



# EP3260: Machine Learning Over Networks

## Lecture 2: Centralized Convex ML

(part 1)

Hossein S. Ghadikolaei, Hadi Ghauch, and Carlo Fischione

Division of Network and Systems Engineering  
School of Electrical Engineering and Computer Science  
KTH Royal Institute of Technology, Stockholm, Sweden

<https://sites.google.com/view/mlons/home>

January 2019

# Learning outcomes

- Basic definitions of smooth and convex functions
- Important properties of smooth and convex functions
- Main (deterministic) iterative algorithms for convex problems
- Connections among them
- Pros and cons of them
- Convergence analysis

# Outline

1. Student groups
2. Basic definitions and properties
3. Iterative solution approaches
4. Supplements

# Outline

1. Student groups
2. Basic definitions and properties
3. Iterative solution approaches
4. Supplements

## Student groups

Any question?

# Outline

1. Student groups
2. Basic definitions and properties
3. Iterative solution approaches
4. Supplements

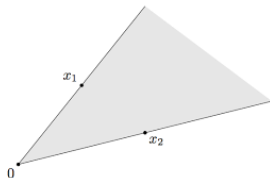
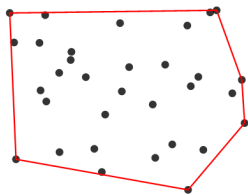
# Basic definitions

**Convex combination** of points  $\mathcal{X} = \{x_i\}_{i \in [n]}$  is  $\sum_{i \in [n]} \theta_i x_i$  where  $\theta = [\theta_1, \dots, \theta_n]$  form a probability simplex;  $\theta \geq 0, \theta^T \mathbf{1} = 1$

**Conic combination** of  $x_1$  and  $x_2$  is  $\theta_1 x_1 + \theta_2 x_2$  where  $\theta_i \geq 0$

**Convex hull of  $\mathcal{X}$** : set of all convex combinations of points in  $\mathcal{X}$

**Convex cone of  $\mathcal{X}$** : set of all conic combinations of points in  $\mathcal{X}$



# Convexity: basic definitions

**Convex set**  $\mathcal{X}$ :  $\forall x_1, x_2 \in \mathcal{X}$  and  $\theta \in [0, 1]$ ,  $\theta x_1 + (1 - \theta)x_2 \in \mathcal{X}$ .

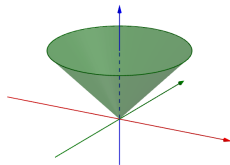
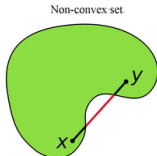
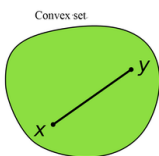
Euclidean/norm ball with radius  $r$  centered at  $x_c$ :

$$\{x \mid \|x - x_c\|_2 \leq r\} = \{x_c + ru \mid \|u\|_2 \leq 1\}$$

Norm cone:  $\{(x, r) \mid \|x\| \leq r\}$

Ellipsoid centered at  $x_c$ : for symmetric positive definite  $P$  and square nonsingular  $A$ :

$$\{(x - x_c)^T P^{-1}(x - x_c) \leq 1\} \equiv \{x_c + Au \mid \|u\|_2 \leq 1\}$$





# Convexity: basic definitions

**Convex function**  $f : \mathcal{X} \rightarrow \mathbb{R}$ : if  $\mathcal{X}$  is convex and  $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$  and  $\theta \in [0, 1]$ :

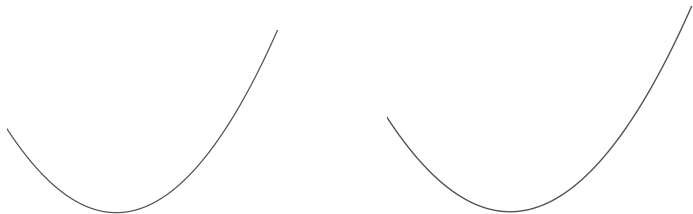
$$f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) \leq \theta f(\mathbf{x}_1) + (1 - \theta) f(\mathbf{x}_2) \quad (1)$$

Assuming differentiability of  $f$ , (1) is equivalent to

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) \quad (2)$$

Local information (gradient) determines a global lower bound

For twice differentiable  $f$ ,  $(1) \Leftrightarrow \nabla^2 f(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}$ : PSD Hessian, non-negative curvature everywhere.



