# Fast and Scalable Learning of Sparse Changes in High-Dimensional Gaussian Graphical Model Structure

Beilun Wang[1]    Arshdeep Sekhon[1]    Yanjun Qi[1]

[1]University of Virginia
http://jointggm.org/

# Outline

# Motivation: Structure Difference Learning from two Datasets

- Two Datasets $\mathbf{X}_c$, $\mathbf{X}_d$ $\rightarrow$ Differential Network $\Delta$.
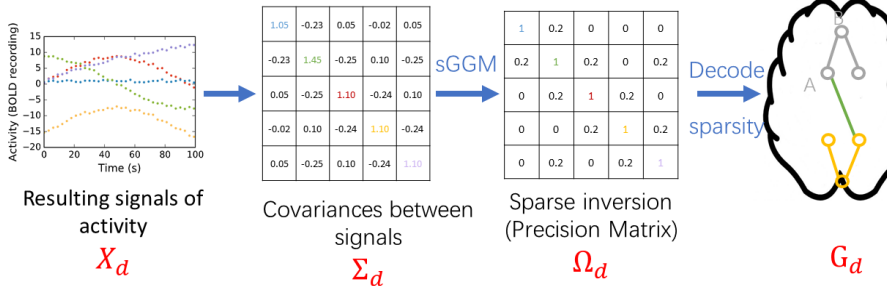  - Case vs. Control;
  - Autism vs. Normal;



Data

Differential Network

# Motivation: Estimating Graph from Dataset via sparse Gaussian Graphical Model:

- A pipeline to infer Graph from one homogeneous dataset $X_d$.



Resulting signals of activity
$X_d$

Covariances between signals
$\Sigma_d$

Sparse inversion (Precision Matrix)
$\Omega_d$

$G_d$

# Motivation: Estimating Graph from Dataset via sparse Gaussian Graphical Model:

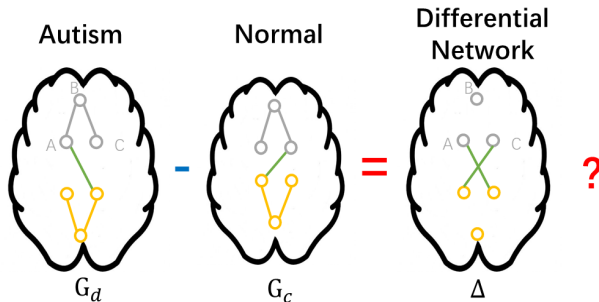- A pipeline to infer Graph from one homogeneous dataset $X_d$.



| Resulting signals of activity | Covariances between signals | Sparse inversion (Precision Matrix) | |
|---|---|---|---|
| $X_d$ | $\Sigma_d$ | $\Omega_d$ | $G_d$ |

- $X_c \rightarrow G_c$ is the same.
- We are more interested in the structure changes between two different but related datatsets.
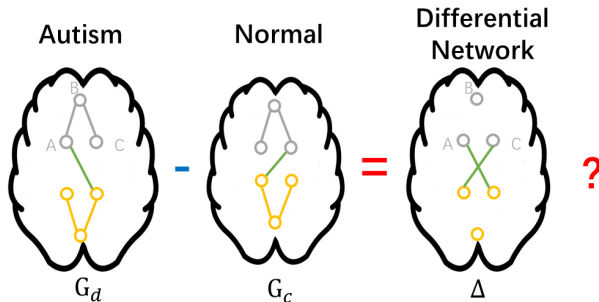
- Sparsity Assumption:
  - If estimating two graphs separately, we need to enforce sparsity assumption on both graphs
  - However, in some real-world applications, $G_c$, $G_d$ are not sparse.



