Optimization Methods for Graph Learning

Prof. Daniel P. Palomar

ELEC5470/IEDA6100A - Convex Optimization The Hong Kong University of Science and Technology (HKUST) Fall 2020-21

Outline

1 Graphs

2 Basics

3 Learning Graphs from Data

- Similarity function based
- Smooth signal based
- i.i.d. model based
- Structured graphs via spectral constraints

Wumerical Experiments

Outline

Graphs

2 Basics

3 Learning Graphs from Data

- Similarity function based
- Smooth signal based
- i.i.d. model based
- Structured graphs via spectral constraints

4 Numerical Experiments

Graphical models

Graphical models are a way to represent knowledge:

- **nodes** correspond to the entities (variables);
- edges encode the relationships between entities (dependencies between the variables).



Examples

Financial graph: represents inter-dependencies of financial companies and the data are the economic indices (stock price, volume, etc.) of each entity.



Examples

Social graph: representa behavioral similarity/influence between people and the data are the online activities (tagging, liking, purchasing).



Graphical model importance

- Graphs are intuitive way of representing and visualising the relationships between entities.
- Graphs allow us to abstract out the **conditional independence** relationships and to answer questions like: "Is x_1 dependent on x_6 given that we know the value of x_8 ?" just by looking at the graph.
- Graph models constitute an effective **representation of data**, available across numerous domains in science and engineering.
- Graphs are widely used in a variety of applications in machine learning, graph CNN, graph signal processing, functional connectivity between brain regions, behavioral influence among groups of people, effect among stocks, etc.
- Graphs offer a **language** through which different **disciplines** can seamlessly interact with each other.
- It captures the actual geometry of data.
- It allows a visualization of high-dimensional data.

- Graphical models are about having a graph representation that can encode **relationships** between entities.
- In many cases, the relationships between entities are straightforward and direct:
 - Are two people friends in a social network?
 - Are two researchers co-authors in a published paper?
- In many other cases, relationships are not known and must be learned:
 - Does one gene regulate the expression of others?
 - Which drug alters the pharmacologic effect of another drug?
- The **choice of graph representation** affects the subsequent analysis and eventually the performance of any graph-based algorithm.

The goal is to learn a graph representation of data with specific properties (e.g., structures).

Graph learning from data

- Given a data matrix $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$, each column $\mathbf{x}_i \in \mathbb{R}^p$ is a graph signal (one observation) and there are *n* observations.
- The goal is to obtain a graph representation of the data.



- A graph is a simple **mathematical structure** described by $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{W})$, where
 - \mathcal{V} contains the set of nodes $\mathcal{V} = \{1, 2, 3, \dots, p\};$
 - $\mathcal{E} = \{(1,2), (1,3), ..., (i,j), ..., (p, p-1)\}$ is the set of edges between pair of nodes (i, j);
 - ${\scriptstyle \bullet}\,$ the weight matrix ${\bf W}$ encodes the strength of the relationships.

Graph learning example

Example of transforming a dataset with points in \mathbb{R}^2 into a graph:





• Introductory references: (Mateos et al. 2019)¹ and (Dong et al. 2019)².

¹G. Mateos, S. Segarra, A. G. Marques, and A. Ribeiro, "Connecting the dots," *IEEE Signal Processing Magazine*, vol. 36, no. 3, pp. 16–43, 2019.

²X. Dong, D. Thanou, M. Rabbat, and P. Frossard, "Learning graphs from data: A signal representation perspective," *IEEE Signal Processing Magazine*, vol. 36, no. 3, pp. 44–63, 2019.

D. Palomar (HKUST)

Graph Learning

Outline

1 Graphs

2 Basics

3 Learning Graphs from Data

- Similarity function based
- Smooth signal based
- i.i.d. model based
- Structured graphs via spectral constraints

With a set of the set

Graph and its matrix representation

• Connectivity matrix C, Adjacency matrix W, and Laplacian matrix L:

$$[\mathbf{C}]_{ij} = \begin{cases} 1 \text{ if } (i,j) \in \mathcal{E} \\ 0 \text{ if } (i,j) \notin \mathcal{E} \\ 0 \text{ if } (i,j) \notin \mathcal{E} \end{cases}, \ [\mathbf{W}]_{ij} = \begin{cases} w_{ij} \text{ if } (i,j) \in \mathcal{E} \\ 0 \text{ if } (i,j) \notin \mathcal{E} \\ 0 \text{ if } i = j \end{cases} \begin{cases} -w_{ij} \text{ if } (i,j) \in \mathcal{E} \\ 0 \text{ if } (i,j) \notin \mathcal{E} \\ \sum_{j=1}^{p} w_{ij} \text{ if } i = j \end{cases}$$

• Example: $\mathcal{V} = \{1, 2, 3, 4\}$, $\mathcal{E} = \{(1, 2), (1, 3), (2, 3), (2, 4)\}$, and $\mathbf{W} = \{2, 2, 3, 1\}$:



$$[\mathbf{C}]_{ij} = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}, \ [\mathbf{W}]_{ij} = \begin{bmatrix} 0 & 2 & 2 & 0 \\ 2 & 0 & 3 & 1 \\ 2 & 3 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}, \ [\mathbf{L}]_{ij} = \begin{bmatrix} 4 & -2 & -2 & 0 \\ -2 & 6 & -3 & -1 \\ -2 & -3 & 5 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix}$$

D. Palomar (HKUST)

Graph Learning

Graph terminology

Graph matrices for a weighted graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{W})$:

- Weighted adjacency matrix of a graph (or simply adjacency matrix), W, is defined as $W_{ij} = w_{ij} \ge 0$. It is a symmetric matrix (undirected graph) and $W_{ii} = 0$ (no self-loops).
- Connectivity matrix: C is a particular case of the adjacency matrix with 0-1 elements.
- The **degree matrix** is the diagonal matrix **D** that contains the degrees d_1, \ldots, d_p along the diagonal:

$$\mathsf{D} = \mathsf{Diag}(\mathsf{W1})$$
 .

The degree d_i of the vertex *i* is defined as the row sum $d_i = \sum_{j \neq i} W_{i,j} = \sum_j w_{ij}$.

• Laplacian of a graph:

$$L = D - W.$$

• Normalized Laplacian of a graph:

$$L_{norm} = D^{1/2}LD^{1/2} = I - D^{1/2}WD^{1/2}.$$

The Laplacian graph matrix

- The adjacency matrix **W** already contains all the graph information, so why do we need other graph matrices like the Laplacian matrix?
- The Laplacian L = D W, where D = Diag (W1), has a nice interesting physical meaning as well as many nice properties:
 - L is a symmetric and positive semidefinite matrix $\textbf{L} \succeq \textbf{0};$
 - the number of zero eigenvalues denotes the number of connected components of the graph;
 - L is singular with eigenvector 1: L1 = 0.
- **Physical interpretation**: Denote with the vector $\mathbf{x} = (x_1, \dots, x_p)$ the value of some quantity on all the *p* nodes. The Laplacian measures the **smoothness** or variance of that vector weighted with the graph weights:

$$\mathbf{x}^T \mathbf{L} \mathbf{x} = \frac{1}{2} \sum_{i,j} w_{ij} (x_i - x_j)^2.$$

Proof:

$$\mathbf{x}^{T}\mathbf{L}\mathbf{x} = \mathbf{x}^{T}\mathbf{D}\mathbf{x} - \mathbf{x}^{T}\mathbf{W}\mathbf{x} = \sum_{i} d_{i}x_{i}^{2} - \sum_{i,j} w_{ij}x_{i}x_{j} = \sum_{i,j} w_{ij}x_{i}^{2} - \sum_{i,j} w_{ij}x_{i}x_{j}.$$

Outline

Graphs

2 Basics

3 Learning Graphs from Data

- Similarity function based
- Smooth signal based
- i.i.d. model based
- Structured graphs via spectral constraints

With a set of the set

We will consider different ways to learn the graph **adjacency matrix**, **W**, or **Laplacian matrix**, **L**, with input either the data matrix **X** or the sample covariance matrix **S**.

- Similarity function based methods for adjacency **W** (with input the data matrix **X**). Basically, two nodes *i* and *j* are connected based on some similarity function, leading to:
 - thresholded Euclidean distance graph;
 - Gaussian graph;
 - thresholded Gaussian graph;
 - k-nearest neighbors (k-NN) graph;
 - feature correlation graph;
 - self-tuned Gaussian graph.

• Smooth signal based methods for Laplacian L:

- graphs from smooth signals;
- closest *k*-connected graph.

• i.i.d. model based methods for Laplacian L (with input the sample covariance matrix S):

- covariance matrix (with and without market factor, which applies to all the methods);
- correlation matrix;
- precision matrix;
- graphical LASSO (GLASSO);
- Laplacian-structured GLASSO;
- Laplacian with spectral constraints (e.g., for *k*-connected graph);
- shift operator based.

Outline

1 Graphs

2 Basics

3 Learning Graphs from Data

- Similarity function based
- Smooth signal based
- i.i.d. model based
- Structured graphs via spectral constraints

4 Numerical Experiments

Basic similarity measures

Write the data matrix in terms of the node signals $\mathbf{x}_i \in \mathbb{R}^n$ along the rows: $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p]^T \in \mathbb{R}^{p \times n}$.

To choose the edges and weights to form the **adjacency matrix W**, either weighted or 0-1 connectivity matrix, we can use different methods:

- Thresholded Euclidean distance graph: nodes *i* and *j* are connected (*w_{ij}* = 1) if the corresponding signals satisfy ||**x**_i − **x**_j||² ≤ γ (where γ is a threshold); otherwise not connected (*w_{ij}* = 0).
- **Gaussian graph**: set every pair of points $i \neq j$ as connected with the following weights:

$$w_{ij} = \exp\left(-rac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}
ight),$$

where σ^2 controls the size of the neighborhood.

• Thresholded Gaussian graph: we can combine the previous two graph constructions, i.e., set every pair of points $i \neq j$ satisfying $\|\mathbf{x}_i - \mathbf{x}_j\|^2 \leq \gamma$ with Gaussian weights; otherwise not connected $(w_{ij} = 0)$.

Basic similarity measures

- k-nearest neighbors (k-NN) graph: nodes i and j are connected (w_{ij} = 1) if x_i is one of the k closest points (Euclidean distance) to x_j or viceversa; otherwise not connected (w_{ij} = 0).
- Feature correlation graph: simply use the pairwise feature correlation for $i \neq j$:

$$w_{ij} = \mathbf{x}_i^T \mathbf{x}_j.$$

• Self-tuned Gaussian graph: the Gaussian graph tends to be too densely connected, this approach automatically chooses a different σ for each point; in particular, it normalizes the distance from each node *i* to the other nodes with the distance to its *k*-NN:

$$w_{ij} = \exp\left(-rac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_i\sigma_j}
ight),$$

where σ_i denotes the distance between the *i*th node and its *k*-NN. Basic reference: (Manor and Perona 2004)³.

