

# Distributed Autonomous Systems M

Dante Piotto

spring semester 2024



# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction and scenarios</b>                                   | <b>7</b>  |
| 1.1      | Distributed Autonomous System . . . . .                             | 7         |
| 1.2      | Scenarios and applications of distributed systems . . . . .         | 7         |
| 1.3      | Measurement filtering in wireless sensor networks . . . . .         | 7         |
| 1.4      | Parameter Estimation in Wireless Sensor Networks . . . . .          | 8         |
| 1.4.1    | Opinion Dynamics in Social Influence Networks . . . . .             | 8         |
| 1.5      | Main questions in averaging algorithms . . . . .                    | 8         |
| 1.6      | Distributed control in cooperative robotics . . . . .               | 8         |
| 1.6.1    | Main questions in cooperative robotics . . . . .                    | 9         |
| 1.6.2    | Distributed optimal control . . . . .                               | 9         |
| <b>2</b> | <b>Preliminaries on Algebraic Graph Theory</b>                      | <b>11</b> |
| <b>3</b> | <b>Averaging Systems</b>  | <b>13</b> |
| 3.1      | Distributed algorithm . . . . .                                     | 13        |
| 3.2      | Discrete-time averaging systems . . . . .                           | 13        |
| 3.3      | Stochastic matrices . . . . .                                       | 14        |
| 3.4      | Example: Metropolis-Hastings weights . . . . .                      | 14        |
| 3.5      | Time-varying digraphs . . . . .                                     | 15        |
| 3.5.1    | Averaging distributed algorithms over time-varying graphs . . . . . | 15        |
| 3.5.2    | Discrete-time consensus over time-varying graphs . . . . .          | 15        |
| 3.6      | Laplacian dynamics . . . . .  | 16        |
| 3.6.1    | Properties of the Laplacian matrix . . . . .                        | 17        |
| 3.6.2    | Consensus for Laplacian dynamics . . . . .                          | 17        |
| <b>4</b> | <b>Optimization basics</b>  | <b>19</b> |
| 4.1      | Optimization algorithms . . . . .                                   | 20        |
| 4.1.1    | A system theoretical perspective to the gradient method . . . . .   | 20        |
| 4.2      | steady-state analysis of the gradient method . . . . .              | 20        |
| 4.2.1    | Gradient method for quadratic programs . . . . .                    | 21        |
| 4.2.2    | Gradient flow . . . . .   | 21        |
| 4.2.3    | Nesterov accelerated gradient method . . . . .                      | 21        |
| <b>5</b> | <b>Parallel Optimization and Federated Learning</b>                 | <b>23</b> |
| 5.1      | Incremental gradient method . . . . .                               | 23        |
| 5.2      | Stochastic Gradient Descent . . . . .                               | 24        |
| 5.2.1    | Convergence . . . . .   | 24        |
| 5.2.2    | Stochastic mini-batch gradient method . . . . .                     | 24        |
| 5.3      | Beyond SGD: Adaptive Momentum Estimation (Adam) . . . . .           | 25        |
| 5.4      | Federated learning . . . . .  | 25        |
| 5.5      | Distributed learning . . . . .                                      | 25        |
| <b>6</b> | <b>Leader-Follower networks: Containment</b>                        | <b>27</b> |
| 6.1      | Definition . . . . .  | 27        |
| 6.2      | Control law . . . . .   | 27        |

|           |   |           |
|-----------|---|-----------|
| <b>7</b>  | <b>Leader Follower networks: Formation control</b>                            | <b>29</b> |
| 7.1       | analogy with mass-spring systems . . . . .                                    | 29        |
| 7.2       | Formation Control . . . . .   | 30        |
| 7.2.1     | Laplacian dynamics: potential function interpretation . . . . .               | 30        |
| 7.2.2     | Distributed control based on potential functions . . . . .                    | 30        |
| 7.2.3     | Inter-agent collision avoidance . . . . .                                     | 31        |
| 7.2.4     | Obstacle avoidance . . . . .  | 31        |
| <b>8</b>  | <b>Distributed Cost-Coupled Optimization</b>                                  | <b>33</b> |
| 8.1       | Cost-Coupled Optimization Set-up . . . . .                                    | 33        |
| 8.2       | Distributed Gradient Algorithm . . . . .                                      | 33        |
| 8.3       | Assumptions . . . . .   | 34        |
| 8.4       | Distributed Projected Subgradient Algorithm . . . . .                         | 34        |
| 8.5       | Improving the Distributed Gradient Algorithm . . . . .                        | 34        |
| 8.5.1     | Dynamic Average Consensus . . . . .   | 35        |
| 8.6       | The Gradient Tracking Algorithm . . . . .                                     | 35        |
| 8.6.1     | Assumptions . . . . .   | 36        |
| 8.6.2     | Convergence Result . . . . .  | 36        |
| <b>9</b>  | <b>Distributed Aggregative Optimization</b>                                   | <b>37</b> |
| 9.1       | Aggregative optimization . . . . .  | 37        |
| 9.2       | Distributed aggregative optimization . . . . .                                | 38        |
| 9.2.1     | Extension to online aggregative optimization . . . . .                        | 39        |
| <b>10</b> | <b>Learning with Neural Networks</b>  | <b>41</b> |
| 10.1      | Basic element: the neuron model . . . . .                                     | 41        |
| <b>11</b> | <b>Multi-Robot Safety Controllers</b>   | <b>43</b> |
| 11.1      | Control Barrier Functions . . . . .   | 43        |
| 11.1.1    | Admissible Control Space . . . . .  | 43        |
| 11.1.2    | Safety Filters via Control Barrier Certificates . . . . .                     | 43        |
| 11.2      | Single-Robot Obstacle Avoidance: Single Integrator Model . . . . .            | 44        |
| 11.3      | Multi-Robot Collision Avoidance: Single Integrators . . . . .                 | 44        |
| 11.3.1    | Control Barrier Function for Multi-Robot Collision Avoidance . . . . .        | 44        |
| 11.3.2    | Centralized Safety Controller for Multi-Robot Collision Avoidance . . . . .   | 45        |
| 11.3.3    | Decentralized Safety Controller for Multi-Robot Collision Avoidance . . . . . | 45        |
| <b>12</b> | <b>Stability Analysis for Linear Averaging</b>                                | <b>47</b> |
| 12.1      | Consensus result . . . . .  | 47        |
| 12.2      | Irreducible and primitive matrices . . . . .                                  | 48        |
| 12.3      | Spectral properties of square matrices . . . . .                              | 48        |
| 12.3.1    | Spectral properties of row-stochastic matrices . . . . .                      | 49        |
| 12.4      | Application of the Perron-Frobenius theorem . . . . .                         | 49        |
| 12.4.1    | Equilibrium manifold of linear averaging . . . . .                            | 49        |
| 12.5      | Convergence proof (I): Jordan form approach . . . . .                         | 49        |
| 12.5.1    | The dissensus dynamics . . . . .  | 50        |
| 12.6      | Convergence proof (II): Lyapunov approach . . . . .                           | 51        |
| 12.7      | Convergence proof (III): alternative Lyapunov approach . . . . .              | 51        |
| <b>13</b> | <b>Gradient tracking analysis</b>   | <b>53</b> |
| 13.1      | Equilibrium manifold for the Gradient Tracking . . . . .                      | 53        |
| 13.2      | Causal coordinates for the Gradient Tracking . . . . .                        | 54        |
| 13.2.1    | Block diagram of the Gradient Tracking Algorithm . . . . .                    | 54        |
| 13.3      | Local perspective of the causal GT . . . . .                                  | 55        |
| 13.4      | Error coordinates for the Gradient Tracking . . . . .                         | 55        |
| 13.4.1    | Block diagram of the Gradient Tracking in error coordinates . . . . .         | 55        |
| 13.5      | Analysis sketch . . . . .   | 56        |

|  |           |
|--|-----------|
| 13.6 Gradient Tracking: extensions . . . . .                           | 57        |
| 13.7 System theoretical approach to distributed optimization . . . . . | 57        |
| <b>14 Constraint-coupled Distributed Optimization</b>                  | <b>59</b> |
| 14.1 Constraint-coupled Setup . . . . .                                | 59        |
| 14.2 Multi-robot scenario . . . . .                                    | 59        |
| 14.3 Distributed optimal control . . . . .                             | 60        |
| 14.4 Duality for constraint-coupled optimization . . . . .             | 61        |
| 14.5 Centralized dual subgradient . . . . .                            | 61        |
| 14.5.1 Distributed dual subgradient convergence . . . . .              | 62        |



# Chapter 1

## Introduction and scenarios

### 1.1 Distributed Autonomous System

Each agent  $i \in \{1, \dots, N\}$  has

- local physical and/or cyber state  $x_i$
- computational and sensing capabilities
- communication capability: exchange messages with "neighbours"

### 1.2 Scenarios and applications of distributed systems

- Averaging: distributed estimation, opinion dynamics
- Distributed control in cooperative robotics
- Distributed optimization
  - distributed machine learning
  - distributed decision-making in cooperative robotics
  - distributed optimal control in energy systems and cooperative robotics

### 1.3 Measurement filtering in wireless sensor networks

Consider a network of  $N$  sensors with local sensing, computation and communication. Agent  $i, i \in \{1, \dots, N\}$ , takes a local measurement from the environment (temperature, pressure, etc.). Let  $x_{i0} \in \mathbb{R}$  be the scalar local measurement. Agents are interested in agreeing on the average of the measurements,

$$x_{\text{avg}} = \frac{1}{N} \sum_{i=1}^N x_{i0}$$

to have a better estimate of the environment quantity

Consider the following "distributed algorithm" based on "local" linear averaging, for each  $i \in \{1, \dots, N\}$

$$\begin{aligned} x_i^0 &= x_{i0} \\ x_i^{k+1} &= \text{average}(x_i^k, \{x_j^k, j \text{ "neighbour" of } i\}), \quad k \in \mathbb{N} \end{aligned}$$

generalizing coefficients of the update:

$$\begin{aligned} x_i^0 &= x_{i0} \\ x_i^{k+1} &= \sum_{j=1}^N a_{ij} x_j^k \quad k \in \mathbb{N} \end{aligned}$$

*Remark.*  $a_{ij} \geq 0$  and  $\sum_{j=1}^N a_{ij} = 1$

*Remark.*  $a_{ij} = 0$ , for some  $j \in \{1, \dots, N\}$ , i.e.  $a_{ij} = 0$  if  $i$  does not have access to the estimate of  $j$

## 1.4 Parameter Estimation in Wireless Sensor Networks

Consider a network of  $N$  sensors with local sensing, computation and communication aiming at estimating a common parameter  $\theta^* \in \mathbb{R}$ . Each sensor  $i$  measures

$$y_i = B_i \theta^* + v_i$$

with  $y_i \in \mathbb{R}^{m_i}$ ,  $B_i$  known matrix and  $v_i$  a random measurement noise. Assume  $v_1, \dots, v_N$  independent and Gaussian, with zero mean and covariance  $E[v_i v_i^T] = \Sigma_i$ . Assume  $\sum_{i=1}^N m_i \geq m$  and  $\begin{bmatrix} B_1 \\ \vdots \\ B_N \end{bmatrix}$  full rank. Compute a least-squares estimate

$$\hat{\theta} = \arg \min_{\theta} \sum_{i=1}^N (y_i - B_i \theta)^T \Sigma_i^{-1} (y_i - B_i \theta)$$

The optimal solution is

$$\begin{aligned} \hat{\theta} &= \left( \sum_{i=1}^N B_i^T \Sigma_i^{-1} B_i \right)^{-1} \sum_{i=1}^N B_i^T \Sigma_i^{-1} y_i \\ &= \left( \frac{1}{N} \sum_{i=1}^N B_i^T \Sigma_i^{-1} B_i \right)^{-1} \frac{1}{N} \sum_{i=1}^N B_i^T \Sigma_i^{-1} y_i \end{aligned}$$

The optimal solution can be obtained by computing two averages  $\frac{1}{N} \sum_{i=1}^N \beta_i$  and  $\frac{1}{N} \sum_{i=1}^N \beta_i$

### 1.4.1 Opinion Dynamics in Social Influence Networks

Group of  $N$  individuals, with  $x_i^k$  being the opinion of individual  $i$  at time  $k$ . Opinions are updated according to

$$x_i^{k+1} = \sum_{j=1}^N a_{ij} x_j^k$$

## 1.5 Main questions in averaging algorithms

- Do node estimates converge? Do they converge to a common value ("reach consensus")?
- Do they reach consensus to the average ("average consensus")?
- How can we model communication in general networks?
- Can we answer the above questions for general networks and communication protocols?
- What assumptions do we need on the communication network?

## 1.6 Distributed control in cooperative robotics

Team of  $N$  (mobile) robots aiming to execute complex tasks



**Basic tasks**

- rendezvous, containment
- formation, flocking
- coverage

**Complex tasks**

- pickup and delivery
- surveillance and patrolling
- exploration
- satellite constellation

**1.6.1 Main questions in cooperative robotics**

- Do robot states asymptotically converge?
- Do the asymptotic states satisfy the global, desired task?
- How can we model communication in (general) robotic networks?
- What assumptions do we need on the communication network?
- Can we answer the above questions for general networks and communication protocols?

**1.6.2 Distributed optimal control**

$$\begin{aligned}
& \min_{\substack{x_1, \dots, x_N \\ u_1, \dots, u_N}} \sum_{i=1}^N \left( \sum_{\tau=0}^{T-1} \ell_i(z_{i,\tau}, u_{i,\tau}) + m_i(z_{i,T}) \right) \\
& \text{subj to } \sum_{i=1}^N H_i z_{i,\tau} \leq h, & \tau \in [0, T] \\
& z_{i,\tau+1} = A_i z_{i,\tau} + B_i u_{i,\tau} & \forall i, \tau \in [0, T] \\
& z_{i,\tau} \in Z_i, \quad u_{i,\tau} \in U_i, & \forall i, \tau \in [0, T]
\end{aligned}$$



## Chapter 2

# Preliminaries on Algebraic Graph Theory

**Definition 2.1** (Digraph)

A digraph is a pair  $G = (I, E)$  where  $I = 1, \dots, N$  is a set of elements called *nodes* and  $E \subset I \times I$  is a set of ordered node pairs called *edges*

*Edge*: the pair  $(i, j)$  denotes an edge from  $i$  to  $j$

*Self-loop*: edge from a node to itself, i.e.  $(i, i)$

**Definition 2.2** (Undirected (di)graph)

if for any  $(i, j) \in E$  then  $(j, i) \in E$

**Definition 2.3** (Subgraph)

$(I', E')$  subgraph of  $(I, E)$  if  $I' \subset I$  and  $E' \subset E$ . Spanning subgraph if  $I' = I$

**Definition 2.4** (In-neighbours of  $i$ )

$j \in I$  is an in-neighbour of  $i \in I$  if  $(j, i) \in E$

**Definition 2.5** (Set of in-neighbours of  $i$ )

$\mathcal{N}_i^{\text{IN}} = \{j \in \{1, \dots, N\} | (j, i) \in E\}$

**Definition 2.6** (Out-neighbours of  $i$ )

$j \in I$  is an out-neighbour of  $i \in I$  if  $(i, j) \in E$

**Definition 2.7** (Set of out-neighbours of  $i$ )

$\mathcal{N}_i^{\text{OUT}} = \{j \in \{1, \dots, N\} | (i, j) \in E\}$

**Definition 2.8** (In-degree  $\deg_i^{\text{IN}}$ )

number of in-neighbours, i.e. cardinality of  $\mathcal{N}_i^{\text{IN}}$  ( $\deg_i^{\text{IN}} = |\mathcal{N}_i^{\text{IN}}|$ )

Out-degree analogous

**Definition 2.9** (Balanced digraph)

A digraph  $G$  is balanced if  $\deg_i^{\text{IN}} = \deg_i^{\text{OUT}}$  for all  $i \in \{1, \dots, N\}$

**Definition 2.10** (Directed path)

ordered sequence of nodes s.t. any pair of consecutive nodes is a directed edge. A path is *simple* if no node appears more than once

**Definition 2.11** (Directed cycle)

simple directed path that starts and ends at the same node. Cycles of length one are called self-loops. Acyclic digraph if there are no cycles.

**Definition 2.12** (Directed tree)

Acyclic digraph s.t. there exists a node, root, s.t. any node can be reached by only one directed path starting at the root.

**Definition 2.13** (Spanning tree)

a spanning subgraph that is a (directed) tree

**Definition 2.14** (Periodic graph)

if there exists a  $k > 1$  period, that divides the length of every cycle.

Note: a graph with self-loops is aperiodic.

