



# SOD314. Cooperative Optimization for data science

**Andrea Simonetto** 

Academic year 2022/2023

#### General info

- 5 classes (2h class, 1h30 project)
- All the info are on moodle
- 2h written Exam: 28/03, with material
- Grade: a function of Exam, TP, and optional homework along the way (it's
  the average between exam and TP plus extra points due to homework, but
  the result is cropped so that it cannot be more or less than 3 points of the
  exam grade)
- Extra points due to homework: max 3 points in total
- Contacts
   Andrea SIMONETTO
   andrea.simonetto@ensta-paris.fr,

#### Content

- Introduction
- ② Distributed Optimization
- Intermezzo: Stochastic gradient
- Federated Learning
- Differential Privacy

Material: Slides, and suggested references to books/articles as we go along.

# Part I Introduction

## Some recap of useful notions

 We look at continuous optimization problems (like in OPT201, OPT202), i.e.,

$$\min_{\mathbf{x} \in X \subseteq \mathbb{R}^n} f(\mathbf{x}), \tag{1}$$

where the set X is closed and convex,  $f: \mathbf{R}^n \to \mathbf{R}$  is a convex function and the minimum is attained for a  $\mathbf{x} \in X$ .

#### Some recap of useful notions

 We look at continuous optimization problems (like in OPT201, OPT202), i.e.,

$$\min_{\mathbf{x} \in X \subset \mathbb{R}^n} f(\mathbf{x}),\tag{1}$$

where the set X is closed and convex,  $f: \mathbf{R}^n \to \mathbf{R}$  is a convex function and the minimum is attained for a  $\mathbf{x} \in X$ .

- What you need to remember from past courses will be reviewed as we go along, but mainly:
  - Optimality conditions (KKT), which are here necessary and sufficient (under a Slater's constraint qualification assumption)
  - ► How to derive dual problems (Lagrangian function, dual function, problems)
  - Subgradients
  - Algorithms: How to derive first-order algorithms (e.g., gradient descent) from the optimality conditions, and how to prove their convergence and convergence rate
  - ► Some basic linear algebra: eigenvalues, singular value decomposition, solution of linear systems, etc..

### The optimization problem of interest

• Said so, the problem we will look at in this course has the form,

(P) 
$$\min_{\mathbf{x} \in X \subseteq \mathbb{R}^n} \sum_{i=1}^N f_i(\mathbf{x}), \tag{2}$$

where the set X is closed and convex,  $f_i: \mathbf{R}^n \to \mathbf{R}$  is a convex function for all i's and the minimum is attained for a  $\mathbf{x} \in X$ .

Here N are the number of agents/players/nodes/sub-systems/etc.. that cooperate to solve the optimization problem.

• Cooperative least-squares. Imagine N users collect noisy measurements  $y_i \in \mathbf{R}^n$  about a quantity  $\mathbf{x} \in X$ . A way to estimate the true value for  $\mathbf{x}$  is to set up a least-squares problem as

$$\min_{\mathbf{x} \in X} \sum_{i=1}^{N} \|\mathbf{x} - \mathbf{y}_i\|^2 = \sum_{i=1}^{N} f_i(\mathbf{x})$$

• Cooperative least-squares. Imagine N users collect noisy measurements  $y_i \in \mathbf{R}^n$  about a quantity  $\mathbf{x} \in X$ . A way to estimate the true value for  $\mathbf{x}$  is to set up a least-squares problem as

$$\min_{\mathbf{x} \in X} \sum_{i=1}^{N} \|\mathbf{x} - \mathbf{y}_i\|^2 = \sum_{i=1}^{N} f_i(\mathbf{x})$$

• Cooperative linear model training. Imagine N users collect input-output pairs  $\mathbf{w}_i \in \mathbf{R}^{n-1}$ ,  $\mathbf{y}_i \in \mathbf{R}$  (e.g., features and labels), and they want to train a global model with weights  $\mathbf{x} \in X \subseteq \mathbf{R}^n$ ,  $\mathbf{x} = [\theta \in \mathbf{R}^{n-1}, c \in \mathbf{R}]$ . Let the local input-output mapping being affine,

$$\mathbf{y}_i = \theta^{\top} \mathbf{w}_i + c, \quad \forall i,$$

then the training problem can be written as

$$\min_{\mathbf{x} \in X} \sum_{i=1}^{N} \|\mathbf{y}_i - (\theta^\top \mathbf{w}_i + c)\|^2 = \sum_{i=1}^{N} f_i(\mathbf{x})$$

Cooperative network problems. Take a sensor network of N sensors. you
want to compute the localization of the whole network based on pair-wise
distance measurements (e.g., sensors can be cars, or people with their
phones). Then each measurement is

$$m_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\| + \text{noise}, \quad \mathbf{x}_i \in \mathbf{R}^2 \text{ is the position of node } i.$$

You have E pair-wise measurements. Then you can write the problem as

$$\min_{\mathbf{x} \in X \subset \mathsf{R}^{2E}} \sum_{i,j}^{E} (m_{ij} - \|\mathbf{x}_i - \mathbf{x}_j\|)^2 = \sum_{k=1}^{E} f_k(\mathbf{x})$$

Cooperative network problems. Take a sensor network of N sensors. you
want to compute the localization of the whole network based on pair-wise
distance measurements (e.g., sensors can be cars, or people with their
phones). Then each measurement is

$$m_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\| + \text{noise}, \quad \mathbf{x}_i \in \mathbf{R}^2 \text{ is the position of node } i.$$

You have E pair-wise measurements. Then you can write the problem as

$$\min_{\mathbf{x} \in X \subset \mathsf{R}^{2E}} \sum_{i,j}^{E} (m_{ij} - \|\mathbf{x}_i - \mathbf{x}_j\|)^2 = \sum_{k=1}^{E} f_k(\mathbf{x})$$

• Cooperative consensus. You have N robots at different locations  $s_i$  and moving at different speeds  $v_i$ , and you want to find the best position in space for the fastest rendez-vous:

$$\min_{\mathbf{x} \in X \subset \mathbb{R}^3} \sum_{i}^{N} \frac{\|\mathbf{s}_i - \mathbf{x}\|}{v_i} = \sum_{i=1}^{N} f_i(\mathbf{x})$$

- Large-scale problems: *N* is big or *n* is big and the problem cannot be solved (not stored) locally. You need to solve it on separate machines, or iteratively.
  - ▶ If n is big, usually we talk about parallel methods;
  - ightharpoonup If N is big, usually we talk about distributed methods.

- Large-scale problems: *N* is big or *n* is big and the problem cannot be solved (not stored) locally. You need to solve it on separate machines, or iteratively.
  - ▶ If *n* is big, usually we talk about parallel methods;
  - ▶ If *N* is big, usually we talk about distributed methods.
- Distributed localization: the data that builds  $f_i$  is stored in separate machines (think: mobile phones). Then, you have communication issues and latencies..
  - Are we communicating to a server? (Cloud-based)
  - Are we communicating to each other (Peer-to-peer)
  - Common issues: asynchronicity, bi-directionality, packet-losses, communication overhead,...
  - Here: graph theory to the rescue!

The SneakerNet paradox. What is the fastest way to send large data sets?

- Large-scale problems: *N* is big or *n* is big and the problem cannot be solved (not stored) locally. You need to solve it on separate machines, or iteratively.
  - ▶ If *n* is big, usually we talk about parallel methods;
  - ▶ If *N* is big, usually we talk about distributed methods.
- Distributed localization: the data that builds  $f_i$  is stored in separate machines (think: mobile phones). Then, you have communication issues and latencies..
  - Are we communicating to a server? (Cloud-based)
  - ► Are we communicating to each other (Peer-to-peer)
  - Common issues: asynchronicity, bi-directionality, packet-losses, communication overhead,...
  - ▶ Here: graph theory to the rescue!

**The SneakerNet paradox.** What is the fastest way to send large data sets? FedEx

- Large-scale problems: *N* is big or *n* is big and the problem cannot be solved (not stored) locally. You need to solve it on separate machines, or iteratively.
  - ▶ If *n* is big, usually we talk about parallel methods;
  - ▶ If *N* is big, usually we talk about distributed methods.
- Distributed localization: the data that builds  $f_i$  is stored in separate machines (think: mobile phones). Then, you have communication issues and latencies..
  - Are we communicating to a server? (Cloud-based)
  - Are we communicating to each other (Peer-to-peer)
  - Common issues: asynchronicity, bi-directionality, packet-losses, communication overhead,...
  - ▶ Here: graph theory to the rescue!
- Data and functions are private: you don't want to disclose  $f_i$  or the local decision to the other players. How do you do that? (Example: training language models based on private text messages on your phone)

 In this course, we will look at algorithms to tackle the challenges above while solving problem (P)

- In this course, we will look at algorithms to tackle the challenges above while solving problem (P)
- The algorithms will be of the first-order kind (more later)

- In this course, we will look at algorithms to tackle the challenges above while solving problem (P)
- The algorithms will be of the first-order kind (more later)
- We divide the course in three parts (distributed, federated, private)

# Part II Basics

#### Some definitions

#### Definition 1 (Convex functions)

A function  $f: X \subseteq \mathbf{R}^n \to \mathbf{R}$  is convex iff X is convex and

(C1) 
$$\forall \mathbf{x}, \mathbf{y} \in X, \lambda \in [0, 1]: f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \le \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}).$$

Multiple definitions exists, for example:

$$(C1) + f \in \mathcal{C}^{1}(X) \quad \Longleftrightarrow \quad \forall \mathbf{x}, \mathbf{y} \in X, \ f(\mathbf{x}) \ge f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \quad (3)$$

$$(C1) + f \in \mathcal{C}^{2}(X) \quad \Longleftrightarrow \quad \forall \mathbf{x}, \mathbf{y} \in X, \, \nabla^{2} f(\mathbf{x}) \succeq 0$$

$$(4)$$

**Proof.** Homework.

#### Some definitions

#### Definition 1 (Strongly convex functions)

A convex function  $f: X \subseteq \mathbf{R}^n \to \mathbf{R}$  is *m*-strongly convex iff

(SC) 
$$f(\mathbf{x}) - \frac{m}{2} \|\mathbf{x}\|^2$$
 is convex.

f doesn't need to be differentiable! Multiple definitions exists, for example:

$$(SC) + f \in \mathcal{C}^{1}(X) \iff \forall \mathbf{x}, \mathbf{y} \in X, \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \ge m \|\mathbf{x} - \mathbf{y}\|^{2}$$
(3)  
$$(SC) + f \in \mathcal{C}^{2}(X) \iff \forall \mathbf{x}, \mathbf{y} \in X, \nabla^{2} f(\mathbf{x}) \succeq m I_{n}$$
(4)

**Proof.** Homework

#### Some definitions

#### Definition 1 (Smooth functions)

A convex function  $f: X \subseteq \mathbf{R}^n \to \mathbf{R}$  is L-smooth iff

(LC) 
$$\frac{L}{2} \|\mathbf{x}\|^2 - f(\mathbf{x})$$
 is convex.

Important (LC)  $\implies f \in C^1(X)$ !

Multiple definitions exists, for example:

(LC) 
$$\iff \forall \mathbf{x}, \mathbf{y} \in X, \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \le L\|\mathbf{x} - \mathbf{y}\|$$
 (3)

(LC) 
$$\iff$$
  $\forall x, y \in X$ ,

$$\frac{1}{l} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 \le \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \quad (4)$$

$$(LC) + f \in \mathcal{C}^{2}(X) \quad \Longleftrightarrow \quad \forall \mathbf{x}, \mathbf{y} \in X, \ 0 \leq \nabla^{2} f(\mathbf{x}) \leq L I_{n}$$
 (5)

#### Proof. Homework.

# Gradient algorithm

Consider solving  $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ .

Rem: the gradient is an ascent direction, and a first-order method

# Gradient algorithm

Consider solving  $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ .

Rem: the gradient is an ascent direction, and a first-order method The simplest scheme, let  $\alpha_k > 0$ :

- Start with  $\mathbf{x}_0 \in \mathbf{R}^n$
- Iterate  $\mathbf{x}_{k+1} = \mathbf{x}_k \alpha_k \nabla f(\mathbf{x}_k)$ ,  $k = 0, 1, \dots$

There are different methods to choose  $\alpha_k$  (either a priori or online):

- Constant:  $\alpha_k = \alpha$
- Vanishing:  $\alpha_k = \frac{\alpha}{\sqrt{k+1}}$
- ...

Why the difference? Convergence

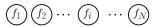
# Gradient algorithm: convergence recap

- Recap from OPT201-202: gradient converges in a well-defined sense provided the step size  $\alpha_k$  is chosen appropriately.
- Depending on the functional class, the gradient algorithm behaves differently with the number of iterations *t*.
- The table below recaps the main results that one can expect:

Туре	just convex	smooth	smooth+strongly convex
		First-order	
Convergence guarantee	$\ \nabla f\ , f - f^*$	$f - f^*$	$  x-x^*  $
Convex	$O(1/\sqrt{t}), O(1/t)$	$O(1/t)$ , $O(1/t^2)$	$O( ho^t)$

**Homework:** Revise your theory

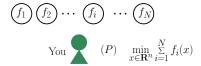
• We want to solve problem (P) but N is too big. How do we do?





You 
$$(P) \quad \min_{x \in \mathbf{R}^n} \sum_{i=1}^{N} f_i(x)$$

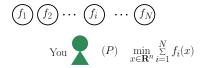
• We want to solve problem (P) but N is too big. How do we do?



• First possibility: you ask one device to apply a local gradient, then you get back the result and you ask a second device..

#### Incremental/ Gauss-Seidel Gradient

• We want to solve problem (P) but N is too big. How do we do?

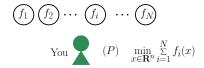


• First possibility: you ask one device to apply a local gradient, then you get back the result and you ask a second device..

#### Incremental/ Gauss-Seidel Gradient

▶ Start with  $\mathbf{x}_0 \in \mathbf{R}^n$ 

• We want to solve problem (P) but N is too big. How do we do?

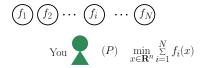


• First possibility: you ask one device to apply a local gradient, then you get back the result and you ask a second device..

#### Incremental/ Gauss-Seidel Gradient

- ▶ Start with  $\mathbf{x}_0 \in \mathbf{R}^n$
- ▶ Iterate: ask device  $i: \mathbf{x}_{k+1} = \mathbf{x}_k \alpha_k \nabla f_i(\mathbf{x}_k), \quad k = 0, 1, \dots$

• We want to solve problem (P) but N is too big. How do we do?