³L. Z. Manor and P. Perona, "Self-tuning spectral clustering," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2004.

D. Palomar (HKUST)

Graphs from basic constructions

Thresholded distance 0-1 graph (two-moon data)



Thresholded distance 0-1 graph (three-circle data)





Gaussian graph (three-circle data)



Thresholded Gaussian graph (two-moon data)



Thresholded Gaussian graph (three-circle data)



Graphs from basic constructions



Outline

Graphs

2 Basics

3 Learning Graphs from Data

- Similarity function based
- Smooth signal based
- i.i.d. model based
- Structured graphs via spectral constraints

With a set of the set

Learning graphs from data: Smooth signals

- Previously, we have learned different basic ways to construct the adjacency matrix W, either weighted or a 0-1 connectivity matrix, from a set of data points {x_i}ⁿ_{i=1}, where each *p*-dimensional vector x_i = (x_{1i},..., x_{pi}) contains the signal from the *p* nodes of the graph.
- In particular, we have considered the following similarity function based methods:
 - thresholded Euclidean distance graph;
 - Gaussian graph;
 - thresholded Gaussian graph;
 - k-nearest neighbors (k-NN) graph;
 - feature correlation graph;
 - self-tuned Gaussian graph.
- We will now consider more sophisticated **smooth signal based** methods to learn the graph Laplacian matrix L:
 - graphs from smooth signals;
 - closest *k*-connected graph.

• Recall that a measure of smoothness or variance of a graph-signal $\mathbf{x} \in \mathbb{R}^{p}$ is

$$\mathbf{x}^T \mathbf{L} \mathbf{x} = \frac{1}{2} \sum_{i,j} w_{ij} (x_i - x_j)^2.$$

- Given *p* signals, each containing *n* observations along the rows, in the data matrix $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p]^T \in \mathbb{R}^{p \times n}$, we can measure its smoothness over the graph \mathcal{G} as $Tr(\mathbf{X}^T \mathbf{L} \mathbf{X}) = \frac{1}{2} \sum_{i,j} w_{ij} ||\mathbf{x}_i \mathbf{x}_j||^2$.
- When the graph \mathcal{G} is not available, we can learn it from the data **X** by finding the graph weights that minimize the variance term combined with some regularization term:

minimize
$$Tr(\mathbf{X}^T \mathbf{L} \mathbf{X}) + \gamma h(\mathbf{L}).$$

- Smaller distance $\|\mathbf{x}_i \mathbf{x}_j\|^2$ between data points \mathbf{x}_i and \mathbf{x}_j will force to learn a graph with larger affinity value w_{ij} , and vice versa.
- Higher values of weight w_{ij} will imply the signals \mathbf{x}_i and \mathbf{x}_j are similar and, hence, strongly connected.
- $h(\mathbf{L})$ is a regularization function, e.g., $\|\mathbf{L}\|_1$, $\|\mathbf{L}\|_F^2$, and $-\log\det(\mathbf{L})$.

D. Palomar (HKUST)

- The graph learning formulation from smooth signals can be formulated in convex form in terms of either the adjacency matrix **W** or the Laplacian matrix $\mathbf{L} = \mathbf{D} \mathbf{W}$, where $\mathbf{D} = \text{Diag}(\mathbf{d})$ and $\mathbf{d} = \mathbf{W}\mathbf{1}$ is the degree vector of the nodes.
- Formulation in terms of Laplacian matrix:

$$\begin{array}{ll} \underset{\textbf{L} \succeq \textbf{0}}{\text{minimize}} & \mathsf{Tr}(\textbf{X}^{T}\textbf{L}\textbf{X}) + \frac{\gamma}{2} \|\textbf{L}\|_{F, \text{off}}^{2} \\ \text{subject to} & \mathsf{diag}(\textbf{L}) = \textbf{1} \\ & \textbf{L}\textbf{1} = \textbf{0}, \quad L_{ij} = L_{ji} \leq 0, \quad \forall i \neq j \end{array}$$

where

- the regularization term $\|\mathbf{L}\|_{F,off}^2$ controls the energy of the off-diagonal elements;
- the constraint d = diag(L) = 1 controls the degrees of the nodes (to get a balanced graph); and
- we have the usual Laplacian constraints L1 = 0 and $L_{ij} = L_{ji} \le 0$ for $i \ne j$.

• Formulation in terms of adjacency matrix:

$$\begin{array}{ll} \underset{\mathbf{W}}{\text{minimize}} & \sum_{i,j=1}^{p} w_{ij} \|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2} + \gamma \sum_{i,j=1}^{p} w_{ij}^{2} \\ \text{subject to} & \sum_{j=1}^{p} w_{ij} = 1, \quad \forall i \\ & w_{ii} = 0, \quad w_{ij} = w_{ji} \geq 0, \quad \forall i \neq j. \end{array}$$

• Defining $Z_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|^2$, we can rewrite the problem compactly:

$$\begin{array}{ll} \underset{\mathbf{W}}{\text{minimize}} & \operatorname{Tr}(\mathbf{WZ}) + \gamma \|\mathbf{W}\|_{F}^{2} \\ \text{subject to} & \mathbf{W1} = \mathbf{1} \\ & \operatorname{diag}(\mathbf{W}) = \mathbf{0}, \quad \mathbf{W} = \mathbf{W}^{T} \geq \mathbf{0} \end{array}$$

- With $\gamma=0$ we would get just one connection per node!
- With $\gamma \to \infty$ we get instead a fully dense graph (Cauchy-Schwartz).
- So the term $\|\mathbf{W}\|_F^2$ makes the solution nonsparse!
- Both formulations (in terms of adjacency matrix or Laplacian matrix) are convex and can be easily solved with a solver. Interestingly, we can develop a simple tailored algorithm.

- If we ignore the symmetry constraint in the adjacency matrix $\mathbf{W} = \mathbf{W}^{T}$, then the problem can be nicely solved separately for each row/column \mathbf{w}_i as in (Nie et al. 2016)⁴.
- Note that at a later stage we will have to symmetrize it:

$$\mathbf{W} \leftarrow rac{1}{2}(\mathbf{W} + \mathbf{W}^{\mathcal{T}})$$

• The problem of the *i*th row is

$$\begin{array}{ll} \underset{\mathbf{w}_{i}}{\text{minimize}} & \sum_{j=1}^{p} w_{ij} z_{ji} + \gamma \sum_{j=1}^{p} w_{ij}^{2} \\ \text{subject to} & \sum_{j=1}^{p} w_{ij} = 1, \ w_{ii} = 0, \ w_{ij} \geq 0 \end{array}$$

More compactly

$$\begin{array}{ll} \underset{\mathbf{w}_{i}}{\text{minimize}} & \mathbf{z}_{i}^{T}\mathbf{w}_{i} + \gamma \|\mathbf{w}_{i}\|^{2} \\ \text{subject to} & \mathbf{w}_{i}^{T}\mathbf{1} = 1, \ w_{ii} = 0, \ \mathbf{w}_{i} \geq \mathbf{0}. \end{array}$$

⁴F. Nie, X. Wang, M. Jordan, and H. Huang, "The constrained Laplacian rank algorithm for graph-based clustering," in *Proc. of the Thirtieth Conference on Artificial Intelligence (AAAI)*, Phoenix, Arizona, USA, 2016, pp. 1969–1976.

• We can finally write our problem as

$$\begin{array}{ll} \underset{\mathbf{w}_{i}}{\text{minimize}} & \frac{1}{2} \|\mathbf{w}_{i} + \frac{\mathbf{z}_{i}}{2\gamma}\|^{2} \\ \text{subject to} & \mathbf{w}_{i}^{T}\mathbf{1} = 1, \ w_{ii} = 0, \ \mathbf{w}_{i} \geq \mathbf{0}. \end{array}$$

• The Lagrangian is

$$\mathcal{L}(\mathbf{w}_i; \eta_i, \beta_i) = \frac{1}{2} \left\| \mathbf{w}_i + \frac{\mathbf{z}_i}{2\gamma} \right\|^2 - \eta_i (\mathbf{w}_i^T \mathbf{1} - 1) - \beta_i^T \mathbf{w}_i$$

where η_i and $\beta_i \in \mathbb{R}^p$ are the Lagrangian multipliers with $\beta_{ji} \ge 0, \ \forall j \neq i$.

• From the KKT optimality conditions, the optimal solution (defining $z_{ii} = \infty$ so that $w_{ii} = 0$) is

$$\mathbf{w}_i = \left(\eta_i - \frac{\mathbf{z}_i}{2\gamma}\right)^+$$

where η is found so that the constraint $\mathbf{w}_i^T \mathbf{1} = 1$ is satisfied.

Technical details: Derivation of solution

• To derive the optimal solution, setting the gradient of the Lagrangian w.r.t. **w**_i to zero gives

$$\mathbf{w}_i + rac{\mathbf{z}_i}{2\gamma} - \eta_i \mathbf{1} - oldsymbol{eta}_i = \mathbf{0}.$$

• Now, from the complementary slackness condition $w_{ji}\beta_{ji} = 0$, we get that, for $j \neq i$:

• if
$$w_{ji} > 0$$
, then $w_{ji} = -\frac{z_{ji}}{2\gamma} + \eta_i$ and $-\frac{z_{ji}}{2\gamma} + \eta_i > 0$

- if $w_{ji} = 0$, then $-\frac{2ji}{2\gamma} + \eta_i = -\beta_{ji} \le 0$
- In other words:

• if
$$\eta_i - \frac{z_{ji}}{2\gamma} > 0$$
, then $w_{ji} = -\frac{z_i}{2\gamma} + \eta$
• if $\eta_i - \frac{z_{ji}}{2\gamma} \le 0$, then $w_{ji} = 0$.

