

# SDSC6015 Stochastic Optimization for Machine Learning

Lu Yu

Department of Data Science, City University of Hong Kong

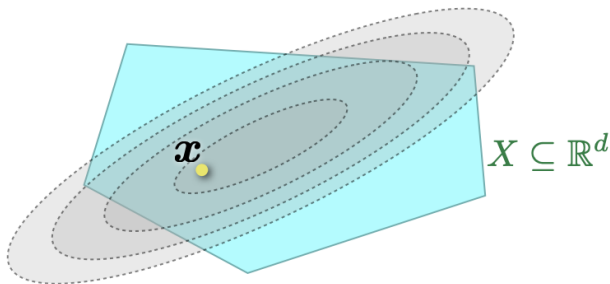
September 25, 2025

# Projected Gradient Descent

# Constrained Optimization

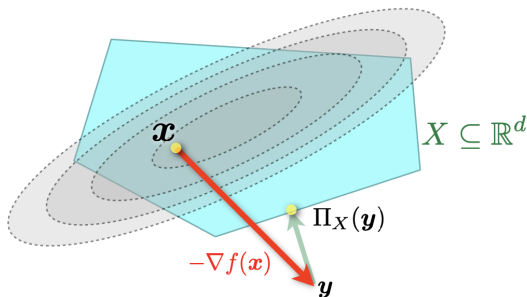
## Constrained Optimization Problem

$$\begin{array}{ll}\text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in X\end{array}$$



# Projected Gradient Descent

**Idea:** project onto  $X$  after every step:  $\Pi_X(\mathbf{y}) := \arg \min_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|$



Projected gradient descent:  $\mathbf{x}_{t+1} = \Pi_X[\mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)]$

# Projected Gradient Descent

Projected gradient descent:

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

$$\mathbf{x}_{t+1} := \Pi_X(\mathbf{y}_{t+1}) = \arg \min_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|$$

for **stepsize**  $\eta_t > 0$  and **timesteps**  $t = 0, 1, \dots$

# Convergence Rate of Projected Gradient Descent

The same number of steps as a gradient over  $\mathbb{R}^d$ !

- ▶ Lipschitz convex functions over  $X$  :  $\mathcal{O}(1/\varepsilon^2)$  steps
- ▶ Smooth convex functions over  $X$  :  $\mathcal{O}(1/\varepsilon)$  steps
- ▶ Smooth and strongly convex functions over  $X$  :  $\mathcal{O}(\log(1/\varepsilon))$  steps

We will adapt the previous proofs for gradient descent.

BUT:

- ▶ Each step involves a projection onto  $X$
- ▶ may or may not be efficient...

# Smooth functions over $X$

Recall:

$f$  is called smooth (with parameter  $L$ ) over  $X$  if

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in X.$$

## Lemma 1

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be differentiable and smooth with parameter  $L$  over  $X$ .  
Choosing stepsize

$$\eta = \frac{1}{L},$$

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

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



# Projected Gradient Descent on Smooth Functions

## Theorem 1

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 smooth over  $X$  with parameter  $L$ .  
Choosing stepsize

$$\eta = \frac{1}{L},$$

projected gradient descent yields

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

# Projected Gradient Descent on Smooth Functions

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

## Proof.

As before, use sufficient decrease to bound sum of squared gradients in vanilla analysis:

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

But now: **extra** term  $\frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2$ .

# Projected Gradient Descent on Smooth Functions

- Replace  $\mathbf{x}_{t+1}$  in the vanilla analysis with  $\mathbf{y}_{t+1}$  (the unprojected gradient step):

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{1}{2\eta} (\eta^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}^*\|^2).$$

- Use Fact (ii):  $\|\mathbf{x} - \Pi_X(\mathbf{y})\|^2 + \|\mathbf{y} - \Pi_X(\mathbf{y})\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2$ .
- With  $\mathbf{x} = \mathbf{x}^*$ ,  $\mathbf{y} = \mathbf{y}_{t+1}$ , we have  $\Pi_X(\mathbf{y}) = \mathbf{x}_{t+1}$ , and hence

$$\|\mathbf{x}^* - \mathbf{x}_{t+1}\|^2 + \underbrace{\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2}_{\leq \|\mathbf{x}^* - \mathbf{y}_{t+1}\|^2} \leq \|\mathbf{x}^* - \mathbf{y}_{t+1}\|^2$$

- We get back to the vanilla analysis...but with a saving!

