

# Where does ADMM come from: Lagrangian point of view

Andersen Ang

Dept. Combinatorics & Optimization, University of Waterloo, Canada

[msxang@uwaterloo.ca](mailto:msxang@uwaterloo.ca), where  $x = \lfloor \pi \rfloor$  Homepage: [angms.science](http://angms.science)

First draft: November 1, 2022    Last update: November 9, 2022

## The starting problem: a linearly constrained minimization

$$(\mathcal{P}) : \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{b}.$$

- ▶  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is
  - ▶ convex (will specify strictly convexity later)
  - ▶ closed (epigraph of  $f$  is a closed set, this is a common assumption)
  - ▶ possibly nondifferentiable ( $\nabla f$  is not continuous)
- ▶  $\mathbf{x} \in \mathbb{R}^n$  is the optimization variable.
- ▶  $\mathbf{A} \in \mathbb{R}^{m \times n}$ 
  - ▶ in general need not to be full rank
  - ▶ for studying theoretical property  $\mathbf{A}$  has to be full rank
- ▶  $\mathbf{b} \in \mathbb{R}^m$ .

# Lagrangian

$$(\mathcal{P}) : \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{Ax} = \mathbf{b}$$

- ▶ Joseph-Louis Lagrange: **add**<sup>1</sup> the constraint into the objective function!

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle.$$

- ▶ Lagrangian  $\mathcal{L}$  is a function of two variables:  $\mathbf{x}$  and  $\boldsymbol{\lambda} \in \mathbb{R}^m$ : Lagrangian multiplier
- ▶ Now we solve

$$(\mathcal{S}) : \min_{\mathbf{x}} \max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}). \quad \text{(Lagrangian saddle point problem)}$$

Assuming strong duality<sup>2</sup>, solving  $\mathcal{S}$  gives the same result as solving the original problem  $\implies$  that's why Lagrangian approach is useful.

---

<sup>1</sup>Actually you can subtract:  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) - \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle$ , the sign of  $\boldsymbol{\lambda}$  doesn't matter.

<sup>2</sup>Slater's condition:  $\mathcal{P}$  is a convex problem and  $\mathbf{Ax} = \mathbf{b}$  has a solution, which translate to  $\mathbf{A}$  is full rank.

# Why Lagrangian?

$$(\mathcal{P}) : \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{Ax} = \mathbf{b}$$

- ▶  $\mathcal{P}$  is **constrained**.
  - ▶ Some constraint is hard to solve.
  - ▶ To find  $\mathbf{x}^*$  by an iterative algorithm we may need to keep all iterates  $\mathbf{x}_k$  to be **feasible**.
  - ▶ You need to worry about both the feasibility and the optimality of  $\mathbf{x}$ .
- ▶  $\mathcal{S}$  has no constraint.

$$(\mathcal{S}) : \min_{\mathbf{x}} \max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}). \quad (\text{Lagrangian saddle point problem})$$

- ▶ All  $\mathbf{x}, \boldsymbol{\lambda}$  are feasible, no need to worry about them being infeasible, only need to worry about them being “not good enough” (not close to optimality).
- ▶ No feasibility restriction seems nice, but you are paying a trade off by introducing an extra variable  $\boldsymbol{\lambda}$ .

## Simplifying the Lagrangian: dual

$$(\mathcal{S}) : \min_x \max_{\lambda} \mathcal{L}(x, \lambda). \quad (\text{Lagrangian saddle point problem})$$

- ▶ We have two variables  $x, \lambda$ .
- ▶ It is natural to think of eliminating one variable.
  - ▶ The idea: we just introduced  $\lambda$  in  $\mathcal{L}$ , so eliminating  $\lambda$  means we go back to the original problem (making a U-turn), so we better eliminate  $x$  instead.
  - ▶ Generally speaking: in some cases eliminating  $\lambda$  is better.
- ▶ How to eliminate  $x$  from  $\mathcal{L}$ : treat  $\lambda$  as a constant and solve  $\min_x \mathcal{L}(x)$ .
- ▶ The remaining of  $\mathcal{L}$  after eliminating  $x$  is called the **dual**.