Autism     Normal     Differential Network

$G_d$     $G_c$     $\Delta$

# Motivation: Estimating the Difference by separately Learning Two Graphs from two datasets has Limitations
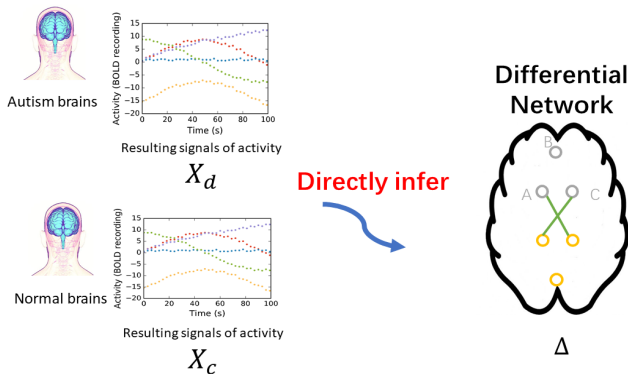
- Sparsity Assumption:
  - If estimating two graphs separately, we need to enforce sparsity assumption on both graphs
  - However, in some real-world applications, $G_c$, $G_d$ are not sparse.



Autism · Normal = Differential Network ?

$G_d$     $G_c$     $\Delta$

- Difficulty in the computation: Current methods can not scale-up. In applications like neuroscience, the number of regions (nodes) $p$ for connectivity analysis in the human brain ranges from 160 to 800,000.

- Our focus: How to directly estimate / learn Differential Network (Δ) from Two datasets ($\mathbf{X}_c$, $\mathbf{X}_d$) about the same set of features in a large scale.



Autism brains

Resulting signals of activity

$X_d$

Normal brains

Resulting signals of activity

$X_c$

Directly infer

**Differential Network**

Δ

$X_c, X_d$  Data matrix.

$\Sigma_{c,d}$  Covariance matrix.

$\Omega_c, \Omega_d$  Inverse of covariance matrix (precision matrix).

$p$  The total number of feature variables.

$n_c, n_d$  The number of samples.

$\Delta$  The Differential Network.
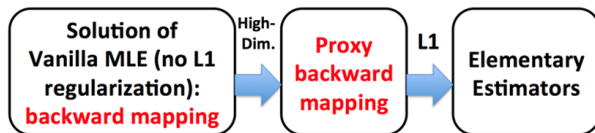
# Outline

# Background: Elementary Estimator for Exponential Family

- The canonical parameter $\theta$ of an exponential family distribution can be learned by the following equation.

## Elementary Estimator

$$\operatorname*{argmin}_{\theta} ||\theta||_1$$

$$\text{Subject to: } ||\theta - \mathcal{B}^*(\widehat{\phi})||_\infty \leq \lambda_n$$
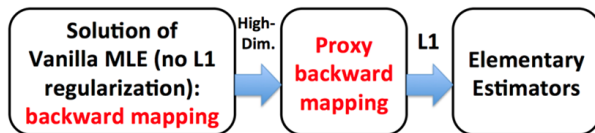
(1.1)

# Background: Elementary Estimator for Exponential Family

- The canonical parameter $\theta$ of an exponential family distribution can be learned by the following equation.

### Elementary Estimator

$$\underset{\theta}{\operatorname{argmin}} ||\theta||_1$$

$$\text{Subject to: } ||\theta - \mathcal{B}^*(\widehat{\phi})||_\infty \le \lambda_n$$

(1.1)

- For high-dimensional case, Vanilla MLE solutions are mostly not available. Therefore, we choose Proxy backward mapping.

# Background: Elementary Estimator for sGGM

- $\theta$ is the canonical parameter of the exponential distribution.
- $\mathcal{B}^*(\widehat{\phi})$ is the backward mapping of $\theta$. Normally, it is the solution of Vanilla MLE.
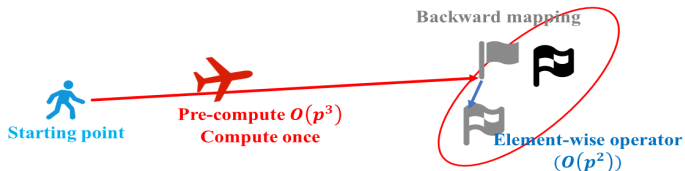- For example, for sGGM:

| EE | $\theta$ | $\mathcal{B}^*$ | $\widehat{\phi}$ |
|---------|-----------|----------------------------------|------------------|
| EE-sGGM | $\Omega$ | $[T_v(\widehat{\Sigma})]^{-1}$ | $\widehat{\Sigma}$ |