• More compactly (defining $z_{ii} = \infty$ so that $w_{ii} = 0$):

$$\mathbf{w}_i = \left(\eta_i - \frac{\mathbf{z}_i}{2\gamma}\right)^+$$

Technical details: Obtaining sparsity via γ

- Let's find a choice of γ so that ||w_i||₀ = m with m ≪ p (i.e., so that each node has exactly m neighbors). In fact, we will need a different γ for each w_i.
- Without loss of generality, suppose $z_{1i}, z_{2i}, \ldots, z_{pi}$ are ordered in increasing order.
- The requirement $\|\mathbf{w}_i\|_0 = m$ implies $w_{mi} > 0$ and $w_{m+1,i} = 0$. Therefore, we have

$$\eta_i - rac{z_{mi}}{2\gamma} > 0 \quad ext{and} \quad \eta_i - rac{z_{m+1,i}}{2\gamma} \leq 0$$

• Combining the solution of \mathbf{w}_i with the constraint $\mathbf{w}_i^T \mathbf{1} = 1$, we get

$$\sum_{j=1}^m \left(\eta_i - \frac{z_{ji}}{2\gamma}\right) = 1 \implies \eta_i = \frac{1}{m} + \frac{1}{2m\gamma} \sum_{j=1}^m z_{ji}.$$

• This leads to following inequality for $\gamma:$

$$\frac{m}{2}z_{mi} - \frac{1}{2}\sum_{j=1}^{m} z_{ji} < \gamma \le \frac{m}{2}z_{m+1,i} - \frac{1}{2}\sum_{j=1}^{m} z_{ji}$$

Technical details: Obtaining sparsity via γ

• Therefore, to obtain an optimal solution \mathbf{w}_i with exactly m nonzero values ($\|\mathbf{w}_i\|_0 = m$), the maximal γ is

$$\gamma = \frac{m}{2} z_{m+1,i} - \frac{1}{2} \sum_{j=1}^{m} z_{ji}$$

• Combining the previous results, we get

$$w_{ji} = \begin{cases} \frac{z_{m+1,i}-z_{ji}}{mz_{m+1,i}-\sum_{h=1}^{m} z_{hi}}, & j \le m \\ 0, & j > m. \end{cases}$$

- We now drop the increasing ordering assumption on z_{ji} and we denote by $z_{(m+1),i}$ the (m+1)-th smallest element of $z_{1,i}, \ldots, z_{p,i}$.
- The final solution is written as:

$$\tilde{\mathbf{w}}_i = \left(1 - \frac{\mathbf{z}_i}{z_{(m+1),i}}\right)^+$$
$$\mathbf{w}_i = \tilde{\mathbf{w}}_i / \mathbf{1}^T \tilde{\mathbf{w}}_i.$$

Learning graphs from smooth signals revisited

• Recall the formulation for the graph learning from smooth signals:

```
 \begin{array}{ll} \underset{\mathbf{W}}{\text{minimize}} & \frac{1}{2} \text{Tr}(\mathbf{W}\mathbf{Z}) + \frac{\gamma}{2} \|\mathbf{W}\|_{F}^{2} \\ \text{subject to} & \mathbf{W}\mathbf{1} = \mathbf{1} \\ & \text{diag}(\mathbf{W}) = \mathbf{0}, \quad \mathbf{W} = \mathbf{W}^{T} \geq \mathbf{0}. \end{array}
```

• If we move the hard constraint on the degrees of the nodes W1 = 1 to the objective as a penalty term to avoid the trivial solution W = 0 (with vanishing degrees), we obtain (Kalofolias 2016)⁵:

$$\begin{array}{ll} \underset{\mathbf{W}}{\text{minimize}} & \frac{1}{2}\mathsf{Tr}(\mathbf{W}\mathbf{Z}) - \alpha \mathbf{1}^{\mathcal{T}}\mathsf{log}(\mathbf{W}\mathbf{1}) + \frac{\beta}{2} \|\mathbf{W}\|_{\mathcal{F}}^{2} \\ \text{subject to} & \mathsf{diag}(\mathbf{W}) = \mathbf{0}, \quad \mathbf{W} = \mathbf{W}^{\mathcal{T}} \geq \mathbf{0}. \end{array}$$

• This formulation is convex and can be easily solved with a solver. Interestingly, we can <u>again develop a simple tailored algorithm</u>.

⁵V. Kalofolias, "How to learn a graph from smooth signals," in *Proc. Int. Conf. Artif. Intell. Statist.*, 2016, pp. 920–929.

Learning graphs from smooth signals revisited

If we ignore the symmetry constraint in the adjacency matrix W = W^T, then the problem can be nicely solved separately for each row/column w_i:

$$\begin{array}{ll} \underset{\mathbf{w}_{i}}{\text{minimize}} & \frac{1}{2}\mathbf{z}_{i}^{T}\mathbf{w}_{i} - \alpha \, \log(\mathbf{1}^{T}\mathbf{w}_{i}) + \frac{\beta}{2} \|\mathbf{w}_{i}\|^{2} \\ \text{subject to} & w_{ii} = 0, \ \mathbf{w}_{i} \geq \mathbf{0}. \end{array}$$

• Note that at a later stage we will have to symmetrize it:

$$\mathbf{W} \leftarrow rac{1}{2}(\mathbf{W} + \mathbf{W}^{\mathcal{T}})$$

• From the KKT optimality conditions, the optimal solution (defining $z_{ii} = \infty$ so that $w_{ii} = 0$) is

$$\mathbf{w}_i = \left(\frac{\alpha/\beta}{t_i} - \frac{\mathbf{z}_i}{2\beta}\right)^+$$

where t_i is found so that $\mathbf{1}^T \mathbf{w}_i = t_i$.

Technical details: Derivation of solution

• The Lagrangian of the problem is

$$\mathcal{L}(\mathbf{w}_i; \lambda_i) = \frac{1}{2} \mathbf{z}_i^T \mathbf{w}_i - \alpha \log(\mathbf{1}^T \mathbf{w}_i) + \frac{\beta}{2} \|\mathbf{w}_i\|^2 - \boldsymbol{\lambda}_i^T \mathbf{w}_i$$

where $\lambda_i \in \mathbb{R}^p$ is the Lagrangian multipliers with $\lambda_{ji} \ge 0, \forall j \neq i$.

- Setting the gradient of the Lagrangian w.r.t. \mathbf{w}_i to zero gives $\frac{1}{2}\mathbf{z}_i \frac{\alpha}{\mathbf{1}^T \mathbf{w}_i} \mathbf{1} + \beta \mathbf{w}_i \lambda_i = \mathbf{0}$.
- Now, from the complementary slackness condition $w_{ji}\lambda_{ji} = 0$, we get that, for $j \neq i$.

• if
$$w_{ji} > 0$$
, then $w_{ji} = \frac{\alpha/\beta}{1^{T}w_i} - \frac{z_{ji}}{2\beta}$ and $\frac{\alpha/\beta}{1^{T}w_i} > \frac{z_{ji}}{2\beta}$
• if $w_{ji} = 0$, then $\frac{z_{ji}}{2} - \frac{\alpha}{1^{T}w_i} = \lambda_{ji} \ge 0$

In other words:

• if
$$\frac{\alpha/\beta}{\mathbf{1}^T \mathbf{w}_i} > \frac{z_{ji}}{2\beta}$$
, then $w_{ji} = \frac{\alpha/\beta}{\mathbf{1}^T \mathbf{w}_i} - \frac{z_{ji}}{2\beta}$
• if $\frac{\alpha/\beta}{\mathbf{1}^T \mathbf{w}_i} \le \frac{z_{ji}}{2\beta}$, then $w_{ji} = 0$.

• More compactly (defining $z_{ii} = \infty$ so that $w_{ii} = 0$): $\mathbf{w}_i = \left(\frac{\alpha/\beta}{1^T \mathbf{w}_i} - \frac{\mathbf{z}_i}{2\beta}\right)^+$.

Technical details: Obtaining sparsity via α and β

- Let's find a choice of α and β so that $\|\mathbf{w}_i\|_0 = m$ with $m \ll p$ (i.e., so that each node has exactly *m* neighbors). In fact, we will need a different α and β for each \mathbf{w}_i .
- Without loss of generality, suppose $z_{1i}, z_{2i}, \ldots, z_{pi}$ are ordered in increasing order.
- The requirement $\|\mathbf{w}_i\|_0 = m$ implies $w_{m,i} > 0$ and $w_{m+1,i} = 0$. Therefore, we have

$$rac{lpha/eta}{\mathbf{1}^T \mathbf{w}_i} > rac{\mathbf{z}_{m,i}}{2eta} \quad ext{and} \quad rac{lpha/eta}{\mathbf{1}^T \mathbf{w}_i} \leq rac{\mathbf{z}_{m+1,i}}{2eta}.$$

• Denoting $\bar{\alpha}=\alpha/\beta,$ we can write it as

$$\frac{z_{m,i}}{2\bar{\alpha}}\mathbf{1}^{\mathsf{T}}\mathbf{w}_{i} < \beta \leq \frac{z_{m+1,i}}{2\bar{\alpha}}\mathbf{1}^{\mathsf{T}}\mathbf{w}_{i}.$$

• Therefore, to obtain an optimal solution \mathbf{w}_i with exactly *m* nonzero values ($\|\mathbf{w}_i\|_0 = m$), the maximal β is

$$\beta = \frac{z_{m+1,i}}{2\bar{\alpha}} \mathbf{1}^T \mathbf{w}_i.$$
Technical details: Obtaining sparsity via α and β

• Combining the previous results, the optimal solution is

$$w_{ji} = rac{ar{lpha}}{\sum_{j=1}^m w_{ji}} \left(1 - rac{z_{ji}}{z_{m+1,i}}
ight)^+$$

• Now we denote $t_i = \sum_{j=1}^m w_{ji}$ and write

$$w_{ji} = \frac{\bar{\alpha}}{t_i} \left(1 - \frac{z_{ji}}{z_{m+1,i}} \right)^+$$

which leads to

$$t_i = \sum_{j=1}^m w_{ji} = \frac{\bar{\alpha}}{t_i} \sum_{j=1}^p \left(1 - \frac{z_{ji}}{z_{m+1,i}} \right)^+$$

or

$$t_i = \sqrt{\bar{\alpha} \sum_{j=1}^p \left(1 - \frac{z_{ji}}{z_{m+1,i}}\right)^+}.$$

D. Palomar (HKUST)

Graph Learning

Technical details: Obtaining sparsity via α and β

• The final solution for **w**_i can be compactly written as

$$\mathbf{w}_i = rac{ar{lpha}}{t_i} \left(1 - rac{\mathbf{z}_i}{z_{m+1,i}}
ight)^+$$

with

$$t_i = \sqrt{\bar{\alpha} \sum_{j=1}^{p} \left(1 - \frac{z_{ji}}{z_{m+1,i}}\right)^+}$$