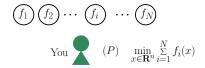
• First possibility: you ask one device to apply a local gradient, then you get back the result and you ask a second device..

#### Incremental/ Gauss-Seidel Gradient

- ▶ Start with  $\mathbf{x}_0 \in \mathbf{R}^n$
- ▶ Iterate: ask device *i*:  $\mathbf{x}_{k+1} = \mathbf{x}_k \alpha_k \nabla f_i(\mathbf{x}_k)$ , k = 0, 1, ...
- Second possibility: you ask every device to apply a local gradient, then you
  get back the result and ??

#### Parallel/ Jacobi Gradient

• We want to solve problem (P) but N is too big. How do we do?



• First possibility: you ask one device to apply a local gradient, then you get back the result and you ask a second device..

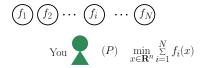
#### Incremental/ Gauss-Seidel Gradient

- ▶ Start with  $\mathbf{x}_0 \in \mathbf{R}^n$
- ▶ Iterate: ask device *i*:  $\mathbf{x}_{k+1} = \mathbf{x}_k \alpha_k \nabla f_i(\mathbf{x}_k)$ , k = 0, 1, ...
- Second possibility: you ask every device to apply a local gradient, then you
  get back the result and ??

#### Parallel/ Jacobi Gradient

▶ Start with  $\mathbf{x}_0 \in \mathbf{R}^n$ 

• We want to solve problem (P) but N is too big. How do we do?



• First possibility: you ask one device to apply a local gradient, then you get back the result and you ask a second device..

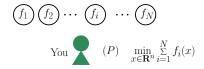
#### Incremental/ Gauss-Seidel Gradient

- ▶ Start with  $\mathbf{x}_0 \in \mathbf{R}^n$
- ▶ Iterate: ask device i:  $\mathbf{x}_{k+1} = \mathbf{x}_k \alpha_k \nabla f_i(\mathbf{x}_k)$ , k = 0, 1, ...
- Second possibility: you ask every device to apply a local gradient, then you
  get back the result and ??

#### Parallel/ Jacobi Gradient

- ▶ Start with  $\mathbf{x}_0 \in \mathbf{R}^n$
- ▶ Iterate: ask every devices:  $\mathbf{x}_{k+1}^i = \mathbf{x}_k \alpha_k \nabla f_i(\mathbf{x}_k)$ ,  $\forall i$

• We want to solve problem (P) but N is too big. How do we do?



• First possibility: you ask one device to apply a local gradient, then you get back the result and you ask a second device..

#### Incremental/ Gauss-Seidel Gradient

- ▶ Start with  $\mathbf{x}_0 \in \mathbf{R}^n$
- ▶ Iterate: ask device  $i: \mathbf{x}_{k+1} = \mathbf{x}_k \alpha_k \nabla f_i(\mathbf{x}_k), \quad k = 0, 1, \dots$
- Second possibility: you ask every device to apply a local gradient, then you
  get back the result and ??

#### Parallel/ Jacobi Gradient

- ▶ Start with  $x_0 \in \mathbf{R}^n$
- ▶ Iterate: ask every devices:  $\mathbf{x}_{k+1}^i = \mathbf{x}_k \alpha_k \nabla f_i(\mathbf{x}_k)$ ,  $\forall i$
- $x_{k+1} = ?$

# Properties and convergence of incremental gradient

• The data that generates  $f_i$  stays private, and if you have enough time you don't care about latencies (etc..): the scheme is quite robust

# Properties and convergence of incremental gradient

- The data that generates  $f_i$  stays private, and if you have enough time you don't care about latencies (etc..): the scheme is quite robust
- Also known in the ML community as: online back-propagation, finite sum stochastic gradient descent (careful: step size is called the learning rate)

#### Theorem 2

Consider problem (P) for a m-strongly convex and L-smooth function f. Consider the incremental gradient method with step size  $\alpha_k$ . Assume that  $\|\nabla f_i(\mathbf{x}) - \sum_i \nabla f_i(\mathbf{x})\| \le G$  for all  $\mathbf{x} \in \mathbf{R}^n$ , choose  $\alpha_k < 2/L$ , and define the quantity,  $\rho_k = \max\{|1 - \alpha_k m|, |1 - \alpha_k L|\} < 1$ .

Then convergence of the incremental gradient goes as

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \rho_k \|\mathbf{x}_k - \mathbf{x}^*\| + \alpha_k G.$$

#### Corollary 3

If 
$$\alpha_k = 1/k^s$$
,  $0 < s < 1$  then,  $\|\mathbf{x}_k - \mathbf{x}^*\| \le O(1/k^s)$ .

#### Proof

• Define  $f(\mathbf{x}) = \sum_i f_i(\mathbf{x})$ . Write

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \le \|\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) - \mathbf{x}^* + \alpha_k \nabla f(\mathbf{x}^*)\| + \alpha_k \|\nabla f_i(\mathbf{x}_k) - \nabla f(\mathbf{x}_k)\|$$

## Proof

• Define  $f(\mathbf{x}) = \sum_i f_i(\mathbf{x})$ . Write

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \|\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) - \mathbf{x}^* + \alpha_k \nabla f(\mathbf{x}^*)\| + \alpha_k \|\nabla f_i(\mathbf{x}_k) - \nabla f(\mathbf{x}_k)\|$$

Use strongly convex and smoothness property to say

$$\|\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) - \mathbf{x}^* + \alpha_k \nabla f(\mathbf{x}^*)\| \le \rho_k \|\mathbf{x}_k - \mathbf{x}^*\|$$

## Proof

• Define  $f(\mathbf{x}) = \sum_i f_i(\mathbf{x})$ . Write

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \|\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) - \mathbf{x}^* + \alpha_k \nabla f(\mathbf{x}^*)\| + \alpha_k \|\nabla f_i(\mathbf{x}_k) - \nabla f(\mathbf{x}_k)\|$$

Use strongly convex and smoothness property to say

$$\|\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) - \mathbf{x}^* + \alpha_k \nabla f(\mathbf{x}^*)\| \le \rho_k \|\mathbf{x}_k - \mathbf{x}^*\|$$

• The thesis follows. The corollary is a bit more involved but not impossible, see for example: arXiv:1510.08562

### Proof

• Define  $f(\mathbf{x}) = \sum_i f_i(\mathbf{x})$ . Write

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \|\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) - \mathbf{x}^* + \alpha_k \nabla f(\mathbf{x}^*)\| + \alpha_k \|\nabla f_i(\mathbf{x}_k) - \nabla f(\mathbf{x}_k)\|$$

Use strongly convex and smoothness property to say

$$\|\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k) - \mathbf{x}^* + \alpha_k \nabla f(\mathbf{x}^*)\| \le \rho_k \|\mathbf{x}_k - \mathbf{x}^*\|$$

- The thesis follows. The corollary is a bit more involved but not impossible, see for example: arXiv:1510.08562
- A lot of research is devoted in lifting the assumptions (strong convexity, smoothness), and relaxing the assumption on *G*, but the basic ideas are still valid: you have a trade-off between speed and error.

## Convergence of parallel gradient?

• It seems that it could be a faster algorithm, but how can we combine different updates  $x_{k+1}^i$ ? Can we just average them?

# Convergence of parallel gradient?

- It seems that it could be a faster algorithm, but how can we combine different updates  $x_{k+1}^i$ ? Can we just average them?
- The scheme is less robust to asynchronicity, package drops, etc...

## Convergence of parallel gradient?

- It seems that it could be a faster algorithm, but how can we combine different updates  $x_{k+1}^i$ ? Can we just average them?
- The scheme is less robust to asynchronicity, package drops, etc..
- The scheme is the starting point of distributed optimization (next!)

# Part III Distributed optimization

• We will need some basic graph theory to devise methods to combine different updates in the parallel scheme and analyze communication issues

- We will need some basic graph theory to devise methods to combine different updates in the parallel scheme and analyze communication issues
- $\bullet$  Consider a graph  ${\cal G}$  as a collection of vertices  ${\cal V},$  edges  ${\cal E},$  and edge weights  ${\cal W}.$

- We will need some basic graph theory to devise methods to combine different updates in the parallel scheme and analyze communication issues
- ullet Consider a graph  ${\cal G}$  as a collection of vertices  ${\cal V}$ , edges  ${\cal E}$ , and edge weights  ${\cal W}$ .
- We define undirected a graph for which an edge can carry information in both directions. Otherwise we call it directed.

- We will need some basic graph theory to devise methods to combine different updates in the parallel scheme and analyze communication issues
- ullet Consider a graph  ${\cal G}$  as a collection of vertices  ${\cal V}$ , edges  ${\cal E}$ , and edge weights  ${\cal W}$ .
- We define undirected a graph for which an edge can carry information in both directions. Otherwise we call it directed.
- We define the Adjacency matrix as  $A \in \{0,1\}^{|\mathcal{V}| \times |\mathcal{V}|}$ , with 1 entries if node i and j share an edge.

- We will need some basic graph theory to devise methods to combine different updates in the parallel scheme and analyze communication issues
- $\bullet$  Consider a graph  ${\cal G}$  as a collection of vertices  ${\cal V},$  edges  ${\cal E},$  and edge weights  ${\cal W}.$
- We define undirected a graph for which an edge can carry information in both directions. Otherwise we call it directed.
- We define the Adjacency matrix as  $A \in \{0,1\}^{|\mathcal{V}| \times |\mathcal{V}|}$ , with 1 entries if node i and j share an edge.
- $\bullet$  The Laplacian matrix of a undirected graph  ${\cal G}$  is an  $|{\cal V}|\times |{\cal V}|$  symmetric matrix defined by

$$\mathcal{L}_{\mathcal{G}} = \mathcal{D} - \mathcal{A}$$
,

where  $\mathcal{D} = \text{diag}(d_1, \ldots, d_n)$  is the degree matrix, which is the diagonal matrix formed from the vertex degrees. Note that  $\mathcal{L}_{\mathcal{G}} \mathbf{1} = \mathbf{0}$ .

- We will need some basic graph theory to devise methods to combine different updates in the parallel scheme and analyze communication issues
- ullet Consider a graph  ${\cal G}$  as a collection of vertices  ${\cal V}$ , edges  ${\cal E}$ , and edge weights  ${\cal W}$ .
- We define undirected a graph for which an edge can carry information in both directions. Otherwise we call it directed.
- We define the Adjacency matrix as  $A \in \{0,1\}^{|\mathcal{V}| \times |\mathcal{V}|}$ , with 1 entries if node i and j share an edge.
- $\bullet$  The Laplacian matrix of a undirected graph  ${\cal G}$  is an  $|{\cal V}|\times |{\cal V}|$  symmetric matrix defined by

$$\mathcal{L}_{\mathcal{G}} = \mathcal{D} - \mathcal{A}$$
,

where  $\mathcal{D} = \text{diag}(d_1, \ldots, d_n)$  is the degree matrix, which is the diagonal matrix formed from the vertex degrees. Note that  $\mathcal{L}_{\mathcal{G}}\mathbf{1} = \mathbf{0}$ .

• A doubly stochastic matrix W is a matrix for which,  $W\mathbf{1} = \mathbf{1}$ , and  $\mathbf{1}^{\top}W = \mathbf{1}^{\top}$ .

 We are ready to get back to our parallel update method. Consider the following update:

 We are ready to get back to our parallel update method. Consider the following update:

#### **Decentralized Gradient Descent**

- ▶ Start with  $\mathbf{x}_0^i \in \mathbf{R}^n$  at each device i
- ▶ Iterate: ask every devices:  $\mathbf{x}_{k+1}^i = \sum_{i=1}^N w_{ij} \mathbf{x}_k^i \alpha_k \nabla f_i(\mathbf{x}_k^i)$ ,  $\forall i$

Here  $w_{ij}$  are not zero only if device i and j are not connected.

 We are ready to get back to our parallel update method. Consider the following update:

#### **Decentralized Gradient Descent**

- ▶ Start with  $\mathbf{x}_0^i \in \mathbf{R}^n$  at each device i
- ▶ Iterate: ask every devices:  $\mathbf{x}_{k+1}^i = \sum_{i=1}^N w_{ij} \mathbf{x}_k^i \alpha_k \nabla f_i(\mathbf{x}_k^i)$ ,  $\forall i$

Here  $w_{ij}$  are not zero only if device i and j are not connected.

• Effectively we are giving a copy of x to each device and looking at the iterates  $y = [x^1; x^2; ...; x^N]$ 

$$y_{k+1} = Wy_k - \alpha_k \nabla_y F(y_k), \qquad F(y) := \sum_{i=1}^N f_i(\mathbf{x}^i),$$

where W is assumed doubly stochastic and it is the mixing matrix.

 We are ready to get back to our parallel update method. Consider the following update:

#### **Decentralized Gradient Descent**

- ▶ Start with  $\mathbf{x}_0^i \in \mathbf{R}^n$  at each device i
- ▶ Iterate: ask every devices:  $\mathbf{x}_{k+1}^i = \sum_{j=1}^N w_{ij} \mathbf{x}_k^j \alpha_k \nabla f_i(\mathbf{x}_k^i)$ ,  $\forall i$

Here  $w_{ij}$  are not zero only if device i and j are not connected.

• Effectively we are giving a copy of x to each device and looking at the iterates  $y = [x^1; x^2; ...; x^N]$ 

$$y_{k+1} = Wy_k - \alpha_k \nabla_y F(y_k), \qquad F(y) := \sum_{i=1}^N f_i(\mathbf{x}^i),$$

where W is assumed doubly stochastic and it is the mixing matrix.

- DGD is peer-to-peer:
  - ► Each device gets and sends  $x_k^j/x_k^i$  from/to the neighbors
  - ► Each device computes:  $\mathbf{x}_{k+1}^i = \sum_{i=1}^N w_{ij} \mathbf{x}_k^i \alpha_k \nabla f_i(\mathbf{x}_k^i), \quad \forall i$

 DGD is the first distributed algorithm that we see in this course and it does not need a server. Each device can compute its version of the optimal decision.