## Elementary Estimator for sGGM

$$\underset{\Omega}{\operatorname{argmin}} ||\Omega||_1$$

$$\text{Subject to: } ||\Omega - [T_v(\widehat{\Sigma})]^{-1}||_\infty \leq \lambda_n \tag{1.2}$$

- Elementary Estimator:



Backward mapping

Starting point

Pre-compute $O(p^3)$
Compute once

Element-wise operator
($O(p^2)$)

- Closed-form solution (non-iterative algorithm)

- Closed-form solution (non-iterative algorithm)

- Fast computation

- Closed-form solution (non-iterative algorithm)

- Fast computation

- Easy to prove the theoretical error bound

- Learn an exponential family distribution $\Longleftrightarrow$ To learn the canonical parameter $\theta$
  - e.g., $\Omega$ is the canonical parameter of the sparse Gaussian Graphical Model

- Learn an exponential family distribution $\iff$ To learn the canonical parameter $\theta$
  - e.g., $\Omega$ is the canonical parameter of the sparse Gaussian Graphical Model

- The density ratio of two Gaussian distributions is naturally an exponential family distribution.
  - We prove $\Delta$ is the canonical parameter of the density ratio distribution.
  - Therefore we can apply Elementary Estimator to estimate $\Delta$.

# Outline

# Goals

Our model aims to have the following properties:

- closed-form solution.

Our model aims to have the following properties:

- closed-form solution.
- Fast and scalable algorithm.

Our model aims to have the following properties:

- closed-form solution.
- Fast and scalable algorithm.
- It provides a strong theoretical guarantee.

# Proposed Method: estimating DIFFerential networks via an Elementary Estimator (DIFFEE) when high-dimensional

We model the differential network $\Delta$ as:

$$\Delta = \Omega_d - \Omega_c \tag{2.1}$$

We apply the elementary estimator to the differential network.

### DIFFEE

$$\underset{\Delta}{\operatorname{argmin}} ||\Delta||_1$$

Subject to: $||\Delta - \mathcal{B}^*(\widehat{\Sigma}_d, \widehat{\Sigma}_c)||_\infty \leq \lambda_n$ $\tag{2.2}$

| EE | $\theta$ | $\mathcal{B}^*$ | $\widehat{\phi}$ |
|---|---|---|---|
| EE-sGGM | $\Omega$ | $[T_v(\widehat{\Sigma})]^{-1}$ | $\widehat{\Sigma}$ |
| DIFFEE | $\Delta$ | $[T_v(\widehat{\Sigma}_d)]^{-1} - [T_v(\widehat{\Sigma}_c)]^{-1}$ | $\widehat{\Sigma}_d, \widehat{\Sigma}_c$ |

- We choose $[T_v(\widehat{\Sigma}_d)]^{-1} - [T_v(\widehat{\Sigma}_c)]^{-1}$ as the proxy backward mapping for $\Delta$.

| EE | $\theta$ | $\mathcal{B}^*$ | $\widehat{\phi}$ |
|---|---|---|---|
| EE-sGGM | $\Omega$ | $[T_v(\widehat{\Sigma})]^{-1}$ | $\widehat{\Sigma}$ |
| DIFFEE | $\Delta$ | $[T_v(\widehat{\Sigma}_d)]^{-1} - [T_v(\widehat{\Sigma}_c)]^{-1}$ | $\widehat{\Sigma}_d, \widehat{\Sigma}_c$ |

- We choose $[T_v(\widehat{\Sigma}_d)]^{-1} - [T_v(\widehat{\Sigma}_c)]^{-1}$ as the proxy backward mapping for $\Delta$.
- It is theoretical guaranteed (will talk later).