- We now drop the increasing ordering assumption on z_{ji} and we denote by $z_{(m+1),i}$ the (m+1)-th smallest element of $z_{1,i}, \ldots, z_{p,i}$.
- The final solution is written as

$$\begin{split} \tilde{\mathbf{w}}_i &= \bar{\alpha} \left(1 - \frac{\mathbf{z}_i}{z_{(m+1),i}} \right)^+ \\ \mathbf{w}_i &= \tilde{\mathbf{w}}_i / \sqrt{\mathbf{1}^T \tilde{\mathbf{w}}_i}. \end{split}$$

• Recall the original convex formulation:

$$\begin{array}{ll} \underset{\mathbf{W}}{\text{minimize}} & \frac{1}{2} \mathsf{Tr}(\mathbf{W}\mathbf{Z}) - \alpha \mathbf{1}^{\mathcal{T}} \mathsf{log}(\mathbf{W}\mathbf{1}) + \frac{\beta}{2} \|\mathbf{W}\|_{\mathcal{F}}^{2} \\ \text{subject to} & \mathsf{diag}(\mathbf{W}) = \mathbf{0}, \quad \mathbf{W} = \mathbf{W}^{\mathcal{T}} \geq \mathbf{0}. \end{array}$$

- Instead of solving it with a solver or with the previously derived columnwise closed-form solution (which ignores the symmetry constraint), we will attempt to derive a joint closed-form solution.
- \bullet Instead of dealing with matrix ${\bf W},$ denote with ${\bf w}$ the off-diagonal elements stacked in a vector.
- Define \mathbf{z} such that $Tr(\mathbf{W}\mathbf{Z}) = \mathbf{z}^T \mathbf{w}$.
- Also, define **G** to extract the row sums of **W**: $W1 = G^T w$
- Note that $\frac{1}{2} \| \mathbf{W} \|_{F}^{2} = \| \mathbf{w} \|^{2}$.
- The constraints $diag(\mathbf{W}) = \mathbf{0}$ and $\mathbf{W} = \mathbf{W}^T$ are unnecessary.

• The problem can be rewritten in terms of \mathbf{w} as (for aesthetics we change two 1/2 factors)

$$\begin{array}{ll} \underset{\mathbf{w}}{\text{minimize}} & \mathbf{z}^{T}\mathbf{w} - \alpha \mathbf{1}^{T} \text{log}(\mathbf{G}^{T}\mathbf{w}) + \frac{\beta}{2} \|\mathbf{w}\|^{2} \\ \text{subject to} & \mathbf{w} \geq \mathbf{0}. \end{array}$$

• The Lagrangian of the problem is

$$\mathcal{L}(\mathbf{w}; \lambda) = \mathbf{z}^T \mathbf{w} - \alpha \mathbf{1}^T \log(\mathbf{G}^T \mathbf{w}) + \frac{\beta}{2} \|\mathbf{w}\|^2 - \boldsymbol{\lambda}^T \mathbf{w}$$

where $\lambda \in \mathbb{R}^{p}_{+}$.

• Setting the gradient of the Lagrangian w.r.t. w to zero gives

$$\mathbf{z} - \alpha \sum_{j=1}^{p} \frac{\mathbf{g}_j}{\mathbf{g}_j^T \mathbf{w}} + \beta \mathbf{w} - \boldsymbol{\lambda} = \mathbf{0}.$$

- Now, from the complementary slackness condition $w_i \lambda_i = 0$, we get:
 - if $w_i > 0$, then $w_i = \frac{\alpha}{\beta} \sum_{j=1}^{p} \frac{g_{ij}}{\mathbf{g}_i^T \mathbf{w}} \frac{z_i}{\beta}$ and $\frac{\alpha}{\beta} \sum_{j=1}^{p} \frac{g_{ij}}{\mathbf{g}_i^T \mathbf{w}} > \frac{z_i}{\beta}$

• if
$$w_i = 0$$
, then $z_i - \alpha \sum_{j=1}^{p} \frac{g_{ij}}{\mathbf{g}_j' \mathbf{w}} = \lambda_i \ge 0$.

• In other words:

• if
$$\frac{\alpha}{\beta} \sum_{j=1}^{p} \frac{\mathbf{g}_{ij}}{\mathbf{g}_{j}^{T}\mathbf{w}} > \frac{\mathbf{z}_{i}}{\beta}$$
, then $w_{i} = \frac{\alpha}{\beta} \sum_{j=1}^{p} \frac{\mathbf{g}_{ij}}{\mathbf{g}_{j}^{T}\mathbf{w}} - \frac{\mathbf{z}_{i}}{\beta}$
• if $\frac{\alpha}{\beta} \sum_{j=1}^{p} \frac{\mathbf{g}_{ij}}{\mathbf{g}_{i}^{T}\mathbf{w}} \le \frac{\mathbf{z}_{i}}{\beta}$, then $w_{ji} = 0$.

• More compactly (defining $z_{ii} = \infty$ so that $w_{ii} = 0$):

$$\mathbf{w} = \left(rac{lpha}{eta} \sum_{j=1}^{p} rac{\mathbf{g}_{j}}{\mathbf{g}_{j}^{T} \mathbf{w}} - rac{\mathbf{z}}{eta}
ight)^{+}$$

• This expression is a fixed-point equation in **w**, which is not easy to compute in practice, although algorithms can be devised.

Denoting explicitly the degrees of the nodes by d_i = g_i^Tw, we can write a double fixed-point equation in (d, w):

$$\begin{split} \mathbf{d} &= \mathbf{G}^T \mathbf{w} \\ \mathbf{w} &= \left(\frac{\alpha}{\beta} \mathbf{G} \mathbf{d}^{-1} - \frac{\mathbf{z}}{\beta}\right)^+ \end{split}$$

• We can also write the fixed-point equation in terms of **d** only:

$$\mathbf{d} = \mathbf{G}^{\mathcal{T}} \left(\frac{\alpha}{\beta} \mathbf{G} \mathbf{d}^{-1} - \frac{\mathbf{z}}{\beta} \right)^+$$

• Alternatively, one would devise a gradient projection method:

$$\mathbf{w}^{k+1} = \left(\mathbf{w}^{k} - \mu^{k} \left(\mathbf{z} - \alpha \mathbf{G} (\mathbf{G}^{\mathsf{T}} \mathbf{w}^{k})^{-1} + \beta \mathbf{w}^{k} \right)\right)^{+}$$

where μ^k is the step size.

• To avoid numerical issues in some degree d_i becoming zero, one can use the heuristic: $\mathbf{w}^{k+1} = \left(\mathbf{w}^k - \mu^k \left(\mathbf{z} - \alpha \mathbf{G} (\mathbf{G}^T \mathbf{w}^k + \epsilon)^{-1} + \beta \mathbf{w}^k\right)\right)^+$

D. Palomar (HKUST)

Graph Learning

Learn a clean graph from a noisy graph by imposing k-connected structure:



Eigenvalue property of Laplacian matrix

• Consider the eigenvalue decomposition of the Laplacian matrix:

$$\mathbf{L} = \mathbf{U} \mathsf{D} \mathsf{iag}(\lambda_1, \lambda_2, \dots, \lambda_p) \mathbf{U}^T,$$

where **U** contains the eigenvectors columnwise and $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_p$ are the eigenvalues in increasing order.

• For a k-connected graph, the k smallest eigenvalues are zero:



- **Goal**: Given an initial noisy adjacency matrix \mathbf{W}_0 , infer a *k*-connected graph.
- In other words, we want to learn **W** as close as possible to **W**₀ but satisfying the property of a *k*-connected graph: zero *k* smallest eigenvalues of the corresponding Laplacian matrix.
- Define the Laplacian operator: L(W) = Diag(W1) W.
- Rank constrained nonconvex formulation (Nie et al. 2016)⁶:

$$\begin{array}{ll} \underset{\mathbf{W}}{\text{minimize}} & \|\mathbf{W} - \mathbf{W}_0\|_F^2 \\ \text{subject to} & \mathbf{W}\mathbf{1} = \mathbf{1}, \ \text{diag}(\mathbf{W}) = \mathbf{0}, \ \mathbf{W} = \mathbf{W}^T \geq \mathbf{0} \\ & \text{rank}(\mathbf{L}(\mathbf{W})) = p - k. \end{array}$$

⁶F. Nie, X. Wang, M. Jordan, and H. Huang, "The constrained Laplacian rank algorithm for graph-based clustering," in *Proc. of the Thirtieth Conference on Artificial Intelligence (AAAI)*, Phoenix, Arizona, USA, 2016, pp. 1969–1976.

If we relax the low-rank constraint λ₁(L(W)) = ··· = λ_k(L(W)) = 0, we can rewrite the problem as

$$\begin{array}{ll} \underset{\mathbf{W}}{\text{minimize}} & \|\mathbf{W} - \mathbf{W}_0\|_F^2 + \beta \sum_{i=1}^k \lambda_i(\mathbf{L}(\mathbf{W})) \\ \text{subject to} & \mathbf{W}\mathbf{1} = \mathbf{1}, \text{ diag}(\mathbf{W}) = \mathbf{0}, \ \mathbf{W} = \mathbf{W}^T \ge \mathbf{0} \end{array}$$

• We now use the variational interpretation of the smallest eigenvalues known as Ky Fan's theorem (Fan 1949):

$$\sum_{i=1}^{k} \lambda_i(\mathbf{X}) = \min_{\mathbf{F} \in \mathbb{R}^{p \times k}} \quad \text{Tr}(\mathbf{F}^T \mathbf{X} \mathbf{F}) \quad \text{s.t.} \quad \mathbf{F}^T \mathbf{F} = \mathbf{I}.$$

• The problem can then be rewritten in a more manageable form (still nonconvex) as

$$\begin{split} \underset{\mathbf{W},\mathbf{F}}{\text{minimize}} & \|\mathbf{W}-\mathbf{W}_0\|_F^2 + \beta \mathsf{Tr}(\mathbf{F}^T \mathbf{L}(\mathbf{W})\mathbf{F}) \\ \text{subject to} & \mathbf{W}\mathbf{1} = \mathbf{1}, \ \mathsf{diag}(\mathbf{W}) = \mathbf{0}, \ \mathbf{W} = \mathbf{W}^T \geq \mathbf{0} \\ & \mathbf{F}^T \mathbf{F} = \mathbf{I}. \end{split}$$

- To solve this nonconvex problem in the variables (**W**, **F**), we will use the block coordinate descent (BCD) method, which in this particular case with two blocks it becomes a simple sequential optimization with respect to **W** and **F** alternatively.
- Optimization w.r.t. **F**:

minimize
$$Tr(F^{T}L(W)F)$$

subject to $F^{T}F = I$.

whose solution is trivially given by the eigenvectors corresponding to the k smallest eigenvalues of L(W).