# Convexity: basic definitions

**Convex function**  $f : \mathcal{X} \rightarrow \mathbb{R}$ : if  $\mathcal{X}$  is convex and  $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$  and  $\theta \in [0, 1]$ :

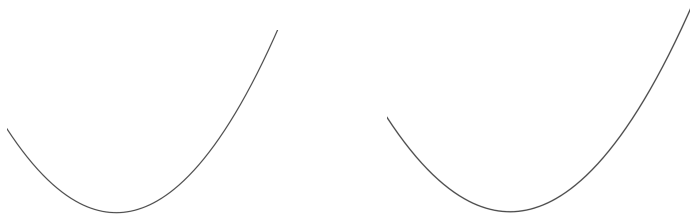
$$f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) \leq \theta f(\mathbf{x}_1) + (1 - \theta) f(\mathbf{x}_2) \quad (1)$$

Assuming differentiability of  $f$ , (1) is equivalent to

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) \quad (2)$$

Local information (gradient) determines a global lower bound

For twice differentiable  $f$ ,  $(1) \Leftrightarrow \nabla^2 f(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}$ : PSD Hessian, non-negative curvature everywhere.



# Convexity: basic definitions

**Convex function**  $f : \mathcal{X} \rightarrow \mathbb{R}$ : if  $\mathcal{X}$  is convex and  $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$  and  $\theta \in [0, 1]$ :

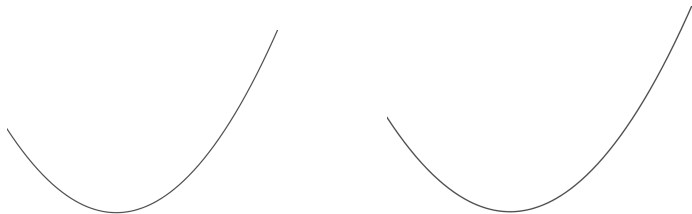
$$f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) \leq \theta f(\mathbf{x}_1) + (1 - \theta) f(\mathbf{x}_2) \quad (1)$$

Assuming differentiability of  $f$ , (1) is equivalent to

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) \quad (2)$$

Local information (gradient) determines a global lower bound

For twice differentiable  $f$ , (1)  $\Leftrightarrow \nabla^2 f(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}$ : PSD Hessian, non-negative curvature everywhere.



# Convexity: some examples

- $\|x\|_p$  for any  $p \geq 1$
- Quadratic function  $f(x) = x^T A x + b^T x + c$  for symmetric matrix  $A$

$$\nabla^2 f(x) = 2A \quad \text{convex iff } A \geq 0$$

- $f(x) = \|Ax - b\|_2^2$  is convex for any  $A$  (observe  $\nabla^2 f(x) = 2A^T A$ )
- $\|Ax - b\|_p$  for any  $p \geq 1$
- $\max_{i \in [n]} \|A_i^T x - b_i\|_p$  for any  $p \geq 1$
- $\lambda_{\max} = \sup_{\|y\|_2=1} y^T X y$  for any symmetric  $X$
- $\text{proj}(x, \mathcal{C}) := \inf_{y \in \mathcal{C}} \|y - x\|$  for convex  $\mathcal{C}$
- Check Boyd and Vandenberghe (2003) for more examples

Nonnegative weighted sum, composition with affine function, pointwise maximum and supremum, composition, minimization, and perspective

## Convexity: more definitions

- $f : \mathcal{X} \rightarrow \mathbb{R}$  is quasi-convex if  $\mathcal{X}$  is convex and sub-level sets

$$\{\mathbf{x} \in \mathcal{X} \mid f(\mathbf{x}) \leq \alpha\}$$

are convex for all  $\alpha$ .

