

# Multi-Agent Optimization and Learning

## Lecture 3: Graphs and their Role in Decentralized Processing

Stefan Vlaski<sup>†</sup> and Ali H. Sayed<sup>\*</sup>

<sup>†</sup>Department of Electrical and Electronic Engineering, Imperial College London, UK  
<sup>\*</sup>Adaptive Systems Laboratory, École Polytechnique Fédérale de Lausanne, Switzerland

IEEE ICASSP 2024 Short Course

**Imperial College  
London**

**EPFL**

## Objective of this lecture

- In Lecture 2 we presented a number of fusion-center based algorithms for multi-agent learning optimization and learning.
- We also developed performance bounds, and investigated the dependence of these bounds on different factors including:
  - ▶ The quality of local gradient approximations.
  - ▶ The level of heterogeneity in the system.
  - ▶ Other algorithm parameters such as participation rate and number of local updates.
- For the remainder of the course, we will develop fully decentralized architectures, where interactions occur over a graph.
- In this lecture, we will introduce graphs and their properties, and illustrate some key dynamics in a simplified setting, namely the problem of computing averages over a graph.
- In the next lecture, we will generalize to optimization and learning over graphs.

# Graphs

An directed, unweighted graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$  consists of a pair of sets  $\mathcal{N}$  and  $\mathcal{E}$ , where  $\mathcal{N}$  denotes a set of *vertices* or *nodes*, and  $\mathcal{E}$  denotes a set of directed *edges* linking pairs of vertices in  $\mathcal{N}$ . Each edge in  $\mathcal{E}$  is represented as a pair (i.e., a 2-tuple) of vertices in  $\mathcal{N}$ . In principle, the elements of  $\mathcal{N}$  can be arbitrary objects. For simplicity, it is generally sufficient number the vertices in  $\mathcal{N}$  beginning at one through  $|\mathcal{N}|$  and represent the elements through their respective index. In most lectures, vertices will be “agents”, or “learners”, terms which we use interchangeably in our presentation. We generally denote the number of agents  $|\mathcal{N}|$  by  $K$ .

## Example

We consider a collection of 4 nodes, numbered 1 through 4. The nodes are collected in the set:

$$\mathcal{N} = \{1, 2, 3, 4\} \quad (1)$$

We construct a sample graph using the following set of edges:

$$\mathcal{E} = \{(1, 1), (1, 2), (2, 1), (2, 3), (3, 2), (3, 4), (4, 3), (2, 4), (4, 2)\} \quad (2)$$

The resulting graph  $\mathcal{G} \triangleq (\mathcal{N}, \mathcal{E})$  is displayed below.

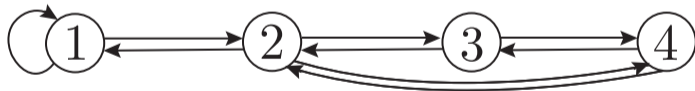


Figure: A simple graph consisting of four vertices.

# Adjacency Matrix

Instead of describing a graph via its set of nodes  $\mathcal{N}$  and vertices  $\mathcal{E}$ , we can equivalently describe it through its **adjacency matrix**  $C \in \{0, 1\}^{K \times K}$ . Here, the element  $c_{\ell k} = [C]_{\ell, k}$  is set to one when there is an edge from node  $\ell$  to node  $k$ , and to zero otherwise. Formally:

$$c_{\ell k} = 1 \iff (\ell, k) \in \mathcal{E} \quad (3)$$

$$c_{\ell k} = 0 \iff (\ell, k) \notin \mathcal{E} \quad (4)$$

In this manner, we can find for the graph from the earlier example:

$$C = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix} \quad (5)$$

## Weighted Adjacency Matrix

In many situations, it will be useful to associate a weight with edges in a graph to quantify the *strength* of the link between a pair of nodes  $\ell$  and  $k$ . We can readily capture these weights by allowing  $c_{\ell k}$  to take non-negative real values, resulting in an adjacency matrix  $C \in \mathbb{R}^{K \times K}$ . We allow only *non-zero* weights to be associated with an edge to avoid degenerate situations where nodes are connected by an edge, but the associated weight is zero. In this manner, it follows from  $c_{\ell k} > 0$  that there is an edge from agent  $\ell$  to agent  $k$ , while  $c_{\ell k} = 0$  implies that there is no edge from  $\ell$  to  $k$ .

$$C = \begin{pmatrix} 1 & 2 & 0 & 0 \\ 0.5 & 0 & 4 & 2 \\ 0 & 2 & 0 & 1 \\ 0 & 0.2 & 1 & 0 \end{pmatrix} \quad (6)$$

## Neighborhoods

For a node  $k$ , we define the **in-neighborhood** as the set of nodes  $\mathcal{N}_k^{\text{in}}$ , for which an edge *from*  $\ell$  *to*  $k$  exists. Formally:

$$\mathcal{N}_k^{\text{in}} = \{\ell : (\ell, k) \in \mathcal{E}\} = \{\ell : c_{\ell k} > 0\} \quad (7)$$

Alternatively, we can define the **out-neighborhood**  $\mathcal{N}_k^{\text{out}}$  as the set of nodes, for which an edge *to*  $\ell$  *from*  $k$  exists:

$$\mathcal{N}_k^{\text{out}} = \{\ell : (k, \ell) \in \mathcal{E}\} = \{\ell : c_{k\ell} > 0\} \quad (8)$$

For graphs where an edge  $c_{\ell k}$  describes the influence of agent  $\ell$  on agent  $k$ , we can interpret  $\mathcal{N}_k^{\text{in}}$  as the set of agents that  $k$  is influenced by, while  $\mathcal{N}_k^{\text{out}}$  denotes the set of nodes that are influenced by  $k$ .

An edge that connects an agent to itself is called a **self-loop**. This means that an agent  $k$  can be contained in its own neighborhood, provided that it has a self-loop.

