

# Joint Gaussian Graphical Model Series – VII

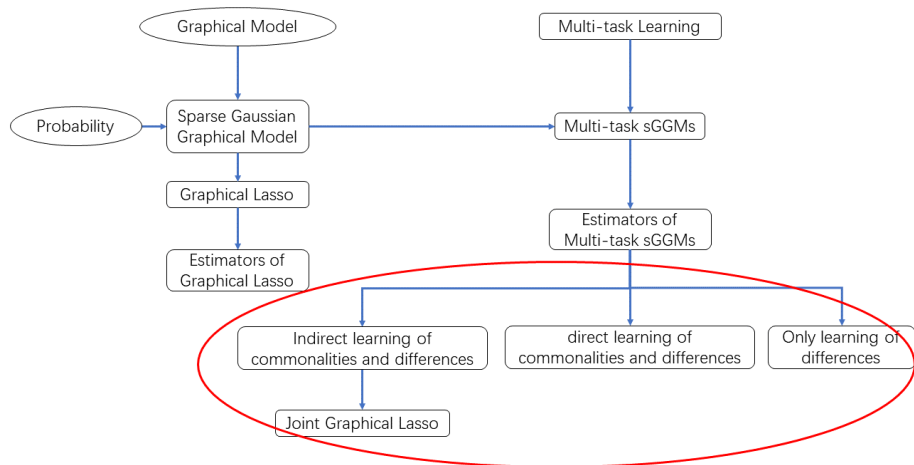
## Multi-task sGGMs estimators

Beilun Wang

<sup>1</sup>Department of Computer Science, University of Virginia  
<http://jointggm.org/>

August 18th, 2017

# Road Map



# Outline

- 1 Notation
- 2 Review
- 3 Multi-task Learning
- 4 Multi-task sGGMs
- 5 Multi-task sGGMs estimators
  - Joint Graphical Lasso
  - Directly learn the commonalities and differences among tasks
  - Directly learn the differences between case and control

# Notation

# Notation

$X^{(i)}$  The  $i$ -th data matrix

$\Sigma^{(i)}$  The  $i$ -th covariance matrix.

$\Omega^{(i)}$  The  $i$ -th precision matrix.

$p$  The number of features.

$n_i$  The number of samples in the  $i$ -th data matrix.

$K$  The number of tasks.

# Review

# Review from last talk

- We introduce multi-task learning sparse Gaussian Graphical Models (sGGMs).
- We introduce the optimization challenges in the multi-task sGGMs.
- We introduce the ADMM method and the solution of Joint Graphical Lasso.

# Review of Gaussian Graphical Model

Suppose the precision matrix  $\Omega = \Sigma^{-1}$ .

The log-likelihood of  $\Omega$  equals to  $\ln \det(\Omega) - \text{tr}(\Omega \hat{S})$



# Multi-task Learning

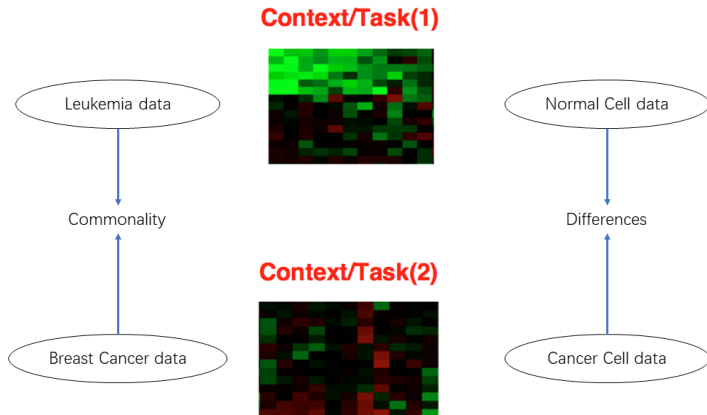
# Multi-task Learning

## Multi-task Learning

Multi-task learning (MTL) is a subfield of machine learning in which multiple learning tasks are solved at the same time, while exploiting commonalities and differences across tasks.

This can result in improved learning efficiency and prediction accuracy for the task-specific models, when compared to training the models separately.

