

GLAD: Learning Sparse Graph Recovery

Le Song

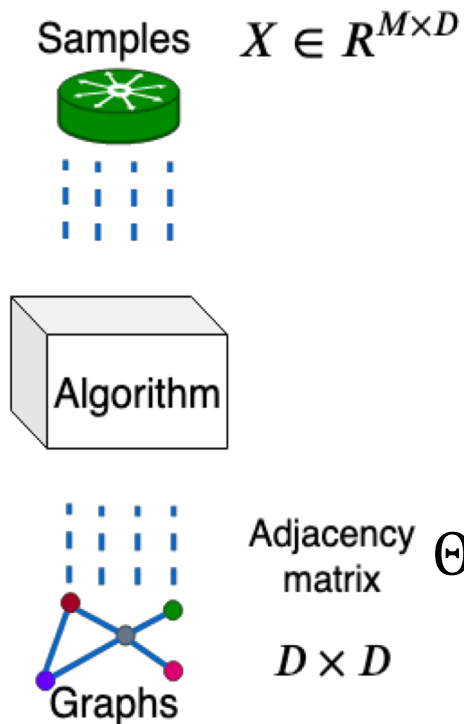
Georgia Tech

Joint work with Harsh Shrivastava, Xinshi Chen, Binghong Chen, Guanghui Lan

Objective

Recovering sparse conditional independence graph G from data

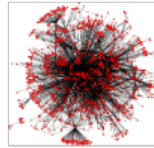
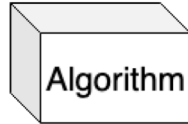
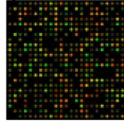
$$\Theta_{ij} = 0 \Leftrightarrow X_i \perp X_j \mid \text{other variables}$$



Applications

Biology

Gene Expression
data - Microarray
experiments



Gene regulatory
network

Finance

Time-series
features



Relationship
between assets

Convex Formulation

- Given M samples from a distribution: $X \in \mathbb{R}^{M \times D}$
- Estimate the matrix ' Θ ' corresponding to the sparse graph

Objective function: L1 regularized log-determinant estimation

$$\hat{\Theta} = \arg \min_{\Theta \in \mathcal{S}_{++}^d} -\log(\det \Theta) + \text{tr}(\hat{\Sigma}\Theta) + \rho \|\Theta\|_{1,\text{off}}$$

Covariance matrix

$$\hat{\Sigma} = \frac{X^T X}{M}$$

Regularization
Parameter

Existing Optimization Algorithms

G-ISTA

Proximal
gradient
method

$$\Theta_{k+1} \leftarrow \eta_{\xi_k \rho} (\Theta_k - \xi_k (\hat{\Sigma} - \Theta_k^{-1})), \quad \text{where } [\eta_\rho(X)]_{ij} := \text{sign}(X_{ij})(|X_{ij}| - \rho)_+$$

Existing Optimization Algorithms

G-ISTA

Proximal
gradient
method

Glasso

Block
coordinate
descent
method

Updates each column (and the corresponding row) of the precision matrix iteratively by solving a sequence of lasso problems

Existing Optimization Algorithms

G-ISTA

Proximal
gradient
method

Glasso

Block
coordinate
descent
method

ADMM

Alternating
direction
method of
multipliers

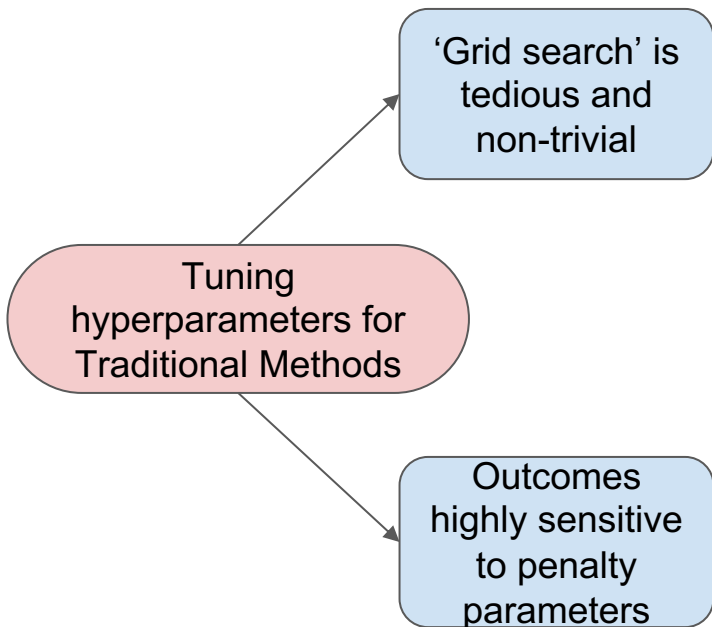
$$-\log(\det \Theta) + \text{tr}(\widehat{\Sigma}\Theta) + \rho \|Z\|_1 + \langle \lambda, \Theta - Z \rangle + \frac{1}{2}\beta \|Z - \Theta\|_F^2.$$