| EE | $\theta$ | $\mathcal{B}^*$ | $\widehat{\phi}$ |
|---|---|---|---|
| EE-sGGM | $\Omega$ | $[T_v(\widehat{\Sigma})]^{-1}$ | $\widehat{\Sigma}$ |
| DIFFEE | $\Delta$ | $[T_v(\widehat{\Sigma}_d)]^{-1} - [T_v(\widehat{\Sigma}_c)]^{-1}$ | $\widehat{\Sigma}_d, \widehat{\Sigma}_c$ |

- We choose $[T_v(\widehat{\Sigma}_d)]^{-1} - [T_v(\widehat{\Sigma}_c)]^{-1}$ as the proxy backward mapping for $\Delta$.
- It is theoretical guaranteed (will talk later).

- Closed-form solution:

$$\widehat{\Delta} = S_{\lambda_n}([T_v(\widehat{\Sigma}_d)]^{-1} - [T_v(\widehat{\Sigma}_c)]^{-1}) \qquad (2.3)$$

Here $[S_\lambda(A)]_{ij} = \text{sign}(A_{ij}) \max(|A_{ij}| - \lambda, 0)$.

- It has closed-form solution.

- It has closed-form solution.
- It is faster than the previous studies:

| DIFFEE | FusedGLasso | Density Ratio | Diff-CLIME |
|--------|-------------|---------------|------------|
| $O(p^3)$ | $O(T*p^3)$ | $O((n_c + p^2)^3)$ | $O(p^8)$ |

# Why DIFFEE is better

- It has closed-form solution.
- It is faster than the previous studies:

| DIFFEE | FusedGLasso | Density Ratio | Diff-CLIME |
|--------|-------------|---------------|------------|
| $O(p^3)$ | $O(T * p^3)$ | $O((n_c + p^2)^3)$ | $O(p^8)$ |

- $O(p^2)$ to tune different $\lambda_n$

# Why DIFFEE is better

- It has closed-form solution.
- It is faster than the previous studies:

| DIFFEE | FusedGLasso | Density Ratio | Diff-CLIME |
|--------|-------------|---------------|------------|
| $O(p^3)$ | $O(T * p^3)$ | $O((n_c + p^2)^3)$ | $O(p^8)$ |

- $O(p^2)$ to tune different $\lambda_n$
- Theoretical guaranteed

# Previous Methods: FusedGLasso for estimating differential network

- Traditionally, we estimate differential network from penalized likelihood formulation.

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

$$\operatorname*{argmin}_{\Omega_c,\Omega_d \succ 0, \Delta} n_c(-\log\det(\Omega_c) + <\Omega_c, \widehat{\Sigma}_c>)$$
$$+ n_d(-\log\det(\Omega_d) + <\Omega_d, \widehat{\Sigma}_d>)$$
$$+ \lambda_2(||\Omega_c||_1 + ||\Omega_d||_1) + \lambda_n||\Delta||_1 \tag{2.4}$$

# Previous Methods: FusedGLasso for estimating differential network

- Traditionally, we estimate differential network from penalized likelihood formulation.
- FusedGLasso adds a second penalty function fused norm $||\Delta||_1$ into the penalized likelihood formulation.

---

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

$$\operatorname*{argmin}_{\Omega_c, \Omega_d \succ 0, \Delta} n_c(-\log\det(\Omega_c) + <\Omega_c, \widehat{\Sigma}_c>)$$
$$+ n_d(-\log\det(\Omega_d) + <\Omega_d, \widehat{\Sigma}_d>)$$
$$+ \lambda_2(||\Omega_c||_1 + ||\Omega_d||_1) + \lambda_n||\Delta||_1$$

(2.4)

# Previous Methods: FusedGLasso for estimating differential network

- Traditionally, we estimate differential network from penalized likelihood formulation.
- FusedGLasso adds a second penalty function fused norm $||\Delta||_1$ into the penalized likelihood formulation.
- $||\Delta||_1$ enforces a sparse difference structure between two graphs.