**Definition 2.15** (Strongly connected digraph)

if there exists a directed path from any node to any other node

**Definition 2.16** (Connected undirected graph)

if there exists a path from any node to any other node

**Definition 2.17** (Globally reachable node)

if one of its nodes can be reached from any other node by traversing a directed path

**Definition 2.18** (Complete graph)

Unweighted graph such that  $\forall i, j \exists (i, j), (j, i) \in E$

**Definition 2.19** (Weighted digraph)

triplet  $G = (I, E, \{a_{(i,j)}\}_{(i,j) \in E})$ , where  $(I, E)$  is a digraph and each  $a_{(i,j)}$  is a strictly positive weight for the edge  $(i, j)$

**Definition 2.20** (Weighted in-degree)

$$\deg_i^{\text{IN}} = \sum_{j=1}^N a_{ji}$$

**Definition 2.21** (Weighted out-degree)

$$\deg_i^{\text{OUT}} = \sum_{j=1}^N a_{ij}$$

A graph is said to be weight balanced if  $\deg_i^{\text{IN}} = \deg_i^{\text{OUT}}$  for all nodes  $i \in \{1, \dots, N\}$

**Definition 2.22** (Weighted adjacency matrix)

non-negative matrix  $A$  s.t. the entry  $(i, j)$  of  $A$ , denoted  $a_{ij}$ , satisfies:

- $a_{ij} > 0$  if  $(i, j) \in E$
- $a_{ij} = 0$  otherwise

**Definition 2.23** (Adjacency matrix)

matrix  $A$  with entries  $a_{ij} = 1$  if  $(i, j) \in E$  and  $a_{ij} = 0$  otherwise

**Definition 2.24** (Weighted in/out degree)

$$\deg_i^{\text{IN}} = \sum_{j=1}^N a_{ji} \text{ and } \deg_i^{\text{OUT}} = \sum_{j=1}^N a_{ij}$$

**Definition 2.25** (Laplacian matrix)

$$L = D^{\text{OUT}} - A$$

## Chapter 3

# Averaging Systems

### 3.1 Distributed algorithm

Given a network of  $N$  agents communicating according to a fixed digraph  $G$ , i.e. each agent  $i$  can receive messages only from in-neighbours in the graph, i.e. from  $j \in \mathcal{N}_i^{\text{IN}}$ . We start by considering a fixed graph, thus, each agent communicates with the same neighbours at each iteration  $k \in \mathbb{N}$

$$x_i^{k+1} = \text{stf}_i(x_i^k, \{x_j^k\}_{j \in \mathcal{N}_i^{\text{IN}}}), \quad i \in \{1, \dots, N\}$$

where  $\text{stf}_i$  is a function depending only on state  $x_i$  and states  $x_j, j \in \mathcal{N}_i^{\text{IN}}$ .

Alternative version with out-neighbours:

$$x_i^{k+1} = \text{stf}_i(\{x_j\}_{j \in \mathcal{N}_i^{\text{OUT}}})$$

### 3.2 Discrete-time averaging systems

Let  $G^{\text{comm}} = (I, E)$  be a fixed (communication) digraph (self loops included). A linear averaging distributed algorithm can be written as:

$$x_i^{k+1} = \sum_{j \in \mathcal{N}_i^{\text{IN}}} a_{ij} x_j^k \quad i \in \{1, \dots, N\}$$

where  $x_i^k \in \mathbb{R}$  is the state of agent  $i$  at  $k$  and  $a_{ij} > 0$  are positive weights.

*Remark.* The weights  $a_{ij}$  are defined only for  $(i, j) \in E$

Each  $i$  uses only the states of neighbours  $j \in \mathcal{N}_i^{\text{IN}}$ , thus this is a distributed algorithm.

For analysis purposes, let us define weights  $a_{ij} = 0$  for  $(j, i) \notin E$ . Thus we can rewrite the distributed algorithm as

$$x_i^{k+1} = \sum_{j=1}^N a_{ij} x_j^k \quad i \in \{1, \dots, N\}$$

This is an LTI autonomous system

$$\begin{bmatrix} x_1^{k+1} \\ \vdots \\ x_N^{k+1} \end{bmatrix} = \begin{bmatrix} a_{11} & \cdots & a_{1N} \\ \vdots & \ddots & \vdots \\ a_{N1} & \cdots & a_{NN} \end{bmatrix} \begin{bmatrix} x_1^k \\ \vdots \\ x_N^k \end{bmatrix}$$

Which can be compactly written as

$$x^{k+1} = Ax^k$$

*Remark.* The matrix  $A$  can be seen as the weighted adjacency matrix of the reverse digraph  $G^{\text{comm,rev}}$  of the digraph  $G^{\text{comm}}$

If instead of in-neighbours we use out-neighbours, we call the digraph a sensing digraph  $G^{\text{sens}}$ . In this case the notation becomes consistent with graph theory, so we get

$$x^{k+1} = Ax^k$$

where  $A$  can be seen as the weighted adjacency matrix of the sensing digraph  $G^{\text{sens}}$

### 3.3 Stochastic matrices

The non-negative square matrix  $A \in \mathbb{R}^{N \times N}$  is

- row stochastic if  $A\mathbf{1} = \mathbf{1}$  (each row sums to 1)
- column stochastic if  $A^\top \mathbf{1} = \mathbf{1}$  (each column sums to 1)
- doubly stochastic if both row and column stochastic.

**Lemma.** Let  $A$  be a row-stochastic matrix and  $G$  the associated digraph. If  $G$  is strongly connected and aperiodic, then

1. the eigenvalue  $\lambda = 1$  is simple;
2. all the other eigenvalues  $\mu$  satisfy  $|\mu| < 1$

*Remark.* The condition “ $G$  contains a globally reachable node and the subgraph of globally reachable nodes is aperiodic” is necessary and sufficient

**Theorem 3.1** (Consensus)

Consider a (discrete-time) averaging system with associated digraph  $G$  and weighted adjacency matrix  $A$ . Assume  $G$  is strongly connected and aperiodic, and  $A$  is row stochastic. Then

1. there exists a left eigenvector  $w \in \mathbb{R}^N, w > 0$  (i.e. with positive components  $w_i > 0$  for all  $i = 1, \dots, N$ ) such that

$$\lim_{k \rightarrow \infty} x^k = \mathbf{1} \frac{w^\top x^0}{w^\top \mathbf{1}} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \frac{\sum_{i=1}^N w_i x_i^0}{\sum_{i=1}^N w_i}$$

i.e., consensus is reached to  $\frac{\sum_{i=1}^N w_i x_i^0}{\sum_{i=1}^N w_i}$

2. if additionally  $A$  is doubly stochastic, then

$$\lim_{k \rightarrow \infty} x^k = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \frac{\sum_{i=1}^N x_i^0}{N}$$

i.e., average consensus is reached

### 3.4 Example: Metropolis-Hastings weights

Given an undirected unweighted graph  $G$  with edge set  $E$  and degrees  $d_1, \dots, d_n$

$$a_{ij} = \begin{cases} \frac{1}{1 + \max\{d_i, d_j\}} & \text{if } (i, j) \in E \text{ and } i \neq j \\ 1 - \sum_{h \in \mathcal{N}_i \setminus \{i\}} a_{ih} & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

Result: the matrix  $A$  is symmetric and doubly-stochastic.

### 3.5 Time-varying digraphs

A time-varying digraph is a sequence of digraphs  $\{G(k)\}_{k \geq 0}$ .

*Remark.* The main definitions of in/out neighbours, in/out degree, adjacency matrix can be generalized by considering time-varying versions, i.e.  $\mathcal{N}_i^{\text{IN}}(k)$ ,  $\mathcal{N}_i^{\text{OUT}}(k)$ ,  $\deg_i^{\text{IN}}(k)$ ,  $\deg_i^{\text{OUT}}(k)$ ,  $A(k)$  associated to each graph  $G(k)$ . Connectivity requires new definitions as assuming each  $G(k)$  to be connected is too conservative.

**Definition 3.1** (Jointly strongly connected digraph)

if  $\bigcup_{\tau=k}^{+\infty} G(\tau)$  is strongly connected  $\forall k \geq 0$

**Definition 3.2** (Uniformly jointly strongly connected (or  $B$ -strongly connected) digraph)

if there exists  $B \in \mathbb{N}$  such that  $\bigcup_{\tau=k}^{k+B} G(\tau)$  is strongly connected  $\forall k \geq 0$

*Remark.* The graph can be disconnected at some time  $k$ .

#### 3.5.1 Averaging distributed algorithms over time-varying graphs

Let  $\{G(k)\}_{k \geq 0}$  be a time-varying digraph (with self loops for each  $G(k)$ ). Consider the distributed algorithm

$$x_i^{k+1} = \sum_{j \in \mathcal{N}_i^{\text{IN}}(k)} a_{ij}(k) x_j^k \quad \forall i \in \{1, \dots, N\}$$

or the out-neighbours version

$$x_i^{k+1} = \sum_{j \in \mathcal{N}_i^{\text{OUT}}(k)} a_{ij}(k) x_j^k \quad \forall i \in \{1, \dots, N\}$$

where  $x_i^k \in \mathbb{R}$  is the state of agent  $i$  at  $k$  and  $a_{ij}(k) > 0$ .

For analysis purposes, let us define weights  $a_{ij}(k) = 0$  for  $(i, j) \notin E(k)$ . Thus we can rewrite the distributed algorithm as

$$x_i^{k+1} = \sum_{j=1}^N a_{ij}(k) x_j^k \quad i \in \{1, \dots, N\}$$

This is a Linear Time-Varying system

$$x^{k+1} = A(k)x^k$$

with state  $x := [x_1, \dots, x_N]^\top$  and state matrix

$$A(k) := \begin{bmatrix} a_{11}(k) & \cdots & a_{1N}(k) \\ \vdots & \ddots & \vdots \\ a_{N1}(k) & \cdots & a_{NN}(k) \end{bmatrix}$$

being a weighted adjacency matrix associated to the digraph  $G(k)$ .

#### 3.5.2 Discrete-time consensus over time-varying graphs

**Theorem 3.2**

Let  $\{A(k)\}_{k \geq 0}$  be a sequence of row-stochastic matrices with associated digraphs  $\{G(k)\}_{k \geq 0}$ . Assume

1. each digraph  $G(k)$  has a self-loop at each node;
2. each non-zero edge weight  $a_{ij}(k)$ , including the self-loop weights  $a_{ii}(k)$ , is larger than a constant  $\epsilon > 0$ ;
3. there exists  $B \in \mathbb{N}$  such that, for all times  $k \geq 0$ , the union digraph  $G(k) \cup \dots \cup G(k+B)$  is strongly connected.

Then

1. there exists a non-negative vector  $w \in \mathbb{R}^N$  such that the solution to  $x^{k+1} = A(k)x^k$  converges (exponentially) to  $\mathbf{1} \frac{w^\top x^0}{w^\top \mathbf{1}}$ , i.e.

$$\lim_{k \rightarrow \infty} x^k = \mathbf{1} \left( \frac{w^\top x^0}{w^\top \mathbf{1}} \right)$$

2. if additionally each matrix in the sequence is doubly-stochastic, then

$$\lim_{k \rightarrow \infty} x^k = \mathbf{1} \frac{1}{N} \sum_{i=1}^N x_i^0$$

i.e., average consensus is achieved

### 3.6 Laplacian dynamics

Consider a network of dynamical systems with dynamics

$$\dot{x}(t) = u_i(t) \quad i \in \{1, \dots, N\}$$

with states  $x_i \in \mathbb{R}$  and inputs  $u_i \in \mathbb{R}$ , communicating (or interacting) according to a digraph  $G$ . Consider a (distributed) “proportional” feedback control

$$u_i(t) = - \sum_{j \in \mathcal{N}_i^{\text{IN}}} a_{ij}(x_i(t) - x_j(t))$$

or the out-neighbour version

$$u_i(t) = - \sum_{j \in \mathcal{N}_i^{\text{OUT}}} a_{ij}(x_i(t) - x_j(t))$$

For analysis purposes, let us define weights  $a_{ij}(k) = 0$  for  $(i, j) \notin E(k)$ . Thus we can rewrite the distributed control systems as

$$\dot{x}_i(t) = - \sum_{j=1}^N a_{ij}(x_i(t) - x_j(t)) \quad \forall i \in \{1, \dots, N\}$$

Defining  $x := [x_1 \cdots x_N]^\top$ , it can be shown that it can be rewritten as the following Linear Time Invariant continuous-time system

$$\dot{x}(t) = -Lx(t)$$

where  $L$  is the (weighted) Laplacian associated to the digraph  $G$  with (weighted) adjacency matrix  $A$

Let

$$\dot{x}_i(t) = - \sum_{j=1}^N a_{ij}(x_i(t) - x_j(t)) \quad \forall i \in \{1, \dots, N\}$$

rearranging terms

$$\dot{x}_i(t) = - \left( \sum_{j=1}^N a_{ij} \right) x_i(t) + \sum_{j=1}^N a_{ij} x_j(t) = -\deg_i^{\text{OUT}} x_i(t) + (Ax(t))_i$$

where  $(Ax(t))_i$  is the  $i$ -th element of  $Ax(t)$ . Writing the previous dynamics in a compact form

$$\dot{x}(t) = -(D^{\text{OUT}} - A)x(t)$$

where we recall that  $D^{\text{OUT}}$  is the (weighted) out-degree matrix. Recalling that  $L = D^{\text{OUT}} - A$ , it holds that

$$\dot{x}(t) = -Lx(t)$$

*Remark.* if the in-neighbours version is considered, then  $\dot{x}(t) = -L^{\text{IN}}x(t)$ , where  $L^{\text{IN}} = D^{\text{IN}} - A^T$  is the in-degree Laplacian (i.e. the Laplacian of the reverse graph of  $G$ )



### 3.6.1 Properties of the Laplacian matrix

It can be easily verified that

$$L\mathbf{1} = D^{\text{OUT}}\mathbf{1} - A\mathbf{1} = \begin{bmatrix} \deg_1^{\text{OUT}} \\ \vdots \\ \deg_i^{\text{OUT}} \end{bmatrix} - \begin{bmatrix} \deg_1^{\text{OUT}} \\ \vdots \\ \deg_i^{\text{OUT}} \end{bmatrix} = 0$$

i.e.,  $\lambda = 0$  is an eigenvalue of  $L$  and  $\mathbf{1}$  is an associated eigenvector.

**Lemma.** Given a weighted digraph with Laplacian  $L$ , then all eigenvalues of  $L$  different from zero have strictly positive real part

**Lemma.** Given a weighted digraph with Laplacian  $L$ , the following statements are equivalent:

1.  $G$  is weight-balanced, i.e.  $D^{\text{IN}} = D^{\text{OUT}}$
2.  $\mathbf{1}L = 0$

#### Theorem 3.3

A weighted digraph with Laplacian  $L$  contains a globally reachable node if and only if  $\lambda = 0$  is simple.