# Undirected Graphs

We say that a graph is **undirected**, if the presence of an edge from agent  $k$  agent to  $\ell$ , implies that there is an edge from agent  $\ell$  to agent  $k$ . In other words:

$$k \in \mathcal{N}_\ell^{\text{in}} \iff \ell \in \mathcal{N}_k^{\text{in}} \quad (9)$$

It can be readily verified that this is equivalent to requiring the in-neighborhoods and out-neighborhoods of each node  $k$  to coincide:

$$\mathcal{N}_k^{\text{in}} = \mathcal{N}_k^{\text{out}} \triangleq \mathcal{N}_k \quad (10)$$

This means that for undirected graphs, there is no need to distinguish between its in-degrees and out-degrees, and we can just refer to the neighborhood  $\mathcal{N}_k$ .

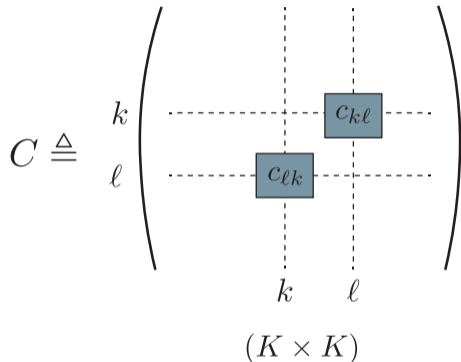
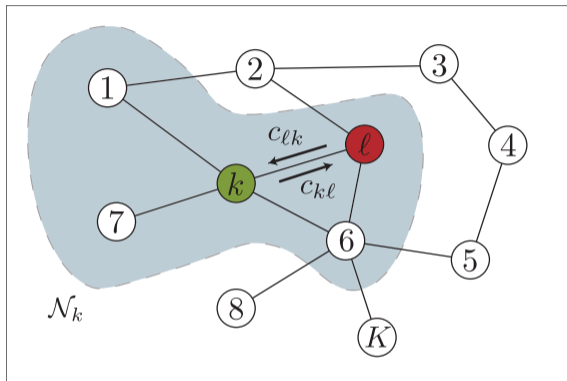


# Symmetric Graphs

We say that a graph is **symmetric**, if the weights linking pairs of agents are symmetric, i.e., when the adjacency matrix is symmetric.

$$C = C^T \quad (11)$$

# Adjacency Matrix



**Figure:** An undirected graph with the neighborhood  $\mathcal{N}_k$  of node highlighted. On the right, we see the corresponding adjacency matrix.

# Strongly-Connected Graphs

We say that a graph is strongly-connected if:

- It is connected.
- There exists at least one self-loop, i.e.,  $c_{kk} > 0$  for at least one  $k$ .

Strong-connectivity of a graph implies that its adjacency matrix is **primitive**, which means that there exists some finite integer  $n_o > 0$  such that all entries of  $C^{n_o}$  are strictly positive:

$$[C^{n_o}]_{\ell,k} > 0, \quad \text{uniformly for all } (\ell, k) \quad (12)$$

## Perron-Frobenius Theorem

One important consequence of the primitiveness of  $C$  is that a famous result in matrix theory, known as the *Perron-Frobenius* theorem, characterizes the eigen-structure of  $C$ . In particular, it will hold that:

- (a) The matrix  $C$  has a *single* leading eigenvalue  $\lambda_1 \in \mathbb{R}$ , which is positive and real.
- (b) The magnitudes of all other (potentially complex) eigenvalues of  $C$  are *strictly less* than  $\lambda_1$ , so that  $\rho(A) = \lambda_1$ .
- (c) The leading eigenvalue  $\lambda_1$  has algebraic multiplicity one, which implies that there is a single eigenvector  $p$  associated with  $\lambda_1$ . Furthermore,  $p$  can be normalized so that all of its entries are strictly positive and add up to one, i.e.,

$$Cp = \lambda_1 p, \quad \mathbf{1}^T p = 1, \quad p_k > 0, \quad k = 1, 2, \dots, K \quad (13)$$

We refer to  $p$  as the *Perron vector* of  $C$ .

## Jordan Decomposition of Primitive Matrices

Owing to the multiplicity of the leading eigenvalue, the Jordan decomposition  $C = V_\epsilon J V_\epsilon^{-1}$  has a very particular structure:

$$V_\epsilon = [ p \quad V_R ], \quad J = \begin{bmatrix} \lambda_1 & 0 \\ 0 & J_\epsilon \end{bmatrix}, \quad V_\epsilon^{-1} = \begin{bmatrix} q^\top \\ V_L^\top \end{bmatrix} \quad (14)$$

where  $J_\epsilon$  is a block Jordan matrix with the eigenvalues  $\lambda_2(C)$  through  $\lambda_K(C)$  on the diagonal and  $\epsilon$  on the first lower sub-diagonal. In particular,  $\rho(J_\epsilon) < \lambda_1$ . One consequence of this decomposition is that:

$$\rho(A - pq^\top) < \lambda_1 \quad (15)$$

and

$$\lim_{i \rightarrow \infty} \frac{C^i}{\lambda_1^i} = pq^\top \quad (16)$$

## Graph Laplacian

An alternative description of a graph is sometimes given through its **graph Laplacian**:

$$L \triangleq \text{diag} \{C\mathbf{1}\} - C \quad (17)$$

When  $C$  is symmetric, several useful facts about a graph can be extracted from its Laplacian. Let  $\theta_1 \geq \theta_2 \geq \dots \geq \theta_K$  denote the ordered eigenvalues of  $L$ . Then:

- (a)  $L$  is symmetric nonnegative-definite so that  $\theta_k \geq 0$ , for  $k = 1, 2, \dots, K$ .
- (b) The entries on each row of  $L$  add up to zero so that  $L\mathbf{1} = 0$ . This means that  $\mathbf{1}$  is a right eigenvector of  $L$  corresponding to the eigenvalue at zero.
- (c) The smallest eigenvalue is always zero, i.e.,  $\theta_K = 0$ . The second smallest eigenvalue,  $\theta_{K-1}$ , is called the algebraic connectivity of the graph.
- (d) The number of times that zero is an eigenvalue of  $L$  (i.e., its multiplicity) is equal to the number of connected subgraphs.
- (e) The algebraic connectivity of a connected graph is nonzero, i.e.,  $\theta_{K-1} \neq 0$ . In other words, a graph is connected if, and only if, its algebraic connectivity is nonzero.