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) \leq \frac{1}{2\eta} \left( \eta^2 \|\mathbf{g}_t\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 - \underbrace{\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2}_{\leq \|\mathbf{x}^* - \mathbf{y}_{t+1}\|^2} \right)$$

# Projected Gradient Descent on Smooth Functions

- Using  $f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*)$  (convexity), vanilla analysis with saving,  $\eta = 1/L$ :

$$\begin{aligned} \sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}^*)) &\leq \sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) \\ &\leq \frac{1}{2L} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2. \end{aligned}$$

- Using sufficient decrease to bound  $\frac{1}{2L} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2$  by

$$\begin{aligned} \sum_{t=0}^{T-1} \left( f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}) + \frac{L}{2} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \right) \\ = f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 \end{aligned}$$

# Projected Gradient Descent on Smooth Functions

- ▶ Putting it together: extra terms cancel, and as in unconstrained case, we get

$$\sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

- ▶ By the definitions of  $\mathbf{x}_{t+1}$  and  $\mathbf{y}_{t+1}$ ,

$$\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\| \leq \|\mathbf{y}_{t+1} - \mathbf{x}_t\| = \eta \|\nabla f(\mathbf{x}_t)\|.$$

Combining this with “Succifiect Decrease ”

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t).$$

Again, we make progress at every step!

- ▶ Hence,

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

# Strongly Convex Functions over $X$

Recall:

$f$  is **strongly convex** (with parameter  $\mu$ ) over  $X$  if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in X.$$

**Lemma 2**

A strongly convex function has a unique minimizer  $\mathbf{x}^*$  of  $f$  over  $X$ .

We prove that projected gradient descent converges to  $\mathbf{x}^*$ .

# Projected GD on Strongly Convex and Smooth Functions

## Theorem 2

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and differentiable. Let  $X \subset \mathbb{R}^d$  be a nonempty 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  $\eta = 1/L$ , **projected** gradient descent with arbitrary  $\mathbf{x}_0$  satisfies the following two properties.

- Squared distances to  $\mathbf{x}^*$  are geometrically decreasing:

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

- The absolute error after  $T$  iterations is exponentially small in  $T$ :

$$\begin{aligned} f(\mathbf{x}_T) - f(\mathbf{x}^*) &\leq \|\nabla f(\mathbf{x}^*)\| \left(1 - \frac{\mu}{L}\right)^{T/2} \|\mathbf{x}_0 - \mathbf{x}^*\| \leftarrow \text{usually, } \nabla f(\mathbf{x}^*) \neq \mathbf{0}! \\ &+ \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2, \quad T > 0. \leftarrow \text{as unconstrained case} \end{aligned}$$

# Projected GD on Strongly Convex and Smooth Functions

Proof.

(i) Geometric decrease plus noise:  $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq \dots$

► unconstrained case:

$$2\eta(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \eta^2 \|\nabla f(\mathbf{x}_t)\|^2 + \underline{(1 - \mu\eta)\|\mathbf{x}_t - \mathbf{x}^*\|^2}.$$

► constrained case (vanilla analysis with a saving):

$$2\eta(f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \eta^2 \|\nabla f(\mathbf{x}_t)\|^2 - \|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2 + \underline{(1 - \mu\eta)\|\mathbf{x}_t - \mathbf{x}^*\|^2}.$$



# Projected GD on Strongly Convex and Smooth Functions

To bound the noise, we use sufficient decrease.

► unconstrained case:

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

► constrained case:

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

Putting it together, the terms  $\|\mathbf{y}_{t+1} - \mathbf{x}_{t+1}\|^2$  cancel, and we get

$$\|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 \leq (1 - \mu\eta) \|\mathbf{x}_t - \mathbf{x}^*\|^2 = \left(1 - \frac{\mu}{L}\right) \|\mathbf{x}_t - \mathbf{x}^*\|^2$$

in both cases.