# Computing the dual

$$(\mathcal{P}) : \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{Ax} = \mathbf{b}$$

$$\min_{\mathbf{x}} \max_{\boldsymbol{\lambda}} \left\{ \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle \right\}$$

$$\begin{aligned} d(\boldsymbol{\lambda}) := \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}) &= \min_{\mathbf{x}} f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle \\ &= \min_{\mathbf{x}} f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} \rangle + \underbrace{\langle \boldsymbol{\lambda}, -\mathbf{b} \rangle}_{\text{constant}} \\ &= \min_{\mathbf{x}} f(\mathbf{x}) + \langle \mathbf{A}^\top \boldsymbol{\lambda}, \mathbf{x} \rangle - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle \\ &= - \left( \max_{\mathbf{x}} -f(\mathbf{x}) - \langle \mathbf{A}^\top \boldsymbol{\lambda}, \mathbf{x} \rangle \right) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle \\ &= - \left( \max_{\mathbf{x}} \langle -\mathbf{A}^\top \boldsymbol{\lambda}, \mathbf{x} \rangle - f(\mathbf{x}) \right) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle \\ &= -f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle \end{aligned}$$

$$\langle \boldsymbol{\lambda}, \mathbf{Ax} \rangle = \langle \mathbf{A}^\top \boldsymbol{\lambda}, \mathbf{x} \rangle$$

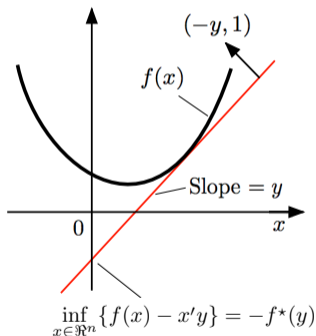
$$\min \phi(\mathbf{x}) = - \max -\phi(\mathbf{x})$$

$$f^*(\boldsymbol{\xi}) := \max_{\zeta} \langle \boldsymbol{\xi}, \zeta \rangle - f(\zeta)$$

- ▶ The dual function, only on  $\boldsymbol{\lambda}$  is  $d(\boldsymbol{\lambda}) = -f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle$ .
- ▶ If you are unfamiliar with convex analysis, the WTF-moment here is the conjugate  $f^*(\boldsymbol{\xi}) := \max_{\zeta} \left( \langle \boldsymbol{\xi}, \zeta \rangle - f(\zeta) \right)$ .

## Revision: what is conjugate (Bertsekas's explanations)

- ▶ Define a closed convex function by its epigraph.
- ▶ Describe the epigraph by supporting hyperplanes.
- ▶ Conjugate function = crossing point of the hyperplanes



Primal Description

Values  $f(x)$

Dual Description

Crossing points  $f^*(y)$

## Solving for $\lambda^*$ on the dual

$$(\mathcal{P}) : \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{Ax} = \mathbf{b}$$

$$\min_{\mathbf{x}} \max_{\boldsymbol{\lambda}} \left\{ \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle \right\} \xrightarrow[\mathbf{x} \text{ is gone}]{\min_{\mathbf{x}}} \max_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda}).$$

► The dual function  $d(\boldsymbol{\lambda}) = -f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle$  is only on  $\boldsymbol{\lambda}$ .

► If we solved for the optimal  $\boldsymbol{\lambda}$

$$\boldsymbol{\lambda}^* \in \operatorname{argmax}_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda}),$$

then we can solve for our original  $\mathbf{x}$  using  $\boldsymbol{\lambda}^*$  as

$$\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*).$$

► Why this works: strong duality<sup>3</sup>.

► Use  $\in$  instead of  $=$  since the solution to  $\operatorname{argmax}_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda})$ ,  $\operatorname{argmin}_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*)$  is possibly nonunique!

---

<sup>3</sup>Revision: weak duality is  $d^* < p^*$ , and strong duality is  $d^* = p^*$ .

Strong duality for  $\mathcal{P}$  holds if Slater's condition is true:  $f$  is convex and  $\mathbf{Ax} = \mathbf{b}$  has a solution.



# Algorithm based on Lagrangian for solving $\mathcal{P}$

$$(\mathcal{P}) : \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{Ax} = \mathbf{b}$$

---

## Algorithm 1: Algorithm based on Lagrangian

---

1 $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle$	write down the Lagrangian
2 $d(\boldsymbol{\lambda}) = -f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle$	derive the dual using conjugate
3 $\boldsymbol{\lambda}^* \in \underset{\boldsymbol{\lambda}}{\operatorname{argmax}} d(\boldsymbol{\lambda})$	solve the dual problem
4 $\mathbf{x}^* \in \underset{\mathbf{x}}{\operatorname{argmin}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*)$	solve the primal solution

---

Drawbacks: it assumes steps 3,4 are **have closed-form solution**.

- ▶ **Exactly solving**  $\operatorname{argmax} d(\boldsymbol{\lambda})$  in general can be hard.
- ▶ **Exactly solving**  $\operatorname{argmin} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*)$  in general can be hard.
- ▶ If we cannot solve step 3, then we cannot even proceed to step 4.