## Example

If we consider again the graph from our earlier example, we find:

$$\begin{aligned} L &= \text{diag} \{C\mathbf{1}\} - C \\ &= \text{diag} \left\{ \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \right\} - \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix} \\ &= \text{diag} \left\{ \begin{pmatrix} 2 \\ 3 \\ 2 \\ 2 \end{pmatrix} \right\} - \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix} \\ &= \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 2 \end{pmatrix} - \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & -1 & 0 & 0 \\ -1 & 3 & -1 & -1 \\ 0 & -1 & 2 & -1 \\ 0 & -1 & -1 & 2 \end{pmatrix} \end{aligned} \quad (18)$$

# Interpretation as a Differential Operator

The Laplacian carries a useful interpretation because it allows us to compute efficiently the variation of a signal over a graph. Specifically, for a symmetric  $C = C^T$  and a signal  $x \in \mathbb{R}^K$ , it holds that:

$$\frac{1}{2} \sum_{k=1}^K \sum_{\ell=1}^K c_{\ell k} \|x_k - x_{\ell}\|^2 = x^T L x \quad (19)$$



## Combination/Weight Matrices

Adjacency matrices will appear frequently in future lectures as a means for describing weighting assigned by agents to information received from their neighbors. In this context, we will impose additional conditions on the entries of the adjacency matrix, and we will refer to matrices satisfying these conditions as *combination matrices* or *weight matrices*. Combination weights  $\{a_{\ell k}\}$  will be nonnegative scalars that satisfy the following condition for each agent  $k = 1, 2, \dots, K$ :

$$a_{\ell k} \geq 0, \quad \sum_{\ell=1}^K a_{\ell k} = 1, \quad \text{and} \quad a_{\ell k} = 0 \text{ if, and only if, } \ell \notin \mathcal{N}_k \quad (20)$$

Condition (20) implies that the combination matrix  $A = [a_{\ell k}]$  is *left-stochastic*, which means that:

$$\mathbf{1}^\top A = \mathbf{1}^\top \iff A^\top \mathbf{1} = \mathbf{1}, \quad (\text{left-stochastic}) \quad (21)$$

## Spectral Properties of Combination Matrices

Since the spectral radius of a matrix is bounded by any of its norms, it follows that:

$$\rho(A) \leq \|A\|_1 = 1 \quad (22)$$

From (21) we observe that  $A$  has an eigenvalue at one and hence  $\rho(A) = 1$ . We also recognize that  $A$  is an adjacency matrix. Hence, if the underlying graph is strongly-connected, it follows that  $A$  is primitive, and we can conclude from Perron-Frobenius, that the eigenvalue at one has multiplicity one with corresponding Perron vector  $p$ :

$$Ap = p \quad (23)$$

If  $A$  is additionally symmetric, it follows that:

$$A^T \mathbf{1} = \mathbf{1}, \quad A = A^T, \quad (\text{doubly-stochastic}) \quad (24)$$

where  $\mathbf{1}$  denotes the vector with all entries equal to one. We say that  $A$  is *doubly-stochastic* since the entries on each of its rows and on each of its columns add up to one. It is useful to note that when  $A$  is doubly-stochastic, its Perron vector  $p$  defined in (22) is given by

$$p = \frac{1}{K} \mathbf{1} \quad (25)$$

# Constructing Combination Matrices

We will frequently be provided with an  $C$ , whose entries are positive but do not necessarily add up to one, and seek to transform it into a valid left- or doubly-stochastic combination matrix  $A$ . There are some popular constructions, which we list here. These constructions start with an undirected, unweighted graph with associated adjacency matrix  $C$  and Laplacian  $L$ . We also denote the degrees of the various agents by  $n_k$  and refer to the maximum degree across the graph by  $n_{\max}$ .

# Constructing Combination Matrices

- Averaging rule (left-stochastic):

$$a_{\ell k} = \begin{cases} 1/n_k, & \text{if } k \neq \ell \text{ are neighbors or } k = \ell \\ 0, & \text{otherwise} \end{cases} \quad (26)$$

- Laplacian rule (doubly-stochastic):

$$A = I_K - \beta L, \quad \beta > 0 \quad (27)$$

- Metropolis rule (doubly-stochastic):

$$a_{\ell k} = \begin{cases} 1/\max\{n_k, n_\ell\}, & \text{if } k \neq \ell \text{ are neighbors} \\ 1 - \sum_{\ell \in \mathcal{N}_k \setminus \{k\}} a_{\ell k}, & \text{when } k = \ell \\ 0, & \text{otherwise} \end{cases} \quad (28)$$

## Averaging over Graphs

We consider a collection  $\mathcal{N}$  of  $K$  agents, indexed by  $k$ , where each agent is assigned some vector-valued quantity  $g_k \in \mathbb{R}^M$  as its “signal”. The agents are arranged in a connected graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ . Our objective is to, in a decentralized manner, compute:

$$\{g_1, \dots, g_K\} \implies \bar{g} \triangleq \sum_{k=1}^K p_k g_k \quad (29)$$

- Here,  $\{p_k\}_{k=1}^K$  denote scalar convex combination weights, such that  $\sum_{k=1}^K p_k = 1$ . It is most common to let  $p_k = \frac{1}{K}$ , though we allow for arbitrary convex weights for generality.
- If the collection of signals  $\{g_k\}_{k=1}^K$  are jointly available at some central location, and in the absence of computational constraints, implementing (29) is straightforward by simply averaging the available realizations.

# Recap and Outlook

- We have introduced various concepts relating to graphs, as well as their properties.
- One key theme was that properties of the graph in many cases imply spectral properties on the adjacency or Laplacian matrix.
  - ▶ Hence, we can study graphs through linear algebra.
- We will now see how we can apply these concepts by performing a simple task over networks, namely computing averages.