---

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

$$\operatorname*{argmin}_{\Omega_c, \Omega_d \succ 0, \Delta} n_c(-\log \det(\Omega_c) + <\Omega_c, \widehat{\Sigma}_c>)$$
$$+ n_d(-\log \det(\Omega_d) + <\Omega_d, \widehat{\Sigma}_d>)$$
$$+ \lambda_2(||\Omega_c||_1 + ||\Omega_d||_1) + \lambda_n||\Delta||_1 \tag{2.4}$$

# Previous Methods: Diff-CLIME

- Another study to learn the $\Delta$ is through a constrained optimization formulation.

## Diff-CLIME [Zhao et al.(2014)Zhao, Cai, and Li]

$$\underset{\Delta}{\mathrm{argmin}} \, ||\Delta||_1$$

Subject to: $||\widehat{\Sigma}_c \Delta \widehat{\Sigma}_d - (\widehat{\Sigma}_c - \widehat{\Sigma}_d)||_\infty \leq \lambda_n$

(2.5)

# Previous Methods: Diff-CLIME

- Another study to learn the $\Delta$ is through a constrained optimization formulation.
- It reduces the estimation to solve multiple linear programming problems.

---

**Diff-CLIME [Zhao et al.(2014)Zhao, Cai, and Li]**

$$\underset{\Delta}{\operatorname{argmin}} ||\Delta||_1$$

$$\text{Subject to: } ||\widehat{\Sigma}_c \Delta \widehat{\Sigma}_d - (\widehat{\Sigma}_c - \widehat{\Sigma}_d)||_\infty \leq \lambda_n \tag{2.5}$$

# Previous Methods: Density Ratio

- Directly model the sparse differential network with density ratio function $r(x; \Delta)$

**Density Ratio [Liu et al.(2013)Liu, Yamada, Collier, and Sugiyama]**

$$\operatorname*{argmax}_{\Delta} \mathcal{L}_{\mathsf{KLIEP}}(\Delta) - \lambda_n \parallel \Delta \parallel_1 - \lambda_2 \parallel \Delta \parallel_2 \qquad (2.6)$$

# Previous Methods: Density Ratio

- Directly model the sparse differential network with density ratio function $r(x; \Delta)$
- Minimizes the KL divergence between $p_d(x)$ and $\widehat{p}_d(x) = r(x; \Delta)p_c(x)$.

### Density Ratio [Liu et al.(2013)Liu, Yamada, Collier, and Sugiyama]

$$\operatorname*{argmax}_{\Delta} \mathcal{L}_{\mathsf{KLIEP}}(\Delta) - \lambda_n \parallel \Delta \parallel_1 - \lambda_2 \parallel \Delta \parallel_2 \qquad (2.6)$$

# Previous Studies: Drawbacks

- The time comparison table:

| FusedGLasso | Density Ratio | Diff-CLIME |
|:---:|:---:|:---:|
| $O(T * p^3)$ | $O((n_c + p^2)^3)$ | $O(p^8)$ |

- The time comparison table:

| FusedGLasso | Density Ratio | Diff-CLIME |
|---|---|---|
| $O(T * p^3)$ | $O((n_c + p^2)^3)$ | $O(p^8)$ |

- **Drawbacks:**
  - **I:** all of them are slow when $p$ is large.

  - **II:** Need terative algorithm solution.
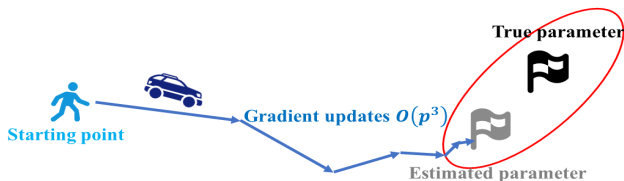
  - **III:** No theoretical analysis in the previous FusedGLasso studies.

