inference in large-scale networks







neighborhood $ne(i) = \{j \in V \setminus \{i\} : (i, j) \in E\}$



neighborhood $ne(i) = \{j \in V \setminus \{i\} : (i, j) \in E\}$

 $ne(1) = \{2, 3, 4\}$ $ne(2) = \{1\}$ $ne(3) = \{1\}$ $ne(4) = \{1\}$



neighborhood $ne(i) = \{j \in V \setminus \{i\} : (i, j) \in E\}$

 $ne(1) = \{2, 3, 4\}$ $ne(1) = \{2, 3, 4\}$ $ne(2) = \{1\}$ $ne(2) = \{1, 3, 4\}$ $ne(3) = \{1\}$ $ne(3) = \{1, 2, 4\}$ $ne(4) = \{1\}$ $ne(4) = \{1, 2, 3\}$

graphical models and networks



description in terms of cliques

graphical models and networks



conditional independencies of graph

Ø

 $X_2 \perp X_3 \mid X_1$ $X_2 \perp X_4 \mid X_1$ $X_3 \perp X_4 \mid X_1$

Gaussian graphical models 2 4 3 true graph $\nabla -1$ $\boldsymbol{\nu}$

$$\begin{array}{ccccc} \mathbf{K} &= \mathbf{2} &= & \\ \begin{pmatrix} 1 & 1/2 & 1/2 & 1/2 \\ 1/2 & 1 & 0 & 0 \\ 1/2 & 0 & 1 & 0 \\ 1/2 & 0 & 0 & 1 \\ \end{array} \right)$$







combine neighborhoods



combine neighborhoods

• and node X_3 is in ne(1) if $\beta_{13} \neq 0$ and $\beta_{31} \neq 0$



combine neighborhoods

- and node X_3 is in ne(1) if $\beta_{13} \neq 0$ and $\beta_{31} \neq 0$
- or node X_3 is in ne(1) if $\beta_{13} \neq 0$ or $\beta_{31} \neq 0$



$$K = \Sigma^{-1} = \begin{pmatrix} 1 & 1/2 & 1/2 & 1/2 \\ 1/2 & 1 & 0 & 0 \\ 1/2 & 0 & 1 & 0 \\ 1/2 & 0 & 0 & 1 \end{pmatrix}$$

• let $K = \Sigma^{-1}$ then the neighborhood of X_1 is

 $\mathsf{ne}(1) = \{i \in V \setminus \{1\} : K_{1i} \neq 0\} = \{2, 3, 4\}$



• let $K = \Sigma^{-1}$ then the neighborhood of X_1 is

 $\mathsf{ne}(1) = \{i \in V \setminus \{1\} : K_{1i} \neq 0\} = \{2, 3, 4\}$

• but the regression coefficient is $\beta_{13} = -K_{13}/K_{11}$ [Lauritzen, 1996, chap. 5], and so

$$\mathsf{ne}(1) = \{i \in V \setminus \{1\} : \beta_{1i} \neq 0\} = \{2, 3, 4\}$$

nodewise estimation n > p



• estimate regression coefficients β_{ij} with LS

$$L_{\beta} := ||Y - X\beta||_{2}^{2} = (Y - X\beta)'(Y - X\beta)$$

 do this for each node separately and combine coefficients (and or or)

algorithm Meinshausen-Bühlmann

Algorithm 1 network $\hat{\Gamma}$ with MB lasso

- 1: for $i \in V$ do
- 2: split data Z in $Y = Z_i$ and $X = Z_{V \setminus i}$
- 3: obtain lasso $\hat{\beta}_L$ for $Y = X\beta + \varepsilon$

4: let
$$\hat{\tau}_{i}^{2} = ||Y - X\hat{\beta}_{L}||_{2}^{2} + \lambda ||\hat{\beta}_{L}||_{1}$$
 and

$$\hat{\gamma}_{ij} = -\hat{\beta}_{Lj}/\hat{\tau}_i^2 \qquad \text{for } i \neq j \hat{\gamma}_{ij} = 1/\hat{\tau}_i^2 \qquad \text{for } i = j$$

5: **end for**

6: output $\hat{\Gamma}$ (perhaps using and or or rule)

relevance of selection/regularisation

linear regression model

$$Y = \beta_{0,0} + X_1 \beta_{0,1} + X_2 \beta_{0,2} + \dots + X_p \beta_{0,p} + \epsilon$$

