

Selection for discrete graphical models: Generalized covariance matrices and their inverses

Martin Wainwright

UC Berkeley
Departments of Statistics, and EECS

Based on joint work with:

Po-Ling Loh (UC Berkeley)

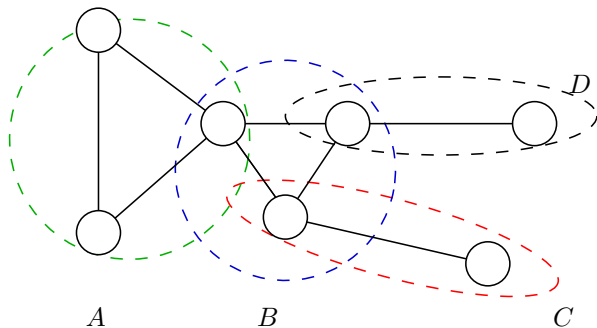
Graphical models for high-dimensional data

Used in a variety of applications:

- bioinformatics, genetics, proteomics
- natural language processing, speech processing
- social networks
- image processing, computer vision
- satisfiability problems, channel communications, data compression

Graph structure and factorization

- Markov random field: random vector (X_1, \dots, X_p) with distribution factoring according to a graph $G = (V, E)$:



- Hammersley-Clifford theorem: factorization over cliques

$$\mathbb{Q}(x_1, \dots, x_p; \theta) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{C \in \mathcal{C}} \theta_C(x_C) \right\}$$

Fitting graphical models to data

- drawn n samples from

$$Q_{\theta}(x_1, \dots, x_p) = \frac{1}{Z(\Theta)} \exp \left\{ \sum_{C \in \mathcal{C}} \theta_C(x_C) \right\}$$

- graph $G = (V, E)$ and clique potentials $\{\theta_C\}$ are unknown

Fitting graphical models to data

- drawn n samples from

$$\mathbb{Q}_\theta(x_1, \dots, x_p) = \frac{1}{Z(\Theta)} \exp \left\{ \sum_{C \in \mathcal{C}} \theta_C(x_C) \right\}$$

- graph $G = (V, E)$ and clique potentials $\{\theta_C\}$ are unknown
- use n samples to form data matrix $\mathbf{X} \in \mathcal{X}^{n \times p}$
- estimator $\mathbf{X} \mapsto (\hat{E}, \{\hat{\theta}_C\})$

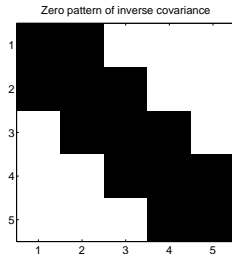
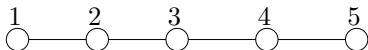
Fitting graphical models to data

- drawn n samples from

$$\mathbb{Q}_\theta(x_1, \dots, x_p) = \frac{1}{Z(\Theta)} \exp \left\{ \sum_{C \in \mathcal{C}} \theta_C(x_C) \right\}$$

- graph $G = (V, E)$ and clique potentials $\{\theta_C\}$ are unknown
- use n samples to form data matrix $\mathbf{X} \in \mathcal{X}^{n \times p}$
- estimator $\mathbf{X} \mapsto (\hat{E}, \{\hat{\theta}_C\})$
- various loss functions are possible:
 - ▶ exact graph selection: $\hat{E} = E?$
 - ▶ Hamming distance between \hat{E} and E
 - ▶ parameter estimation $\|\hat{\theta} - \theta\|$
 - ▶ bounds on Kullback-Leibler divergence $D(\mathbb{Q}_{\hat{\theta}} \parallel \mathbb{Q}_\theta)$

Illustration: Graphs and Gaussian inverse covariance



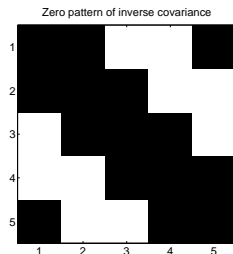
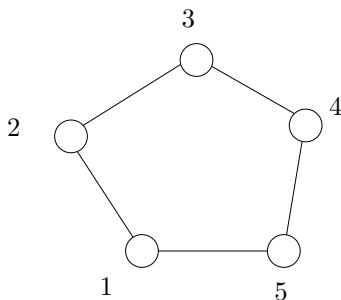
Multivariate Gaussian $(X_1, \dots, X_p) \sim N(0, \Theta^{-1})$:

$$\mathbb{Q}(x_1, \dots, x_p; \Theta) = \frac{\det(\Theta)}{(2\pi)^{p/2}} \exp\left(-\sum_{s \in V} \theta_{ss} x_s^2 - \sum_{(s,t) \in E} \theta_{st} x_s x_t\right).$$

Classical fact:

Inverse covariance matrix Θ of any multivariate Gaussian is graph-structured. (Consequence of Hammersley-Clifford theorem).