- Previous studies:



True parameter

**Starting point**

**Gradient updates $O(p^3)$**

Estimated parameter

- Previous studies:



- DIFFEE:

# Outline

- error bound: $||\Delta^* - \widehat{\Delta}||$
- DIFFEE achieves similar error bound as the previous studies.

| DIFFEE | FusedGLasso | Density Ratio | Diff-CLIME |
|---|---|---|---|
| $\frac{\log p}{\min(n_c, n_d)}$ | $N/A$ | $\frac{\log p}{\min(n_c, n_d)}$ | $\frac{\log p}{\min(n_c, n_d)}$ |

- Sharp convergence rate as the state-of-art

$$||\widehat{\Delta} - \Delta^*||_\infty \leq \frac{16\kappa_1 a}{\kappa_2}\sqrt{\frac{\log p}{\min(n_c, n_d)}}$$

$$||\widehat{\Delta} - \Delta^*||_F \leq \frac{32\kappa_1 a}{\kappa_2}\sqrt{\frac{k\log p}{\min(n_c, n_d)}} \qquad (3.1)$$

$$||\widehat{\Delta} - \Delta^*||_1 \leq \frac{64\kappa_1 a}{\kappa_2}k\sqrt{\frac{\log p}{\min(n_c, n_d)}}$$

# Outline

- Comparison with the best baseline – FusedGLasso:



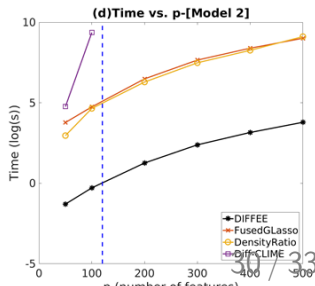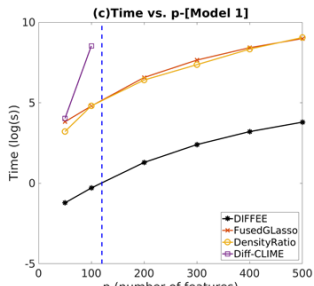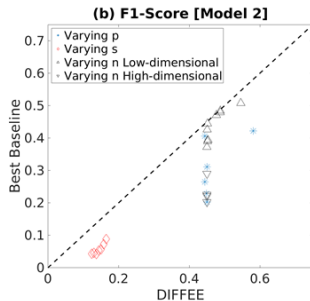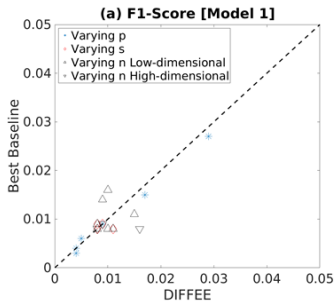(a) Precision@k [Model 1]    (b) Time (log(s)) [Model 1]

# Results on Synthetic Datasets: Accuracy and Computational time When Varying $p$

- Compare the estimation accuracy and scalabilty of different methods



(a) Precision@k vs p [Model 1]
(b) Time vs p [Model 1]
(c) F1-score vs p [Model 2]
(d) Time vs p [Model 2]

# Results on fMRI Datasets: the Classification Accuracy

- (1) ABIDE dataset
- (2) Train the differential network and use it as the parameter of a LDA classifier

| Method | DIFFEE | FusedGLasso | Diff-CLIME |
|---|---|---|---|
| Accuracy (%) | **57.58%** | 56.90% | 53.79% |

# R Package is Available !!!

- The project website: `http://jointggm.org/`

- R package "diffee":
  - `install.packages("diffee")`
  - `demo(diffeeDemo)` !
  - `https://cran.r-project.org/web/packages/diffee/index.html`

## References

📄 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.

📄 S. Liu, M. Yamada, N. Collier, and M. Sugiyama.
Change-point detection in time-series data by relative density-ratio estimation.
*Neural Networks*, 43:72–83, 2013.

📄 S. D. Zhao, T. T. Cai, and H. Li.
Direct estimation of differential networks.
*Biometrika*, 101(2):253–268, 2014.