Equivalently, if for all  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$  and  $\theta \in [0, 1]$ ,

$$f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) \leq \max\{f(\mathbf{x}_1), f(\mathbf{x}_2)\}$$

- A positive function is log-concave if  $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$  and  $\theta \in [0, 1]$

$$\log f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) \geq \theta \log f(\mathbf{x}_1) + (1 - \theta) \log f(\mathbf{x}_2)$$

Most of probability densities are log-concave

For convex  $\mathcal{S}$  and random variable  $\mathbf{y}$  with log-concave PDF,  
 $f(\mathbf{x}) = \Pr(\mathbf{x} + \mathbf{y} \in \mathcal{S})$  is log-concave and  $\{\mathbf{x} \mid f(\mathbf{x}) \geq t\}$  is convex

- Majorization ( $\mathbf{a} \succ \mathbf{b}$ ) and Schur convexity ( $\mathbf{a} \succ \mathbf{b}$  implies  $f(\mathbf{a}) \geq f(\mathbf{b})$ )

# Standard forms

$$\begin{aligned} &\text{minimize} && f_0(\mathbf{x}) \\ &\text{s.t.} && f_i(\mathbf{x}) \leq 0, i \in [m] \\ &&& h_i(\mathbf{x}) = 0, i \in [p] \end{aligned} \tag{3}$$

**Linear program:** affine objective over a (open/closed) polyhedron

**Quadratic program:** convex quadratic objective over a polyhedron

**Quadratically constrained quadratic program**

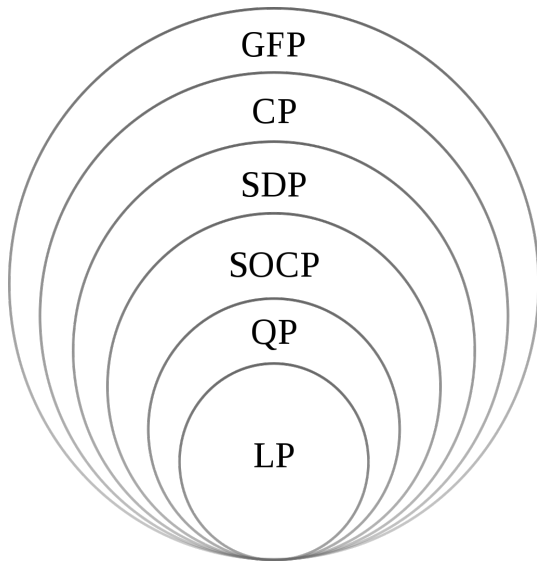
**Second-order cone program:**

$$\begin{aligned} &\text{minimize} && \mathbf{c}^T \mathbf{x} \\ &\text{s.t.} && \|\mathbf{A}_i \mathbf{x} + \mathbf{b}_i\|_2 \leq \mathbf{d}_i^T \mathbf{x} + \mathbf{e}_i, i \in [m], \quad \mathbf{F} \mathbf{x} = \mathbf{g} \end{aligned}$$

**Semidefinite program** for symmetric  $\mathbf{A}$  and  $\mathbf{B}$ :

$$\begin{aligned} &\text{minimize} && \mathbf{c}^T \mathbf{x} \\ &\text{s.t} && \mathbf{A} + \sum_{i \in [n]} x_i \mathbf{B}_i \leq 0, \quad \mathbf{D} \mathbf{x} = \mathbf{e} \end{aligned}$$

# Standard forms



# Duality

Consider (3) with objective  $f_0$ , inequality and equality constraints  $f_i$  and  $h_i$ , and optimal solution  $(x^*, f^* = f_0(x^*))$

**Lagrange dual function:**  $g : \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$

$$g(\lambda, \nu) = \inf_{x \in \mathcal{X}} L(x, \lambda, \nu) := f_0(x) + \sum_{i \in [m]} \lambda_i f_i(x) + \sum_{i \in [p]} \nu_i h_i(x)$$

if  $\lambda \geq 0$ , then  $g(\lambda, \nu) \leq f^*$ .

