# JOINT ESTIMATION OF MULTIPLE GRANGER GRAPHICAL MODELS USING NON-CONVEX PENALTIES

Thesis Proposal

Parinthorn Manomaisaowapak

Advisor: Assist. Prof. Dr. Jitkomut Songsiri

Department of Electrical Engineering, Faculty of Engineering

Chulalongkorn University

1

# OUTLINE

- Introduction
- Overview
- Related works
- Work plan
- Background
- Methodology
- Preliminary results
- Future works

## Graphical representation

**How to study relationship of variables?**

Causality network    Causality matrix



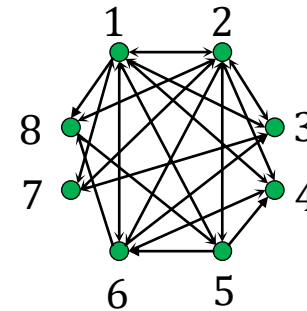**Causality analysis**

*A strength of evidence*

**Granger graphical model**

**Granger causality(GC)**

→ Based on dynamical models

→ Directly related to sparsity of model coefficient

Graphical representation
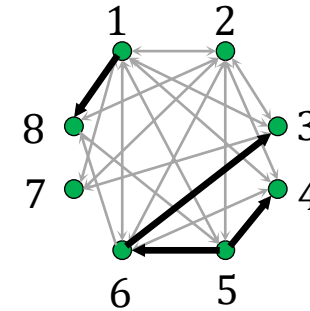
# High dimensional GC network

- GC network has **large amount of connections**

- We aim to extract only **significant connections**

**Sparse representation** of GC network
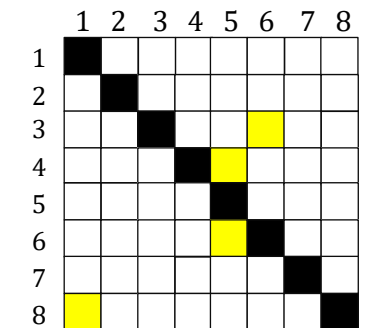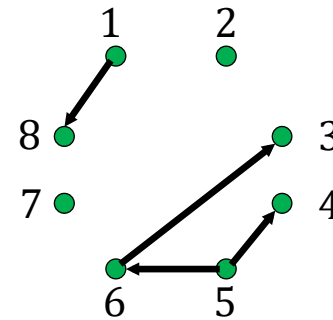
Causality network    Causality matrix

**Sparse estimation formulation in general form.**

Model parameter

$$\min_{\theta} f(\theta) + \lambda g(\theta)$$

Fitting term

Sparsity inducing penalty

$\lambda \uparrow$

**Dense**

**Sparse**

consider when the same multivariate time-series are measured in different settings



analyze each **setting** separately → the relations will show if the samples are large enough

jointly analyze **all setting** → allow low sample size estimation By adding **prior information**

channel

setting

**Joint estimation of multiple models**

# Goal: Find important connections of multiple networks with prior knowledge

$$\min_{\theta_1,\dots,\theta_K} \sum_i [f_i(\theta_i) + \lambda_1 h_i(\theta_i)] + \lambda_2 g(\theta_1,\dots,\theta_K)$$

where $h_i$ aims to promote differential sparsity in each model. ●

$\boxed{g}$ aims to promote $\boxed{\text{common sparsity}}$ across all models. ●
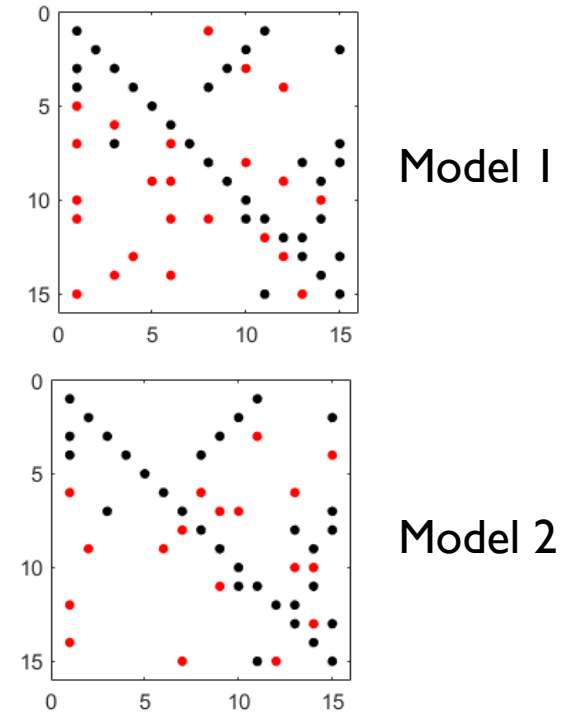
Require **definition of similarity**
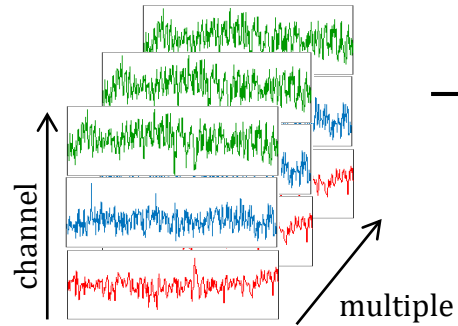
→ Sparsity inducing function