# Projected GD on Strongly Convex and Smooth Functions

(ii) Error bound from smoothness:

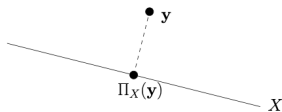
$$\begin{aligned} f(\mathbf{x}_T) - f(\mathbf{x}^*) &\leq \nabla f(\mathbf{x}^*)^\top (\mathbf{x}_T - \mathbf{x}^*) + \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_T\|^2 \\ &\leq \|\nabla f(\mathbf{x}^*)\| \|\mathbf{x}_T - \mathbf{x}^*\| + \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_T\|^2 \quad (\text{Cauchy-Schwarz}) \\ &\leq \|\nabla f(\mathbf{x}^*)\| \left(1 - \frac{\mu}{L}\right)^{T/2} \|\mathbf{x}_0 - \mathbf{x}^*\| + \frac{L}{2} \left(1 - \frac{\mu}{L}\right)^T \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \quad (i) \end{aligned}$$

constrained error bound  $\approx \sqrt{\text{unconstrained error bound}}$   
required number of steps roughly doubles.

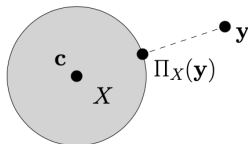
# The Projection Step

Computing  $\Pi_X(\mathbf{y}) := \arg \min_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{y}\|$  is an optimization problem itself. It can efficiently be solved in relevant cases:

- Projecting onto an affine subspace (leads to system of linear equations, similar to least squares)



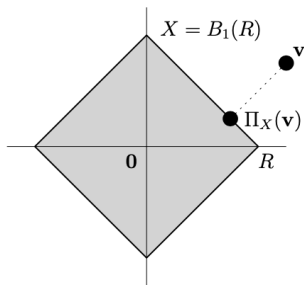
- Projecting onto a Euclidean ball with center  $\mathbf{c}$  (simply scale the vector  $\mathbf{y} - \mathbf{c}$ )



# Projecting onto $\ell_1$ -balls (needed in Lasso)

W.l.o.g. restrict to center at  $\mathbf{0}$

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



►  $B_1(R)$  is the cross polytope ( $2d$  vertices,  $2^d$  facets)

# Projecting onto $\ell_1$ -balls

- This problem can be reduced to a projection onto a simplex set

$$\Pi_X(\mathbf{v}) = \arg \min_{\mathbf{x} \in \Delta_d} \|\mathbf{x} - \mathbf{v}\|^2,$$

where  $\Delta_d := \{\mathbf{x} \in \mathbb{R}^d : \sum_{i=1}^d x_i = 1, x_i \geq 0, \forall i\}$  is called the **standard simplex**.

- Projection onto a simplex can be computed in  $\mathcal{O}(d \log d)$  time (can be improved to  $\mathcal{O}(d)$ ) [DSSSC08]

# Questions?

# Proximal Gradient Descent

# Composite Optimization Problems

Consider objective functions composed as

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

- ▶  $g$  is a “nice” function
- ▶  $h$  is a “simple” additional term, which however doesn’t satisfy the assumptions of niceness which we used in the convergence analysis so far.
- ▶ Proximal Gradient Descent is useful for solving nonsmooth, constrained, or structured optimization problems.



The classical gradient step for minimizing a differentiable function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is typically written as:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \nabla g(\mathbf{x}_t)$$

An equivalent way to express this step is by [minimizing a local quadratic approximation of  \$g\(\mathbf{x}\)\$](#) , given by:

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x}} \left\{ g(\mathbf{x}_t) + \nabla g(\mathbf{x}_t)(\mathbf{x} - \mathbf{x}_t) + \frac{1}{2\eta} \|\mathbf{x} - \mathbf{x}_t\|^2 \right\}$$

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

$$\begin{aligned}\mathbf{x}_{t+1} &:= \arg \min_{\mathbf{y}} \left\{ g(\mathbf{x}_t) + \nabla g(\mathbf{x}_t)^\top (\mathbf{y} - \mathbf{x}_t) + \frac{1}{2\eta} \|\mathbf{y} - \mathbf{x}_t\|^2 + h(\mathbf{y}) \right\} \\ &= \arg \min_{\mathbf{y}} \left\{ \frac{1}{2\eta} \|\mathbf{y} - (\mathbf{x}_t - \eta \nabla g(\mathbf{x}_t))\|^2 + h(\mathbf{y}) \right\},\end{aligned}$$

the proximal gradient descent update.

