

Hub Graphical Lasso

Kean Ming Tan

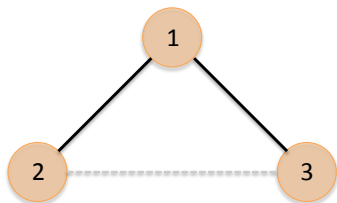
June 18, 2013

Joint work with Karthik Mohan, Palma London, Maryam Fazel,
Su-In Lee, and Daniela Witten
University of Washington

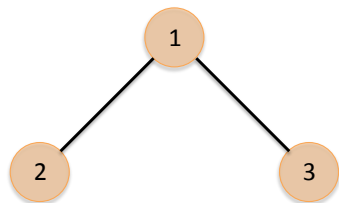
Arabidopsis Thaliana data set

- ▶ $n = 118$ samples and $p = 39$ genes
- ▶ Two isoprenoid biosynthesis pathways, Mevalonate Acid (MVA) and **Methylerythritol Phosphate (MEP)**
- ▶ Understand the conditional dependence relationship among the genes

Marginal versus conditional



Marginal



Conditional

Gaussian Graphical Model

- ▶ Suppose $\mathbf{x}_1, \dots, \mathbf{x}_n \stackrel{iid}{\sim} N(\mathbf{0}, \boldsymbol{\Sigma})$, where $\boldsymbol{\Sigma}$ is a $p \times p$ matrix
- ▶ We know $\boldsymbol{\Sigma}_{ij}^{-1} = 0$ implies the i th and j th variables are **conditionally independent**, given the other variables
- ▶ We can construct a conditional independence graph by estimating $\boldsymbol{\Sigma}^{-1}$
- ▶ The MLE for $\boldsymbol{\Sigma}^{-1}$ is \mathbf{S}^{-1} , where $\mathbf{S} = \frac{1}{n} \mathbf{X} \mathbf{X}^T$
- ▶ Not invertible and all entries are non-zero

The Graphical Lasso

- ▶ ℓ_1 penalized log-likelihood for **sparse estimation** of Σ^{-1}
- ▶ Let $\Theta = \Sigma^{-1}$ and $\hat{\Theta}$ be an estimate of Σ^{-1}

$$\underset{\Theta}{\text{minimize}} \{-\log \det \Theta + \text{trace}(\mathbf{S}\Theta) + \lambda \|\Theta\|_1\}$$

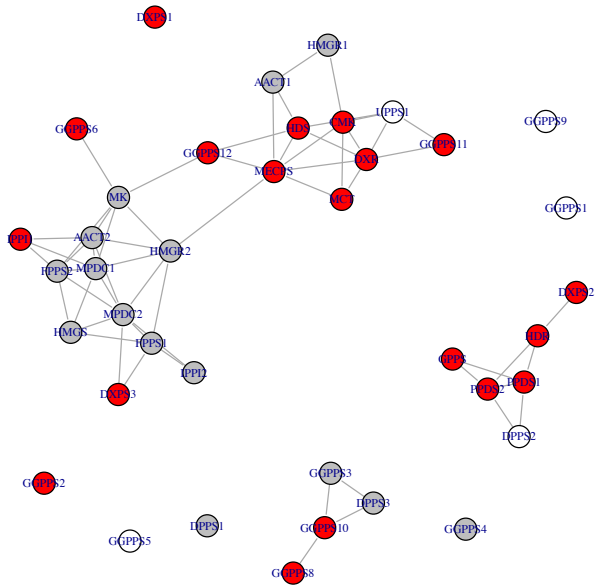
- ▶ λ controls sparsity of $\hat{\Theta}$

Citation: Friedman et al. 2007, Yuan and Lin 2006, Banerjee et al. 2008

Introduction

Hub Graphical Lasso

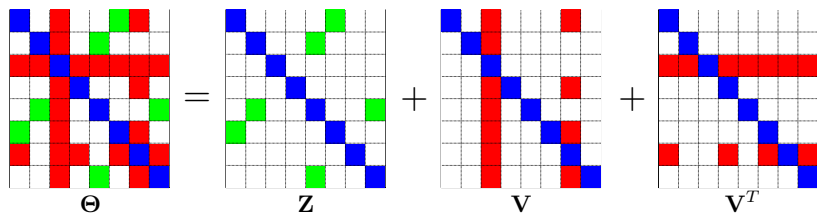
Simulation and real data



Problems with the Graphical Lasso

- ▶ Each edge is equally likely to appear, independent of all other edges
- ▶ Unrealistic in many settings
 - ▶ **Gene regulatory network** - transcription factors are connected to a number of genes
 - ▶ **World Wide Web** - Google is connected to most webpages
 - ▶ **Network of collaborations among scientists** - Paul Erdos collaborated with 511 scientists
- ▶ Want some **columns** of $\hat{\Theta}$ to be **dense**