Example

group lasso ⟶ $\theta_1,\dots,\theta_K$ has **same** non-zero **pattern**

fused lasso ⟶ $\theta_m - \theta_\ell$ is sparse ⟶ **Some** model **coefficients are identical**

Model 1

Model 2

7

**Objectives**



Multiple multivariate time-series

Formulation C → identical GC networks

Formulation D → common GC network & differential network

Formulation S → common GC network with identical strength & differential network

- To propose three formulations. The formulations are
  - Formulation C: The estimated networks have an **identical sparsity pattern**
  - Formulation D: The estimated networks have some common parts and some different parts.
  - Formulation S: The estimated networks have some common parts and some different parts. The common parts also **share model parameters**.

- To provide efficient numerical methods for solving the proposed estimation methods in a large-scale setting.

8

## Scope of work

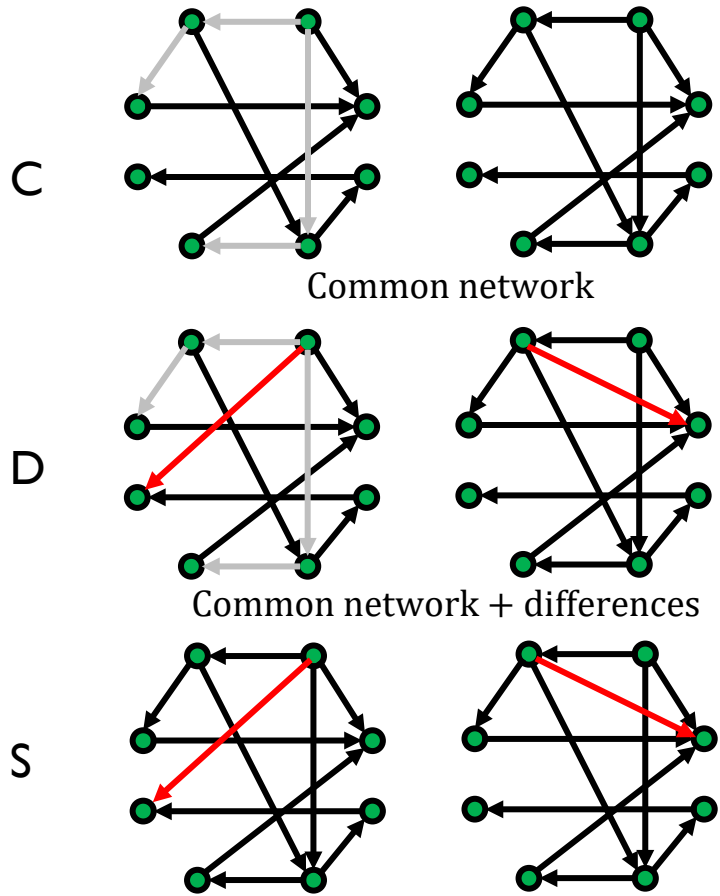- The proposed framework will be verified intensively in a simulated data sets and one real-world data set
- The usefulness of the methods will be illustrated on brain network application

## Expected outcome

- Estimation formulations of multiple Granger graphical models
- A computer program that has input as a set of multivariate time-series and return group and individual Granger graphical model of the multiple time-series

**RELATED WORKS**

**non-convex group penalties [Our work]**

**[Songsiri, 2017]**  group lasso, fused lasso

**extension**

C — Common network

D — Common network + differences

S — Common network has identical strength + differences

[Wilms, 18]  group lasso

[Gregorova, 2015]  group lasso+Tikhonov

[Skripnikov, 2019]  group lasso, two-stages

**[Guo, 2011],**
**[Chun, 2015]**  **non-convex penalty** Gaussian graphical model

[Skripnikov, 2019]  sparse fused-lasso

[Tuck, 2020]  sparse fused-lasso Gaussian graphical model

**[Bore, 2020]**  **non-convex penalty** single Granger model

**non-convex group norm penalty**
$$\sum \|x_{G_i}\|_p^q \quad p \geq 1, 0 < q < 1$$
**[Hu, 2017]**

**non-convex penalty**
$$\sum |x_i|^q \qquad 0 < q < 1$$
**[Chartrand, 2008]**

**Vector autoregressive model (VAR)**

$$y(t) = \sum_{r=1}^{p} A_r y(t-r) + \eta(t)$$

$$A_r \in \mathbf{R}^{n \times n} \qquad y = (y_1, \dots, y_n) \in \mathbf{R}^n$$



$A_2$
$A_1$

**Granger causality on VAR models**

- Granger causality(GC, $F_{ij}$) is a strength of evidence

- Absence of GC connection can be investigated by the relation

$$F_{ij} = 0 \Leftrightarrow (A_r)_{ij} = 0; r = 1, 2, \dots p \quad \text{[Granger, 1980]}$$

Stacked lags

1
2
3

1　2　3

**Sparsity inducing penalty can be designed using this prior knowledge**

12

$K$ multivariate time-series

Joint estimation formulation

Model selection

Estimated networks

vary sparsity

13

Expected outcome

We proposed three formulations,

Formulation **C**    **Common**

- **Common** pattern.