# How to improve **Algorithm based on Lagrangian**

Instead of exactly solving  $\operatorname{argmax}_{\lambda} d(\lambda)$ , we approximately solve it

---

**Algorithm 2:** An iterative algorithm

---

- 1 Write down  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{b} \rangle$  and  $d(\boldsymbol{\lambda}) = -f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle$
  - 2 Initialize  $\mathbf{x}_0, \boldsymbol{\lambda}_0$
  - 3 **for**  $k = 1, 2, \dots$  **do**
  - 4     Get  $\boldsymbol{\lambda}_{k+1}$  from  $\mathbf{x}_k, \boldsymbol{\lambda}_k$  that approximately solve  $\operatorname{argmax}_{\lambda} d(\lambda)$  better
  - 5     Solve  $\mathbf{x}_{k+1} \in \operatorname{argmin}_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_{k+1})$  using  $\boldsymbol{\lambda}_{k+1}$
  - 6 **end**
- 

How to approximately solve  $\operatorname{argmax}_{\lambda} d(\lambda)$ : gradient ascent

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \alpha_k \nabla_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda}) \Big|_{\boldsymbol{\lambda}=\boldsymbol{\lambda}_k}.$$

Is  $d(\boldsymbol{\lambda})$  even differentiable? Is  $\nabla_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda})$  unique? How to compute  $\nabla_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda})$ ?

# Differentiation of the dual

$$d(\boldsymbol{\lambda}) = -f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle$$

Some facts<sup>4</sup>:

1. Fact:  $d(\boldsymbol{\lambda})$  is concave<sup>5</sup> so gradient ascent solves  $\max d(\boldsymbol{\lambda})$  eventually if we can compute  $\nabla d$
2. Fact: if  $\mathbf{x}_k \in \operatorname{argmin} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k)$  is unique (only one minimizer), then
  - ▶  $d(\boldsymbol{\lambda})$  is differentiable at  $\boldsymbol{\lambda}_k$
  - ▶  $\nabla_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda}) \Big|_{\boldsymbol{\lambda}=\boldsymbol{\lambda}_k} = \mathbf{A}\mathbf{x}_k - \mathbf{b}$
3. Fact: if  $\mathbf{x}_k \in \operatorname{argmin} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k)$  is not unique, then
  - ▶  $d(\boldsymbol{\lambda})$  is not classically differentiable at  $\boldsymbol{\lambda}_k$
  - ▶ we have to use subgradient for  $d(\boldsymbol{\lambda})$  at  $\boldsymbol{\lambda}_k$
  - ▶  $\mathbf{A}\mathbf{x}_k - \mathbf{b}$  is only a subgradient of  $d(\boldsymbol{\lambda})$  at  $\boldsymbol{\lambda}_k$
4. Fact: if  $f$  is strictly convex then (2) is true and we have

$$\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \alpha_k(\mathbf{A}\mathbf{x}_k - \mathbf{b}).$$

---

<sup>4</sup>For proof see any convex analysis text.

<sup>5</sup>conjugate function is always convex, negative of convex is concave

## Dual ascent algorithm for solving Lagrangian saddle point problem

► Problem:  $\min_{\mathbf{x}} f(\mathbf{x})$  s.t.  $\mathbf{Ax} = \mathbf{b}$ ,  $f$  strictly convex and  $\mathbf{A}$  full rank.

► Lagrangian saddle point problem:  $\min_{\mathbf{x}} \max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$

---

### Algorithm 3: Dual ascent Lagrangian algorithm

---

```
1 Write down  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle$  and  $d(\boldsymbol{\lambda}) = -f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle$ 
2 Initialize  $\mathbf{x}_0, \boldsymbol{\lambda}_0$ 
3 for  $k = 1, 2, \dots$  do
4      $\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k)$ 
5      $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \alpha_k (\mathbf{Ax}_{k+1} - \mathbf{b})$ 
6 end
```

---

- Note: here the x-step and the y-step have swapped order.
- Need to determine dual ascent stepsize  $\alpha_k$
- This algorithm is already looking similar to ADMM!

# Drawback of the dual ascent algorithm

---

**Algorithm 4:** Dual ascent Lagrangian algorithm

---

```
1  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{b} \rangle$ 
2  $d(\boldsymbol{\lambda}) = -f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle$ 
3 Initialize  $\mathbf{x}_0, \boldsymbol{\lambda}_0$ 
4 for  $k = 1, 2, \dots$  do
5    $\mathbf{x}_{k+1} = \underset{\mathbf{x}}{\operatorname{argmin}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k)$ 
6    $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \alpha_k(\mathbf{A}\mathbf{x}_{k+1} - \mathbf{b})$ 
7 end
```

---

- ▶ We have equality  $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \alpha_k(\mathbf{A}\mathbf{x}_{k+1} - \mathbf{b})$  based on the assumption that  $d(\boldsymbol{\lambda})$  is differentiable  $\Leftrightarrow \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k)$  has unique minimizer  $\Leftrightarrow f$  is strictly convex
- ▶  $f$  is strictly convex is a strong assumption: that means only problems with strictly convex  $f$  can use dual ascent
- ▶ To make the algorithm works for more problems, we relax such condition.