Illustration: Graphs and Gaussian inverse covariance



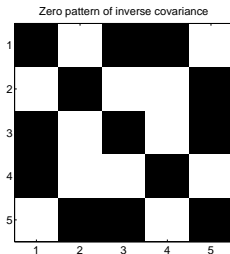
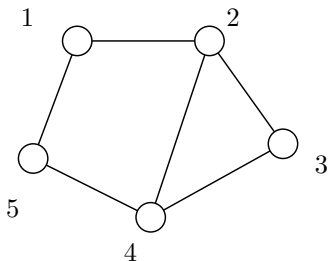
Multivariate Gaussian $(X_1, \dots, X_p) \sim N(0, \Theta^{-1})$:

$$\mathbb{Q}(x_1, \dots, x_p; \Theta) = \frac{\det(\Theta)}{(2\pi)^{p/2}} \exp \left(- \sum_{s \in V} \theta_{ss} x_s^2 - \sum_{(s,t) \in E} \theta_{st} x_s x_t \right).$$

Classical fact:

Inverse covariance matrix Θ of any multivariate Gaussian is graph-structured. (Consequence of Hammersley-Clifford theorem).

Illustration: Graphs and Gaussian inverse covariance



Multivariate Gaussian $(X_1, \dots, X_p) \sim N(0, \Theta^{-1})$:

$$Q(x_1, \dots, x_p; \Theta) = \frac{\det(\Theta)}{(2\pi)^{p/2}} \exp\left(-\sum_{s \in V} \theta_{ss} x_s^2 - \sum_{(s,t) \in E} \theta_{st} x_s x_t\right).$$

Classical fact:

Inverse covariance matrix Θ of any multivariate Gaussian is graph-structured. (Consequence of Hammersley-Clifford theorem).

Utility for Gaussian models

Sparsity of inverse covariance can be exploited in many ways:

- global Gaussian likelihood (log-determinant) with ℓ_1 -regularization (d'Aspremont et al., 2007; Friedman et al., 2008; Ravikumar et al., 2009)
- various neighborhood-based approaches:
 - ▶ neighborhood Lasso method (Meinshausen & Bühlmann, 2006)
 - ▶ neighborhood Dantzig selector (Yuan, 2010)
- CLIME estimator for inverse covariance Θ (Cai et al., 2011)
- non-convex Lasso for Gaussians with missing data (Loh & Wainwright, 2011)
- copula methods: univariate transformations of Gaussians (Liu, Lafferty & Wasserman, 2011)

Utility for Gaussian models

Sparsity of inverse covariance can be exploited in many ways:

- global Gaussian likelihood (log-determinant) with ℓ_1 -regularization (d'Aspremont et al., 2007; Friedman et al., 2008; Ravikumar et al., 2009)
- various neighborhood-based approaches:
 - ▶ neighborhood Lasso method (Meinshausen & Bühlmann, 2006)
 - ▶ neighborhood Dantzig selector (Yuan, 2010)
- CLIME estimator for inverse covariance Θ (Cai et al., 2011)
- non-convex Lasso for Gaussians with missing data (Loh & Wainwright, 2011)
- copula methods: univariate transformations of Gaussians (Liu, Lafferty & Wasserman, 2011)

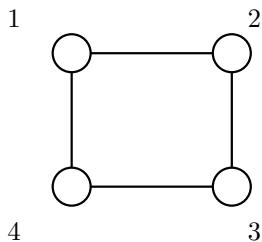
A natural question:

Is there any analog of this type of correspondence for non-Gaussian models?

