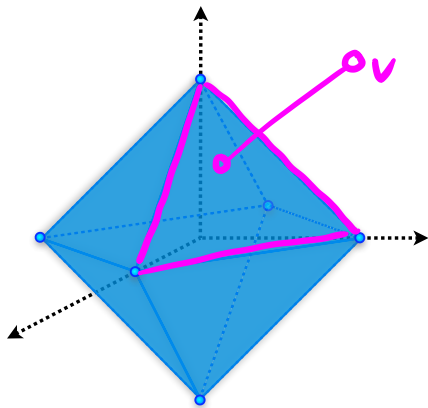
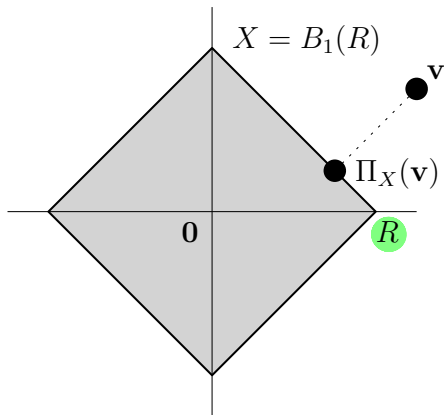
using strong convexity.

Then proceed as in the unconstrained theorem. □

*continue ... use property of projection*

# Projecting onto $\ell_1$ -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 $\ell_1$ -balls

→ project onto simplex

w.l.o.g.

- ▶  $R = 1,$  ✓
  - ▶  $v_i \geq 0$  for all  $i,$  ✓
  - ▶  $\sum_{i=1}^d v_i > 1.$  ✓
- (\*)

And using this,  
*corollary*

$\mathbf{x} = \Pi_X(\mathbf{v})$  satisfies  $x_i \geq 0$  for all  $i$  and  $\sum_{i=1}^d x_i = 1.$

*simplex, see next*

proof:  $x_i < 0$  ? No!

could flip sign  $x_i \leftarrow -x_i$  (Exercise)

# Projecting onto $\ell_1$ -balls

## Corollary

Under our assumption (\*),

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

where

$$\Delta_d := \left\{ \mathbf{x} \in \mathbb{R}^d : \sum_{i=1}^d x_i = 1, x_i \geq 0 \ \forall i \right\}$$

is the standard simplex.

Also, w.l.o.g. assume that  $v$  is ordered decreasingly,

$$v_1 \geq v_2 \geq \dots \geq v_d.$$

$\geq 0$

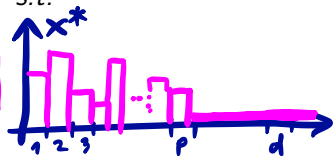


# Projecting onto $\ell_1$ -balls

## Lemma

Let  $\mathbf{x}^* := \operatorname{argmin}_{\mathbf{x} \in \Delta_d} \|\mathbf{x} - \mathbf{v}\|^2$ , and  $\mathbf{v}$  ordered decreasingly.  
There exists (a unique) index  $p \in \{1, \dots, d\}$  s.t.

$$\begin{aligned} x_i^* &> 0, & i \leq p, \\ x_i^* &= 0, & i > p. \end{aligned}$$



Proof.

$$d_{\mathbf{v}}(\mathbf{x}) := \|\mathbf{x} - \mathbf{v}\|^2$$

Optimality criterion for constrained optimization:

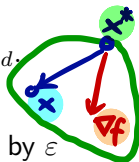
$$\nabla d_{\mathbf{v}}(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) = 2(\mathbf{x}^* - \mathbf{v})^\top (\mathbf{x} - \mathbf{x}^*) \geq 0, \quad \forall \mathbf{x} \in \Delta_d.$$

$\exists$  a positive entry in  $\mathbf{x}^*$  (because  $\sum_{i=1}^d x_i^* = 1$ ).

Why not  $x_i^* = 0$  and  $x_{i+1}^* > 0$ ? If so, we could decrease  $x_{i+1}^*$  by  $\varepsilon$  and increase  $x_i^*$  to  $\varepsilon$  to obtain  $\mathbf{x} \in \Delta_d$  s.t.

$$(\mathbf{x}^* - \mathbf{v})^\top (\mathbf{x} - \mathbf{x}^*) = (0 - v_i)\varepsilon - (x_{i+1}^* - v_{i+1})\varepsilon = \varepsilon \underbrace{(v_{i+1} - v_i)}_{\leq 0} - \underbrace{(x_{i+1}^*)}_{> 0} < 0,$$

contradicting the optimality.  $\square$



# Projecting onto $\ell_1$ -balls

Can say more about  $\mathbf{x}^*$ :

## Lemma

With  $p$  as in the above Lemma, and  $\mathbf{v}$  ordered decreasingly, we have

$$x_i^* = v_i - \Theta_p, \quad i \leq p,$$

constant !

where

$$\Theta_p = \frac{1}{p} \left( \sum_{i=1}^p v_i - 1 \right).$$

Proof.