## Revision: strong and strict convexity

- ▶ Fact: strong convexity  $\implies$  strict convexity.
- ▶ Fact: a function  $f$  is  $\rho$ -strongly convex for  $\rho > 0$  if  $f(\mathbf{x}) - \frac{\rho}{2}\|\mathbf{x}\|_2^2$  is convex.
- ▶ To relax the strict convexity condition in dual ascent, consider the following augmented Lagrangian: let  $\rho \geq 0$

$$\mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle + \frac{\rho}{2}\|\mathbf{x}\|_2^2.$$

Now  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k)$  is strongly convex ( $\implies$  strictly convex) so  $\min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k)$  has unique solution.

- ▶ But adding  $\frac{\rho}{2}\|\mathbf{x}\|_2^2$  will force  $\mathbf{x}$  to be unnecessarily small (bad for some applications). A better augmentation is  $\frac{\rho}{2}\|\mathbf{Ax} - \mathbf{b}\|_2^2$ , i.e.

$$\mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle + \frac{\rho}{2}\|\mathbf{Ax} - \mathbf{b}\|_2^2.$$

Now it will not force  $\mathbf{x}$  to be unnecessarily small but encourage  $\mathbf{Ax} = \mathbf{b}$ , exactly what we want!

- ▶ If  $\mathbf{A}$  is full rank, then  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}_k)$  is  $\rho\gamma$ -strongly convex, where  $\gamma$  is the smallest eigenvalue value of  $\mathbf{A}^\top \mathbf{A}$ .

## Augmented Lagrangian is $\rho\gamma$ -strongly convex

$$\mathcal{L}_\rho(\mathbf{x}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} \rangle + \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$$

- ▶ Let  $\gamma$  be the smallest eigenvalue of  $\mathbf{A}^\top \mathbf{A}$ .

$$\begin{aligned} \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2 &= \frac{\rho}{2} \langle \mathbf{A}^\top \mathbf{Ax}, \mathbf{x} \rangle + \text{other terms} \\ &\geq \frac{\rho}{2} \langle \gamma \mathbf{I} \mathbf{x}, \mathbf{x} \rangle + \text{other terms} \\ &= \frac{\rho\gamma}{2} \|\mathbf{x}\|_2^2 + \text{other terms} \end{aligned}$$

- ▶ Hence  $\mathcal{L}_\rho(\mathbf{x})$  is  $\rho\gamma$ -strongly convex.
- ▶ This also means that  $\mathbf{A}$  need to be full rank, otherwise  $\gamma = 0$ .
- ▶ For comparison, dual  $d(\boldsymbol{\lambda}) = -f^*(-\mathbf{A}^\top \boldsymbol{\lambda}) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle$  is concave but not strongly concave.

# Augmented Lagrangian Algorithm

- ▶ Problem:  $\min_{\mathbf{x}} f(\mathbf{x})$  s.t.  $\mathbf{Ax} = \mathbf{b}$ ,  $f$  convex<sup>6</sup> and  $\mathbf{A}$  full rank.
- ▶ Lagrangian saddle point problem:  $\min_{\mathbf{x}} \max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle$ .
- ▶ Augmented Lagrangian  $\mathcal{L}_{\rho}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} - \mathbf{b} \rangle + \frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$  is  $\rho\gamma$ -strongly convex.

---

## Algorithm 5: Augmented Lagrangian algorithm

---

```
1 Initialize  $\mathbf{x}_0, \boldsymbol{\lambda}_0$ 
2 for  $k = 1, 2, \dots$  do
3    $\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \mathcal{L}_{\rho}(\mathbf{x}, \boldsymbol{\lambda}_k)$ 
4    $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \rho(\mathbf{Ax}_{k+1} - \mathbf{b})$ 
5 end
```

---

### Remarks

- ▶ We use  $\rho$  as the gradient ascent stepsize.
- ▶ We do not see the dual  $d(\boldsymbol{\lambda}) = -f^*(-\mathbf{A}^{\top} \boldsymbol{\lambda}) - \langle \boldsymbol{\lambda}, \mathbf{b} \rangle$  because of dual ascent.
- ▶ Fact: after adding  $\frac{\rho}{2} \|\mathbf{Ax} - \mathbf{b}\|_2^2$ , the dual gradient is still  $\mathbf{Ax}_{k+1} - \mathbf{b}$ .