Proof: given non-negative  $\lambda_i$ ,  $h_i(x) = 0$ , and  $\lambda_i f_i(x) \leq 0$  for any feasible point  $x$ . Therefore,

$$f^* = f_0(x^*) \geq L(x^*, \lambda, \nu) \geq \inf_{x \in \mathcal{X}} L(x, \lambda, \nu) = g(\lambda, \nu)$$

**Lagrange dual problem:** with solution  $d^*(\leq f^*, = \text{with strong duality})$

$$\begin{aligned} &\text{maximize} && g(\lambda, \nu) \\ &\text{s.t.} && \lambda \geq 0 \end{aligned}$$



# Duality

**Primal:** minimize  $\mathbf{c}^T \mathbf{x}$

s.t.  $\mathbf{x} \geq 0$

$\mathbf{A}\mathbf{x} = \mathbf{b}$

**Dual:** maximize  $-\mathbf{b}^T \boldsymbol{\nu}$

s.t.  $\mathbf{A}^T \boldsymbol{\nu} + \mathbf{c} = 0$

In a network with one unit of communication per constraint, dual is more communication-efficient for tall  $\mathbf{A}$

- Check Boyd and Vandenberghe (2003) for more details

Weak and strong duality

Constraint qualifications

Slater's constraint qualification

Complementary slackness

Karush-Kuhn-Tucker (KKT) conditions

# Strong convexity

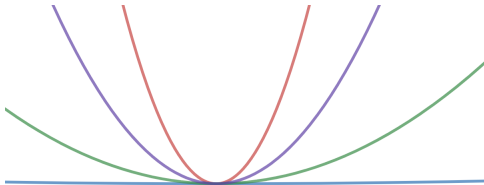
Differentiable function  $f$  is  $\mu$ -strongly convex iff  $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}, \mu > 0$

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) + \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2$$

Gradient can be replaced by sub-gradient for non-smooth functions

**Main intuition:** linear lower bound with convexity, quadratic lower bound with strong convexity

Global definitions not local ( $\forall \mathbf{x} \in \mathcal{X}$ )



## Strong convexity

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) + \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2 \quad (4)$$

- (4) is equivalent to a minimum positive curvature  $\nabla^2 f(\mathbf{x}) \geq \mu \mathbf{I}_d, \forall \mathbf{x} \in \mathcal{X}$

- (4) is equivalent to  $(\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1) \geq \mu \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2$

- (4) implies

$$(a) \quad f(\mathbf{x}) - f^* \leq \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2, \forall \mathbf{x}$$

$$(b) \quad \|\mathbf{x}_2 - \mathbf{x}_1\|_2 \leq \frac{1}{\mu} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2, \forall \mathbf{x}_1, \mathbf{x}_2$$

$$(c) \quad (\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1) \leq \frac{1}{\mu} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2^2, \forall \mathbf{x}_1, \mathbf{x}_2$$

(d)  $f(\mathbf{x}) + r(\mathbf{x})$  is strongly convex for any convex  $f$  and strongly convex  $r$

**HW1(a):** prove all the statements of this slide

# Smoothness

A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L$ -smooth iff it is differentiable and its gradient is  $L$ -Lipschitz-continuous (usually w.r.t. norm-2):

$$\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d, \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2 \leq L\|\mathbf{x}_2 - \mathbf{x}_1\|_2 \quad (5)$$

Recall strong convexity result:  $\|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2 \geq \mu\|\mathbf{x}_2 - \mathbf{x}_1\|_2$

For twice differentiable  $f$ ,  $(5) \leftrightarrow \nabla^2 f(\mathbf{x}) \leq L\mathbf{I}_d$