## Averaging over Graphs

We consider a collection  $\mathcal{N}$  of  $K$  agents, indexed by  $k$ , where each agent is assigned some vector-valued quantity  $g_k \in \mathbb{R}^M$  as its “signal”. The agents are arranged in a connected graph  $\mathcal{G} = (\mathcal{N}, \mathcal{E})$ . The objective of a consensus algorithm is to, in a decentralized manner, compute:

$$\{g_1, \dots, g_K\} \implies \bar{g} \triangleq \sum_{k=1}^K p_k g_k \quad (30)$$

Here,  $\{p_k\}_{k=1}^K$  denote scalar convex combination weights, such that  $\sum_{k=1}^K p_k = 1$ . It is most common to let  $p_k = \frac{1}{K}$ , though we allow for arbitrary convex weights for generality.

## The Average Consensus Algorithm

We claim that we can pursue the average value  $\bar{g} = \sum_{k=1}^K p_k g_k$  by initializing  $w_{k,0} = g_k$  and then iterating:

$$w_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} w_{\ell,i-1} \quad (31)$$

where the weights  $a_{\ell k}$  are elements of a left-stochastic combination matrix  $A \in \mathbb{R}^{K \times K}$  with:

$$a_{\ell k} \begin{cases} > 0, & \text{if } \ell \in \mathcal{N}_k, \\ = 0, & \text{otherwise.} \end{cases} \quad \text{and} \quad \sum_{\ell=1}^K a_{\ell k} = 1 \quad (32)$$

As long as  $A$  is primitive, it then follows from the Perron-Frobenius Theorem, that:

$$Ap = p, \quad A^T \mathbf{1} = \mathbf{1}, \quad \lim_{i \rightarrow \infty} A^i = p \mathbf{1}^T \quad (33)$$

where  $p$  denotes the Perron eigenvector of  $A$ . Assume for now that the Perron entries  $p_k$  coincide with the desired averaging weights (more on design later).



## Network Quantities

Whenever studying the evolution of a collection of quantities over a network, it will be useful to introduce “network quantities” that summarize the state of the network in a compact form. We collect:

$$\mathcal{W}_i \triangleq \begin{pmatrix} w_{1,i} \\ w_{2,i} \\ \vdots \\ w_{K,i} \end{pmatrix} \quad (34)$$

and define

$$\mathcal{A} \triangleq A \otimes I_M \quad (35)$$

We can then write the consensus averaging recursion compactly as:

$$\mathcal{W}_i = \mathcal{A}^T \mathcal{W}_{i-1} \quad (36)$$

## Limiting Dynamics

Iterating the network recursion:

$$\lim_{i \rightarrow \infty} w_i = \lim_{i \rightarrow \infty} (\mathcal{A}^\top)^i w_0 = \left( \lim_{i \rightarrow \infty} (A^\top)^i \otimes I_M \right) w_0 \quad (37)$$

But we know from Perron-Frobenius that:

$$\lim_{i \rightarrow \infty} A^i = p \mathbf{1}^\top \implies \lim_{i \rightarrow \infty} (A^\top)^i = \mathbf{1} p^\top \quad (38)$$

and hence:

$$\lim_{i \rightarrow \infty} w_i = \left( \mathbf{1} p^\top \otimes I_M \right) w_0 = \begin{pmatrix} p^\top \otimes I_M \\ \vdots \\ p^\top \otimes I_M \end{pmatrix} \begin{pmatrix} g_1 \\ \vdots \\ g_K \end{pmatrix} = \begin{pmatrix} \sum_{k=1}^K p_k g_k \\ \vdots \\ \sum_{k=1}^K p_k g_k \end{pmatrix} \quad (39)$$

We conclude that  $\lim_{i \rightarrow \infty} w_{k,i} = \bar{g}$  for all  $k$ .

# Convergence of the consensus averaging algorithm

## Linear convergence of the consensus averaging algorithm

Each agent  $w_{k,i}$  following the consensus protocol (31) converges linearly to the weighted mean  $\bar{g} \triangleq \sum_{k=1}^K p_k g_k$ :

$$\sum_{k=1}^K \|\bar{g} - w_{k,i}\|^2 \leq \left\| p\mathbf{1}^\top - A \right\|^{2i} \times \sum_{k=1}^K \|\bar{g} - w_{k,0}\|^2, \quad (40)$$

where the rate of convergence is dictated by  $\|p\mathbf{1}^\top - A\|$ , and we are free to choose the norm.

## Sketch of Proof: Step 1 – Centroid is Stationary

Let us first examine the evolution of the weighted mean of iterates across the network

$\bar{w}_i = \sum_{k=1}^K p_k w_{k,i}$ . We have:

$$\begin{aligned}\bar{w}_i &= \sum_{k=1}^K p_k w_{k,i} = \left( p^\top \otimes I_M \right) w_i = \left( p^\top \otimes I_M \right) \mathcal{A}^\top w_{i-1} = \left( (Ap)^\top \otimes I_M \right) w_{i-1} \\ &= \left( p^\top \otimes I_M \right) w_{i-1} \\ &= \sum_{k=1}^K p_k w_{k,i-1} \\ &= \bar{w}_{i-1}\end{aligned}\tag{41}$$

We observe that, as a consequence of the fact that  $p$  is an eigenvector for the combination matrix  $A$ , the weighted mean  $\bar{w}_i$  is constant over time. We can conclude by iterating that:

$$\bar{w}_i = \bar{w}_{i-1} = \dots = \bar{w}_0 = \sum_{k=1}^K p_k g_k\tag{42}$$

## Sketch of Proof: Step 2 – Deviations around the Centroid Contract

It follows that the weighted mean, which we also refer to as a *centroid*, is fixed, while the local iterates  $w_{k,i}$  are scattered around it. We can then use the centroid as a reference point, and study to evolution of the network around it. Subtracting the consensus recursion from the augmented mean  $\mathbb{1}_K \otimes \bar{w}_i$ , we find:

$$\begin{aligned}\mathbb{1}_K \otimes \bar{w}_i - w_i &= \mathbb{1}_K \otimes \bar{w}_i - \mathcal{A}^\top w_{i-1} \\ &= \mathbb{1}_K \otimes \bar{w}_{i-1} - \mathcal{A}^\top w_{i-1} \\ &= \mathbb{1}_K \otimes \left( \left( p^\top \otimes I_M \right) w_{i-1} \right) - \mathcal{A}^\top w_{i-1} \\ &= \left( \mathbb{1} p^\top \otimes I_M \right) w_{i-1} - \mathcal{A}^\top w_{i-1} \\ &= \left( \mathbb{1} p^\top \otimes I_M - \mathcal{A}^\top \right) w_{i-1} \\ &= \left( \mathbb{1} p^\top \otimes I_M - \mathcal{A}^\top \right) (w_{i-1} - \mathbb{1} \otimes \bar{w}_{i-1})\end{aligned}\tag{43}$$

The result follows after taking norms and applying the sub-multiplicative property of norms.