---

<sup>6</sup>Only convexity is assumed,  $f$  needs no strictly convex.



## Look at the gradient of dual again

- ▶ For  $\mathcal{L}$  with strictly convex  $f$ ,  
we have  $\nabla d(\boldsymbol{\lambda}) = \mathbf{A}\mathbf{x}_{k+1} - \mathbf{b}$ , where  $d(\boldsymbol{\lambda}) = \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$
- ▶ For  $\mathcal{L}_\rho$  with convex  $f$  and  $\mathbf{A}$  full rank,  
we have  $\nabla d_\rho(\boldsymbol{\lambda}) = \mathbf{A}\mathbf{x}_{k+1} - \mathbf{b}$ , where  $d_\rho(\boldsymbol{\lambda}) = \min_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\lambda})$
- ▶  $\nabla d(\boldsymbol{\lambda}) = \nabla d_\rho(\boldsymbol{\lambda}) = \mathbf{A}\mathbf{x}_{k+1} - \mathbf{b}$  and being  $\|\mathbf{A}\|_2$ -Lipschitz
- ▶ For  $\nabla d_\rho(\boldsymbol{\lambda})$

$$\|\nabla d_\rho(\boldsymbol{\lambda}) - \nabla d_\rho(\boldsymbol{\lambda}')\|_2 \leq \|\mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{x}'\|_2 \leq \frac{1}{\rho} \|\boldsymbol{\lambda} - \boldsymbol{\lambda}'\|_2.$$

Prove  $\|\nabla d_\rho(\boldsymbol{\lambda}) - \nabla d_\rho(\boldsymbol{\lambda}')\|_2 \leq \|\mathbf{Ax} - \mathbf{Ax}'\|_2 \leq \frac{1}{\rho} \|\boldsymbol{\lambda} - \boldsymbol{\lambda}'\|_2$

► Let  $\mathbf{x} = \underset{\mathbf{x}}{\operatorname{argmin}} \mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\lambda})$  and  $\mathbf{x}' = \underset{\mathbf{x}}{\operatorname{argmin}} \mathcal{L}_\rho(\mathbf{x}, \boldsymbol{\lambda}')$ .

► By subgradient 1st-order optimality,

$$\begin{aligned} \mathbf{0} &\in \partial f(\mathbf{x}) + \mathbf{A}^\top \boldsymbol{\lambda} + \rho \mathbf{A}^\top (\mathbf{Ax} - \mathbf{b}), & \mathbf{0} &\in \partial f(\mathbf{x}') + \mathbf{A}^\top \boldsymbol{\lambda}' + \rho \mathbf{A}^\top (\mathbf{Ax}' - \mathbf{b}) \\ \iff -\mathbf{A}^\top \boldsymbol{\lambda} - \rho \mathbf{A}^\top (\mathbf{Ax} - \mathbf{b}) &\in \partial f(\mathbf{x}), & \mathbf{A}^\top \boldsymbol{\lambda}' - \rho \mathbf{A}^\top (\mathbf{Ax}' - \mathbf{b}) &\in \partial f(\mathbf{x}') \end{aligned}$$

► By the fact that subgradient  $\partial f$  is a monotone operator:  $\langle \partial f(\mathbf{x}) - \partial f(\mathbf{x}'), \mathbf{x} - \mathbf{x}' \rangle \geq 0$ , we have

$$\left\langle -\mathbf{A}^\top \boldsymbol{\lambda} - \rho \mathbf{A}^\top (\mathbf{Ax} - \mathbf{b}) + \mathbf{A}^\top \boldsymbol{\lambda}' + \rho \mathbf{A}^\top (\mathbf{Ax}' - \mathbf{b}), \mathbf{x} - \mathbf{x}' \right\rangle \geq 0.$$

$$\left\langle -(\boldsymbol{\lambda} - \boldsymbol{\lambda}') - \rho(\mathbf{Ax} - \mathbf{Ax}'), \mathbf{Ax} - \mathbf{Ax}' \right\rangle \geq 0.$$

$$\left\langle -(\boldsymbol{\lambda} - \boldsymbol{\lambda}'), \mathbf{Ax} - \mathbf{Ax}' \right\rangle \geq \rho \|\mathbf{Ax} - \mathbf{Ax}'\|_2^2$$