• Optimization w.r.t. W is a quadratic problem (QP):

minimize
W
$$\|\mathbf{W} - \mathbf{W}_0\|_F^2 + \beta \operatorname{Tr}(\mathbf{F}^T \mathbf{L}(\mathbf{W})\mathbf{F})$$
subject to
$$\mathbf{W} \mathbf{1} = \mathbf{1}, \operatorname{diag}(\mathbf{W}) = \mathbf{0}, \ \mathbf{W} = \mathbf{W}^T \ge \mathbf{0}.$$

• From the property of the Laplacian matrix we can write:

$$\operatorname{Tr}(\mathbf{F}^{T}\mathbf{L}(\mathbf{W})\mathbf{F}) = \frac{1}{2}\sum_{i,j}w_{ij}\|\mathbf{f}_{i}-\mathbf{f}_{j}\|^{2}$$

• The problem w.r.t. \boldsymbol{W} can then be rewritten as

$$\begin{array}{ll} \underset{\mathbf{W}}{\text{minimize}} & \|\mathbf{W} - \mathbf{W}_0\|_F^2 + \frac{\beta}{2} \sum_{i,j} w_{ij} \|\mathbf{f}_i - \mathbf{f}_j\|^2\\ \text{subject to} & \mathbf{W}\mathbf{1} = \mathbf{1}, \text{ diag}(\mathbf{W}) = \mathbf{0}, \ \mathbf{W} = \mathbf{W}^T \geq \mathbf{0}. \end{array}$$

• If we ignore the symmetry constraint in the adjacency matrix $\mathbf{W} = \mathbf{W}^{T}$, then the problem can be nicely solved separately for each row/column \mathbf{w}_{i} as in (Nie et al. 2016)⁷:

$$\begin{array}{ll} \underset{\mathbf{w}_{i}}{\text{minimize}} & \|\mathbf{w}_{i} - \mathbf{w}_{0,i}\|^{2} + \frac{\beta}{2}\mathbf{w}_{i}^{T}\mathbf{v}_{i} \\ \text{subject to} & \mathbf{w}_{i}^{T}\mathbf{1} = 1, \ w_{ii} = 0, \ \mathbf{w}_{i} \geq \mathbf{0}, \end{array}$$

⁷F. Nie, X. Wang, M. Jordan, and H. Huang, "The constrained Laplacian rank algorithm for graph-based clustering," in *Proc. of the Thirtieth Conference on Artificial Intelligence (AAAI)*, Phoenix, Arizona, USA, 2016, pp. 1969–1976.

D. Palomar (HKUST)

where $v_{ii} = \|\mathbf{f}_{i} - \mathbf{f}_{i}\|^{2}$.

• The problem can be finally rewritten as

$$\begin{array}{ll} \underset{\mathbf{w}_{i}}{\text{minimize}} & \left\|\mathbf{w}_{i} - \left(\mathbf{w}_{0,i} - \frac{\beta}{2}\mathbf{v}_{i}\right)\right\|^{2} \\ \text{subject to} & \mathbf{w}_{i}^{T}\mathbf{1} = 1, \ w_{ii} = 0, \ \mathbf{w}_{i} \geq \mathbf{0}, \end{array}$$

where $v_{ij} = \|\mathbf{f}_i - \mathbf{f}_j\|^2$.

• We have seen this problem before and it has the closed-form solution for $j \neq i$ (recall $w_{ii} = 0$)

$$w_{ji} = \left(\mathbf{w}_{0,i} - \frac{\beta}{2}\mathbf{v}_i + \eta_i\right)^{-1}$$

where η is found so that the constraint $\mathbf{w}_i^T \mathbf{1} = 1$ is satisfied.

Graphs based on smooth signals



Outline

Graphs

2 Basics

3 Learning Graphs from Data

- Similarity function based
- Smooth signal based
- i.i.d. model based
- Structured graphs via spectral constraints

With a set of the set

Notation for time series

- Our convention for graphs is to form the data matrix as X = [x₁,..., x_n] ∈ ℝ^{p×n}, where p is the number of nodes and n is the number of observations per node. Note that each x_i denotes one graph observation.
- In the context of time series, we usually denote the number of series by N and the number of time observations by T (with observation along rows instead of columns). So the correspondence of the two conventions is
 - p = N: number of nodes/series
 - *n* = *T*: number of observations/features
- The way to construct our graph data matrix X from a time series is by transposing the usual time series matrix (where each column represents the T observations of one series). In our graph notation, we want each series along a row of X and each column is one vector observation, so that X is N × T (p × n):

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{N \times T}.$$

Modeling as i.i.d. Random Variables

- Until now the data matrix **X** was composed of some (columnwise) observations for each of the nodes without any statistical modeling.
- Examples of observations **x**_i:
 - the x-y coordinates of some points (e.g., two-moon dataset, three-circle dataset)
 - feature vector with arbitrary features (e.g., the animal dataset, where the features include attributed such as whether the animal is a mammal, has wings, has feathers, weight, size, etc.)
- Now we will assume a **statistical model** where the observations are multivariate i.i.d. random variables and each multivariate observation **x**_i is distributed according to some multivariate distribution.
- Example: Gaussian distribution

$$\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

where μ is the mean and $\pmb{\Sigma}$ the covariance matrix of the observations.

Graph from covariance matrix

- If our dataset is a collection of multivariate i.i.d. random variables {x_i}, we may look for more interesting similarity measures.
- Let Σ = E[x_ix_i^T] be the covariance matrix of the multivariate random i.i.d. variables. The correlation matrix is a normalized version (equivalent to normalizing the variance of each stock):

$$\mathbf{C} = \mathbf{D}^{-1/2} \mathbf{\Sigma} \mathbf{D}^{-1/2}$$

where $\mathbf{D} = \text{Diag}(\mathbf{\Sigma})$.

- The correlation matrix can be used for graph construction (Heimo et al. $2007)^8$.
- We can consider a modification of the correlation matrix as an adjacency matrix:

$$W = |C| - I$$

so that \mathbf{W} has nonnegative elements and zero along the diagonal.

⁸T. Heimo, J. Saramaki, J.-P. Onnela, and K. Kaski, "Spectral and network methods in the analysis of correlation matrices of stock returns," *Physica A: Statistical Mechanics and its Applications*, vol. 383, no. 1, pp. 147–151, 2007.

Graph from precision matrix

- The covariance/correlation matrix measures the direct dependency between two nodes but ignores the other nodes.
- Can we improve that? Answer: conditional dependence graph.
- Definition: We say that two random variables X and Y are conditionally independent given a third variable Z if their conditional probability distributions given Z are independent. We denote it by X⊥ Y | Z.
 - **Example**: Height and vocabulary of kids are not independent, but they are conditionally independent if you also consider age.
- **Definition**: The **precision matrix** is defined as $\Theta = \Sigma^{-1}$.
- **Theorem**: Suppose our random variables are Gaussian with zero mean and covariance **Σ**. Then

$$\boldsymbol{\Theta}_{ij} = 0 \Longleftrightarrow X_i \perp X_j \mid \{X_l\}_{l \neq i,j}.$$

• In words: the non-zero entries of Θ indicate **conditional dependence** between the two random variables given all the other random variables.

Historical timeline of Markov graphical models

- Suppose the columns of the data matrix **X** follow $x_i \sim \mathcal{N}(\mu, \Sigma)$ and let **S** be the sample covariance matrix.
- Covariance selection (Dempster 1972): graph from the elements of **S**⁻¹ inverse sample covariance matrix.
- Neighborhood regression (Meinshausen and Bühlmann 2006):

$$\arg\min_{\boldsymbol{\beta}_1} |\mathbf{x}^{(1)} - \boldsymbol{\beta}_1^T \mathbf{X}_{/\mathbf{x}^{(1)}}|^2 + \alpha \|\boldsymbol{\beta}\|_1.$$

• ℓ_1 -regularized MLE (Banerjee et al. 2008; Friedman et al. 2008):

$$\underset{\boldsymbol{\Theta}\succ\mathbf{0}}{\operatorname{maximize}} \log \det(\boldsymbol{\Theta}) - \operatorname{Tr}(\boldsymbol{\Theta}\mathbf{S}) - \alpha \|\boldsymbol{\Theta}\|_{1}.$$

- Ising model: ℓ_1 -regularized logistic regression (Ravikumar et al. 2010).
- Attractive IGMRF (Slawski and Hein 2015).
- Laplacian structure in Θ (Lake and Tenenbaum 2010).
- ℓ_1 -regularized MLE with Laplacian structure (Egilmez et al. 2017; Zhao et al. 2019).

Graphical LASSO (GLASSO)

- In real life, interactions are typically local and sparse.
- Normally the ideal precision matrix will be sparse.
- But the sample covariance matrix is **noisy**.
- Graphical LASSO tries to learn a sparse precision matrix:

$$\max_{\substack{\boldsymbol{\Theta}\succ\boldsymbol{0}\\\boldsymbol{\Theta}\succ\boldsymbol{0}}} \operatorname{log} \det\left(\boldsymbol{\Theta}\right) - \mathsf{Tr}\left(\boldsymbol{S}\boldsymbol{\Theta}\right) - \rho \|\boldsymbol{\Theta}\|_{1, \mathsf{off}}$$

where **S** is the sample covariance matrix and $\|\cdot\|_{1,off}$ denotes the ℓ_1 elementwise of the off-diagonal elements.

← If $\rho = 0$, the solution is given precisely by the inverse sample covariance matrix $\Theta = \mathbf{S}^{-1}$.

• Seminal reference (Friedman et al. 2008)⁹.

⁹J. Friedman, T. Hastie, and R. Tibshirani, "Sparse inverse covariance estimation with the graphical lasso," *Biostatistics*, vol. 9, no. 3, pp. 432–441, 2008.

Derivation of GLASSO algorithm

• The optimality condition for the problem

$$\begin{array}{l} \underset{\boldsymbol{\Theta} \succ \boldsymbol{0}}{\text{maximize}} \quad \log \det \left(\boldsymbol{\Theta} \right) - \mathsf{Tr} \left(\boldsymbol{S} \boldsymbol{\Theta} \right) - \rho \| \boldsymbol{\Theta} \|_{1, \text{off}} \end{array}$$

is

$$\Theta^{-1} - \mathbf{S} - \alpha \mathbf{\Gamma} = \mathbf{0},$$

where $\mathbf{\Gamma}$ is a matrix of component-wise signs of $\boldsymbol{\Theta}$:

$$[\mathbf{\Gamma}]_{jk} = \begin{cases} \gamma_{jk} = \operatorname{sign}(\Theta_{jk}), & \text{if } \Theta_{jk} \neq 0\\ \gamma_{jk} \in [-1, 1], & \text{if } \Theta_{jk} \neq 0. \end{cases}$$

• The equation for optimality condition is also known as the normal equation.

