### Mathematics for Machine Learning

— Continuous Optimization: Preliminary Convex Optimization

Joseph Chuang-Chieh Lin

Department of Computer Science & Information Engineering, Tamkang University

Fall 2023

#### Credits for the resource

- The slides are based on the textbooks:
  - Marc Peter Deisenroth, A. Aldo Faisal, and Cheng Soon Ong: Mathematics for Machine Learning. Cambridge University Press. 2020.
  - Howard Anton, Chris Rorres, Anton Kaul: Elementary Linear Algebra. Wiley. 2019.
- We could partially refer to the monograph: Francesco Orabona: A Modern Introduction to Online Learning. https://arxiv.org/abs/1912.13213

#### Outline

- Convex Programming
- 2 Linear Programming
- Quadratic Programming

#### Outline

- Convex Programming
- 2 Linear Programming
- Quadratic Programming

#### Our Focus & Motivation

#### Convex Optimization.

- A class of optimization problems where we can guarantee global optimality.
  - $f(\cdot)$  is a convex function.

The constraints  $g(\cdot)$  and  $h(\cdot)$  form convex sets.

#### Convex Sets & Functions

#### Convex set

A set C is convex if for any  $\mathbf{x}, \mathbf{y} \in C$ , we have

$$\forall \alpha \in [0,1], \alpha \mathbf{x} + (1-\alpha)\mathbf{y} \in \mathcal{C}.$$

#### Convex function

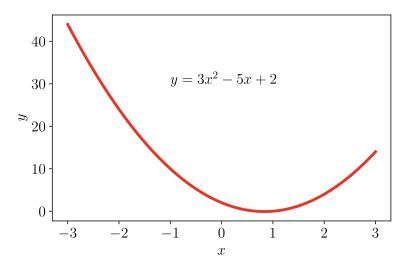
A function  $f: \mathcal{C} \subseteq \mathbb{R}^D \mapsto \mathbb{R}$  is convex if for any  $\mathbf{x}, \mathbf{y} \in \mathcal{C}$ ,

$$\forall \alpha \in [0,1], f((1-\alpha)\mathbf{x} + \alpha\mathbf{y}) \leq (1-\alpha)f(\mathbf{x}) + \alpha f(\mathbf{y}).$$

Equivalently, if f is differentiable (i.e.,  $\nabla f(\mathbf{x})$  exists for all  $\mathbf{x} \in \mathcal{C}$ ), then f is convex if and only if for all  $\mathbf{x}, \mathbf{y} \in \mathcal{C}$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla_{\mathbf{x}} f(\mathbf{x})^{\top} (\mathbf{y} - \mathbf{x}).$$

#### An Example of Convex Functions



#### Remark

• If  $f(\mathbf{x})$  is twice differentiable (i.e., the Hessian exists for all  $\mathbf{x} \in \mathcal{C}$ ), then

 $f(\mathbf{x})$  is convex  $\iff \nabla_{\mathbf{x}}^2 f(\mathbf{x})$  is positive semidefinite.

#### Example

Show that  $f(x) = x \lg x$  is convex for x > 0.

#### Example

Show that  $f(x) = x \lg x$  is convex for x > 0.

• Note:  $\lg x := \log_2 x$  and  $\ln x := \log_e x$ .

#### Example

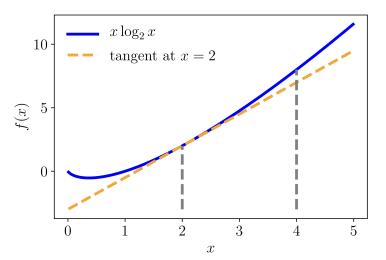
Show that  $f(x) = x \lg x$  is convex for x > 0.

- Note:  $\lg x := \log_2 x$  and  $\ln x := \log_e x$ .
- Compute  $\nabla_x f(x)$ .

#### Example

Show that  $f(x) = x \lg x$  is convex for x > 0.

- Note:  $\lg x := \log_2 x$  and  $\ln x := \log_e x$ .
- Compute  $\nabla_x f(x)$ .
- Say given x = 2, y = 4, compute  $f(x) + \nabla_x f(x) \top (y x)$ .