**Smoothness:**  $f(\mathbf{x}_2) - f(\mathbf{x}_1)$  can be over-estimated by a quadratic function

- (5) implies for all  $\mathbf{x}_1, \mathbf{x}_2$  (**HW1(b)**: prove them. Assume convexity if needed)

$$(a) \quad f(\mathbf{x}_2) \leq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T(\mathbf{x}_2 - \mathbf{x}_1) + \frac{L}{2}\|\mathbf{x}_2 - \mathbf{x}_1\|_2^2$$

$$(b) \quad f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T(\mathbf{x}_2 - \mathbf{x}_1) + \frac{1}{2L}\|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2^2$$

$$(c) \quad (\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T(\mathbf{x}_2 - \mathbf{x}_1) \geq \frac{1}{L}\|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2^2$$

# Outline

1. Student groups
2. Basic definitions and properties
3. Iterative solution approaches
4. Supplements

# Gradient descent

- Problem: minimize  $f(x)$  for some differentiable  $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\pm\infty\}$

**Gradient descent (GD)**, also called batch GD, full GD, ...

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) \quad (6)$$

for some sequence of non-negative step sizes  $(\alpha_k)_{k \in \mathbb{N}}$ .

## Theorem 1: Convergence of GD with constant step size

**Convex and  $L$ -smooth**  $f$  with  $\alpha \leq 1/L$  satisfies  $f(\mathbf{x}_k) - f^* \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2k\alpha}$ .

**$\mu$ -strongly convex and  $L$ -smooth**  $f$  with  $\alpha \leq 2/(\mu + L)$  satisfies  $f(\mathbf{x}_k) - f^* \leq e^{-ck} L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 / 2$  for  $c = 2\alpha\mu L / (\mu + L)$ . With  $\alpha = 2/(\mu + L)$ , we have  $\|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \leq \left(1 - \frac{2}{1+L/\mu}\right)^{2k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$ .

Smooth convex:  $\mathcal{O}(1/\epsilon)$  iterations for  $\epsilon$ -optimality

Smooth strongly-convex:  $\mathcal{O}(\log(1/\epsilon))$  iterations for  $\epsilon$ -optimality

# Foods for thought

1. Define the conjugate function as

$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^d} \mathbf{y}^T \mathbf{x} - f(\mathbf{x}).$$

Observe that  $f^*$  is convex even when  $f$  is not (why?). When  $f$  is  $\mu$ -strongly convex and  $L$ -smooth,  $f^*$  is  $\frac{1}{L}$ -strongly convex and  $\frac{1}{\mu}$ -smooth (why?).

2. Define projection operator for convex set  $\mathcal{X}$  as

$$\text{proj}(\mathbf{x}, \mathcal{X}) := \operatorname{argmin}_{\mathbf{y} \in \mathcal{X}} \|\mathbf{y} - \mathbf{x}\|.$$

Observe  $\|\text{proj}(\mathbf{y}, \mathcal{X}) - \mathbf{x}\|^2 \leq \|\mathbf{y} - \mathbf{x}\|^2$  for any  $\mathbf{x} \in \mathcal{X}$  and any  $\mathbf{y}$ . Modify (6) to solve  $\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ . What is the convergence of this “projected GD” algorithm?

3. Define the set of subgradients of  $f : \mathcal{X} \rightarrow \mathbb{R}$  at  $\mathbf{x} \in \mathcal{X}$  as

$$\partial f(\mathbf{x}) := \{s \mid f(\mathbf{x}) - f(\mathbf{y}) \leq s^T(\mathbf{x} - \mathbf{y}) \ \forall \mathbf{y} \in \mathcal{X}\}.$$

Modify (6) to solve  $\min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$  for non-smooth (non-differentiable) functions that are Lipschitz ( $|f(\mathbf{x}) - f(\mathbf{y})| \leq \beta \|\mathbf{x} - \mathbf{y}\|_2$ ). What is the convergence of this “projected sub-GD” for Lipschitz functions? What is the optimal step size?

