

*Inria*



ENS DE LYON

**Learning graphs with precision matrices:  
statistical estimators, compressive approaches,  
unrolled neural networks**

**Titouan Vayer**

Can Pouliquen Etienne Lasalle Mathurin Massias Rémi Gribonval Paulo Gonçalves

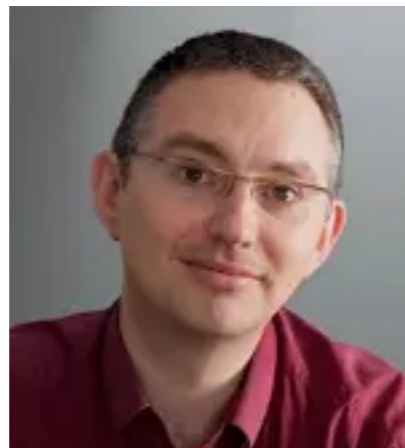


# | Overview of the talk

■ Part I: **Finding graphs from unstructured data**

■ Part II: **Schur's Positive-Definite Network**

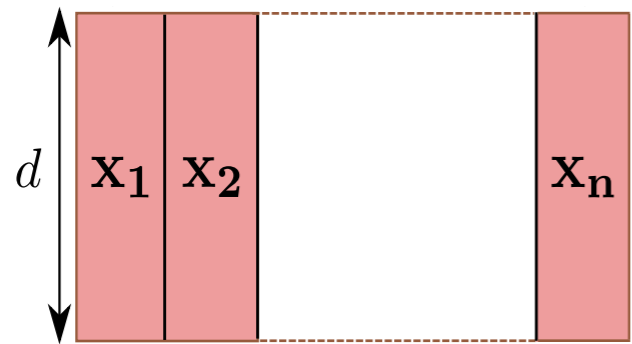
■ Part III: **The sketching approach**



# | Graph Learning

■ Input: a dataset

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$$

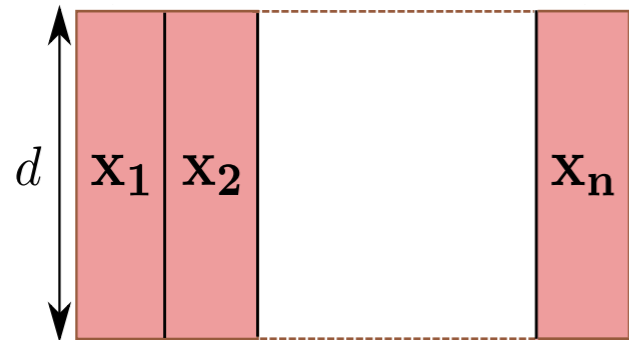


$$\mathbf{x}_i \in \mathbb{R}^d \sim \mu$$

# Graph Learning

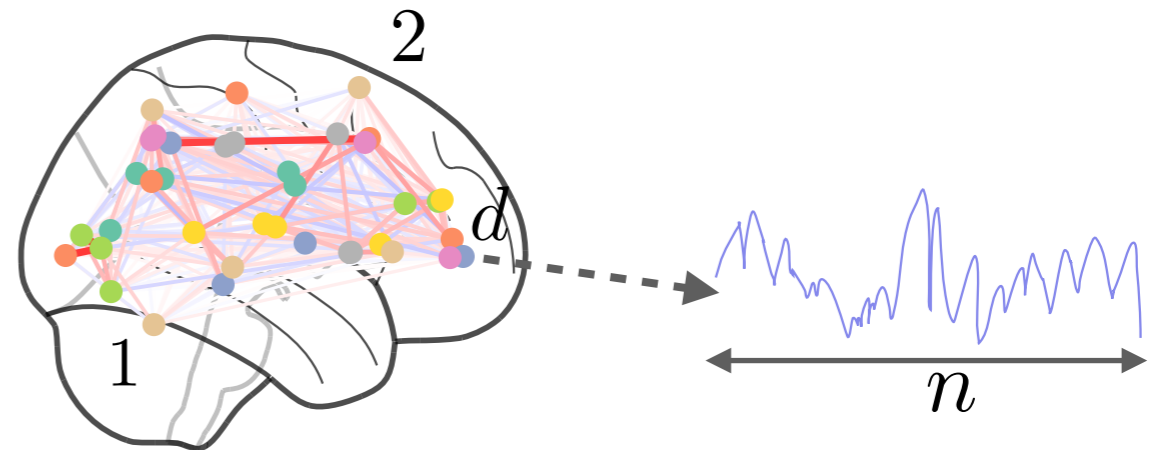
Input: a dataset

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$$



$$\mathbf{x}_i \in \mathbb{R}^d \sim \mu$$

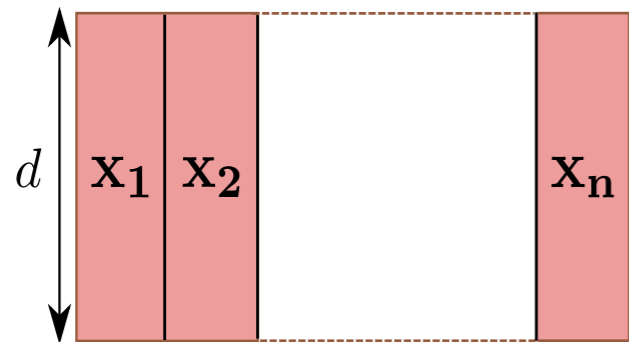
Output: graph of relations between the  $d$  variables



# Graph Learning

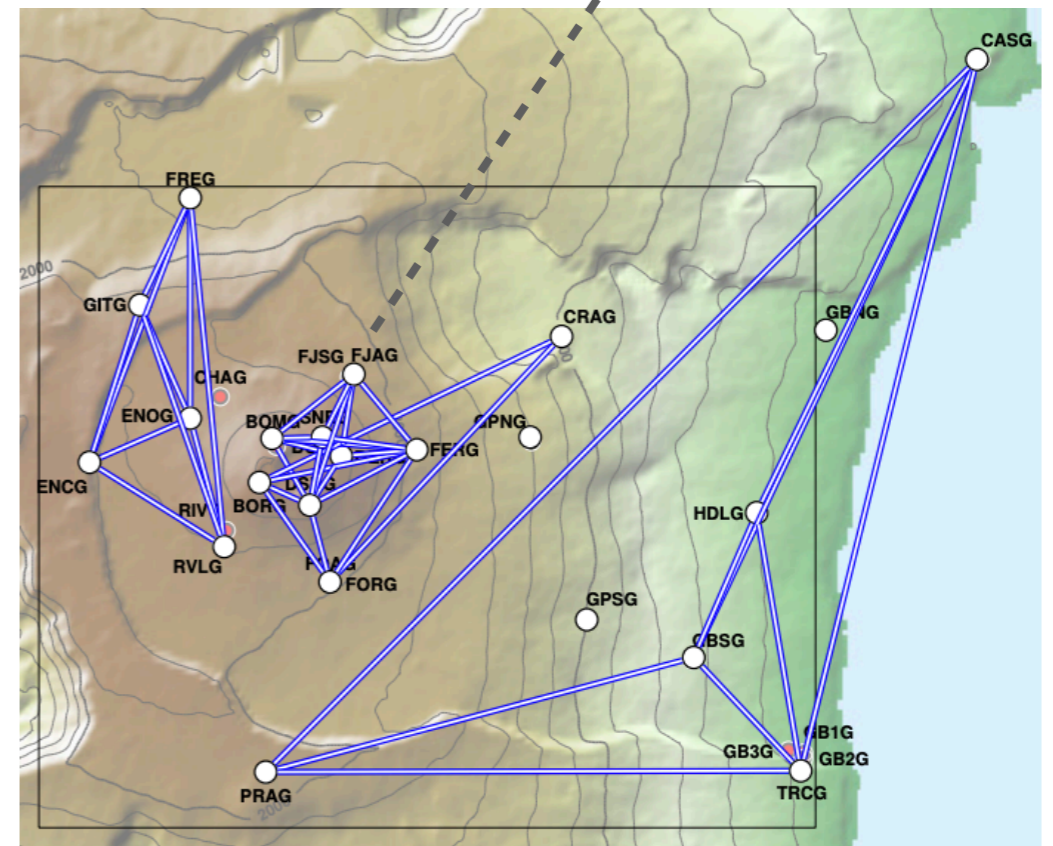
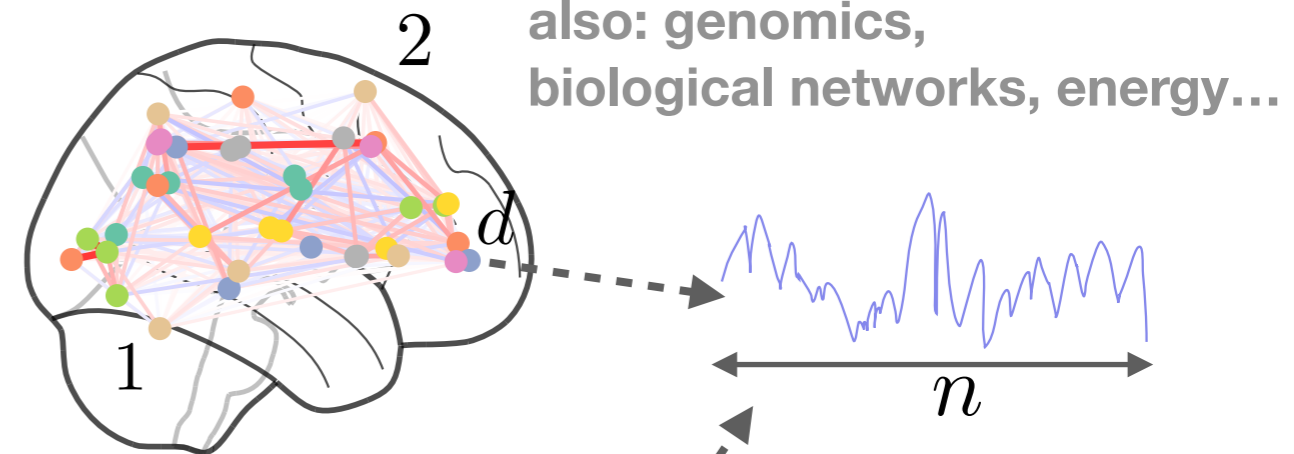
Input: a dataset

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$$



$$\mathbf{x}_i \in \mathbb{R}^d \sim \mu$$

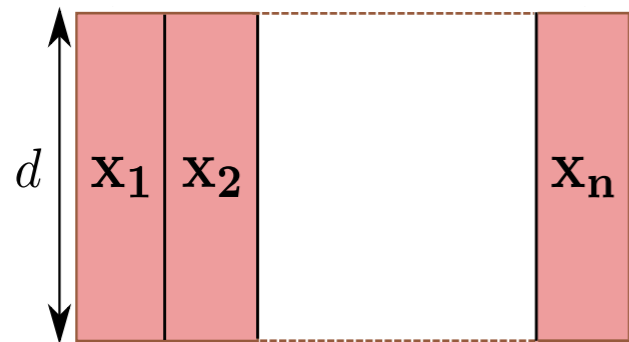
Output: graph of relations between the  $d$  variables



# Graph Learning

Input: a dataset

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$$



$$\mathbf{x}_i \in \mathbb{R}^d \sim \mu$$

Graph modeled as a matrix:

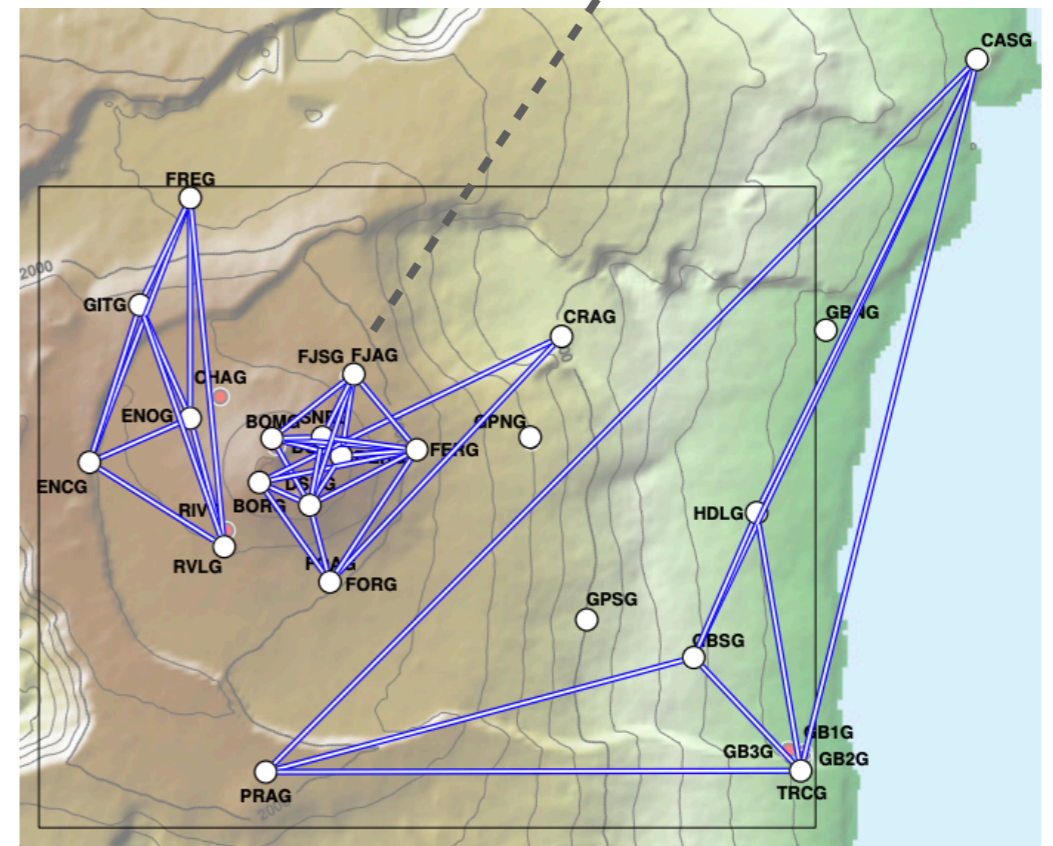
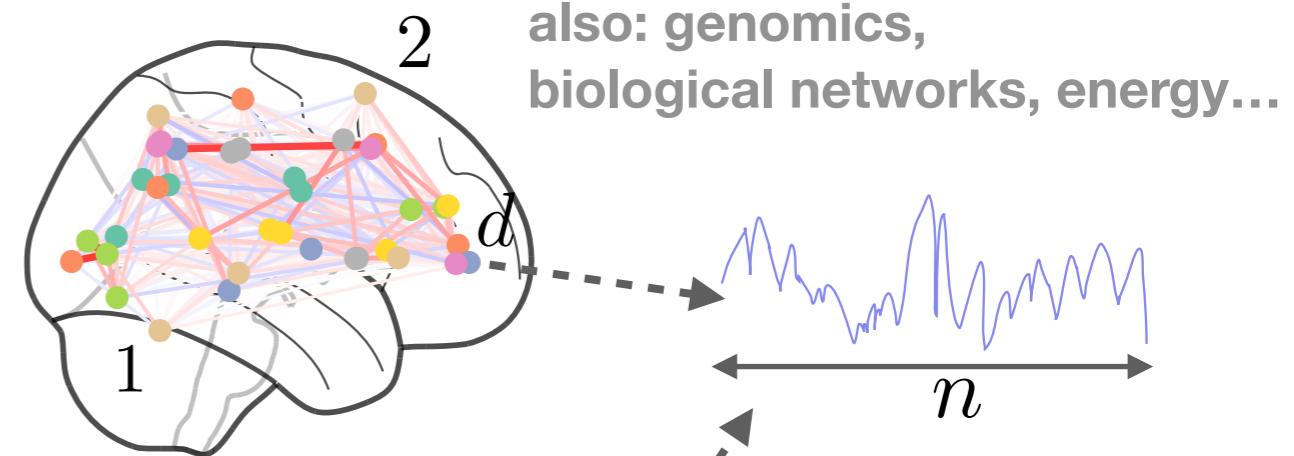
$$\Theta \in \mathbb{R}^{d \times d}$$

$\Theta_{ij}$ : **interaction** between variable  $i$  and  $j$

| statistical correlations

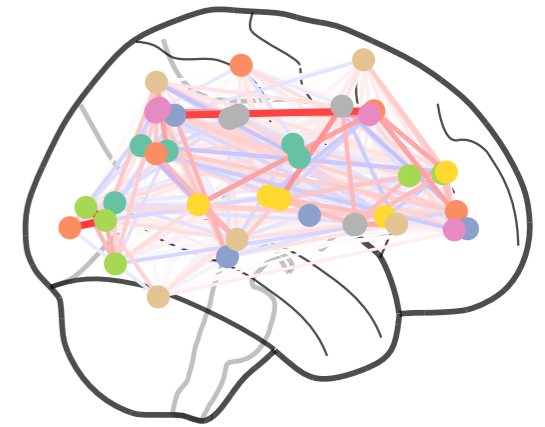
| statistical dependencies

Output: graph of relations between the  $d$  variables



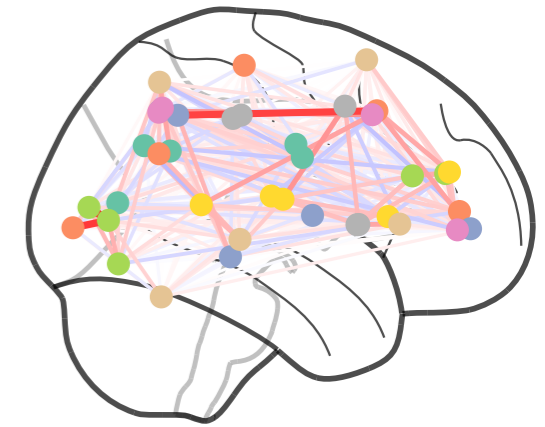
# | Epilepsy: the big picture

- One of the **most common neurological disorder**, affecting 1% of the global population
- Nearly 30% of patients are **drug-resistant**
- Surgical solution: remove the epileptic onset area (**resection**)

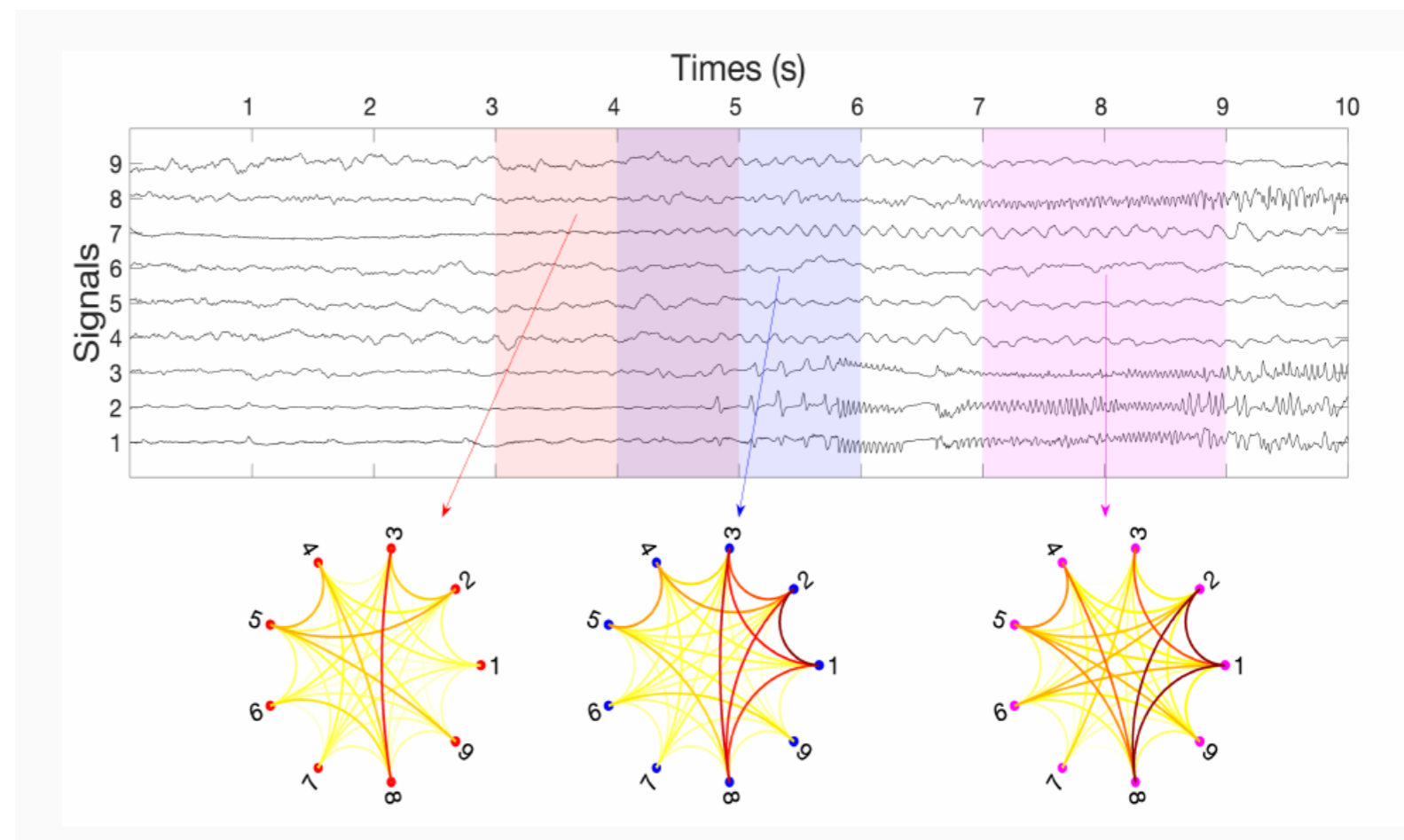


# Epilepsy: the big picture

- One of the **most common neurological disorder**, affecting 1% of the global population
- Nearly 30% of patients are **drug-resistant**
- Surgical solution: remove the epileptic onset area (**resection**)



- Sparse dynamic functional connectivity graphs from EEG signals





# | Graphical LASSO

Side note

- Input:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$
- $\mathbf{x}_i \in \mathbb{R}^d \sim \mu$
- Output:  $\Theta \in \mathbb{R}^{d \times d}$

# | Graphical LASSO

## ■ Gaussian Graphical Model

Gaussian assumption  $\mu = \mathcal{N}(0, \Sigma = \Theta^{-1})$

■  $\Theta_{ij} = 0 \iff$  variable  $i$  is independent of  $j$  conditionally to the others

Side note

■ Input:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$   
↓  $\mathbf{x}_i \in \mathbb{R}^d \sim \mu$   
■ Output:  $\Theta \in \mathbb{R}^{d \times d}$

# Graphical LASSO

## Gaussian Graphical Model

Gaussian assumption  $\mu = \mathcal{N}(0, \Sigma = \Theta^{-1})$

■  $\Theta_{ij} = 0 \iff$  variable  $i$  is independent of  $j$  conditionally to the others

## Maximum Likelihood estimator

Emp. cov.  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$

$$\Theta_{\text{MLE}} = \arg \min_{\Theta \succ 0} -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle_F$$

■ When  $\hat{\Sigma}$  is invertible  $\Theta_{\text{MLE}} = (\hat{\Sigma})^{-1}$

Side note

■ Input:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$   
    ↓  $\mathbf{x}_i \in \mathbb{R}^d \sim \mu$   
■ Output:  $\Theta \in \mathbb{R}^{d \times d}$

# Graphical LASSO

## Gaussian Graphical Model

Gaussian assumption  $\mu = \mathcal{N}(0, \Sigma = \Theta^{-1})$

■  $\Theta_{ij} = 0 \iff$  variable  $i$  is independent of  $j$  conditionally to the others

## Maximum Likelihood estimator

Emp. cov.  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$

$$\Theta_{\text{MLE}} = \arg \min_{\Theta \succ 0} -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle_F$$

■ When  $\hat{\Sigma}$  is invertible  $\Theta_{\text{MLE}} = (\hat{\Sigma})^{-1}$

May not be true  
in high dim  $n < d$

Side note

■ Input:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$   
    ↓  $\mathbf{x}_i \in \mathbb{R}^d \sim \mu$   
■ Output:  $\Theta \in \mathbb{R}^{d \times d}$

# Graphical LASSO

## Gaussian Graphical Model

Gaussian assumption  $\mu = \mathcal{N}(0, \Sigma = \Theta^{-1})$

$\Theta_{ij} = 0 \iff$  variable  $i$  is independent of  $j$  conditionally to the others

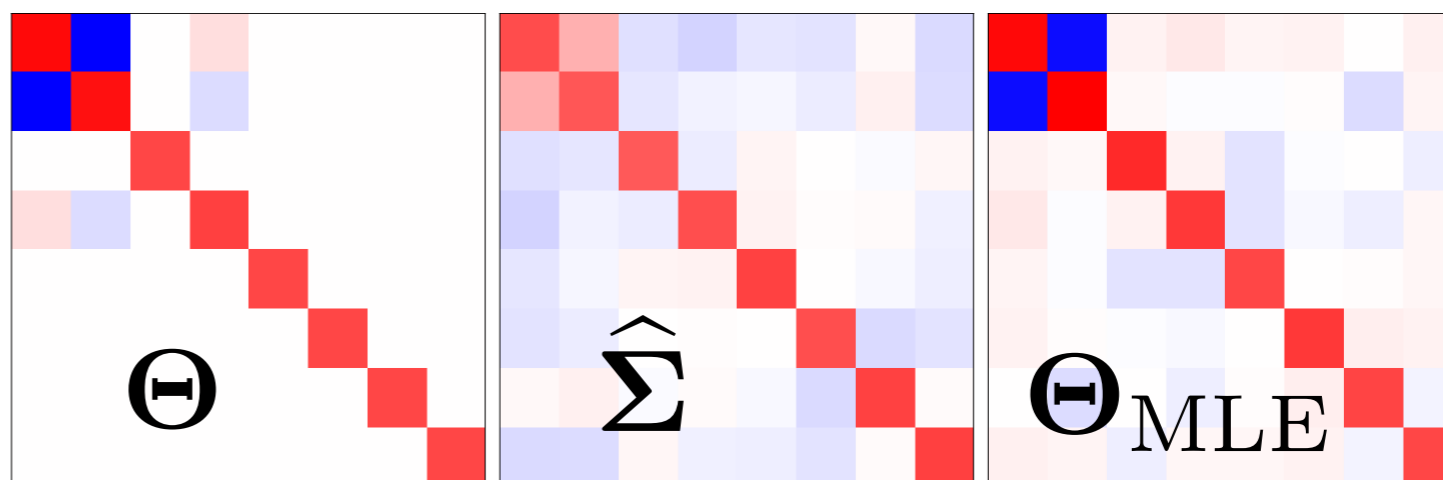
## Maximum Likelihood estimator

Emp. cov.  $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$

$$\Theta_{\text{MLE}} = \arg \min_{\Theta \succ 0} -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle_F$$