Taking $U := \lambda/\beta$ as the scaled dual variable, the update rules for the ADMM algorithm are

$$\Theta_{k+1} \leftarrow (-Y + \sqrt{Y^\top Y + (4/\beta)I})/2, \text{ where } Y = \widehat{\Sigma}/\beta - Z_k + U_k$$

$$Z_{k+1} \leftarrow \eta_{\rho/\beta}(\Theta_{k+1} + U_k), \quad U_{k+1} \leftarrow U_k + \Theta_{k+1} - Z_{k+1}$$

Hard to Tune Hyperparameters

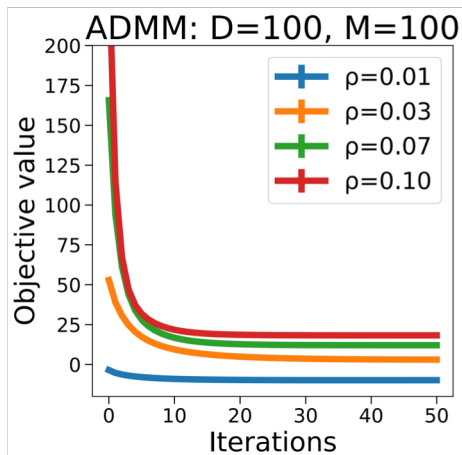


$\rho \backslash \beta$	5	1	0.5	0.1	0.01
0.01	-2.51	-2.25	-2.06	-2.06	-2.69
0.03	-5.59	-9.05	9.48	-9.61	-9.41
0.07	-9.53	-7.58	-7.42	-7.38	-7.46
0.1	-9.38	-6.51	-6.43	-6.41	-6.50
0.2	-6.76	-4.68	-4.55	-4.47	-4.80

Errors of different parameter combinations

$$-\log(\det \Theta) + \text{tr}(\hat{\Sigma}\Theta) + \rho \|Z\|_1 + \langle \lambda, \Theta - Z \rangle + \frac{1}{2}\beta \|Z - \Theta\|_F^2.$$

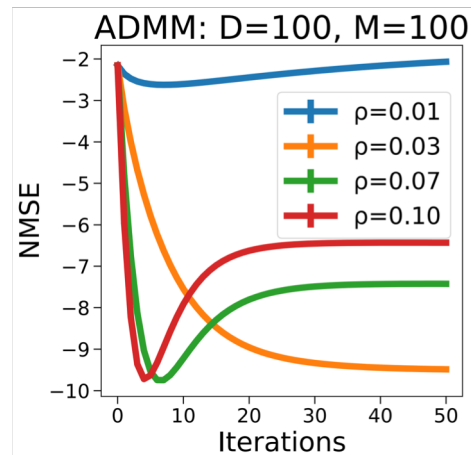
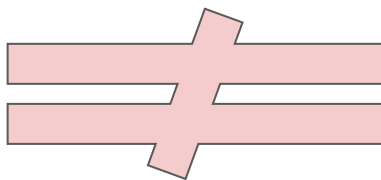
Mismatch in Objectives



Log-determinant estimator

$$-\log(\det \Theta) + \text{tr}(\hat{\Sigma}\Theta) + \rho \|\Theta\|_{1,\text{off}}$$

mismatch!