# Proximal Gradient Descent Algorithm

An iteration of **proximal gradient descent** is defined as

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

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

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

The update step can be equivalently written as

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

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

# A Generalization of Gradient Descent?

- ▶  $h = 0$  : recover **gradient descent**
- ▶  $h = \mathbf{1}_X$  : recover **projected gradient descent**

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

$$\mathbf{1}_X : \mathbb{R}^d \rightarrow \mathbb{R} \cup +\infty$$
$$\mathbf{x} \mapsto \mathbf{1}_X(\mathbf{x}) := \begin{cases} 0 & \text{if } \mathbf{x} \in X \\ +\infty & \text{otherwise.} \end{cases}$$

Proximal mapping becomes

$$\text{prox}_{h,\eta}(\mathbf{x}) = \arg \min_{\mathbf{y}} \left\{ \frac{1}{2\eta} \|\mathbf{y} - \mathbf{z}\|^2 + \mathbf{1}_X(\mathbf{y}) \right\} = \arg \min_{\mathbf{y} \in X} \|\mathbf{y} - \mathbf{z}\|^2$$

# Proximal Gradient Descent

Aim to minimize

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

Proximal Gradient Descent iteration:

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

where

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

- ▶ Convergence of proximal gradient can be as fast as classic gradient descent
- ▶ In every iteration, we have to additionally compute the proximal mapping  $h$ .

# Convergence of Proximal Gradient Descent

## Theorem 3

Let  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  be convex and smooth with parameter  $L$ ,  $h$  convex, and  $\text{prox}_{h,\eta}(\mathbf{x}) = \arg \min_{\mathbf{z}} \{\|\mathbf{x} - \mathbf{z}\|^2 / (2\eta) + h(\mathbf{z})\}$  can be computed. Choosing the fixed stepsize

$$\eta = \frac{1}{L},$$

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

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

- Intuitively, this means that proximal method only “sees” the nice smooth part  $g$  of the objective, and is not impacted by the additional  $h$ , which it treats separately in each step.

# Convergence of Proximal Gradient Descent

## Proof.

- Recall that the proximal step could be written as

$$\begin{aligned}\mathbf{x}_{t+1} &= \arg \min_{\mathbf{y} \in \mathbb{R}^d} \{g(\mathbf{x}_t) + \nabla g(\mathbf{x}_t)^\top (\mathbf{y} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}_t\|^2 + h(\mathbf{y})\} \\ &= \arg \min_{\mathbf{y} \in \mathbb{R}^d} \{\psi(\mathbf{y})\},\end{aligned}$$

where the function

$$\psi(\mathbf{y}) = g(\mathbf{x}_t) + \nabla g(\mathbf{x}_t)^\top (\mathbf{y} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}_t\|^2 + h(\mathbf{y})$$

is strongly convex with  $L$ .

# Convergence of Proximal Gradient Descent

- By the definition of  $\mathbf{x}_{t+1}$  and strong convexity of  $\psi$

$$\psi(\mathbf{y}) \geq \psi(\mathbf{x}_{t+1}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}_{t+1}\|^2,$$

which is equivalent to

$$\begin{aligned} & \nabla g(\mathbf{x}_t)^\top (\mathbf{y} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}_t\|^2 + h(\mathbf{y}) \\ & \geq \nabla g(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + h(\mathbf{x}_{t+1}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}_{t+1}\|^2 \end{aligned}$$

- Rearranging terms and subtracting  $h(\mathbf{x}_t)$  from both sides,

$$\begin{aligned} & \nabla g(\mathbf{x}_t)^\top (\mathbf{y} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}_t\|^2 - \frac{L}{2} \|\mathbf{y} - \mathbf{x}_{t+1}\|^2 + h(\mathbf{y}) - h(\mathbf{x}_t) \\ & \geq \nabla g(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 + h(\mathbf{x}_{t+1}) - h(\mathbf{x}_t) \end{aligned}$$



# Convergence of Proximal Gradient Descent

- ▶ As the function  $g$  is  $L$ -smooth, we can estimate the right side as

$$\nabla g(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \geq g(\mathbf{x}_{t+1}) - g(\mathbf{x}_t).$$

- ▶ Since  $g$  is convex, on the left side we estimate

$$\nabla g(\mathbf{x}_t)^\top (\mathbf{y} - \mathbf{x}_t) \leq g(\mathbf{y}) - g(\mathbf{x}_t).$$

- ▶ Putting this together

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

This holds for any  $\mathbf{y} \in \mathbb{R}^d$ .