### Example (Theorem)

#### **Theorem**

Given a nonnegative real  $\alpha \geq 0$  and two convex functions  $f_1$  and  $f_2$ , then  $\alpha \cdot f_1 + (1 - \alpha)f_2$  is still convex.

### Example (Theorem)

#### **Theorem**

Given a nonnegative real  $\alpha \geq 0$  and two convex functions  $f_1$  and  $f_2$ , then  $\alpha \cdot f_1 + (1 - \alpha)f_2$  is still convex.

By definition,

$$f_1(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f_1(\mathbf{x}) + (1 - \alpha)f_1(\mathbf{y})$$
  
 $f_2(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f_2(\mathbf{x}) + (1 - \alpha)f_2(\mathbf{y})$ 

• Summing up:

$$f_1(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) + f_2(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y})$$

$$\leq \alpha f_1(\mathbf{x}) + (1 - \alpha)f_1(\mathbf{y}) + \alpha f_2(\mathbf{x}) + (1 - \alpha)f_2(\mathbf{y})$$

$$\alpha(f_1(\mathbf{x}) + f_2(\mathbf{x})) + (1 - \alpha)(f_1(\mathbf{y}) + f_2(\mathbf{y}))$$

### Example (Theorem)

#### **Theorem**

Given a nonnegative real  $\alpha \geq 0$  and two convex functions  $f_1$  and  $f_2$ , then  $\alpha \cdot f_1 + (1 - \alpha)f_2$  is still convex.

By definition,

$$f_1(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f_1(\mathbf{x}) + (1 - \alpha)f_1(\mathbf{y})$$
  
 $f_2(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f_2(\mathbf{x}) + (1 - \alpha)f_2(\mathbf{y})$ 

Summing up:

$$f_1(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) + f_2(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y})$$

$$\leq \alpha f_1(\mathbf{x}) + (1 - \alpha)f_1(\mathbf{y}) + \alpha f_2(\mathbf{x}) + (1 - \alpha)f_2(\mathbf{y})$$

$$\alpha (f_1(\mathbf{x}) + f_2(\mathbf{x})) + (1 - \alpha)(f_1(\mathbf{y}) + f_2(\mathbf{y}))$$

#### Outline

- Convex Programming
- 2 Linear Programming
- 3 Quadratic Programming

### Linear Programming

• Consider the special case that all the preceding functions are linear.

$$\label{eq:continuity} \begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^d} \quad & \mathbf{c}^\top \mathbf{x} \\ \text{subject to} \quad & \mathbf{A} \mathbf{x} \leq \mathbf{b}. \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times d}$  and  $\mathbf{b} \in \mathbb{R}^m$ .

• The Lagrangian:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{c}^{\top} \mathbf{x} + \boldsymbol{\lambda}^{\top} (\boldsymbol{A} \mathbf{x} - \mathbf{b})$$

where  $\lambda \in \mathbb{R}^m$  is the vector of non-negative Lagrange multipliers.

• The Lagrangian:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{c}^{\top}\mathbf{x} + \boldsymbol{\lambda}^{\top}(\boldsymbol{A}\mathbf{x} - \mathbf{b})$$

where  $\lambda \in \mathbb{R}^m$  is the vector of non-negative Lagrange multipliers.

Rearranging the terms:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = (\mathbf{c} + \boldsymbol{A}^{\top} \boldsymbol{\lambda})^{\top} \mathbf{x} - \boldsymbol{\lambda}^{\top} \mathbf{b}.$$

• The Lagrangian:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{c}^{\top} \mathbf{x} + \boldsymbol{\lambda}^{\top} (\boldsymbol{A} \mathbf{x} - \mathbf{b})$$

where  $\lambda \in \mathbb{R}^m$  is the vector of non-negative Lagrange multipliers.

Rearranging the terms:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = (\mathbf{c} + \mathbf{A}^{\top} \boldsymbol{\lambda})^{\top} \mathbf{x} - \boldsymbol{\lambda}^{\top} \mathbf{b}.$$

Taking the derivate w.r.t. x:

$$\mathbf{c} + \mathbf{A}^{\mathsf{T}} \boldsymbol{\lambda} = \mathbf{0}.$$

• The Lagrangian:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{c}^{\top}\mathbf{x} + \boldsymbol{\lambda}^{\top}(\boldsymbol{A}\mathbf{x} - \mathbf{b})$$

where  $\pmb{\lambda} \in \mathbb{R}^m$  is the vector of non-negative Lagrange multipliers.