- DGD is the first distributed algorithm that we see in this course and it does not need a server. Each device can compute its version of the optimal decision.
- Communication needs to happen synchronously and bi-directionally

- DGD is the first distributed algorithm that we see in this course and it does not need a server. Each device can compute its version of the optimal decision.
- Communication needs to happen synchronously and bi-directionally
- Other versions exist with the mixing matrix being applied before/after the local gradient, but the idea stays the same

- DGD is the first distributed algorithm that we see in this course and it does not need a server. Each device can compute its version of the optimal decision.
- Communication needs to happen synchronously and bi-directionally
- Other versions exist with the mixing matrix being applied before/after the local gradient, but the idea stays the same
- What can we say about convergence? For a constant step size  $\alpha$ , it turns out that it is not so difficult, since

$$y_{k+1} = Wy_k - \alpha \nabla_y F(y_k) = y_k - \alpha \left[ \nabla_y F(y_k) - \frac{1}{\alpha} (W - I) y_k \right]$$

so it's the same as optimizing,

$$\min_{\mathbf{y}\in\mathsf{R}^{Nn}}\alpha F(\mathbf{y})+\frac{1}{2}\mathbf{y}^{\top}(I-W)\mathbf{y}.$$

- DGD is the first distributed algorithm that we see in this course and it does not need a server. Each device can compute its version of the optimal decision.
- Communication needs to happen synchronously and bi-directionally
- Other versions exist with the mixing matrix being applied before/after the local gradient, but the idea stays the same
- What can we say about convergence? For a constant step size  $\alpha$ , it turns out that it is not so difficult, since

$$y_{k+1} = Wy_k - \alpha \nabla_y F(y_k) = y_k - \alpha \left[ \nabla_y F(y_k) - \frac{1}{\alpha} (W - I) y_k \right]$$

so it's the same as optimizing,

$$\min_{\mathbf{y}\in\mathsf{R}^{Nn}}\alpha F(\mathbf{y})+\frac{1}{2}\mathbf{y}^{\top}(I-W)\mathbf{y}.$$

ullet For a general diminishing  $lpha_k$  sequence, it is a bit more complex but doable

## DGD convergence

• We assume that the mixing matrix  $W = [w_{ij}]$  is symmetric and doubly stochastic. The eigenvalues of W are real and sorted in a nonincreasing order

$$1 = \lambda_1(W) \ge \lambda_2(W) \ge \ldots \ge \lambda_N(W) \ge -1.$$

Let the second largest magnitude of the eigenvalues of W be denoted as

$$\gamma = \max\{|\lambda_2(W)|, |\lambda_N(W)|\}$$

# DGD convergence

• We assume that the mixing matrix  $W = [w_{ij}]$  is symmetric and doubly stochastic. The eigenvalues of W are real and sorted in a nonincreasing order

$$1 = \lambda_1(W) \ge \lambda_2(W) \ge \ldots \ge \lambda_N(W) \ge -1.$$

Let the second largest magnitude of the eigenvalues of W be denoted as

$$\gamma = \max\{|\lambda_2(W)|, |\lambda_N(W)|\}$$

• Some basic questions in the decentralized/distributed setting arise: (1) When does  $\mathbf{x}_k^i$  converge? (2) Does it converge to  $\mathbf{x}^*$  (3) If not does consensus (i.e.,  $\mathbf{x}_k^i = \mathbf{x}_k^i$ ) hold asymptotically? (4) How do the properties of  $f_i$  and the network affect convergence?

## DGD convergence: statement

## Theorem 4 (DGD convergence)

Consider Problem (P) and its solution via a decentralized gradient descent algorithm with constant step size  $\alpha$ , and doubly stochastic communication matrix W with  $\gamma < 1$ .

Let convex functions  $f_i$  be  $L_i$ -smooth and let  $L = \max_i \{L_i\}$ . Let the mean value be

$$\bar{\mathbf{x}}_k = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_k^i.$$

If  $\alpha$  is chosen small enough and in particular  $\leq O(1/L)$ , then

Consensus:

$$\|\mathbf{x}_k^i - \bar{\mathbf{x}}_k\| \to O(\frac{\alpha}{1-\gamma});$$

2 Convergence

$$f(\bar{\mathbf{x}}_k) - f^* \to O(\frac{\alpha}{1 - \gamma}).$$

 In the statement, we observe two things: convergence of the mean values and consensus around the mean, this is a general characteristic of distributed optimization

- In the statement, we observe two things: convergence of the mean values and consensus around the mean, this is a general characteristic of distributed optimization
- We also see a dependency on the graph entering via its spectrum  $\gamma$ , and in particular (here) the second largest eigenvalue (in modulus).

- In the statement, we observe two things: convergence of the mean values and consensus around the mean, this is a general characteristic of distributed optimization
- We also see a dependency on the graph entering via its spectrum  $\gamma$ , and in particular (here) the second largest eigenvalue (in modulus).
- The proofs can be extended to cases in which  $f_i$  are strongly convex (easier to handle), but it is much harder to work with non-doubly stochastic matrices (i.e., communication issues). See arXiv:1310.7063

- In the statement, we observe two things: convergence of the mean values and consensus around the mean, this is a general characteristic of distributed optimization
- We also see a dependency on the graph entering via its spectrum  $\gamma$ , and in particular (here) the second largest eigenvalue (in modulus).
- The proofs can be extended to cases in which  $f_i$  are strongly convex (easier to handle), but it is much harder to work with non-doubly stochastic matrices (i.e., communication issues). See arXiv:1310.7063
- We can see how for a constant step size we obtain a constant error bound (not so surprising, since we are changing the cost function!)

- In the statement, we observe two things: convergence of the mean values and consensus around the mean, this is a general characteristic of distributed optimization
- We also see a dependency on the graph entering via its spectrum  $\gamma$ , and in particular (here) the second largest eigenvalue (in modulus).
- The proofs can be extended to cases in which  $f_i$  are strongly convex (easier to handle), but it is much harder to work with non-doubly stochastic matrices (i.e., communication issues). See arXiv:1310.7063
- We can see how for a constant step size we obtain a constant error bound (not so surprising, since we are changing the cost function!)
- When  $\alpha$  is diminishing, you can obtain a zero error bound, but the analysis is more complicated, and you require typically extra assumptions.

• The parameter  $\gamma$  depends on the graph. For a full graph with weights 1/N,  $\gamma=0$ , because ..

- The parameter  $\gamma$  depends on the graph. For a full graph with weights 1/N,  $\gamma=0$ , because ..
- $\gamma$  will grow depending on how the graph is connected (in general, the more connected the lower is the  $\gamma$ ).

- The parameter  $\gamma$  depends on the graph. For a full graph with weights 1/N,  $\gamma=0$ , because ..
- $\gamma$  will grow depending on how the graph is connected (in general, the more connected the lower is the  $\gamma$ ).
- Communication issues

- The parameter  $\gamma$  depends on the graph. For a full graph with weights 1/N,  $\gamma=0$ , because ..
- $\gamma$  will grow depending on how the graph is connected (in general, the more connected the lower is the  $\gamma$ ).
- Communication issues
- Directed communication: this is hard, since we loose double-stochasticity. Also you can have cycles in the graph and counter-intuitive behaviors. Usually solved with algorithms of the type "push-sum"

- The parameter  $\gamma$  depends on the graph. For a full graph with weights 1/N,  $\gamma=0$ , because ..
- $\gamma$  will grow depending on how the graph is connected (in general, the more connected the lower is the  $\gamma$ ).

- Directed communication: this is hard, since we loose double-stochasticity.
   Also you can have cycles in the graph and counter-intuitive behaviors.
   Usually solved with algorithms of the type "push-sum"
- Package drops, asynchronous communication. This is often less hard to deal with, if the instantaneous W stay symmetric and doubly stochastic. Typically, the idea is to define a finite set of communication of matrices {W<sup>r</sup>}<sub>r</sub> and at each time k we randomly pick from them. If each W<sup>r</sup> has the same structural properties of W and the eigenvalues of a restricted sequence of W<sup>r</sup> have the same properties of the eigenvalues of W, then we can extend convergence results.

- The parameter  $\gamma$  depends on the graph. For a full graph with weights 1/N,  $\gamma=0$ , because ..
- $\gamma$  will grow depending on how the graph is connected (in general, the more connected the lower is the  $\gamma$ ).

- Directed communication: this is hard, since we loose double-stochasticity.
   Also you can have cycles in the graph and counter-intuitive behaviors.
   Usually solved with algorithms of the type "push-sum"
- Package drops, asynchronous communication. This is often less hard to deal with, if the instantaneous W stay symmetric and doubly stochastic.
   Typically, the idea is to define a finite set of communication of matrices {W<sup>r</sup>}<sub>r</sub> and at each time k we randomly pick from them. If each W<sup>r</sup> has the same structural properties of W and the eigenvalues of a restricted sequence of W<sup>r</sup> have the same properties of the eigenvalues of W, then we can extend convergence results.
- Gossip protocols: asynchronous and directed. A node wakes up and transmits to some of its neighbors, it also analyses what it has received, then it sleeps again. This is very realistic but quite hard to do. Here latencies...

• We have seen that DGD has fundamental trade-off between accuracy and performance. Can we do better?

- We have seen that DGD has fundamental trade-off between accuracy and performance. Can we do better?
- Consider the following two-communication rounds algorithm

- We have seen that DGD has fundamental trade-off between accuracy and performance. Can we do better?
- Consider the following two-communication rounds algorithm
  - Decentralized Gradient Tracking

- We have seen that DGD has fundamental trade-off between accuracy and performance. Can we do better?
- Consider the following two-communication rounds algorithm
  - Decentralized Gradient Tracking
  - ▶ Initialize  $\mathbf{x}_k^i$ ,  $\mathbf{g}_k^i$

- We have seen that DGD has fundamental trade-off between accuracy and performance. Can we do better?
- Consider the following two-communication rounds algorithm
  - Decentralized Gradient Tracking
  - ▶ Initialize  $\mathbf{x}_k^i$ ,  $\mathbf{g}_k^i$
  - Each device gets and sends  $x_k^j/x_k^i$  from/to the neighbors

- We have seen that DGD has fundamental trade-off between accuracy and performance. Can we do better?
- Consider the following two-communication rounds algorithm

#### Decentralized Gradient Tracking

- ▶ Initialize  $\mathbf{x}_k^i$ ,  $\mathbf{g}_k^i$
- ▶ Each device gets and sends  $x_k^j/x_k^i$  from/to the neighbors
- ► Each device computes:

$$\mathbf{x}_{k+1}^{i} = \sum_{j=1}^{N} w_{ij} \mathbf{x}_{k}^{j} - \alpha \mathbf{g}_{k}^{i}, \quad \forall i$$

- We have seen that DGD has fundamental trade-off between accuracy and performance. Can we do better?
- Consider the following two-communication rounds algorithm

#### Decentralized Gradient Tracking

- ▶ Initialize  $\mathbf{x}_k^i$ ,  $\mathbf{g}_k^i$
- ► Each device gets and sends  $x_k^j/x_k^i$  from/to the neighbors
- ► Each device computes:

$$\mathbf{x}_{k+1}^{i} = \sum_{i=1}^{N} w_{ij} \mathbf{x}_{k}^{j} - \alpha \mathbf{g}_{k}^{i}, \quad \forall i$$

• Each device gets and sends  $g_k^j/g_k^i$  from/to the neighbors

- We have seen that DGD has fundamental trade-off between accuracy and performance. Can we do better?
- Consider the following two-communication rounds algorithm

#### Decentralized Gradient Tracking

- ▶ Initialize  $\mathbf{x}_k^i$ ,  $\mathbf{g}_k^i$
- ► Each device gets and sends  $x_k^j/x_k^i$  from/to the neighbors
- ► Each device computes:

$$\mathbf{x}_{k+1}^{i} = \sum_{j=1}^{N} w_{ij} \mathbf{x}_{k}^{j} - \alpha \mathbf{g}_{k}^{i}, \quad \forall i$$

- Each device gets and sends  $g_k^j/g_k^i$  from/to the neighbors
- Each device computes:

$$\boldsymbol{g}_{k+1}^{i} = \sum_{i=1}^{N} w_{ij} \boldsymbol{g}_{k}^{j} + (\nabla f_{i}(\boldsymbol{x}_{k+1}^{i}) - \nabla f_{i}(\boldsymbol{x}_{k}^{i})), \quad \forall$$

- We have seen that DGD has fundamental trade-off between accuracy and performance. Can we do better?
- Consider the following two-communication rounds algorithm

#### Decentralized Gradient Tracking

- ▶ Initialize  $\mathbf{x}_k^i$ ,  $\mathbf{g}_k^i$
- ► Each device gets and sends  $x_k^j/x_k^i$  from/to the neighbors
- ► Each device computes:

$$\mathbf{x}_{k+1}^{i} = \sum_{j=1}^{N} w_{ij} \mathbf{x}_{k}^{j} - \alpha \mathbf{g}_{k}^{i}, \quad \forall i$$

- Each device gets and sends  $g_k^j/g_k^i$  from/to the neighbors
- ► Each device computes:

$$\boldsymbol{g}_{k+1}^{i} = \sum_{i=1}^{N} w_{ij} \boldsymbol{g}_{k}^{j} + (\nabla f_{i}(\boldsymbol{x}_{k+1}^{i}) - \nabla f_{i}(\boldsymbol{x}_{k}^{i})), \quad \forall$$

 We have now a term that tracks the gradient values and "integrates" the errors

# Gradient tracking convergence

 $\bullet$  Start by defining the mean quantity,  $\bar{\mathbf{x}}_k = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_k^i$ 

# Gradient tracking convergence

 $\bullet$  Start by defining the mean quantity,  $\overline{\mathbf{x}}_k = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_k^i$ 

## Theorem 5 (Gradient Tracking convergence)

Consider Problem (P) and its solution via a gradient tracking algorithm with constant step size  $\alpha$ , and doubly stochastic communication matrix W with  $w_{ij} \geq 0$ . Let convex functions  $f_i$  be L-smooth and strongly convex. If  $\alpha$  is chosen small enough and in particular  $\leq O(1/L)$ , then

Consensus:

$$\|\mathbf{x}_k^i - \bar{\mathbf{x}}_k\| \to 0;$$

2 Convergence

$$\|\bar{\mathbf{x}}_k - \mathbf{x}^*\| \to 0.$$

And convergence is linear.

Proof: arXiv:1906.10760

• The proof construct a recursion of the form  $\mathbf{v}_{k+1} \leq J(\alpha)\mathbf{v}_k$ , and the condition on  $\alpha$  guarantees that  $||J(\alpha)|| < 1$ .