• Furthermore, the constraint requires Θ_{jj} to be positive, this implies that

$$\hat{\Sigma}_{ii}=S_{ii}+\alpha, \quad i=1,\ldots,p,$$

where $\hat{\boldsymbol{\Sigma}} = \boldsymbol{\Theta}^{-1}$

D. Palomar (HKUST)

Derivation of GLASSO algorithm

- GLASSO uses a block-coordinate method for solving the problem.
- Consider a partitioning of Θ and $\hat{\Sigma}:$

$$\boldsymbol{\Theta} = \left(\begin{array}{cc} \boldsymbol{\Theta}_{11} & \boldsymbol{\theta}_{12} \\ \boldsymbol{\theta}_{12}^{\mathsf{T}} & \boldsymbol{\theta}_{22} \end{array} \right), \quad \hat{\boldsymbol{\Sigma}} = \left(\begin{array}{cc} \hat{\boldsymbol{\Sigma}}_{11} & \hat{\boldsymbol{\sigma}}_{12} \\ \hat{\boldsymbol{\sigma}}_{12}^{\mathsf{T}} & \hat{\boldsymbol{\sigma}}_{22} \end{array} \right),$$

where $\Theta_{11} \in \mathbb{R}^{(p-1) \times (p-1)}$, $\hat{\theta}_{12} \in \mathbb{R}^{p-1}$ and θ_{22} is a scalar, and similarly for the other partitions.

 \bullet Then, $\hat{\Sigma} = \Theta^{-1} \; (\Theta \hat{\Sigma} = I)$ can be expressed as

$$\hat{\boldsymbol{\Sigma}} = \begin{pmatrix} \boldsymbol{\Theta}_{11}^{-1} + \frac{\boldsymbol{\Theta}_{11}^{-1} \boldsymbol{\theta}_{12} \boldsymbol{\theta}_{12}^{T} \boldsymbol{\Theta}_{11}^{-1}}{\boldsymbol{\theta}_{22} - \boldsymbol{\theta}_{12}^{T} \boldsymbol{\Theta}_{11}^{-1} \boldsymbol{\theta}_{12}} & \frac{\boldsymbol{\Theta}_{11}^{-1} \boldsymbol{\theta}_{12}}{\boldsymbol{\theta}_{22} - \boldsymbol{\theta}_{12}^{T} \boldsymbol{\Theta}_{11}^{-1} \boldsymbol{\theta}_{12}} \\ \cdot & \frac{1}{\boldsymbol{\theta}_{22} - \boldsymbol{\theta}_{12}^{T} \boldsymbol{\Theta}_{11}^{-1} \boldsymbol{\theta}_{12}} \end{pmatrix}$$

• GLASSO solves for a row/column of Θ at a time, holding the rest fixed. Considering the *p*th column of the normal equation, we get

$$-\hat{\boldsymbol{\sigma}}_{12}+\mathbf{s}_{12}+\alpha\boldsymbol{\gamma}_{12}=\mathbf{0}.$$

D. Palomar (HKUST)

Graph Learning

Derivation of GLASSO algorithm

• Consider reading off $heta_{12}$ from the partitioned expression:

$$\frac{\boldsymbol{\Theta}_{11}^{-1}\boldsymbol{\theta}_{12}}{\boldsymbol{\theta}_{22} - \boldsymbol{\theta}_{12}^{\mathsf{T}}\boldsymbol{\Theta}_{11}^{-1}\boldsymbol{\theta}_{12}} + \mathbf{s}_{12} + \alpha \boldsymbol{\gamma}_{12} = \mathbf{0}.$$

• The above also simplifies to

$$\boldsymbol{\Theta}_{11}^{-1}\boldsymbol{\theta}_{12}\hat{\sigma}_{22} + \mathbf{s}_{12} + \alpha\boldsymbol{\gamma}_{12} = \mathbf{0}$$

with $\nu = \theta_{12}\hat{\sigma}_{22}$ (with $\hat{\sigma}$ fixed), $\Theta_{11} \succ \mathbf{0}$ is equivalent to the stationary condition for (Mazumder and Hastie 2012)

$$\underset{\boldsymbol{\nu} \in \mathbb{R}^{p-1}}{\text{minimize}} \quad \frac{1}{2} \boldsymbol{\nu}^{T} \boldsymbol{\Theta}_{11}^{-1} \boldsymbol{\nu} + \boldsymbol{\nu}^{T} \mathbf{s}_{12} + \alpha \|\boldsymbol{\nu}\|_{1}$$

• Let u^{\star} be the minimizer, then

$$egin{aligned} oldsymbol{ heta}_{12}^{\star} &= oldsymbol{
u}^{\star} \hat{\sigma}_{22} \ oldsymbol{ heta}_{22}^{\star} &= rac{1}{\hat{\sigma}_{22}} + (oldsymbol{ heta}_{12})^T oldsymbol{\Theta}_{11}^{-1} oldsymbol{ heta}_{12}^{\star} \end{aligned}$$

Algorithm Graphical LASSO

initialize: $\hat{\boldsymbol{\Sigma}} = \text{Diag}(\boldsymbol{S}) + \alpha \boldsymbol{I}$ and $\boldsymbol{\Theta} = \hat{\boldsymbol{\Sigma}}^{-1}$. repeat

• Rearrange rows and columns such that the target column is the last.

• Compute
$$oldsymbol{\Theta}_{11}^{-1} = \hat{oldsymbol{\Sigma}}_{11} - rac{\hat{\sigma}_{12}\hat{\sigma}_{12}^{ op}}{\hat{\sigma}_{22}}$$

- Obtain $\boldsymbol{\nu}$ and update $\boldsymbol{\theta}_{12}^{\star}$ and $\boldsymbol{\theta}_{22}^{\star}.$
- \bullet Update Θ and $\hat{\Sigma}$ using the second partition function, ensuring $\Theta\hat{\Sigma}=I.$

until convergence

return Θ

Laplacian-structured GLASSO

• We can compute a precision-like matrix with the **Laplacian structure** by adding the following constraints:

$$\mathbf{L} \succeq \mathbf{0}, \ \mathbf{L}\mathbf{1} = \mathbf{0}, \ L_{ij} = L_{ji} \leq 0, \quad \forall i \neq j.$$

• Laplacian-structured GLASSO: precision estimation with Laplacian constraints and ℓ_1 -norm regularization:

$$\begin{array}{ll} \underset{\mathbf{L}\succeq\mathbf{0}}{\text{maximize}} & \log \operatorname{gdet}\left(\mathbf{L}\right) - \operatorname{Tr}\left(\mathbf{SL}\right) - \rho \|\mathbf{L}\|_{1,\operatorname{off}} \\ \text{subject to} & \mathbf{L}\mathbf{1} = \mathbf{0}, \ L_{ij} = L_{ji} \leq 0, \quad \forall i \neq j \end{array}$$

• Basic references (Lake and Tenenbaum 2010)¹⁰, (Egilmez et al. 2017)¹¹, and (Zhao et al. 2019)¹².

¹⁰B. Lake and J. Tenenbaum, "Discovering structure by learning sparse graphs," in *Proc. the 33rd Annual Cognitive Science Conference*, 2010.

¹¹H. E. Egilmez, E. Pavez, and A. Ortega, "Graph learning from data under Laplacian and structural constraints," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 6, pp. 825–841, 2017. ¹²L. Zhao, Y. Wang, S. Kumar, and D. P. Palomar, "Optimization algorithms for graph Laplacian estimation via ADMM and MM," *IEEE Trans. on Signal Processing*, vol. 67, no. 16, pp. 4231–4244, 2019.

D. Palomar (HKUST)

Graphs based on the i.i.d. model





GLASSO graph (three-circle data)



Laplacian-structured GLASSO graph (two-moon data)



Laplacian-structured GLASSO graph (three-circle data)



Outline

1 Graphs

2 Basics

3 Learning Graphs from Data

- Similarity function based
- Smooth signal based
- i.i.d. model based
- Structured graphs via spectral constraints

4 Numerical Experiments

Structured graphs



D. Palomar (HKUST)

Graph Learning

Structured graphs: Importance and challenges

Useful graph structures:

- Multi-component: for clustering, classification.
- Bipartite: for matching and constructing two-channel filter banks.
- Multi-component bipartite: for co-clustering.
- Tree: for sampling algorithms.
- Modular: for social network analysis.
- Connected sparse: for graph signal processing applications.

Structured graph learning from data

- involves both the estimation of structure (graph connectivity) and parameters (graph weights),
- parameter estimation is well explored (e.g., maximum likelihood),
- but structure is a combinatorial property which makes structure estimation very challenging.

Structure learning is NP-hard for a general class of graphical models (Bogdanov et al. 2008).

Eigenvalue property of Laplacian matrix

• Consider the eigenvalue decomposition of the Laplacian matrix:

$$\mathbf{L} = \mathbf{U} \mathsf{D} \mathsf{iag}(\lambda_1, \lambda_2, \dots, \lambda_p) \mathbf{U}^T,$$

where **U** contains the eigenvectors columnwise and $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_p$ are the eigenvalues in increasing order.

• For a **multi-component graph** (aka *k*-connected graph), the *k* smallest eigenvalues are zero:

$$\lambda_1 = \cdots = \lambda_k = 0.$$



Eigenvalue property of adjacency matrix

• Consider the eigenvalue decomposition of the adjacency matrix:

$$\mathbf{W} = \mathbf{V} \mathsf{Diag}(\psi_1, \psi_2, \dots, \psi_p) \mathbf{V}^T,$$

where **V** contains the eigenvectors columnwise and $\psi_1, \psi_2, \ldots, \psi_p$ are the eigenvalues in increasing order.