## Centralities and Mixing Rate

For a graph  $\mathcal{G}$  with combination matrix  $A$ , we refer to the elements of the Perron eigenvector:

$$Ap = p \quad (44)$$

as the *agent centrality*, and the quantity:

$$\lambda_2 \triangleq \left\| A - p\mathbf{1}^\top \right\| \quad (45)$$

as the *mixing rate*. Note that for a symmetric combination matrix  $A = A^\top$ , by choosing  $\|\cdot\|$  to be the spectral norm, it follows that  $p = \frac{1}{K}\mathbf{1}$  and:

$$\left\| A - \frac{1}{K}\mathbf{1}\mathbf{1}^\top \right\| = \max \{ |\lambda_2(A)|, |\lambda_K(A)| \} \quad (46)$$

## Designing Perron Vectors and Mixing Rates

We consider an undirected graph with  $k \in \mathcal{N}_\ell \Leftrightarrow \ell \in \mathcal{N}_k$ . Then, the following choice of combination weights  $a_{\ell k}$  yields any desired Perron vector  $p = Ap$ :

$$a_{\ell k} = \begin{cases} 0, & \text{if } \ell \notin \mathcal{N}_k, \\ p_\ell, & \text{if } \ell \in \mathcal{N}_k \setminus \ell, \\ 1 - \sum_{m \in \mathcal{N}_k \setminus \ell} a_{mk}, & \text{when } \ell = k. \end{cases} \quad (47)$$

Combination matrices can also be designed to optimize  $\rho\left(A - \frac{1}{K}\mathbb{1}\mathbb{1}^\top\right)$  [Boyd, Diaconis, and Xiao 2004].

# Dynamic Consensus

We consider now a collection of  $K$  vector-valued and potentially random and time-varying signals  $\mathbf{g}_{k,i} \in \mathbb{R}^M$ , and wish to efficiently compute the average:

$$\bar{\mathbf{g}}_i \triangleq \sum_{k=1}^K p_k \mathbf{g}_{k,i} \quad (48)$$

where  $\{p_k\}_{k=1}^K$  denote scalar convex combination weights, such that  $\sum_{k=1}^K p_k = 1$ .



## Dynamic Consensus

To this end, it will be useful to formulate an optimization problem:

$$\bar{\mathbf{g}}_i = \arg \min_{\mathbf{w}} \sum_{k=1}^K p_k \|\mathbf{w}_k - \mathbf{g}_{k,i}\|^2, \text{ subject to } \mathbf{w}_k = \mathbf{w}_\ell \quad \forall k, \ell. \quad (49)$$

To verify that (49) indeed defines the weighted average  $\bar{\mathbf{g}}_i$ , we note that under the constraint  $\mathbf{w}_k = \mathbf{w}_\ell \quad \forall k, \ell$ , (49) is equivalent to:

$$\bar{\mathbf{g}}_i = \arg \min_{\mathbf{w}} \sum_{k=1}^K p_k \|\mathbf{w} - \mathbf{g}_{k,i}\|^2 \quad (50)$$

By differentiating, we find:

$$\sum_{k=1}^K p_k (\mathbf{w} - \mathbf{g}_{k,i}) = 0 \iff \mathbf{w} = \sum_{k=1}^K p_k \mathbf{g}_{k,i} \quad (51)$$

# Dynamic Consensus

We will now examine algorithms for dynamic consensus by reformulating (49) in a manner which makes it more amenable to decentralized computations. To this end, we note that, as long as the graph  $\mathcal{G}$  is connected, the constraint  $\mathbf{w}_k = \mathbf{w}_\ell \forall k, \ell$  can be simplified to the equivalent constraint  $\mathbf{w}_k = \mathbf{w}_\ell \forall \ell \in \mathcal{N}_k$ , where  $\mathcal{N}_k$  denotes the neighborhood of node  $k$ . This is because for any connected graph, there exists a sequence of edges leading from any node  $k$ , to an arbitrary node  $\ell$ , yielding a string of equality constraints leading from node  $k$  all the way to node  $\ell$ . It then follows that (49) is equivalent to:

$$\bar{\mathbf{g}}_i = \arg \min_{\mathbf{w}} \sum_{k=1}^K p_k \|\mathbf{w}_k - \mathbf{g}_{k,i}\|^2, \text{ subject to } \mathbf{w}_k = \mathbf{w}_\ell \forall \ell \in \mathcal{N}_k. \quad (52)$$

# Dynamic Consensus

Furthermore, we can write equivalently:

$$\bar{\mathbf{g}}_i = \arg \min_{\mathbf{w}} \sum_{k=1}^K p_k \|\mathbf{w}_k - \mathbf{w}_{k,i}\|^2, \text{ subject to } \sum_{k=1}^K \sum_{\ell=1}^K c_{\ell k} \|\mathbf{w}_k - \mathbf{w}_\ell\|^2 = 0. \quad (53)$$

where  $c_{\ell k} = c_{k\ell}$  are symmetric, non-negative weights with:

$$c_{\ell k} = \begin{cases} > 0 & \text{if } \ell \in \mathcal{N}_k, \\ = 0 & \text{otherwise.} \end{cases} \quad (54)$$