# Multi-task Learning



# Multi-task Learning–Linear Classifier Example

## Linear Classifier

$$f(x) = \text{sgn}(w^T x + b) \quad (3.1)$$

## Multi-task Linear Classifiers

For the  $i$ -th task,

$$f_i(x) = \text{sgn}((w_S^T + w_i^T)x + b) \quad (3.2)$$

## Multi-task sGGMs

## Problem

- Input:  $\{X^{(i)}\}$
- Output:  $\{\Omega^{(i)}\}$
- Assumption I: Sparsity
- Assumption II: Commonalities and Differences

# Multi-task sGGMs

## Likelihood

$$\sum_i n_i (\ln \det(\Omega^{(i)}) - \text{tr}(\Omega^{(i)} \widehat{\mathcal{S}}^{(i)})) \quad (4.1)$$

## Likelihood with sparsity assumption

$$\operatorname{argmax}_{\Omega^{(i)}} \sum_i n_i (\ln \det(\Omega^{(i)}) - \text{tr}(\Omega^{(i)} \widehat{\mathcal{S}}^{(i)})) \quad (4.2)$$

$$\text{Subject to: } \|\Omega^{(i)}\|_1 \leq t \quad (4.3)$$

# Multi-task sGGMs

## Likelihood with multi-task setting

$$\operatorname{argmax}_{\Omega^{(i)}} \sum_i n_i (\ln \det(\Omega^{(i)}) - \operatorname{tr}(\Omega^{(i)} \widehat{\mathcal{S}}^{(i)})) \quad (4.4)$$

$$\text{Subject to: } \|\Omega^{(i)}\|_1 \leq t \quad (4.5)$$

$$P(\Omega^{(1)}, \Omega^{(2)}, \dots, \Omega^{(K)}) \leq t_2 \quad (4.6)$$

## Joint Graphical Lasso

[Danaher et al.(2013) Danaher, Wang, and Witten]

$$\operatorname{argmin}_{\Omega^{(i)}} - \sum_i n_i (\ln \det(\Omega^{(i)}) + \operatorname{tr}(\Omega^{(i)} \widehat{\mathcal{S}}^{(i)})) + \lambda_1 \|\Omega^{(i)}\|_1 + \lambda_2 P(\Omega^{(1)}, \Omega^{(2)}, \dots, \Omega^{(K)}) \quad (4.7)$$



## Multi-task sGGMs estimators

# Multi-task sGGMs estimators

- Joint Graphical Lasso type estimators
- Directly learn the commonalities and differences among tasks
- Directly learn the differences between case and control

# Joint Graphical Lasso estimators

## Different Joint Graphical Lasso

In the end, different multi-task sGGMs estimators choose different  $P(\Omega^{(1)}, \Omega^{(2)}, \dots, \Omega^{(K)})$ .

## Solutions

Most methods use ADMM as the solution of the estimators.

## JGL:Problem

- Input:  $\{X^{(i)}\}$
- Output:  $\{\Omega^{(i)}\}$
- Assumption I: Sparsity
- Assumption II: Commonalities and Differences

# Multi-task sGGMs estimators

Group Lasso [Danaher et al. (2013) Danaher, Wang, and Witten]

$$P(\Omega^{(1)}, \Omega^{(2)}, \dots, \Omega^{(K)}) = \|\Omega^{(1)}, \Omega^{(2)}, \dots, \Omega^{(K)}\|_{\mathcal{G}, 2}.$$

SIMONE [Chiquet et al. (2011) Chiquet, Grandvalet, and Ambroise]

$$P(\Omega^{(1)}, \Omega^{(2)}, \dots, \Omega^{(K)}) = \sum_{i \neq j} \left( \left( \sum_{k=1}^T (\Omega_{ij}^{(k)})_+^2 \right) \right)^{\frac{1}{2}} + \left( \left( \sum_{k=1}^K (-\Omega_{ij}^{(k)})_+^2 \right) \right)^{\frac{1}{2}}.$$

Node

JGL [Mohan et al. (2013) Mohan, London, Fazel, Lee, and Witten]

$$P(\Omega^{(1)}, \Omega^{(2)}, \dots, \Omega^{(K)}) = \sum_{ij, i > j} RCON(\Omega^{(i)} - \Omega^{(j)}).$$

**Definition 1** The row-column overlap norm (RCON) induced by a matrix norm  $\|\cdot\|$  is defined as

$$\Omega(\Theta^1, \Theta^2, \dots, \Theta^K) = \min_{V^1, V^2, \dots, V^K} \left\| \begin{bmatrix} V^1 \\ V^2 \\ \vdots \\ V^K \end{bmatrix} \right\|$$

subject to  $\Theta^k = V^k + (V^k)^T$  for  $k = 1, \dots, K$ .

## Directly learn the commonalities and differences among tasks: Problem

- Input:  $\{X^{(i)}\}$
- Output:  $\{\Omega_I^{(i)}, \Omega_S\}$
- Assumption I: Sparsity
- Assumption II: Commonalities and Differences

## Multi-task sGGMs estimators – Direct modeling

The second penalty function is still an indirect way to model the commonality and differences among tasks. Some works try to directly model this relationship.