- increase prediction accuracy (decrease prediction error) by selecting a subset of predictors
- increase model interpretability by selecting only 'relevant' (correlated to Y) predictors





regularization

lasso variants



$$Y = \beta_{0,0} + X_1 \beta_{0,1} + X_2 \beta_{0,2} + \dots + X_p \beta_{0,p} + \epsilon$$

$$Y = \beta_{0,0} + X_1 \beta_{0,1} + X_2 \beta_{0,2} + \dots + X_p \beta_{0,p} + \epsilon$$

• best subset selection: for each number of predictors k, test $\binom{p}{k}$ models, select the best from all

- best subset selection: for each number of predictors k, test $\binom{p}{k}$ models, select the best from all
- forward selection: start with no predictors, add the best predictor one at a time

- best subset selection: for each number of predictors k, test $\binom{p}{k}$ models, select the best from all
- forward selection: start with no predictors, add the best predictor one at a time
- backward selection: start with all predictors, remove the worst predictor one at a time

$$Y = \beta_{0,0} + X_1 \beta_{0,1} + X_2 \beta_{0,2} + \dots + X_p \beta_{0,p} + \epsilon$$

$$Y = \beta_{0,0} + X_1 \beta_{0,1} + X_2 \beta_{0,2} + \dots + X_p \beta_{0,p} + \epsilon$$

 best subset selection: complexity is 2^p, e.g., 20 predictors gives 1.048.576 models

- best subset selection: complexity is 2^p, e.g., 20 predictors gives 1.048.576 models
- forward selection: complexity is p(p + 1)/2 + 1, e.g., 20 predictors gives 211 models

- best subset selection: complexity is 2^p, e.g., 20 predictors gives 1.048.576 models
- forward selection: complexity is p(p + 1)/2 + 1, e.g., 20 predictors gives 211 models
- backward selection: complexity is p(p+1)/2 + 1

- best subset selection: complexity is 2^p, e.g., 20 predictors gives 1.048.576 models
- forward selection: complexity is p(p + 1)/2 + 1, e.g., 20 predictors gives 211 models
- backward selection: complexity is p(p+1)/2 + 1
- no guarantees that the true model is selected (asymptotically)

 $Y = \beta_{0,0} + X_1 \beta_{0,1} + X_2 \beta_{0,2} + \dots + X_p \beta_{0,p} + \epsilon = X \beta_0 + \epsilon$

 $Y = \beta_{0,0} + X_1 \beta_{0,1} + X_2 \beta_{0,2} + \dots + X_p \beta_{0,p} + \epsilon = X \beta_0 + \epsilon$

problem when p > n

• design matrix X is not of full rank

$$Y = \beta_{0,0} + X_1 \beta_{0,1} + X_2 \beta_{0,2} + \dots + X_p \beta_{0,p} + \epsilon = X \beta_0 + \epsilon$$

problem when p > n

$$Y = X(\beta_0 + \mathbf{u}) + \varepsilon$$

- design matrix X is not of full rank
- where *u* is in the null-space of *X*, i.e.

$$\mathsf{null}(X) = \{ u \in \mathbb{R}^p : Xu = 0 \}$$

$$Y = \beta_{0,0} + X_1 \beta_{0,1} + X_2 \beta_{0,2} + \dots + X_p \beta_{0,p} + \epsilon = X \beta_0 + \epsilon$$



- design matrix X is not of full rank
- where *u* is in the null-space of *X*, i.e.

$$\mathsf{null}(X) = \{u \in \mathbb{R}^p : Xu = 0\}$$

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \psi(\beta)$$

LS fit penalty

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \psi(\beta)$$

LS fit penalty

types	of penaliz	zation	
		ridge ℓ_2	
	ψ	$\sum_{i=1} \beta_i^2$	
	bias	all β_i	
	treats β s	unequally	

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \psi(\beta)$$

LS fit penalty

types	of penaliz	zation		
		ridge ℓ_2	lasso ℓ_1	
	ψ bias treats β s	$\sum_{i=1} eta_i^2$ all eta_i unequally	$\begin{array}{l} \sum_{i=1} \beta_i \\ \text{small } \beta_i \\ \text{equally} \end{array}$	-
				-

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \psi(\beta)$$

LS fit penalty



solution when p > n



types of penalization

	ridge ℓ_2	lasso ℓ_1	ℓ_0
ψ	$\sum_{i=1} \beta_i^2$	$\sum_{i=1} \beta_i $	р
bias	all β_i	small β_i	all β_i
treats β s	unequally	equally	equally
adaptive lasso

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \sum_{i} \frac{|\beta_i|}{|\hat{\beta}_i^{\mathsf{init}}|}$$

penalty

where β_{j}^{init} is an initial estimate (prediction, with CV λ).