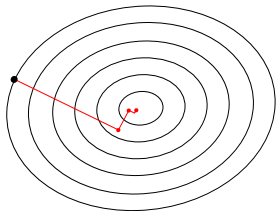
4. Check <https://ee227c.github.io/code/lecture4.html>

# GD for smooth and strongly convex functions

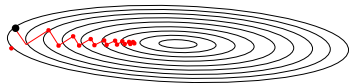
Linear convergence rate  $1 - \frac{2}{1 + L/\mu}$

GD may need many iterations to converge

Preconditioning to change the space geometry, to make sub-levels similar in all coordinates



small  $L/\mu$



large  $L/\mu$



# Descent methods

**Descent methods:**  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}(\mathbf{x}_k)$  s.t.  $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$  (7)

for some sequence of non-negative step sizes  $(\alpha_k)_{k \in \mathbb{N}}$  and decent direction  $\mathbf{d}$   
 $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$  implies  $-\nabla f(\mathbf{x}_k)^T \mathbf{d}(\mathbf{x}_k) > 0$ , same half-space as the negative gradient

**Steepest descent:**  $\mathbf{d}(\mathbf{x}) = \operatorname{argmin}\{\nabla f(\mathbf{x})^T \nu \mid \|\nu\| = 1\}$  in some norm  $\|\cdot\|$

\* Note that we need to unnormalize the descent direction using the dual-norm

Maximizes the first order prediction of decrease (for small  $\nu$ ):

$$f(\mathbf{x} + \nu) - f(\mathbf{x}) \approx \nabla f(\mathbf{x})^T \nu$$

Define  $\|\mathbf{x}\|_P = (\mathbf{x}^T \mathbf{P} \mathbf{x})^{1/2}$  for positive definite  $\mathbf{P}$  (this is called Mahalanobis distance):  $\mathbf{d}(\mathbf{x}) = -\mathbf{P}^{-1} \nabla f(\mathbf{x})$

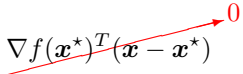
Reduces to GD on Euclidian norm ( $\mathbf{P} = \mathbf{I}$ ,  $\mathbf{d}(\mathbf{x}) = -\nabla f(\mathbf{x}_k)^T$ )

Good norm: should be consistent with the geometry of sublevel sets

Same theoretical convergence as of GD, much better in practice

# Newton methods

Around optimal point:

$$f(\mathbf{x}) \approx f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) + \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^T \nabla^2 f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*)$$


Sublevel sets are like ellipsoids (determined by Hessian) near the minimum

Recall  $\|\mathbf{x}\|_{\mathbf{P}} = (\mathbf{x}^T \mathbf{P} \mathbf{x})^{1/2}$  and its descent direction  $\mathbf{d} = \mathbf{P}^{-1} \nabla f(\mathbf{x})$

How about steepest descent on the norm induced by Hessian  $\nabla^2 f(\mathbf{x}^*)$  ?

Oops! we do not know  $\mathbf{x}^*$

**Newton method:**

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k). \quad (8)$$

# Newton methods

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla^2 f(\mathbf{x}_k)^{-1} \nabla f(\mathbf{x}_k)$$

## Theorem 2: Quadratic convergence of Newton's method

Assume  $f$  is a twice continuously differentiable and set  $\alpha_k = 1$ . If  $\|\mathbf{x}_k - \mathbf{x}^*\|$  is small enough, there exist a positive constant  $c$  such that  $\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq c(\|\mathbf{x}_k - \mathbf{x}^*\|_2)^2$ .

constant +  $\mathcal{O}(\log \log(1/\epsilon))$  iterations for  $\epsilon$ -optimality

**Expensive iterations** due to  $\nabla^2 f(\mathbf{x}_k)$

- **Affine invariance:** apply coordinate change for non-singular matrix  $\mathbf{A}$ . Newton's method have same iterations for  $\min_{\mathbf{x}} f(\mathbf{x})$  and  $\min_{\mathbf{y}} f(\mathbf{A}\mathbf{y})$  (namely  $\mathbf{x}_k = \mathbf{A}\mathbf{y}_k$ ), whereas GD has  $\nabla f(\mathbf{x}) = \mathbf{A}^T \nabla f(\mathbf{A}\mathbf{y})$ .