- The proof construct a recursion of the form  $\mathbf{v}_{k+1} \leq J(\alpha)\mathbf{v}_k$ , and the condition on  $\alpha$  guarantees that  $||J(\alpha)|| < 1$ .
- Extensions to non-strongly convex or to other communication patters are possible, but very involved

- The proof construct a recursion of the form  $\mathbf{v}_{k+1} \leq J(\alpha)\mathbf{v}_k$ , and the condition on  $\alpha$  guarantees that  $||J(\alpha)|| < 1$ .
- Extensions to non-strongly convex or to other communication patters are possible, but very involved
- On the good side, we have a first-order primal method that is fully distributed and allows us to reach zero consensus and residual error.

- The proof construct a recursion of the form  $\mathbf{v}_{k+1} \leq J(\alpha)\mathbf{v}_k$ , and the condition on  $\alpha$  guarantees that  $||J(\alpha)|| < 1$ .
- Extensions to non-strongly convex or to other communication patters are possible, but very involved
- On the good side, we have a first-order primal method that is fully distributed and allows us to reach zero consensus and residual error.
- On the minus side, we communicate quite a lot in terms of gradients, so we lose in privacy

• The general problem of interest is minimizing  $\sum_i f_i(\mathbf{x})$ 

- The general problem of interest is minimizing  $\sum_i f_i(\mathbf{x})$
- When N is big or  $f_i$  cannot be shared (efficiently or at all), decentralized schemes play an important role

- The general problem of interest is minimizing  $\sum_i f_i(\mathbf{x})$
- When N is big or  $f_i$  cannot be shared (efficiently or at all), decentralized schemes play an important role
- Communication and computation and their trade-off are then key: and of course, first-order algorithms

- The general problem of interest is minimizing  $\sum_i f_i(\mathbf{x})$
- When N is big or  $f_i$  cannot be shared (efficiently or at all), decentralized schemes play an important role
- Communication and computation and their trade-off are then key: and of course, first-order algorithms
- Two type of schemes: Incremental or parallel

- The general problem of interest is minimizing  $\sum_i f_i(\mathbf{x})$
- When N is big or  $f_i$  cannot be shared (efficiently or at all), decentralized schemes play an important role
- Communication and computation and their trade-off are then key: and of course, first-order algorithms
- Two type of schemes: Incremental or parallel
- For the former, we have seen a convergence result. The latter yields "distributed optimization" and, in the primal space, we have seen two example of algorithms: DGD and Gradient Tracking in a simplified "controlled" setting

- The general problem of interest is minimizing  $\sum_i f_i(\mathbf{x})$
- When N is big or  $f_i$  cannot be shared (efficiently or at all), decentralized schemes play an important role
- Communication and computation and their trade-off are then key: and of course, first-order algorithms
- Two type of schemes: Incremental or parallel
- For the former, we have seen a convergence result. The latter yields "distributed optimization" and, in the primal space, we have seen two example of algorithms: DGD and Gradient Tracking in a simplified "controlled" setting
- Convergence proofs in distributed settings are harder since the communication plays an important role!

- The general problem of interest is minimizing  $\sum_i f_i(\mathbf{x})$
- When N is big or  $f_i$  cannot be shared (efficiently or at all), decentralized schemes play an important role
- Communication and computation and their trade-off are then key: and of course, first-order algorithms
- Two type of schemes: Incremental or parallel
- For the former, we have seen a convergence result. The latter yields "distributed optimization" and, in the primal space, we have seen two example of algorithms: DGD and Gradient Tracking in a simplified "controlled" setting
- Convergence proofs in distributed settings are harder since the communication plays an important role!
- Why would I care about communication or distributed when we can all upload on the cloud?

- The general problem of interest is minimizing  $\sum_i f_i(\mathbf{x})$
- When N is big or  $f_i$  cannot be shared (efficiently or at all), decentralized schemes play an important role
- Communication and computation and their trade-off are then key: and of course, first-order algorithms
- Two type of schemes: Incremental or parallel
- For the former, we have seen a convergence result. The latter yields "distributed optimization" and, in the primal space, we have seen two example of algorithms: DGD and Gradient Tracking in a simplified "controlled" setting
- Convergence proofs in distributed settings are harder since the communication plays an important role!
- Why would I care about communication or distributed when we can all upload on the cloud? (Think BIG: large-data files, server farms, or many sensors, etc..)

### Sample references

- The standard book: Dimitri P. Bertsekas and John N. Tsitsiklis, Parallel and Distributed Computation: Numerical Methods, 1997, https://web.mit.edu/dimitrib/www/pdc.html
- Two recent articles: Kun Yuan, Qing Ling, Wotao Yin, On the Convergence of Decentralized Gradient Descent, arXiv:1310.7063 and SIAM Journal on Optimization, 2016
  - Giuseppe Notarstefano, Ivano Notarnicola, Andrea Camisa, **Distributed Optimization for Smart Cyber-Physical Networks**, arXiv:1906.10760 and Foundations and Trends in Systems and Control, 2019
- Many variants out there.

### Class 2

We remind that the problem at hand is

$$(P) \qquad \min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^N f_i(\mathbf{x}).$$

and we proceed as in the first class by endowing each device with a copy of x, so that the problem becomes,

$$(P') \qquad \min_{\mathbf{x}^i \in \mathbb{R}^n, i=1,\dots,N} \sum_{i=1}^N f_i(\mathbf{x}^i) \quad \text{subject to } \mathbf{x}^i = \mathbf{x}^j, \forall i \sim j$$

We remind that the problem at hand is

$$(P) \qquad \min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^N f_i(\mathbf{x}).$$

and we proceed as in the first class by endowing each device with a copy of **x**, so that the problem becomes,

$$(P') \qquad \min_{\mathbf{x}^i \in \mathbb{R}^n, i=1,\dots,N} \sum_{i=1}^N f_i(\mathbf{x}^i) \quad \text{subject to } \mathbf{x}^i = \mathbf{x}^j, \forall i \sim j$$

• We can then write (P') compactly as,

$$(P'') \qquad \min_{y \in \mathbb{R}^{N_n}} F(y) \quad \text{subject to } Ay = 0$$

We remind that the problem at hand is

$$(P) \qquad \min_{\mathbf{x} \in \mathbb{R}^n} \sum_{i=1}^N f_i(\mathbf{x}).$$

and we proceed as in the first class by endowing each device with a copy of  $\mathbf{x}$ , so that the problem becomes,

$$(P') \qquad \min_{\mathbf{x}^i \in \mathbb{R}^n, i=1,\dots,N} \sum_{i=1}^N f_i(\mathbf{x}^i) \quad \text{subject to } \mathbf{x}^i = \mathbf{x}^j, \forall i \sim j$$

• We can then write (P') compactly as,

$$(P'') \qquad \min_{y \in \mathbb{R}^{Nn}} F(y) \quad \text{subject to } Ay = 0$$

• Can we now look at its dual problem?

### Let's revisit duality

• Consider the problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad \text{subject to } A\mathbf{x} = b,$$

for a convex function  $f: \mathbf{R}^n \to \mathbf{R}$ , and matrix  $A \in \mathbf{R}^{p \times n}$ , vector  $b \in \mathbf{R}^p$  (with b in the image of A).

### Let's revisit duality

Consider the problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad \text{subject to } A\mathbf{x} = b,$$

for a convex function  $f: \mathbf{R}^n \to \mathbf{R}$ , and matrix  $A \in \mathbf{R}^{p \times n}$ , vector  $b \in \mathbf{R}^p$  (with b in the image of A).

• The Lagrangian function is defined as,

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

where  $\lambda \in \mathbf{R}^p$  are the Lagrangian multipliers (aka the dual variables)

### Let's revisit duality

• Consider the problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad \text{subject to } A\mathbf{x} = b,$$

for a convex function  $f: \mathbf{R}^n \to \mathbf{R}$ , and matrix  $A \in \mathbf{R}^{p \times n}$ , vector  $b \in \mathbf{R}^p$  (with b in the image of A).

• The Lagrangian function is defined as,

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

where  $\lambda \in \mathbf{R}^p$  are the Lagrangian multipliers (aka the dual variables)

The Lagrangian dual function is then,

$$q(\lambda) = \inf_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \lambda) = -(f^*(-A^\top \lambda) + \lambda^\top b),$$

where  $f^*$  is the conjugate function of f, i.e.,  $f^*(y) = \sup_{\mathbf{x}} \{ y^\top x - f(\mathbf{x}) \}$ . **Homework.** Prove the last equality.

• The dual problem is then,

 $\max_{\lambda \in \mathsf{R}^p} q(\lambda)$ 

The dual problem is then,

$$\max_{\lambda \in \mathsf{R}^p} q(\lambda)$$

 The dual problem is always convex, even if the primal is not and provide a lower bound on the minimum of the primal problem, i.e.,

$$q^* \leq f^*$$
.

The dual problem is then,

$$\max_{\lambda \in \mathbb{R}^p} q(\lambda)$$

 The dual problem is always convex, even if the primal is not and provide a lower bound on the minimum of the primal problem, i.e.,

$$q^* \leq f^*$$
.

 Equality holds if some constraint qualification holds. For convex problems, if Slater's condition is verified (there exists a strictly feasible solution), then equality holds and we say that we have strong duality.

• The dual problem is then,

$$\max_{\lambda \in \mathbb{R}^p} q(\lambda)$$

 The dual problem is always convex, even if the primal is not and provide a lower bound on the minimum of the primal problem, i.e.,

$$q^* \leq f^*$$
.

- Equality holds if some constraint qualification holds. For convex problems, if Slater's condition is verified (there exists a strictly feasible solution), then equality holds and we say that we have strong duality.
- For this course strong duality always holds because I don't look at dualizing inequality constraints.

• We remind that the problem at hand is

$$(P'')$$
  $\min_{y \in \mathbb{R}^{N_n}} F(y)$  subject to  $Ay = 0$ 

whose Lagrangian is

$$L(\mathbf{y},\lambda) = F(\mathbf{y}) + \lambda^{\top} A \mathbf{y}$$

and dual problem,

$$\max_{\lambda} q(\lambda) := \inf_{y} \{ F(y) + \lambda^{\top} A y \} = -F^{\star} (-A^{\top} \lambda).$$

## Let's revisit our problem

We remind that the problem at hand is

$$(P'')$$
  $\min_{y \in \mathbb{R}^{N_n}} F(y)$  subject to  $Ay = 0$ 

whose Lagrangian is

$$L(\mathbf{y},\lambda) = F(\mathbf{y}) + \lambda^{\top} A \mathbf{y}$$

and dual problem,

$$\max_{\lambda} q(\lambda) := \inf_{y} \{ F(y) + \lambda^{\top} A y \} = -F^{\star} (-A^{\top} \lambda).$$

• Can we set up a dual ascent algorithm?

• Start with a  $\lambda_0$  and iterate:

- Start with a  $\lambda_0$  and iterate:
- $\lambda_{k+1} = \lambda_k + \alpha_k \mathbf{v}_k$ , with  $\mathbf{v}_k \in \partial q(\lambda_k)$ , for all  $k = 0, 1, \ldots$ Rem: the dual function may not be differentiable even if f is so, so we need the sub-differential

- Start with a  $\lambda_0$  and iterate:
- $\lambda_{k+1} = \lambda_k + \alpha_k \mathbf{v}_k$ , with  $\mathbf{v}_k \in \partial q(\lambda_k)$ , for all  $k = 0, 1, \ldots$ Rem: the dual function may not be differentiable even if f is so, so we need the sub-differential
- What's  $\partial q(\lambda_k)$ ?

- Start with a  $\lambda_0$  and iterate:
- $\lambda_{k+1} = \lambda_k + \alpha_k \mathbf{v}_k$ , with  $\mathbf{v}_k \in \partial q(\lambda_k)$ , for all  $k = 0, 1, \ldots$ Rem: the dual function may not be differentiable even if f is so, so we need the sub-differential
- What's  $\partial q(\lambda_k)$  ?
- Start with

$$\partial_{\lambda} q(\lambda_k) = -A \partial_{\lambda} [-F^*(-A^{\top} \lambda_k)]$$

Then, we know that for convex functions  $(\partial_v F)^{-1}(v) = \partial_u F^*(u)$ , and in addition let,

$$y^*(\lambda) := \arg\min_{y} L(y, \lambda) \iff y^*(\lambda) = (\partial F)^{-1}(-A^{\top}\lambda)$$

Therefore,

$$\partial q(\lambda_k) = Ay^*(\lambda_k)$$
 residual map!

• Therefore the dual ascent yields, Start with a  $\lambda_0$  and iterate:

- Therefore the dual ascent yields. Start with a  $\lambda_0$  and iterate:
- Solve locally

$$y^*(\lambda) := \arg\min_{\mathbf{y}} L(\mathbf{y}, \lambda_k) = \arg\min_{\mathbf{x}^1, \dots, \mathbf{x}^N} \sum_{i=1}^n f_i(\mathbf{x}^i) + \lambda_k^\top A[\mathbf{x}^1, \dots, \mathbf{x}^N]^\top$$

Upon dividing  $A = [A_1 | \dots | A_N]$ ,

$$\iff \mathbf{x}^{i,*}(\lambda_k) := \arg\min_{\mathbf{x}_i} f_i(\mathbf{x}^i) + \lambda_k^\top A_i \mathbf{x}^i \qquad \forall i$$

A. Simonetto 38 | 61

- Therefore the dual ascent yields, Start with a  $\lambda_0$  and iterate:
- Solve locally

$$y^*(\lambda) := \arg\min_{\mathbf{y}} L(\mathbf{y}, \lambda_k) = \arg\min_{\mathbf{x}^1, \dots, \mathbf{x}^N} \sum_{i=1}^n f_i(\mathbf{x}^i) + \lambda_k^\top A[\mathbf{x}^1, \dots, \mathbf{x}^N]^\top$$

Upon dividing  $A = [A_1 | \dots | A_N]$ ,

$$\iff \mathbf{x}^{i,*}(\lambda_k) := \arg\min_{\mathbf{x}_i} f_i(\mathbf{x}^i) + \lambda_k^{\top} A_i \mathbf{x}^i \qquad \forall i,$$

• Send and update  $\mathbf{x}^{i,*}(\lambda_k)$  to all, and gather,

$$\partial q(\lambda_k) = Ay^*(\lambda_k) = \sum_i A_i \mathbf{x}^{i,*}(\lambda_k)$$

- Therefore the dual ascent yields, Start with a  $\lambda_0$  and iterate:
- Solve locally

$$y^*(\lambda) := \arg\min_{\mathbf{y}} L(\mathbf{y}, \lambda_k) = \arg\min_{\mathbf{x}^1, \dots, \mathbf{x}^N} \sum_{i=1}^n f_i(\mathbf{x}^i) + \lambda_k^\top A[\mathbf{x}^1, \dots, \mathbf{x}^N]^\top$$

Upon dividing  $A = [A_1 | \dots | A_N]$ ,

$$\iff \mathbf{x}^{i,*}(\lambda_k) := \arg\min_{\mathbf{x}_i} f_i(\mathbf{x}^i) + \lambda_k^{\top} A_i \mathbf{x}^i \qquad \forall i,$$

• Send and update  $\mathbf{x}^{i,*}(\lambda_k)$  to all, and gather,

$$\partial q(\lambda_k) = Ay^*(\lambda_k) = \sum_i A_i \mathbf{x}^{i,*}(\lambda_k)$$

• Update then,  $\lambda_{k+1} = \lambda_k + \alpha_k \sum_i A_i x^{i,*}(\lambda_k)$ , for all  $k = 0, 1, \ldots$ 

• The method is called dual decomposition, since you decompose the primal problem via the dual variables;

- The method is called dual decomposition, since you decompose the primal problem via the dual variables;
- The dual variables are updated globally, the primal locally, and this is convenient if  $f_i$  depends on data that you can't share/ don't want to!