Recovery Objective (NMSE)

$$\|\hat{\Theta} - \Theta^*\|_F^2 / \|\Theta^*\|_F^2$$

Limitations of Existing Optimization Algorithms

Limitations of the **convex** formulation

Consistency of estimator

$$\hat{\Theta}$$

- Based on 'carefully chosen conditions' like
1. Lower bound on sample size
 2. Sparsity of Θ
 3. Degree of graph
 4. Magnitude of covariance entries

Specific regularization parameter

1. Highly sensitive parameter
2. Depends on tail behavior of maximum deviation $\max_{i,j} |\hat{\Sigma}_{ij} - \Sigma_{ij}^*|$

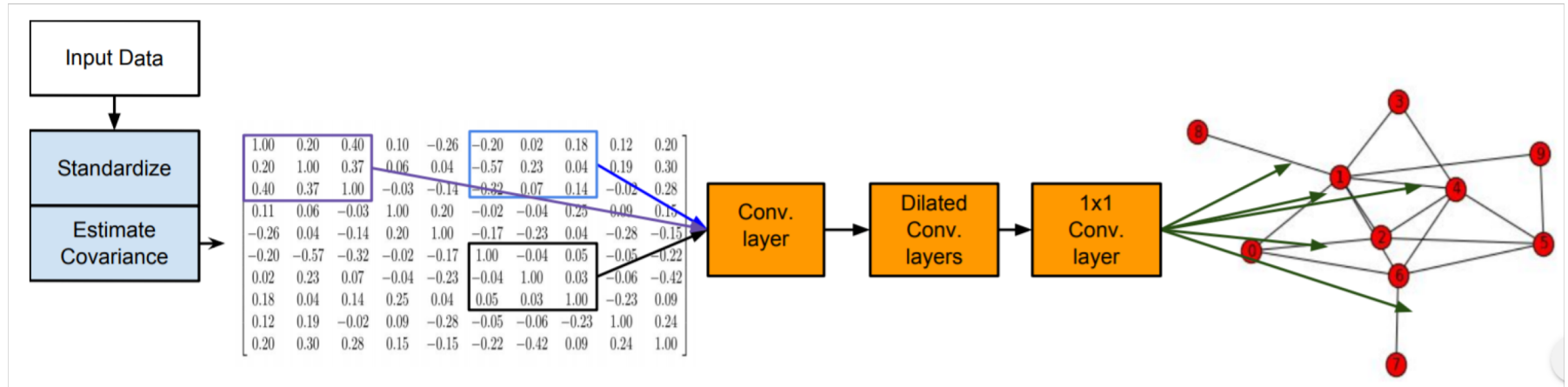
Room for Improvement!

Big Picture Question

- Given a collection of ground truth precision matrix Θ^* , and the corresponding empirical covariance $\hat{\Sigma}$
- Learn an algorithm f which directly produces an estimate of the precision matrix Θ ?

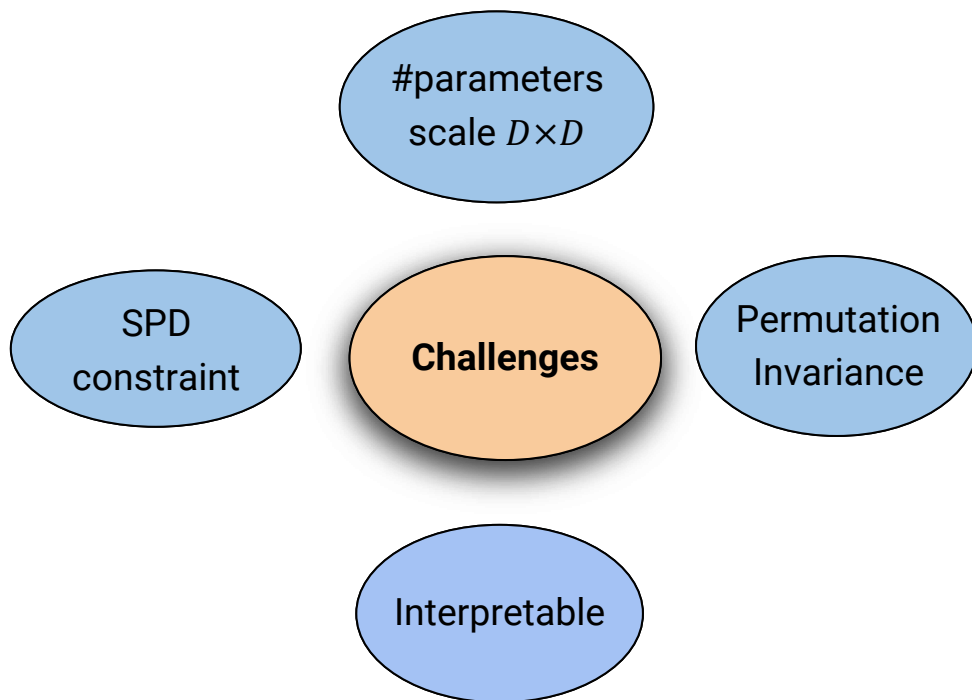
$$\min_f \frac{1}{|\mathcal{D}|} \sum_{(\hat{\Sigma}_i, \Theta_i^*) \in \mathcal{D}} \|\Theta_i - \Theta_i^*\|_F^2, \quad s.t. \Theta_i = f(\hat{\Sigma}_i)$$