► Cauchy-Schwarz inequality

$$\|\boldsymbol{\lambda} - \boldsymbol{\lambda}'\|_2 \|\mathbf{Ax} - \mathbf{Ax}'\|_2 \geq \left| \left\langle -(\boldsymbol{\lambda} - \boldsymbol{\lambda}'), \mathbf{Ax} - \mathbf{Ax}' \right\rangle \right| \geq \left\langle -(\boldsymbol{\lambda} - \boldsymbol{\lambda}'), \mathbf{Ax} - \mathbf{Ax}' \right\rangle \geq \rho \|\mathbf{Ax} - \mathbf{Ax}'\|_2^2.$$

First term and last term gives  $\|\mathbf{Ax} - \mathbf{Ax}'\|_2 \leq \frac{1}{\rho} \|\boldsymbol{\lambda} - \boldsymbol{\lambda}'\|_2$ .

# Augmented Lagrangian for composite function

- ▶ Consider problem

$$(\mathcal{P}_0) : \min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{A}\mathbf{x})$$

where both  $f, g$  are convex functions and  $\mathbf{A}$  is full rank.

- ▶ Rewrite  $\mathcal{P}_0$  as

$$(\mathcal{P}_1) : \min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{y}) \text{ s.t. } \mathbf{A}\mathbf{x} = \mathbf{y}$$

- ▶ Solve  $\mathcal{P}_1$  by Augmented Lagrangian

$$\min_{\mathbf{x}, \mathbf{y}} \max_{\boldsymbol{\lambda}} \mathcal{L}_{\rho}(\mathbf{x}, \mathbf{y}, \boldsymbol{\lambda}) = f(\mathbf{x}) + g(\mathbf{y}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{x} - \mathbf{y} \rangle + \frac{\rho}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2.$$

## Augmented Lagrangian for general problem

- ▶ Consider problem

$$(\mathcal{P}') : \min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{y}) \text{ s.t. } \mathbf{Ax} + \mathbf{By} + \mathbf{c} = \mathbf{0}$$

where both  $f, g$  are convex functions and  $\mathbf{A}, \mathbf{B}$  are full rank.

- ▶ Solve  $\mathcal{P}'$  by Augmented Lagrangian

$$\min_{\mathbf{x}, \mathbf{y}} \max_{\boldsymbol{\lambda}} \mathcal{L}_{\rho}(\mathbf{x}, \mathbf{y}, \boldsymbol{\lambda}) = f(\mathbf{x}) + g(\mathbf{y}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} + \mathbf{By} + \mathbf{c} \rangle + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{By} + \mathbf{c}\|_2^2.$$

---

### Algorithm 6: Augmented Lagrangian algorithm

---

- 1 Initialize  $\mathbf{x}_0, \mathbf{y}_0, \boldsymbol{\lambda}_0$
  - 2 **for**  $k = 1, 2, \dots$  **do**
  - 3      $(\mathbf{x}_{k+1}, \mathbf{y}_{k+1}) = \underset{\mathbf{x}, \mathbf{y}}{\operatorname{argmin}} \mathcal{L}_{\rho}(\mathbf{x}, \mathbf{y}, \boldsymbol{\lambda}_k)$  joint-minimization
  - 4      $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \alpha_k (\mathbf{Ax}_{k+1} + \mathbf{By}_{k+1} + \mathbf{c})$  dual ascent
  - 5 **end**
-

# How ADMM is discovered

---

**Algorithm 7:** Augmented Lagrangian algorithm

---

```
1 Initialize  $\mathbf{x}_0, \mathbf{y}_0, \boldsymbol{\lambda}_0$ 
2 for  $k = 1, 2, \dots$  do
3    $(\mathbf{x}_{k+1}, \mathbf{y}_{k+1}) = \underset{\mathbf{x}, \mathbf{y}}{\operatorname{argmin}} \mathcal{L}_\rho(\mathbf{x}, \mathbf{y}, \boldsymbol{\lambda}_k)$  joint-minimization
4    $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \alpha_k (\mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}\mathbf{y}_{k+1} + \mathbf{c})$  dual ascent
5 end
```

---

- ▶ Glowinski and Marrocco (1976) suggested to split  $\mathbf{x}, \mathbf{y}$  in a coordinate descent fashion, i.e.

$$\begin{aligned}\mathbf{x}_{k+1} &= \underset{\mathbf{x}, \mathbf{y}}{\operatorname{argmin}} \mathcal{L}_\rho(\mathbf{x}, \mathbf{y}_k, \boldsymbol{\lambda}_k) \\ \mathbf{y}_{k+1} &= \underset{\mathbf{x}, \mathbf{y}}{\operatorname{argmin}} \mathcal{L}_\rho(\mathbf{x}_{k+1}, \mathbf{y}, \boldsymbol{\lambda}_k)\end{aligned}$$

they found that this gives the best performance.