Simplest version is for the *Ising model*:

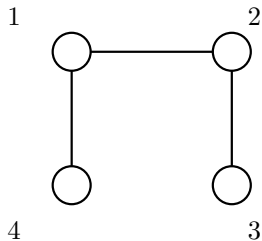
$$\mathbb{Q}_\theta(x_1, \dots, x_p) \propto \exp \left\{ \sum \theta_s x_s + \sum \theta_{st} x_s x_t \right\} \quad \text{where } X_s \in \{-1, +1\}.$$

A counter example



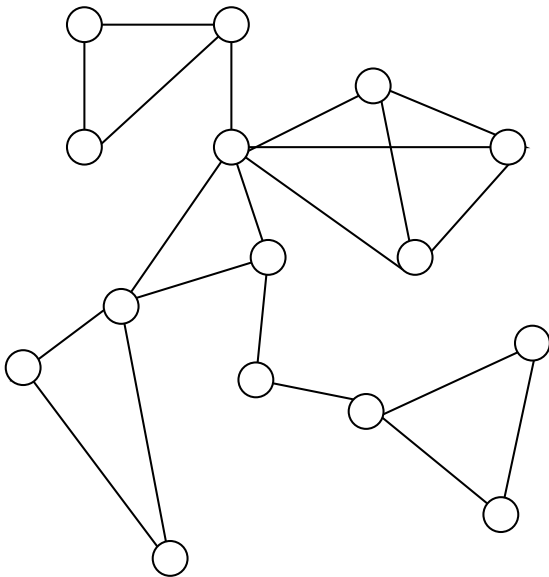
	X_1	X_2	X_3	X_4
X_1				
X_2				
X_3				
X_4				

A success

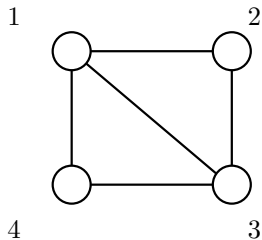


	X_1	X_2	X_3	X_4
X_1	Shaded	Shaded	White	Shaded
X_2	Shaded	Shaded	Shaded	White
X_3	White	Shaded	Shaded	White
X_4	Shaded	White	White	Shaded

A bigger success: Dinosaur graph

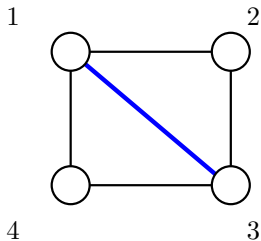


A triangulated counter example



	X_1	X_2	X_3	X_4
X_1				
X_2				
X_3				
X_4				

Generalized moment matrices: a success



	X_1	X_2	X_3	X_4	X_1X_3
X_1	■	■	■	■	■
X_2	■	■	■	□	■
X_3	■	■	■	■	■
X_4	■	□	■	■	■
X_1X_3	■	■	■	■	■

Generalized moment matrices

- discrete graphical model with $x_j \in \mathcal{X} = \{0, 1, \dots, m-1\}$:
- for each clique C , define subset of possible configurations

$$\mathcal{X}_0^{|C|} := \{j_C = (j_s, s \in C) \mid j_s \neq 0 \text{ for all } s \in C\}.$$

- parameterization as a minimal exponential family:

$$\mathbb{Q}_\theta(x_1, \dots, x_p) = \exp \left\{ \sum_{C \in \mathcal{C}} \sum_{J \in \mathcal{X}_0^{|C|}} \theta_{C;J} \mathbb{I}[x_C = J] - A(\theta) \right\},$$

- for any set \mathcal{S} of subsets of vertices, define random vector:

$$\Phi(X; \mathcal{S}) = \{\mathbb{I}[X_C = J], J \in \mathcal{X}_0^{|C|}, S \in \mathcal{S}\}.$$

Generalized inverse covariances and graph structure

- for any set \mathcal{S} of subsets of vertices, define random vector:

$$\Phi(X; \mathcal{S}) = \{\mathbb{I}[X_C = J], J \in \mathcal{X}_0^{|C|}, S \in \mathcal{S}\}.$$

- for any set of subsets \mathcal{S} , define

$$\text{pow}(\mathcal{S}) = \{\text{pow}(S), | S \in \mathcal{S}\}.$$

Theorem (Loh & Wainwright, 2012)

For any clique set \mathcal{T} in some triangulation of G , the inverse Γ of the augmented covariance matrix $\text{cov}(\Phi(X; \text{pow}(\mathcal{T})))$ is block graph-structured in the following sense:

- (a)** *For any two subsets A and B which are not subsets of the same maximal clique, the block $\Gamma(\text{pow}(A), \text{pow}(B))$ is zero.*
- (b)** *For a generic parameter vector θ , the block $\Gamma(\text{pow}(A), \text{pow}(B))$ is non-zero whenever A and B belong to a common maximal clique.*

Separator sets are sufficient

- consider any triangulated version $\tilde{G} = (V, \tilde{E})$
- let \mathcal{S} be set of separator sets in \tilde{G}

Corollary

Let Γ be the inverse of $\text{cov}(\Phi(X; V \cup \text{pow}(\mathcal{S})))$. Then $\Gamma_{V,V}$ is block edge-structured, meaning that

$$(s, t) \notin \tilde{E} \iff \Gamma_{s,t} = 0.$$

Corollary

For any graph with singleton separator sets, the inverse covariance matrix is graph-structured. (Includes forests/trees as a special case).

Proof ideas

Key ingredients:

- Legendre duality and minimal exponential families (e.g., Brown, 1986)
- Junction tree representation (e.g., Lauritzen & Spiegelhalter, 1988)
- Entropy decompositions and marginal polytopes (Wainwright & Jordan, 2008)

Proof ideas

Key ingredients:

- Legendre duality and minimal exponential families (e.g., Brown, 1986)
- Junction tree representation (e.g., Lauritzen & Spiegelhalter, 1988)
- Entropy decompositions and marginal polytopes (Wainwright & Jordan, 2008)
- Cumulant generating function of our exponential family:

$$A(\theta) := \log \left\{ \sum_x \exp \left(\sum_{C \in \mathcal{C}} \sum_{j_C \in \mathcal{X}_0^{|C|}} \theta_{C; j_C} \mathbb{I}[x_C = j_C] \right) \right\}.$$

Proof ideas

Key ingredients:

- Legendre duality and minimal exponential families (e.g., Brown, 1986)
- Junction tree representation (e.g., Lauritzen & Spiegelhalter, 1988)
- Entropy decompositions and marginal polytopes (Wainwright & Jordan, 2008)
- Cumulant generating function of our exponential family:

$$A(\theta) := \log \left\{ \sum_x \exp \left(\sum_{C \in \mathcal{C}} \sum_{j_C \in \mathcal{X}_0^{C|c}} \theta_{C;J_C} \mathbb{I}[x_C = J_C] \right) \right\}.$$

- strictly convex with second partial derivatives

$$\frac{\partial A}{\partial \theta_{C;J_C} \theta_{D;K_D}}(\theta) = \text{cov} \left\{ \mathbb{I}[X_C = J_C], \mathbb{I}[X_D = K_D] \right\}.$$

- result concerns the structure of $(\nabla^2 A(\theta))^{-1}$: connect this to the Legendre dual A^* of A .

Some statistical consequences

- regularized log-determinant methods for discrete graphical models (Banerjee, d'Aspremont & El Ghaoui, 2009):
 - ▶ consistent for “good” graphs (i.e., singleton separator sets, including trees)
 - ▶ inconsistent in general; partial consistency results are possible

Some statistical consequences

- regularized log-determinant methods for discrete graphical models (Banerjee, d'Aspremont & El Ghaoui, 2009):
 - ▶ consistent for “good” graphs (i.e., singleton separator sets, including trees)
 - ▶ inconsistent in general; partial consistency results are possible
- local methods for discrete graphical selection (Loh & W., 2012)
 - ▶ for “good” graph, ordinary nodewise linear regression is consistent:

$$\widehat{\theta}[s] \in \arg \min_{\theta \in \mathbb{R}^{p-1}} \left\{ \frac{1}{2} \|X_s - X_{\setminus\{s\}} \theta\|_2^2 + \lambda_n^2 \|\theta\|_1^2 \right\}.$$

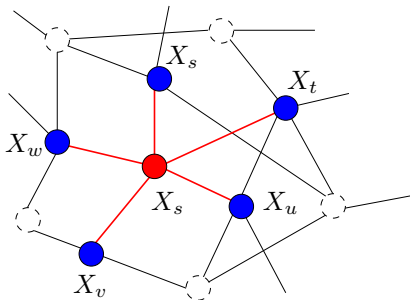
- ▶ for graphs with maximum degree d , nodewise regression with lifted moment matrices of order d are consistent