**Corollary.** If a weighted digraph is strongly connected, then  $\lambda = 0$  is simple

### 3.6.2 Consensus for Laplacian dynamics

#### Theorem 3.4

let  $L$  be a (weighted) Laplacian matrix with associated strongly connected (weighted) digraph  $G$ . Consider the Laplacian dynamics  $\dot{x}(t) = -Lx(t)$ ,  $t \geq 0$ , then

- 1.

$$\lim_{t \rightarrow \infty} x(t) = \mathbf{1} \left( \frac{w^\top x(0)}{w^\top \mathbf{1}} \right)$$

with  $w^\top L = 0$ , i.e.  $w$  is a left eigenvector for the eigenvalue  $\lambda = 0$ ;

2. if additionally  $G$  is weight-balanced then

$$\lim_{t \rightarrow \infty} x(t) = \mathbf{1} \frac{\sum_{i=1}^N x_i(0)}{N}$$



## Chapter 4

# Optimization basics

### Convexity and gradient monotonicity

If a convex function  $\ell$  is also differentiable, then its gradient  $\nabla\ell : \mathbb{R}^d \rightarrow \mathbb{R}^d$  satisfies

$$(\nabla\ell(z_A) - \nabla\ell(z_B))^\top (z_A - z_B) \geq 0$$

for all  $z_A, z_B$ . That is, the gradient  $\nabla\ell$  is a monotone operator

### Strict convexity and gradient monotonicity

A function  $\ell$  is strictly convex if for  $z_A \neq z_B$  and  $\theta \in (0, 1)$

$$\ell(\theta z_A + (1 - \theta)z_B) < \theta\ell(z_A) + (1 - \theta)\ell(z_B)$$

If the strictly convex function  $\ell$  is also differentiable, then its gradient satisfies

$$(\nabla\ell(z_A) - \nabla\ell(z_B))^\top (z_A - z_B) > 0$$

for all  $z_A, z_B$ . That is, the gradient  $\nabla\ell$  is a strictly monotone operator

### Strong convexity and gradient monotonicity

A function  $\ell$  is strongly convex with parameter  $\mu > 0$  if for  $z_A \neq z_B$  and  $\theta \in (0, 1)$

$$\ell(\theta z_A + (1 - \theta)z_B) < \theta\ell(z_A) + (1 - \theta)\ell(z_B) - \mu\theta(1 - \theta)\|z_A - z_B\|^2$$

The gradient of a differentiable strongly convex function satisfies

$$(\nabla\ell(z_A) - \nabla\ell(z_B))^\top (z_A - z_B) \geq \mu\|z_A - z_B\|^2$$

for all  $z_A, z_B$ . That is, the gradient  $\nabla\ell$  is a strongly monotone operator

### Convexity and Lipschitz continuity of the gradient

Consider a differentiable convex function  $\ell$  with a Lipschitz continuous gradient with parameter  $L > 0$ , i.e.

$$\|\nabla\ell(z_A) - \nabla\ell(z_B)\| \leq L\|z_A - z_B\|$$

for all  $z_A, z_B$ . Then, the following characterization holds

$$(\nabla\ell(z_A) - \nabla\ell(z_B))^\top (z_A - z_B) \geq \frac{1}{L}\|\nabla\ell(z_A) - \nabla\ell(z_B)\|^2$$

for all  $z_A, z_B$ . That is, the gradient  $\nabla\ell$  is a co-coercive operator

### Strong convexity and Lipschitz continuity of the gradient

Consider a strongly convex (with parameter  $\mu > 0$ ) function  $\ell$  with Lipschitz continuous gradient (with parameter  $L > 0$ ). The the followin characterization holds

$$(\nabla\ell(z_A) - \nabla\ell(z_B))^\top (z_A - z_B) \geq \frac{\mu L}{\mu + L} \|z_A - z_B\|^2 + \frac{1}{\mu + L} \|\nabla\ell(z_A) - \nabla\ell(z_B)\|^2$$

for all  $z_A, z_B$ .

## 4.1 Optimization algorithms

We consider optimization algorithms based on iterative descent.

*Notation.* We denote by  $z^k \in \mathbb{R}^d$  the estimate at iteration  $k \in \mathbb{N}$  of a local minimum.

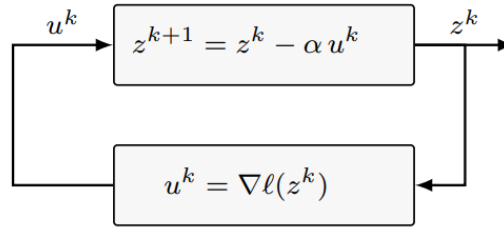
The algorithm starts at a given initial guess ecc ecc we know iterative descent

### 4.1.1 A system theoretical perspective to the gradient method

Let  $\ell$  be  $\mu$ -strongly convex and with  $L$ -Lipschitz continuous gradient. The gradient method is a discrete-time integrator in feedback interconnection with a static map

$$\begin{aligned} z^{k+1} &= z^k - \alpha u^k, & z^0 \text{ given} \\ u^k &= \nabla\ell(z^k) \end{aligned}$$

where  $u^k$  is a "control input" obtained through a static (nonlinear) state feedback. The block diagram is



This is known as the *Lur'e problem*

## 4.2 steady-state analysis of the gradient method

Assume that the state  $z^k$  converges to some value  $z_{eq}$ . Then, such an equilibrium must satisfy:

$$\begin{aligned} z_{eq} &= z_{eq} - \alpha \nabla\ell(z_{eq}) \implies & z_{eq} &: \nabla\ell(z_{eq}) \\ &\implies & z_{eq} &= z^* \end{aligned}$$

Consider the change of coordinates  $z^k \rightarrow \tilde{z}^k - z_{eq} = z^k - z^*$ . Then, the error dynamics is

$$\begin{aligned} \tilde{z}^{k+1} &= \tilde{z}^k - \alpha u^k \\ u^k &= \nabla\ell(\tilde{z}^k + z^*) - \nabla\ell(z^*) \end{aligned}$$

where  $u^k$  and  $\tilde{z}^k$  satisfy, in light of the assumption on  $\ell$ , the following inequality<sup>1</sup>

$$-(u^k)^T \tilde{z}^k \leq -\gamma_1 \|\tilde{z}^k\|^2 - \gamma_2 \|u^k\|^2$$

---

<sup>1</sup>For all  $z_A, z_B$  it holds that

$$(\nabla\ell(z_A) - \nabla\ell(z_B))^\top (z_A - z_B) \geq \frac{\mu L}{\mu + L} \|z_A - z_B\|^2 + \frac{1}{\mu + L} \|\nabla\ell(z_A) - \nabla\ell(z_B)\|^2$$

Consider a Lyapunov function  $V : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$  given by  $V(\tilde{z}) = \|\tilde{z}\|^2$ . Then

$$\begin{aligned} V(\tilde{z}^{k+1}) - V(\tilde{z}^k) &= \|\tilde{z}^{k+1}\|^2 - \|\tilde{z}^k\|^2 \\ &= \|\tilde{z}^k\|^2 - 2\alpha(u^k)^\top \tilde{z}^k + \alpha^2 \|u^k\|^2 - \|\tilde{z}^k\|^2 \\ &\leq -2\alpha\gamma_1 \|\tilde{z}^k\|^2 + \alpha(\alpha - 2\gamma_2) \|u^k\|^2 \end{aligned}$$

For a small enough stepsize  $\alpha$  (i.e.,  $\alpha \leq 2\gamma_2$ ), we can write

$$V(\tilde{z}^{k+1}) - V(\tilde{z}^k) < -2\alpha\gamma_1 \|\tilde{z}^k\|^2 \implies \|\tilde{z}^{k+1}\|^2 \leq (1 - 2\alpha\gamma_1) \|\tilde{z}^k\|^2 \leq (1 - 2\alpha\gamma_1)^k \|\tilde{z}^0\|^2$$

Therefore  $\{\tilde{z}^k\}_{k \in \mathbb{N}}$  goes exponentially/geometrically fast to zero

#### 4.2.1 Gradient method for quadratic programs

Consider a quadratic program

$$\min_z \frac{1}{2} z^\top Q z + r^\top z$$

With  $Q = Q^\top > 0$  The gradient method is an affine linear system

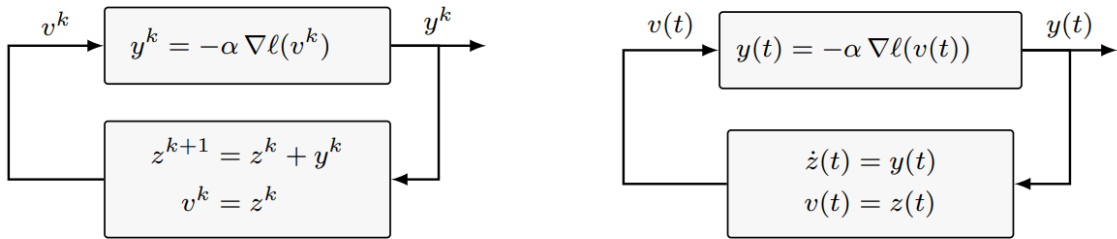
$$\begin{aligned} z^{k+1} &= z^k - \alpha(Qz^k + r) \quad z^k \text{ given} \\ &= (I - \alpha Q)z^k - \alpha r \end{aligned}$$

For a sufficiently small  $\alpha$ , the state matrix  $(I - \alpha Q)$  is Schur. Hence, the state trajectory is <sup>2</sup>

$$z^k = (I - \alpha Q)^k z^0 - \alpha \sum_{i=0}^{k-1} (I - \alpha Q)^i r \xrightarrow{k \rightarrow \infty} -\alpha \left( \sum_{i=0}^{\infty} (I - \alpha Q)^i \right) r = -Q^{-1}r$$

#### 4.2.2 Gradient flow

Let us swap the roles of the plant (the static nonlinearity) and the controller (the integrator) We obtain the



so-called gradient flow (continuous-time dynamics)

$$\dot{z}(t) = -\nabla \ell(z(t)) \quad z(0) = z_0$$

*Remark.* A solution to the ODE exists if the vector field is Lipschitz continuous

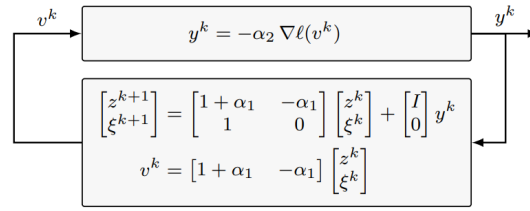
#### 4.2.3 Nesterov accelerated gradient method

Consider the following two-step algorithm: for all  $k \in \mathbb{N}$

$$\zeta^{k+1} = \zeta^k + \alpha_1(\zeta^k - \zeta^{k+1}) - \alpha_2 \nabla \ell(\zeta^k + \alpha_1(\zeta^k - \zeta^{k+1})), \quad \zeta^0, \zeta^{-1} \text{ given}$$

for some  $\alpha_1, \alpha_2 > 0$ . It admits the state-space representation

<sup>2</sup>The geometric series  $\sum_{i=0}^{\infty} \rho^i$  is equal to  $(1 - \rho)^{-1}$  for all  $\rho < 1$



More general updates can also be considered: for all  $k \in \mathbb{N}$

$$\zeta^{k+1} = \zeta^k + \alpha_1(\zeta^k - \zeta^{k+1}) - \alpha_2 \nabla \ell(\zeta^k + \alpha_3(\zeta^k - \zeta^{k-1})), \quad \zeta^0, \zeta^{-1} \text{ given}$$

for some  $\alpha_1, \alpha_2, \alpha_3 > 0$

## Chapter 5

# Parallel Optimization and Federated Learning

### Cost-coupled optimization for learning

In learning applications, we usually consider optimization problems in the form

$$\min_{z \in \mathbb{R}^d} \sum_{i=1}^N \ell_i(z)$$

where, for all  $i = 1, \dots, N$ , the cost function  $\ell_i : \mathbb{R}^d \rightarrow \mathbb{R}$  is local and private

### (Batch) gradient method for learning

Consider the optimization problem

$$\min_z \sum_{i=1}^N \ell_i(z)$$

The (batch) gradient method is: for each iteration  $k \in \mathbb{N}$

$$z^{k+1} = z^k - \alpha \sum_{i=1}^N \nabla \ell_i(z^k)$$

*Remark.* computation can be expensive

## 5.1 Incremental gradient method

Consider the optimization problem

$$\min_z \sum_{i=1}^N \ell_i(z)$$

Idea: rather than using the whole batch gradient at each  $k \in \mathbb{N}$ , just select one single "sample" per iteration.

The *incremental gradient method* is: for each iteration  $k \in \mathbb{N}$

$$z^{k+1} = z^k - \alpha \nabla \ell_{i^k}(z^k)$$

where  $i^k \in \{1, \dots, N\}$ .

Two rules for choosing index  $i^k$  at iteration  $k$ :

- Cyclic rule: choose  $i^k = 1, 2, \dots, N, 1, 2, \dots, N, \dots$
- Randomized rule: draw  $i^k \in \{1, \dots, N\}$  at random<sup>1</sup>, e.g according to a discrete uniform distribution

---

<sup>1</sup>Each  $i^k$  is a realization of a discrete random variable  $\mathcal{X} \sim \mathcal{U}\{1, N\}$  in which every one of the  $N$  possible outcome values has an equal probability

## 5.2 Stochastic Gradient Descent

consider the stochastic optimization problem

$$\min_z \mathbb{E}_{\mathcal{W}}[\ell(z, \mathcal{W})]$$

where  $\mathbb{E}_{\mathcal{W}}[\cdot]$  denotes the expected value with respect to the random variable  $\mathcal{W}$  (possibly having an unknown probability distribution  $p_{\mathcal{W}}(w)$ )

*Remark.* for all  $z$ , also  $\ell(z, \mathcal{W})$  is a random variable, whose probability distribution depends on  $p_{\mathcal{W}}$  and  $\ell$ . Moreover, the gradient  $\nabla \ell(z, \mathcal{W})$  at each  $z$  is a random quantity

Assumption: There exists an oracle that, given a realization  $\bar{w}$  of  $\mathcal{W}$ , returns the corresponding realization of the gradient  $\nabla \ell(\bar{z}, \bar{w})$  at any query point  $\bar{z}$

The stochastic gradient descent is: for each iteration  $k \in \mathbb{N}$  draw a realization  $w^k$  of  $\mathcal{W}$  and update

$$z^{k+1} = z^k - \alpha \nabla \ell(z^k, w^k)$$

### 5.2.1 Convergence

**Proposition 5.1** (convergence with constant step-size)

Assume:

- $\ell$  is a  $\mu$ -strongly convex function with  $L$ -Lipschitz continuous gradient (uniformly in its second argument)
- $\nabla \ell(z, \mathcal{W})$  is an unbiased estimatee of  $\nabla_z \mathbb{E}_{\mathcal{W}}[\ell(z, \mathcal{W})]$
- $\|\nabla \ell(z, \mathcal{W})\| \leq M$  almost surely<sup>2</sup> for some  $M > 0$

Let  $\{z^k\}_{k \in \mathbb{N}}$  be a realized sequence of the SGD with stepsize  $\alpha \leq 1/2\mu$ . Then

$$\|z^k - z^*\|^2 \leq (1 - 2\mu\alpha)^k \left( \|z^0 - z^*\|^2 - \frac{\alpha M^2}{2\mu} \right) + \frac{\alpha M^2}{2\mu}$$

*Analysis idea:* Use the sequence of realizations  $w^0, w^1, \dots, w^k$  of  $\mathcal{W}$  to approximate the original stochastic problem as

$$\mathbb{E}_{\mathcal{W}}[\ell(z, \mathcal{W})] \approx \frac{1}{K} \sum_{k=1}^L \ell(z, w^k)$$

### 5.2.2 Stochastic mini-batch gradient method

At each iteration  $k \in \mathbb{N}$ , one can use a set of realizations of  $\mathcal{W}$  to statistically approximate the (random) gradient of  $\ell$  at a given point  $z^k$

*Remark.* The set of realizations is denoted by  $\mathcal{I}^k$  and is usually called minibatch

For each  $k \in \mathbb{N}$ , choose  $\mathcal{I}^k \subset \{1, \dots, N\}$  with  $|\mathcal{I}^k| \ll N$ . The stochastic minibatch gradient method is

$$z^{k+1} = z^k - \alpha \sum_{i \in \mathcal{I}^k} \nabla \ell(z^k, w^i)$$

*Remark.* The minibatch can be chosen randomly or deterministically

*Remark.* This approach can also be applied to a "deterministic" optimization problem

---

<sup>2</sup>A sequence of random variables  $\{\mathcal{X}_k\}_{k \in \mathbb{N}}$  converges almost surely to the rv  $\mathcal{X}$  if  $\mathbb{P}(\lim_{k \rightarrow \infty} \mathcal{X}_k = \mathcal{X}) = 1$



### 5.3 Beyond SGD: Adaptive Momentum Estimation (Adam)

The ADAM algorithm reads as follows

- Mean and Variance (first momentum and second momentum)

$$\begin{aligned} m^{k+1} &= \beta_1 m^k + (1 - \beta_1) \nabla \ell(z^k, w^k), & \text{for some } \beta \in (0, 1) \\ v^{k+1} &= \beta_2 v^k + (1 - \beta_2) [\nabla \ell(z^k, w^k)]^2, & \text{for some } \beta_2 \in (0, 1) \end{aligned}$$

where the square operation  $[\cdot]^2$  is meant component-wise

- Construct the descent direction

$$\begin{aligned} \hat{m} &= \frac{1}{1 - \beta_1^{k+1}} m^{k+1} \\ \hat{v} &= \frac{1}{1 - \beta_2^{k+1}} v^{k+1} \\ d^k &= -\frac{\hat{m}}{\sqrt{\hat{v}} + \epsilon}, \quad \text{for some } \epsilon > 0 \end{aligned}$$

where the division in the last equation is meant element-wise

- Update of the solution estimate

$$z^{k+1} = z^k + \alpha d^k$$

### 5.4 Federated learning

Consider the optimization problem

$$\min_z \sum_{i=1}^N \ell(z; \mathcal{D}^i, p^i)$$

Paradigm:

- local private data  $\mathcal{D}^i = ([\mathcal{D}^i]_1, [\mathcal{D}^i]_2, \dots, [\mathcal{D}^i]_d)$  and  $p^i$
- learn common parameters  $z^* \in \mathbb{R}^d$  (common neural network)
- communication with a parameter server

### 5.5 Distributed learning

Consider the optimization problem

$$\min_z \sum_{i=1}^N \ell(z; \mathcal{D}^i, p^i)$$

Paradigm:

- local private data  $\mathcal{D}^i = ([\mathcal{D}^i]_1, [\mathcal{D}^i]_2, \dots, [\mathcal{D}^i]_d)$  and  $p^i$
- learn common parameters  $z^* \in \mathbb{R}^d$  (common neural network)
- communication with neighbours only



## Chapter 6

# Leader-Follower networks: Containment

Consider a network of  $N$  agents communicating/interacting according to a fixed, undirected graph  $G = (\{1, \dots, N\}, E)$ . Let  $x_i(t) \in \mathbb{R}^d$  be the state of agent  $i$  and for simplicity suppose  $d = 1$ , i.e.  $x_i(t) \in \mathbb{R}$ .