Induce block sparsity by group norm penalty

$F_{12}^{(1)}$  $F_{12}^{(2)}$

$F_{13}^{(2)}$

$F_{13}^{(1)}$

$F_{21}^{(1)}$

Stacked VAR coefficient matrix

Formulation **D**    **Differential**

- **Common** pattern
- **Different** pattern

Formulation **S**    **Similar**

- **Shared VAR coefficients** value in all models
- **Different** pattern

**C**

**D**

**S**



Model 1            Model 2

14

We used BIC criteria to find optimal tuning-parameters

$$BIC(\lambda_1, \lambda_2) = -2\,\mathcal{L}(\lambda_1, \lambda_2) + \log(N) \cdot \mathrm{df}(\lambda_1, \lambda_2)$$

Log-likelihood of VAR model.
(Fitness of models)

Effective degree of freedom
(Complexity of models).

# off-diagonal nonzero estimated parameters

There are other choices.

15

Problem properties (Formulation C, D, S)

- The problem is in the form of $\min_{x} f(x) + g(x)$

- $\nabla f$ is Lipschitz-continuous.
- Function $g$ is not differentiable at zero while we prefer sparse solutions
- We aim to solve high-dimensional problem or in a large-scale setting.

$\downarrow$

**First order** algorithm should be considered first

$\downarrow$

**Proximal gradient methods** **unify the framework that solve this problem**

ISTA    GISA

FISTA    GIST

Iterative hard-thresholding algorithm    Half-thresholding algorithm

## Proximal algorithms

- require evaluation of **<u>proximal operator</u>**

  Definition: proximal operator of function $g$

  $$\text{prox}_{\alpha g}(v) = \underset{x}{\text{argmin}} \; g(x) + \frac{1}{2\alpha} \|x - v\|_2^2$$

- are widely used in sparse estimation using lasso, group lasso for a convex case

- proximal operator has a closed-form expression for some functions, such as

$\ell_1$ norm $\qquad (\text{prox}_{\lambda\|x\|_1}(v))_i = \text{sign}(v_i)\max\{0, v_i - \lambda\}$ $\qquad$ Soft thresholding operator

$\ell_2$ norm $\qquad \text{prox}_{\lambda\|x\|_2}(v) = \max\left\{0, 1 - \frac{\lambda}{\|v\|_2}\right\}v$ $\qquad$ Block-soft thresholding operator

proximal gradient methods

Alternating direction methods of multipliers(ADMM)

solve $\min\limits_{x} F(x) = f(x) + g(x)$

solve $\min\limits_{x,z} F(x) = f(x) + g(z) + \rho\|Ax + Bz - c\|_2^2$
subjected to
$Ax+Bz=c$

require proximal operator of
$g(x) = h_1(L_1 x) + h_2(L_2 x)$

set $A = \begin{bmatrix} L_1 \\ L_2 \end{bmatrix}, B = -I, c = 0$

related work: **SDMM**
[Combettes, 2011]

set $A = \begin{bmatrix} I \\ L_1 \\ L_2 \end{bmatrix}, B = -I, c = 0$

Higher
computation complexity
than ours

only $h_1(x), h_2(x)$
have closed-form proximal operator
[Hu, 2017]

Evaluation of $\text{prox}_{\alpha g}(v)$
Reduce to $\text{prox}_{\alpha h_1}(v), \text{prox}_{\alpha h_2}(v)$

**no convergence guaranteed in non-convex formulations**

**a convergence to critical point** can be obtained
by **selecting a proper penalty parameter $\rho$**

adaptive $\rho$ may solve the problem ⟶ **Spectral ADMM [Xu, 2017]**

19

# PRELIMINARY RESULTS

**Problem formulation**

**Literature review**

**Implementation**

| convex | non-convex |
|---|---|
| **C** | **C** |
| **D** | nmAPG |
| **S** | **D** |
| ADMM | **S** |
| | Spectral ADMM |

**Thesis writing** & **Publication**

**Experiments**

- Algorithm hyperparameters selection and testing

  Experiment 1: initial point selection for nmAPG

- Efficiency evaluation of numerical methods

  Experiment 2: efficiency in linear regression model

  Experiment 3: VAR time-series generation

- Effectiveness of formulations

  Experiment 4: Common Granger network extraction

  Experiment 5: Classification

  Experiment 6: Effectiveness of formulation D

  Experiment 7: Effectiveness of formulation S

- Brain network application

  Experiment 8: Application on fMRI time-series

# Performance index

- Area under ROC curve

$$FPR = \frac{\text{False positive}}{\text{\# Negative}}$$

$$TPR = \frac{\text{True positive}}{\text{\# Positive}}$$

- Relative parameter bias

  - $\dfrac{\|\hat{x} - x_{\text{true}}\|_2}{\|x_{\text{true}}\|_2}$



$\lambda = 0$

$(FPR, TPR)_{\lambda=0}$

$\lambda \uparrow$

$(FPR, TPR)_{\lambda\uparrow}$

TP

TN

FN

FP

$\lambda_c$

Estimation

$(FPR, TPR)_{\lambda_c}$

Ground-truth

Estimated GC matrix

AUC

TPR

FPR

## Experiment 3: VAR time-series with prespecified GC patterns generation

Objective: To test the formulations with known given structure

We randomized stable VAR coefficients that the Granger causality patterns are

1. **Common type** ground truth        2. **Differential type** ground truth        3. **Similar type** ground truth