Motivation of Hub Graphical Lasso (HGL)



- ▶ Z is sparse, something like the graphical lasso solution
- ▶ Non-zero columns of V corresponds to the hub nodes

Optimization problem

$$\begin{aligned} \underset{\Theta, \mathbf{V}, \mathbf{Z}}{\text{minimize}} \quad & \left\{ -\log\det\Theta + \text{trace}(\mathbf{S}\Theta) + \lambda_1\|\mathbf{Z} - \text{diag}(\mathbf{Z})\|_1 \right. \\ & \left. + \lambda_2\|\mathbf{V} - \text{diag}(\mathbf{V})\|_1 + \lambda_3\sum_{j=1}^p\|(\mathbf{V} - \text{diag}(\mathbf{V}))_j\|_2 \right\} \\ \text{subject to} \quad & \Theta = \mathbf{Z} + \mathbf{V} + \mathbf{V}^T \end{aligned}$$

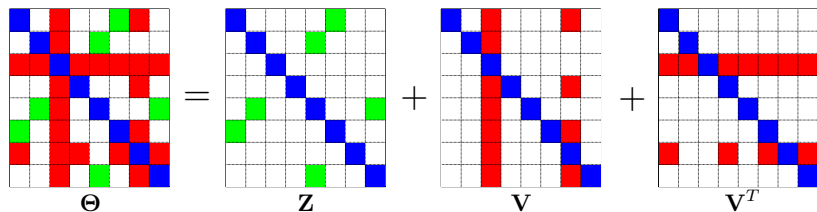
- ▶ λ_1 controls the **sparsity of the matrix \mathbf{Z}**
- ▶ λ_2 controls the **sparsity of the hub nodes**
- ▶ λ_3 controls the **selection of hub nodes**

Convex problem - solve using Alternating Direction Method of Multiplier

Measures of performance

- ▶ Correctly estimated edges
- ▶ Proportion of correctly estimated **hub edges**
- ▶ Proportion of correctly identified **hub nodes**
 - ▶ nodes that have at least **r edges** and also estimated to be **\mathcal{H} most highly-connected nodes**

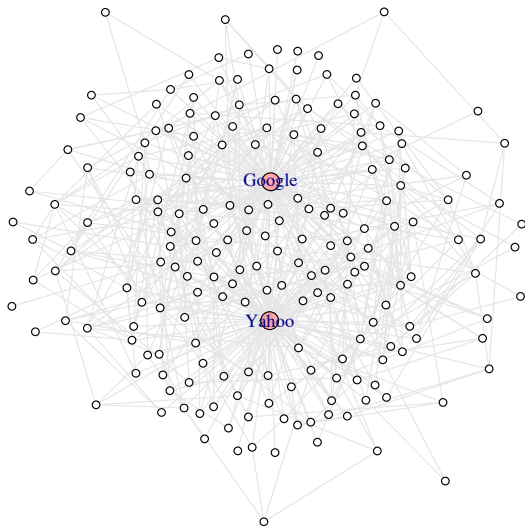
Simulation: Network with several hub nodes



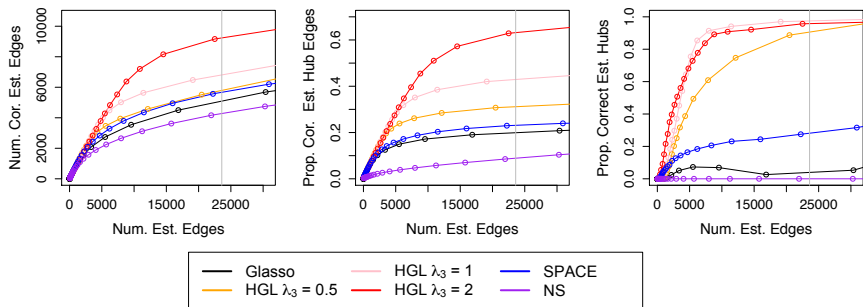
Generate $\Theta = \Sigma^{-1}$ as follows

- ▶ Matrix \mathbf{Z} is 98% sparse
- ▶ There are $\mathcal{H} = 20$ hub nodes (30% sparse) in the matrix \mathbf{V}
- ▶ Simulate data $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} N(0, \Sigma)$

Realization of such network for $p = 200$



Simulation: $n = 500$ $p = 1000$



Citation: Meinshausen and Bühlmann 2006, Friedman et al. 2007, Peng et al. 2009

Summary

- ▶ Propose HGL to identify hub nodes
- ▶ ADMM algorithm guarantees solutions to converge to the global minimum
- ▶ Future direction of research
 - ▶ Modeling binary network with hub nodes
 - ▶ Study the theoretical properties of the estimator

Question?

