

# Review of sparse reduced-rank regression model

报告人: 尚尔淦

数学科学学院 概率统计系

2021年10月29日





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- 7 wavelet-based sparse reduce-rank regression





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- 7 wavelet-based sparse reduce-rank regression

## definition of graph



A graph G=(V,E), V: set of vertices, E: set of edges; XY is adiacent if there's an edge joining them: $X \sim Y$ Suppose V represents a set of random variables having jointly distribution P

The graph gives a visual way of understanding the joint distribution of the entire set of random variables. The absence of an edge between 2 vertices has a special meaning: the corresponding random variables are conditionally independent given other variables.

We assume that the observation has a multivariate Gaussian distribution with mean  $\mu$  and covariance matrix  $\Sigma$ . Then all conditional distributions are also Gaussian.  $\Sigma^{-1}$  contains partial covariance: covariance between i and i conditioned on all other variables

### stability approach



Our major work is to determine which features of the system are conditionally independent i.e. Estimating the inverse covariance matrix  $\Sigma^{-1}$ 

We can put forward a regularization parameter  $\lambda$  that controls the sparsity of the graph; A new approach to model selection based on **model stability**.

We start with a large regularization which corresponds to an empty and hence highly stable graph. We gradually reduce the amount of regularization until there's a small but acceptable amount of variability of the graph across subsamples.



 $X = (X(x), \dots, X(p))'$  is the random vector with distribution P, G=(V,E) with vertices  $V=\{X(1),\cdots,X(p)\}$ . We use E to denote the adjacency matrix and edges.

Our aim is to infer E from i.i.d observed data  $X_1, \dots, X_n$ where  $X_i = (X_i(1), \cdots, X_i(p))'$ 

Suppose P is Gaussian with mean vector  $\mu$  and covariance matrix  $\Sigma$ . let  $\Omega = \Sigma^{-1}$ 



 $X = (X(x), \dots, X(p))'$  is the random vector with distribution P, G=(V,E) with vertices  $V = \{X(1), \dots, X(p)\}$ . We use E to denote the adjacency matrix and edges.

Our aim is to infer E from i.i.d observed data  $X_1, \dots, X_n$ where  $X_i = (X_i(1), \cdots, X_i(p))'$ 

Suppose P is Gaussian with mean vector  $\mu$  and covariance matrix  $\Sigma$ . let  $\Omega = \Sigma^{-1}$ 

 $\Omega_{ik} = 0 \iff \text{no edge corresponding to } X(i) \text{ and } X(k)$ 



 $X = (X(x), \dots, X(p))'$  is the random vector with distribution P, G=(V,E) with vertices  $V = \{X(1), \dots, X(p)\}$ . We use E to denote the adjacency matrix and edges.

Our aim is to infer E from i.i.d observed data  $X_1, \dots, X_n$ where  $X_i = (X_i(1), \cdots, X_i(p))'$ 

Suppose P is Gaussian with mean vector  $\mu$  and covariance matrix  $\Sigma$ . let  $\Omega = \Sigma^{-1}$ 

$$\Omega_{jk} = 0 \iff$$
 no edge corresponding to  $X(j)$  and  $X(k)$ 

We come to estimate the sparsity pattern of  $\Omega$ , if ignoring the constants, the log-likelihood:

$$\ell(\Omega) = \log |\Omega| - \operatorname{trace}(\hat{\Sigma}\Omega)$$



 $\dot{\Sigma}$ :the sample covariance matrix

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^{N} (x_i - \bar{x})(x_i - \bar{x})'$$

 $\hat{\Omega}(\lambda) = \arg\min_{\Omega>0} (-\ell(\Omega) + \lambda ||\Omega||_1)$  where  $||\Omega||_1 = \sum_{i,k} |\Omega_{ik}|$ is the  $\ell_1$  norm

Therefore we get the estimator

$$\hat{G}(\lambda) = (V, \hat{E}(\lambda))$$

## Regularization selection



Obviously, larger values of  $\lambda$  tend to yield sparse graphs, we define  $\Lambda = \frac{1}{\lambda}$  so that small  $\Lambda$  corresponds to sparse graph. We get Grid of regularization parameters  $G_n = \{\lambda_1 \cdots, \Lambda_K\}$ , to choose one  $\hat{\Lambda} \in G_n$  such that the true graph E is contained in  $\hat{E}(\hat{\lambda})$  with high probability.