Examples of generated GC matrix topology



**Common network** density and **differential network** density can be set.

Experiment 4: Group level Granger network extraction

Objective: To extract common GC network with a presence of heterogeneous connections



channel

Formulation C

Extracted common network

$K$ models

Groundtruth networks

23

## Experiment 4: Group level Granger network extraction

Objective: To extract common GC network with a presence of heterogeneous connections

- 4 **sets of** 15-dimensional 2nd-order-VAR models
- Common density : 10%, 20%
- Differential density : 5%



Common density = 10%

Common density = 20%

Generate time-series
with unit variance Gaussian noise.

Experiment 4: Group level Granger network extraction



Common density : **0. 1**

Common density : **0. 2**

# Experiment 5: Supervised-classification using learned common Granger network

Objective: To illustrate the application of common Granger network extraction



Class #1 GC network extraction

testing time series

channel

channel

Class #1's time-series

GC network templates learned from joint estimation

Likelihood

Class 1    Class 2    Class 3    Class 4

Experiment 5: Supervised-classification using learned common Granger network

Setting

- 10 GC networks defined on $2^{nd}$ order 15-dimensional VAR models

- The GC matrix of classifying time-series has sparsity pattern same as one of classes

- Common network density is set to 20%

- vary VAR lag order to test the performance when model order is wrongly chosen

## Experiment 5: Supervised-classification using learned common Granger network

Result



- Near perfect classification rate in non-convex case
- Non-convex case did not deteriorate much when model order is wrong compared to convex case.

28

## Experiment 6: Performance of differential priors

Objective: To illustrate the performance of formulation D

Setting

- 4 sets of 15-dimensional 2nd-order-VAR models
- Common network density is set to 20%
- Differential network density is set to 5%
- The ground-truth types are **common type**, **differential type**, **similar type**



common type                    differential type                    similar type

29

## Experiment 6: Performance of differential priors



Ground truth type

C        D        S

non-convex formulation D
using
Group norm penalty

convex formulation D
using
Group lasso

Overall ROC

**Experiments**

- Effectiveness of formulations
    Experiment 4: Common Granger network extraction
    Experiment 5: Classification
    Experiment 6: Effectiveness of formulation D
    Experiment 7: Effectiveness of formulation S
- Brain network application
    Experiment 8: Application on fMRI time-series

**Thesis writing** & **Publication**

"Learning A Common Granger Causality Network Using A Non−Convex Regularization", ICASSP-2020

Formulation C (non-convex)

## **Goal**

- Control the convergence of ADMM algorithm to solve formulation D, S
- Increase performance of algorithms
- Apply formulations on fMRI data

# Q&A

# REFERENCES

[Bore20]       J. C. Bore, P. Li, D. J. Harmah, F. Li, D. Yao, P. Xu, Directed EEG neural network analysis by LAPPS (p≤1) Penalized sparse Granger approach, Neural Networks, Volume 124, 2020, Pages 213-222,

[Boyd11]       S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers", Foundation and Trends in Machine Learning, vol. 3, no. 1, pp. 1-122, Jan. 2011.

[Chartrand08]R. Chartrand, V. Staneva,. (2008). Restricted isometry properties and nonconvex compressive sensing. Inverse Problems 24(3).

[Chun15]       H. Chun, X. Zhang, and H. Zhao, "Gene regulation network inference with joint sparse Gaussian graphical models," Journal of Computational and Graphical Statistics, vol. 24, no. 4, pp. 954–974, 2015.

[Combettes11]P. L. Combettes and J. C. Pesquet. Proximal Splitting Methods in Signal Processing, pages 185-212. Springer New York, New York, NY, 2011.

[Granger1980]C.W.J. Granger, Testing for causality: A personal viewpoint, Journal of Economic Dynamics and Control, Volume 2, 1980, Pages 329-352, ISSN 0165-1889,

[Gregorova15]M. Gregorova, A. Kalousis, and S. Marchand-Maillet. Learning coherent Granger causality in panel vector autoregressive models. In Proceedings of the Demand Forecasting Workshop of the 32nd International Conference on Machine Learning. ICML, 2015.

[Guo11]       J. Guo, E. Levina, G. Michailidis, and J. Zhu, "Joint estimation of multiple graphical models," Biometrika, vol. 98, no. 1, pp. 1–15, 2011.

[Hu17]       Hu, C. Li, K. Meng, J. Qin, and X. Yang, "Group sparse optimization via $\ell_{p,q}$ regularization," Journal of Machine Learning Research, vol. 18, no. 30, pp. 1–52, 2017.