To verify that (52) and (53) are equivalent, we observe that:

$$\sum_{k=1}^K \sum_{\ell=1}^K c_{\ell k} \|\mathbf{w}_k - \mathbf{w}_\ell\|^2 = 0 \iff \|\mathbf{w}_k - \mathbf{w}_\ell\|^2 = 0 \forall \ell \in \mathcal{N}_k \iff \mathbf{w}_k = \mathbf{w}_\ell, \forall \ell \in \mathcal{N}_k. \quad (55)$$

## Penalty-Based Dynamic Consensus

Rather than solve (53) with the exact constraint  $\sum_{k=1}^K \sum_{\ell=1}^K c_{\ell k} \|\mathbf{w}_k - \mathbf{w}_\ell\|^2 = 0$ , we penalize deviation from the constraint as in:

$$\bar{\mathbf{g}}_i^\eta = \arg \min_{\mathbf{w}} \sum_{k=1}^K p_k \|\mathbf{w}_k - \mathbf{g}_{k,i}\|^2 + \frac{\eta}{4} \sum_{\ell=1}^K c_{\ell k} \|\mathbf{w}_k - \mathbf{w}_\ell\|^2 \quad (56)$$

where  $\eta > 0$  denotes a new penalty parameter. It is important to note that, while (49) through (53) were *equivalent* reformulations of the consensus problem, relation (56) is only an approximation for finite  $\eta$ . Letting  $\eta \rightarrow \infty$  ensures  $\sum_{\ell=1}^K c_{\ell k} \|\mathbf{w}_k - \mathbf{w}_\ell\|^2 = 0$ , and hence exact consensus.

## Penalty-Based Dynamic Consensus

In the sequel, we will quantify this difference, and present an algorithm for pursuing  $\bar{\mathbf{g}}_i^\eta$ . To this end, it will be useful to write the penalized cost (56) in terms of network level quantities. We introduce:

$$P = \begin{pmatrix} p_1 & 0 & \cdots & 0 \\ 0 & p_2 & \cdots & 0 \\ 0 & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & p_K \end{pmatrix} = \text{diag}\{p\} \quad (57)$$

$$\mathcal{P} = P \otimes I_M \quad (58)$$

Then:

$$\sum_{k=1}^K p_k \|\mathbf{w}_k - \mathbf{g}_{k,i}\|^2 = \|\mathbf{w} - \mathbf{g}_i\|_{\mathcal{P}}^2 \quad (59)$$

# Penalty-Based Dynamic Consensus

For the second term in (56), we can verify, that:

$$\frac{\eta}{2} \sum_{k=1}^K \sum_{\ell \in \mathcal{N}_k} c_{\ell k} \|\mathbf{w}_k - \mathbf{w}_\ell\|^2 = \eta \mathbf{w}^\top \mathcal{L} \mathbf{w} \quad (60)$$

where  $\mathcal{L} = L \otimes I_M$  and we defined:

$$L = \text{diag}\{C\mathbf{1}\} - C \quad (61)$$

Hence, problem (56) is equivalent to:

$$\bar{\mathbf{g}}_i^\eta = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{g}_i\|_{\mathcal{P}}^2 + \frac{\eta}{2} \mathbf{w}^\top \mathcal{L} \mathbf{w} \quad (62)$$

## Penalty-Based Dynamic Consensus

Note that this objective function is time-varying. Nevertheless, we can pursue its minimizing argument via (online) gradient-descent, pre-conditioned by the positive-definite matrix  $\mathcal{P}^{-1}$ :

$$\begin{aligned}\mathbf{w}_i &= \mathbf{w}_{i-1} - \mu \mathcal{P}^{-1} (\mathcal{P} (\mathbf{w}_{i-1} - \mathbf{g}_i) + \mu \eta \mathcal{L} \mathbf{w}_{i-1}) \\ &= \mathbf{w}_{i-1} - \mu (\mathbf{w}_{i-1} - \mathbf{g}_i) - \mu \eta \mathcal{P}^{-1} \mathcal{L} \mathbf{w}_{i-1} \\ &= ((1 - \mu)I - \mu \eta \mathcal{P}^{-1} \mathcal{L}) \mathbf{w}_{i-1} + \mu \mathbf{g}_i \\ &= (I - \mu \eta \mathcal{P}^{-1} \mathcal{L}) \mathbf{w}_{i-1} + \mu (\mathbf{g}_i - \mathbf{w}_{i-1})\end{aligned}\tag{63}$$

We set  $\mu = \eta^{-1}$  and define:

$$\mathcal{A}^\top \triangleq I - \mathcal{P}^{-1} \mathcal{L}\tag{64}$$

# Penalty-Based Dynamic Consensus

We can then equivalently write:

$$\mathbf{w}_i = \mathcal{A}^\top \mathbf{w}_{i-1} + \mu (\mathbf{g}_i - \mathbf{w}_{i-1}) \quad (65)$$

It is instructive to examine the spectral structure of  $\mathcal{A}^\top$  defined as in (64). Since  $\mathcal{L}\mathbf{1} = 0$ , we have:

$$A^\top \mathbf{1} = (I - \mu\eta P^{-1}L) \mathbf{1} = \mathbf{1} - \mu\eta P^{-1}L\mathbf{1} = \mathbf{1} \quad (66)$$

$$Ap = (I - \mu\eta LP^{-1}) p = p - \mu\eta LP^{-1}p = p - \mu\eta L\mathbf{1} = p \quad (67)$$

We conclude that  $\mathbf{1}$  and  $p$  are left and right eigenvectors corresponding to an eigenvalue at one respectively. Returning to node level quantities, we arrive at:

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \mathbf{w}_{\ell,i-1} + \mu (\mathbf{g}_{k,i} - \mathbf{w}_{k,i-1}) \quad (68)$$



## Bias of Penalty-Based Algorithms

The dynamic consensus algorithm we derived relies on a penalty-reformulation of the networked averaging problem, which is merely an approximation unless  $\eta \rightarrow \infty$  which implies  $\mu \rightarrow 0$ . We can verify that even when applying this algorithm to a constant set of signal, this will result in a bias compared to the classical consensus averaging algorithm. Thus, consider the case when  $\mathbf{g}_{k,i} \triangleq \mathbf{g}_k$  is deterministic and constant. We then have:

$$\mathbf{w}_i = \left( \mathcal{A}^T - \mu I \right) \mathbf{w}_{i-1} + \mu \mathbf{g} \quad (69)$$