# Convergence of Proximal Gradient Descent

- Let  $\mathbf{y} = \mathbf{x}^*$  and sum up the inequality over  $t = 0$  to  $t = T - 1$

$$\sum_{t=0}^{T-1} (f(\mathbf{x}^*) - f(\mathbf{x}_t)) + \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_0\|^2 - \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}_T\|^2 \geq f(\mathbf{x}_T) - f(\mathbf{x}_0).$$

- Note that  $f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t)$ , as  $\psi(\mathbf{x}_{t+1}) \leq \psi(\mathbf{x}_t)$  for each  $0 \leq t \leq T$

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{1}{T} \sum_{t=1}^T (f(\mathbf{x}_t) - f(\mathbf{x}^*)) \leq \frac{L}{2T} \|\mathbf{x}^* - \mathbf{x}_0\|^2.$$

# Convergence of Proximal Gradient Descent

$$R^2 = \|\mathbf{x}_0 - \mathbf{x}^*\|^2.$$

$$T \geq \frac{R^2 L}{2\varepsilon} \quad \Rightarrow \quad \text{error} \leq \frac{L}{2T} R^2 \leq \varepsilon.$$

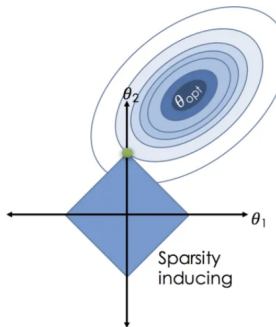
- The convergence rate  $\mathcal{O}(\frac{1}{\varepsilon})$  is the same as [gradient descent](#) on convex smooth functions (Theorem 2 from lecture 2), but now for any  $h$  for which we can compute the proximal mapping.

# Example: Iterative Soft-Thresholding Algorithm (ISTA)

## Lasso regression:

$$\min_{\beta} f(\beta) = \frac{1}{2} \|\mathbf{y} - A\beta\|_2^2 + \lambda \|\beta\|_1$$

- ▶  $A \in \mathbb{R}^{n \times d}$  is the feature matrix,  $\beta \in \mathbb{R}^d$  is the coefficient vector,  $\mathbf{y} \in \mathbb{R}^n$  is the response variable
- ▶  $\lambda > 0$  is a tuning parameter that controls sparsity (the  $\ell_1$ -norm forces coefficients  $\beta_j$  to be exactly zero, performing feature selection)



## Example: Iterative Soft-Thresholding Algorithm (ISTA)

$$f(\boldsymbol{\beta}) = \underbrace{\frac{1}{2}\|\mathbf{y} - A\boldsymbol{\beta}\|_2^2}_{g(\boldsymbol{\beta})} + \underbrace{\lambda\|\boldsymbol{\beta}\|_1}_{h(\boldsymbol{\beta})}$$

Proximal mapping is now

$$\text{prox}_{h,\eta}(\mathbf{z}) = \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2\eta} \|\boldsymbol{\beta} - \mathbf{z}\|^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\} =: S_{\lambda,\eta}(\mathbf{z}).$$

Here,  $S_{\lambda,\eta}(\mathbf{z})$  is the soft-thresholding operator

$$[S_{\lambda,\eta}(\mathbf{z})]_i = \begin{cases} z_i - \eta\lambda & \text{if } z_i > \eta\lambda \\ 0 & \text{if } |z_i| < \eta\lambda \\ z_i + \eta\lambda & \text{if } z_i < -\eta\lambda \end{cases}, \quad i = 1, \dots, d$$

# Example: Iterative Soft-Thresholding Algorithm (ISTA)

Recall  $\nabla g(\beta) = -A^\top(\mathbf{y} - A\beta)$ , hence proximal gradient update is:

$$\mathbf{x}_{t+1} = S_{\lambda, \eta}(\mathbf{x}_t + \eta A^\top(\mathbf{y} - A\mathbf{x}_t)).$$

Often called the **iterative soft-thresholding algorithm (ISTA)**

# Questions?

# Mirror Descent



# Mirror Descent: Motivation

Consider the simplex-constrained optimization problem

$$\min_{\mathbf{x} \in \Delta_d} f(\mathbf{x}),$$

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

Now, we assume  $\|\nabla f(\mathbf{x})\|_\infty = \max_{i=1, \dots, d} |[\nabla f(\mathbf{x})]_i| \leq 1, \forall \mathbf{x} \in \Delta_d$ .