- The method is called dual decomposition, since you decompose the primal problem via the dual variables;
- The dual variables are updated globally, the primal locally, and this is convenient if f<sub>i</sub> depends on data that you can't share/ don't want to!
- $\lambda$  can also live locally on the edges! Let A be the collection of  $\mathbf{x}^i - \mathbf{x}^j$  for the edge set. Then the associated dual variable is  $\lambda^{ij}$  and that can live on the edge i,j.

- The method is called dual decomposition, since you decompose the primal problem via the dual variables;
- The dual variables are updated globally, the primal locally, and this is convenient if f<sub>i</sub> depends on data that you can't share/ don't want to!
- $\lambda$  can also live locally on the edges! Let A be the collection of  $\mathbf{x}^i - \mathbf{x}^j$  for the edge set. Then the associated dual variable is  $\lambda^{ij}$  and that can live on the edge i,j.
- The last means that also  $\lambda^{ij}$  can also be done locally!  $\lambda^{ij}_{k+1} = \lambda^{ij}_k + \alpha_k [\mathbf{x}^{i,*}(\lambda_k) \mathbf{x}^{j,*}(\lambda_k)]$ , for all  $k = 0, 1, \dots$

• Consider a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , undirected. We indicate  $i \sim j$  if there is an edge between i and j. We also define A as

$$A_{ij} = [1_i, -1_j]$$
  $\forall i \sim j, j < i.$ 

• Consider a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , undirected. We indicate  $i \sim j$  if there is an edge between i and j. We also define A as

$$A_{ij} = [1_i, -1_j] \qquad \forall i \sim j, j < i.$$

• Initialize  $\lambda_0^{ij}$  for  $(i,j) \in \mathcal{E}, j < i$ 

• Consider a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , undirected. We indicate  $i \sim j$  if there is an edge between i and j. We also define A as

$$A_{ij} = [1_i, -1_j] \qquad \forall i \sim j, j < i.$$

- Initialize  $\lambda_0^{ij}$  for  $(i,j) \in \mathcal{E}, j < i$
- For all nodes *i* do.

$$\mathbf{x}^{i,*}(\lambda_k) := \arg\min_{\mathbf{x}_i} f_i(\mathbf{x}^i) + \sum_{i \sim j} \lambda_k^{ij,\top} A_{ij} \mathbf{x}^i$$

• Consider a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , undirected. We indicate  $i \sim j$  if there is an edge between i and j. We also define A as

$$A_{ij} = [1_i, -1_j] \quad \forall i \sim j, j < i.$$

- Initialize  $\lambda_0^{ij}$  for  $(i,j) \in \mathcal{E}, j < i$
- For all nodes i do,

$$\mathbf{x}^{i,*}(\lambda_k) := \arg\min_{\mathbf{x}_i} f_i(\mathbf{x}^i) + \sum_{i \sim j} \lambda_k^{ij,\top} A_{ij} \mathbf{x}^i$$

• Send  $\mathbf{x}^{i,*}(\lambda_k)$  to neighbors

• Consider a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , undirected. We indicate  $i \sim j$  if there is an edge between i and j. We also define A as

$$A_{ij} = [1_i, -1_j] \qquad \forall i \sim j, j < i.$$

- Initialize  $\lambda_0^{ij}$  for  $(i,j) \in \mathcal{E}, j < i$
- For all nodes i do,

$$\mathbf{x}^{i,*}(\lambda_k) := \arg\min_{\mathbf{x}_i} f_i(\mathbf{x}^i) + \sum_{i \sim j} \lambda_k^{ij,\top} A_{ij} \mathbf{x}^i$$

- Send  $\mathbf{x}^{i,*}(\lambda_k)$  to neighbors
- Update for node *i*, *j*:

$$\lambda_{k+1}^{ij} = \lambda_k^{ij} + \alpha_k [\mathbf{x}^{i,*}(\lambda_k) - \mathbf{x}^{j,*}(\lambda_k)]$$

• Consider a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , undirected. We indicate  $i \sim j$  if there is an edge between i and j. We also define A as

$$A_{ij} = [1_i, -1_j] \qquad \forall i \sim j, j < i.$$

- Initialize  $\lambda_0^{ij}$  for  $(i,j) \in \mathcal{E}, j < i$
- For all nodes *i* do.

$$\mathbf{x}^{i,*}(\lambda_k) := \arg\min_{\mathbf{x}_i} f_i(\mathbf{x}^i) + \sum_{i \sim j} \lambda_k^{ij,\top} A_{ij} \mathbf{x}^i$$

- Send  $\mathbf{x}^{i,*}(\lambda_k)$  to neighbors
- Update for node *i*, *j*:

$$\lambda_{k+1}^{ij} = \lambda_k^{ij} + \alpha_k[\mathbf{x}^{i,*}(\lambda_k) - \mathbf{x}^{j,*}(\lambda_k)]$$

 This has two local updates and one communication round (synchronous, undirected)

• We look now at the convergence of dual decomposition. We need to recall some of the properties of the dual function. In particular, let  $\sigma(A)$  be the singular values of A.

If f is m-strongly convex then -q is  $\sigma_{\max}^2(A)/m$  smooth; If f is L-smooth then -q is  $\sigma_{\min}^2(A)/L$  strongly convex; **Homework.** Remind yourself of why it is so.

• We look now at the convergence of dual decomposition. We need to recall some of the properties of the dual function. In particular, let  $\sigma(A)$  be the singular values of A.

If f is m-strongly convex then -q is  $\sigma_{\max}^2(A)/m$  smooth; If f is L-smooth then -q is  $\sigma_{\min}^2(A)/L$  strongly convex; **Homework.** Remind yourself of why it is so.

• Then we have the following

• We look now at the convergence of dual decomposition. We need to recall some of the properties of the dual function. In particular, let  $\sigma(A)$  be the singular values of A.

If f is m-strongly convex then -q is  $\sigma_{\max}^2(A)/m$  smooth; If f is L-smooth then -q is  $\sigma_{\min}^2(A)/L$  strongly convex; **Homework.** Remind yourself of why it is so.

Then we have the following

#### Theorem 6 (Dual decomposition convergence)

Consider problem (P) and the dual decomposition approach for a certain connection matrix A. If  $f_i$ 's are m-strongly convex and L-smooth, then we can select  $\alpha < 2m/\sigma_{\max}^2(A)$  and obtaining a linear rate of convergence as,

$$\|\lambda_k - \lambda^*\| \leq \left[ \max\{|1 - \alpha \frac{\sigma_{\max}^2(A)}{m}|, |1 - \alpha \frac{\sigma_{\min}^2(A)}{L}|\} \right]^k \|\lambda_0 - \lambda^*\|$$

and

$$\|\mathbf{x}_k^i - \mathbf{x}^*\| \leq \frac{\sigma_{\mathsf{max}}(A)}{m} \|\lambda_k - \lambda^*\|.$$

• Proof follows from OPT201/202. Homework. Work it out.

- Proof follows from OPT201/202. Homework. Work it out.
- Convergence is proven when A is full row rank (so that  $\sigma_{\min} > 0$ ). But why would it be?

- Proof follows from OPT201/202. Homework. Work it out.
- Convergence is proven when A is full row rank (so that  $\sigma_{\min} > 0$ ). But why would it be?
- In the case A is not full row rank?  $\lambda^*$  is not unique and -q is not strongly convex. But, we can look at restricted strong convexity instead!

- Proof follows from OPT201/202. Homework. Work it out.
- Convergence is proven when A is full row rank (so that  $\sigma_{\min} > 0$ ). But why would it be?
- In the case A is not full row rank?  $\lambda^*$  is not unique and -q is not strongly convex. But, we can look at restricted strong convexity instead!

Let's see how to do it

- Proof follows from OPT201/202. Homework. Work it out.
- Convergence is proven when A is full row rank (so that  $\sigma_{\min} > 0$ ). But why would it be?
- In the case A is not full row rank?  $\lambda^*$  is not unique and -q is not strongly convex. But, we can look at restricted strong convexity instead!

#### Let's see how to do it

• First  $\lambda^* = \lambda_0^* + \lambda_1^*$ , with  $\lambda_0^* \in \operatorname{im}(A)$  and  $\lambda_1^* \in \operatorname{null}(A^\top)$ . One can show that  $\lambda_0^*$  is unique and we concentrate on that one, since  $\lambda_1^*$  is redundant  $A^\top \lambda_0^* = 0$ .

- Proof follows from OPT201/202. **Homework.** Work it out.
- Convergence is proven when A is full row rank (so that  $\sigma_{\min} > 0$ ). But why would it be?
- In the case A is not full row rank?  $\lambda^*$  is not unique and -q is not strongly convex. But, we can look at restricted strong convexity instead!

#### Let's see how to do it

- First  $\lambda^* = \lambda_0^* + \lambda_1^*$ , with  $\lambda_0^* \in \operatorname{im}(A)$  and  $\lambda_1^* \in \operatorname{null}(A^\top)$ . One can show that  $\lambda_0^*$  is unique and we concentrate on that one, since  $\lambda_1^*$  is redundant  $A^\top \lambda_0^* = 0$ .
- Start with  $\lambda_0 \in \operatorname{im}(A)$ ; then, it is direct that:

$$\lambda_{k+1} = \lambda_k + \alpha A \mathbf{x} \in \mathrm{im}(A)$$

- Proof follows from OPT201/202. Homework. Work it out.
- Convergence is proven when A is full row rank (so that  $\sigma_{\min} > 0$ ). But why would it be?
- In the case A is not full row rank?  $\lambda^*$  is not unique and -q is not strongly convex. But, we can look at restricted strong convexity instead!

#### Let's see how to do it

- First  $\lambda^* = \lambda_0^* + \lambda_1^*$ , with  $\lambda_0^* \in \operatorname{im}(A)$  and  $\lambda_1^* \in \operatorname{null}(A^\top)$ . One can show that  $\lambda_0^*$  is unique and we concentrate on that one, since  $\lambda_1^*$  is redundant  $A^\top \lambda_0^* = 0$ .
- Start with  $\lambda_0 \in \operatorname{im}(A)$ ; then, it is direct that:

$$\lambda_{k+1} = \lambda_k + \alpha A \mathbf{x} \in \mathrm{im}(A)$$

Then, can we say that (restricted strong convexity)

$$(\partial q(\lambda) - \partial q(\lambda'))^{\top}(\lambda' - \lambda) \ge \sigma_{\min}^2 / L \|\lambda' - \lambda\|^2, \quad \forall \lambda, \lambda' \in \operatorname{im}(A)$$
?

with  $\sigma_{\overline{\min}}$  the minimum non-zero singular value of A?

• Yes, substitute  $\partial q(\lambda) = A \partial F^*(-A^T \lambda)$ , then,

$$(\partial F^{*}(-A^{\top}\lambda) - \partial F^{*}(-A^{\top}\lambda'))^{\top}A^{\top}(\lambda' - \lambda) =$$

$$(\partial F^{*}(-A^{\top}\lambda) - \partial F^{*}(-A^{\top}\lambda'))^{\top}(-A^{\top}\lambda + A^{\top}\lambda') \geq$$

$$1/L\|A^{\top}(\lambda' - \lambda)\|^{2} \geq \sigma_{\min}^{2}/L\|\lambda' - \lambda)\|^{2}.$$

• Yes, substitute  $\partial q(\lambda) = A\partial F^*(-A^T\lambda)$ , then,

$$(\partial F^{*}(-A^{\top}\lambda) - \partial F^{*}(-A^{\top}\lambda'))^{\top}A^{\top}(\lambda' - \lambda) =$$

$$(\partial F^{*}(-A^{\top}\lambda) - \partial F^{*}(-A^{\top}\lambda'))^{\top}(-A^{\top}\lambda + A^{\top}\lambda') \geq$$

$$1/L\|A^{\top}(\lambda' - \lambda)\|^{2} \geq \sigma_{\min}^{2}/L\|\lambda' - \lambda)\|^{2}.$$

• The latest step is due to the fact that  $A^{T}(\lambda - \lambda') = 0$  iff  $\lambda = \lambda'$ .