• Rearranging the terms:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = (\mathbf{c} + \mathbf{A}^{\top} \boldsymbol{\lambda})^{\top} \mathbf{x} - \boldsymbol{\lambda}^{\top} \mathbf{b}.$$

Taking the derivate w.r.t. x:

$$\mathbf{c} + \mathbf{A}^{\mathsf{T}} \boldsymbol{\lambda} = \mathbf{0}.$$

• Thus, the dual Lagrangian is  $\mathcal{D}(\boldsymbol{\lambda}) = -\lambda^{\top} \mathbf{b}$ .

- Recall that we would like to maximize  $\mathcal{D}(\lambda)$  and the constraint that  $\lambda \geq \mathbf{0}$ .
- The dual optimization problem is

$$egin{array}{ll} \max_{m{\lambda} \in \mathbb{R}^m} & -\mathbf{b}^ op m{\lambda} \ & ext{subject to} & \mathbf{c} + m{A}^ op m{\lambda} = m{0} \ & m{\lambda} \geq m{0} \end{array}$$

which is also a linear program but with m variables.

- Recall that we would like to maximize  $\mathcal{D}(\lambda)$  and the constraint that  $\lambda \geq \mathbf{0}$ .
- The dual optimization problem is

$$\max_{\pmb{\lambda} \in \mathbb{R}^m} \qquad -\pmb{b}^\top \pmb{\lambda}$$
 subject to 
$$\pmb{c} + \pmb{A}^\top \pmb{\lambda} = \pmb{0}$$
 
$$\pmb{\lambda} \geq \pmb{0}$$

which is also a linear program but with m variables.

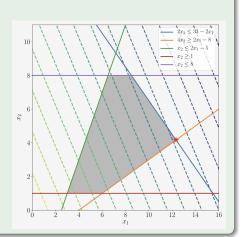
\* Solve the primal or the dual program depending on whether m (i.e., # constraints) or d (i.e., # variables) is larger.

#### Consider the linear program

$$\min_{\mathbf{x} \in \mathbb{R}^2} - \begin{bmatrix} 5 \\ 3 \end{bmatrix}^{\top} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

#### subject to

$$\begin{bmatrix} 2 & 2 \\ 2 & -4 \\ -2 & 1 \\ 0 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \le \begin{bmatrix} 33 \\ 8 \\ 5 \\ -1 \\ 8 \end{bmatrix}$$



#### Outline

- Convex Programming
- 2 Linear Programming
- Quadratic Programming

### Quadratic Programming

Consider the case of a convex quadratic objective function, where the constraints are affine:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \qquad \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{c}^\top \mathbf{x}$$
 subject to 
$$\mathbf{A} \mathbf{x} \leq \mathbf{b},$$

#### where

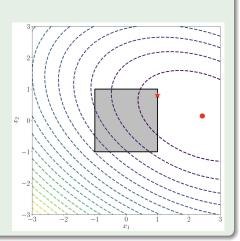
- $\mathbf{A} \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$  and  $\mathbf{c} \in \mathbb{R}^d$ .
- $Q \in \mathbb{R}^{d \times d}$ : a positive definite matrix. d variables and m linear constraints.

#### Consider the quadratic program

$$\min_{\mathbf{x} \in \mathbb{R}^2} \frac{1}{2} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^{\top} \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix}^{\top} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 5 \\ 3 \end{bmatrix}^{\top} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

#### subject to

$$\left[ egin{array}{ccc} 1 & 0 \ -1 & 0 \ 0 & 1 \ 0 & -1 \end{array} 
ight] \left[ egin{array}{c} x_1 \ x_2 \end{array} 
ight] \leq \left[ egin{array}{c} 1 \ 1 \ 1 \ 1 \end{array} 
ight].$$



Consider the case of a convex quadratic objective function, where the constraints are affine:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \qquad \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{c}^\top \mathbf{x}$$
  
subject to 
$$\mathbf{A} \mathbf{x} \leq \mathbf{b},$$

Consider the case of a convex quadratic objective function, where the constraints are affine:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \quad \frac{1}{2} \mathbf{x}^{\top} \mathbf{Q} \mathbf{x} + \mathbf{c}^{\top} \mathbf{x}$$
subject to  $\mathbf{A} \mathbf{x} \leq \mathbf{b}$ ,