When  $\hat{\Sigma}$  is invertible  $\Theta_{\text{MLE}} = (\hat{\Sigma})^{-1} \dashrightarrow$  usually not sparse

May not be true  
in high dim  $n < d$



Side note

Input:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$   
 $\mathbf{x}_i \in \mathbb{R}^d \sim \mu$   
Output:  $\Theta \in \mathbb{R}^{d \times d}$

# Graphical LASSO

Side note

Input:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$   
 $\mathbf{x}_i \in \mathbb{R}^d \sim \mu$   
Output:  $\Theta \in \mathbb{R}^{d \times d}$

## Gaussian Graphical Model

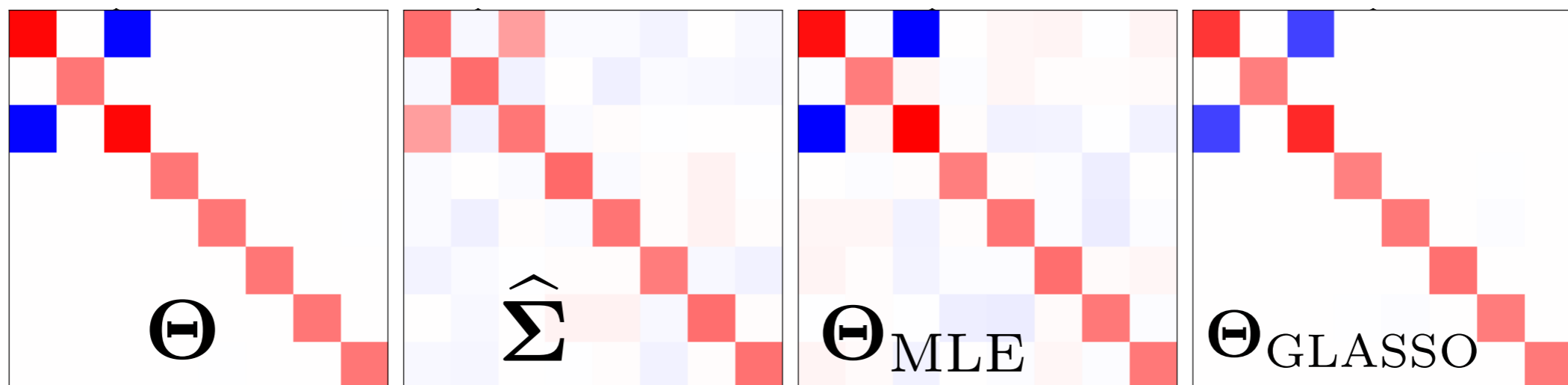
Gaussian assumption  $\mu = \mathcal{N}(0, \Sigma = \Theta^{-1})$

$\Theta_{ij} = 0 \iff$  variable  $i$  is independent of  $j$  conditionally to the others

## Penalized Maximum Likelihood estimator [Friedman-Hastie-Tibshirani, 2007]

$$\Theta_{\text{GLASSO}} = \arg \min_{\Theta \succ 0} -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle_F + \lambda \|\Theta\|_{1,\text{off}}$$

$\|\Theta\|_{1,\text{off}} = \sum_{i < j} |\Theta_{ij}|$  promotes sparsity for the output graph

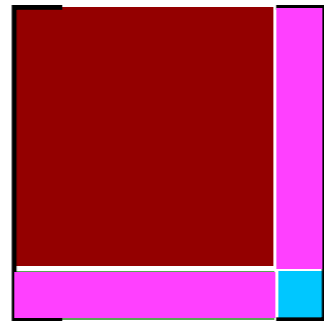


# | Graphical LASSO

## ■ Solving for GLASSO

$$\operatorname{argmin}_{\Theta > 0} F(\Theta) = -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle + \lambda \|\Theta\|_{1, \text{off}}$$

$$\Theta = \begin{pmatrix} \Theta_{11} & \theta_{12} \\ \theta_{12}^\top & \theta_{22} \end{pmatrix}$$



■ Schur complement

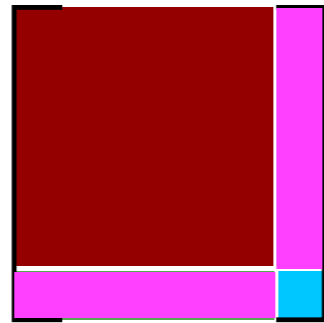
$$\det(\Theta) = \det(\Theta_{11}) \times \det(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12})$$

# Graphical LASSO

## Solving for GLASSO

$$\operatorname{argmin}_{\Theta > 0} F(\Theta) = -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle + \lambda \|\Theta\|_{1, \text{off}}$$

$$\Theta = \begin{pmatrix} \Theta_{11} & \theta_{12} \\ \theta_{12}^\top & \theta_{22} \end{pmatrix}$$



Schur complement

$$\det(\Theta) = \det(\Theta_{11}) \times \det(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12})$$

$$F(\Theta) = -\log \det(\Theta_{11}) - \log(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12}) + \hat{\sigma}_{22} \theta_{22} + \langle \hat{\Sigma}_{11}, \Theta_{11} \rangle + 2 \langle \hat{\sigma}_{12}, \theta_{12} \rangle + \lambda \|\Theta_{11}\|_{1, \text{off}} + \lambda \|\theta_{12}\|_1$$

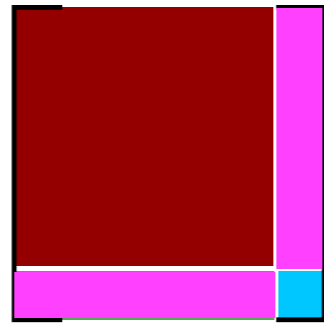


# Graphical LASSO

## Solving for GLASSO

$$\operatorname{argmin}_{\Theta > 0} F(\Theta) = -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle + \lambda \|\Theta\|_{1, \text{off}}$$

$$\Theta = \begin{pmatrix} \Theta_{11} & \theta_{12} \\ \theta_{12}^\top & \theta_{22} \end{pmatrix}$$



Schur complement

$$\det(\Theta) = \det(\Theta_{11}) \times \det(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12})$$

$$F(\Theta) = -\log \det(\Theta_{11}) - \log(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12}) + \hat{\sigma}_{22} \theta_{22} + \langle \hat{\Sigma}_{11}, \Theta_{11} \rangle + 2 \langle \hat{\sigma}_{12}, \theta_{12} \rangle + \lambda \|\Theta_{11}\|_{1, \text{off}} + \lambda \|\theta_{12}\|_1$$

BCD algorithm:  $\operatorname{argmin}_{\theta_{22}} F(\Theta_{11}, \theta_{12}, \theta_{22}) = \frac{1}{\hat{\sigma}_{22}} + \theta_{12}^\top \Theta_{11}^{-1} \theta_{12} := \theta_{22}^*$

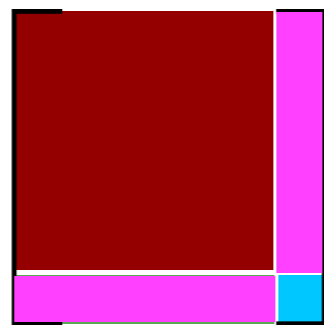
$$\operatorname{argmin}_{\theta_{12}} F(\Theta_{11}, \theta_{12}, \theta_{22}^*) = 2 \langle \hat{\sigma}_{12}, \theta_{12} \rangle + \lambda \|\theta_{12}\|_1 + \hat{\sigma}_{22} \theta_{12}^\top \Theta_{11}^{-1} \theta_{12}$$

# Graphical LASSO

## Solving for GLASSO

$$\operatorname{argmin}_{\Theta > 0} F(\Theta) = -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle + \lambda \|\Theta\|_{1, \text{off}}$$

$$\Theta = \begin{pmatrix} \Theta_{11} & \theta_{12} \\ \theta_{12}^\top & \theta_{22} \end{pmatrix}$$



Schur complement

$$\det(\Theta) = \det(\Theta_{11}) \times \det(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12})$$

$$F(\Theta) = -\log \det(\Theta_{11}) - \log(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12}) + \hat{\sigma}_{22} \theta_{22} + \langle \hat{\Sigma}_{11}, \Theta_{11} \rangle + 2 \langle \hat{\sigma}_{12}, \theta_{12} \rangle + \lambda \|\Theta_{11}\|_{1, \text{off}} + \lambda \|\theta_{12}\|_1$$

BCD algorithm:  $\operatorname{argmin}_{\theta_{22}} F(\Theta_{11}, \theta_{12}, \theta_{22}) = \frac{1}{\hat{\sigma}_{22}} + \theta_{12}^\top \Theta_{11}^{-1} \theta_{12} := \theta_{22}^*$

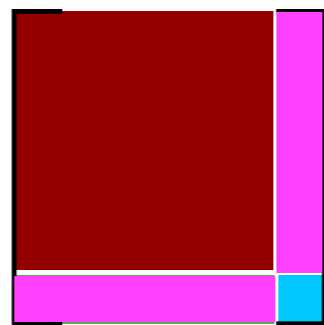
$$\text{LASSO} \sim \operatorname{argmin}_{\theta_{12}} F(\Theta_{11}, \theta_{12}, \theta_{22}^*) = 2 \langle \hat{\sigma}_{12}, \theta_{12} \rangle + \lambda \|\theta_{12}\|_1 + \hat{\sigma}_{22} \theta_{12}^\top \Theta_{11}^{-1} \theta_{12}$$

# Graphical LASSO

## Solving for GLASSO

$$\operatorname{argmin}_{\Theta > 0} F(\Theta) = -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle + \lambda \|\Theta\|_{1, \text{off}}$$

$$\Theta = \begin{pmatrix} \Theta_{11} & \theta_{12} \\ \theta_{12}^\top & \theta_{22} \end{pmatrix}$$



Schur complement

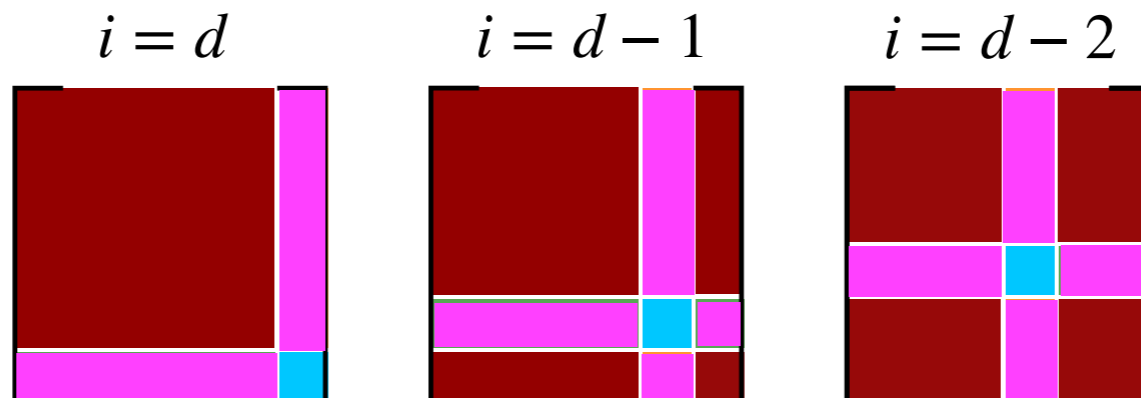
$$\det(\Theta) = \det(\Theta_{11}) \times \det(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12})$$

$$F(\Theta) = -\log \det(\Theta_{11}) - \log(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12}) + \hat{\sigma}_{22} \theta_{22} + \langle \hat{\Sigma}_{11}, \Theta_{11} \rangle + 2 \langle \hat{\sigma}_{12}, \theta_{12} \rangle + \lambda \|\Theta_{11}\|_{1, \text{off}} + \lambda \|\theta_{12}\|_1$$

BCD algorithm:  $\operatorname{argmin}_{\theta_{22}} F(\Theta_{11}, \theta_{12}, \theta_{22}) = \frac{1}{\hat{\sigma}_{22}} + \theta_{12}^\top \Theta_{11}^{-1} \theta_{12} := \theta_{22}^*$

$$\text{LASSO} \sim \operatorname{argmin}_{\theta_{12}} F(\Theta_{11}, \theta_{12}, \theta_{22}^*) = 2 \langle \hat{\sigma}_{12}, \theta_{12} \rangle + \lambda \|\theta_{12}\|_1 + \hat{\sigma}_{22} \theta_{12}^\top \Theta_{11}^{-1} \theta_{12}$$