# REFERENCES

[Huang15]    F. Huang and S. Chen, "Joint learning of multiple sparse matrix Gaussian graphical models," IEEE Transactions on Neural Networks and Learning Systems, vol. 26, no. 11, pp. 2606–2620, 2015.

[Li15]    H. Li and Z. Lin, "Accelerated proximal gradient methods for nonconvex programming," in Advances in Neural Information Processing Systems 28, pp. 379–387, 2015.

[Skrip19]    A. Skripnikov and G. Michailidis, "Joint estimation of multiple network Granger causal models," Econometrics and Statistics, vol. 10, pp. 120–133, 2019.

[Skrip19]    A. Skripnikov and G. Michailidis, "Regularized joint estimation of related vector autoregressive models," Computational Statistics & Data Analysis, vol. 139, pp. 164–177, 2019.

[Songsiri 17]    J. Songsiri. Estimations in Learning Granger Graphical Models with Application to fMRI Time Series. Technical report, Chulalongkorn University, Department of Electrical engineering, July 2017.

[Teboulle18]    M. Teboulle. A simplified view of first order methods for optimization. Math. Program. 170, 1 (2018), 67-96.

[Tuck20]    J. Tuck and S. Boyd. Fitting Laplacian regularized stratified Gaussian models. ArXiv, abs/2005.01752, 2020.

[Wang19]    Y. Wang, W. Yin, J. Zeng. Global Convergence of ADMM in Nonconvex Nonsmooth Optimization. Journal of Scientific Computing 78, 29–63 (2019).

[Wilms18]    I. Wilms, L. Barbaglia, and C. Croux, "Multiclass vector auto-regressive models for multistore sales data," Journal of the Royal Statistical Society: Series C (Applied Statistics), vol. 67, no. 2, pp. 435–452, 2018.

[Xu17]    Z. Xu, M. Figueiredo, T. Goldstein, "Adaptive ADMM with Spectral Penalty Parameter Selection," Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, PMLR 54:718-727, 2017.

**PHASE ONE – PHASE TWO – PHASE THREE (June 2019 – May 2020)**

| Objective | TASK | Scheduled period (shaded weeks) |
|---|---|---|
| **1** | **Literature review on sparse formulation in graphical models** | |
| 1.1 | Joint sparse estimation of Gaussian graphical models | June 2019 |
| 1.2 | Joint sparse estimation of Granger graphical models | June–July 2019 |
| 1.3 | Survey of non-smooth, non-convex optimization algorithms | August 2019; January–May 2020 |
| 1.5 | Survey of non-convex regularization | July–August 2019; January–May 2020 |
| **2** | **Algorithm implementation & Experiment design** | |
| 2.1 | ADMM algorithm implementation for non-convex optimization | September 2019 |
| 2.2 | nmAPG algorithm implementation | September 2019 |
| 2.3 | Coding optimization | September 2019 – May 2020 |
| 2.4 | Experiments planning | October 2019 |
| **3** | **Perform simulated data experiments** | |
| 3.1 | Experiment 1: initial point selection in linear regression model | September 2019 |
| 3.2 | Experiment 2: non-convex group norm regularization in linear regression models | September 2019 |
| 3.3 | Experiment 3: VAR time-series with pre-specified GC patterns generation | June 2019 |
| 3.4 | Experiment 4: Common Granger network extraction | September–October 2019 |
| 3.5 | Experiment 5: Supervised-classification using learned common Granger network | September–October 2019 |
| 3.6 | Submit part of the proposal work to ICASSP | October 2019 |
| 3.7 | Experiment 6: Effectiveness of differential priors | January–February 2020 |
| 3.8 | Preparing proposal | February–May 2020 |

**PHASE FOUR – PHASE FIVE – PHASE SIX (June 2020 – January 2021)**

| Objective | TASK | Scheduled period (shaded weeks) |
|---|---|---|
| **4** | **Future works** | |
| 4.1 | Spectral ADMM implementation & Optimization | June 2020 |
| 4.2 | Literature review on fMRI data preprocessing | June 2020 |
| 4.3 | Perform real data experiments | July–November 2020 |
| 4.4 | Conclude result & preparing thesis defense | December 2020 – January 2021 |

# SUPPLEMENTARY: FORMULATION COST FUNCTION

$$(1/2) \left\| Y^{(k)} - A^{(k)} H^{(k)} \right\|_2^2$$ Least square (individual)

$$\sum_{k=1}^{K} (1/2) \left\| Y^{(k)} - A^{(k)} H^{(k)} \right\|_2^2$$ Least square (joint)