Let  $d(\lambda)$  denote the degree of freedom of the corresponding Gaussian model, we have these existing criterion

$$\mathit{AIC}: \quad \hat{\Lambda} = \arg\min_{\lambda \in \mathit{G}_{\mathit{n}}} (-2\ell(\Omega(\Lambda)) + 2\mathit{d}(\Lambda))$$

$$\mathit{BIC}: \quad \hat{\Lambda} = \arg\min_{\lambda \in \mathit{G}_{\mathit{n}}} (-2\ell(\Omega(\Lambda)) + \mathit{d}(\Lambda) \cdot \mathit{logn})$$

A common practice is to calculate  $d(\Lambda) = m(\Lambda)(m(\Lambda) - 1)/2 + p$ , where  $m(\Lambda)$  denotes the number of nonzero elements of  $\hat{\Omega}(\Lambda)$ 



When  $\Lambda = 0$ , the graph is empty. As we increase  $\Lambda$ , the variability of graph increases and the stability decreases. Now we put forward the concrete rule for choosing  $\Lambda$ 



When  $\Lambda = 0$ , the graph is empty. As we increase  $\Lambda$ , the variability of graph increases and the stability decreases. Now we put forward the concrete rule for choosing  $\Lambda$ Let b = b(n) be the parameter such that 1 < b(n) < n, we draw N random subsamples  $S_1, \dots, S_N$  from  $X_1, \dots, X_n$  each of size b. We choose a large number N of subsamples at random. Each subsample is drawn without replacement



When  $\Lambda = 0$ , the graph is empty. As we increase  $\Lambda$ , the variability of graph increases and the stability decreases. Now we put forward the concrete rule for choosing  $\Lambda$ Let b = b(n) be the parameter such that 1 < b(n) < n, we draw N random subsamples  $S_1, \dots, S_N$  from  $X_1, \dots, X_n$  each of size b. We choose a large number N of subsamples at random. Each subsample is drawn without replacement For each  $\Lambda \in G_n$ , we construct a graph using glasso for each subsample. This results in N estimated edge matrices  $\hat{E}_1^b(\Lambda), \cdots, \hat{E}_N^b(\Lambda)$ 



When  $\Lambda = 0$ , the graph is empty. As we increase  $\Lambda$ , the variability of graph increases and the stability decreases. Now we put forward the concrete rule for choosing  $\Lambda$ Let b = b(n) be the parameter such that 1 < b(n) < n, we draw N random subsamples  $S_1, \dots, S_N$  from  $X_1, \dots, X_n$  each of size b. We choose a large number N of subsamples at random. Each subsample is drawn without replacement For each  $\Lambda \in G_n$ , we construct a graph using glasso for each subsample. This results in N estimated edge matrices  $\hat{E}_1^b(\Lambda), \cdots, \hat{E}_N^b(\Lambda)$ 

Now we focus on one edge (s,t) and one value of  $\Lambda$ . Let  $\psi^{\lambda}(\cdot)$ denotes the glasso algorithm with  $\Lambda$ 



For any subsample  $S_j$ , let  $\psi_{st}^{\Lambda}(S_j)=1$  if the algorithm puts an edge between s and t, otherwise  $\psi_{st}^{\Lambda}(S_i) = 0$ 



For any subsample  $S_i$ , let  $\psi_{st}^{\Lambda}(S_i) = 1$  if the algorithm puts an edge between s and t, otherwise  $\psi_{st}^{\Lambda}(S_i) = 0$ Define

$$\theta_{st}^b(\Lambda) = P(\psi_{st}^{\Lambda}(X_1, \cdots, X_b) = 1)$$



For any subsample  $S_i$ , let  $\psi_{st}^{\Lambda}(S_i) = 1$  if the algorithm puts an edge between s and t, otherwise  $\psi_{st}^{\Lambda}(S_i) = 0$ Define

$$\theta_{st}^b(\Lambda) = P(\psi_{st}^{\Lambda}(X_1, \cdots, X_b) = 1)$$

To estimate  $\theta_{st}^b(\Lambda)$ , define  $\hat{\theta}_{st}^b(\Lambda) = \frac{1}{N} \sum_{i=1}^N \psi_{st}^{\Lambda}(S_i)$ 



For any subsample  $S_i$ , let  $\psi_{st}^{\Lambda}(S_i) = 1$  if the algorithm puts an edge between s and t, otherwise  $\psi_{ct}^{\Lambda}(S_i) = 0$ Define

$$\theta_{st}^b(\Lambda) = P(\psi_{st}^{\Lambda}(X_1, \cdots, X_b) = 1)$$

To estimate  $\theta_{st}^b(\Lambda)$ , define  $\hat{\theta}_{st}^b(\Lambda) = \frac{1}{N} \sum_{i=1}^N \psi_{st}^{\Lambda}(S_i)$ Now define the parameter

$$\xi_{\rm st}^b(\Lambda) = 2\theta_{\rm st}^b(\Lambda)(1 - \theta_{\rm st}^b(\Lambda))$$



For any subsample  $S_i$ , let  $\psi_{st}^{\Lambda}(S_i) = 1$  if the algorithm puts an edge between s and t, otherwise  $\psi_{st}^{\Lambda}(S_i) = 0$ Define

$$\theta_{st}^b(\Lambda) = P(\psi_{st}^{\Lambda}(X_1, \cdots, X_b) = 1)$$

To estimate  $\theta_{st}^b(\Lambda)$ , define  $\hat{\theta}_{st}^b(\Lambda) = \frac{1}{N} \sum_{i=1}^N \psi_{st}^{\Lambda}(S_i)$ Now define the parameter

$$\xi_{st}^b(\Lambda) = 2\theta_{st}^b(\Lambda)(1 - \theta_{st}^b(\Lambda))$$

let  $\hat{\mathcal{E}}_{\mathfrak{s}}^{b}(\Lambda) = 2\hat{\theta}_{\mathfrak{s}}^{b}(\Lambda)(1-\hat{\theta}_{\mathfrak{s}}^{b}(\Lambda))$  be its estimate. We can regard  $\xi_{st}^b(\Lambda)$  as being twice the variance of the Bernoulli indicator of the edge (s,t) and as a measure of instability of the edge across subsample with  $0 \leq \xi_{st}^b(\Lambda) \leq \frac{1}{2}$ 