• For a bipartite graph, the eigenvalues are symmetric:

$$\psi_i = -\psi_{p-i} \quad \forall i$$



D. Palomar (HKUST)

Graph Learning

Laplacian with spectral constraints

- The Laplacian-structured GLASSO is useful, but in many practical situations, graphs of more complex structures need to be estimated, e.g., *k*-component graph, bipartite graph, etc.
- $\bullet\,$ Such classes of graphs can be enforced via spectral constraints on the graph matrices W and L.
- The (nonconvex) formulation of the Laplacian estimation with spectral constraints is

$$\begin{array}{ll} \underset{\mathbf{L}\succeq\mathbf{0}}{\text{maximize}} & \log \operatorname{gdet}\left(\mathbf{L}\right) - \operatorname{Tr}\left(\mathbf{SL}\right) - \rho h(\mathbf{L}) \\ \text{subject to} & \mathbf{L}\mathbf{1} = \mathbf{0}, \ L_{ij} = L_{ji} \leq 0, \quad \forall i \neq j \\ & \boldsymbol{\lambda}(\mathbf{L}) \in \mathcal{S}_{\boldsymbol{\lambda}}. \end{array}$$

• References: (Kumar et al. 2019)¹³ and (Kumar et al. 2020)¹⁴.

¹³S. Kumar, J. Ying, J. V. de M. Cardoso, and D. P. Palomar, "Structured graph learning via laplacian spectral constraints," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, Vancouver, Canada, 2019.

¹⁴S. Kumar, J. Ying, J. V. de M. Cardoso, and D. P. Palomar, "A unified framework for structured graph learning via spectral constraints," *Journal of Machine Learning Research (JMLR)*, pp. 1–60, 2020.

D. Palomar (HKUST)

Graphs with spectral constraints

Consider learning a graph with k = 3 clusters:

Smooth-optimized graph (three-circle data)

- without imposing the clusters, we get a fully connected graph;
- further approximating the previous graph with a low-rank one, we don't get the best results;

Low-rank approximated graph (three-circle data)

Low-rank ML graph (three-circle data)

• learning directly a low-rank graph is the best option.



Derivation of Structured Graph Learning (SGL) algorithm

• Recall the original problem

$$\begin{array}{ll} \underset{\mathbf{L}\succeq\mathbf{0}}{\text{maximize}} & \log \operatorname{gdet}\left(\mathbf{L}\right) - \operatorname{Tr}\left(\mathbf{SL}\right) - \rho \|\mathbf{L}\|_{1,\operatorname{off}} \\ \text{subject to} & \mathbf{L}\mathbf{1} = \mathbf{0}, \ L_{ij} = L_{ji} \leq 0, \quad \forall i \neq j \\ & \boldsymbol{\lambda}(\mathbf{L}) \in \mathcal{S}_{\boldsymbol{\lambda}}. \end{array}$$

 \bullet To properly control the eigenvalues of $\boldsymbol{\mathsf{L}},$ let's include them explicitly as variables:

$$\begin{array}{ll} \underset{\textbf{L}\succeq \textbf{0}, \lambda, \textbf{U}}{\text{maximize}} & \log \operatorname{gdet} (\textbf{L}) - \operatorname{Tr} (\textbf{SL}) - \rho \| \textbf{L} \|_{1, \operatorname{off}} \\ \text{subject to} & \textbf{L1} = \textbf{0}, \ L_{ij} = L_{ji} \leq 0, \quad \forall i \neq j \\ & \textbf{L} = \textbf{U} \operatorname{Diag} (\boldsymbol{\lambda}) \textbf{U}^{\mathsf{T}}, \quad \boldsymbol{\lambda} \in \mathcal{S}_{\boldsymbol{\lambda}}, \quad \textbf{U}^{\mathsf{T}} \textbf{U} = \textbf{I}, \end{array}$$

where λ denote the eigenvalues and ${f U}$ the eigenvectors columnwise.

• This formulation seems intractable with so many nonconvex constraints coupling the different variables.

 \checkmark To simplify it, we will introduct a linear operator \mathcal{L} that transforms the Laplacian structural constraints into simple algebraic constraints.

D. Palomar (HKUST)

Derivation of SGL algorithm

- We will now derive the linear operator ${\cal L}$ to characterize the Laplacian constraints.
- Recall that the valid Laplacian are defined in the set

$$\mathcal{S}_{\mathsf{L}} = \{\mathsf{L} \succeq \mathbf{0} \mid \mathsf{L}\mathbf{1} = \mathbf{0}, \ L_{ij} = L_{ji} \leq 0, \quad \forall i \neq j\},\$$

- It's not difficult to see that due to the constraints L1 = 0 and L_{ij} = L_{ji} the degrees of freedom of the Laplacian L ∈ ℝ^{p×p} are actually p(p-1)/2.
- We can define a linear operator $\mathcal{L} : \mathbf{w} \in \mathbb{R}^{p(p-1)/2}_+ \to \mathcal{L}\mathbf{w} \in \mathbb{R}^{p \times p}$ that maps a weight vector \mathbf{w} to a valid Laplacian matrix satisfying all the required properties.
- Eaxmple for p = 4: $\mathbf{w} = [w_1, w_2, \dots, w_6] \in \mathbb{R}^6$ and

$$\mathcal{L}\mathbf{w} = \begin{bmatrix} \sum_{i=1,2,3} w_i & -w_1 & -w_2 & -w_3 \\ -w_1 & \sum_{i=1,4,5} w_i & -w_4 & -w_5 \\ -w_2 & -w_4 & \sum_{i=2,4,6} w_i & -w_6 \\ -w_3 & -w_5 & -w_6 & \sum_{i=3,5,6} w_i \end{bmatrix}$$
• Using $\mathbf{L} = \mathcal{L} \mathbf{w}$ and $\operatorname{gdet}(\mathbf{L}) = \operatorname{gdet}(\operatorname{Diag}(\boldsymbol{\lambda}))$, we can rewrite the problem as

$$\begin{array}{ll} \underset{\mathbf{w},\lambda,\mathbf{U}}{\operatorname{maximize}} & \log \operatorname{gdet}\left(\operatorname{Diag}(\lambda)\right) - \operatorname{Tr}\left(\mathbf{S}\mathcal{L}\mathbf{w}\right) - \rho \|\mathcal{L}\mathbf{w}\|_{1,\operatorname{off}} \\ \operatorname{subject to} & \mathcal{L}\mathbf{w} = \mathbf{U}\operatorname{Diag}(\lambda)\mathbf{U}^{\mathcal{T}}, \quad \lambda \in \mathcal{S}_{\lambda}, \quad \mathbf{U}^{\mathcal{T}}\mathbf{U} = \mathbf{I}, \end{array}$$

- Noting that $\|\mathcal{L}\mathbf{w}\|_{1,\text{off}} = \frac{1}{2}\|\mathbf{w}\|_1 = \frac{1}{2}\mathbf{1}^T\mathbf{w}$, we can group together the two linear terms on \mathbf{w} (i.e., $\text{Tr}(\mathbf{S}\mathcal{L}\mathbf{w})$ and $\|\mathcal{L}\mathbf{w}\|_{1,\text{off}}$) into the single linear term $\mathbf{k}^T\mathbf{w}$ with a properly defined \mathbf{k} .
- Finally, relaxing the hard constraint *L*w = UDiag(λ)U^T to the objective, we get the approximate formulation:

$$\begin{array}{ll} \underset{\mathbf{w},\lambda,\mathbf{U}}{\text{minimize}} & \mathbf{k}^{T}\mathbf{w} - \log \operatorname{gdet}\left(\operatorname{Diag}(\lambda)\right) + \frac{\beta}{2} \|\mathcal{L}\mathbf{w} - \mathbf{U}\operatorname{Diag}(\lambda)\mathbf{U}^{T}\|_{F}^{2} \\ \text{subject to} & \mathbf{w} \geq \mathbf{0}, \quad \lambda \in \mathcal{S}_{\lambda}, \quad \mathbf{U}^{T}\mathbf{U} = \mathbf{I}. \end{array}$$

- To solve the problem in the three variables (w, λ, U) we will use the block majorization-minimization (MM) method, which updates each variable sequentially while keeping the others fixed based on simple surrogate functions (Sun et al. 2017)¹⁵, (Razaviyayn et al. 2013)¹⁶.
- For illustration purposes we will give the algorithm for a k-connected graph.
- The constraints for each of the three variables are:
 - Spectral constraint: $S_{\lambda} = \{\{\lambda_j = 0\}_{j=1}^k, c_1 \leq \lambda_{k+1} \leq \cdots \leq \lambda_p \leq c_2\}.$
 - Nonnegativity constraint: $\mathbf{w} \geq \mathbf{0}$.
 - Orthogonality constraint: $\mathbf{U}^{T}\mathbf{U} = \mathbf{I}$.

¹⁵Y. Sun, P. Babu, and D. P. Palomar, "Majorization-minimization algorithms in signal processing, communications, and machine learning," *IEEE Trans. Signal Processing*, vol. 65, no. 3, pp. 794–816, 2017.
¹⁶M. Razaviyayn, M. Hong, and Z. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization," *SIAM J. Optim.*, vol. 23, no. 2, pp. 1126–1153, 2013.

Derivation of SGL algorithm: Update for w

• Sub-problem for **w**:

$$\begin{array}{ll} \underset{\mathsf{w}\geq \mathsf{0}}{\text{minimize}} \quad \mathsf{k}^{\mathsf{T}}\mathsf{w} + \frac{\beta}{2} \|\mathcal{L}\mathsf{w} - \mathsf{U}\mathsf{D}\mathsf{iag}(\boldsymbol{\lambda})\mathsf{U}^{\mathsf{T}}\|_{\mathsf{F}}^2. \end{array}$$

• It can be rewritten as

$$\begin{array}{ll} \underset{\mathbf{w}\geq\mathbf{0}}{\text{minimize}} \quad f(\mathbf{w}) \triangleq \|\mathcal{L}\mathbf{w}\|_{F}^{2} - \mathbf{c}^{T}\mathbf{w} \end{array}$$

which is clearly a convex quadratic program (QP).