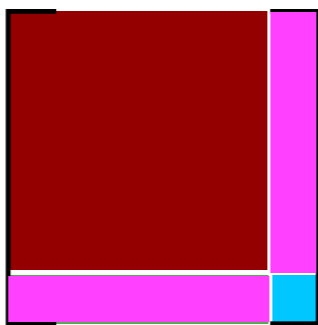
Then iterates on the columns



# Graphical LASSO

## Solving for GLASSO

$$\operatorname{argmin}_{\Theta > 0} F(\Theta) = -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle + \lambda \|\Theta\|_{1, \text{off}}$$

$$\Theta = \begin{pmatrix} \Theta_{11} & \theta_{12} \\ \theta_{12}^\top & \theta_{22} \end{pmatrix}$$


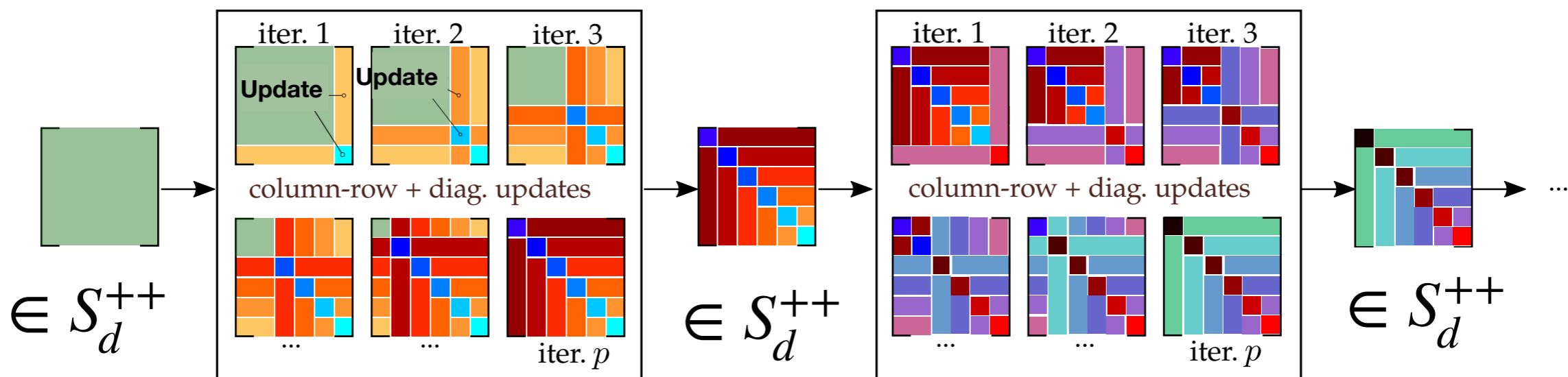
Schur complement

$$\det(\Theta) = \det(\Theta_{11}) \times \det(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12})$$

$$F(\Theta) = -\log \det(\Theta_{11}) - \log(\theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12}) + \hat{\sigma}_{22} \theta_{22} + \langle \hat{\Sigma}_{11}, \Theta_{11} \rangle + 2 \langle \hat{\sigma}_{12}, \theta_{12} \rangle + \lambda \|\Theta_{11}\|_{1, \text{off}} + \lambda \|\theta_{12}\|_1$$

BCD algorithm:  $\operatorname{argmin}_{\theta_{22}} F(\Theta_{11}, \theta_{12}, \theta_{22}) = \frac{1}{\hat{\sigma}_{22}} + \theta_{12}^\top \Theta_{11}^{-1} \theta_{12} := \theta_{22}^*$

$$\text{LASSO} \sim \operatorname{argmin}_{\theta_{12}} F(\Theta_{11}, \theta_{12}, \theta_{22}^*) = 2 \langle \hat{\sigma}_{12}, \theta_{12} \rangle + \lambda \|\theta_{12}\|_1 + \hat{\sigma}_{22} \theta_{12}^\top \Theta_{11}^{-1} \theta_{12}$$



# Graphical LASSO

## Gaussian Graphical Model

Gaussian assumption  $\mu = \mathcal{N}(0, \Sigma = \Theta^{-1})$

■  $\Theta_{ij} = 0 \iff$  variable  $i$  is independent of  $j$  conditionally to the others

## Penalized Maximum Likelihood estimator

$$\Theta_{\text{GLASSO}} = \arg \min_{\Theta \succ 0} -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle_F + \lambda \|\Theta\|_{1, \text{off}}$$

■ Optimization: convex problem

Coordinate descent

Involves LASSO steps (on the rows)

Side note

■ Input:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$   
    ↓  $\mathbf{x}_i \in \mathbb{R}^d \sim \mu$   
■ Output:  $\Theta \in \mathbb{R}^{d \times d}$

# Graphical LASSO

## Gaussian Graphical Model

Gaussian assumption  $\mu = \mathcal{N}(0, \Sigma = \Theta^{-1})$

■  $\Theta_{ij} = 0 \iff$  variable  $i$  is independent of  $j$  conditionally to the others

## Penalized Maximum Likelihood estimator

$$\Theta_{\text{GLASSO}} = \arg \min_{\Theta \succ 0} -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle_F + \lambda \|\Theta\|_{1,\text{off}}$$

■ Optimization: convex problem

Coordinate descent

Involves LASSO steps (on the rows)

■ Many large scale variants:

QUIC, Big & QUIC [Hsieh & al, 2013-2014]

SQUIC [Bollhöfer, 2019] + other estimators...

Side note

■ Input:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$   
↓  $\mathbf{x}_i \in \mathbb{R}^d \sim \mu$   
■ Output:  $\Theta \in \mathbb{R}^{d \times d}$

# Graphical LASSO

Side note

Input:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$   
 $\mathbf{x}_i \in \mathbb{R}^d \sim \mu$   
Output:  $\Theta \in \mathbb{R}^{d \times d}$

## Gaussian Graphical Model

Gaussian assumption  $\mu = \mathcal{N}(0, \Sigma = \Theta^{-1})$

$\Theta_{ij} = 0 \iff$  variable  $i$  is independent of  $j$  conditionally to the others

## Penalized Maximum Likelihood estimator

$$\Theta_{\text{GLASSO}} = \arg \min_{\Theta \succ 0} -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle_F + \lambda \|\Theta\|_{1, \text{off}}$$

Optimization: convex problem

Coordinate descent

Involves LASSO steps (on the rows)

Many modelisation variants:

$\Theta = \mathcal{L}(\mathcal{G})$  is a Laplacian matrix of a graph  
[Kumar, 2020]

Many large scale variants:

QUIC, Big & QUIC [Hsieh & al, 2013-2014]

SQUIC [Bollhöfer, 2019] + other estimators...

# Graphical LASSO

Side note

Input:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$   
 $\mathbf{x}_i \in \mathbb{R}^d \sim \mu$   
Output:  $\Theta \in \mathbb{R}^{d \times d}$

## Gaussian Graphical Model

Gaussian assumption  $\mu = \mathcal{N}(0, \Sigma = \Theta^{-1})$

$\Theta_{ij} = 0 \iff$  variable  $i$  is independent of  $j$  conditionally to the others

## Penalized Maximum Likelihood estimator

$$\Theta_{\text{GLASSO}} = \arg \min_{\Theta \succ 0} -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle_F + \lambda \|\Theta\|_{1, \text{off}}$$

Optimization: convex problem

Coordinate descent

Involves LASSO steps (on the rows)

Many large scale variants:

QUIC, Big & QUIC [Hsieh & al, 2013-2014]

SQUIC [Bollhöfer, 2019] + other estimators...

Many modelisation variants:

$\Theta = \mathcal{L}(\mathcal{G})$  is a Laplacian matrix of a graph  
[Kumar, 2020]

Complexity of GLASSO:

	In memory	In time
$\hat{\Sigma}$	$\mathcal{O}(d^2)$	$\mathcal{O}(d^3)$



# | Overview of the talk

- **Part I: Finding graphs from unstructured data**
- **Part II: Schur's Positive-Definite Network**
- **Part III: The sketching approach**

# | SpodNet: Schur's Positive-Definite Network

## ■ Motivations

- Model-based vs learning-based approach

# | SpodNet: Schur's Positive-Definite Network

## ■ Motivations

■ Model-based vs learning-based approach

■ Neural network architecture for SDP learning **with structure**

NN :  $\mathcal{S}_d^{+++} \rightarrow \mathcal{S}_d^{+++}$  **focus on element-wise sparsity of output**

■ Inspired by the GLASSO solver: **unrolled architecture**

# SpodNet: Schur's Positive-Definite Network

## Motivations

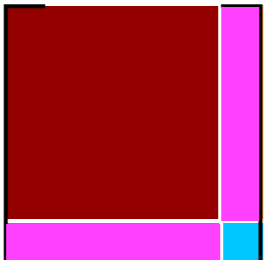
Model-based vs learning-based approach

Neural network architecture for SDP learning **with structure**

$\text{NN} : \mathcal{S}_d^{+++} \rightarrow \mathcal{S}_d^{+++}$  **focus on element-wise sparsity of output**

Inspired by the GLASSO solver: **unrolled architecture**

## Key ingredient: Schur's condition for PSDness

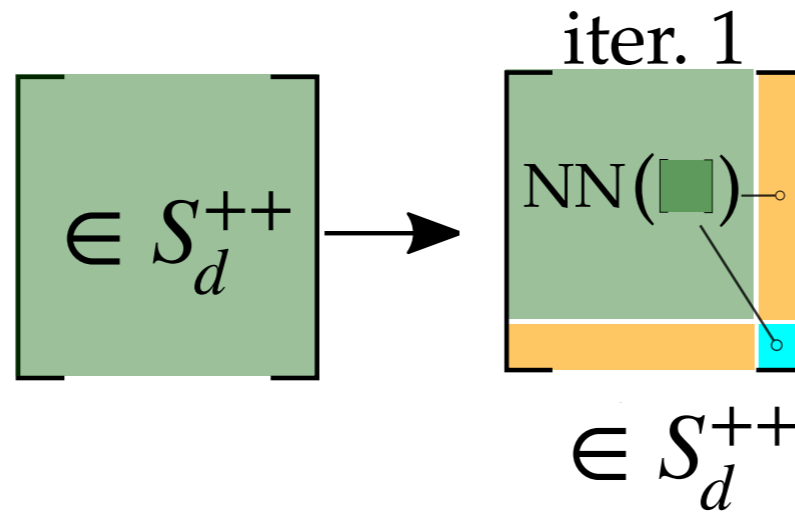
$$\Theta = \begin{pmatrix} \Theta_{11} & \theta_{12} \\ \theta_{12}^\top & \theta_{22} \end{pmatrix} \succ 0 \iff \begin{matrix} \Theta_{11} \succ 0 \\ \& \\ \theta_{22} - \theta_{12}^\top \Theta_{11}^{-1} \theta_{12} > 0 \end{matrix}$$


Holds for any value of the column !

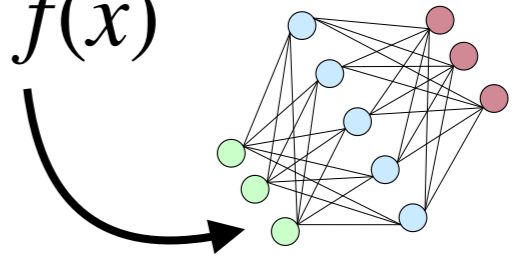
# SpodNet: Schur's Positive-Definite Network

## ■ The architecture

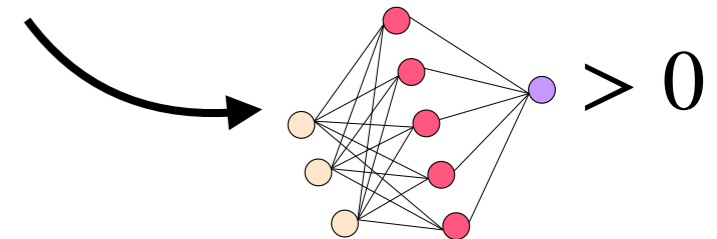
■ One layer



$$\theta_{12}^+ = f(x)$$



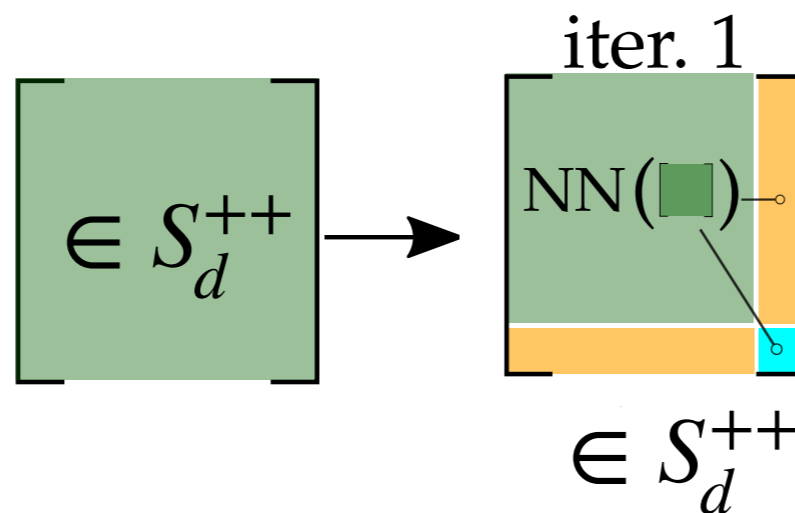
$$\theta_{22}^+ = g(y) + \theta_{12}^{+\top} \Theta_{11}^{-1} \theta_{12}^+ > 0$$