### 6.1 Definition

Suppose that agents are divided into two groups:

- $N_f$  followers
- $N - N_f$  leaders

Define  $x = \begin{bmatrix} x_f \\ x_l \end{bmatrix}$  with  $x_f \in \mathbb{R}^{N_f}$  the followers' state and  $x_l \in \mathbb{R}^{N-N_f}$  the leaders' state. Consistently, partition the Laplacian matrix as  $L = \begin{bmatrix} L_f & L_{fl} \\ L_{fl}^+ & L_l \end{bmatrix}$ . If leaders and followers run the same Laplacian-based distributed control law, then the Laplacian dynamics can be written as

$$\begin{bmatrix} \dot{x}_f(t) \\ \dot{x}_l(t) \end{bmatrix} = - \begin{bmatrix} L_f & L_{fl} \\ L_{fl}^+ & L_l \end{bmatrix} \begin{bmatrix} x_f(t) \\ x_l(t) \end{bmatrix}$$

### 6.2 Control law

Now, suppose that leaders and followers have different roles:

- $N_f$  followers running Laplacian dynamics
- $N - N_f$  leaders being stationary

The distributed control system is given by

$$\begin{aligned} \dot{x}_i(t) &= - \sum_{j \in \mathcal{N}_i} a_{ij}(x_i(t) - x_j(t)) & \forall i \in \{1, \dots, N_f\} \\ \dot{x}_i(t) &= 0 & \forall i \in \{N_f + 1, \dots, N\} \end{aligned}$$

and in compact form

$$\begin{bmatrix} \dot{x}_f(t) \\ \dot{x}_l(t) \end{bmatrix} = - \begin{bmatrix} L_f & L_{fl} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_f(t) \\ x_l(t) \end{bmatrix}$$

The dynamics of the followers is given by

$$\dot{x}_f(t) = -L_f x_f(t) - L_{fl} x_l$$

where  $x_l(t) = x_l(0) = x_l$



## Chapter 7

# Leader Follower networks: Formation control

### 7.1 analogy with mass-spring systems

Consider a platoon of  $N$  masses such that each mass  $i$  is connected with mass  $i - 1$  and  $i + 1$  through a spring with elastic constants respectively  $a_{i-1,i} = a_{i,i-1} > 0$  and  $a_{i+1,i} = a_{i,i+1} > 0$ . Let  $x_i \in \mathbb{R}$  be the position of mass  $i$

The elastic force at mass  $i$ ,  $F_{e,i}(x)$  is given by

$$F_{e,i}(x) = -a_{i,i-1}(x_i - x_{i-1}) - a_{i,i+1}(x_i - x_{i+1})$$

For each spring, we can write the associated elastic force as the negative gradient of the elastic (potential) energy, so that

$$F_{e,i}(x) = -\frac{\partial}{\partial x_i} \left( \frac{1}{2} a_{i,i-1} \|x_i - x_{i-1}\|^2 + \frac{1}{2} a_{i,i+1} \|x_i - x_{i+1}\|^2 \right)$$

Let us suppose that each spring can be written as the parallel of two springs with elastic constants respectively  $\frac{1}{2} a_{i,i-1} > 0$  and  $\frac{1}{2} a_{i,i+1} > 0$ .

Let  $x_i \in \mathbb{R}$  be the position of mass  $i$ . The elastic force at mass  $i$ ,  $F_{e,i}(x)$  is given by

$$F_{e,i} = -\left( \frac{1}{2} a_{i,i-1} + \frac{1}{2} a_{i,i-1} \right) (x_i - x_{i-1}) - \left( \frac{1}{2} a_{i,i+1} + \frac{1}{2} a_{i,i+1} \right) (x_i - x_{i+1})$$

As before, we can write the elastic force as the negative gradient of the elastic (potential) energy, i.e.

$$F_{e,i}(x) = -\frac{\partial}{\partial x_i} \left( \frac{1}{2} \frac{a_{i,i-1}}{2} \|x_i - x_{i-1}\|^2 + \frac{1}{2} \frac{a_{i,i-1}}{2} \|x_i - x_{i-1}\|^2 + \frac{1}{2} \frac{a_{i,i+1}}{2} \|x_i - x_{i+1}\|^2 + \frac{1}{2} \frac{a_{i,i+1}}{2} \|x_i - x_{i+1}\|^2 \right)$$

The total elastic (potential) energy of the mass-spring system can be written as

$$\begin{aligned} V(x) &= \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} \frac{1}{2} \frac{a_{i,j}}{2} \|x_i - x_j\|^2 \\ &= \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} V_{ij}(x_i, x_j) \end{aligned}$$

where we have defined  $\mathcal{N}_i := \{i - 1, i + 1\}$  and  $V_{ij}(x_i, x_j) := \frac{1}{2} \frac{a_{i,j}}{2} \|x_i - x_j\|^2$ .

Thus, the elastic force at mass  $i$  can be seen as the negative gradient of the energy, i.e.

$$\begin{aligned} F_{e,i}(x) &= -\frac{\partial}{\partial x_i} \sum_{j \in \mathcal{N}_i} (V_{ij}(x_i, x_j) + V_{ji}(x_j, x_i)) \\ &= -\frac{\partial}{\partial x_i} V(x) \end{aligned}$$

This formulation can be extended to more general systems in which masses are interconnected according to a topology described by an undirected graph  $G = (\{1, \dots, N\}, E)$

By adding a damping term on each mass, the system dynamics can be written as

$$\begin{aligned}\dot{x}_i &= v_i \\ m_i \dot{v}_i &= -v_i - \frac{\partial}{\partial x_i} V(x) \quad \forall i \in \{1, \dots, N\}\end{aligned}$$

where we have considered the damping coefficient equal to one.

If we assume that masses are small, we may write

$$v_i \approx -\frac{\partial}{\partial x_i} V(x)$$

so that the dynamics may be approximated by the following first order dynamics

$$\dot{x}_i = -\frac{\partial}{\partial x_i} V(x) \quad \forall i \in \{1, \dots, N\}$$

## 7.2 Formation Control

### 7.2.1 Laplacian dynamics: potential function interpretation

Consider a network of  $N$  agents communicating/interacting according to a fixed, undirected graph  $G$ . Let  $x_i(t) \in \mathbb{R}^d$  be the state of agent  $i$ . Let agents run a Laplacian dynamics

$$\dot{x}_i = -\sum_{j \in \mathcal{N}_i} a_{ij}(x_i - x_j) \quad \forall i \in \{1, \dots, N\}$$

We can rewrite it as

$$\dot{x}_i(t) = -\sum_{j \in \mathcal{N}_i} \frac{\partial}{\partial x_i} (V_{ij}(x_i, x_j) + V_{ji}(x_j, x_i)) \quad \forall i \in \{1, \dots, N\}$$

with  $V_{ij}(x) = \frac{1}{2} \frac{a_{i,j}}{2} \|x_i - x_j\|^2$ .

By recalling the definition of the total energy

$$V(x) = \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} V_{ij}(x_i, x_j)$$

the Laplacian dynamics

$$\dot{x} = -Lx$$

can be seen as a "gradient flow", i.e.

$$\dot{x} = -\nabla V(x)$$

Thus, the consensus configuration can be seen as a stationary point of  $V$ . This idea can be extended to general potential functions and applied to distributed control systems.

### 7.2.2 Distributed control based on potential functions

Consider a network of  $N$  autonomous agents communicating/interacting according to a fixed, undirected graph. Let  $x_i(t) \in \mathbb{R}^d$  be the state of agent  $i$ . Consider a global potential function defined as

$$V(x) = \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} V_{ij}(x_i, x_j)$$

such that (local) minima of the potential correspond to desired configurations of the team. The gradient flow dynamics  $\dot{x} = -\nabla V(x)$  turns out to be distributed. That is,

$$\dot{x}_i(t) = \sum_{j \in \mathcal{N}_i} \frac{\partial}{\partial x_i} (V_{ij}(x_i, x_j) + V_{ji}(x_j, x_i)) \quad \forall i \in \{1, \dots, N\}$$

We can define a desired formation by assigning a set of distances,  $d_{ij}$ , between neighbouring agents  $i$  and  $j$  in a suitable graph

The main idea for formation control is to define a potential function matching the sparsity of  $G$ ,  $V^{\text{form}}(x) = \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} V_{ij}^{\text{form}}(i, j)$ , such that a configuration  $x^{\text{form}}$  satisfying

$$\|x_i^{\text{form}} - x_j^{\text{form}}\| = d_{ij} \quad \forall (i, j) \in E$$

is a minimum of  $V$ .

In order to reach a formation with assigned distances  $d_{ij}$ , let us define

$$V_{ij}^{\text{form}}(x) = \frac{1}{8} (\|x_i - x_j\|^2 - d_{ij}^2)^2$$

with corresponding (global) potential function  $V^{\text{form}}(x) = \sum_{i=1}^N \sum_{j \in \mathcal{N}_i} V_{ij}^{\text{form}}(x_i, x_j)$ .

The gradient flow dynamics of each agent  $i$  is given by

$$\dot{x}_i = - \sum_{j \in \mathcal{N}_i} \frac{\partial}{\partial x_i} (V_{ij}^{\text{form}}(x_i, x_j) + V_{ji}^{\text{form}}(x_j, x_i)) \quad \forall i \in \{1, \dots, N\}$$

which reads as

$$\dot{x}_i = - \sum_{j \in \mathcal{N}_i} (\|x_i - x_j\|^2 - d_{ij}^2) (x_i - x_j) \quad \forall i \in \{1, \dots, N\}$$

### 7.2.3 Inter-agent collision avoidance

This dynamics has multiple equilibrium points, including the desired formation in which the agents are at the assigned distances. In particular, the consensual solution  $x_1 = x_2 = \dots = x_N$  is an (undesired) equilibrium.

Such a "degenerate" equilibrium can be avoided by means of additional "collision avoidance" potential functions  $V_{ij}^{\text{ca}}(x_i, x_j)$  such that

$$\lim_{\|x_i - x_j\| \rightarrow 0} V_{ij}^{\text{ca}}(x_i, x_j) = +\infty$$

A possible solution is a barrier function given

$$V_{ij}^{\text{ca}} = -\log(\|x_i - x_j\|)$$

### 7.2.4 Obstacle avoidance

Similarly, barrier potential functions  $V^{\text{obs}}(x_i)$ , depending only on the state of agent  $x_i$ , can be used to avoid obstacles.

The formation control dynamics becomes

$$\dot{x}_i = - \frac{\partial V^{\text{form}}}{\partial x_i} - \frac{\partial V^{\text{ca}}(x)}{\partial x_i} - \nabla V^{\text{obs}}(x) \quad \forall i \in \{1, \dots, N\}$$





## Chapter 8

# Distributed Cost-Coupled Optimization

Scenario: Network of  $N$  agents communicating according to a digraph  $G = (\{1, \dots, N\}E)$

Goal: define a distributed algorithm able to solve an optimization problem

### 8.1 Cost-Coupled Optimization Set-up

Consider the optimization problem

$$\min_{z_i \in Z} \sum_{i=1}^N \ell_i(z)$$

with  $z \in \mathbb{R}^d$

- $N$  agents communicate over graph  $G$
- agent  $i$  knows  $\ell_i$  and  $Z$  only
- $z_i^k$  solution estimate of  $i$

### 8.2 Distributed Gradient Algorithm

The centralized gradient algorithm reads

$$z^{k+1} = z^k - \alpha^k \sum_{h=1}^N \nabla \ell_h(z^k)$$

- One single estimate  $z^k \in \mathbb{R}^d$  initialized at some  $z^0$
- $d^k = - \sum_{h=1}^N \nabla \ell_h(z^k)$
- $\sum_{h=1}^N \nabla \ell_h(z^k)$  not available to each agent

The distributed gradient algorithm reads

$$\begin{aligned} v_i^{k+1} &= \sum_{j \in \mathcal{N}_i} a_{ij} z_j^k \\ z_i^{k+1} &= v_i^{k+1} - \alpha^k \nabla \ell_i(v_i^{k+1}) \end{aligned}$$

with  $z_i^0 \in \mathbb{R}^d$

### 8.3 Assumptions

- Assumption 1. Let  $a_{ij}$ ,  $i, j \in \{1, \dots, N\}$  be nonnegative entries of a weighted adjacency matrix  $A$  associated to the undirected and connected graph  $G$ , with  $a_{ij} > 0$  and  $A$  doubly stochastic.
- Assumption 2. The step-size sequence  $\{\alpha^k\}_{k \in \mathbb{N}}$  satisfies the conditions

$$\alpha^k \geq 0, \quad \sum_{k=0}^{\infty} \alpha^k = \infty, \quad \sum_{k=0}^{\infty} (\alpha^k)^2 < \infty$$

- Assumption 3a. Let the following conditions hold
  - each  $\ell_i : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex and has bounded gradients, i.e., there exists a scalar  $C_i > 0$  such that  $\|\nabla \ell_i(z)\| \leq C_i$  at any  $z \in \mathbb{R}^d$
  - the optimization problem has at least one optimal solution, i.e. the optimal solution set  $Z^* = \{z \in \mathbb{R}^d \mid \ell(z) = \ell^*\}$  is nonempty, where  $\ell^*$  denotes the optimal value of the problem

#### Theorem 8.1

Let assumptions 1,2 and 3a hold true. Then, the sequences of local solution estimates  $\{z_i^k\}_{k \in \mathbb{N}, i=1, \dots, N}$ , generated by the Distributed Gradient Algorithm, asymptotically converge to a common solution  $z^*$  of the problem, i.e. for all  $i = 1, \dots, N$  it holds

$$\lim_{k \rightarrow \infty} \|z_i^k - z^*\| = 0$$

for some  $z^*$  in the optimal solution set  $Z^*$

### 8.4 Distributed Projected Subgradient Algorithm

The previous result can be extended to nonsmooth convex functions  $\ell_i$  and  $z$  constrained to a closed, convex set  $Z \subseteq \mathbb{R}^d$  by implementing the following distributed step

$$\begin{aligned} v_i^{k+1} &= \sum_{j \in \mathcal{N}_i} a_{ij} z_j^k \\ z_i^{k+1} &= P_Z \left( v_i^{k+1} - \alpha^k \tilde{\nabla} \ell_i(v_i^{k+1}) \right) \end{aligned}$$

where  $P_Z(\cdot)$  denotes the Euclidean projection onto  $Z$  while  $\tilde{\nabla} \ell_i$  denotes a subgradient of  $\ell_i$

Convergence to an optimal solution can be achieved by substituting Assumption 3a with Assumption 3b. Assume  $Z \subseteq \mathbb{R}^d$  be a closed and bounded set. Let the following conditions hold

- each  $\ell_i : Z \rightarrow \mathbb{R}$  is convex and has bounded subgradients, i.e. there exists a scalar  $C_i > 0$  such that  $\|\tilde{\nabla} \ell_i(z)\| \leq C_i$  for any subgradient  $\tilde{\nabla} \ell_i(z)$  of  $\ell_i$  at any  $z \in \mathbb{R}^d$
- the optimization problem has at least one optimal solution

### 8.5 Improving the Distributed Gradient Algorithm

Recall the Distributed Gradient Algorithm

$$z_i^{k+1} = \sum_{j \in \mathcal{N}_i} a_{ij} z_j^k + \alpha d_i^k$$

We need to replace  $d_i^k = -\nabla \ell_i(z_i^k)$  with a better descent direction. The idea is to use a local descent direction  $d_i^k$  that aims at reconstructing the correct gradient

$$d_i^k \xrightarrow[k \rightarrow \infty]{} \frac{1}{N} \sum_{h=1}^N \nabla \ell_h(z_h^k)$$