- ▶ **This is the moment ADMM was born.** It was discovered as an empirical heuristic with no theory support.
- ▶ Jonathan Eckstein: **ADMM does not approximate the ALM.**

# ADMM algorithm

► For problem

$$(\mathcal{P}') : \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + g(\mathbf{y}) \text{ s.t. } \mathbf{Ax} + \mathbf{By} + \mathbf{c} = \mathbf{0}$$

where both  $f, g$  are convex functions and  $\mathbf{A}, \mathbf{B}$  are full rank.

---

## Algorithm 8: ADMM

---

```
1 Initialize  $\mathbf{x}_0, \mathbf{y}_0, \boldsymbol{\lambda}_0, \rho$ 
2 for  $k = 1, 2, \dots$  do
3    $\mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x}} \mathcal{L}_\rho(\mathbf{x}, \mathbf{y}_k, \boldsymbol{\lambda}_k)$ 
4    $\mathbf{y}_{k+1} = \operatorname{argmin}_{\mathbf{y}} \mathcal{L}_\rho(\mathbf{x}_{k+1}, \mathbf{y}, \boldsymbol{\lambda}_k)$ 
5    $\boldsymbol{\lambda}_{k+1} = \boldsymbol{\lambda}_k + \rho(\mathbf{Ax}_{k+1} + \mathbf{By}_{k+1} + \mathbf{c})$ 
6 end
```

---

Extra slide 1: what's the dual for  $\min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + g(\mathbf{y})$  s.t.  $\mathbf{Ax} + \mathbf{By} + \mathbf{c} = \mathbf{0}$ ?

- ▶ Original problem

$$\min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{y}) \text{ s.t. } \mathbf{Ax} + \mathbf{By} + \mathbf{c} = \mathbf{0} \quad f, g \text{ convex and } \mathbf{A}, \mathbf{B} \text{ full rank.}$$

- ▶ Augmented Lagrangian saddle point problem

$$\min_{\mathbf{x}, \mathbf{y}} \max_{\boldsymbol{\lambda}} \mathcal{L}_{\rho}(\mathbf{x}, \mathbf{y}, \boldsymbol{\lambda}) = f(\mathbf{x}) + g(\mathbf{y}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} + \mathbf{By} + \mathbf{c} \rangle + \frac{\rho}{2} \|\mathbf{Ax} + \mathbf{By} + \mathbf{c}\|_2^2.$$

- ▶ Un-augmented Lagrangian saddle point problem ( $\rho = 0$ )

$$\min_{\mathbf{x}, \mathbf{y}} \max_{\boldsymbol{\lambda}} \mathcal{L}(\mathbf{x}, \mathbf{y}, \boldsymbol{\lambda}) = f(\mathbf{x}) + g(\mathbf{y}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} + \mathbf{By} + \mathbf{c} \rangle.$$

- ▶ Dual on  $\mathcal{L}$  (not  $\mathcal{L}_{\rho}$ )

$$\begin{aligned} d(\boldsymbol{\lambda}) &:= \min_{\mathbf{x}, \mathbf{y}} \mathcal{L}(\mathbf{x}, \mathbf{y}, \boldsymbol{\lambda}) \\ &= \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + g(\mathbf{y}) + \langle \boldsymbol{\lambda}, \mathbf{Ax} \rangle + \langle \boldsymbol{\lambda}, \mathbf{By} \rangle + \langle \boldsymbol{\lambda}, \mathbf{c} \rangle \\ &= - \left( \max_{\mathbf{x}, \mathbf{y}} -f(\mathbf{x}) - g(\mathbf{y}) - \langle \boldsymbol{\lambda}, \mathbf{Ax} \rangle - \langle \boldsymbol{\lambda}, \mathbf{By} \rangle - \langle \boldsymbol{\lambda}, \mathbf{c} \rangle \right) \\ &= - \left( \left[ \max_{\mathbf{x}} -f(\mathbf{x}) - \langle \boldsymbol{\lambda}, \mathbf{Ax} \rangle \right] + \left[ \max_{\mathbf{y}} -g(\mathbf{y}) - \langle \boldsymbol{\lambda}, \mathbf{By} \rangle \right] + \left[ -\langle \boldsymbol{\lambda}, \mathbf{c} \rangle \right] \right) \\ &= - \left( \left[ \max_{\mathbf{x}} -f(\mathbf{x}) + \langle -\mathbf{A}^{\top} \boldsymbol{\lambda}, \mathbf{x} \rangle \right] + \left[ \max_{\mathbf{y}} -g(\mathbf{y}) + \langle -\mathbf{B}^{\top} \boldsymbol{\lambda}, \mathbf{y} \rangle \right] + \left[ -\langle \boldsymbol{\lambda}, \mathbf{c} \rangle \right] \right) \\ &= - \left( \left[ f^*(-\mathbf{A}^{\top} \boldsymbol{\lambda}) \right] + \left[ g^*(\mathbf{B}^{\top} \boldsymbol{\lambda}) \right] + \left[ -\langle \boldsymbol{\lambda}, \mathbf{c} \rangle \right] \right) \\ &= -f^*(-\mathbf{A}^{\top} \boldsymbol{\lambda}) - g^*(\mathbf{B}^{\top} \boldsymbol{\lambda}) + \langle \boldsymbol{\lambda}, \mathbf{c} \rangle \end{aligned}$$