$$A^{(k)} = \left[ \hat{A}_1^{(k)} ..., \hat{A}_p^{(k)} \right]$$

$$B_{ij}^{(k)} = [(A_1^{(k)})_{ij} ... (A_p^{(k)})_{ij}]$$

$$C_{ij} = [B_{ij}^{(1)} ... B_{ij}^{(K)}]$$

Regularization

$$\sum_{k=1}^{K} \sum_{i \neq j} \left\| C_{ij} \right\|_2$$

Formulation C

$$\sum_{k=1}^{K} \sum_{i \neq j} \left\| B_{ij}^{(k)} \right\|_2$$

$$\sum_{k=1}^{K} \sum_{i \neq j} \left\| C_{ij} \right\|_2$$

Formulation D

$$\sum_{k=1}^{K} \sum_{i \neq j} \left\| B_{ij}^{(k)} \right\|_2$$

$$\sum_{k < k'} \sum_{i \neq j} \left\| B_{ij}^{k} - B_{ij}^{k'} \right\|_2$$

Formulation S

38

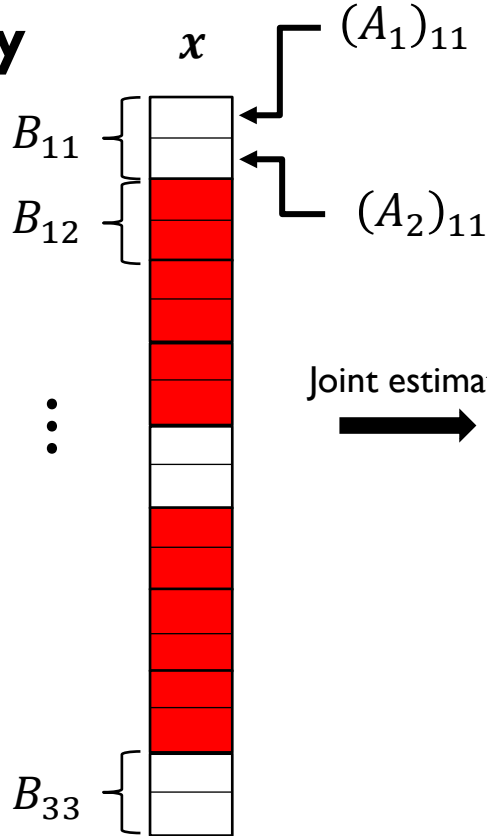**Group norm penalty**

Stacked VAR coefficient matrix

Knowing that the sparsity must be a block of size $p$

Penalize $\sum_{i \neq j} \left\| B_{ij} \right\|_2^q \longrightarrow \left\| Px \right\|_{2,q}^{(p)}$

$q = 1$, Group lasso
$q = 1/2$, **Our non-convex extension**

$F_{ij}^{(k)} = 0 \Leftrightarrow \left( A_r^{(k)} \right)_{ij} = 0; r = 1,2,...p$

Penalize $\sum_{i \neq j} \left\| C_{ij} \right\|_2^q \longrightarrow \left\| Px \right\|_{2,q}^{(pK)}$

Penalize $\sum_{m \neq l} \sum_{i \neq j} \left\| B_{ij}^{(m)} - B_{ij}^{(l)} \right\|_2^q \longrightarrow \left\| Dx \right\|_{2,q}^{(p)}$

Penalize $\sum_{k} \sum_{i \neq j} \left\| B_{ij}^{(k)} \right\|_2^q \longrightarrow \left\| Px \right\|_{2,q}^{(p)}$

39

C.  $\min_{x} \|y - Gx\|_2^2 + \lambda\|Px\|_{2,q}^{(pK)}$

D.  $\min_{x} \|y - Gx\|_2^2 + \lambda_1\|Px\|_{2,q}^{(p)} + \lambda_2\|Px\|_{2,q}^{(pK)}$

S.  $\min_{x} \|y - Gx\|_2^2 + \underbrace{\lambda_1\|Px\|_{2,q}^{(p)} + \lambda_2\|Dx\|_{2,q}^{(p)}}_{g(z)}$

$$\begin{bmatrix} T_1 \\ T_2 \end{bmatrix} x - z = 0$$

# SUPPLEMENTARY : ADMM STEP

$$x^+ = \underset{x}{\mathrm{argmin}} \|Gx - b\|_2^2 + \frac{\rho}{2} \left\| \begin{bmatrix} L_1 \\ L_2 \end{bmatrix} x - z + \frac{y}{\rho} \right\|_2^2$$

$$z^+ = \underset{z_1, z_2}{\min} \lambda_1 \|z_1\|_{2,q}^{(M_1)} + \lambda_2 \|z_2\|_{2,q}^{(M_2)} + \frac{\rho}{2} \left\| \begin{bmatrix} L_1 \\ L_2 \end{bmatrix} x - z + \frac{y}{\rho} \right\|_2^2$$

$$y^+ = y + \rho \left( \begin{bmatrix} L_1 \\ L_2 \end{bmatrix} x - z \right)$$

**Monotone accelerated proximal gradient (mAPG)**

Beck & Teboulle

Descent

$$y_k = x_k + \frac{t_{k-1} - 1}{t_k}(x_k - x_{k-1}) + \frac{t_{k-1}}{t_k}(z_k - x_k)$$

$$t_{k+1} = 0.5(1 + \sqrt{1 + 4t_k^2})$$

$$z_{k+1} = \text{prox}_{\lambda g}(y_k - \lambda\nabla f(y_k))$$

$$x_{k+1} = \text{argmin}\{F(x_k), F(z_{k+1})\}$$ **Monitoring step**

Does not generate sufficient decreasing sequence.