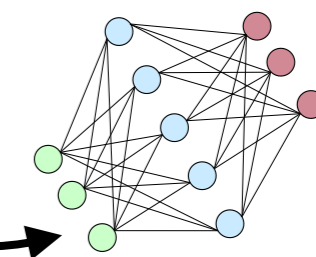
# SpodNet: Schur's Positive-Definite Network

## The architecture

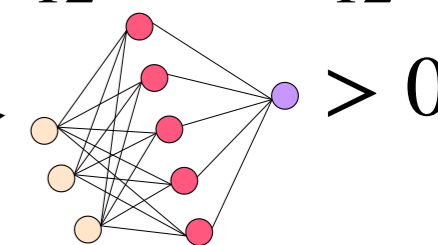
■ One layer



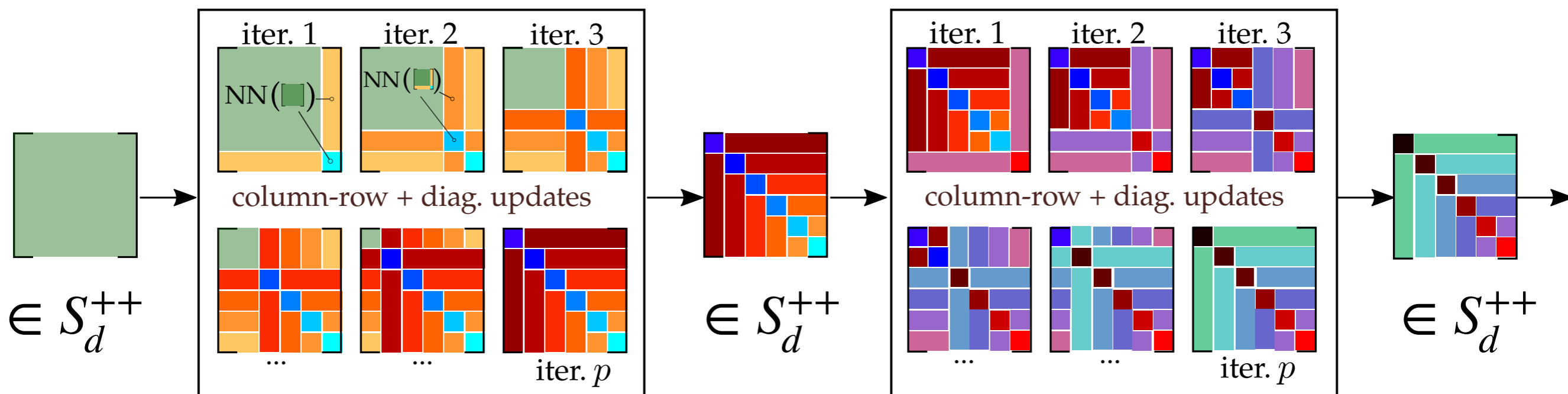
$$\theta_{12}^+ = f(x)$$



$$\theta_{22}^+ = g(y) + \theta_{12}^{+\top} \Theta_{11}^{-1} \theta_{12}^+ > 0$$



■ Multiple layers:



# | SpodNet: Schur's Positive-Definite Network

## ■ Efficient updates

■ Each iter. requires the computation of  $\theta_{12}^{+T} \Theta_{11}^{-1} \theta_{12}^{+} \longrightarrow \mathcal{O}(d^3)$

# SpodNet: Schur's Positive-Definite Network

## Efficient updates

Each iter. requires the computation of  $\theta_{12}^{+\top} \Theta_{11}^{-1} \theta_{12}^+ \longrightarrow \mathcal{O}(d^3)$

We can improve by keeping in memory  $W = \Theta^{-1}$

$$\Theta_{11}^{-1} = W_{11} - \frac{1}{w_{22}} w_{12} w_{12}^{\top}$$

$$\theta_{12}^{+\top} \Theta_{11}^{-1} \theta_{12}^+ \longrightarrow \mathcal{O}(d^2)$$



# SpodNet: Schur's Positive-Definite Network

## Efficient updates

Each iter. requires the computation of  $\theta_{12}^{+\top} \Theta_{11}^{-1} \theta_{12}^+ \longrightarrow \mathcal{O}(d^3)$

We can improve by keeping in memory  $W = \Theta^{-1}$

$$\Theta_{11}^{-1} = W_{11} - \frac{1}{w_{22}} w_{12} w_{12}^\top$$

$$\theta_{12}^{+\top} \Theta_{11}^{-1} \theta_{12}^+ \longrightarrow \mathcal{O}(d^2)$$

## Maintaining the inverse

Using Schur's inversion theorem if  $\Theta^+ = \begin{pmatrix} \Theta_{11} & \theta_{12}^+ \\ \theta_{12}^{+\top} & \theta_{22}^+ \end{pmatrix}$  is the update

$$W^+ = \begin{pmatrix} [\Theta_{11}]^{-1} + \frac{[\Theta_{11}]^{-1} \theta_{12}^+ \theta_{12}^{+\top} [\Theta_{11}]^{-1}}{g(y)} & -\frac{[\Theta_{11}]^{-1} \theta_{12}^+}{g(y)} \\ \left( -\frac{[\Theta_{11}]^{-1} \theta_{12}^+}{g(y)} \right)^\top & \frac{1}{g(y)} \end{pmatrix} \text{ satisfies } W^+ = [\Theta^+]^{-1}$$

# | SpodNet: Schur's Positive-Definite Network

- Examples of implementations

# | SpodNet: Schur's Positive-Definite Network

## ■ Examples of implementations

### ■ Unrolled Block-Graphical ISTA (UBG)

$$f : (\theta_{12}, \hat{\sigma}_{12}, w_{12}) \rightarrow \text{ST}_{\lambda^+} (\theta_{12} - \gamma^+ \times (\hat{\sigma}_{12} - w_{12}))$$

◆  $\lambda^+, \gamma^+$  are small perceptrons that learn regularization and step-size

# | SpodNet: Schur's Positive-Definite Network

## ■ Examples of implementations

### ■ Unrolled Block-Graphical ISTA (UBG)

$$f : (\theta_{12}, \hat{\sigma}_{12}, w_{12}) \rightarrow \text{ST}_{\lambda^+} (\theta_{12} - \gamma^+ \times (\hat{\sigma}_{12} - w_{12}))$$

◆  $\lambda^+, \gamma^+$  are small perceptrons that learn regularization and step-size

◆ Inspired by a proximal coordinate gradient descent step on the GLASSO objective

Expressivity +

Sparsity +++

Interpretability +++

# | SpodNet: Schur's Positive-Definite Network

## ■ Examples of implementations

### ■ Unrolled Block-Graphical ISTA (UBG)

$$f : (\theta_{12}, \hat{\sigma}_{12}, w_{12}) \rightarrow \text{ST}_{\lambda^+} (\theta_{12} - \gamma^+ \times (\hat{\sigma}_{12} - w_{12}))$$

### ■ Plug and play

$$f : (\theta_{12}, \hat{\sigma}_{12}, w_{12}) \rightarrow \Psi (\theta_{12} - \gamma^+ \times (\hat{\sigma}_{12} - w_{12}))$$

◆  $\Psi : \mathbb{R}^{d-1} \rightarrow \mathbb{R}^{d-1}$  is a learned operator (proximal?)

Expressivity ++

Sparsity ++

Interpretability ++

# SpodNet: Schur's Positive-Definite Network

## Examples of implementations

### Unrolled Block-Graphical ISTA (UBG)

$$f : (\theta_{12}, \hat{\sigma}_{12}, w_{12}) \rightarrow \text{ST}_{\lambda^+} (\theta_{12} - \gamma^+ \times (\hat{\sigma}_{12} - w_{12}))$$

### Plug and play

$$f : (\theta_{12}, \hat{\sigma}_{12}, w_{12}) \rightarrow \Psi (\theta_{12} - \gamma^+ \times (\hat{\sigma}_{12} - w_{12}))$$

### End to end

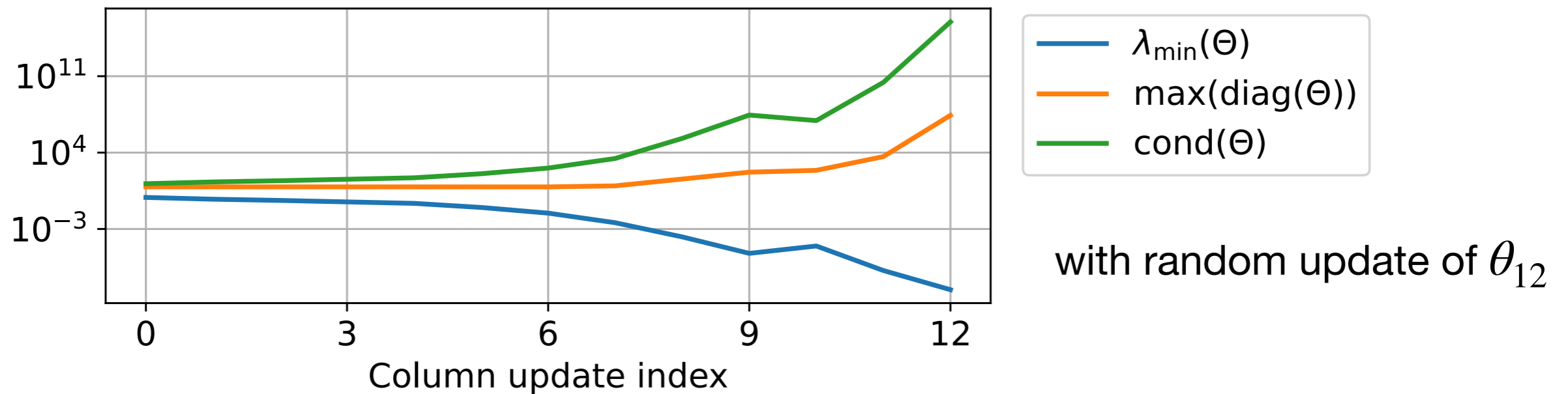
$$f : \theta_{12} \rightarrow \Phi (\theta_{12})$$

### For all architectures

$$g = \text{NN}(\theta_{22}, \hat{\sigma}_{22}, \text{schur}) > 0$$

# SpodNet: Schur's Positive-Definite Network

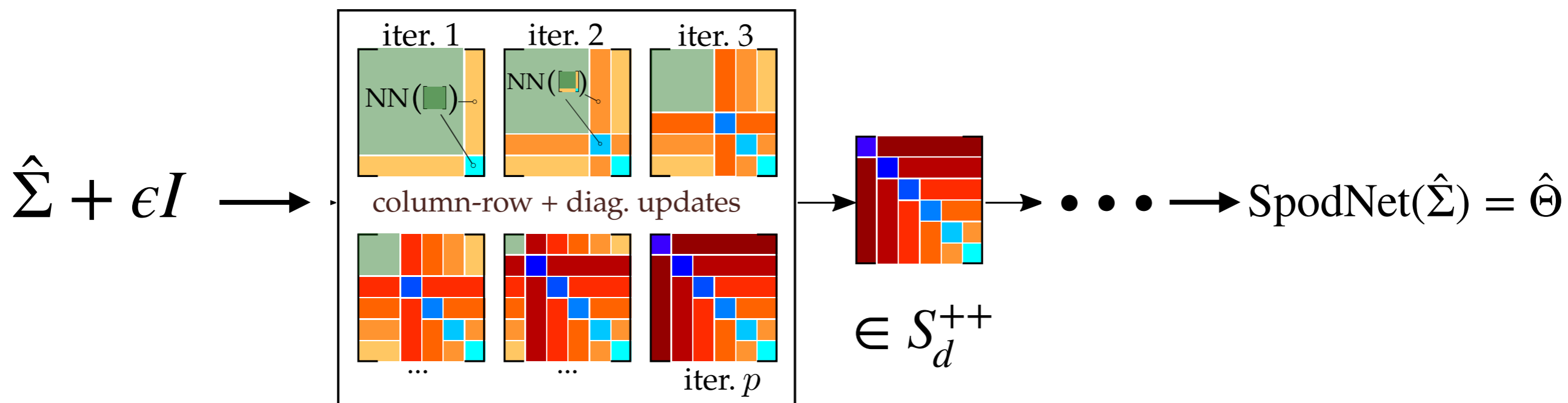
■ Nice but ...