### 8.5.1 Dynamic Average Consensus

- $N$  agents measure only a local, time-varying signal  $r_i^k : \mathbb{N} \rightarrow \mathbb{R}$ , for all  $i = 1, \dots, N$
- Each agent  $i$  wants to estimate/track the average signal

$$\bar{r}^k = \frac{1}{N} \sum_{i=1}^N r_i^k$$

Each agent  $i$  maintains and updates a local estimate  $s_i^k$  so that  $\lim_{k \rightarrow \infty} \|s_i^k - \bar{r}^k\| = 0$ . A distributed algorithm based on a (perturbed) consensus is the following

$$s_i^{k+1} = \sum_{j \in \mathcal{N}_i} a_{ij} s_j^k + (r_i^{k+1} - r_i^k)$$

with  $s_i^0 = r_{i,0}$ , where the first term represents a consensus term, and the second is a local innovation

#### Convergence

Let the first-order differences be relatively bounded, i.e.

$$\|r_i^{k+1} - r_i^k\| \leq C_1$$

for some  $C_1 > 0$ . Then, the local sequence generated by the dynamic average consensus satisfies

$$\lim_{k \rightarrow \infty} \|s_i^k - \bar{r}^k\| \leq C_2$$

for some  $C_2 > 0$  and for all  $i = 1, \dots, N$

*Remark.* The local states achieve a nonzero steady-state error, which is uniformly upper bounded

*Remark.* If  $\|r_i^{k+1} - r_i^k\| \rightarrow 0$ , then the steady-state error is zero, i.e. perfect tracking is achieved

## 8.6 The Gradient Tracking Algorithm

The strategy is to use the previous signal tracking mechanism to reconstruct the gradient of  $\sum_{i=1}^N \ell_i$ . The local signal is

$$r_i^k = \nabla \ell_i(z_i^k)$$

The local estimate  $s_i^k$  is updated as

$$s_i^{k+1} = \sum_{j \in \mathcal{N}_i} a_{ij} s_j^k + (\nabla \ell_i(z_i^{k+1}) - \nabla \ell_i(z_i^k)), \quad s_i^0 = \nabla \ell_i(z_i^0)$$

and use its opposite as local descent direction  $d_i^k = -s_i^k$

$$d_i^k \xrightarrow[k \rightarrow \infty]{-} \frac{1}{N} \sum_{h=1}^N \nabla \ell_h(z_h^k)$$

*Remark.* If each  $\nabla \ell_i(z_i^k)$  converges to a constant value then perfect tracking is achieved

*Remark.* The convergence of  $\nabla \ell_i(z_i^k)$  is intertwined with the convergence of  $z_i^k$ , e.g. to  $z^*$

The Gradient Tracking Algorithm is

$$\begin{aligned} z_i^{k+1} &= \sum_{j \in \mathcal{N}_i} a_{ij} z_j^k - \alpha s_i^k & z_i^0 &\in \mathbb{R} \\ s_i^{k+1} &= \sum_{j \in \mathcal{N}_i} a_{ij} s_j^k + \nabla \ell_i(z_i^{k+1}) - \nabla \ell_i(z_i^k) & s_i^0 &= \nabla \ell_i(z_i^0) \end{aligned}$$

### 8.6.1 Assumptions

- Assumption 1. Let  $a_{ij}$ ,  $i, j \in \{1, \dots, N\}$  be non-negative entries of a weighted adjacency matrix  $A$  associated to the undirected and connected graph  $G$ , with  $a_{ii} > 0$  and  $A$  doubly stochastic.
- Assumption 2. For all  $i \in \{1, \dots, N\}$  each cost function  $\ell_i : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfies the following conditions
  - it is strongly convex with coefficient  $\mu > 0$
  - it has Lipschitz continuous gradient with constant  $L > 0$

*Remark.* Strong convexity implies uniqueness of the optimal solution  $z^* = \arg \min_x \sum_{i=1}^N \ell_i(z)$

### 8.6.2 Convergence Result

#### Theorem 8.2

Let assumptions 1 and 2 hold true. Then, there exists  $\alpha^* > 0$  such that for all  $\alpha \in (0, \alpha^*)$  the sequences of local solution estimates  $\{z_i^k\}_{k \in \mathbb{N}}$ ,  $i = 1, \dots, N$  generated by the Gradient Tracking algorithm asymptotically converge to a (consensual) solution  $z^*$  of the problem, i.e. for all  $i = 1, \dots, N$  it holds

$$\lim_{k \rightarrow \infty} \|z_i^k - z^*\| = 0$$

Moreover, the convergence rate is linear/the stability is exponential, i.e., for all  $i = 1, \dots, N$  there exists  $\rho \in (0, 1)$  such that

$$\|z_i^k - z^*\| \leq \rho \|z_i^{k-1} - z^*\| \leq \rho^k \|z_i^0 - z^*\|$$

for all  $k \in \mathbb{N}$

## Chapter 9

# Distributed Aggregative Optimization

Consider  $N$  robots in the plane that want to optimize their positions  $z_i \in \mathbb{R}^2$ , for all  $i = 1, \dots, N$  to perform multi-robot surveillance. Let:

- $r_0 \in \mathbb{R}^2$  be a target to protect
- $r_i \in \mathbb{R}^2$  be the intruder associated to robot  $i$
- $\sigma(z) = \frac{1}{N} \sum_{i=1}^N z_i$  is the barycenter of the robots
- Local cost function of robot  $i$

$$\ell_i(z_i, \sigma(z)) = \gamma_i \|z_i - r_i\|^2 + \|\sigma(z) - r_0\|^2$$

with  $z \in \mathbb{R}^{2N}$  the stack of  $z_1, \dots, z_N$  and  $\gamma_i > 0$  being a tradeoff parameter.

### 9.1 Aggregative optimization

Let us consider aggregative optimization problems in the form

$$\min_{z_1, \dots, z_N} \sum_{i=1}^N \ell_i(z_i, \sigma(z))$$

where the aggregative variable  $\sigma(z)$  is defined as

$$\sigma(z) = \frac{1}{N} \sum_{i=1}^N \phi_i(z_i)$$

where

- $z = (z_1, \dots, z_N)$ , with each  $z_i \in \mathbb{R}^{n_i}$ , for all  $i = 1, \dots, N$
- $\ell_i : \mathbb{R}^{n_i} \times \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\phi_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^d$ , for all  $i = 1, \dots, N$

For scalar states,  $z_i \in \mathbb{R}$ , the *centralized gradient* method at iteration  $k$  reads as

$$z_i^{k+1} = z_i^k - \alpha \frac{\partial}{\partial z_i} \left( \sum_{j=1}^N \ell_j(z_j, \sigma(z_1, \dots, z_N)) \right) \Bigg|_{z_1=z_1^k, \dots, z_N=z_N^k}$$

for all  $i = 1, \dots, N$  where  $\alpha > 0$  is the stepsize

**Gradient computation (scalar case)**

Since the cost function is a composite function, we need the chain rule to compute its derivative with respect to  $z_i$

$$\begin{aligned} & \left. \frac{\partial}{\partial z_i} \left( \sum_{j=1}^N \ell_j(z_j, \sigma(z_1, \dots, z_N)) \right) \right|_{z_1=z_1^k, \dots, z_N=z_N^k} \\ &= \left. \frac{\partial}{\partial z_i} \ell_i(z_i, \sigma) \right|_{z_i=z_i^k, \sigma=\frac{1}{N} \sum_{j=1}^N \phi_j(z_j^k)} \\ &+ \sum_{j=1}^N \left. \frac{\partial}{\partial \sigma} \ell_j(z_j, \sigma) \right|_{z_j=z_j^k, \sigma=\frac{1}{N} \sum_{j=1}^N \phi_j(z_j^k)} \cdot \left. \frac{\partial \sigma(z_1, \dots, z_N)}{\partial z_i} \right|_{z_1=z_1^k, \dots, z_N=z_N^k} \end{aligned}$$

Notice that  $\frac{\partial \sigma(z_1, \dots, z_N)}{\partial z_i} = \frac{1}{N} \frac{d}{dz_i} \phi_i(z_i)$  can be computed locally

**Gradient computation (vector case)**

As in the scalar case, we use the chain rule to compute the gradient of the composite function. The  $i$ -th block of the gradient, denoted as  $\left[ \nabla \left( \sum_{j=1}^N \ell_j(z_j, \sigma(z_1, \dots, z_N)) \right) \right]_i \in \mathbb{R}^{n_i}$ , is given by

$$\begin{aligned} & \left[ \nabla \left( \sum_{j=1}^N \ell_j(z_j, \sigma(z_1, \dots, z_N)) \right) \right]_i \\ &= \nabla_1 \ell_i(z_i, \sigma) \Big|_{z_j=z_j^k, \sigma=\frac{1}{N} \sum_{j=1}^N \phi_j(z_j^k)} \\ &+ \frac{1}{N} \nabla \phi_i(z_i) \Big|_{z_i=z_i^k} \cdot \sum_{j=1}^N \nabla_2 \ell_j(z_j, \sigma) \Big|_{z_j=z_j^k, \sigma=\frac{1}{N} \sum_{j=1}^N \phi_j(z_j^k)} \end{aligned}$$

**9.2 Distributed aggregative optimization**

In a distributed context, each agent  $i$

- knows only  $\ell_i$  and  $\phi_i$
- maintains an estimate  $z_i^k$  of  $z_i^*$
- maintains an estimate  $s_i^k$  of  $\phi(z^k) = \frac{1}{N} \sum_{j=1}^N \phi_j(z_j^k)$
- maintains an estimate  $v_i^k$  of  $\sum_{j=1}^N \nabla_2 \ell_j(z_j^k, \sigma(z^k))$

The "tracking" idea of gradient tracking algorithm is applied to aggregative optimization

$$\begin{aligned} z_i^{k+1} &= z_i^k - \alpha (\nabla_1 \ell_i(z_i^k, s_i^k) + \nabla \phi_i(z_i^k) v_i^k) & z_i^0 &\in \mathbb{R}^{n_i} \\ s_i^{k+1} &= \sum_{j \in \mathcal{N}_i} a_{ij} s_j^k + \phi_i(z_i^{k+1}) - \phi_i(z_i^k) & s_i^0 &= \phi_i(z_i^0) \\ v_i^{k+1} &= \sum_{j \in \mathcal{N}_i} a_{ij} v_j^k + \nabla_2 \ell_i(z_i^{k+1}, s_i^{k+1}) - \nabla_2 \ell_i(z_i^k, s_i^k) & v_i^0 &= \nabla_2 \ell_i(z_i^0, s_i^0) \end{aligned}$$

**Theorem 9.1** (aggregative tracking distributed optimization algorithm: convergence)

Assume  $G$  is a strongly connected and aperiodic digraph, and  $A$  is doubly stochastic. Assume that each function  $\ell_i$  is strongly convex, the gradients  $\nabla_1 \ell_i$  and  $\nabla_2 \ell_i$  are Lipschitz continuous, and  $\phi_i$  is differentiable and Lipschitz continuous.

Then, there exists  $\alpha^*$  such that for all  $\alpha \in (0, \alpha^*)$  the sequences of local solution estimates  $\{z_1^k, \dots, z_N^k\}_{k \in \mathbb{N}}$  generated by the aggregative tracking distributed optimization algorithm satisfy

$$\lim_{k \rightarrow \infty} \|z_i^k - z_i^*\| = 0$$

at a linear rate, for all  $i = 1, \dots, N$

### 9.2.1 Extension to online aggregative optimization

Consider a time-varying instance of the problem

$$\begin{aligned} \min_z \quad & \sum_{i=1}^N \ell_i^k(z_i, \sigma^k(z)) \\ \text{subj. to } & z_i \in Z_i^k, \quad \forall i = 1, \dots, N \end{aligned}$$

where  $Z_i^k \subset \mathbb{R}^{n_i}$  are local constraint sets.

The goal is to design an algorithm generating a sequence  $\{z_i^k\}_{k \in \mathbb{N}}$  that "tracks" the solution  $z^{k,*} = (z_1^{k,*}, \dots, z_N^{k,*})$  of the  $k$ -th problem instance.

*Remark.* A regret  $R_K := \sum_{k=1}^K \sum_{i=1}^N \left( \ell_i^k(z_i^k, \sigma^k(z^k)) - \ell_i^k(z_i^{k,*}, \sigma^k(z^{k,*})) \right)$  can be introduced for the analysis. Under suitable assumptions, it can be proven that  $R_K \leq C_1 + C_2 K$  for some constants  $C_1$  and  $C_2$





## Chapter 10

# Learning with Neural Networks

In supervised learning we have labeled data whose generic element(sample) is composed by

- an input point  $\mathcal{D} \in \mathbb{R}^d$ , i.e. a vector  $\mathcal{D} = ([\mathcal{D}]_1, \dots, [\mathcal{D}]_d)$
- a label  $p \in \mathbb{R}$  associated to the input point

A dataset is usually made by  $\mathcal{M}$  samples, i.e.  $((\mathcal{D}^1, p^1), \dots, (\mathcal{D}^{\mathcal{M}}, p^{\mathcal{M}}))$

Goal: reconstruct the unknown input-output map between inputs and labels

Strategy: approximate the map as a nonlinear function  $\phi(\cong; \cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$  parametrized by  $\cong$

Problem: find the best parameters  $\cong^*$  based on data

### 10.1 Basic element: the neuron model

A generic neuron  $h$  is a computational unit that

- has a set of weights  $u_h \in \mathbb{R}^d$
- elaborates a vector  $x \in \mathbb{R}^d$
- computes a scalar quantity

$$x_h^+ = \phi(x^\top u_h)$$

with  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  being the *activation function*



# Chapter 11

## Multi-Robot Safety Controllers

Consider a nonlinear system

$$\dot{x}(t) = f(x(t), u(t)), \quad x \in \mathbb{R}^n \quad u \in U \subset \mathbb{R}^m$$

Let us define a safe (state) set

$$X^s := \{x \in \mathbb{R}^n | V^s(x) \geq 0\}$$

for some sufficiently regular function  $V^s : \mathbb{R}^n \rightarrow \mathbb{R}$

Goal: design a feedback control law  $\kappa(x), \kappa : \mathbb{R}^n \rightarrow \mathbb{R}^m$  such that the set  $X^s$  is forward invariant

### 11.1 Control Barrier Functions

Consider the time derivative of  $V^s(x(t))$  along the system trajectories

$$\begin{aligned} \frac{d}{dt} V^s(x(t)) &= \nabla V^s(x(t))^\top f(x(t)) + \sum_{h=1}^m \nabla V^s(x(t))^\top g_h(x(t)) u_h(x(t)) \\ &= L_f V^s(x(t)) + L_g V^s(x(t)) u(t) \end{aligned}$$

where  $L_f$  and  $L_g$  represent Lie derivatives. We say that  $V^s$  is a Control Barrier Function (CBF) if there exists a continuous, strictly increasing function  $\gamma : \mathbb{R} \rightarrow \mathbb{R}$ , with  $\gamma(0) = 0$ , such that

$$\sup_{u \in U} \{L_f V^s(x) + L_g V^s(x) u + \gamma(V^s(x))\} \geq 0$$

for all  $x \in \mathbb{R}^n$ .

The above inequality induced by a CBF is called a *control barrier certificate*.

#### 11.1.1 Admissible Control Space

We can define the set of admissible (safe) controllers for a given state  $x$  as

$$U^s(x) = \{u \in \mathbb{R}^m | L_f V^s(x) + L_g V^s(x) u + \gamma(V^s(x))\} \geq 0$$

i.e., the set of inputs satisfying the control barrier certificates.

*Remark.* It can be shown that (sufficiently regular) feedback controllers satisfying the above condition render  $X^s$  forward invariant and asymptotically stable.

#### 11.1.2 Safety Filters via Control Barrier Certificates

Let  $u^{\text{ref}} \in \mathbb{R}^m$  be a (possibly unsafe) reference input, e.g., obtained by a higher-level controller. The safety controller is designed to be "minimally invasive", i.e. it aims at altering the reference controller as little as possible. A safe policy  $u = \kappa(x)$  can be designed as

$$\begin{aligned} \kappa(x) &= \arg \min_{u \in \mathbb{R}^m} \|u - u^{\text{ref}}(x)\|^2 \\ \text{subj. to } &-L_f V^s(x) - L_g V^s(x) u - \gamma(V^s(x)) \leq 0 \end{aligned}$$

## 11.2 Single-Robot Obstacle Avoidance: Single Integrator Model

Consider a robot modeled as a single integrator, i.e.

$$\dot{x} = u$$

with  $x \in \mathbb{R}^d$  being the state of the robot and  $u \in \mathbb{R}^d$  the control input. The goal is to keep the robot at a safety distance  $d$  from an obstacle.