Deep Learning Model Example

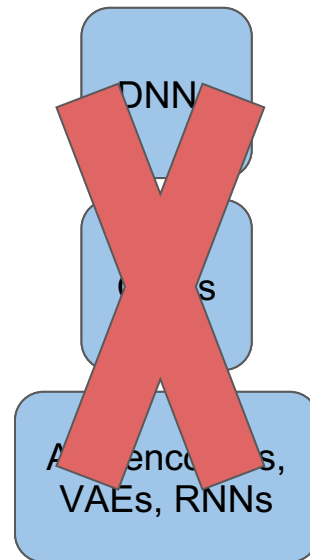


DeepGraph (DG) architecture. The input is first standardized and then the sample covariance matrix is estimated. A neural network consisting of multiple dilated convolutions (Yu & Koltun, 2015) and a final 1×1 convolution layer is used to predict edges corresponding to non-zero entries in the precision matrix.

Challenges in Designing Learning Models



Traditional Approaches



GLAD: DL model based on Unrolled Algorithm

Alternating Minimization (AM) algorithm: Objective function

$$\hat{\Theta}_\lambda, \hat{Z}_\lambda := \arg \min_{\Theta, Z \in \mathcal{S}_{++}^d} -\log(\det \Theta) + \text{tr}(\hat{\Sigma}\Theta) + \rho \|Z\|_1 + \frac{1}{2}\lambda \|Z - \Theta\|_F^2$$

AM: Update Equations (Nice closed form updates!)

$$\Theta_{k+1}^{\text{AM}} \leftarrow \frac{1}{2} \left(-Y + \sqrt{Y^\top Y + \frac{4}{\lambda} I} \right), \text{ where } Y = \frac{1}{\lambda} \hat{\Sigma} - Z_k^{\text{AM}}$$

$$Z_{k+1}^{\text{AM}} \leftarrow \eta_{\rho/\lambda}(\Theta_{k+1}^{\text{AM}}), \quad \text{where } \eta_{\rho/\lambda}(\theta) := \text{sign}(\theta) \max(|\theta| - \rho/\lambda, 0)$$

- Unroll to fixed #iterations 'K'.
- Treat it as a deep model

GLAD: Training

Loss function: Frobenius norm with discounted cumulative reward

$$\min_f \text{loss}_f := \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \gamma^{K-k} \left\| \Theta_k^{(i)} - \Theta^{*(i)} \right\|_F^2$$

Optimizer for training: 'Adam'.

Learning rate chosen between [0.01, 0.1] in conjunction with Multi-step LR scheduler.

Gradient Computation through matrix square root in the GLADcell:

For any SPD matrix X : $X = X^{1/2} X^{1/2}$

Solve Sylvester's equation for $d(X^{1/2})$:

$$dX = d(X^{1/2})X^{1/2} + X^{1/2}d(X^{1/2})$$

Use Neural Networks for (ρ, λ)

$$\lambda \leftarrow \Lambda_{nn}(\|Z - \Theta\|_F^2, \lambda)$$

$$\rho_{ij} = \rho_{nn}(\Theta_{ij}, \hat{\Sigma}_{ij}, Z_{ij})$$

of layers
= 2
Hidden unit
size = 3

**Minimalist
designing of Neural
Networks**

of layers
= 4
Hidden unit
size = 3

Non-Linearity:
Hidden layers = 'tanh'
Final layer = 'sigmoid'

GLAD

Algorithm 1: GLAD

Function GLADcell($\hat{\Sigma}, \Theta, Z, \lambda$):

$\lambda \leftarrow \Lambda_{nn}(\|Z - \Theta\|_F^2, \lambda)$

$Y \leftarrow \lambda^{-1} \hat{\Sigma} - Z$

$\Theta \leftarrow \frac{1}{2}(-Y + \sqrt{Y^T Y + \frac{4}{\lambda} I})$

For all i, j **do**

$\rho_{ij} = \rho_{nn}(\Theta_{ij}, \hat{\Sigma}_{ij}, Z_{ij})$

$Z_{ij} \leftarrow \eta_{\rho_{ij}}(\Theta_{ij})$

return Θ, Z, λ

Function GLAD($\hat{\Sigma}$):