Extra slide 2: why  $\nabla_{\lambda} d(\lambda) = \mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}\mathbf{y}_{k+1} + \mathbf{c}$ ?

► Recall subgradient for **concave** function:

$$\mathbf{q} \in \partial d(\lambda) \iff d(\gamma) \leq d(\lambda) + \langle \mathbf{q}, \gamma - \lambda \rangle \quad \forall \gamma.$$

► By definition of dual

$$\begin{aligned} d(\gamma) &= \min_{\mathbf{x}, \mathbf{y}} \mathcal{L}(\mathbf{x}, \mathbf{y}, \gamma) \\ &= \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + g(\mathbf{y}) + \langle \gamma, \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{c} \rangle \\ &= \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + g(\mathbf{y}) + \langle \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{c}, \gamma \rangle \\ &= \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + g(\mathbf{y}) + \langle \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{c}, \gamma + \lambda - \lambda \rangle \\ &= \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{x}) + g(\mathbf{y}) + \langle \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{c}, \lambda \rangle + \langle \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \mathbf{c}, \gamma - \lambda \rangle \\ &\stackrel{*}{\leq} f(\mathbf{x}_{k+1}) + g(\mathbf{y}_{k+1}) + \langle \mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}\mathbf{y}_{k+1} + \mathbf{c}, \lambda \rangle + \langle \mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}\mathbf{y}_{k+1} + \mathbf{c}, \gamma - \lambda \rangle \\ &\stackrel{**}{\leq} d(\lambda) + \langle \mathbf{A}\mathbf{x}_{k+1} + \mathbf{B}\mathbf{y}_{k+1} + \mathbf{c}, \gamma - \lambda \rangle \end{aligned}$$

\*: by definition of argmin.

\*\* : by definition of  $\lambda$  in the ADMM algorithm and by the definition of dual.



# Early history of ADMM

- ▶ Douglas-Rachford splitting<sup>7</sup> 1956
  - ▶ Another way to derive ADMM, much earlier than ADMM.
  - ▶ Requires knowledge of convex analysis and monotone operator theory.
  - ▶ It was first used to solve heat equation.
- ▶ Glowinski and Marrocco suggested to split  $x, y$  in a coordinate descent fashion 1976
  - ▶ The moment ADMM was born.
  - ▶ Pure heuristic, no theory.
- ▶ Lions and Mercier 1979
- ▶ Gabay 1983
- ▶ Eckstein and Bertsekas 1992

---

<sup>7</sup>J. Douglas and H.H. Rachford, "On the numerical solution of heat conduction problems in two and three space variables," Transactions of the American Mathematical Society 82 (1956) 421-439.

# A slide from Jonathan Eckstein (early contributor of ADMM)

## Getting the History Right

- Douglas and Rachford had a different representation of the operations in their method, but equivalent
- However, the original Douglas-Rachford publication was only for linear operators
  - Applied to very specific linear operators related to the discretized 2-D heat equation
- Lions and Mercier (1979) generalized the idea from linear maps to general monotone set-valued maps (but kept the name)
- Gabay (1983) showed that the ADMM is just this idea applied to the dual of  $\min f(x) + g(Mx)$
- The composition-of-nonexpansive maps interpretation may first be found (as an aside) in Lawrence & Spingarn (1987)
- E & Bertsekas (1992) contains some equivalent analysis and exploits the relationship with the proximal point algorithm to derive approximate and over-relaxed versions

## Last page - summary

- ▶ The history of ADMM
  - ▶ Lagrangian multiplier
  - ▶ Augmented Lagrangian
  - ▶ ADMM: splitting primal variables
- ▶ Not discussed / discuss next
  - ▶ Douglas-Rachford splitting
  - ▶ Convergence of ADMM

to handle constrained program  
to relax  $f$  strictly convex  
discovered empirically

End of document