- ▶ The largest element of any gradient is bounded by 1.
- ▶ All the elements of any gradient are bounded by 1.
- ▶ The extreme cases here are the following two vectors taken as the gradient.

$$(\text{the minimal vector}) \mathbf{0}_d = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}, \quad (\text{the maximal vector}) \mathbf{1}_d = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

# Mirror Descent: Motivation

- ▶ For the vector  $\mathbf{1}_d$ , it has  $\ell_2$ -norm  $\|\mathbf{1}_d\|_2 = \sqrt{d}$
- ▶ In other words,  $\|\nabla f(\mathbf{x})\|_\infty \leq 1$  gives  $\|\nabla f(\mathbf{x})\|_2 \leq L = \sqrt{d}$
- ▶ Convergence of GD (on convex and  $L$ -Lipschitz functions):

$$f(\mathbf{x}_{\text{best}}^{(T)}) - f(\mathbf{x}^*) \leq R\sqrt{\frac{d}{T}}$$

- ▶ It turns out the rate  $\mathcal{O}(\sqrt{\frac{d}{T}})$  is not optimal, mirror descent can do better as  $\mathcal{O}(\sqrt{\frac{\log d}{T}})$

# Mirror Descent: Preliminary

- Fix an arbitrary norm  $\|\cdot\|$  on  $\mathbb{R}^d$ , and a compact convex set  $X \subseteq \mathbb{R}^d$ . The dual norm  $\|\cdot\|_*$  is defined as

$$\|\mathbf{g}\|_* = \sup_{\|\mathbf{x}\| \leq 1} \mathbf{g}^\top \mathbf{x}.$$

- We say that a convex function  $f : X \rightarrow \mathbb{R}$  is
  - $L$ -Lipschitz w.r.t.  $\|\cdot\|$  if  $\forall \mathbf{x} \in X, \mathbf{g} \in \partial f(\mathbf{x}), \|\mathbf{g}\|_* \leq L$
  - $\beta$ -smooth w.r.t.  $\|\cdot\|$  if
$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_* \leq \beta \|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x}, \mathbf{y} \in X$$
  - $\mu$ -strongly convex w.r.t.  $\|\cdot\|$  if
$$f(\mathbf{x}) - f(\mathbf{y}) \leq \mathbf{g}^\top (\mathbf{x} - \mathbf{y}) - \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|^2, \forall \mathbf{x}, \mathbf{y} \in X, \mathbf{g} \in \partial f(\mathbf{x})$$

# Mirror Descent

Consider the **mirror descent** [Nemirovski and Yudin (1983)] iteration

$$\mathbf{y}_{t+1} = (\nabla\Phi)^{-1}(\nabla\Phi(\mathbf{x}_t) - \eta_t \mathbf{g}_t) \quad \text{and} \quad \mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in X} D_\Phi(\mathbf{x}, \mathbf{y}_{t+1}),$$

with  $\mathbf{g}_t \in \partial f(\mathbf{x}_t)$ .

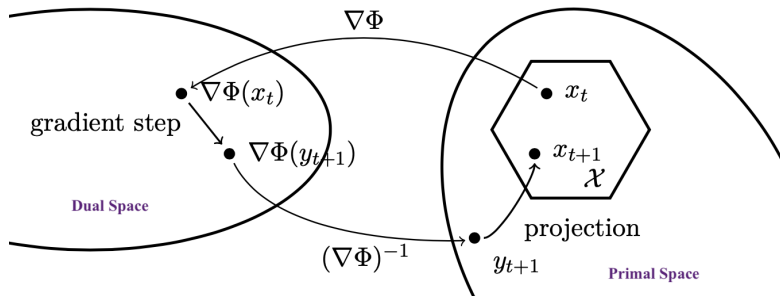
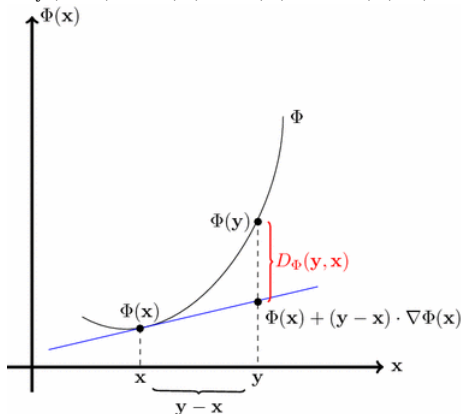


Figure: The “geometry” of mirror descent from [Bubeck 2015].