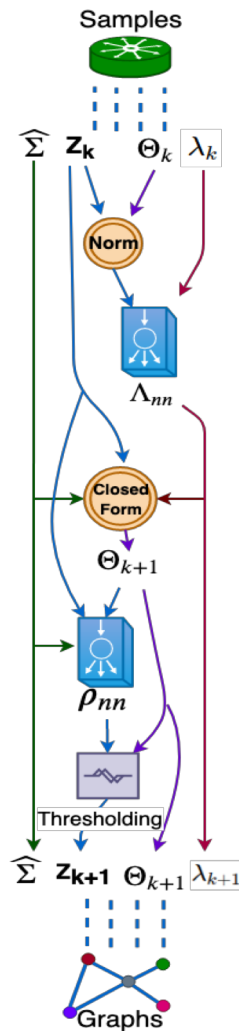
$\Theta_0 \leftarrow (\hat{\Sigma} + tI)^{-1}, \lambda_0 \leftarrow 1$

For $k = 0$ **to** $K - 1$ **do**

$\Theta_{k+1}, Z_{k+1}, \lambda_{k+1}$

\leftarrow GLADcell($\hat{\Sigma}, \Theta_k, Z_k, \lambda_k$)

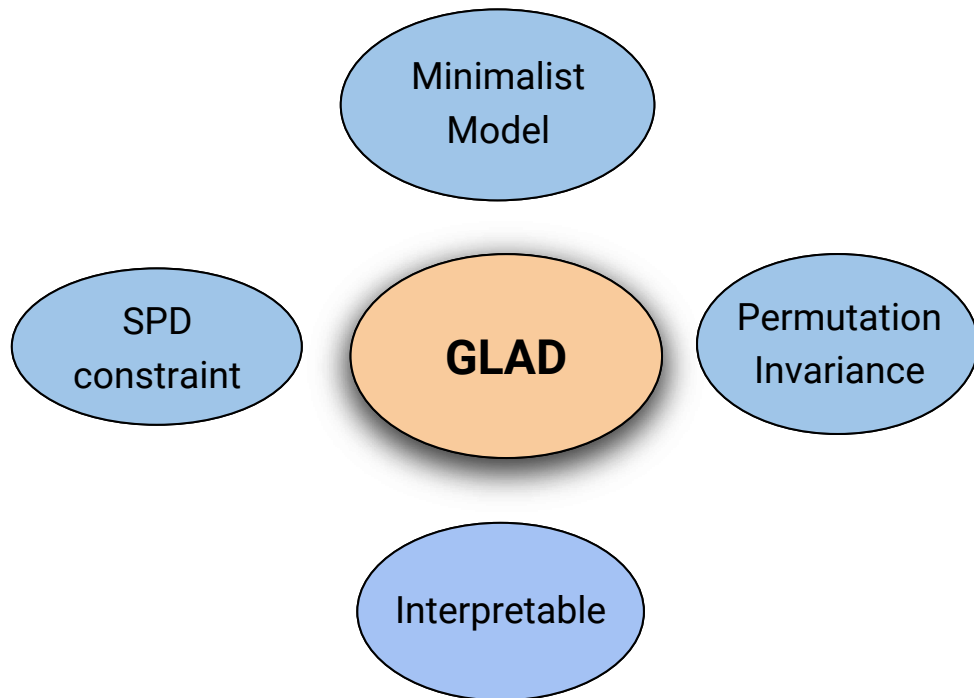
return Θ_K, Z_K



GLAD

Using algorithm structure as inductive bias for designing unrolled DL architectures

Desiderata for GLAD



Algorithm 1: GLAD

Function GLADcell($\hat{\Sigma}, \Theta, Z, \lambda$):

$\lambda \leftarrow \Lambda_{nn}(\|Z - \Theta\|_F^2, \lambda)$

$Y \leftarrow \lambda^{-1} \hat{\Sigma} - Z$

$\Theta \leftarrow \frac{1}{2}(-Y + \sqrt{Y^\top Y + \frac{4}{\lambda} I})$

For all i, j **do**

$\rho_{ij} = \rho_{nn}(\Theta_{ij}, \hat{\Sigma}_{ij}, Z_{ij})$
 $Z_{ij} \leftarrow \eta_{\rho_{ij}}(\Theta_{ij})$

return Θ, Z, λ

Function GLAD($\hat{\Sigma}$):