■ One workaround: « adaptive » normalization of the columns

# SpodNet: Schur's Positive-Definite Network

## Learning with SpodNet

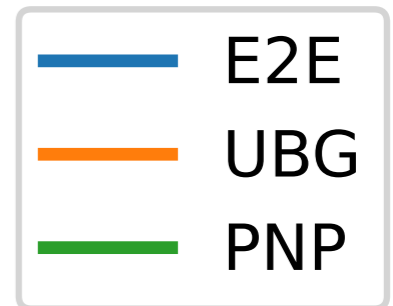
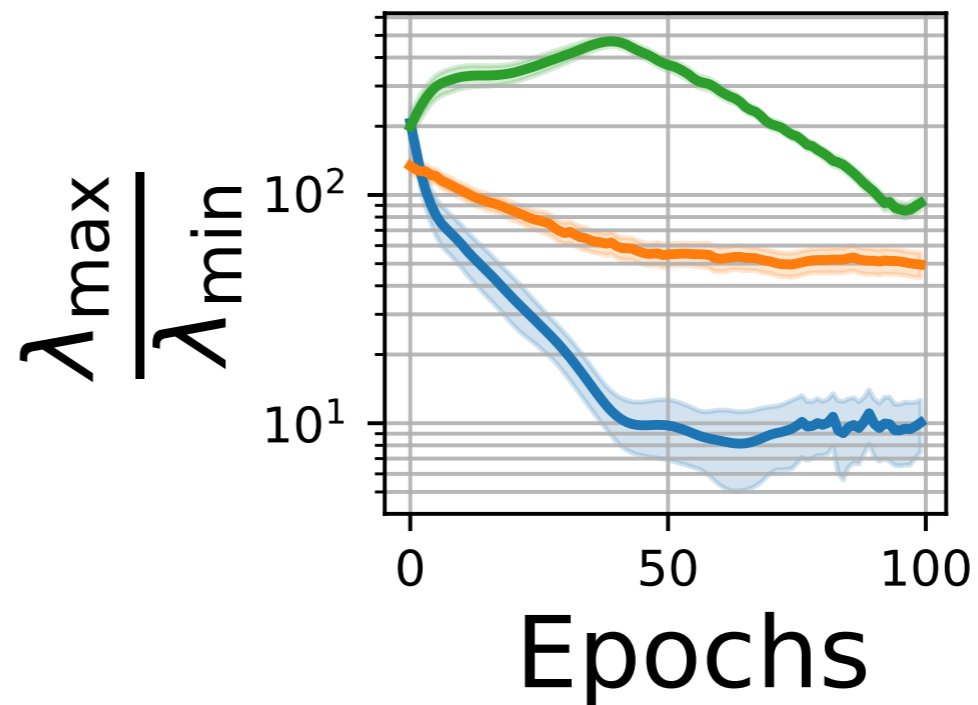
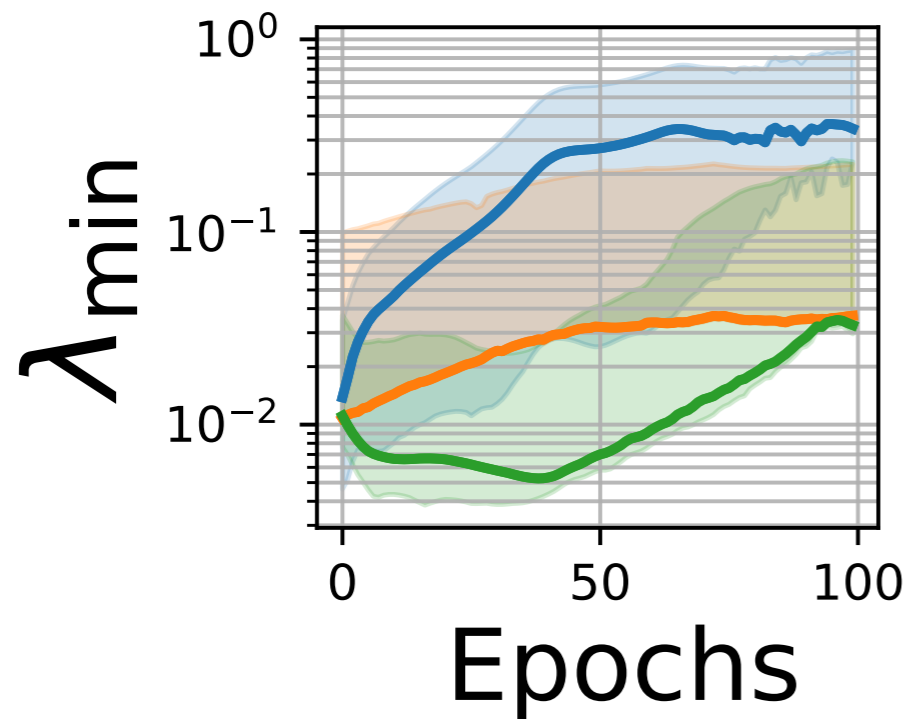


- **Data driven approach**
  - Generate  $N$  sparse PSD matrices  $(\Theta_i)_i$
  - Minimize  $L_{MSE} = \sum_{i=1}^N \|\Theta_i - \hat{\Theta}_i\|_F^2$



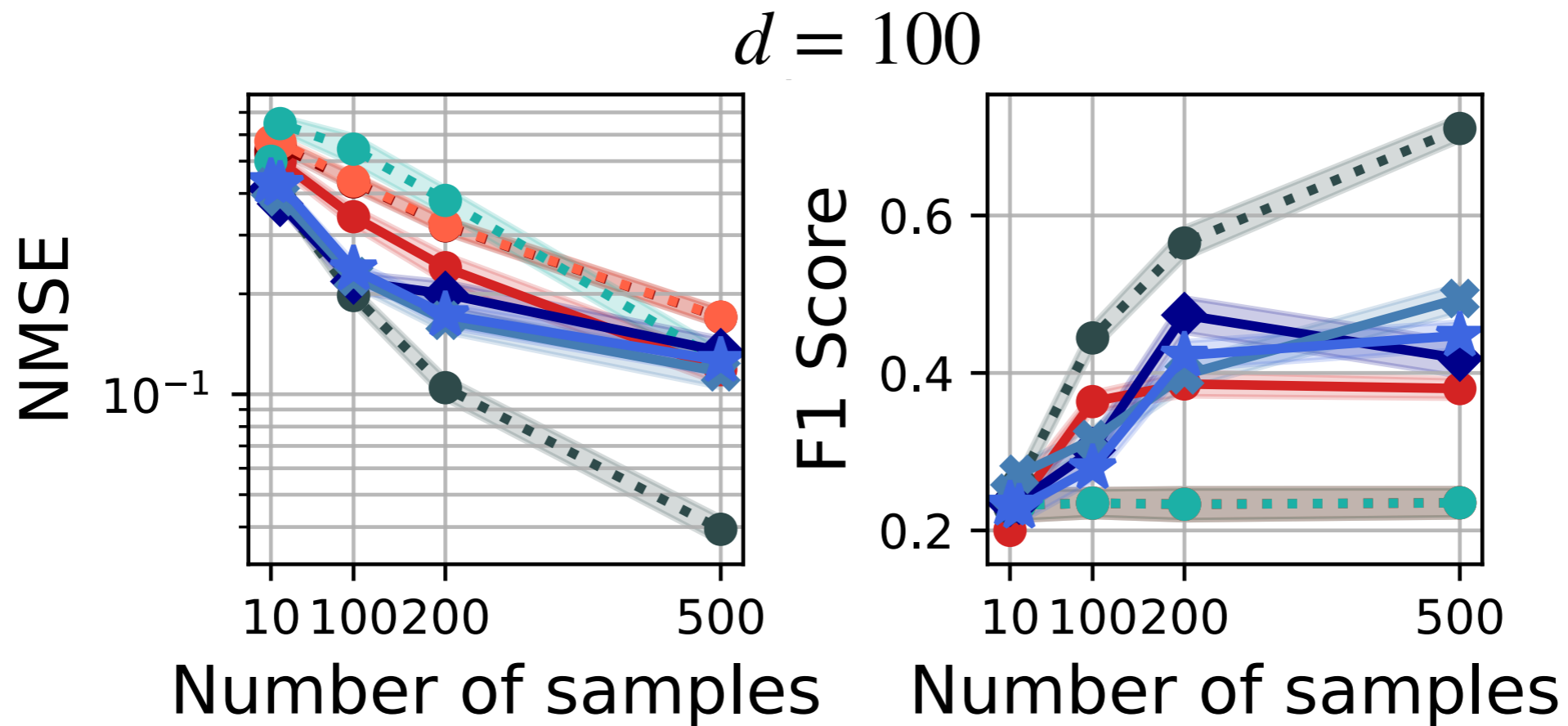
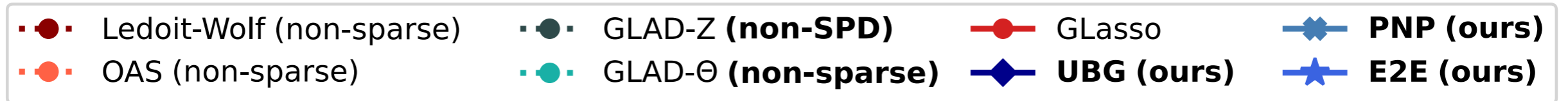
# SpodNet: Schur's Positive-Definite Network

■ Stability



# SpodNet: Schur's Positive-Definite Network

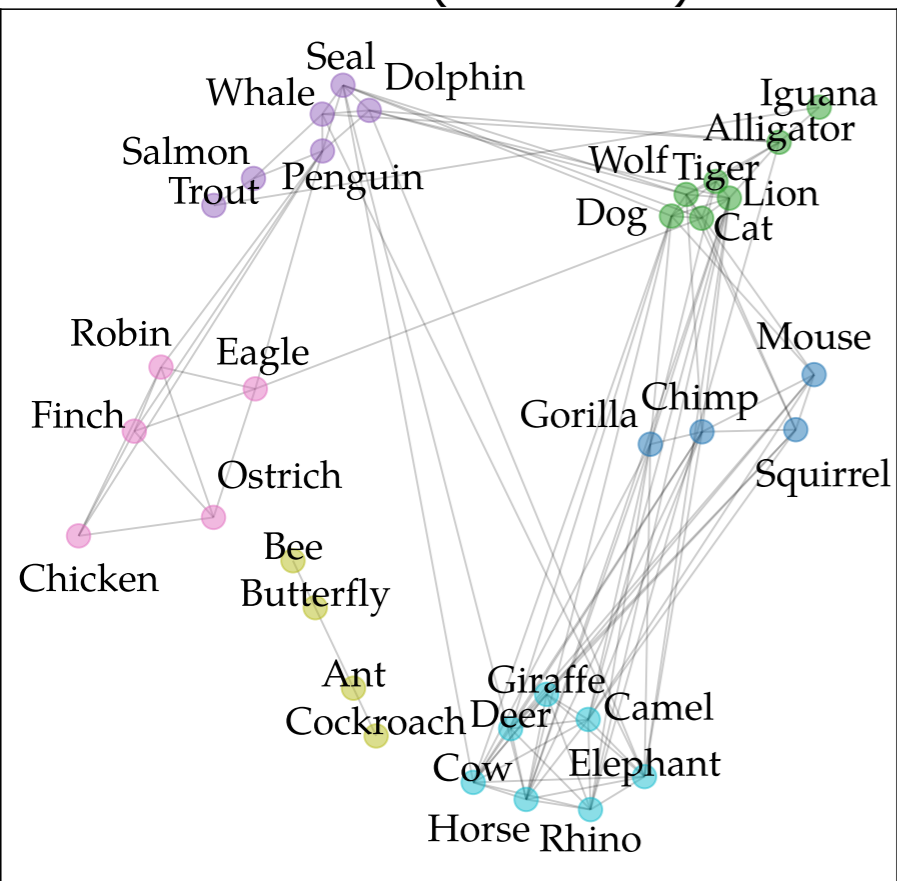
## Performances on generated data



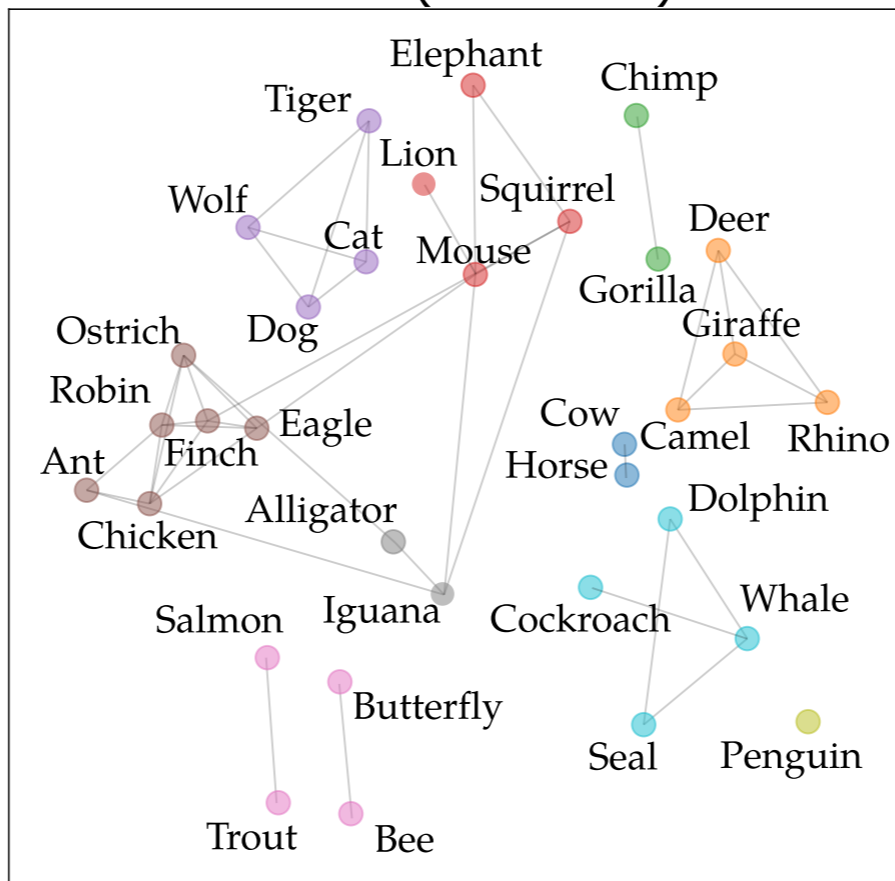
# SpodNet: Schur's Positive-Definite Network