### Mixed Neighborhood Selection

(MSN)[Monti et al.(2015)Monti, Anagnostopoulos, and Montana]

the neighborhood edges of a given node  $v$  in the  $i$ -task is modeled as  $\beta^v + \tilde{\mathbf{b}}^{(i),v}$ . Here  $\tilde{\mathbf{b}}^{(i),v} \sim N(0, \Phi^v)$ .

Consider the CLIME estimator, we can directly model the graphs as the sum of commonality and differences

### SIMULE

$$\Omega^{(i)} = \epsilon\Omega_S + \Omega_I^{(i)}.$$



## SIMULE

$$\hat{\Omega}_I^{(1)}, \hat{\Omega}_I^{(2)}, \dots, \hat{\Omega}_I^{(K)}, \hat{\Omega}_S = \operatorname{argmin}_{\Omega_I^{(i)}, \Omega_S} \sum_i \|\Omega_I^{(i)}\|_1 + \epsilon K \|\Omega_S\|_1$$

Subject to:  $\|\Sigma^{(i)}(\Omega_I^{(i)} + \Omega_S) - I\|_\infty \leq \lambda_n, i = 1, \dots, K$

## Multi-task sGGMs estimators – Direct modeling the differential networks: Problem

- Input:  $\{X^{(i)}\}$
- Output:  $\{\Delta\}$
- Assumption I: Sparse Differential networks

# Multi-task sGGMs estimators – Direct modeling the differential networks I

## Fused GLasso

By adding a regularization to enforce the sparsity of  $\Delta = \Omega_c - \Omega_d$ , we have the following formulation:

$$\underset{\Omega_c, \Omega_d \succ 0, \Delta}{\operatorname{argmin}} \mathcal{L}(\Omega_c) + \mathcal{L}(\Omega_d) \lambda_n (\|\Omega_c\|_1 + \|\Omega_d\|_1) + \lambda_2 \|\Delta\|_1 \quad (5.1)$$

The Fused Lasso assumes  $\Omega_{case}, \Omega_{control}, \Delta$ . However, many real world applications, like brain imaging data, only assume the differential network  $\Delta$  is sparse.

# Direct modeling the differential networks II: Differential CLIME

A recent study proposes the following model, which only assume the sparsity of  $\Delta$ .

## Differential CLIME

$$\begin{aligned} & \underset{\Delta}{\operatorname{argmin}} \|\Delta\|_1 \\ \text{Subject to: } & \|\widehat{\Sigma}_c \Delta \widehat{\Sigma}_d - (\widehat{\Sigma}_c - \widehat{\Sigma}_d)\|_\infty \leq \lambda_n \end{aligned} \quad (5.2)$$

However, this method is solved by a linear programming. It has  $p^2$  variables in this method. Therefore, the time complexity is at least  $O(p^8)$ . In practice, it takes more than 2 days to finish running the method when  $p = 120$ .

# Direct modeling the differential networks III: Density Ratio

The above methods all make the Gaussian assumption. This method relaxes the model to the exponential family distribution.

## Density Ratio

$$\frac{p_c(x, \theta_c)}{p_d(x, \theta_d)} \propto \exp\left(\sum_t \Delta_t f_t(x)\right) \quad (5.3)$$

Here  $\Delta_t$  encodes the difference between two Networks for factor  $f_t$ .

## Density Ratio

$$r(x; \theta) = \frac{1}{N(\theta)} \exp\left(\sum_t \Delta_t f_t(x)\right) \quad (5.4)$$

Here  $\Delta_t$  encodes the difference between two Networks for factor  $f_t$ .  $N(\theta)$  is a normalization term.

## Density Ratio for Markov Random Field

$$\begin{aligned} \hat{p}(x) &= p_d(x)r(x; \theta) \\ \text{KL}[p_c || \hat{p}] &= \text{Const.} - \int p_c(x) \log r(x; \theta) dx. \end{aligned} \tag{5.5}$$

# Summary

- We introduce the multi-task sGGMs estimators.
- We introduce the multi-task sGGMs estimators, which directly model the commonalities and differences.
- We introduce the multi-task sGGMs estimators, which directly model the differences.

# References I

 J. Chiquet, Y. Grandvalet, and C. Ambroise.

Inferring multiple graphical structures.

*Statistics and Computing*, 21(4):537–553, 2011.

 P. Danaher, P. Wang, and D. M. Witten.

The joint graphical lasso for inverse covariance estimation across multiple classes.

*Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2013.

 K. Mohan, P. London, M. Fazel, S.-I. Lee, and D. Witten.

Node-based learning of multiple gaussian graphical models.

*arXiv preprint arXiv:1303.5145*, 2013.



## References II



R. P. Monti, C. Anagnostopoulos, and G. Montana.

Learning population and subject-specific brain connectivity networks via mixed neighborhood selection.

*arXiv preprint arXiv:1512.01947*, 2015.