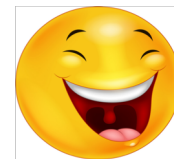
$\Theta_0 \leftarrow (\hat{\Sigma} + tI)^{-1}, \lambda_0 \leftarrow 1$

For $k = 0$ **to** $K - 1$ **do**

$\Theta_{k+1}, Z_{k+1}, \lambda_{k+1}$
 \leftarrow GLADcell($\hat{\Sigma}, \Theta_k, Z_k, \lambda_k$)

return Θ_K, Z_K

GLAD: Graph recovery **L**earning **A**lgorithm using **D**ata-driven training



Experiments: Convergence

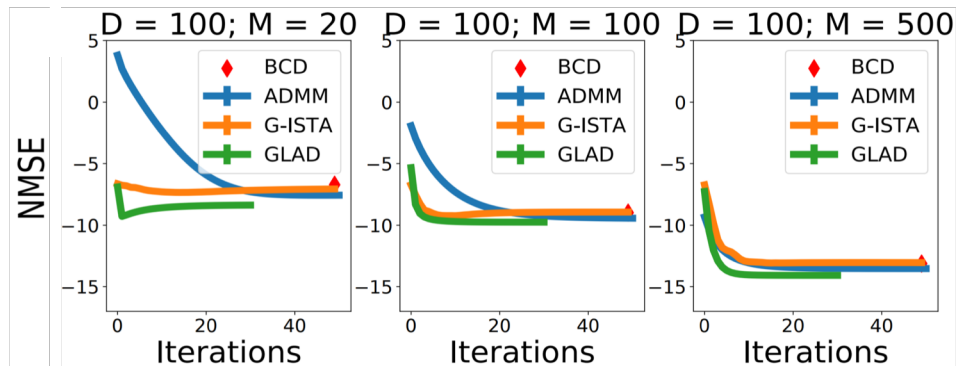
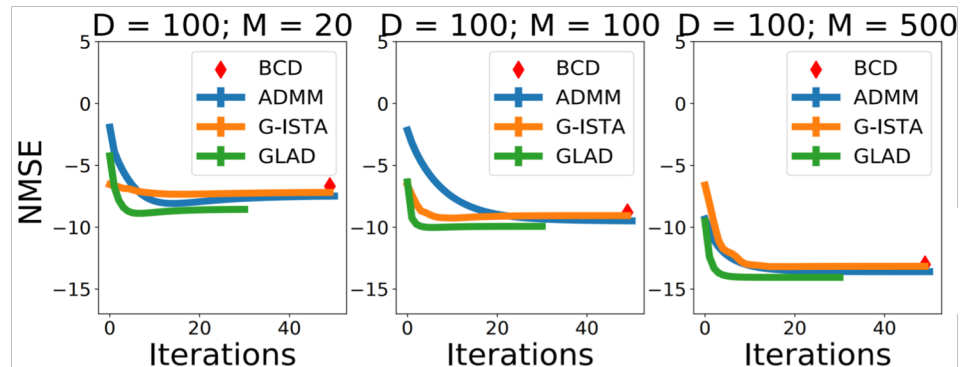
Train/finetuning
using 10 random
graphs