# Mirror Descent: Key elements

- Mirror potential  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}$  is strictly convex, continuously differentiable with  $\lim_{\|x\|_2 \rightarrow \infty} \|\nabla \Phi(x)\| = \infty$ .
- Define the Bregman divergence associated to  $f$  as

$$D_f(x, y) = f(x) - f(y) - \nabla f(y)^\top (x - y).$$



# Mirror Descent: Key elements

- The projection via Bregman divergence associated to  $\Phi$

$$\Pi_X^\Phi(\mathbf{y}) = \arg \min_{\mathbf{x} \in X} D_\Phi(\mathbf{x}, \mathbf{y}), \quad \forall \mathbf{y} \in X.$$

- Properties of  $\Phi$  ensures the existence and uniqueness of this projection  $\Pi_X^\Phi$ .

# Convergence of Mirror Descent

Let  $\mathbf{x}_1 \in \arg \min_{\mathbf{x} \in X} \Phi(\mathbf{x})$ . For  $t \geq 1$ , let  $\mathbf{y}_{t+1} \in \mathbb{R}^d$  such that

$$\nabla \Phi(\mathbf{y}_{t+1}) = \nabla \Phi(\mathbf{x}_t) - \eta \mathbf{g}_t, \text{ where } \mathbf{g}_t \in \partial f(\mathbf{x}_t),$$

and

$$\mathbf{x}_{t+1} \in \Pi_X^\Phi(\mathbf{y}_{t+1}).$$

## Theorem 4

Let

- ▶  $\Phi$  be a mirror map  $\rho$ -strongly convex on  $X$  w.r.t  $\|\cdot\|$ .
- ▶  $R^2 = \sup_{\mathbf{x} \in X} \Phi(\mathbf{x}) - \Phi(\mathbf{x}_1)$ .
- ▶  $f$  be convex and  $L$ -Lipschitz w.r.t  $\|\cdot\|$ .

Mirror descent with  $\eta = \frac{2R}{L} \sqrt{\frac{\rho}{T}}$  satisfies

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

# Convergence of Mirror Descent

To prove Theorem 4, we need some auxiliary arguments.

- Given the Bregman divergence associated to  $\Phi$ , it holds that

$$(\nabla\Phi(\mathbf{x}) - \nabla\Phi(\mathbf{y}))^\top(\mathbf{x} - \mathbf{z}) = D_\Phi(\mathbf{x}, \mathbf{y}) + D_\Phi(\mathbf{z}, \mathbf{x}) - D_\Phi(\mathbf{z}, \mathbf{y}).$$

- Moreover,

$$D_\Phi(\mathbf{x}, \Pi_X^\Phi(\mathbf{y})) + D_\Phi(\Pi_X^\Phi(\mathbf{y}), \mathbf{y}) \leq D_\Phi(\mathbf{x}, \mathbf{y}).$$



# Convergence of Mirror Descent

Proof.

$$\begin{aligned} & f(\mathbf{x}_t) - f(\mathbf{x}) \\ & \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}) \\ & = \frac{1}{\eta} (\nabla \Phi(\mathbf{x}_t) - \nabla \Phi(\mathbf{y}_{t+1}))^\top (\mathbf{x}_t - \mathbf{x}) \\ & = \frac{1}{\eta} \left( D_\Phi(\mathbf{x}, \mathbf{x}_t) + D_\Phi(\mathbf{x}_t, \mathbf{y}_{t+1}) - D_\Phi(\mathbf{x}, \mathbf{y}_{t+1}) \right) \\ & \leq \frac{1}{\eta} \left( D_\Phi(\mathbf{x}, \mathbf{x}_t) + D_\Phi(\mathbf{x}_t, \mathbf{y}_{t+1}) - D_\Phi(\mathbf{x}, \mathbf{x}_{t+1}) - D_\Phi(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) \right). \end{aligned}$$

The term  $D_\Phi(\mathbf{x}, \mathbf{x}_t) - D_\Phi(\mathbf{x}, \mathbf{x}_{t+1})$  will lead to a telescopic sum when summing over  $t = 1$  to  $t = T$ . It remains to bound the other term...