Markov property and neighborhood structure

- Markov properties encode neighborhood structure:

$$\underbrace{(X_s \mid X_{V \setminus s})}_{\text{Condition on full graph}} \stackrel{d}{=} \underbrace{(X_s \mid X_{N(s)})}_{\text{Condition on Markov blanket}}$$

$N(s) = \{s, t, u, v, w\}$



- pseudolikelihood method (Besag, 1974)
- basis of many graph learning algorithms (Spirites et al., 2000; Abeel et al., 2006; Meinshausen & Buhlmann, 2006)

Two methods for discrete binary models

Method 1: ℓ_1 -regularized logistic regression

- perform **logistic** regression of each variable X_s on $X_{V \setminus \{s\}}$
 - ℓ_1 -regularization to enforce sparsity
 - known to succeed w.h.p. with sample size $n \gtrsim d^2 \log p$, where d is maximum degree (Ravikumar et al., 2010)
-

Two methods for discrete binary models

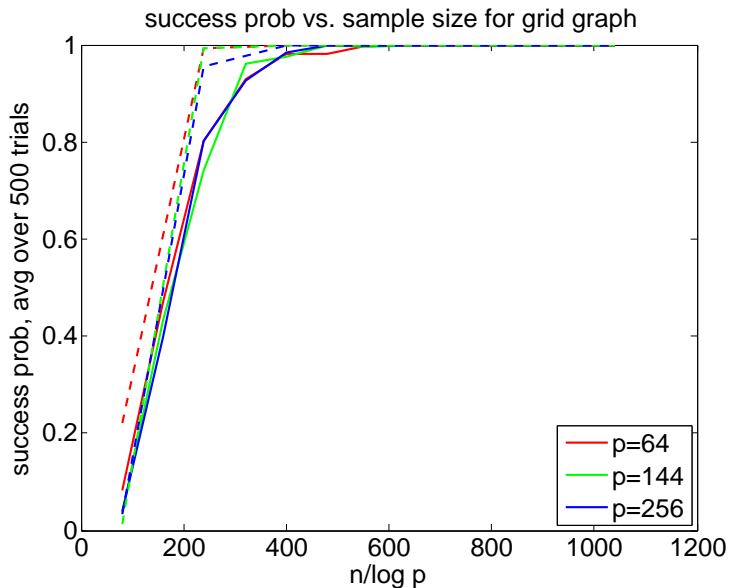
Method 1: ℓ_1 -regularized logistic regression

- perform **logistic** regression of each variable X_s on $X_{V \setminus \{s\}}$
 - ℓ_1 -regularization to enforce sparsity
 - known to succeed w.h.p. with sample size $n \gtrsim d^2 \log p$, where d is maximum degree (Ravikumar et al., 2010)
-

Method 2: ℓ_1 -regularized **linear** regression

- perform linear regression of each variable X_s on the d -extended version of $X_{V \setminus \{s\}}$
- ℓ_1 -regularization to enforce sparsity
- involves $\mathcal{O}(p^d)$ covariates, so suitable only for bounded degree graphs
- extends naturally to missing/corrupted covariates

Logistic versus linear regression for grid



Straightforward extension to missing data

- observe corrupted version $\tilde{Z} \in \mathbb{R}^{n \times (p-1)}$ of matrix $X_{\setminus \{s\}} \in \mathbb{R}^{n \times (p-1)}$:

$$\tilde{Z}_{ij} = \begin{cases} X_{ij} & \text{with probability } 1 - \alpha \\ \star & \text{with probability } \alpha. \end{cases}$$

Straightforward extension to missing data

- observe corrupted version $\tilde{Z} \in \mathbb{R}^{n \times (p-1)}$ of matrix $X_{\setminus\{s\}} \in \mathbb{R}^{n \times (p-1)}$:

$$\tilde{Z}_{ij} = \begin{cases} X_{ij} & \text{with probability } 1 - \alpha \\ \star & \text{with probability } \alpha. \end{cases}$$

- set $\star \equiv 0$ and $\hat{Z} = \frac{\tilde{Z}}{1-\alpha}$, and form quantities

$$\hat{\Gamma} = \frac{\hat{Z}^T \hat{Z}}{n} - \alpha \text{diag} \left(\frac{\hat{Z}^T \hat{Z}}{n} \right), \quad \text{and} \quad \hat{\gamma} = \frac{\hat{Z}^T y}{n},$$