Test on 100
random graphs

Fixed Sparsity level
 $s=0.1$

GLAD vs
traditional
methods

Mixed Sparsity level
 $s \sim U(0.05, 0.15)$

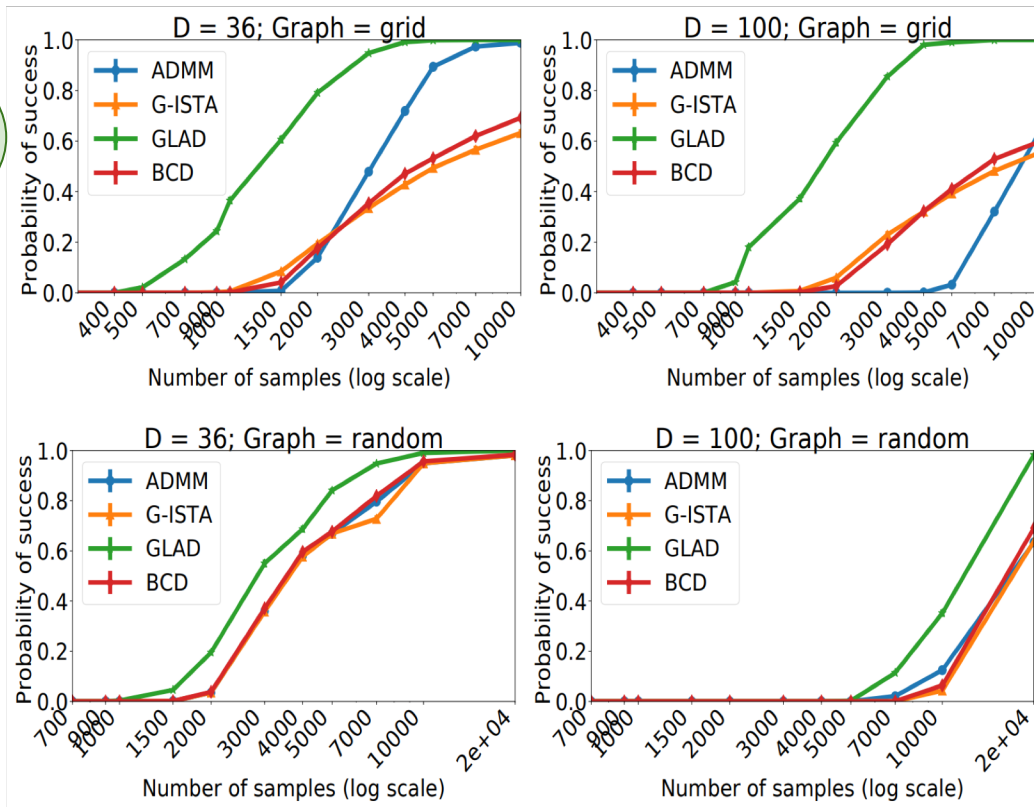


Experiments: Recovery probability

Sample complexity for model selection consistency

PS is non-zero if all graph edges are recovered with correct signs

GLAD able to recover true edges with considerably fewer samples



Experiments: Data Efficiency (cont...)

GLAD vs CNN*

Training graphs
100 vs 100,000

of parameters
<25 vs >>>25

Runtime
< 30 mins vs
several hours

Methods	M=15	M=35	M=100
BCD	0.578±0.006	0.639±0.007	0.704±0.006
CNN	0.664±0.008	0.738±0.006	0.759±0.006
CNN+P	0.672±0.008	0.740±0.007	0.771±0.006
GLAD	0.788±0.003	0.811±0.003	0.878±0.003

AUC on 100 test graphs, Gaussian random graph sparsity=0.05 and edge values sampled from $\sim U(-1, 1)$.

* DeepGraph-39 model from “Learning to Discover Sparse Graphical Models” by Belilovsky et. al.

Table 1. of Belilovsky et. al.

Gene Regulation Data: SynTReN details

Synthetic gene expression data generator creating biologically plausible networks

Models biological & correlation noises

SynTReN

The topological characteristics of generated networks closely resemble transcriptional networks

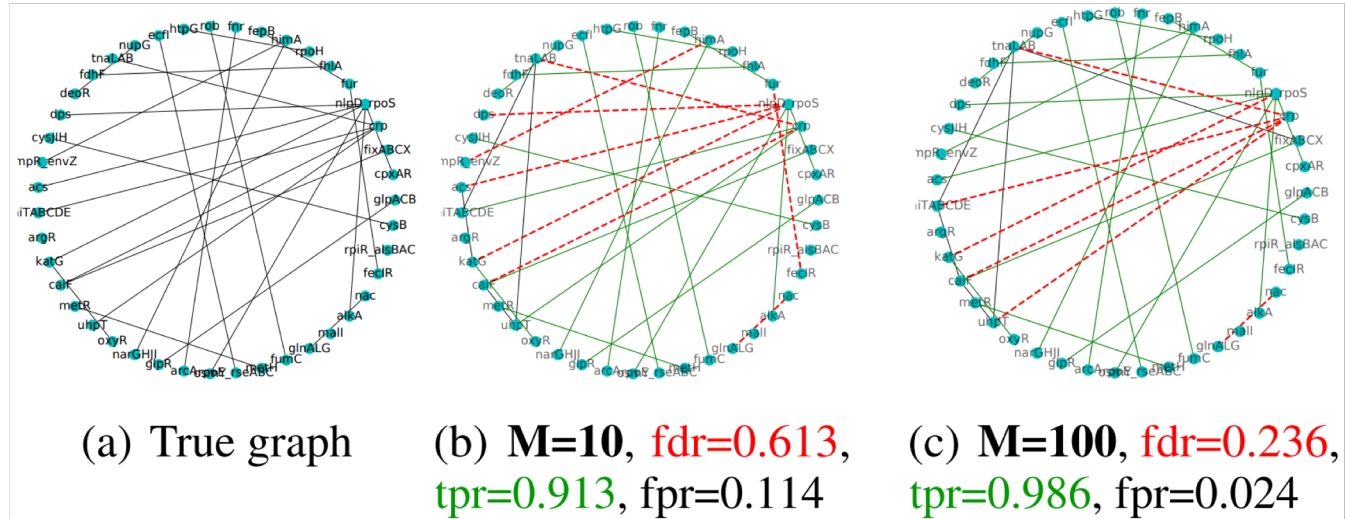
Contains instances of Ecoli bacteria and other true interaction networks