• Yes, substitute  $\partial q(\lambda) = A \partial F^*(-A^\top \lambda)$ , then,

$$(\partial F^*(-A^{\top}\lambda) - \partial F^*(-A^{\top}\lambda'))^{\top}A^{\top}(\lambda' - \lambda) =$$

$$(\partial F^*(-A^{\top}\lambda) - \partial F^*(-A^{\top}\lambda'))^{\top}(-A^{\top}\lambda + A^{\top}\lambda') \ge$$

$$1/L\|A^{\top}(\lambda' - \lambda)\|^2 \ge \sigma_{\min}^2/L\|\lambda' - \lambda)\|^2.$$

- The latest step is due to the fact that  $A^{\top}(\lambda \lambda') = 0$  iff  $\lambda = \lambda'$ .
- So the Theorem becomes: we can select  $\alpha < 2m/\sigma_{\max}^2(A)$  and  $\lambda_0 \in \operatorname{im}(A)$ , and obtaining a linear rate of convergence as,

$$\|\lambda_k - \lambda^*\| \leq \left\lceil \max\{|1 - \alpha \frac{\sigma_{\max}^2(A)}{m}|, |1 - \alpha \frac{\sigma_{\min}^2(A)}{L}|\} \right\rceil^k \|\lambda_0 - \lambda^*\|$$

• Yes, substitute  $\partial q(\lambda) = A \partial F^*(-A^\top \lambda)$ , then,

$$(\partial F^*(-A^{\top}\lambda) - \partial F^*(-A^{\top}\lambda'))^{\top}A^{\top}(\lambda' - \lambda) =$$

$$(\partial F^*(-A^{\top}\lambda) - \partial F^*(-A^{\top}\lambda'))^{\top}(-A^{\top}\lambda + A^{\top}\lambda') \ge$$

$$1/L\|A^{\top}(\lambda' - \lambda)\|^2 \ge \sigma_{\min}^2/L\|\lambda' - \lambda)\|^2.$$

- The latest step is due to the fact that  $A^{\top}(\lambda \lambda') = 0$  iff  $\lambda = \lambda'$ .
- So the Theorem becomes: we can select  $\alpha < 2m/\sigma_{\max}^2(A)$  and  $\lambda_0 \in \operatorname{im}(A)$ , and obtaining a linear rate of convergence as,

$$\|\lambda_k - \lambda^*\| \leq \left\lceil \max\{|1 - \alpha \frac{\sigma_{\max}^2(A)}{m}|, |1 - \alpha \frac{\sigma_{\min}^2(A)}{L}|\} \right\rceil^k \|\lambda_0 - \lambda^*\|$$

• Note:  $\lambda_0 = \mathbf{0} \in \operatorname{im}(A)$ 

• The step size can be chosen constant and you still converge linearly! (As in the gradient tracking case); but the step size has to be chosen carefully

- The step size can be chosen constant and you still converge linearly! (As in the gradient tracking case); but the step size has to be chosen carefully
- Convergence depends on the graph, and in particular on  $\sigma_{min}$ ,  $\sigma_{max}$  It is not only a matter of function properties, but also graph properties!

- The step size can be chosen constant and you still converge linearly! (As in the gradient tracking case); but the step size has to be chosen carefully
- Convergence depends on the graph, and in particular on  $\sigma_{\min}$ ,  $\sigma_{\max}$  It is not only a matter of function properties, but also graph properties!
- Convergence is more tricky when F is not strongly convex, since q will not be differentiable, so it also hard to recover the primal solution here

## Dual decomposition: convergence III

- The step size can be chosen constant and you still converge linearly! (As in the gradient tracking case); but the step size has to be chosen carefully
- Convergence depends on the graph, and in particular on  $\sigma_{\min}$ ,  $\sigma_{\max}$  It is not only a matter of function properties, but also graph properties!
- Convergence is more tricky when F is not strongly convex, since q will not be differentiable, so it also hard to recover the primal solution here
- Convergence may be complicated when communication is directed, or asynchronous, or you have latencies, as in the first class

# The Alternating Direction Method of Multipliers (ADMM)

 To overcome the drawbacks of dual decomposition, we introduce another iterative method: the ADMM. We are going to develop it for a general case and then focus it on our distributed problem

# The Alternating Direction Method of Multipliers (ADMM)

- To overcome the drawbacks of dual decomposition, we introduce another iterative method: the ADMM. We are going to develop it for a general case and then focus it on our distributed problem
- Consider the problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^p} f(\mathbf{x}) + g(\mathbf{y})$$
subject to 
$$A\mathbf{x} + B\mathbf{y} = c.$$

For well-defined matrices and vectors, as well as convex functions f, g.

# The Alternating Direction Method of Multipliers (ADMM)

- To overcome the drawbacks of dual decomposition, we introduce another iterative method: the ADMM. We are going to develop it for a general case and then focus it on our distributed problem
- Consider the problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^p} \qquad f(\mathbf{x}) + g(\mathbf{y})$$
  
subject to 
$$A\mathbf{x} + B\mathbf{y} = c.$$

For well-defined matrices and vectors, as well as convex functions f, g.

• Define the augmented Lagrangian,

$$L(\mathbf{x}, \mathbf{y}, \lambda) := f(\mathbf{x}) + g(\mathbf{y}) + \lambda^{\top} (A\mathbf{x} + B\mathbf{y} - c) + \frac{\beta}{2} ||A\mathbf{x} + B\mathbf{y} - c||^{2},$$

for any scalar  $\beta > 0$ .

• Then the ADMM does: start with an initial value for  $x_0$ ,  $y_0$ ,  $\lambda_0$ ,

- Then the ADMM does: start with an initial value for  $x_0$ ,  $y_0$ ,  $\lambda_0$ ,
- Iterate:

$$\begin{aligned} \mathbf{x}_{k+1} &= & \arg\min_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \mathbf{y}_k, \lambda_k) \\ \mathbf{y}_{k+1} &= & \arg\min_{\mathbf{y} \in \mathbb{R}^p} L(\mathbf{x}_{k+1}, \mathbf{y}, \lambda_k) \\ \lambda_{k+1} &= & \lambda_{k+1} + \beta(A\mathbf{x}_{k+1} + B\mathbf{y}_{k+1} - c) \end{aligned}$$

- Then the ADMM does: start with an initial value for  $x_0$ ,  $y_0$ ,  $\lambda_0$ ,
- Iterate:

$$\begin{aligned} \mathbf{x}_{k+1} &= & \arg\min_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \mathbf{y}_k, \lambda_k) \\ \mathbf{y}_{k+1} &= & \arg\min_{\mathbf{y} \in \mathbb{R}^p} L(\mathbf{x}_{k+1}, \mathbf{y}, \lambda_k) \\ \lambda_{k+1} &= & \lambda_{k+1} + \beta(A\mathbf{x}_{k+1} + B\mathbf{y}_{k+1} - c) \end{aligned}$$

 ADMM is similar to an incremental approach, but on two variables. It is also similar to dual decomposition, but on two variables

- Then the ADMM does: start with an initial value for  $x_0$ ,  $y_0$ ,  $\lambda_0$ ,
- Iterate:

```
\begin{aligned} \mathbf{x}_{k+1} &= \arg\min_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \mathbf{y}_k, \lambda_k) \\ \mathbf{y}_{k+1} &= \arg\min_{\mathbf{y} \in \mathbb{R}^p} L(\mathbf{x}_{k+1}, \mathbf{y}, \lambda_k) \\ \lambda_{k+1} &= \lambda_{k+1} + \beta(A\mathbf{x}_{k+1} + B\mathbf{y}_{k+1} - c) \end{aligned}
```

- ADMM is similar to an incremental approach, but on two variables. It is also similar to dual decomposition, but on two variables
- It looks more complicated than dual ascent, is it better in some sense?

- Then the ADMM does: start with an initial value for  $x_0$ ,  $y_0$ ,  $\lambda_0$ ,
- Iterate:

$$\mathbf{x}_{k+1} = \arg\min_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \mathbf{y}_k, \lambda_k)$$

$$\mathbf{y}_{k+1} = \arg\min_{\mathbf{y} \in \mathbb{R}^p} L(\mathbf{x}_{k+1}, \mathbf{y}, \lambda_k)$$

$$\lambda_{k+1} = \lambda_{k+1} + \beta(A\mathbf{x}_{k+1} + B\mathbf{y}_{k+1} - c)$$

- ADMM is similar to an incremental approach, but on two variables. It is also similar to dual decomposition, but on two variables
- It looks more complicated than dual ascent, is it better in some sense?
- Assume strong duality holds and primal and dual solutions exist, then convergence is ensure for any  $\beta > 0$  in a weak sense.

• Assume f to be m-strongly convex and L-smooth and that a primal-dual solution exists (g can be a generic convex function). This is a special case and ADMM generally converges with minimal assumptions, but in a much weaker sense, so we keep here the stronger assumptions.

- Assume f to be m-strongly convex and L-smooth and that a primal-dual solution exists (g can be a generic convex function). This is a special case and ADMM generally converges with minimal assumptions, but in a much weaker sense, so we keep here the stronger assumptions.
- Assume that A is full row rank.

- Assume f to be m-strongly convex and L-smooth and that a primal-dual solution exists (g can be a generic convex function). This is a special case and ADMM generally converges with minimal assumptions, but in a much weaker sense, so we keep here the stronger assumptions.
- Assume that A is full row rank.
- Then ADMM converges linearly for any step size  $\beta > 0$ .

- Assume f to be m-strongly convex and L-smooth and that a primal-dual solution exists (g can be a generic convex function). This is a special case and ADMM generally converges with minimal assumptions, but in a much weaker sense, so we keep here the stronger assumptions.
- Assume that A is full row rank.
- Then ADMM converges linearly for any step size  $\beta > 0$ .

#### Theorem 7 (ADMM convergence)

With the assumptions in place,

$$\|\lambda_k - \lambda^*\| \leq \varrho^k \|\lambda_0 - \lambda^*\| \qquad \varrho = \max\left\{ \left| \frac{1 - \beta \frac{\sigma_{\max}^2(A)}{m}}{1 + \beta \frac{\sigma_{\max}^2(A)}{m}} \right|, \left| \frac{1 - \beta \frac{\sigma_{\min}^2(A)}{L}}{1 + \beta \frac{\sigma_{\min}^2(A)}{L}} \right| \right\}$$

and.

$$\|\mathbf{x}_k - \mathbf{x}^*\| \le \frac{\|A\|}{m} \|\lambda_k - \lambda^*\|.$$

 ADMM is extremely popular, for its loose assumptions and the fact that it works in many contexts, even when we have no idea why (non-convex, combinatorial, etc..)

- ADMM is extremely popular, for its loose assumptions and the fact that it works in many contexts, even when we have no idea why (non-convex, combinatorial, etc..)
- Hence, we have many ways to prove its convergence. Here, I want to give a glimpse of an operator-splitting technique.

- ADMM is extremely popular, for its loose assumptions and the fact that it works in many contexts, even when we have no idea why (non-convex, combinatorial, etc..)
- Hence, we have many ways to prove its convergence. Here, I want to give a glimpse of an operator-splitting technique.
- The proof goes as follows: ADMM is an application of a special algorithm (the Douglas-Rachford splitting) applied to the dual of our initial problem.
   The Douglas-Rachford splitting converges in a certain way given functional properties. We derive the dual of those functional properties and (as in the dual decomposition case) the convergence of the dual algorithm (ADMM).

- ADMM is extremely popular, for its loose assumptions and the fact that it works in many contexts, even when we have no idea why (non-convex, combinatorial, etc..)
- Hence, we have many ways to prove its convergence. Here, I want to give a glimpse of an operator-splitting technique.
- The proof goes as follows: ADMM is an application of a special algorithm (the Douglas-Rachford splitting) applied to the dual of our initial problem.
   The Douglas-Rachford splitting converges in a certain way given functional properties. We derive the dual of those functional properties and (as in the dual decomposition case) the convergence of the dual algorithm (ADMM).
- This is why the result looks very similar to the result of the dual decomposition. ADMM is a dual algorithm.

## Proof: step I, Douglas-Rachford splitting

• Consider the problem,

$$\min_{x \in \mathbb{R}^n} f(x) + g(x),$$

with f and g convex closed and proper (CCP). Consider now the following method to find a solution  $x^*$ .

# Proof: step I, Douglas-Rachford splitting

• Consider the problem,

$$\min_{x \in \mathbb{R}^n} f(x) + g(x),$$

with f and g convex closed and proper (CCP). Consider now the following method to find a solution  $x^*$ .

• Start at a certain  $z_0$  and iterate for all  $k \in \mathbb{N}$ :

$$x_k = \operatorname{prox}_{\beta f}(z_k) \tag{6a}$$

$$z_{k+1} = z_k + \text{prox}_{\beta q}(2x_k - z_k) - x_k,$$
 (6b)

where  $prox_{\beta\phi}$  is the usual prox operator:

$$\operatorname{prox}_{\beta\phi}(u) = \arg\min_{v} \{\phi(v) + \frac{1}{2\beta} \|v - u\|^2\}$$

The method is called the Douglas-Rachford splitting and it converges in some defined sense.

## Proof: step I, Douglas-Rachford splitting

• Consider the problem,

$$\min_{x \in \mathbb{R}^n} f(x) + g(x),$$

with f and g convex closed and proper (CCP). Consider now the following method to find a solution  $x^*$ .

• Start at a certain  $z_0$  and iterate for all  $k \in \mathbb{N}$ :

$$x_k = \operatorname{prox}_{\mathcal{B}f}(z_k) \tag{6a}$$

$$z_{k+1} = z_k + \text{prox}_{\beta g}(2x_k - z_k) - x_k,$$
 (6b)

where  $prox_{\beta\phi}$  is the usual prox operator:

$$\operatorname{prox}_{\beta\phi}(u) = \arg\min_{v} \{\phi(v) + \frac{1}{2\beta} \|v - u\|^2\}$$

The method is called the Douglas-Rachford splitting and it converges in some defined sense.