## Performances on animals

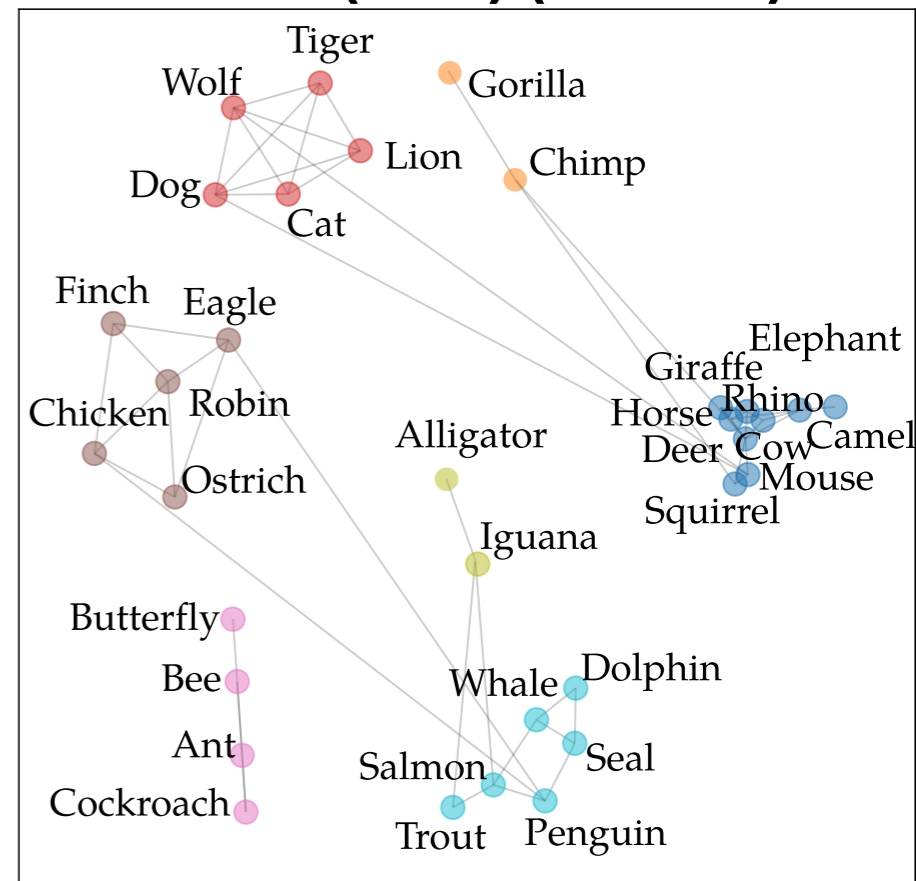
GLasso ( $m = 0.61$ )



EGFM ( $m = 0.86$ )



UBG (ours) ( $m=0.78$ )



# | SpodNet: Schur's Positive-Definite Network

## ■ Takeaways

- Architecture for SPD + XXX is hard
- SpodNet: clever column/row update maintains SPD
- Can plug any column value: additional structure (e.g. sparse)
- Numerical stability is still a pain

# | Overview of the talk

- **Part I: Finding graphs from unstructured data**
- **Part II: Schur's Positive-Definite Network**
- **Part III: The sketching approach**

# Graphical LASSO

## Penalized Maximum Likelihood estimator

$$\Theta_{\text{GLASSO}} = \arg \min_{\Theta \succ 0} -\log \det(\Theta) + \langle \hat{\Sigma}, \Theta \rangle_F + \lambda \|\Theta\|_{1,\text{off}}$$

### Optimization: convex problem

Coordinate descent

Involves LASSO steps (on the rows)

### Many large scale variants:

QUIC, Big & QUIC [Hsieh & al, 2013-2014]

SQUIC [Bollhöfer, 2019] + other estimators...

### Many modelisation variants:

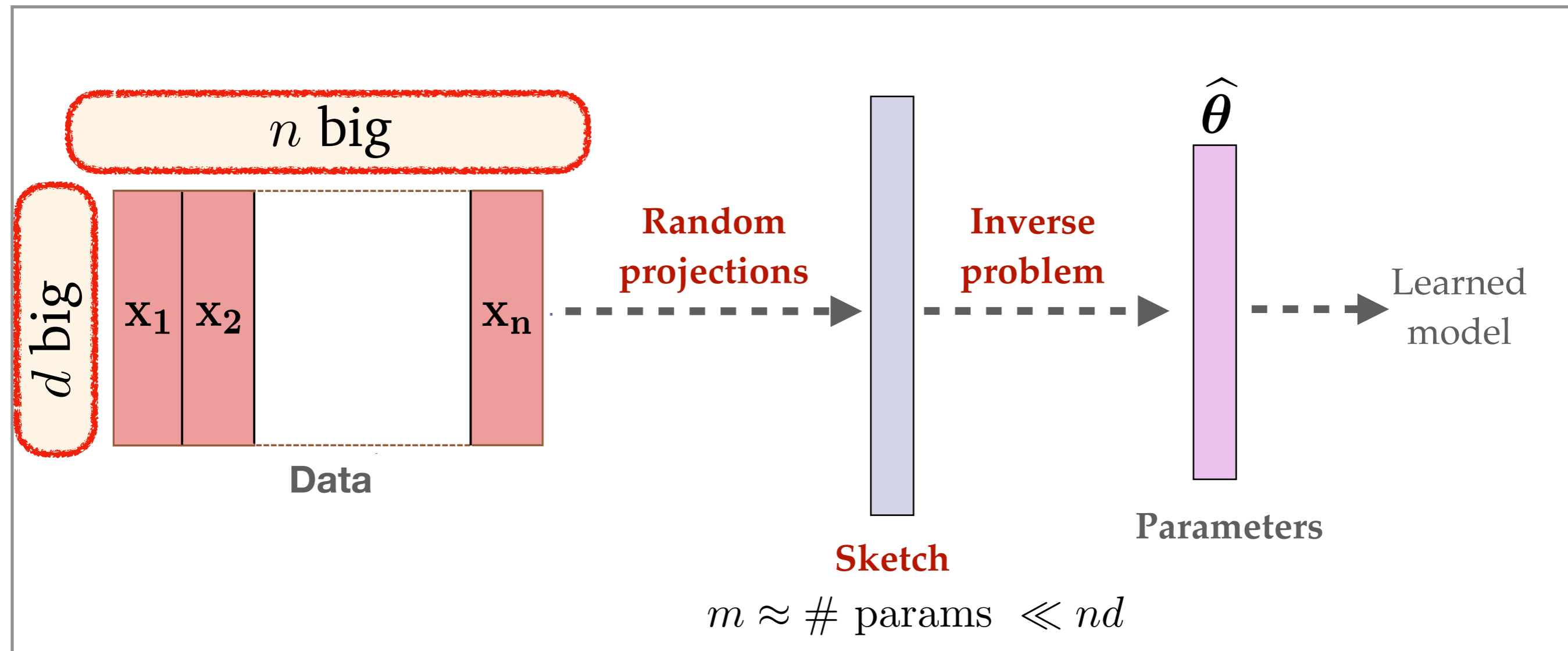
$\Theta = \mathcal{L}(\mathcal{G})$  is a **Laplacian matrix** of a graph  
[Kumar, 2020]

### Complexity of GLASSO:

	In memory	In time
$\hat{\Sigma}$	$\mathcal{O}(d^2)$	$\mathcal{O}(d^3)$

# The sketching approach

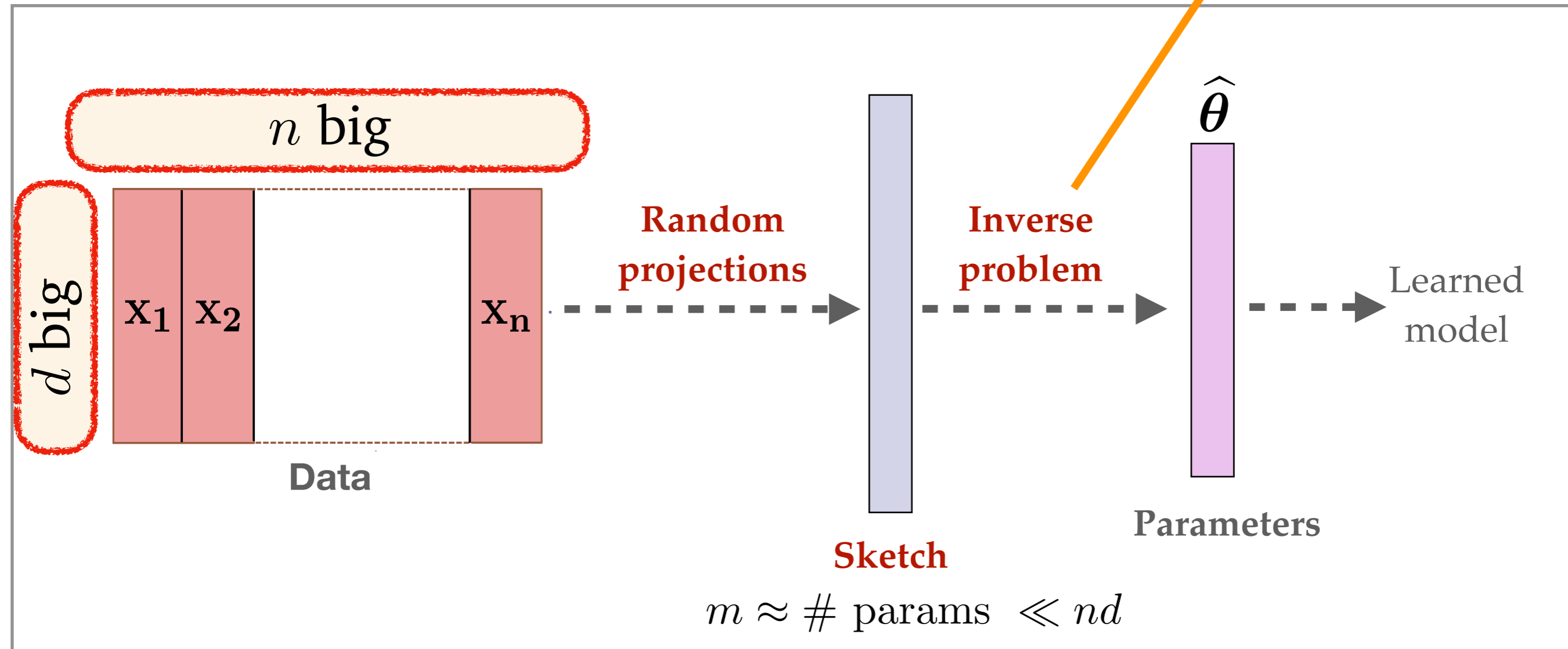
## High overview:



# The sketching approach

generalizes the principles of **compressed sensing**

High overview:



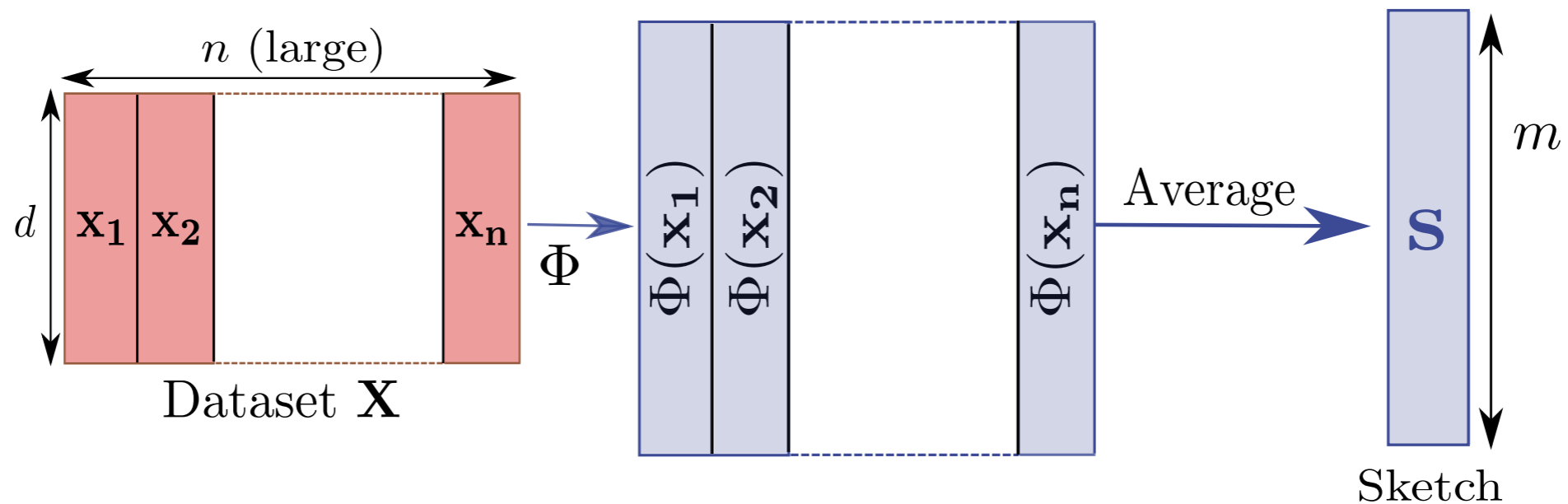


# The sketching approach

## Obtaining the sketch

■ A function called **feature operator**  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$