Iterating we find:

$$\mathbf{w}_i = \left( \mathcal{A}^T - \mu I \right)^i \mathbf{w}_0 + \mu \left( \sum_{j=0}^{i-1} \left( \mathcal{A}^T - \mu I \right)^j \right) \mathbf{g} \quad (70)$$

# Bias of Penalty-Based Algorithms

In the limit, we then have:

$$\lim_{i \rightarrow \infty} w_i = \mu \left( \sum_{j=0}^{i-1} (\mathcal{A}^\top - \mu I)^j \right) g = \mu \left( I - (\mathcal{A}^\top - \mu I) \right)^{-1} g = \mu \left( (1 + \mu)I - \mathcal{A}^\top \right)^{-1} g \quad (71)$$

By exploiting the Jordan decomposition of the primitive matrix  $A$ , we can expand this expression and conclude:

$$\sum_{k=1}^K \left\| \lim_{i \rightarrow \infty} w_{k,i} - g_k \right\|^2 = O(\mu^2) \quad (72)$$

Which means that we pay for the ability to track time-varying signals with a bias that is  $O(\mu^2)$  when the signals are constant.

# Exact Dynamic Consensus Algorithm

The previous discussion motivates the question, whether it is possible to design a consensus algorithm over a graph, which is able to average a constant signal perfectly, while having some ability to track slowly time varying signals. It turns out that this is possible, and this algorithm can be derived using control-theoretic tools [Kia et al 2019] or primal-dual techniques we will develop in the next lecture. For the time being, we simply list the exact dynamic consensus algorithm, which takes the form:

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \mathbf{w}_{\ell,i-1} + \mathbf{g}_{k,i} - \mathbf{g}_{k,i-1} \quad (73)$$

or in network level quantities:

$$\mathbf{w}_i = \mathcal{A}^T \mathbf{w}_{i-1} + \mathbf{g}_i - \mathbf{g}_{i-1} \quad (74)$$

## Constant Signals

We begin by establishing the ability of the dynamic consensus (74) to compute exactly the mean of a constant signal over the graph. To this end, let  $\mathbf{g}_i \triangleq \mathbf{g}$ , and we find:

$$\mathbf{w}_i = \mathcal{A}^T \mathbf{w}_{i-1} + \mathbf{g} - \mathbf{g} = \mathcal{A}^T \mathbf{w}_{i-1} \quad (75)$$

which corresponds precisely to the static consensus algorithm (36). Under the initialization  $\mathbf{w}_0 \triangleq \mathbf{g}$ , we conclude that the evolution of the dynamic consensus algorithm (74) matches that of the static consensus algorithm (36), provided that the signal of interest  $\mathbf{g}_{i-1}$  is constant.

# Tracking Random Walks

The question now is whether, and how well, the dynamic variant of the consensus algorithms is able to track the mean of the signal of interest when it is time-varying. To quantify the tracking performance, we will employ the simple random-walk model:

$$\mathbf{g}_i = \mathbf{g}_{i-1} + \mathbf{v}_i \quad (76)$$

where  $\mathbf{v}_i \in \mathbb{R}^{MK}$  denotes a random driving variable with mean  $\mu_v = \mathbb{E}\mathbf{v}_i$  and variance  $\sigma_v^2 = \mathbb{E}\|\mathbf{v}_i - \mu_v\|^2$ . We note that, while the driving term  $\mathbf{v}_i$  is presumed stationary, the resulting dynamics of  $\mathbf{g}_i$  are not. Under this model, we have:

$$\mathbf{w}_i = \mathcal{A}^\top \mathbf{w}_{i-1} + \mathbf{g}_i - \mathbf{g}_{i-1} = \mathcal{A}^\top \mathbf{w}_{i-1} + \mathbf{v}_i \quad (77)$$

## Tracking Analysis: Centroid

We proceed similarly to the analysis of the static consensus algorithm, and decompose the evolution of the local estimates into the network centroid  $\mathbf{w}_{c,i} \triangleq \sum_{k=1}^K p_k \mathbf{w}_{k,i}$  and the deviations from the centroid  $\mathbf{w}_{k,i} - \mathbf{w}_{c,i}$ . We begin with the centroid, which evolves according to:

$$\mathbf{w}_{c,i} = \mathbf{w}_{c,i-1} + \sum_{k=1}^K p_k \mathbf{v}_{k,i} \quad (78)$$

## Tracking Analysis: Centroid

In comparison, for the spatial average of the time varying signal  $\mathbf{g}_i$ , we have:

$$\begin{aligned}\mathbf{g}_{c,i} &\triangleq \sum_{k=1}^K p_k \mathbf{g}_{k,i} = \left( p^\top \otimes I_M \right) \mathbf{g}_i \\ &\stackrel{(76)}{=} \left( p^\top \otimes I_M \right) (\mathbf{g}_{i-1} + \mathbf{v}_i) \\ &= \sum_{k=1}^K p_k \mathbf{g}_{k,i-1} + \sum_{k=1}^K p_k \mathbf{v}_{k,i} \\ &= \mathbf{g}_{c,i-1} + \sum_{k=1}^K p_k \mathbf{v}_{k,i}\end{aligned}\tag{79}$$

Comparing the expression (78) for  $\mathbf{w}_{c,i}$  with (79) for  $\mathbf{g}_{c,i}$ , we find that the recursions match exactly.

## Tracking Analysis: Centroid

We can hence conclude that, as long as the  $\{\mathbf{w}_{k,0}\}_{k=1}^K$  are initialized appropriately, namely to  $\mathbf{w}_{k,0} = \mathbf{g}_{k,0}$ , we have with probability one:

$$\mathbf{w}_{c,i} = \mathbf{g}_{c,i} \iff \sum_{k=1}^K p_k \mathbf{w}_{k,i} = \sum_{k=1}^K p_k \mathbf{g}_{k,i} \quad (80)$$

This fact is analogous to relation (41) for the static consensus algorithm, with the important advantage that the correction term  $\mathbf{g}_{i-1} - \mathbf{g}_{i-2}$  in (74) allows the dynamic algorithm to track a *time-varying* centroid. Perhaps surprisingly, the network centroid  $\mathbf{w}_{c,i}$  under this construction tracks the signal centroid  $\mathbf{g}_{c,i}$  perfectly, and with probability one, independently of the power or variance of the random variable  $\mathbf{v}_i$  driving the random walk.