**Monotone accelerated proximal gradient (mAPG)**

Beck & Teboulle ➡ Li & Lin

Descent ➡ Sufficient descent

$$y_k = x_k + \frac{t_{k-1} - 1}{t_k}(x_k - x_{k-1}) + \frac{t_{k-1}}{t_k}(z_k - x_k)$$

$$t_{k+1} = 0.5(1 + \sqrt{1 + 4t_k^2})$$

Compute original proximal gradient step

$$z_{k+1} = \text{prox}_{\lambda g}(y_k - \lambda \nabla f(y_k)) \implies v_{k+1} = \text{prox}_{\lambda g}(x_k - \lambda \nabla f(x_k))$$

$$x_{k+1} = \text{argmin}\{F(\boldsymbol{v_{k+1}}), F(z_{k+1})\}$$ **Monitoring step** ➡ Sufficient descent

**Monotone APG** is proved to converge in some **non-convex** problems.

However, monitoring step is **too conservative**.

**Monotone accelerated proximal gradient (mAPG)**

Beck & Teboulle $\Rightarrow$ Li & Lin

Descent $\Rightarrow$ Sufficient descent

$$y_k = x_k + \frac{t_{k-1} - 1}{t_k}(x_k - x_{k-1}) + \frac{t_{k-1}}{t_k}(z_k - x_k)$$

$$t_{k+1} = 0.5(1 + \sqrt{1 + 4t_k^2})$$

$$z_{k+1} = \text{prox}_{\lambda g}(y_k - \lambda \nabla f(y_k))$$

$$F(z_{k+1}) \leq F(x_k) - \delta \|z_{k+1} - x_k\|_2^2?$$

YES          NO

$x_{k+1} = z_{k+1}$          $x_{k+1} = \text{prox}_{\lambda g}(x_k - \lambda \nabla f(x_k))$
No proximal step

Can **sufficient descent property** can be dropped in non-convex setting ?

There is a trick.

## Non-monotone accelerated proximal gradient (nmAPG)

In objective function

Beck & Teboulle ⟹ Li & Lin

Descent ⟹ Sufficient descent

$$y_k = x_k + \frac{t_{k-1} - 1}{t_k}(x_k - x_{k-1}) + \frac{t_{k-1}}{t_k}(z_k - x_k)$$