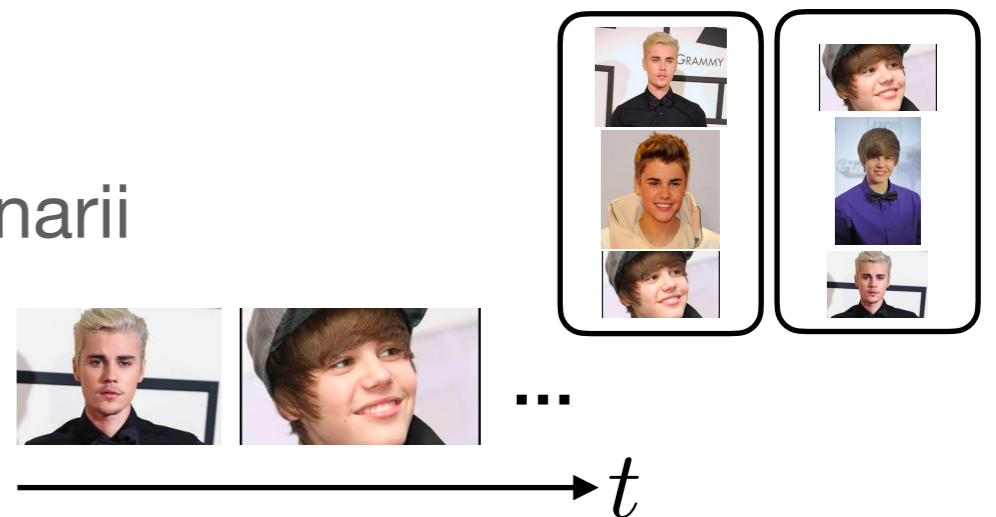
■ Averaging **n points**  $\rightarrow \mathbf{s} := \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)$



## Average is a simple idea but ...

■ Suitable for **distributed / streaming** scenarios

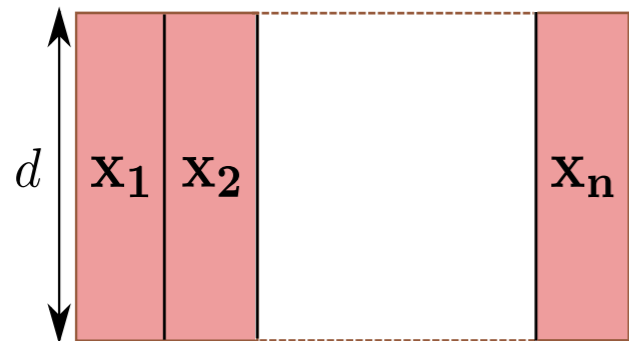
■ It can be calculated in **parallel**



# Goal of this talk

Input: a dataset

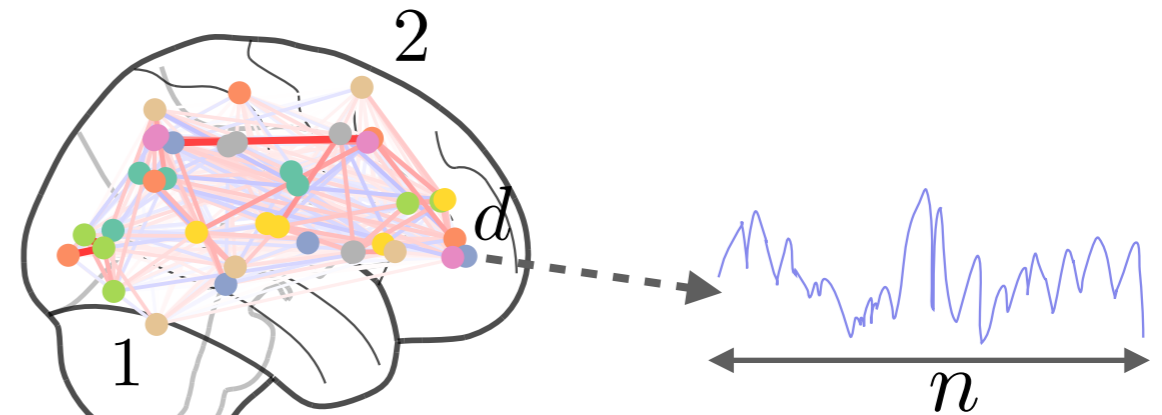
$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$$



$$\mathbf{x}_i \in \mathbb{R}^d \sim \mu$$

GLASSO

Output: graph of relations between the  $d$  variables  $\Theta \in \mathbb{R}^{d \times d}$

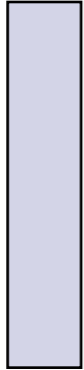


GLASSO

	In memory	In time
$\hat{\Sigma}$	$\mathcal{O}(d^2)$	$\mathcal{O}(d^3)$

# Goal of this talk

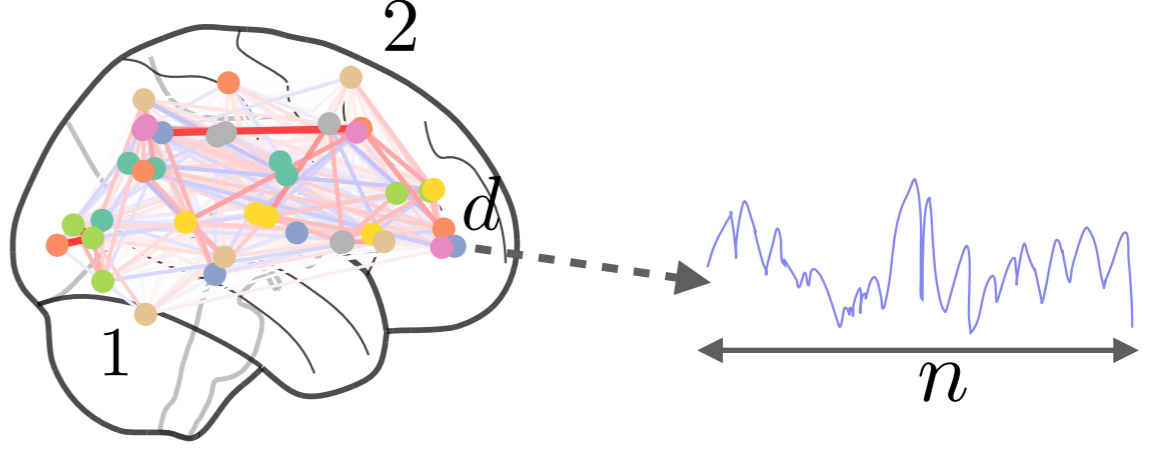
Input: a dataset

$$\mathbf{s} = \text{Sketch}(\mathbf{X}) \in \mathbb{R}^m$$


$m \approx \# \text{ edges}$

**Sketch**

Output: graph of relations between the  $d$  variables  $\Theta \in \mathbb{R}^{d \times d}$



## GLASSO

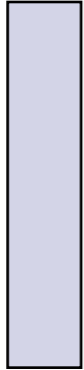
	In memory	In time
$\hat{\Sigma}$	$\mathcal{O}(d^2)$	$\mathcal{O}(d^3)$

## Sketching

	In memory	In time
$\mathbf{s}$	$\mathcal{O}(m) \ll d^2$	$\mathcal{O}(?)$

# Goal of this talk

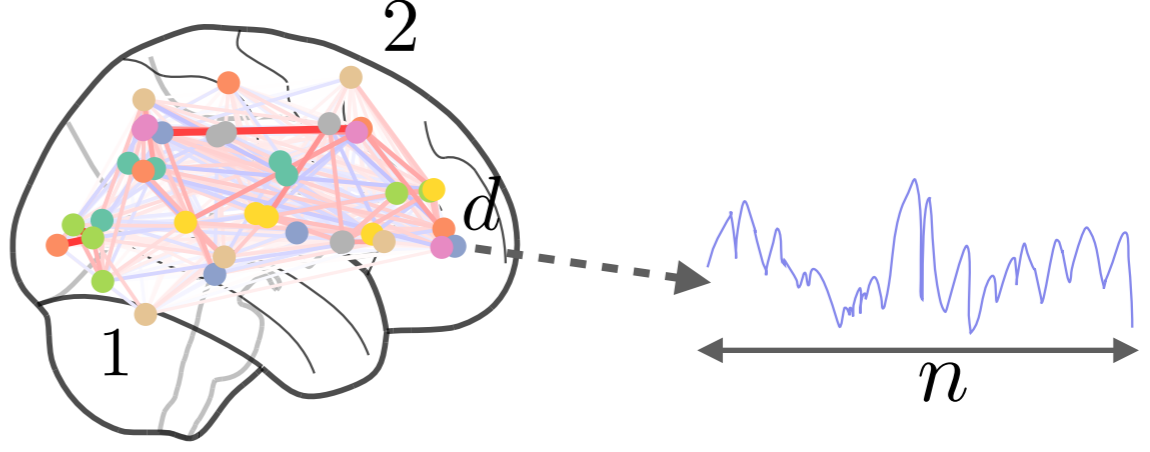
Input: a dataset

$$\mathbf{s} = \text{Sketch}(\mathbf{X}) \in \mathbb{R}^m$$


$m \approx \# \text{ edges}$

**Sketch**

Output: graph of relations between the  $d$  variables  $\Theta \in \mathbb{R}^{d \times d}$



## GLASSO

In memory	In time
$\hat{\Sigma} \dashrightarrow \mathcal{O}(d^2)$	$\mathcal{O}(d^3)$

## Sketching

In memory	In time
$\mathbf{s} \dashrightarrow \mathcal{O}(m) \ll d^2$	$\mathcal{O}(?)$

## Why should it work ?

- The underlying graph is **sparse**
- Keep only what we need through the sketch

# | Towards theoretical compressive recovery

■ The feature operator  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$

# Towards theoretical compressive recovery

■ The feature operator  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$

■ In this talk: quadratic measurements  $\mathbf{A}_j \sim \Lambda$  is a random matrix

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{m}} (\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_m \mathbf{x})^\top$$

# Towards theoretical compressive recovery

■ The feature operator  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$

■ In this talk: quadratic measurements  $\mathbf{A}_j \sim \Lambda$  is a random matrix

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{m}} (\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_m \mathbf{x})^\top$$

Gaussian measurements

$$\mathbf{A}_j \underset{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_{d \times d})$$

# Towards theoretical compressive recovery

■ The feature operator  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$

■ In this talk: quadratic measurements  $\mathbf{A}_j \sim \Lambda$  is a random matrix

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{m}} (\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_m \mathbf{x})^\top$$

Gaussian measurements

$$\mathbf{A}_j \underset{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_{d \times d})$$

Structured rank-one

End of presentation

Rank-one measurements

$$\mathbf{A}_j = \mathbf{a}_j \mathbf{a}_j^\top \quad \mathbf{a}_j \underset{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$$

$$\Phi(\mathbf{x}) = (|\langle \mathbf{a}_j, \mathbf{x} \rangle|^2)_{j \in \llbracket m \rrbracket}$$

■ Inspired by works on low-rank matrix completion



# Towards theoretical compressive recovery

■ The feature operator  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$

■ In this talk: **quadratic measurements**  $\mathbf{A}_j \sim \Lambda$  is a **random matrix**

$$\Phi(\mathbf{x}) = \frac{1}{\sqrt{m}} (\mathbf{x}^\top \mathbf{A}_1 \mathbf{x}, \dots, \mathbf{x}^\top \mathbf{A}_m \mathbf{x})^\top$$

**Gaussian measurements**

$$\mathbf{A}_j \underset{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_{d \times d})$$

**Structured rank-one**

End of presentation

**Rank-one measurements**

$$\mathbf{A}_j = \mathbf{a}_j \mathbf{a}_j^\top \quad \mathbf{a}_j \underset{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_d)$$

$$\Phi(\mathbf{x}) = (|\langle \mathbf{a}_j, \mathbf{x} \rangle|^2)_{j \in \llbracket m \rrbracket}$$

■ Inspired by works on low-rank matrix completion

■ Defined a **linear op.** on **symmetric matrices**  $\mathcal{A} : S_d \rightarrow \mathbb{R}^m$

$$\mathcal{A}\mathbf{S} = \frac{1}{\sqrt{m}} (\langle \mathbf{A}_j, \mathbf{S} \rangle_F)_{j \in \llbracket m \rrbracket}$$

$$\mathbf{s} = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i) = \mathcal{A} \hat{\Sigma} \in \mathbb{R}^m$$

Emp. cov. 

# Summary