- The QP does not have a closed-form solution due to the constraint $\mathbf{w} \geq \mathbf{0}$ and we will use MM.
- The function $f(\mathbf{w})$ is majorized at \mathbf{w}^t by

$$g(\mathbf{w}|\mathbf{w}^{t}) = f(\mathbf{w}^{t}) + (\mathbf{w} - \mathbf{w}^{t})^{T} \nabla f(\mathbf{w}^{t}) + \frac{L}{2} \|\mathbf{w} - \mathbf{w}^{t}\|^{2}$$

• The majorized problem has now a simple closed-form solution:

$$\mathbf{w}^{t+1} = \left(\mathbf{w}^t - \frac{1}{2\rho} \nabla f(\mathbf{w}^t)\right)^+$$

Derivation of SGL algorithm: Update for U

• Sub-problem for **U**:

$$\begin{array}{ll} \underset{\mathbf{U}}{\text{minimize}} & \frac{\beta}{2} \| \mathcal{L} \mathbf{w} - \mathbf{U} \text{Diag}(\boldsymbol{\lambda}) \mathbf{U}^{T} \|_{F}^{2} \\ \text{subject to} & \mathbf{U}^{T} \mathbf{U} = \mathbf{I}. \end{array}$$

• It can be rewritten as

$$\begin{array}{ll} \max \underset{\mathbf{U}}{\operatorname{maximize}} & \operatorname{Tr}(\mathbf{U}^{T}\mathcal{L}\mathbf{w}\mathbf{U}\operatorname{Diag}(\boldsymbol{\lambda}))\\ \text{subject to} & \mathbf{U}^{T}\mathbf{U} = \mathbf{I}. \end{array}$$

- This sub-problem is an optimization on the orthogonal Stiefel manifold (Absil et al. 2009; Benidis et al. 2016).
- From the KKT optimality conditions the solution is given by

$$\mathbf{U}^{t+1} = ext{eigenvectors}(\mathcal{L}\mathbf{w}^{t+1})[k+1:p],$$

that is, the p - k principal eigenvectors of the matrix $\mathcal{L}\mathbf{w}^{t+1}$ in increasing order of eigenvalue magnitude.

Derivation of SGL algorithm: Update for λ

• Sub-problem for λ :

$$\begin{array}{ll} \underset{\boldsymbol{\lambda}\in\mathcal{S}_{\boldsymbol{\lambda}}}{\text{minimize}} & -\text{log gdet}\left(\text{Diag}(\boldsymbol{\lambda})\right) + \frac{\beta}{2}\|\mathcal{L}\mathbf{w} - \mathbf{U}\text{Diag}(\boldsymbol{\lambda})\mathbf{U}^{\mathcal{T}}\|_{F^{1}}^{2} \end{array}$$

• It can be rewritten for the k-component graph as

$$\underset{c_1 \leq \lambda_{k+1} \leq \dots \leq \lambda_p \leq c_2}{\text{minimize}} - \sum_{i=1}^{p-k} \log(\lambda_{k+i}) + \frac{\beta}{2} \|\boldsymbol{\lambda} - \mathbf{d}\|^2$$

- This sub-problem is popularly known as a regularized isotonic regression problem. It is a convex optimization problem and the solution can be obtained from the KKT optimality conditions.
- An efficient algorithm with a fast convergence to the global optimum can be derived with a maximum of p k iterations (Kumar et al. 2020).

Algorithm SGL for *k*-connected graph

Input: S, k, c_1, c_2, β Output: $\mathcal{L}w$ $t \leftarrow 0$

repeat

•
$$\mathbf{w}^{t+1} = \left(\mathbf{w}^t - \frac{1}{2p}\nabla f(\mathbf{w}^t)\right)^{\top}$$

- $\mathbf{U}^{t+1} \leftarrow \text{eigenvectors}(\mathbf{Lw}^{t+1})$, suitably ordered.
- Update $oldsymbol{\lambda}^{t+1}$ via isotonic regression.
- $t \leftarrow t+1$

until convergence return $\mathcal{L}\mathbf{w}^t$

The following gives the convergence of the SGL algorithm to a stationary point (Kumar et al. 2019)¹⁷ and (Kumar et al. 2020)¹⁸.

Theorem

The limit point $(\mathbf{w}^*, \lambda^*, \mathbf{U}^*)$ generated by the SGL algorithm converges to the set of KKT points of the optimization problem.

Furthermore, the worst-case computational complexity of the proposed algorithm is $O(p^3)$.

¹⁷S. Kumar, J. Ying, J. V. de M. Cardoso, and D. P. Palomar, "Structured graph learning via laplacian spectral constraints," in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, Vancouver, Canada, 2019.

¹⁸S. Kumar, J. Ying, J. V. de M. Cardoso, and D. P. Palomar, "A unified framework for structured graph learning via spectral constraints," *Journal of Machine Learning Research (JMLR)*, pp. 1–60, 2020.

D. Palomar (HKUST)

Outline

1 Graphs

2 Basics

3 Learning Graphs from Data

- Similarity function based
- Smooth signal based
- i.i.d. model based
- Structured graphs via spectral constraints

4 Numerical Experiments

Synthetic experiment setup

- Generate a graph with desired structure.
- Generate weights for the graph edges.
- Obtain true Laplacian \mathbf{L}_{true} .
- Sample data $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ with $\mathbf{x}_i \in \mathbb{R}^p \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma} = \mathbf{L}_{true}^{\dagger})$.
- $\mathbf{S} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i \mathbf{x}_i^T$
- \bullet Use \bm{S} as input for the learning and some prior spectral information, if available, to estimate the graph $\hat{\bm{L}}$
- Performance metric:

$$\label{eq:Relative Error} \text{Relative Error} = \frac{\|\hat{\textbf{L}} - \textbf{L}_{\text{true}}\|_{\textit{F}}}{\|\textbf{L}_{\text{true}}\|_{\textit{F}}} \quad \text{and} \quad \text{F-Score} = \frac{2\text{tp}}{2\text{tp} + \text{fp} + \text{fn}},$$

where tp, fp, fn correspond to true positives, false positives, and false negatives, respectively.

• SGL refers to structured graph learning.

Grid graph



D. Palomar (HKUST)

Graph Learning

Noisy multi-component graph



D. Palomar (HKUST)

Graph Learning

Model mismatch



Popular multi-component structures



Real data: cancer dataset (Weinstein et al. 2013)



• Clustering accuracy (ACC): CLR = 0.9862 and SGL = 0.99875.

Animal dataset (Osherson et al. 1991)



D. Palomar (HKUST)

Cockroach

Bipartite structure via adjacency spectral constraints



Multi-component bipartite structure via joint spectral constraints





For more information visit:

https://www.danielppalomar.com



Absil, P. A., Mahony, R., & Sepulchre, R. (2009). *Optimization algorithms on matrix manifolds.* Princeton University Press.

Banerjee, O., Ghaoui, L., & d'Aspremont, A. (2008). Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data. *Journal of Machine Learning Research*, 9(Mar), 485–516.

Benidis, K., Sun, Y., Babu, P., & Palomar, D. P. (2016). Orthogonal sparse PCA and covariance estimation via Procrustes reformulation. *IEEE Trans. Signal Processing*, *64*(23), 6211–6226.

Bogdanov, A., Mossel, E., & Vadhan, S. (2008). The complexity of distinguishing markov random fields. In *Approximation, randomization and combinatorial optimization. Algorithms and techniques* (pp. 331–342). Springer.

Dempster, A. P. (1972). Covariance selection. Biometrics, 157-175.

References II

Dong, X., Thanou, D., Rabbat, M., & Frossard, P. (2019). Learning graphs from data: A signal representation perspective. *IEEE Signal Processing Magazine*, *36*(3), 44–63.

Egilmez, H. E., Pavez, E., & Ortega, A. (2017). Graph learning from data under Laplacian and structural constraints. *IEEE Journal of Selected Topics in Signal Processing*, 11(6), 825–841.

Fan, K. (1949). On a theorem of weyl concerning eigenvalues of linear transformations i. *Proceedings of the National Academy of Sciences of the United States of America*, *35*(11), 652.

Friedman, J., Hastie, T., & Tibshirani, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, *9*(3), 432–441.

Heimo, T., Saramaki, J., Onnela, J.-P., & Kaski, K. (2007). Spectral and network methods in the analysis of correlation matrices of stock returns. *Physica A: Statistical Mechanics and its Applications*, *383*(1), 147–151.

Kalofolias, V. (2016). How to learn a graph from smooth signals. In *Proc. Int. Conf. Artif. Intell. Statist.* (pp. 920–929).

Kumar, S., Ying, J., M. Cardoso, J. V. de, & Palomar, D. P. (2019). Structured graph learning via laplacian spectral constraints. In *Proc. Advances in neural information processing systems (neurips)*. Vancouver, Canada.

Kumar, S., Ying, J., M. Cardoso, J. V. de, & Palomar, D. P. (2020). A unified framework for structured graph learning via spectral constraints. *Journal of Machine Learning Research (JMLR)*, 1–60.

Lake, B., & Tenenbaum, J. (2010). Discovering structure by learning sparse graphs. In *Proc. The 33rd annual cognitive science conference.*

Manor, L. Z., & Perona, P. (2004). Self-tuning spectral clustering. In *Proc. Advances in neural information processing systems (neurips)*.

Mateos, G., Segarra, S., Marques, A. G., & Ribeiro, A. (2019). Connecting the dots. *IEEE Signal Processing Magazine*, *36*(3), 16–43.

Mazumder, R., & Hastie, T. (2012). The graphical lasso: New insights and alternatives. *Electronic journal of statistics*, *6*, 2125.

Meinshausen, N., & Bühlmann, P. (2006). High-dimensional graphs and variable selection with the lasso. *The annals of statistics*, *34*(3), 1436–1462.

Nie, F., Wang, X., Jordan, M., & Huang, H. (2016). The constrained Laplacian rank algorithm for graph-based clustering. In *Proc. Of the thirtieth conference on artificial intelligence (aaai)* (pp. 1969–1976). Phoenix, Arizona, USA.

Osherson, D. N., Stern, J., Wilkie, O., Stob, M., & Smith, E. E. (1991). Default probability. *Cognitive Science*, *15*(2), 251–269.

Ravikumar, P., Wainwright, M., Lafferty, J. D., & others. (2010). High-dimensional ising model selection using ℓ_1 -regularized logistic regression. *The Annals of Statistics*, *38*(3), 1287–1319.

Razaviyayn, M., Hong, M., & Luo, Z. (2013). A unified convergence analysis of block successive minimization methods for nonsmooth optimization. *SIAM J. Optim.*, *23*(2), 1126–1153.

Slawski, M., & Hein, M. (2015). Estimation of positive definite m-matrices and structure learning for attractive gaussian markov random fields. *Linear Algebra and its Applications*, 473, 145–179.

Sun, Y., Babu, P., & Palomar, D. P. (2017). Majorization-minimization algorithms in signal processing, communications, and machine learning. *IEEE Trans. Signal Processing*, *65*(3), 794–816.

Weinstein, J. N., Collisson, E. A., Mills, G. B., Shaw, K. R. M., Ozenberger, B. A., Ellrott, K., et al. (2013). The cancer genome atlas pan-cancer analysis project. *Nature Genetics*, *45*(10), 1113.

Zhao, L., Wang, Y., Kumar, S., & Palomar, D. P. (2019). Optimization algorithms for graph Laplacian estimation via ADMM and MM. *IEEE Trans. on Signal Processing*, *67*(16), 4231–4244.