adaptive lasso

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \sum_{i} \frac{|\beta_i|}{|\hat{\beta}_i^{\mathsf{init}}|}$$

penalty

where β_j^{init} is an initial estimate (prediction, with CV λ). **properties of adaptive lasso** • when $\hat{\beta}_i^{\text{init}} = 0$ then $\hat{\beta}_j^{\text{ad}} = 0$

adaptive lasso

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \sum_{i} \frac{|\beta_i|}{|\hat{\beta}_i^{\mathsf{init}}|}$$

penalty

where β_j^{init} is an initial estimate (prediction, with CV λ). **properties of adaptive lasso** • when $\hat{\beta}_i^{\text{init}} = 0$ then $\hat{\beta}_j^{\text{ad}} = 0$ • when $\hat{\beta}_i^{\text{init}}$ is large, $\hat{\beta}_j^{\text{ad}}$ has small penalty, resulting in less bias

thresholded lasso

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \psi(\beta)$$

LS fit penalty

and select for some constant c > 0

 $\hat{\beta}_j \mathbb{1}\{\hat{\beta}_j > c\}$

[Bühlmann and van der Geer, 2011]

thresholded lasso

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \psi(\beta)$$

LS fit penalty

and select for some constant c > 0

$$\hat{\beta}_j \mathbb{1}\{\hat{\beta}_j > c\}$$

properties of thresholded lasso

• screening property $S_0 \subseteq \hat{S}$

thresholded lasso

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}} (Y - X\beta) + \lambda \psi(\beta)$$

LS fit penalty

and select for some constant c > 0

$$\hat{\beta}_j \mathbb{1}\{\hat{\beta}_j > c\}$$

properties of thresholded lasso

• screening property $S_0 \subseteq \hat{S}$

more 'accurate' than adaptive lasso (smaller predictive and estimation error)

scaled lasso

iterate as long as $L_{\hat{\beta}(\hat{\sigma}\lambda_0)} \leq L_{\hat{\beta},\lambda}$

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \psi(\beta)$$

LS fit penalty
 $\hat{\beta} = \hat{\beta}(\hat{\sigma}\lambda_0)$ and $\hat{\sigma} = \mathsf{LS}^{1/2}/n$

scaled lasso

iterate as long as $L_{\hat{\beta}(\hat{\sigma}\lambda_0)} \leq L_{\hat{\beta},\lambda}$

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \psi(\beta)$$

LS fit penalty
 $\hat{\beta} = \hat{\beta}(\hat{\sigma}\lambda_0)$ and $\hat{\sigma} = \mathsf{LS}^{1/2}/n$

o screening property $S_0 \subseteq \hat{S}$

scaled lasso

iterate as long as $L_{\hat{\beta}(\hat{\sigma}\lambda_0)} \leq L_{\hat{\beta},\lambda}$

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \psi(\beta)$$

LS fit penalty
 $\hat{\beta} = \hat{\beta}(\hat{\sigma}\lambda_0)$ and $\hat{\sigma} = \mathsf{LS}^{1/2}/n$

properties of scaled lasso

- screening property $S_0 \subseteq \hat{S}$
- $\circ\,$ selects automatically a value for λ depending on the noise level $\hat{\sigma}\,$

obtain inverse covariance matrix $K = \Sigma^{-1}$ from

$$L_{\beta,\lambda} = \log |K| - \operatorname{tr} S K - \lambda \sum_{i < j} |K_{ij}|$$

where S is sample covariance

obtain inverse covariance matrix $K = \Sigma^{-1}$ from

$$L_{\beta,\lambda} = \log |K| - \operatorname{tr} S K - \lambda \sum_{i < j} |K_{ij}|$$

where S is sample covariance

• screening property $S_0 \subseteq \hat{S}$

obtain inverse covariance matrix $K = \Sigma^{-1}$ from

$$L_{\beta,\lambda} = \log |K| - \operatorname{tr} S K - \lambda \sum_{i < j} |K_{ij}|$$

where S is sample covariance

properties of glasso

- screening property $S_0 \subseteq \hat{S}$
- more restrictive than MB approach

obtain inverse covariance matrix $K = \Sigma^{-1}$ from

$$L_{\beta,\lambda} = \log |K| - \operatorname{tr} S K - \lambda \sum_{i < j} |K_{ij}|$$

where S is sample covariance

properties of glasso

- screening property $S_0 \subseteq \hat{S}$
- more restrictive than MB approach
- more accurate that MB apraoch (slightly)

dependent measures of effect



dependent measures of effect



dependent measures of effect

correct rejections/true edges (true positive rate) $\mathbf{recall} := \frac{|\hat{E} \cap E_0|}{|E_0|}$ correct rejections/rejections (positive predictive value) precision := $\frac{|\hat{E} \cap E_0|}{|\hat{E}|}$



true graph





true graph

MB lasso





true graph			MB lasso	
	precision	recall	rejection rate	density
	0.26	0.98	0.077	0.010