If we change coordinate/metric, GD iterates change

Finding a good coordinate for GD is usually very hard in high-dimension!

# Proximal methods

Smoothness implies  $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) + \frac{L}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2, \forall \mathbf{x}$

GD iterations to minimize differentiable  $f$  is like *successive quadratic upper-bound minimization*:

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2$$

New objective: minimize  $f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$  for convex differentiable  $g$  and convex (possibly) non-differentiable  $h$

Define **proximal mapping** as

$$\operatorname{prox}_{\alpha h}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{u}} h(\mathbf{u}) + \frac{1}{2\alpha} \|\mathbf{u} - \mathbf{x}\|_2^2$$

**Proximal method:**

$$\mathbf{x}_k = \operatorname{prox}_{\alpha_k h}(\mathbf{x}_{k-1} - \alpha_k \nabla g(\mathbf{x}_k))$$

# Proximal methods

Smoothness implies  $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) + \frac{L}{2}\|\mathbf{x} - \mathbf{x}_k\|_2^2, \forall \mathbf{x}$

GD iterations to minimize differentiable  $f$  is like *successive quadratic upper-bound minimization*:

$$\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2$$

New objective: minimize  $f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$  for convex differentiable  $g$  and convex (possibly) non-differentiable  $h$

Define **proximal mapping** as

$$\operatorname{prox}_{\alpha h}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{u}} h(\mathbf{u}) + \frac{1}{2\alpha} \|\mathbf{u} - \mathbf{x}\|_2^2$$

**Proximal method:**

$$\mathbf{x}_k = \operatorname{prox}_{\alpha_k h}(\mathbf{x}_{k-1} - \alpha_k \nabla g(\mathbf{x}_k))$$

# Proximal methods

- Observe from the definition of the proximal method

$$\begin{aligned}\mathbf{x}_k &= \operatorname{argmin}_{\mathbf{x}} h(\mathbf{x}) + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}_{k-1} + \alpha_k \nabla g(\mathbf{x}_k)\|_2^2 \\ &= \operatorname{argmin}_{\mathbf{x}} h(\mathbf{x}) + g(\mathbf{x}_k) + \nabla g(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2\alpha_k} \|\mathbf{x} - \mathbf{x}_k\|_2^2\end{aligned}$$

GD  $\leftrightarrow$  proximal method with  $h(\mathbf{x}) = 0$  and  $\alpha = 1/L$

Projected GD  $\leftrightarrow$  proximal method with  $h(\mathbf{x}) = \begin{cases} 0, & \text{if } \mathbf{x} \in \mathcal{X} \\ \infty, & \text{otherwise.} \end{cases}$

Soft thresholding for  $\ell_1$  regularization  $\leftrightarrow$  proximal method with  $h(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$ .

$$[\operatorname{prox}_h(\mathbf{x})]_i = \begin{cases} x_i - \lambda, & \text{if } x_i \geq \lambda \\ 0, & \text{if } -\lambda \leq x_i \leq \lambda \\ x_i + \lambda, & \text{if } x_i \leq -\lambda. \end{cases}$$

# Resource allocation

**HW1(c):** Consider

$$\begin{aligned} & \text{minimize} \quad \frac{1}{N} \sum_{i \in [N]} f_i(x_i) \\ & \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{b}. \end{aligned}$$

for  $\mathbf{A} \in \mathbb{R}^{p \times N}$  and  $\mathbf{x} = [x_1, \dots, x_N]^T$ .

- (a) Assume strong-convexity and smoothness on  $f$ . How would you solve this problem when  $N = 1000$ ?
- (b) What if  $N = 10^9$  ?
- (c) Can we use Newton's method for  $N = 10^9$ ? Try efficient method for computing  $\nabla^2 f(\mathbf{x}_k)$  for  $p = 1$  and  $b = 1$  (probability simplex constraint). Extend it to  $1 \leq p \ll N$ .
- (d) Now, add twice differentiable  $r(\mathbf{x})$  to the objective and solve (a)-(c).