The Lagrangian is given by

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{x}^{\top} \boldsymbol{Q} \mathbf{x} + \mathbf{c}^{\top} \mathbf{x} + \boldsymbol{\lambda}^{\top} (\boldsymbol{A} \mathbf{x} - \mathbf{b})$$
$$= \frac{1}{2} \mathbf{x}^{\top} \boldsymbol{Q} \mathbf{x} + (\mathbf{c} + \boldsymbol{A}^{\top} \boldsymbol{\lambda})^{\top} \mathbf{x} - \boldsymbol{\lambda}^{\top} \mathbf{b}.$$

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{x}^{\top} \boldsymbol{Q} \mathbf{x} + (\mathbf{c} + \boldsymbol{A}^{\top} \boldsymbol{\lambda})^{\top} \mathbf{x} - \boldsymbol{\lambda}^{\top} \mathbf{b}.$$

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{x}^{\top} \boldsymbol{Q} \mathbf{x} + (\mathbf{c} + \boldsymbol{A}^{\top} \boldsymbol{\lambda})^{\top} \mathbf{x} - \boldsymbol{\lambda}^{\top} \mathbf{b}.$$

Taking the derivative w.r.t. **x** and setting it to zero:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{x}^{\top} \boldsymbol{Q} \mathbf{x} + (\mathbf{c} + \boldsymbol{A}^{\top} \boldsymbol{\lambda})^{\top} \mathbf{x} - \boldsymbol{\lambda}^{\top} \mathbf{b}.$$

Taking the derivative w.r.t.  $\mathbf{x}$  and setting it to zero:

$$Qx + (c + A^{T}\lambda) = 0.$$

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{x}^{\top} \boldsymbol{Q} \mathbf{x} + (\mathbf{c} + \boldsymbol{A}^{\top} \boldsymbol{\lambda})^{\top} \mathbf{x} - \boldsymbol{\lambda}^{\top} \mathbf{b}.$$

Taking the derivative w.r.t.  $\mathbf{x}$  and setting it to zero:

$$\mathbf{Q}\mathbf{x} + (\mathbf{c} + \mathbf{A}^{\top} \boldsymbol{\lambda}) = \mathbf{0}.$$

Assume that Q is invertible, then

$$\mathbf{x} = -\mathbf{Q}^{-1}(\mathbf{c} + \mathbf{A}^{\top} \lambda).$$

Substituting it back to  $\mathcal{L}(\mathbf{x}, \lambda)$ , we get the dual Lagrangian

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{x}^{\top} \boldsymbol{Q} \mathbf{x} + (\mathbf{c} + \boldsymbol{A}^{\top} \boldsymbol{\lambda})^{\top} \mathbf{x} - \boldsymbol{\lambda}^{\top} \mathbf{b}.$$

Taking the derivative w.r.t.  $\mathbf{x}$  and setting it to zero:

$$Qx + (c + A^{T}\lambda) = 0.$$

Assume that Q is invertible, then

$$\mathbf{x} = -\mathbf{Q}^{-1}(\mathbf{c} + \mathbf{A}^{\top} \lambda).$$

Substituting it back to  $\mathcal{L}(\mathbf{x}, \lambda)$ , we get the dual Lagrangian

$$\mathcal{D}(\boldsymbol{\lambda}) = -rac{1}{2}(\mathbf{c} + \boldsymbol{A}^{ op} \boldsymbol{\lambda})^{ op} \boldsymbol{Q}^{-1}(\mathbf{c} + \boldsymbol{A}^{ op} \boldsymbol{\lambda}) - \boldsymbol{\lambda}^{ op} \mathbf{b}.$$

Therefore, the dual optimization problem is given by

$$\max_{oldsymbol{\lambda} \in \mathbb{R}^m} \quad -\frac{1}{2} (\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda})^\top \mathbf{Q}^{-1} (\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda}) - \boldsymbol{\lambda}^\top \mathbf{b}$$
 subject to  $\boldsymbol{\lambda} \geq \mathbf{0}$ 

• Heads up: Application in Support Vector Machine (SVM).

# **Discussions**