# Convergence of Mirror Descent

Recall that  $\Phi$  is  $\rho$ -strongly convex, and  $az - bz^2 \leq \frac{a^2}{4b}, \forall z \in \mathbb{R}$ .

$$\begin{aligned} & D_{\Phi}(\mathbf{x}_t, \mathbf{y}_{t+1}) - D_{\Phi}(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}) \\ &= \Phi(\mathbf{x}_t) - \Phi(\mathbf{x}_{t+1}) - \nabla \Phi(\mathbf{y}_{t+1})^{\top} (\mathbf{x}_t - \mathbf{x}_{t+1}) \\ &\leq (\nabla \Phi(\mathbf{x}_t) - \nabla \Phi(\mathbf{y}_{t+1}))^{\top} (\mathbf{x}_t - \mathbf{x}_{t+1}) - \frac{\rho}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &= \eta \mathbf{g}_t^{\top} (\mathbf{x}_t - \mathbf{x}_{t+1}) - \frac{\rho}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &\leq \eta L \|\mathbf{x}_t - \mathbf{x}_{t+1}\| - \frac{\rho}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &\leq \frac{(\eta L)^2}{2\rho}. \end{aligned}$$

Thus,

$$\sum_{t=1}^T \left( f(\mathbf{x}_t) - f(\mathbf{x}) \right) \leq \frac{D_{\Phi}(\mathbf{x}, \mathbf{x}_1)}{\eta} + \eta \frac{L^2 T}{2\rho}.$$

- **“Ball setup”**. The mirror potential

$$\Phi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|_2^2, \quad \forall \mathbf{x} \in \mathbb{R}^d.$$

- Associated Bregman divergence  $D_\Phi(\mathbf{x}, \mathbf{y}) = \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|_2^2$ .
- This is exactly equivalent to the projected subgradient descent.

# Standard Setups for Mirror Descent

- **“Simplex setup”**. The mirror potential

$$\Phi(\mathbf{x}) = \sum_{i=1}^d x_i \log(x_i), \quad \mathbf{x} \in \mathbb{R}_{++}^d = \{\mathbf{x} \in \mathbb{R}^d : x_i > 0, i = 1, \dots, d\}.$$

- The gradient update  $\nabla\Phi(\mathbf{y}_{t+1}) = \nabla\Phi(\mathbf{x}_t) - \eta\nabla f(\mathbf{x}_t)$  can be written as

$$[\mathbf{y}_{t+1}]_i = [\mathbf{x}_t]_i \exp(-\eta[\nabla f(\mathbf{x}_t)]_i), \quad i = 1, \dots, d.$$

- The Bregman divergence of  $\Phi$  is

$$D_{\Phi}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^d x_i \log(x_i/y_i) \quad (\text{Kullback-Leibler divergence})$$

- Projection of  $\mathbf{y}$  onto the simplex  $\triangle_d$  under the KL divergence leads to renormalization  $\mathbf{y} \rightarrow \mathbf{y}/\|\mathbf{y}\|_1$  (see notes).
- For  $X = \triangle_d$ ,  $\mathbf{x}_1 = (1/d, \dots, 1/d)$  and  $R^2 = \log d$  (see notes).

# References



Sébastien Bubeck, *Convex optimization: Algorithms and complexity*, Foundations and Trends in Machine Learning **8** (2015), no. 3-4, 231–357.



John Duchi, Shai Shalev-Shwartz, Yoram Singer, and Tushar Chandra, *Efficient projections onto the  $l_1$ -ball for learning in high dimensions*, Proceedings of the 25th international conference on Machine learning, 2008, pp. 272–279.



Mor Harchol-Balter, *Introduction to probability for computing*, Cambridge University Press, 2023.



Laurent Lessard, Benjamin Recht, and Andrew Packard, *Analysis and design of optimization algorithms via integral quadratic constraints*, SIAM Journal on Optimization **26** (2016), no. 1, 57–95.



Arkadij Semenovič Nemirovskij and David Borisovich Yudin, *Problem complexity and method efficiency in optimization*.



Stephen J Wright and Benjamin Recht, *Optimization for data analysis*, Cambridge University Press, 2022.