## Some references

- S. Bubeck, "Convex optimization: Algorithms and complexity," Foundations and Trends in Machine Learning, 2015.
- Y. Nesterov, "Introductory lectures on convex optimization: A basic course, " Springer Science & Business Media, 2004.
- F. Bach, "Large-scale machine learning and convex optimization," Machine Learning Summer School, 2018.
- L. Bottou, F. Curtis, and J. Norcedal, "Optimization methods for large-scale machine learning," SIAM Rev., 2018.
- N. Parikh and S. Boyd, "Proximal algorithms", gradient methods," Foundations and Trends in Optimization, 2013.
- A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," SIAM Journal on Imaging Sciences, 2009.



# Outline

1. Student groups
2. Basic definitions and properties
3. Iterative solution approaches
4. Supplements

# Proof sketch for convex and $L$ -smooth $f$

By convexity of  $f$ ,  $f(\mathbf{x}_i) \leq f(\mathbf{x}^*) + \langle \nabla f(\mathbf{x}_i), \mathbf{x}_i - \mathbf{x}^* \rangle$  (\*)

Use smoothness property  $f(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^T (\mathbf{x}_{i+1} - \mathbf{x}_i) + \frac{L}{2} \|\mathbf{x}_{i+1} - \mathbf{x}_i\|_2^2$   
and  $\alpha \leq 1/L$  to conclude  $f(\mathbf{x}_{i+1}) \leq f(\mathbf{x}_i) - \frac{\alpha}{2} \|\nabla f(\mathbf{x}_i)\|_2^2$  (\*\*)

Substitute (\*) into (\*\*) and note that from (6):  $\nabla f(\mathbf{x}_i) = \frac{\mathbf{x}_i - \mathbf{x}_{i+1}}{\alpha}$

Observe  $f(\mathbf{x}_{i+1}) \leq f^* + \frac{1}{2\alpha} (\|\mathbf{x}_i - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{i+1} - \mathbf{x}^*\|_2^2)$

Summing over iterations:  $\frac{1}{k} \sum_{i \in [k]} (f(\mathbf{x}_i) \leq f^*) \leq \frac{1}{2\alpha k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$

Conclude from the non-increasing property of GD iterates:

$$f(\mathbf{x}_k) - f^* \leq \frac{1}{k} \sum_{i \in [k]} f(\mathbf{x}_i) - f^* \leq \frac{1}{2\alpha k} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$

# Proof sketch for strongly-convex and $L$ -smooth $f$

From smoothness and vanishing gradient of the optimal point, conclude

$$f(\mathbf{x}_i) - f(\mathbf{x}^*) \leq \frac{L}{2} \|\mathbf{x}_i - \mathbf{x}^*\|_2^2 \quad (*)$$

Use the coercivity of the gradient (**HW1(d)**: prove it)

$$(\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}))^T (\mathbf{x} - \mathbf{y}) \geq \frac{\mu L}{\mu + L} \|\mathbf{x} - \mathbf{y}\|_2^2 + \frac{1}{\mu + L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2$$

Use  $\alpha < 2/(L + \mu)$  to obtain  $\|\mathbf{x}_i - \mathbf{x}^*\|^2 \leq \left(1 - 2t\alpha \frac{\mu L}{\mu + L}\right) \|\mathbf{x}_i - \mathbf{x}^*\|^2$

Iterate over  $i$  and use  $(*)$  to obtain

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{L}{2} \prod_{i \in [k]} \left(1 - 2t\alpha \frac{\mu L}{\mu + L}\right) \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$$

Use  $e^{-x} \geq 1 - x$  to conclude  $f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{L}{2} e^{-ck} \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2$  for  $c = \frac{2\alpha\mu L}{\mu + L}$