Consider a CBF defined as

$$V^s(x) = \|x - x_{\text{obs}}\|^2 - d^2$$

where  $d$  denotes the safety distance to keep from the obstacle. The gradient wrt  $x$  is given by

$$\nabla V^s(x) = 2(x - x_{\text{obs}})$$

The safety, CBF-based policy can be computed as

$$\begin{aligned} \kappa_s(x) = \arg \min_{u \in \mathbb{R}^d} & \|u - u^{\text{ref}}(x)\|^2 \\ \text{subj. to} & -2(x - x_{\text{obs}})^\top u - \gamma(\|x - x_{\text{obs}}\| - d^2) \leq 0 \end{aligned}$$

*Remark.* For a given, fixed  $x \in \mathbb{R}^d$ , the optimization problem above is a quadratic program

*Remark.* The safety controller is meant to prevent collisions while being "minimally invasive" in the sense that it aims at altering the reference controller (coming from a higher level) as little as possible

*Remark.* From a practical viewpoint the above problem requires a computation time, while the safety policy should be available in continuous time

## 11.3 Multi-Robot Collision Avoidance: Single Integrators

Consider a team of  $N$  robots modeled as single integrators, i.e.

$$\dot{x}_i = u_i$$

with  $x_i \in \mathbb{R}^d$  being the state of robot  $i$  and  $u_i \in \mathbb{R}^d$  being the control input. The goal is to keep the robots at a safety distance greater than  $d$ .

To keep the notation simpler, we consider a scalar robot state, and input, i.e.  $x_i, u_i \in \mathbb{R}$ .

### 11.3.1 Control Barrier Function for Multi-Robot Collision Avoidance

Consider a "pair-wise", local CBF defined as

$$V_{i,j}^s(x_i, x_j) = \|x_i - x_j\|^2 - d^2$$

where  $d$  denotes the safety distance between robots  $i$  and  $j$ .

The gradient wrt  $x_i$  and  $x_j$  is given by

$$\begin{aligned} \nabla_1 V_{i,j}^s(x_i, x_j) &= 2(x_i - x_j) \\ \nabla_2 V_{i,j}^s(x_i, x_j) &= 2(x_j - x_i) \end{aligned}$$

while the gradient wrt all the other states is zero.

Let  $u \in \mathbb{R}^N$  denote the stack of inputs of the  $N$  robots. The admissible control space for collision avoidance is obtained by considering the intersection of all the constraint sets related to a pair-wise control barrier certificates. Thus,

$$U^s(x) \{u \in \mathbb{R}^N \mid -\nabla_1 V_{i,j}^s(x_i, x_j)^\top u_i - \nabla_2 V_{i,j}^s(x_i, x_j) u_j - \gamma V_{i,j}^s(x_i, x_j) \leq 0, \forall j \in \mathcal{N}_i, \forall i \in \{1, \dots, N\}\}$$

where  $\mathcal{N}_i$  includes all the robots  $j$  that may collide with robot  $i$ .

*Remark.* If there are no apriori restrictions on the motion of the robots and they share the same area, we have  $\mathcal{N}_i = \{1, \dots, N\}$ . However, for some state configurations  $\mathcal{N}_i$  may involve only a restricted number of neighbouring robots.

### 11.3.2 Centralized Safety Controller for Multi-Robot Collision Avoidance

The safety, CBF-based policy is obtained by the following optimization problem

$$\begin{aligned} \kappa(x) = \arg \min_{u \in \mathbb{R}^N} & \sum_{i=1}^N \|u_i - u_i^{\text{ref}}\|^2 \\ \text{subj. to} & \quad -\nabla_1 V_{i,j}^s(x_i, x_j)^\top u_i - \nabla_2 V_{i,j}^s(x_i, x_j) u_j - \gamma V_{i,j}^s(x_i, x_j) \leq 0 \\ & \quad \|u_i\| \leq u_i^{\text{max}} \\ & \quad \forall j \in \mathcal{N}_i, \forall i \in \{1, \dots, N\} \end{aligned}$$

where  $u_1^{\text{ref}}, \dots, u_N^{\text{ref}}$  are reference inputs for each robot coming from a higher control layer, while  $u_1^{\text{max}}, \dots, u_N^{\text{max}}$  are input bounds for each robot.

### 11.3.3 Decentralized Safety Controller for Multi-Robot Collision Avoidance

In cooperative (distributed) multi-robot scenarios it is desirable not to have a central coordinator. A more constrained version of the problem can be formulated to obtain a decentralized solver.

Safety controller for robot  $i$ :

$$\begin{aligned} \kappa(x) = \arg \min_{u \in \mathbb{R}^N} & \sum_{i=1}^N \|u_i - u_i^{\text{ref}}\|^2 \\ \text{subj. to} & \quad -2\nabla_1 V_{i,j}^s(x_i, x_j)^\top u_i - \frac{1}{2}\gamma V_{i,j}^s(x_i, x_j) \leq 0 \\ & \quad \|u_i\| \leq u_i^{\text{max}} \\ & \quad \forall j \in \mathcal{N}_i \end{aligned}$$

*Remark.* The decentralized safety controller requires robot  $i$  to receive information by all other robots that may incur in a collision.

*Remark.* More complex models may be considered to take into account dynamics constraints

*Remark.* When more complex dynamics are considered together with input bounds, the (centralized or decentralized) optimization problem may become infeasible (safety cannot be guaranteed). More sophisticated approaches are required to synthesize safety barrier certificates with guaranteed feasibility.

*Remark.* When the objectives of multiple robots conflict with the safety barrier certificates, robots might incur into a deadlock, i.e., a configuration that is safe but does not allow for task completion.



## Chapter 12

# Stability Analysis for Linear Averaging

Consider  $N$  agents in a network modeled as a graph  $G$  with associated weighted adjacency matrix  $A$   
Linear averaging distributed algorithm

$$x_i^{k+1} = \sum_{j \in \mathcal{N}_i} a_{ij} x_j^k, \quad x^0 \text{ is given} \quad i = 1, \dots, N$$

where  $x_i^k \in \mathbb{R}$  is the state of agent  $i$  at time  $k$  and  $a_{ij} \geq 0$

Each agent  $i$  uses only the states of neighbours  $j \in \mathcal{N}_i$ , thus it is a distributed algorithm. The goal is to reach consensus.

The linear averaging is

$$\begin{bmatrix} x_1^{k+1} \\ \vdots \\ x_N^{k+1} \end{bmatrix} = \begin{bmatrix} a_{11} & \cdots & a_{1N} \\ \vdots & \ddots & \vdots \\ a_{N1} & \cdots & a_{NN} \end{bmatrix} \begin{bmatrix} x_1^k \\ \vdots \\ x_N^k \end{bmatrix}, \quad x^0 = \begin{bmatrix} x_1^0 \\ \vdots \\ x_N^0 \end{bmatrix}$$

and can be compactly written as

$$x^{k+1} = Ax^k, \quad x^0 = \begin{bmatrix} x_1^0 \\ \vdots \\ x_N^0 \end{bmatrix}$$

with  $x^k \in \mathbb{R}^n$  and the matrix  $A \in \mathbb{R}^{N \times N}$  "compatible" with the graph  $G$ . We want to study the stability properties of the equilibrium point(s) of such an autonomous system

The compact version of the linear averaging is an LTI dynamical system

$$x^{k+1} = Ax^k \quad x^0 \text{ given}$$

The equilibrium points can be computed as

$$x_{\text{eq}} = Ax_{\text{eq}} \implies (I - A)x_{\text{eq}} = 0$$

- Is this linear system stable or unstable?
- Spectral properties of  $A$  are crucial for the analysis

### 12.1 Consensus result

#### Theorem 12.1

Consider the discrete-time averaging system with associated digraph  $G$  and weighted adjacency matrix  $A$ . Assume  $G$  is strongly connected and aperiodic, and  $A$  is row stochastic. Then

1. there exists a positive, left eigenvector  $w \in \mathbb{R}^n$  such that

$$\lim_{k \rightarrow \infty} x^k = \mathbf{1} \frac{w^T x^0}{w^T \mathbf{1}} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \frac{\sum_{i=1}^N w_i x_i^0}{\sum_{i=1}^N w_i}$$

i.e., consensus is reached to  $\frac{\sum_{i=1}^N w_i x_i^0}{\sum_{i=1}^N w_i}$

2. if, additionally  $A$  is doubly stochastic, then

$$\lim_{k \rightarrow \infty} x^k = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \frac{\sum_{i=1}^N x_i^0}{N}$$

i.e., *average* consensus is reached

## 12.2 Irreducible and primitive matrices

A square matrix  $A \in \mathbb{R}^{N \times N}$  is

- non-negative if all its entries are greater than or equal to zero and we write  $A \geq 0$
- irreducible if  $\sum_{k=0}^{N-1} A^k > 0$   
Equivalently, there exists no permutation matrix  $T$  such that  $TAT^T$  is block triangular
- primitive if there exists an integer  $k \in \{1, \dots, N\}$  such that  $A^k$  is positive
- positive<sup>1</sup> if all its entries are greater than zero and we write  $A > 0$

## 12.3 Spectral properties of square matrices

The spectrum of a square matrix  $A$ ,  $\text{spec}(A)$ , is the set of eigenvalues of  $A$ . The spectral radius,  $\rho(A)$ , is the maximum norm of the eigenvalues of  $A$ .

**Theorem 12.2** (Gershgorin Theorem)

For any (real or complex) square matrix  $A$ , it holds

$$\text{spec}(A) \subset \bigcup_{i=1}^N \left\{ s \in \mathbb{C} \mid |s - a_{ii}| \leq \sum_{j=1, j \neq i}^N |a_{ij}| \right\}$$

That is, the union of disks in  $\mathbb{C}$  centered at each diagonal entry  $a_{ii}$  with radius  $\sum_{j=1, j \neq i}^N |a_{ij}|$

**Theorem 12.3**

Let  $G$  be a weighted digraph with  $N \geq 2$  nodes and with weighted adjacency matrix  $A$ . Then  $A$  is irreducible iff  $G$  is strongly connected

**Theorem 12.4**

Let  $G$  be a weighted digraph with  $N \geq 2$  nodes and with weighted adjacency matrix  $A$ . Then  $A$  is primitive iff  $G$  is strongly connected and aperiodic

**Theorem 12.5** (Perron-Frobenius theorem)

Let  $A \in \mathbb{R}^{N \times N}$  with  $N \geq 2$  be a non-negative matrix, then

1. there exists a real eigenvalue  $\lambda \geq |\mu| \geq 0$  for all other eigenvalues  $\mu \in \mathbb{C}$  of  $A$

---

<sup>1</sup>Not to be confused with a *positive definite* matrix, denoted as  $A \succ 0$



2. the right eigenvector  $v \in \mathbb{R}^N$  and left eigenvector  $w \in \mathbb{R}^N$  associated to  $\lambda$  can be selected non-negative  
If, additionally  $A \in \mathbb{R}^{N \times N}$  is irreducible, then
3. the eigenvalue  $\lambda \in \mathbb{R}$  is strictly positive and simple
4. the right eigenvector  $v \in \mathbb{R}^N$  and left eigenvector  $w \in \mathbb{R}^N$  associated to  $\lambda$  are unique and positive, possibly up to rescaling  
If, additionally  $A \in \mathbb{R}^{N \times N}$  is primitive, then
5. The eigenvalue  $\lambda > |\mu|$  for all other eigenvalues  $\mu \in \mathbb{C}$  of  $A$

### 12.3.1 Spectral properties of row-stochastic matrices

**Lemma.** For a row-stochastic matrix  $A \in \mathbb{R}^{N \times N}$

1.  $\lambda = 1$  is an eigenvalue of  $A$
2.  $\text{spec}(A)$  is a subset of the unit disk

*Proof.* For the first statement, use  $A\mathbf{1} = \mathbf{1}$ . For the second statement use the Gershgorin Theorem 12.2  $\square$

## 12.4 Application of the Perron-Frobenius theorem

Let  $A \in \mathbb{R}^{N \times N}$  be a row stochastic and primitive matrix. We already know that  $\lambda = 1$  is an eigenvalue of  $A$  and  $\lambda = 1 \geq |\mu|$  for all other eigenvalues  $\mu$  (since they lie in the unit disk). Invoking the Perron-Frobenius Theorem, the eigenvalue  $\lambda = 1$  must be simple and strictly dominant, i.e.  $1 = \lambda > |\mu|$  for all other eigenvalues  $\mu \in \text{spec}(A) \setminus \{1\}$ .

Therefore, the averaging system

$$x^{k+1} = Ax^k, \quad x^0 \text{ given}$$

is marginally stable.

### 12.4.1 Equilibrium manifold of linear averaging

Let us investigate the equilibrium manifold of  $x^{k+1} = Ax^k$ . We have

- $x_{eq} = Ax_{eq} \implies x_{eq} \in \ker(I - A) = \text{span } \mathbf{1} = \mathbf{1}\beta$  for some  $\beta \in \mathbb{R}$
- $w^T x^{k+1} = w^T Ax^k = w^T x^k = w^T x^0$  for all  $k \in \mathbb{N}$   
Hence also  $w^T x^0 = w^T x_{eq} = w^T \mathbf{1}\beta \implies \beta = (w^T x^0)/(w^T \mathbf{1})$
- Thus:  $x_{eq} = \mathbf{1} \frac{w^T x^0}{w^T \mathbf{1}}$

*Remark.* The equilibrium is not necessarily the origin and depends on the initial condition  $x^0$

*Remark.* Wlog, we consider a rescaled version of  $w$ , so that  $w^T \mathbf{1} = 1$

## 12.5 Convergence proof (I): Jordan form approach

Being  $G$  strongly connected and aperiodic, then  $A$  is also primitive (other than being row stochastic). Hence,

- it has one dominant eigenvalue in 1 associated to eigenvectors  $\mathbf{1} \in \mathbb{R}^N$  (right) and  $w \in \mathbb{R}^N$  (left)
- it has  $N - 1$  eigenvalues with norm strictly less than 1

Consider the change of coordinates

$$x \mapsto \tilde{x} = T^{-1}x$$

with the nonsingular matrix  $T \in \mathbb{R}^{N \times N}$  being associated to the Jordan form  $T^{-1}AT$  of  $A$

$$T^{-1}AT = \left[ \begin{array}{c|c} 1 & \\ \hline & J_2 \end{array} \right]$$

with  $J_2 \in \mathbb{R}^{(N-1) \times (N-1)}$  Schur (the empty spaces are zeros)

Let matrix  $T$  be arranged as

$$T = \left[ \begin{array}{c|ccc} \uparrow & \uparrow & & \uparrow \\ \mathbf{1} & v^2 & \cdots & v^N \\ \downarrow & \downarrow & & \downarrow \end{array} \right] = [\mathbf{1} \mid W_L] \quad \text{and} \quad T^{-1} = \left[ \begin{array}{c|ccc} \leftarrow & w^T & \rightarrow \\ \hline \leftarrow & w^{2T} & \rightarrow \\ & \vdots & \\ \leftarrow & w^{NT} & \rightarrow \end{array} \right] = \left[ \frac{w^T}{W_R} \right]$$

with  $W_L \in \mathbb{R}^{N \times (N-1)}$  and  $W_R \in \mathbb{R}^{(N-1) \times N}$

The dynamics of  $x^{k+1} = Ax^k$  in the new coordinates is

$$x \mapsto \tilde{x} = T^{-1}x \quad \implies \quad \tilde{x}^{k+1} = T^{-1}AT\tilde{x}^k = \left[ \begin{array}{c|c} 1 & \\ \hline & J_2 \end{array} \right] \tilde{x}^k, \quad \tilde{x}^0 = T^{-1}x^0$$

with  $J_2$  Schur

From the definition of  $T^{-1}$ , for all  $k \in \mathbb{N}$  we have  $\tilde{x}^k = (\tilde{x}_m^k, \tilde{x}_\perp^k)$  with

- $\tilde{x}_m^k = w^T x^k \in \mathbb{R}$  being the weighted average of  $x^k$
- $\tilde{x}_\perp^k = W_R^T x^k \in \mathbb{R}^{N-1}$  being the remaining components (dispersion)

Then

$$\lim_{k \rightarrow \infty} \tilde{x}^k = \lim_{k \rightarrow \infty} \begin{bmatrix} \tilde{x}_m^k \\ \tilde{x}_\perp^k \end{bmatrix} = \lim_{k \rightarrow \infty} \left[ \begin{array}{c|c} 1 & \\ \hline & J_2 \end{array} \right]^k \begin{bmatrix} \tilde{x}_m^0 \\ \tilde{x}_\perp^0 \end{bmatrix} = \lim_{k \rightarrow \infty} \begin{bmatrix} 1^k \tilde{x}_m^0 \\ J_2^k \tilde{x}_\perp^0 \end{bmatrix} = \begin{bmatrix} \tilde{x}_m^0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

where, by definition,  $\tilde{x}_m^0 = w^T x^0$

Therefore, in the original coordinates, it holds

$$\lim_{k \rightarrow \infty} x^k = \lim_{k \rightarrow \infty} T\tilde{x}^k = \begin{bmatrix} \tilde{x}_m^0 \\ \vdots \\ \tilde{x}_m^0 \end{bmatrix} = \mathbf{1} w^T x^0$$

### 12.5.1 The dissensus dynamics

Consider the so-called deflated matrix defined as

$$A - \mathbf{1}w^T \in \mathbb{R}^{N \times N}$$

that, differently from  $A$ , is Schur. In fact, the simple eigenvalue in 1 of  $A$  has been "moved" to 0

Define the *dissensus* vector in  $\mathbb{R}^N$  as

$$\delta^k := x^k - \mathbf{1}w^T x^0 = (I - \mathbf{1}w^T)x^k$$

for all  $k \in \mathbb{N}$

*Remark.* If the dissensus is zero, it means that consensus has been achieved

The dissensus dynamics is

$$\delta^{k+1} = (A - \mathbf{1}w^T)\delta^k, \quad \delta^0 = (I - \mathbf{1}w^T)x^0$$

## 12.6 Convergence proof (II): Lyapunov approach

Consider the Jordan form of the deflated matrix

$$J = T^{-1}(A - \mathbf{1}w^T)T = \left[ \begin{array}{c|c} 0 & \\ \hline & J_2 \end{array} \right] \implies (A - \mathbf{1}w^T) = TJT^{-1}$$

with  $T \in \mathbb{R}^{N \times N}$  nonsingular.