- unbiased estimates of $\Gamma = \text{cov}(X_{\setminus\{s\}})$ and $\gamma = \text{cov}(X_s, X_{\setminus\{s\}})$

Straightforward extension to missing data

- observe corrupted version $\tilde{Z} \in \mathbb{R}^{n \times (p-1)}$ of matrix $X_{\setminus\{s\}} \in \mathbb{R}^{n \times (p-1)}$:

$$\tilde{Z}_{ij} = \begin{cases} X_{ij} & \text{with probability } 1 - \alpha \\ \star & \text{with probability } \alpha. \end{cases}$$

- set $\star \equiv 0$ and $\hat{Z} = \frac{\tilde{Z}}{1-\alpha}$, and form quantities

$$\hat{\Gamma} = \frac{\hat{Z}^T \hat{Z}}{n} - \alpha \text{diag} \left(\frac{\hat{Z}^T \hat{Z}}{n} \right), \quad \text{and} \quad \hat{\gamma} = \frac{\hat{Z}^T y}{n},$$

- unbiased estimates of $\Gamma = \text{cov}(X_{\setminus\{s\}})$ and $\gamma = \text{cov}(X_s, X_{\setminus\{s\}})$
- solve optimization problem:

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^{p-1}} \left\{ \frac{1}{2} \theta^T \hat{\Gamma} \theta - \langle \hat{\gamma}, \theta \rangle + \lambda_n^2 \|\theta\|_1^2 \right\}$$

Straightforward extension to missing data

- observe corrupted version $\tilde{Z} \in \mathbb{R}^{n \times (p-1)}$ of matrix $X_{\setminus\{s\}} \in \mathbb{R}^{n \times (p-1)}$:

$$\tilde{Z}_{ij} = \begin{cases} X_{ij} & \text{with probability } 1 - \alpha \\ \star & \text{with probability } \alpha. \end{cases}$$

- set $\star \equiv 0$ and $\hat{Z} = \frac{\tilde{Z}}{1-\alpha}$, and form quantities

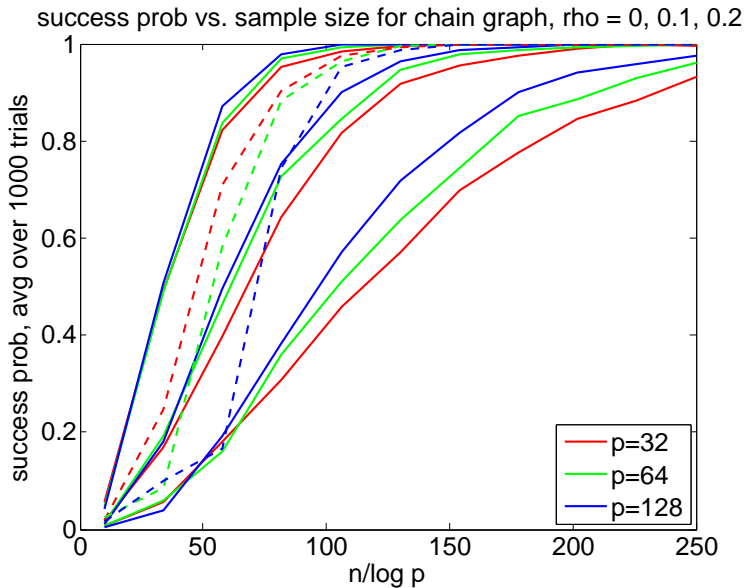
$$\hat{\Gamma} = \frac{\hat{Z}^T \hat{Z}}{n} - \alpha \text{diag} \left(\frac{\hat{Z}^T \hat{Z}}{n} \right), \quad \text{and} \quad \hat{\gamma} = \frac{\hat{Z}^T y}{n},$$

- unbiased estimates of $\Gamma = \text{cov}(X_{\setminus\{s\}})$ and $\gamma = \text{cov}(X_s, X_{\setminus\{s\}})$
- solve optimization problem:

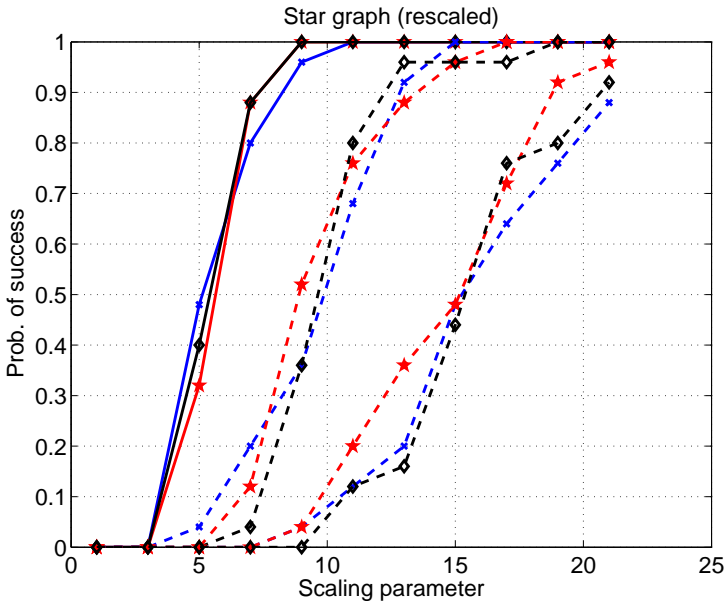
$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^{p-1}} \left\{ \frac{1}{2} \theta^T \hat{\Gamma} \theta - \langle \hat{\gamma}, \theta \rangle + \lambda_n^2 \|\theta\|_1^2 \right\}$$

- “corrected” estimates of this type:
 - Xu & You (2007): corrected Lasso in classical setting
 - Loh & Wainwright (2012): corrected Lasso and algorithms for non-convexity

Performance with missing data: Chain graph



Performance with missing data: Star graph



Statistical guarantees

- take n independent samples $\{X_i\}_{i=1}^n$ from Ising model with zero mean, and inverse covariance matrix $\Theta^* = (\text{cov}(X))^{-1}$
- entries of samples $\{X_i\}_{i=1}^n$ missing independently with probability α
- graph has maximum degree d , and singleton separator sets

Statistical guarantees

- take n independent samples $\{X_i\}_{i=1}^n$ from Ising model with zero mean, and inverse covariance matrix $\Theta^* = (\text{cov}(X))^{-1}$
- entries of samples $\{X_i\}_{i=1}^n$ missing independently with probability α
- graph has maximum degree d , and singleton separator sets

Theorem (Loh & W., 2012)

As long as $n \gtrsim \frac{d^2 \log p}{(1-\alpha)^2}$ and $\lambda_n \asymp \sqrt{\frac{\log p}{n}}$, then with probability greater than $1 - 2 \exp(-cn\lambda_n^2)$, any global optimum $\hat{\theta}$ satisfies the bound

$$\|\hat{\theta} - \theta^*\|_\infty \lesssim \|\Theta^*\|_\infty \sqrt{\frac{\log p}{n}}.$$

Consequently, the method is graph-selection consistent as long as

$$\min_{t \in \mathcal{N}(s)} |\theta_{st}^*| \gtrsim \sqrt{\frac{\log p}{n}}.$$

Summary

- graphical model selection:
 - ▶ wide-range of applications
 - ▶ convex relaxations play an important role
 - ▶ Gaussian data benefits from link to inverse covariance matrix

- general connection between graph structure and inverse covariance
 - ▶ depends on structure of separator sets in a junction tree
 - ▶ inverses of higher-order moment matrices are useful

- various statistical consequences for discrete graph selection

Summary

- graphical model selection:
 - ▶ wide-range of applications
 - ▶ convex relaxations play an important role
 - ▶ Gaussian data benefits from link to inverse covariance matrix

- general connection between graph structure and inverse covariance
 - ▶ depends on structure of separator sets in a junction tree
 - ▶ inverses of higher-order moment matrices are useful

- various statistical consequences for discrete graph selection

- many remaining questions:
 - ▶ how to handle more realistic models of missing data?
 - ▶ provable guarantees for other types of non-convex M -estimators?
 - ▶ discrete graphical models with hidden variables?

Some papers/pre-prints

- Agarwal, Negahban & W. (2012). Fast convergence rates for high-dimensional statistical recovery. *Annals of Statistics*, 40(5): 2452–2482, December.
- Loh & W. (2012). High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity. *Annals of Statistics*, 40(3): 1637–1664, September.
- Loh & W. (2012). Structure estimation for discrete graphical models: Generalized covariance matrices and their inverses *Arxiv pre-print*.