Gene Regulation Data: Ecoli Network Predictions

GLAD trained on
Erdos-Renyi
graphs of
dimension=25.

of train/valid
graphs were
20/20.

M samples were
generated per
graph



Recovered graph structures for a sub-network of the *E. coli* consisting of 43 genes and 30 interactions with increasing samples. All noises sampled $\sim U(0.01, 0.1)$
Increasing the samples reduces the fdr by discovering more true edges.

Theoretical Analysis: Assumptions

Assumption 1 (Tail conditions). *There exists $v_* \in (0, \infty]$ and a function $h : \mathbb{N} \times (0, \infty) \rightarrow (0, \infty)$ such that $\forall i, j, \mathbb{P}[|\widehat{\Sigma}_{ij}^m - (\Theta^*)_{ij}^{-1}| \geq \delta] \leq 1/h(m, \delta)$, $\forall \delta \in (0, 1/v_*]$. Further, assume h is monotonically increasing in sample size m and δ and define $\bar{\delta}_h(m, r) := \arg \max \{\delta | h(m, \delta) \leq r\}$.*

Ensures that sample sizes are large enough for an accurate estimation of the covariance matrix

Assumption 2 (Incoherence condition). *Denote the Hessian by $\Gamma^* := \Theta^{*-1} \otimes \Theta^{*-1}$, the indices of nonzero entries by $S := \{(i, j) | \Theta_{ij}^* \neq 0\}$ and its complement set by S^c . There exists $\alpha \in (0, 1]$ such that $\max_{e \in S^c} \|\Gamma_{eS}^* (\Gamma_{SS}^*)^{-1}\|_1 \leq (1 - \alpha)$.*

Restricts the interaction between edge and non-edge terms in the precision matrix

Consistency Analysis

Recalling AM

$$\Theta_{k+1}^{\text{AM}} \leftarrow \frac{1}{2} \left(\Theta_k^{\text{AM}} + \hat{\Theta}_\lambda \right)$$

$$Z_{k+1}^{\text{AM}} \leftarrow \eta_{\rho/\lambda} \left(Z_k^{\text{AM}} + \left(\frac{1}{\lambda} \hat{\Sigma} - Z_k^{\text{AM}} \right) \right)$$

An adaptive sequence of penalty parameters should achieve a better error bound

$$\frac{1}{\lambda} \hat{\Sigma} - Z_k^{\text{AM}}$$

$$\left(\|\theta\| - \rho/\lambda, 0 \right)$$

Summary

Under assumption 1 and 2, suppose that the sample size is larger than d^τ . Then with probability at least $1 - 1/d^{\tau-2}$,

$$\|\Theta_k^{\text{AM}} - \Theta^*\| \leq C_{\rho, \lambda, \Theta_0^{\text{AM}}, Z_0^{\text{AM}}} k^{-1/2} + \mathcal{O} \left(\delta_h(m, d^\tau) \lambda^{-1/2} \right) + \mathcal{O} \left(\delta_h(m, d^\tau) \right)$$

where $C_{\rho, \lambda, \Theta_0^{\text{AM}}, Z_0^{\text{AM}}}$ depends on the initialization $\|\Theta_0^{\text{AM}} - \hat{\Theta}_\lambda\|_F$ and λ .

Optimal parameter values depends on the tail behavior and the prediction error

Hard to choose these parameters manually

Conclusion

Unrolled DL architecture,
GLAD, for sparse graph
recovery

Empirically, GLAD is able to
reduce sample complexity

Empirical evidence that
learning can improve
graph recovery

Highlighting the potential of
using algorithms as
inductive bias for DL
architectures

Thank you!