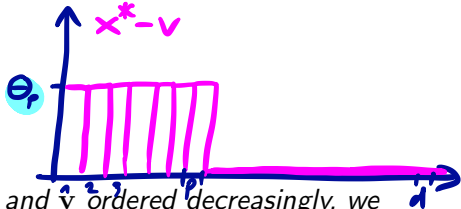
Assume there is  $i, j \leq p$  with  $x_i^* - v_i < x_j^* - v_j$ . As before, we could decrease  $x_i^* > 0$  by  $\varepsilon$  and increase  $x_j^*$  by  $\varepsilon$  to get  $\mathbf{x} \in \Delta_d$  s.t.

$$(\mathbf{x}^* - \mathbf{v})^\top (\mathbf{x} - \mathbf{x}^*) = (x_i^* - v_i)\varepsilon - (x_j^* - v_j)\varepsilon = \varepsilon \underbrace{((x_i^* - v_i) - (x_j^* - v_j))}_{< 0} < 0,$$

$\nabla d_v$

again contradicting optimality of  $\mathbf{x}^*$ .

□



# Projecting onto $\ell_1$ -balls

**Summary:** have  $d$  candidates for  $\mathbf{x}^*$ , namely

$$\mathbf{x}^*(p) := \underbrace{(v_1 - \Theta_p, \dots, v_p - \Theta_p, 0, \dots, 0)}, \quad p \in \{1, \dots, d\},$$

Need to find the right one. In order for candidate  $\mathbf{x}^*(p)$  to comply with our first Lemma, we must have

$$v_p - \Theta_p > 0,$$

and this actually ensures  $\mathbf{x}^*(p)_i > 0$  for all  $i \leq p$  (because  $\mathbf{v}$  is ordered) and therefore  $\mathbf{x}^*(p) \in \Delta_d$ .

But there could still be several choices for  $p$ . Among them, we simply pick the one for which  $\mathbf{x}^*(p)$  minimizes the distance to  $\mathbf{v}$ .

In time  $\mathcal{O}(d \log d)$ , by first sorting  $v$  and checking incrementally.

# Projecting onto $\ell_1$ -balls

## Theorem

Let  $\mathbf{v} \in \mathbb{R}^d$ ,  $R \in \mathbb{R}_+$ ,  $X = B_1(R)$  the  $\ell_1$ -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

„composite objective function“

Consider objective functions composed as

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

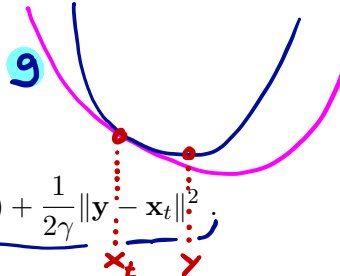
where  $g$  is a “nice” function, where as  $h$  is a “simple” additional term, which however doesn’t satisfy the assumptions of niceness which we used in the convergence analysis so far.

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

# Idea

From unconstrained minimization

The classical gradient step for minimizing  $g$ :

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


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

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

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  $f$ .



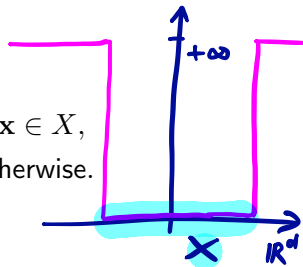
# A generalization of gradient descent?

- ▶  $h \equiv 0$ : 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

$$\iota_X : \mathbb{R}^d \rightarrow \mathbb{R} \cup +\infty$$

$$\mathbf{x} \mapsto \iota_X(\mathbf{x}) := \begin{cases} 0 & \text{if } \mathbf{x} \in X, \\ +\infty & \text{otherwise.} \end{cases}$$



Proximal mapping becomes

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

$$\text{projection!} = \Pi_X(\mathbf{z})$$

# Convergence in $\mathcal{O}(1/\varepsilon)$ steps

Same as vanilla case for smooth functions, but now for any  $h$  for which we can compute the proximal mapping.

- Examples:
- $h(x) = \lambda \|x\|_1$   
prox is soft-thresholding operator.
  - $h(x) = \mathbf{1}_{B_1(R)}$   
prox is projection on to  $\ell_1$ -ball.