BA graph



true graph

BA graph



true graph



MB lasso

BA graph





true graph

MB lasso

precision	recall	rejection rate	density
0.032	1.00	0.99	0.0625

• prediction squared prediction loss (est – truth)²

$$\mathbb{E}[(\hat{\beta} - \beta_0)^\mathsf{T} X_{\mathsf{new}}^\mathsf{T} X_{\mathsf{new}} (\hat{\beta} - \beta_0) \mid X]$$

• prediction squared prediction loss (est – truth)²

$$\mathbb{E}[(\hat{\beta} - \beta_0)^\mathsf{T} X_{\mathsf{new}}^\mathsf{T} X_{\mathsf{new}} (\hat{\beta} - \beta_0) \mid X]$$

• variable screening estimated variables in $\hat{S} = \{j : \hat{\beta}_j \neq 0\}$ wrt true $S_0 = \{j : \beta_{0j} \neq 0\}$

$$\mathbb{P}[S_0 \subseteq \hat{S}] \to 1$$

• prediction squared prediction loss (est – truth)²

$$\mathbb{E}[(\hat{\beta} - \beta_0)^\mathsf{T} X_{\mathsf{new}}^\mathsf{T} X_{\mathsf{new}} (\hat{\beta} - \beta_0) \mid X]$$

• variable screening estimated variables in $\hat{S} = \{j : \hat{\beta}_j \neq 0\}$ wrt true $S_0 = \{j : \beta_{0j} \neq 0\}$

$$\mathbb{P}[S_0 \subseteq \hat{S}] \to 1$$

allow for false positives, no false negatives

• prediction squared prediction loss (est – truth)²

$$\mathbb{E}[(\hat{\beta} - \beta_0)^\mathsf{T} X_{\mathsf{new}}^\mathsf{T} X_{\mathsf{new}} (\hat{\beta} - \beta_0) \mid X]$$

• variable screening estimated variables in $\hat{S} = \{j : \hat{\beta}_j \neq 0\}$ wrt true $S_0 = \{j : \beta_{0j} \neq 0\}$

$$\mathbb{P}[S_0 \subseteq \hat{S}] \to 1$$

allow for false positives, no false negativesvariable selection

$$\mathbb{P}[S_0 = \hat{S}] \to 1$$

assumptions for variable selection (a) sparsity number $s_0 = |\{j : \beta_{0j} \neq 0\}|$ cannot be large

$$|\{j: \beta_{0j} \neq 0\}| \le \sqrt{\frac{n}{\log(p)}}$$

assumptions for variable selection (a) sparsity number $s_0 = |\{j : \beta_{0j} \neq 0\}|$ cannot be large

$$|\{j: \beta_{0j} \neq 0\}| \le \sqrt{\frac{n}{\log(p)}}$$

(b) consistency tuning λ , no collinearity in true set of variables S_0

assumptions for variable selection (a) sparsity number $s_0 = |\{j : \beta_{0j} \neq 0\}|$ cannot be large

$$|\{j: \beta_{0j} \neq 0\}| \le \sqrt{\frac{n}{\log(p)}}$$

- (b) consistency tuning λ , no collinearity in true set of variables S_0
- (c) beta-min the signal (edge weights) cannot be too small

$$\min_{j\in E} |\beta_{0j}| > c$$



Gaussian graphical models

- stepwise procedures
 problem
- 2
- regularization
- Iasso variants



desparsified lasso

the lasso estimate $\hat{\beta}_L$ is obtained by minimizing

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \sum_{i} |\beta_{i}|$$

desparsified lasso

the lasso estimate $\hat{\beta}_L$ is obtained by minimizing

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \sum_{i} |\beta_{i}|$$

Properties lasso estimate $\hat{\beta}_L$

$$\mathbb{E}(\hat{\beta}_L) = \beta_0 - \mathsf{bias}$$

desparsified lasso

the lasso estimate $\hat{\beta}_L$ is obtained by minimizing

$$L_{\beta,\lambda} = (Y - X\beta)^{\mathsf{T}}(Y - X\beta) + \lambda \sum_{i} |\beta_{i}|$$