• In particular if f is m-strongly convex and L-smooth, then, for all  $\beta>0$ 

$$||z_{k+1} - z^*|| \le \varrho^k ||z_k - z^*||, \qquad \varrho = \max \left\{ \left| \frac{1 - \beta L}{1 + \beta L} \right|, \left| \frac{1 - \beta m}{1 + \beta m} \right| \right\}$$

• Start from our problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^p} \qquad f(\mathbf{x}) + g(\mathbf{y})$$
subject to 
$$A\mathbf{x} + B\mathbf{y} = c,$$

form the Lagrangian and take the dual

• Start from our problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^p} \quad f(\mathbf{x}) + g(\mathbf{y})$$
subject to 
$$A\mathbf{x} + B\mathbf{y} = c,$$

form the Lagrangian and take the dual

Dual problem is,

$$\max_{\lambda} \inf_{\mathbf{x}, \mathbf{y}} \{ f(\mathbf{x}) + g(\mathbf{y}) + \lambda^{\top} (A\mathbf{x} + B\mathbf{y} - c) \}$$

and so.

$$\min_{\lambda} f^{\star}(-A^{\top}\lambda) + g^{\star}(-A^{\top}\lambda) + c^{\top}\lambda$$

• Start from our problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^p} \qquad f(\mathbf{x}) + g(\mathbf{y})$$
  
subject to 
$$A\mathbf{x} + B\mathbf{y} = c,$$

form the Lagrangian and take the dual

• Dual problem is,

$$\max_{\lambda} \inf_{\mathbf{x}, \mathbf{y}} \{ f(\mathbf{x}) + g(\mathbf{y}) + \lambda^{\top} (A\mathbf{x} + B\mathbf{y} - c) \}$$

and so,

$$\min_{\lambda} f^{\star}(-A^{\top}\lambda) + g^{\star}(-A^{\top}\lambda) + c^{\top}\lambda$$

• Consider "f" =  $f^*(-A^T\lambda) + c^T\lambda$  and "g" =  $g^*(-B^T\lambda)$ 

• Start from our problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^p} \qquad f(\mathbf{x}) + g(\mathbf{y})$$
subject to 
$$A\mathbf{x} + B\mathbf{y} = c,$$

form the Lagrangian and take the dual

• Dual problem is,

$$\max_{\lambda} \inf_{\mathbf{x}, \mathbf{y}} \{ f(\mathbf{x}) + g(\mathbf{y}) + \lambda^{\top} (A\mathbf{x} + B\mathbf{y} - c) \}$$

and so,

$$\min_{\lambda} f^{\star}(-A^{\top}\lambda) + g^{\star}(-A^{\top}\lambda) + c^{\top}\lambda$$

- Consider "f" =  $f^*(-A^T\lambda) + c^T\lambda$  and "g" =  $g^*(-B^T\lambda)$
- $\bullet$  As before we know that "f" is  $\sigma_{\min}^2/L\text{-strongly}$  convex and  $\sigma_{\max}^2/m\text{-smooth}$

## Proof: step III, applying DR to the dual

• Apply DRS "f" =  $f^*(-A^T\lambda) + c^T\lambda$  and "g" =  $g^*(-B^T\lambda)$  and use its convergence properties.

## Proof: step III, applying DR to the dual

- Apply DRS "f" =  $f^*(-A^T\lambda) + c^T\lambda$  and "g" =  $g^*(-B^T\lambda)$  and use its convergence properties.
- Do some computational gymnastic: check out the references

## Proof: step III, applying DR to the dual

- Apply DRS "f" =  $f^*(-A^T\lambda) + c^T\lambda$  and "g" =  $g^*(-B^T\lambda)$  and use its convergence properties.
- Do some computational gymnastic: check out the references
- Arrive at the ADMM iterations

• Let's go back to:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^p} \quad f(\mathbf{x}) + g(\mathbf{y})$$
subject to 
$$A\mathbf{x} + B\mathbf{y} = c,$$

and our problem,

$$\begin{aligned} & \min_{\mathbf{x}^i \in \mathbb{R}^n, \mathbf{y}^{ij} \in \mathbb{R}^{|\mathcal{E}|}} & & \sum_{i=1}^N f_i(\mathbf{x}^i) \\ & \text{subject to} & & \mathbf{x}^i = \mathbf{y}^{ij}, & \forall i \sim j. \end{aligned}$$

• Let's go back to:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^p} \quad f(\mathbf{x}) + g(\mathbf{y})$$
subject to 
$$A\mathbf{x} + B\mathbf{y} = c,$$

and our problem,

$$\min_{\mathbf{x}^i \in \mathbb{R}^n, \mathbf{y}^{ij} \in \mathbb{R}^{|\mathcal{E}|}} \qquad \sum_{i=1}^N f_i(\mathbf{x}^i) \\
\text{subject to} \qquad \mathbf{x}^i = \mathbf{y}^{ij}, \qquad \forall i \sim j.$$

• Think a bit at what is written. This is of course not the only possibility. Different splittings yield different properties. This will give a distributed algorithm similar to the one we have seen for dual decomposition.

• Let's go back to:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^p} \quad f(\mathbf{x}) + g(\mathbf{y})$$
subject to 
$$A\mathbf{x} + B\mathbf{y} = c,$$

and our problem,

$$\min_{\mathbf{x}^i \in \mathbb{R}^n, y^{ij} \in \mathbb{R}^{|\mathcal{E}|}} \qquad \sum_{i=1}^N f_i(\mathbf{x}^i) \\
\text{subject to} \qquad \mathbf{x}^i = y^{ij}, \qquad \forall i \sim j.$$

- Think a bit at what is written. This is of course not the only possibility. Different splittings yield different properties. This will give a distributed algorithm similar to the one we have seen for dual decomposition.
- Later, we will see a different cloud-based splitting instead.

• Let's go back to:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^p} f(\mathbf{x}) + g(\mathbf{y})$$
  
subject to 
$$A\mathbf{x} + B\mathbf{y} = c,$$

and our problem,

$$\min_{\mathbf{x}^i \in \mathbb{R}^n, \mathbf{y}^{ij} \in \mathbb{R}^{|\mathcal{E}|}} \qquad \sum_{i=1}^N f_i(\mathbf{x}^i) \\
\text{subject to} \qquad \mathbf{x}^i = \mathbf{y}^{ij}, \qquad \forall i \sim j.$$

- Think a bit at what is written. This is of course not the only possibility. Different splittings yield different properties. This will give a distributed algorithm similar to the one we have seen for dual decomposition.
- Later, we will see a different cloud-based splitting instead.
- Compactify  $\mathbf{x}^i = \mathbf{y}^{ij}$  as,  $A\mathbf{x} + B\mathbf{y} = 0$ . Careful here that A is not full row rank, so linear convergence requires more work, but possible.

• Then you apply ADMM and simplify,

- Then you apply ADMM and simplify,
- For all nodes update:

$$\mathbf{x}_{k+1}^{i} = \arg\min_{\mathbf{x}^{i}} \{f_{i}(\mathbf{x}) + (\lambda_{k})^{\top} (A\mathbf{x} + B\mathbf{y}_{k}) + \frac{\beta}{2} \|A\mathbf{x} + B\mathbf{y}_{k}\|^{2} \}$$

which is equivalent to,

$$\mathbf{x}_{k+1}^{i} = \arg\min_{\mathbf{x}^{i}} \{f_{i}(\mathbf{x}^{i}) + \sum_{i \in \mathcal{N}^{i}} \frac{\beta}{2} \|\mathbf{x}^{i} - \mathbf{y}_{k}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} \}$$

- Then you apply ADMM and simplify,
- For all nodes update:

$$\mathbf{x}_{k+1}^{i} = \arg\min_{\mathbf{x}^{i}} \{f_{i}(\mathbf{x}) + (\lambda_{k})^{\top} (A\mathbf{x} + B\mathbf{y}_{k}) + \frac{\beta}{2} \|A\mathbf{x} + B\mathbf{y}_{k}\|^{2} \}$$

which is equivalent to,

$$\mathbf{x}_{k+1}^{i} = \arg\min_{\mathbf{x}^{i}} \{f_{i}(\mathbf{x}^{i}) + \sum_{i \in \mathcal{N}^{i}} \frac{\beta}{2} \|\mathbf{x}^{i} - \mathbf{y}_{k}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} \}$$

• Communicate  $\mathbf{x}_{k+1}^{i}$  to your neighbors

- Then you apply ADMM and simplify,
- For all nodes update:

$$\mathbf{x}_{k+1}^{i} = \arg\min_{\mathbf{x}^{i}} \{f_{i}(\mathbf{x}) + (\lambda_{k})^{\top} (A\mathbf{x} + B\mathbf{y}_{k}) + \frac{\beta}{2} \|A\mathbf{x} + B\mathbf{y}_{k}\|^{2} \}$$

which is equivalent to,

$$\mathbf{x}_{k+1}^{i} = \arg\min_{\mathbf{x}^{i}} \{f_{i}(\mathbf{x}^{i}) + \sum_{j \in \mathcal{N}^{i}} \frac{\beta}{2} \|\mathbf{x}^{i} - \mathbf{y}_{k}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} \}$$

- Communicate  $x_{k+1}^i$  to your neighbors
- Each node,

$$y_{k+1}^{ij} = \arg\min_{y^{ij}} \{ (\lambda_k)^{\top} (A_{k+1} + B_y) + \frac{\beta}{2} ||A_{k+1} + B_y||^2 \}$$

which is equivalent to,

$$y_{k+1}^{ij} = \arg\min_{\mathbf{y}^{ij}} \{ \frac{\beta}{2} \| \mathbf{x}_{k+1}^{i} - \mathbf{y}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} + \frac{\beta}{2} \| \mathbf{x}_{k+1}^{j} - \mathbf{y}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} \}$$

# ADMM: distributed algorithm cont'ed

•

$$\mathbf{y}_{k+1}^{ij} = \arg\min_{\mathbf{y}^{ij}} \{ \frac{\beta}{2} \| \mathbf{x}_{k+1}^{i} - \mathbf{y}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} + \frac{\beta}{2} \| \mathbf{x}_{k+1}^{j} - \mathbf{y}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} \}$$

which is

$$y_{k+1}^{ij} = \frac{\mathbf{x}_{k+1}^{i} + \mathbf{x}_{k+1}^{j}}{2} + \frac{\lambda_{k}^{ij}}{\beta}$$

## ADMM: distributed algorithm cont'ed

•

$$y_{k+1}^{ij} = \arg\min_{y^{ij}} \{ \frac{\beta}{2} \| \mathbf{x}_{k+1}^{i} - \mathbf{y}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} + \frac{\beta}{2} \| \mathbf{x}_{k+1}^{j} - \mathbf{y}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} \}$$

which is

$$y_{k+1}^{ij} = \frac{\mathbf{x}_{k+1}^{i} + \mathbf{x}_{k+1}^{j}}{2} + \frac{\lambda_{k}^{ij}}{\beta}$$

• Each node.

$$\lambda_{k+1}^{ij} = \lambda_k^{ij} + \beta(\mathbf{x}_{k+1}^i + \mathbf{x}_{k+1}^j - 2\mathbf{y}_{k+1}^{ij}) = 0$$

# ADMM: distributed algorithm cont'ed

•

$$y_{k+1}^{ij} = \arg\min_{y^{ij}} \{ \frac{\beta}{2} \| \mathbf{x}_{k+1}^{i} - \mathbf{y}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} + \frac{\beta}{2} \| \mathbf{x}_{k+1}^{j} - \mathbf{y}^{ij} + \frac{\lambda_{k}^{ij}}{\beta} \|^{2} \}$$

which is

$$y_{k+1}^{ij} = \frac{\mathbf{x}_{k+1}^{i} + \mathbf{x}_{k+1}^{j}}{2} + \frac{\lambda_{k}^{ij}}{\beta}$$

• Each node,

$$\lambda_{k+1}^{ij} = \lambda_k^{ij} + \beta(\mathbf{x}_{k+1}^i + \mathbf{x}_{k+1}^j - 2\mathbf{y}_{k+1}^{ij}) = 0$$

- Start with  $\lambda_0^{ij} = 0$ , and  $y_0^{ij}$  then,
- **②** Each node updates,  $\mathbf{x}_{k+1}^i = \arg\min_{\mathbf{x}^i} \{ f_i(\mathbf{x}^i) + \sum_{j \in \mathcal{N}^i} \frac{\beta}{2} \|\mathbf{x}^i \mathbf{y}_k^{ij}\|^2 \}$
- **Solution** Each node communicates  $\mathbf{x}_{k+1}^{i}$  to its neighbors
- **1** Each node updates,  $y_{k+1}^{ij} = \frac{\mathbf{x}_{k+1}^i + \mathbf{x}_{k+1}^j}{2}$

Consider the task of training a model via the convex problem,

$$\min_{\mathbf{x}\in\mathbb{R}^n}\ell(A\mathbf{x}-b)+r(\mathbf{x}),$$

where  $\ell : \mathbf{R}^p \to \mathbf{R}$  is a convex loss function,  $A \in \mathbf{R}^{p \times n}$  is the feature matrix,  $b \in \mathbf{R}^p$  is the output vector,  $\mathbf{x}$  are the parameters of the model, and  $r : \mathbf{R}^n \to \mathbf{R}$  is a convex regularization function.

• Consider the task of training a model via the convex problem,

$$\min_{\mathbf{x}\in\mathbb{R}^n}\ell(A\mathbf{x}-b)+r(\mathbf{x}),$$

where  $\ell: \mathbf{R}^p \to \mathbf{R}$  is a convex loss function,  $A \in \mathbf{R}^{p \times n}$  is the feature matrix,  $b \in \mathbf{R}^p$  is the output vector,  $\mathbf{x}$  are the parameters of the model, and  $r: \mathbf{R}^n \to \mathbf{R}$  is a convex regularization function.

ullet Assume that  $\ell$  is additive across training examples, as

$$\ell(\mathbf{x}) = \sum_{i=1}^{m} \ell_i (a_i^{\top} \mathbf{x} - b_i).$$

• Consider the task of training a model via the convex problem,

$$\min_{\mathbf{x}\in\mathbb{R}^n}\boldsymbol{\ell}(A\mathbf{x}-b)+r(\mathbf{x}),$$

where  $\ell: \mathbf{R}^p \to \mathbf{R}$  is a convex loss function,  $A \in \mathbf{R}^{p \times n}$  is the feature matrix,  $b \in \mathbf{R}^p$  is the output vector,  $\mathbf{x}$  are the parameters of the model, and  $r: \mathbf{R}^n \to \mathbf{R}$  is a convex regularization function.

ullet Assume that  $\ell$  is additive across training examples, as

$$\ell(\mathbf{x}) = \sum_{i=1}^{m} \ell_i (a_i^{\top} \mathbf{x} - b_i).$$

• r is often separable too. For instance if r is the Tikhonov regularization (aka ridge penalty):  $r(\mathbf{x}) = \nu ||\mathbf{x}||_2^2$ , or if r is the lasso penalty:  $r(\mathbf{x}) = \nu ||\mathbf{x}||_1^2$ 

• Consider the task of training a model via the convex problem,

$$\min_{\mathbf{x}\in\mathbb{R}^n}\boldsymbol{\ell}(A\mathbf{x}-b)+r(\mathbf{x}),$$

where  $\ell : \mathbf{R}^p \to \mathbf{R}$  is a convex loss function,  $A \in \mathbf{R}^{p \times n}$  is the feature matrix,  $b \in \mathbf{R}^p$  is the output vector,  $\mathbf{x}$  are the parameters of the model, and  $r : \mathbf{R}^n \to \mathbf{R}$  is a convex regularization function.

ullet Assume that  $\ell$  is additive across training examples, as

$$\ell(\mathbf{x}) = \sum_{i=1}^{m} \ell_i (a_i^{\top} \mathbf{x} - b_i).$$

- r is often separable too. For instance if r is the Tikhonov regularization (aka ridge penalty):  $r(\mathbf{x}) = \nu \|\mathbf{x}\|_2^2$ , or if r is the lasso penalty:  $r(\mathbf{x}) = \nu \|\mathbf{x}\|_1$
- Typically you have a modest number of features but a very large number of training examples (i.e.,  $m \gg n$ ).