For all  $Q_2 = Q_2^T \succ 0 \in \mathbb{R}^{(N-1) \times (N-1)}$ , there exists a matrix  $P_2 = P_2^T \succ 0 \in \mathbb{R}^{(N-1) \times (N-1)}$  and a scalar  $p_1 > 0$  satisfying the following (discrete-time) Lyapunov equation

$$\left[ \begin{array}{c|c} 0 & \\ \hline & J_2 \end{array} \right]^T \left[ \begin{array}{c|c} p_1 & \\ \hline & P_2 \end{array} \right] \left[ \begin{array}{c|c} 0 & \\ \hline & J_2 \end{array} \right] - \left[ \begin{array}{c|c} p_1 & \\ \hline & P_2 \end{array} \right] = \left[ \begin{array}{c|c} -p_1 & \\ \hline & -Q_2 \end{array} \right]$$

Consider the congruence transformation (preserving the sign of the eigenvalues)

$$P \mapsto \Pi := (T^{-1})^T \left[ \begin{array}{c|c} p_1 & \\ \hline & P_2 \end{array} \right] T^{-1} \implies P = T^T \Pi T$$

and choose the Lyapunov function  $V_\Pi : \mathbb{R}^N \rightarrow \mathbb{R}_{\geq 0}$

$$V_\Pi(\delta) := \delta^T \Pi \delta$$

Then its increment along trajectories satisfies

$$\begin{aligned} V_\Pi(\delta^{k+1}) - V_\Pi(\delta^k) &= (\delta^k)^T (A - \mathbf{1}w^T)^T \Pi (A - \mathbf{1}w^T) \delta^k - (\delta^k)^T \Pi \delta^k \\ &= (\delta^k)^T (TJT^{-1})^T \Pi (TJT^{-1}) \delta^k - (\delta^k)^T \Pi \delta^k \\ &= (T^{-1}\delta^k)^T (J^T T^T \Pi T J - T^T \Pi T) T^{-1} \delta^k \\ &= (T^{-1}\delta^k)^T \left[ \begin{array}{c|c} -p_1 & \\ \hline & -Q_2 \end{array} \right] T^{-1} \delta^k \\ &= -\|T^{-1}\delta^k\|_Q^2 < 0 \end{aligned}$$

## 12.7 Convergence proof (III): alternative Lyapunov approach

Consider the so-called max-min disagreement function  $V_{m-m} : \mathbb{R}^N \rightarrow \mathbb{R}_{\geq 0}$

$$V_{m-m}(x) = \left( \max_{i \in \{1, \dots, N\}} x_i \right) - \left( \min_{i \in \{1, \dots, N\}} x_i \right)$$

- it is a non-quadratic candidate Lyapunov function
- Notice that  $V_{m-m}(x) = 0$  iff  $x = \mathbf{1}\beta$  for some  $\beta \in \mathbb{R}$
- its increment along trajectories decreases monotonically, i.e. for all  $k$

$$V_{m-m}(Ax^k) < V_{m-m}(x^k)$$

Invoking the Lyapunov theorem, we can conclude that  $x^k \xrightarrow[k \rightarrow \infty]{} \mathbf{1}\beta$

Since we have shown that the equilibrium is  $\beta = w^T x^0$ , it follows that  $x^k \xrightarrow[k \rightarrow \infty]{} \mathbf{1}w^T x^0$



## Chapter 13

# Gradient tracking analysis

In the gradient tracking each agent  $i$  performs

$$\begin{aligned} z_i^{k+1} &= \sum_{j \in \mathcal{N}_i} a_{ij} z_j^k - \alpha s_i^k & z_i^0 &\in \mathbb{R} \\ s_i^{k+1} &= \sum_{j \in \mathcal{N}_i} a_{ij} s_j^k + \nabla \ell_i(z_i^{k+1}) - \nabla \ell_i(z_i^k) & s_i^0 &= \nabla \ell_i(z_i^0) \end{aligned}$$

### Aggregate reformulation

Let

$$z^k = \begin{bmatrix} z_1^k \\ \vdots \\ z_N^k \end{bmatrix} \quad s^k = \begin{bmatrix} s_1^k \\ \vdots \\ s_N^k \end{bmatrix} \quad \nabla \ell(z^k) = \begin{bmatrix} \nabla \ell_1(z_1^k) \\ \vdots \\ \nabla \ell_N(z_N^k) \end{bmatrix}$$

Then, the gradient tracking can be compactly written as

$$\begin{aligned} z^{k+1} &= Az^k - \alpha s^k & z^0 &\in \mathbb{R}^n \\ s^{k+1} &= As^k - \nabla \ell(z^{k+1}) - \nabla \ell(z^k) & s^0 &= \nabla \ell(z^0) \end{aligned}$$

where  $A \in \mathbb{R}^{N \times N}$  is the matrix collecting the weights  $a_{ij}$

*Remark.* It is not a causal dynamical system

## 13.1 Equilibrium manifold for the Gradient Tracking

The equilibrium of the system  $(z_{eq}, s_{eq})$  must satisfy

$$\begin{aligned} z_{eq} &= Az_{eq} - \alpha s_{eq} \\ s_{eq} &= As_{eq} + \nabla \ell(z_{eq}) - \nabla \ell(z_{eq}) \end{aligned}$$

Therefore

- from the second:  $s_{eq} \in \ker(I - A) = \text{span } \mathbf{1} \implies s_{eq} = \mathbf{1}\beta_1$  for some  $\beta_1 \in \mathbb{R}$
- from the first:  $(I - A)z_{eq} = -\alpha \mathbf{1}\beta_1$ . But since  $\ker(I - A) = \text{span } \mathbf{1}$  it must be  $\beta_1 = 0$  and

$$\begin{aligned} s_{eq} &= 0 \\ z_{eq} &= \mathbf{1}\beta_2 \text{ for some } \beta_2 \in \mathbb{R} \end{aligned}$$

Moreover, by pre-multiplying by  $\mathbf{1}^T$  the second state equation, for all  $k \in \mathbb{N}$  we have

$$\begin{aligned} \mathbf{1}^T s^{k+1} &= \mathbf{1}^T A s^k + \mathbf{1}^T \nabla \ell(z^{k+1}) - \mathbf{1}^T \nabla \ell(z^k) \\ &= \mathbf{1}^T s^k + \mathbf{1}^T \nabla \ell(z^{k+1}) - \mathbf{1}^T \nabla \ell(z^k) \end{aligned}$$

This implies an invariance condition

$$\begin{aligned}
 \mathbf{1}^T s^{k+1} - \mathbf{1}^T \nabla \ell(z^{k+1}) &= s^k - \mathbf{1}^T \nabla \ell(z^k) \\
 &= \mathbf{1}^T s_{eq} - \mathbf{1}^T \nabla \ell(z^{eq}) = 0 - \mathbf{1}^T \nabla \ell(z_{eq}) = -\mathbf{1}^T \nabla \ell(\mathbf{1}\beta_2) \\
 &= \mathbf{1}^T s_0 - \mathbf{1}^T \nabla \ell(z^0) = 0
 \end{aligned}$$

Therefore it must be  $\beta_2 = z^*$ , so that  $z_{eq} = \mathbf{1}z^*$

## 13.2 Causal coordinates for the Gradient Tracking

Consider the nonlinear change of coordinates

$$\begin{bmatrix} z^k \\ s^k \end{bmatrix} \mapsto \begin{bmatrix} z^k \\ \xi^k \end{bmatrix} = \begin{bmatrix} z^k \\ \alpha(\nabla \ell(z^k) - s^k) \end{bmatrix}$$

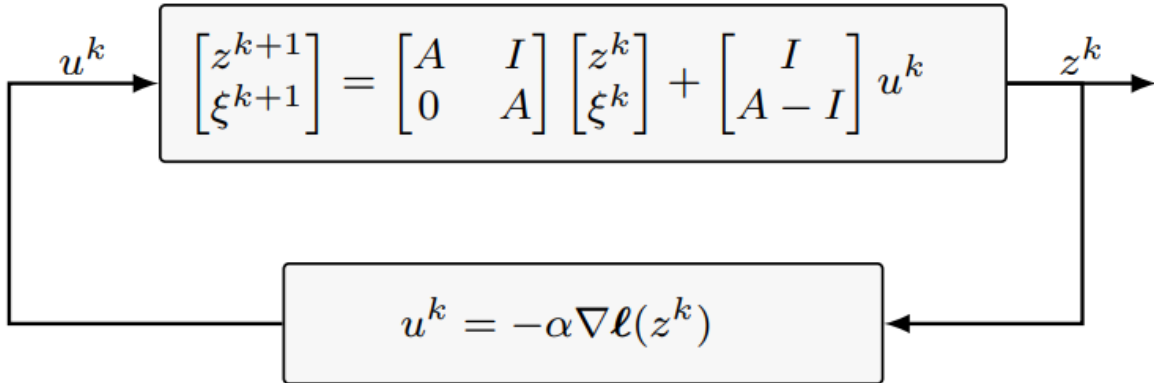
*Remark.* The inverse transformation is  $s^k = -\frac{1}{\alpha}\xi^k + \nabla \ell(z^k)$

The GT dynamics in the new coordinates is

$$\begin{aligned}
 z^{k+1} &= Az^k + \xi^k - \alpha \nabla \ell(z^k) & z^0 &\in \mathbb{R}^N \\
 \xi^{k+1} &= A\xi^k - \alpha(A - I)\nabla \ell(z^k) & \xi^0 &= 0
 \end{aligned}$$

*Remark.* This is a causal system

### 13.2.1 Block diagram of the Gradient Tracking Algorithm



- The upper part is a linear dynamical system with states  $(z, \xi)$  and input  $u$  (and output  $y = z$ )
- The state matrix  $\begin{bmatrix} A & I \\ 0 & A \end{bmatrix}$  is Schur except for the semi-simple eigenvalue in 1 (geom. mult. 2)
- The step-size  $\alpha > 0$  plays the role of a gain
- the lower part is a static nonlinearity (as in the gradient method)

### 13.3 Local perspective of the causal GT

From a local perspective, each agent  $i$  performs

$$\begin{aligned} z_i^{k+1} &= \sum_{j \in \mathcal{N}_i} a_{ij} z_j^k + \xi_i^k - \alpha \nabla \ell_i(z_i^k) & z_i^0 &\in \mathbb{R}^N \\ \xi_i^{k+1} &= \sum_{j \in \mathcal{N}_i} a_{ij} \xi_j^k - \alpha \left( \sum_{j \in \mathcal{N}_i} a_{ij} \nabla \ell_j(z_j^k) - \nabla \ell_i(z_i^k) \right) & \xi_i^0 &= 0 \end{aligned}$$

*Remark.* more communication is needed in these coordinates

Since each cost function  $\ell_i$  is strongly convex (with parameter  $\mu > 0$ ) with Lipschitz continuous gradient (with parameter  $L > 0$ ) then for all  $z_A, z_B \in \mathbb{R}^N$  it holds

$$\begin{aligned} (\nabla \ell(z_A) - \nabla \ell(z_B))^T (z_A - z_B) &= \sum_{i=1}^N (\nabla \ell_i(z_{i,A}) - \nabla \ell_i(z_{i,B}))^T (z_{i,A} - z_{i,B}) \\ &\geq \frac{\mu L}{L + \mu} \|z_A - z_B\|^2 + \frac{1}{L + \mu} \|\nabla \ell(z_A) - \nabla \ell(z_B)\|^2 \end{aligned}$$

### 13.4 Error coordinates for the Gradient Tracking

Consider the error coordinates

$$\begin{bmatrix} z^k \\ \xi^k \end{bmatrix} \mapsto \begin{bmatrix} \tilde{z}^k \\ \tilde{\xi}^k \end{bmatrix} = \begin{bmatrix} z^k - \mathbf{1} \\ \xi^k - \alpha \nabla \ell(\mathbf{1} z^*) \end{bmatrix}$$

The dynamics in the new coordinates is

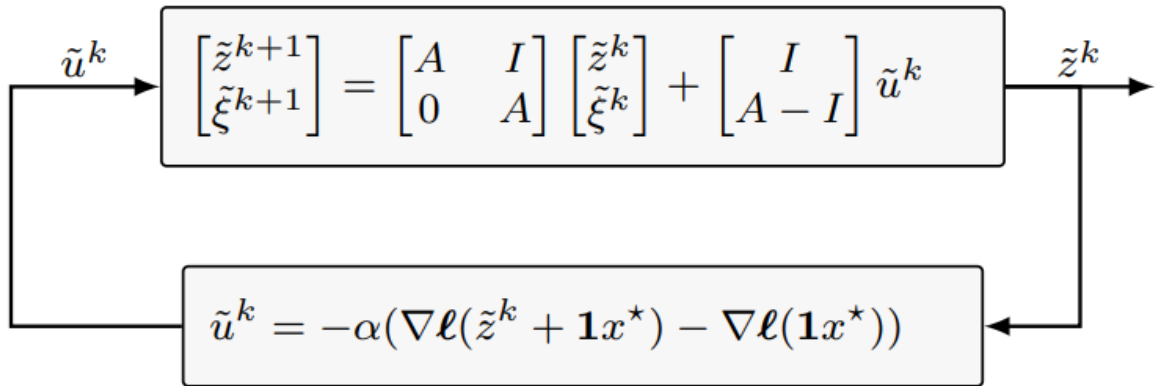
$$\begin{aligned} \tilde{z}^{k+1} &= A \tilde{z}^k + \tilde{\xi}^k - \alpha (\nabla \ell(z^k) - \nabla \ell(\mathbf{1} z^*)) & \tilde{z}^0 &\text{ given} \\ \tilde{\xi}^{k+1} &= A \tilde{\xi}^k - \alpha (A_I) (\nabla \ell(z^k) - \nabla \ell(\mathbf{1} z^*)) & \tilde{\xi}^0 &\text{ given} \end{aligned}$$

that can be written as

$$\begin{bmatrix} \tilde{z}^{k+1} \\ \tilde{\xi}^{k+1} \end{bmatrix} = \begin{bmatrix} A & I \\ 0 & A \end{bmatrix} \begin{bmatrix} \tilde{z}^k \\ \tilde{\xi}^k \end{bmatrix} - \alpha \begin{bmatrix} I \\ A_I \end{bmatrix} (\nabla \ell(\tilde{z}^k + \mathbf{1} z^*) - \nabla \ell(\mathbf{1} z^*))$$

*Remark.* The equilibrium has been shifted to the origin  $(\tilde{z}, \tilde{\xi}) = (0, 0)$

#### 13.4.1 Block diagram of the Gradient Tracking in error coordinates



- The origin  $(\tilde{z}, \tilde{\xi}) = (0, 0)$  is the equilibrium

- The feedback action  $\tilde{u}^k$  is a static map whose characterization depends on the regularity properties of the cost function, i.e.  $(\mu, L)$ , and is given by

$$(\nabla \ell(\tilde{z}^k + \mathbf{1}z^*) - \nabla \ell(\mathbf{1}z^*))^T \tilde{z}^k \geq \frac{\mu L}{L + \mu} \|\tilde{z}^k\|^2 + \frac{1}{L + \mu} \|\nabla \ell(\tilde{z}^k + \mathbf{1}z^*) - \nabla \ell(\mathbf{1}z^*)\|^2$$