## Tracking Analysis: Deviation from the Centroid

This fact of course does not imply that any given agent  $\mathbf{w}_{k,i}$  is necessarily able to track the signal centroid  $\mathbf{g}_{c,i}$  well. We now proceed to study the deviation of local estimates  $\mathbf{w}_{k,i}$  around the network centroid  $\mathbf{w}_{c,i}$  under the random walk model (76) and the dynamic consensus algorithm (74) to establish this stronger result. We have:

$$\begin{aligned} & \mathbf{w}_i - \mathbf{1} \otimes \mathbf{w}_{c,i} \\ &= \mathbf{w}_i - \left( \mathbf{1} p^\top \otimes I_M \right) \mathbf{w}_i \\ &= \mathcal{A}^\top \mathbf{w}_{i-1} + \mathbf{v}_i - \left( \mathbf{1} p^\top \otimes I_M \right) \left( \mathcal{A}^\top \mathbf{w}_{i-1} + \mathbf{v}_i \right) \\ &= \left( \mathcal{A}^\top - \mathbf{1} p^\top \otimes I_M \right) \mathbf{w}_{i-1} + \left( I_{MK} - \mathbf{1} p^\top \otimes I_M \right) \mathbf{v}_i \end{aligned} \quad (81)$$

Upon iterating and taking the limit as  $i \rightarrow \infty$ , we have:

$$\lim_{i \rightarrow \infty} \mathbb{E} \|\mathbf{w}_i - \mathbf{1} \otimes \mathbf{w}_{c,i}\|^2 \leq \frac{1}{1 - \lambda_2} \left( \sum_{j=0}^{\infty} \lambda_2^j \right) \sigma_v^2 = \frac{\sigma_v^2}{(1 - \lambda_2)^2} \quad (82)$$

## Simulation Study

A key take-away of this lecture is that the rate at which average consensus is achieved over a graph is determined by its mixing rate. We illustrate this numerically here. To this end, we generate first an unweighted Erdos-Renyi graph, where  $c_{\ell k} = 1$  with probability  $p_{\text{edge}}$ . To ensure the graph is undirected, we set:

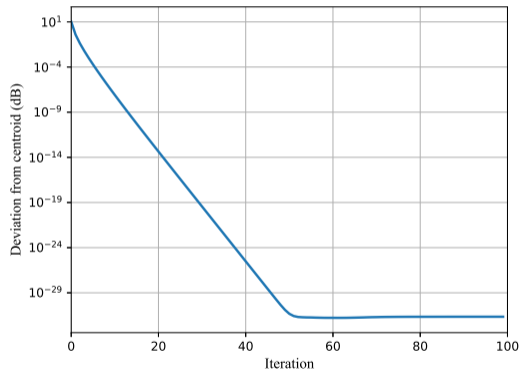
$$C \longleftarrow \frac{1}{2} (C + C^T) \quad (83)$$

Given this symmetric adjacency matrix, we can then construct a symmetric combination matrix  $A = A^T$  using the Metropolis-Rule. The mixing rate in that case is given by:

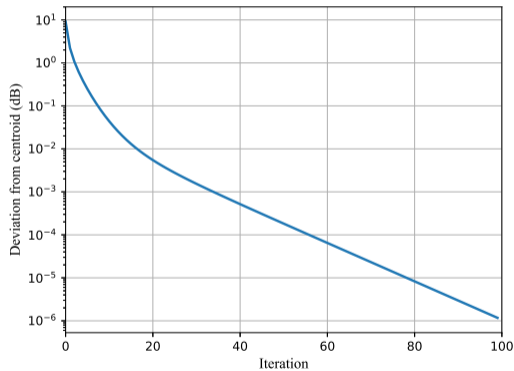
$$\lambda_2 = \rho \left( A - \frac{1}{K} \mathbf{1}\mathbf{1}^T \right) \quad (84)$$

We can control  $\lambda_2$  indirectly by adjusting the edge probability  $p_{\text{edge}}$ .

# Convergence of the Consensus Averaging Algorithm for Differing Network Connectivities



$$p_{\text{edge}} = 0.3, \lambda_2 = 0.498$$



$$p_{\text{edge}} = 0.05, \lambda_2 = 0.950$$

# Conclusion

- We introduced graphs and their (spectral) properties. Two key quantities are:
  - ▶ The mixing rate: Quantifies how well connected the graph is, and how quickly information spreads.
  - ▶ The agent centrality: Quantifies the influence an agent has over the limiting behavior of the algorithm.
- We illustrated these in the context of averaging over graphs, a simple task that nevertheless captures some of the fundamental trade-offs.
- We also developed dynamic consensus algorithms, which allow for the tracking of averages of time-varying signals.

# References and Further Reading

- General references and surveys:
  - ▶ A. H. Sayed, *Inference and Learning from Data*, Cambridge University Press, 2022.
  - ▶ A. H. Sayed, “Adaptation, learning, and optimization over networks,” *Foundations and Trends in Machine Learning*, vol. 7, no. 4-5, pp. 311–801, July 2014.
- Averaging over graphs:
  - ▶ S. Boyd, P. Diaconis, and L. Xiao, “Fastest Mixing Markov Chain on a Graph”, *SIAM Review*, vol. 46, no. 4, 2004.
  - ▶ S. S. Kia, B. Van Scoy, J. Cortes, R. A. Freeman, K. M. Lynch and S. Martinez, “Tutorial on Dynamic Average Consensus: The Problem, Its Applications, and the Algorithms,” in *IEEE Control Systems Magazine*, vol. 39, no. 3, pp. 40-72, June 2019.