• The goal is to solve the problem in a distributed way, with each processor handling a subset of the training data. This is useful either when there are so many training examples that it is inconvenient or impossible to process them on a single machine or when the data is naturally collected or stored in a distributed fashion. This includes, for example, online social network data, webserver access logs, wireless sensor networks, and many cloud computing applications more generally.

- The goal is to solve the problem in a distributed way, with each processor handling a subset of the training data. This is useful either when there are so many training examples that it is inconvenient or impossible to process them on a single machine or when the data is naturally collected or stored in a distributed fashion. This includes, for example, online social network data, webserver access logs, wireless sensor networks, and many cloud computing applications more generally.
- Suppose you associate a number of training examples to a number N of processors (in the extreme N=m). Then, you may solve

$$\min_{\mathbf{x}^{j} \in \mathbb{R}^{n}, \mathbf{y} \in \mathbb{R}^{n}} \qquad \sum_{j=1}^{N} \ell_{j} (A_{j} \mathbf{x}^{j} - b_{j}) + r(\mathbf{y}),$$
subject to 
$$\mathbf{x}^{j} = \mathbf{y}, \quad i = 1, \dots, N.$$

- The goal is to solve the problem in a distributed way, with each processor handling a subset of the training data. This is useful either when there are so many training examples that it is inconvenient or impossible to process them on a single machine or when the data is naturally collected or stored in a distributed fashion. This includes, for example, online social network data, webserver access logs, wireless sensor networks, and many cloud computing applications more generally.
- Suppose you associate a number of training examples to a number N of processors (in the extreme N=m). Then, you may solve

$$\min_{\mathbf{x}^{j} \in \mathbb{R}^{n}, \mathbf{y} \in \mathbb{R}^{n}} \qquad \sum_{j=1}^{N} \ell_{j} (A_{j} \mathbf{x}^{j} - b_{j}) + r(\mathbf{y}),$$
subject to 
$$\mathbf{x}^{j} = \mathbf{y}, \quad i = 1, \dots, N.$$

This can be solved via ADMM

Each processor solves

$$\begin{aligned} \mathbf{x}_{k+1}^{j} &= \arg\min_{\mathbf{x}^{j}} \{ \ell_{j} (A_{j} \mathbf{x}^{j} - b^{j}) + \lambda_{k}^{j} (\mathbf{x}^{j} - \mathbf{y}_{k}) + \frac{\beta}{2} \| \mathbf{x}^{j} - \mathbf{y}_{k} \|^{2} \} \\ &= \arg\min_{\mathbf{x}^{j}} \{ \ell_{j} (A_{j} \mathbf{x}^{j} - b^{j}) + \frac{\beta}{2} \| \mathbf{x}^{j} - \mathbf{y}_{k} + \frac{\lambda_{k}^{j}}{\beta} \|^{2} \} \end{aligned}$$

Each processor solves

$$\begin{aligned} \mathbf{x}_{k+1}^{j} &= \arg\min_{\mathbf{x}^{j}} \{ \ell_{j} (A_{j} \mathbf{x}^{j} - b^{j}) + \lambda_{k}^{j} (\mathbf{x}^{j} - \mathbf{y}_{k}) + \frac{\beta}{2} \| \mathbf{x}^{j} - \mathbf{y}_{k} \|^{2} \} \\ &= \arg\min_{\mathbf{x}^{j}} \{ \ell_{j} (A_{j} \mathbf{x}^{j} - b^{j}) + \frac{\beta}{2} \| \mathbf{x}^{j} - \mathbf{y}_{k} + \frac{\lambda_{k}^{j}}{\beta} \|^{2} \} \end{aligned}$$

• Communication to the cloud of  $x_{k+1}^{j}$ 

Each processor solves

$$\begin{aligned} \mathbf{x}_{k+1}^{j} &= \arg\min_{\mathbf{x}^{j}} \{ \ell_{j} (A_{j} \mathbf{x}^{j} - b^{j}) + \lambda_{k}^{j} (\mathbf{x}^{j} - \mathbf{y}_{k}) + \frac{\beta}{2} \| \mathbf{x}^{j} - \mathbf{y}_{k} \|^{2} \} \\ &= \arg\min_{\mathbf{x}^{j}} \{ \ell_{j} (A_{j} \mathbf{x}^{j} - b^{j}) + \frac{\beta}{2} \| \mathbf{x}^{j} - \mathbf{y}_{k} + \frac{\lambda_{k}^{j}}{\beta} \|^{2} \} \end{aligned}$$

- Communication to the cloud of  $x_{k+1}^{j}$
- On the cloud we solve.

$$y_{k+1} = \arg\min_{\mathbf{y}} \{ r(\mathbf{y}) + \sum_{j=1}^{N} \lambda_{k}^{j} (\mathbf{x}_{k+1}^{j} - \mathbf{y}) + \frac{\beta}{2} \| \mathbf{x}_{k+1}^{j} - \mathbf{y} \|^{2} \}$$

$$= \arg\min_{\mathbf{y}} \{ r(\mathbf{y}) + \sum_{j=1}^{N} \frac{\beta}{2} \| \mathbf{x}_{k+1}^{j} - \mathbf{y} - \frac{\lambda_{k}^{j}}{\beta} \|^{2} \}$$

Each processor solves

$$\begin{aligned} \mathbf{x}_{k+1}^{j} &= \arg\min_{\mathbf{x}^{j}} \{ \ell_{j} (A_{j} \mathbf{x}^{j} - b^{j}) + \lambda_{k}^{j} (\mathbf{x}^{j} - \mathbf{y}_{k}) + \frac{\beta}{2} \| \mathbf{x}^{j} - \mathbf{y}_{k} \|^{2} \} \\ &= \arg\min_{\mathbf{x}^{j}} \{ \ell_{j} (A_{j} \mathbf{x}^{j} - b^{j}) + \frac{\beta}{2} \| \mathbf{x}^{j} - \mathbf{y}_{k} + \frac{\lambda_{k}^{j}}{\beta} \|^{2} \} \end{aligned}$$

- Communication to the cloud of  $x_{k+1}^{j}$
- On the cloud we solve.

$$\begin{aligned} y_{k+1} &= \arg\min_{\mathbf{y}} \{ r(\mathbf{y}) + \sum_{j=1}^{N} \lambda_{k}^{j} (\mathbf{x}_{k+1}^{j} - \mathbf{y}) + \frac{\beta}{2} \| \mathbf{x}_{k+1}^{j} - \mathbf{y} \|^{2} \} \\ &= \arg\min_{\mathbf{y}} \{ r(\mathbf{y}) + \sum_{i=1}^{N} \frac{\beta}{2} \| \mathbf{x}_{k+1}^{j} - \mathbf{y} - \frac{\lambda_{k}^{j}}{\beta} \|^{2} \} \end{aligned}$$

• Update  $\lambda_{k+1}^j = \lambda_k^j + \beta(\mathbf{x}_{k+1}^j - \mathbf{y}_{k+1})$  and communicate back to processors  $\mathbf{y}_{k+1}, \lambda_{k+1}^j$ 

• The data never leaves the processors

- The data never leaves the processors
- Different updates of y depending on r and on x depending on  $\ell$ .

- The data never leaves the processors
- Different updates of y depending on r and on x depending on  $\ell$ .
- Consider for example:  $\ell(\cdot) = \|\cdot\|_2^2$ ,  $r(y) = \nu \|y\|_1$ , then we obtain the algorithm,

- The data never leaves the processors
- Different updates of y depending on r and on x depending on  $\ell$ .
- Consider for example:  $\ell(\cdot) = \|\cdot\|_2^2$ ,  $r(y) = \nu \|y\|_1$ , then we obtain the algorithm,
- Each processor solves

$$\mathbf{x}_{k+1}^{j} = \arg\min_{\mathbf{x}^{j}} \{ \|A_{j}\mathbf{x}^{j} - b_{j}\|^{2} \} + \frac{\beta}{2} \|\mathbf{x}^{j} - \mathbf{y}_{k} + \frac{\lambda_{k}^{j}}{\beta} \|^{2} \}$$
$$= (A_{j}^{T}A_{j} + \beta I)^{-1} (A_{j}^{T}b_{j} - \beta \mathbf{y}_{k} + \lambda_{k}^{j}).$$

- The data never leaves the processors
- Different updates of y depending on r and on x depending on  $\ell$ .
- Consider for example:  $\ell(\cdot) = \|\cdot\|_2^2$ ,  $r(y) = \nu \|y\|_1$ , then we obtain the algorithm,
- Each processor solves

$$\mathbf{x}_{k+1}^{j} = \arg\min_{\mathbf{x}^{j}} \{ \|A_{j}\mathbf{x}^{j} - b_{j}\|^{2} \} + \frac{\beta}{2} \|\mathbf{x}^{j} - \mathbf{y}_{k} + \frac{\lambda_{k}^{j}}{\beta} \|^{2} \}$$
$$= (A_{j}^{T}A_{j} + \beta I)^{-1} (A_{j}^{T}b_{j} - \beta \mathbf{y}_{k} + \lambda_{k}^{j}).$$

ullet Communication to the cloud of  ${oldsymbol x}_{k+1}^j$ 

- The data never leaves the processors
- Different updates of y depending on r and on x depending on  $\ell$ .
- Consider for example:  $\ell(\cdot) = \|\cdot\|_2^2$ ,  $r(y) = \nu \|y\|_1$ , then we obtain the algorithm,
- Each processor solves

$$\mathbf{x}_{k+1}^{j} = \arg\min_{\mathbf{x}^{j}} \{ \|A_{j}\mathbf{x}^{j} - b_{j}\|^{2} \} + \frac{\beta}{2} \|\mathbf{x}^{j} - \mathbf{y}_{k} + \frac{\lambda_{k}^{j}}{\beta} \|^{2} \}$$
$$= (A_{j}^{T}A_{j} + \beta I)^{-1} (A_{j}^{T}b_{j} - \beta \mathbf{y}_{k} + \lambda_{k}^{j}).$$

- Communication to the cloud of  $\mathbf{x}_{k+1}^{j}$
- On the cloud we solve,

$$y_{k+1} = \arg\min_{y} \{\nu \|y\|_{1} + \sum_{i=1}^{N} \frac{\beta}{2} \|\mathbf{x}_{k+1}^{j} - \mathbf{y} + \frac{\lambda_{k}^{j}}{\beta} \|^{2} \} = S_{\nu/\beta}(\overline{\mathbf{x}}_{k+1} + \frac{\overline{\lambda}_{k}}{\beta})$$

- The data never leaves the processors
- Different updates of y depending on r and on x depending on  $\ell$ .
- Consider for example:  $\ell(\cdot) = \|\cdot\|_2^2$ ,  $r(y) = \nu \|y\|_1$ , then we obtain the algorithm,
- Each processor solves

$$\mathbf{x}_{k+1}^{j} = \arg\min_{\mathbf{x}^{j}} \{ \|A_{j}\mathbf{x}^{j} - b_{j}\|^{2} \} + \frac{\beta}{2} \|\mathbf{x}^{j} - \mathbf{y}_{k} + \frac{\lambda_{k}^{j}}{\beta} \|^{2} \}$$
$$= (A_{j}^{T} A_{j} + \beta I)^{-1} (A_{j}^{T} b_{j} - \beta \mathbf{y}_{k} + \lambda_{k}^{j}).$$

- ullet Communication to the cloud of  ${f x}_{k+1}^j$
- On the cloud we solve,

$$y_{k+1} = \arg\min_{y} \{\nu \|y\|_{1} + \sum_{j=1}^{N} \frac{\beta}{2} \|\mathbf{x}_{k+1}^{j} - \mathbf{y} + \frac{\lambda_{k}^{j}}{\beta} \|^{2} \} = S_{\nu/\beta}(\overline{\mathbf{x}}_{k+1} + \frac{\overline{\lambda}_{k}}{\beta})$$

• Update  $\lambda_{k+1}^j = \lambda_k^j + \beta(\mathbf{x}_{k+1}^j - \mathbf{y}_{k+1})$  and communicate back to processors  $\mathbf{y}_{k+1}, \lambda_{k+1}^j$ 

#### Communication issues

• In the description above, we have used bi-direction synchronous communication, this may not be the case in reality

#### Communication issues

- In the description above, we have used bi-direction synchronous communication, this may not be the case in reality
- Asynchronicity, latencies, package losses, all add to the difficulty in proving convergence of the algorithm

#### Communication issues

- In the description above, we have used bi-direction synchronous communication, this may not be the case in reality
- Asynchronicity, latencies, package losses, all add to the difficulty in proving convergence of the algorithm
- Research in this domain is very rich and active!

• We have looked at distributed optimization in the dual domain, with two "classical" algorithms: dual decomposition and the ADMM

- We have looked at distributed optimization in the dual domain, with two "classical" algorithms: dual decomposition and the ADMM
- We have studied their convergence, which seem better (in some sense) than primal methods

- We have looked at distributed optimization in the dual domain, with two "classical" algorithms: dual decomposition and the ADMM
- We have studied their convergence, which seem better (in some sense) than primal methods
- ADMM in particular offers multiple ways to distributed computations across multiple devices and it has received a lot of attention in recent years

- We have looked at distributed optimization in the dual domain, with two "classical" algorithms: dual decomposition and the ADMM
- We have studied their convergence, which seem better (in some sense) than primal methods
- ADMM in particular offers multiple ways to distributed computations across multiple devices and it has received a lot of attention in recent years
- ADMM works by splitting the problem into a part that can be solved locally, and a part that we can afford to solve sharing information

- We have looked at distributed optimization in the dual domain, with two "classical" algorithms: dual decomposition and the ADMM
- We have studied their convergence, which seem better (in some sense) than primal methods
- ADMM in particular offers multiple ways to distributed computations across multiple devices and it has received a lot of attention in recent years
- ADMM works by splitting the problem into a part that can be solved locally, and a part that we can afford to solve sharing information
- ADMM can be applied to many settings (multi-core, cloud-computing, distributed computing, etc..)

#### Sample references

- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein, Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers, Foundations and Trend in Machine Learning, 2010
- Ernest Ryu, Stephen Boyd, A Primer on Monotone Operator Methods, Appl. Comput. Math., 2016
- Many variants out there.