$$t_{k+1} = 0.5(1 + \sqrt{1 + 4t_k^2})$$

$$z_{k+1} = \text{prox}_{\lambda g}(y_k - \lambda \nabla f(y_k))$$

$$F(z_{k+1}) \leq c_k - \delta \|z_{k+1} - y_k\|_2^2 ?$$

YES     NO

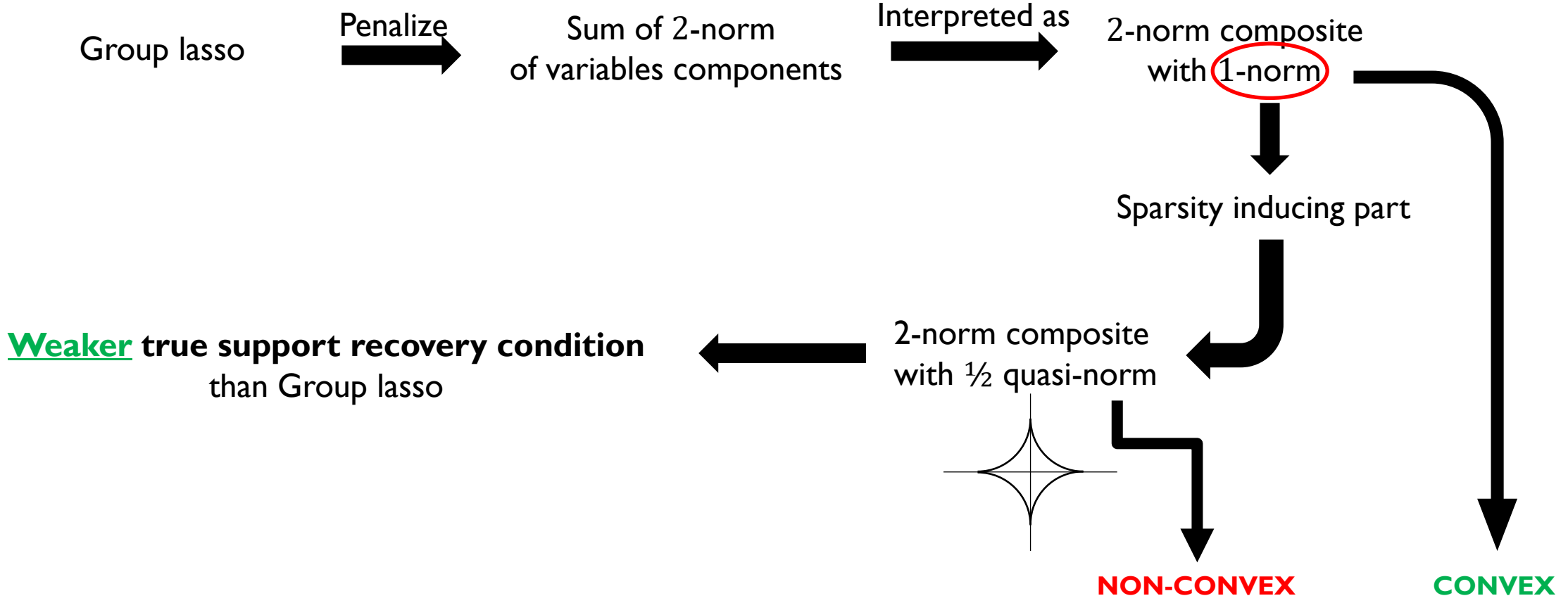$$x_{k+1} = z_{k+1}$$
No proximal step

Monitoring step in mAPG

$$x_{k+1} = \text{argmin}\{F(v_{k+1}), F(z_{k+1})\}$$

$c_k$ is weighted average of objective function in iterations $1, \dots, k$.

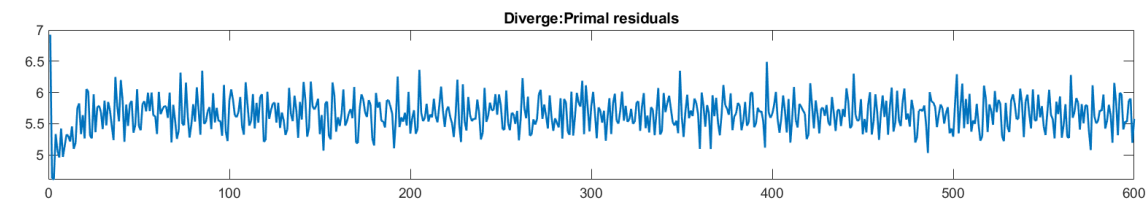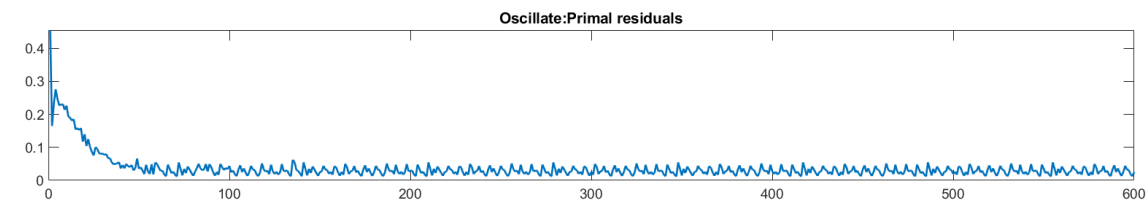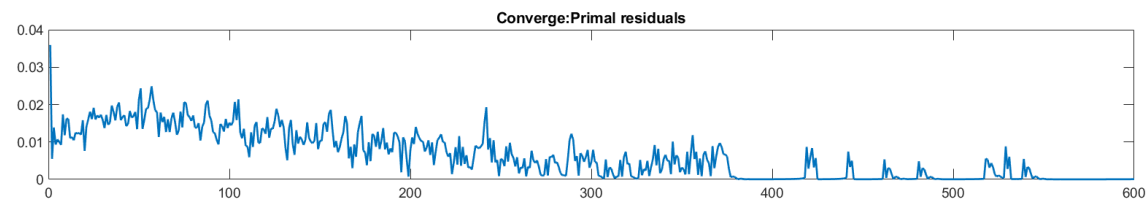Sequence $c_k$ is strictly monotone decreasing while $F(x_k)$ **may not.**

45

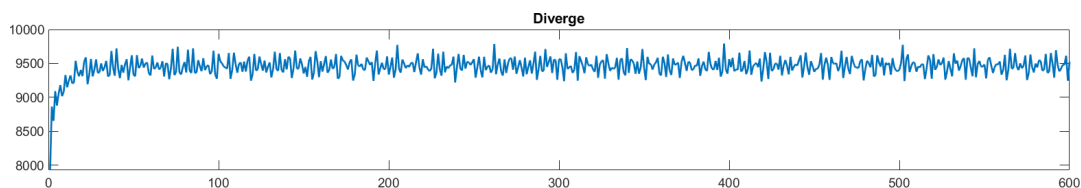# Choices of sparsity inducing-penalty

Group lasso → **Penalize** → Sum of 2-norm of variables components → **Interpreted as** → 2-norm composite with (1-norm)

↓ Sparsity inducing part

2-norm composite with ½ quasi-norm

2-norm composite with ½ quasi-norm → **Weaker** true support recovery condition than Group lasso

**NON-CONVEX**   **CONVEX**

46

ADMM convergence issues

$$\min_{Ax+Bz=c} f(x) + g(z)$$



$f(x) + g(z)$

$\|Ax + Bz - c\|_2$

$$\mathcal{L}(x, z, y, \rho) = f(x) + g(z) + y^T r + \frac{\boldsymbol{\rho}}{2}\|r\|_2^2$$

$$r = (Ax + Bz - c)$$

47