- The feedback must be a stabilizing action for the closed-loop system

### 13.5 Analysis sketch

Let us focus on the upper (open-loop) linear system

$$\begin{bmatrix} \tilde{z}^{k+1} \\ \tilde{\xi}^{k+1} \end{bmatrix} = \begin{bmatrix} A & I \\ 0 & A \end{bmatrix} \begin{bmatrix} \tilde{z}^k \\ \tilde{\xi}^k \end{bmatrix} + \begin{bmatrix} I \\ A_I \end{bmatrix} \tilde{u}^k$$

Let us consider a change of coordinates (with  $T_z = \begin{bmatrix} \mathbf{1} & W \end{bmatrix}$  and  $T_\xi = \begin{bmatrix} W & \mathbf{1} \end{bmatrix}$ , orthonormal) given by

$$\begin{bmatrix} \tilde{z}^k \\ \tilde{\xi}^k \end{bmatrix} \mapsto \begin{bmatrix} \tilde{z}_m^k \\ \tilde{z}_\perp^k \\ \tilde{\xi}_\perp^k \\ \tilde{\xi}_m^k \end{bmatrix} = \left[ \begin{array}{c|c} T_z^T & \\ \hline & T_\xi^T \end{array} \right] \begin{bmatrix} \tilde{z}^k \\ \tilde{\xi}^k \end{bmatrix}$$

The dynamics in the new coordinates reads

$$\begin{bmatrix} \tilde{z}_m^{k+1} \\ \tilde{z}_\perp^{k+1} \\ \tilde{\xi}_\perp^{k+1} \\ \tilde{\xi}_m^{k+1} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & J_2 & I_{N-1} & 0 \\ 0 & 0 & J_2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \tilde{z}_m^k \\ \tilde{z}_\perp^k \\ \tilde{\xi}_\perp^k \\ \tilde{\xi}_m^k \end{bmatrix} + \begin{bmatrix} \mathbf{1}^T \\ W^T \\ W^T(A - I) \\ 0 \end{bmatrix} \tilde{u}^k$$

where  $J_2 := W^T A W$  is Schur

It can be shown that

1. the matrix  $F_{22} = \begin{bmatrix} J_2 & I_{N-1} \\ 0 & J_2 \end{bmatrix}$  is Schur
2.  $\xi_i^0 = 0$  for all  $i$  implies that  $\tilde{x}_m^0 = 0$

Therefore the evolution of  $\tilde{\xi}_m^k$  can be ignored to obtain a "reduced" system

$$\begin{bmatrix} \tilde{z}_m^{k+1} \\ \tilde{z}_\perp^{k+1} \\ \tilde{\xi}_\perp^{k+1} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & F_{22} \end{bmatrix} \begin{bmatrix} \tilde{z}_m^k \\ \tilde{z}_\perp^k \\ \tilde{\xi}_\perp^k \end{bmatrix} + \begin{bmatrix} \mathbf{1}^T \\ G_2 \end{bmatrix} \tilde{u}^k$$

with  $G_2 = \begin{bmatrix} W^T \\ W^T(A - I) \end{bmatrix}$ , where  $\tilde{u}^k = -\alpha(\nabla \ell(\mathbf{1}\tilde{z}_m^k + W\tilde{z}_\perp^k + \mathbf{1}z^*) - \nabla \ell(\mathbf{1}z^*))$  (coupling the dynamics)

Consider the quadratic Lyapunov function

$$\mathcal{V}(\tilde{z}_m^k, \tilde{z}_\perp^k, \tilde{\xi}_\perp^k) = \begin{bmatrix} \tilde{z}_m^k \\ \tilde{z}_\perp^k \\ \tilde{\xi}_\perp^k \end{bmatrix}^T \begin{bmatrix} 1 & 0 \\ 0 & P_2 \end{bmatrix} \begin{bmatrix} \tilde{z}_m^k \\ \tilde{z}_\perp^k \\ \tilde{\xi}_\perp^k \end{bmatrix}$$

with  $P_2 = P_2^T \succ 0$  designed on the matrix  $F_{22}$  (i.e., via  $F_{22}^T P_2 F_{22} - P_2 \prec 0$ )

There exists  $\alpha^* > 0$  such that for all  $\alpha \in (0, \alpha^*)$  it holds

$$\begin{aligned} \mathcal{V}(\tilde{z}_m^{k+1}, \tilde{z}_\perp^{k+1}, \tilde{\xi}_\perp^{k+1}) - \mathcal{V}(\tilde{z}_m^k, \tilde{z}_\perp^k, \tilde{\xi}_\perp^k) &= \begin{bmatrix} \tilde{z}_m^k \\ \tilde{z}_\perp^k \\ \tilde{\xi}_\perp^k \end{bmatrix}^T (F^T P F - P) \begin{bmatrix} \tilde{z}_m^k \\ \tilde{z}_\perp^k \\ \tilde{\xi}_\perp^k \end{bmatrix} + 2 \begin{bmatrix} \tilde{z}_m^k \\ \tilde{z}_\perp^k \\ \tilde{\xi}_\perp^k \end{bmatrix}^T F^T P G \tilde{u}^k + (\tilde{u}^k)^T G^T P G \tilde{u}^k \\ &\leq -q_1 \|\tilde{z}_m^k\|^2 - q_2 \|\tilde{z}_\perp^k\|^2 - q_3 \|\tilde{\xi}_\perp^k\|^2 \end{aligned}$$

for some  $q_1, q_2, q_3 > 0$  (depending on  $\alpha$ )

It follows that the origin is exponentially stable, which translates in the sought linear convergence of the solution estimates generated by the Gradient Tracking



## 13.6 Gradient Tracking: extensions

Gradient Tracking can be extended to handle

- Nonconvex cost functions
- Convex constraints on the optimization variable
- Time-varying and directed communication networks (packet losses)
- Block-wise updates and communication
- Stochastic optimization problems
- Heavy-ball, Nesterov, Adam-like implementations

## 13.7 System theoretical approach to distributed optimization

- Gradient tracking can be shown to implement a *distributed integral action*
- The approach can be generalized to derive more general controllers and/or algorithms

For instance, consider the following algorithm

$$\begin{aligned} z^{k+1} &= Az^k - (I - A)\xi^k - \alpha \nabla \ell(z^k) & z^0 &\in \mathbb{R}^N \\ \xi^{k+1} &= \xi^k + (I - A)z^k & \xi^0 &\in \mathbb{R}^N \end{aligned}$$

*Remark.* The dynamics of the controller state  $\xi^k$  is a pure discrete-time integrator trying to compensate for the consensus error  $(I - A)z^k$



## Chapter 14

# Constraint-coupled Distributed Optimization

Constraint-coupled optimization problem

$$\begin{aligned} \min_{z_1, \dots, z_N} \quad & \sum_{i=1}^N \ell_i(z_i) \\ \text{subj. to} \quad & z_i \in Z_i \subset \mathbb{R}^{n_i} \quad i = 1, \dots, N \\ & \sum_{i=1}^N g_{ij}(z_i) \leq 0 \quad j = 1, \dots, r \end{aligned}$$

with  $\ell_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R}$  convex,  $g_{ij} : \mathbb{R}^{n_i} \rightarrow \mathbb{R}$  convex and  $Z_i \subset \mathbb{R}^{n_i}$  convex and compact.

Constraint-coupled optimization problem (vector form)

$$\begin{aligned} \min_{z_1, \dots, z_N} \quad & \sum_{i=1}^N \ell_i(z_i) \\ \text{subj. to} \quad & z_i \in Z_i \subset \mathbb{R}^{n_i} \quad i = 1, \dots, N \\ & \sum_{i=1}^N g_i(z_i) \leq 0 \end{aligned}$$

with  $\ell_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R}$  convex,  $g_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^r$  convex and  $Z_i \subset \mathbb{R}^{n_i}$  convex and compact

### 14.1 Constraint-coupled Setup

- $N$  agents communicate over graph  $G$
- $(z_1, \dots, z_N)$  decision variables - size grows with  $N$
- agent  $i$  knows  $\ell_i, g_i$  and  $Z_i$

### 14.2 Multi-robot scenario

Let us assume a framework with  $N$  robots that must collectively serve  $N$  tasks. We model the state of the  $i$ -th robot as  $z_i \in \mathbb{R}^{n_i}$  where  $n_i$  is the number of tasks robot  $i$  can serve, and we denote

$$\begin{cases} z_{ik} = 1 & \text{if robot } i \text{ serves task } k \\ z_{ik} = 0 & \text{otherwise} \end{cases}$$

Note that each robot serves only one task, thus imposing a constraint of the form

$$\sum_{k=1}^N z_{ik} = 1$$

The constraint set is thus

$$Z_i = \{z_i \in \{0, 1\}^N \mid \sum_{k=1}^N z_{ik} = 1\}$$

and we say that robot  $i$  serves task  $j$  with some cost  $c_{ij}$ . Suppose each robot can serve all tasks. Then, to ensure is task is carried out by only one robot,

$$\sum_{i=1}^N z_{ik} = 1 \quad k = 1, \dots, N$$

in a more compact way

$$\sum_{i=1}^N z_i = 1$$

if robots can serve only a subset of of tasks we can define some "selection" matrices  $H_i$ ,  $i \in \{1, \dots, N\}$ . We can thus write the constraint as

$$\sum_{i=1}^N H_i z_i = 1$$

To summarize, we can rewrite the problem as

$$\begin{aligned} \min_{z_1, \dots, z_N} \quad & \sum_{i=1}^N c_i^T z_i \\ \text{subj. to } \quad & z_i \in Z_i \quad i = 1, \dots, N \\ & \sum_{i=1}^N H_i z_i = 1 \end{aligned}$$

*Remark.* Task allocation problems satisfy the "unimodularity" property. We can relax to a linear program (in place of an integer program):

$$Z_i = \{z_i \in \mathbb{R}^{n_i} \mid 0 \leq z_i \leq 1, \sum_{k=1}^N z_{ik} = 1\}$$

The solution  $z_i^* \in \{0, 1\}$

### 14.3 Distributed optimal control

$$\begin{aligned} \min_{\substack{x_1, \dots, x_N \\ u_1, \dots, u_N}} \quad & \sum_{i=1}^N \left( \sum_{\tau=0}^{T-1} \ell_i^{SC}(x_i(\tau), u_i(\tau)) + \ell_i^{TC}(x_i(T)) \right) \\ \text{subj. to } \quad & x_i(\tau+1) = A_i x_i(\tau) + B_i u_i(\tau), \quad \forall i, \quad \tau \in [0, T-1] \\ & x_i(\tau) \in X_i \quad u_i(\tau-1) \in U_i, \quad \forall i, \quad \tau \in [1, T] \\ & \sum_{i=1}^N H_i x_i(\tau) \leq \bar{h} \quad \tau \in [1, T] \end{aligned}$$

## 14.4 Duality for constraint-coupled optimization

Let the Lagrangian function  $\mathcal{L} : \mathbb{R}^{\sum_{i=1}^N n_i} \times \mathbb{R}^r \rightarrow \mathbb{R}$  be

$$\mathcal{L}(x, \mu) = \sum_{i=1}^N (\ell_i(z_i) + \mu^T g_i(z_i))$$

Then, the dual function can be computed as

$$\begin{aligned} q(\mu) &= \inf_{\substack{z_i \in Z_i \\ i \in \{1, \dots, N\}}} \sum_{i=1}^N (\ell_i(z_i) + \mu^T g_i(z_i)) \\ &= \sum_{i=1}^N \inf_{z_i \in Z_i} (\ell_i(z_i) + \mu^T g_i(z_i)) \\ &= \sum_{i=1}^N q_i(\mu) \end{aligned}$$

with  $q_i(\mu) := \inf_{z_i \in Z_i} (\ell_i(z_i) + \mu^T g_i(z_i))$ . Because each local cost function depends only on its own variable and the constraint sets are independent we can minimize each cost function separately.

The dual problem is thus

$$\max_{\mu \geq 0} \sum_{i=1}^N q_i(\mu)$$

*Remark.* The dual problem is a cost-coupled optimization problem.

## 14.5 Centralized dual subgradient

The subgradient algorithm non the dual problem reads

$$\begin{aligned} \mu^{k+1} &= \mathcal{P}_{\mu \geq 0} \left( \mu^k + \alpha^k \tilde{\nabla} q(\mu^k) \right) \\ &= \mathcal{P}_{\mu \geq 0} \left( \mu^k + \alpha^k \sum_{i=1}^N \tilde{\nabla} q_i(\mu^k) \right) \end{aligned}$$

where  $\alpha^k$  is the stepsize.

A subgradient of  $q_i$  at  $\mu^k$  can be computed as

$$\begin{aligned} z_i^{k+1} &= \arg \min_{z_i \in Z_i} \ell_i(z_i) + (\mu^k)^T g_i(z_i) \\ \tilde{\nabla} q_i(\mu^k) &= g_i(z_i^{k+1}) \end{aligned}$$

### Centralized dual subgradient algorithm

$$\begin{aligned} z_i^{k+1} &= \arg \min_{z_i \in Z_i} \ell_i(z_i) + (\mu^k)^T g_i(z_i) \\ \mu^{k+1} &= \mathcal{P}_{\mu \geq 0} \left( \mu^k + \alpha^k \sum_{i=1}^N g_i(z_i^{k+1}) \right) \end{aligned}$$

*Remark.* The centralized part is the update of  $\mu$ , which is to be computed by a "master" node, whereas the update of  $z_i$  can be done locally to each agent.

*Remark.* This procedure is known as dual decomposition, as the problem is decomposed and computation can be parallelized

**Distributed dual subgradient algorithm**

$$\begin{aligned}
v_i^{k+1} &= \sum_{j=1}^N a_{ij} \mu_j^k \\
z_i^{k+1} &= \arg \min_{z_i \in Z_i} \ell_i(z_i) + (v_i^{k+1})^T g_i(z_i) \\
\mu_i^{k+1} &= \mathcal{P}_{z \geq 0} (\mu_i^k + \alpha^k g_i(z_i^{k+1}))
\end{aligned}$$

This algorithm was obtained by applying the distributed gradient algorithm to the centralized algorithm, thus obtaining a completely distributed algorithm.

If the minimum in the second step can be proven to be unique at every iteration, then  $z_i^k$  converges to the optimal solution of the primal problem. Otherwise, it can be proven that the running average of  $z_i^k$  converges to the optimal solution.

**14.5.1 Distributed dual subgradient convergence**

- Assumption 1. Let  $a_{ij}$ ,  $i, j \in \{1, \dots, N\}$  be non-negative entries of a weighted adjacency matrix  $A$  associated to the undirected graph  $G$ , with  $a_{ii} > 0$  and  $A$  doubly stochastic.
- Assumption 2. The step-size sequence  $\{\alpha^k\}_{k \geq 0}$  satisfies the conditions

$$\alpha^k \geq 0 \quad \sum_{k=0}^{\infty} \alpha^k = \infty \quad \sum_{k=0}^{\infty} (\alpha^k)^2 < \infty$$

- Assumption 3d. For all  $i \in \{1, \dots, N\}$ : each function  $\ell_i$  is convex, each constraint  $Z_i$  is a non-empty, compact and convex set; each function  $g_i$  is a component-wise convex function. Moreover, there exist  $\bar{z}_i \in Z_1, \dots, \bar{z}_N \in Z_N$  such that  $\sum_{i=1}^N g_i(\bar{z}_i) < 0$

**Theorem 14.1** (Distributed Dual Subgradient)

Let Assumptions 1, 2, and 3d hold. Then, the sequence of dual variables  $\{\mu_1^k, \dots, \mu_N^k\}_{k \geq 0}$  generated by the Distributed Dual Subgradient satisfies

$$\lim_{k \rightarrow \infty} \|\mu_i^k - \mu^*\| = 0, \quad i \in \{1, \dots, N\}$$

where  $\mu^*$  is an optimal solution of the dual problem. Moreover, let for each  $k$

$$\hat{z}_i^k = \frac{1}{k} \sum_{\tau=0}^k z_i^\tau$$

Then, it holds

$$\begin{aligned}
\lim_{k \rightarrow \infty} \sum_{i=1}^N \ell_i(\hat{z}_i^k) &= \ell^* \\
\lim_{k \rightarrow \infty} \|\hat{z}_i^k - z_i^*\| &= 0 \quad i \in \{1, \dots, N\}
\end{aligned}$$

where  $z_i^*$  and  $\ell^*$  denote an optimal solution and the optimal cost of the primal problem