Properties lasso estimate $\hat{\beta}_L$

$$\mathbb{E}(\hat{\beta}_L) = \beta_0 - \text{bias}$$

bias = $\hat{\Theta}X^{\mathsf{T}}(Y - X\hat{\beta}_L)/n$
Gaussian graphical models



desparsified (debiased) lasso

$$\hat{\beta}_{dL} = \underbrace{\hat{\beta}_L}_{\text{normal lasso}} + \underbrace{\hat{\Theta}X^{\mathsf{T}}(Y - X\hat{\beta}_L)/n}_{\text{debias part}}$$

 $\hat{\Theta}$ is a 'relaxed' inverse of $X^{\mathsf{T}}X/n$, such that



 $\hat{\Theta}$ is a 'relaxed' inverse of $X^{\mathsf{T}}X/n$, such that

$$\hat{\beta}_{dL} = \underbrace{\hat{\beta}_L}_{\text{normal lasso}} + \underbrace{\hat{\Theta}X^{\mathsf{T}}(Y - X\hat{\beta}_L)/n}_{\text{debias part}} -\Delta$$
$$\Delta = \sqrt{n}(\hat{\Theta}(X^{\mathsf{T}}X/n) - I)(\hat{\beta}_L - \beta) \quad \text{(small)}$$

desparsified lasso

is asymptotically unbiased

 $\hat{\Theta}$ is a 'relaxed' inverse of $X^{\mathsf{T}}X/n$, such that

$$\hat{\beta}_{dL} = \underbrace{\hat{\beta}_L}_{\text{normal lasso}} + \underbrace{\hat{\Theta}X^{\mathsf{T}}(Y - X\hat{\beta}_L)/n}_{\text{debias part}} -\Delta$$
$$\Delta = \sqrt{n}(\hat{\Theta}(X^{\mathsf{T}}X/n) - I)(\hat{\beta}_L - \beta) \quad \text{(small)}$$

desparsified lasso

- is asymptotically unbiased
- has standard, computable variance

 $\hat{\Theta}$ is a 'relaxed' inverse of $X^{\mathsf{T}}X/n$, such that

$$\hat{\beta}_{dL} = \underbrace{\hat{\beta}_L}_{\text{normal lasso}} + \underbrace{\hat{\Theta}X^{\mathsf{T}}(Y - X\hat{\beta}_L)/n}_{\text{debias part}} -\Delta$$
$$\Delta = \sqrt{n}(\hat{\Theta}(X^{\mathsf{T}}X/n) - I)(\hat{\beta}_L - \beta) \quad \text{(small)}$$

desparsified lasso

- is asymptotically unbiased
- has standard, computable variance
- is amenable to standard analysis

nodewise desparsified lasso



parameter index

tuning parameter λ estimated by 10-fold cv

estimate true value

• 0

simulations

generate data from different networks ER or BA

simulations

- generate data from different networks ER or BA
- estimate using desparsified lasso with MB or GL

simulations

- generate data from different networks ER or BA
- estimate using desparsified lasso with MB or GL

	ER		BA	
	MB	GL	MB	GL
sparse			$\begin{array}{c} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 &$	$\begin{array}{c} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 &$
not sparse				

correct rejections/true edges (true positive rate)

$$\operatorname{recall} := \frac{|\hat{E} \cap E_0|}{|E_0|}$$

correct rejections/true edges (true positive rate)

$$\operatorname{recall} := \frac{|\hat{E} \cap E_0|}{|E_0|}$$

correct rejections/rejections (positive predictive value)

precision :=
$$\frac{|\hat{E} \cap E_0|}{|\hat{E}|}$$

correct rejections/true edges (true positive rate)

$$\operatorname{recall} := \frac{|\hat{E} \cap E_0|}{|E_0|}$$

correct rejections/rejections (positive predictive value)

precision :=
$$\frac{|\hat{E} \cap E_0|}{|\hat{E}|}$$

probability of CI containing true parameter γ_{0ij}

$$\text{coverage} := \frac{1}{|V|(|V|-1)} \sum_{i \neq j \in V} \mathbb{P}[\beta_{0ij} \in \mathsf{Cl}_{ij}]$$

32/38



true graph



true graph



desparsified lasso









true graph



graphical lasso

6

12 (13)

60

59

(5) (2) (5) (5)

35 / 36 / 37

19) (20)

50

(18)

2

62 63

(22) (23

57 68

(49)

(4) (4) (43)

(33)

25

17

9

(1)

true graph

64

56

48

40

32

(24)

(16)

(8)

(47

15







recall and precision for sparse graphs



-MB - SL

– dL – GL

BA

recall and precision for sparse graphs





-MB - SL





– dL – GL

conclusions variable screening



ER



- dL - GL

ER



-MB - SL

- dL - GL







-MB - SL





ER



BA

• coverage is close to nominal level (95%)

conclusions variable selection CI