Define the total instability by averaging over all edges

$$\hat{D}_b(\Lambda) = \frac{\sum_{s < t} \hat{\xi_{st}^b}}{C_p^2}$$



Define the total instability by averaging over all edges

$$\hat{D}_b(\Lambda) = \frac{\sum_{s < t} \hat{\xi_{st}^b}}{C_p^2}$$

Then define

$$\bar{D}_b(\Lambda) = \sup_{0 \leq t \leq \Lambda} \hat{D}_b(t) \quad \hat{\Lambda}_s = \sup\{\Lambda : \bar{D}_b(\Lambda) \leq \beta\}$$

 $\beta$  is the threshold





- 1 Stability approach to regularization selection
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- wavelet-based sparse reduce-rank regression





- 1 Stability approach to regularization selection
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- wavelet-based sparse reduce-rank regression

#### SRRR



The dimension reduction aspect of multivariate regression is taken care of by reduced-rank regression(RRR)

The variable selction aspect is addressed by adding a penalty to the least squares fitting criterion to enforce the sparsity of reduced-rank coefficient matrix.



The dimension reduction aspect of multivariate regression is taken care of by reduced-rank regression(RRR)

The variable selction aspect is addressed by adding a penalty to the least squares fitting criterion to enforce the sparsity of reduced-rank coefficient matrix

The model is

$$Y = XC + E$$

with n observations. Taking advantage of possible interrelationships between response variables is to impose a constraint on the rank of C:rank(C) = r < min(p, q)Then C = BA', where  $B : p \times r, A : q \times r$ , and

$$Y = (XB)A' + E$$



XB is of reduced dimension with only r components, which can be interpreted as unobservable latent factors. By solving the optimization problem

$$\min_{C: rank(C)=r} ||Y - XC||^2$$

Denote  $S_{xx} = \frac{1}{n}X'X$ ,  $S_{xy} = \frac{1}{n}X'Y$ ,  $S_{yx} = \frac{1}{n}Y'X$ , we have the solution

$$\hat{A}^{(r)} = V \quad \hat{B}^{(r)} = S_{xx}^{-1} S_{xy} V$$

where  $V = (v_1, \dots, v_r)$  and  $v_i$  is the eigen vector of  $S_{vx}S_{vx}^{-1}S_{xv}$ The solution satisfies that:  $A'A = I_r$ ,  $B'S_{xx}B$  being diagnonal.

## SRRR through penalized least square 中国神经技术大学

Exclude the redundant predictors when some predictor variables are not useful for prediction ←⇒ set as zero an entire row of B.

## SRRR through penalized least square 中国神经技术大学

Exclude the redundant predictors when some predictor variables are not useful for prediction ←⇒ set as zero an entire row of B. We Consider the following optimization problem

$$\min_{A,B} ||Y - XBA'||^2 + \sum_{i=1}^p \lambda_i ||B^i||$$
 such that  $A'A = I$ 

where  $B^i$  denotes the ith row of B. This is a penalized regression with a grouped lasso penalty.

## SRRR through penalized least square 中国神学技术大学

Exclude the redundant predictors when some predictor variables are not useful for prediction ←⇒ set as zero an entire row of B. We Consider the following optimization problem

$$\min_{A,B} ||Y - XBA'||^2 + \sum_{i=1}^p \lambda_i ||B^i||$$
 such that  $A'A = I$ 

where  $B^i$  denotes the ith row of B. This is a penalized regression with a grouped lasso penalty.

#### Lemma 2.1

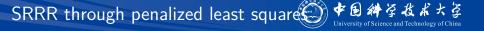
The solution to the optimization problem is unique up to an orthogonal matrix.

### proof of lemma 1



### Proof 2.1

Let  $(\hat{A}, \hat{B})$  is a solution and Q is an orthogonal matrix. Let  $\tilde{A} = \hat{A}Q, \tilde{B} = \hat{B}Q$ , then  $\tilde{B}\tilde{A}' = \hat{B}\hat{A}'$  and  $||\tilde{B}^i|| = ||\hat{B}^i||$  via QQ' = I. As a result,  $(\tilde{A}, \tilde{B})$  is also a solution Moreover, if  $(\tilde{A}, \tilde{B})$  and  $(\hat{A}, \hat{B})$  are both the solution, considering  $rank(\hat{A}) = rank(\hat{A}) = r$ , then there's a non-singular matrix Q of  $r \times r$  such that A = AQ, we reach the conclusion that Q is orthogonal because  $I_r = \tilde{A}'\tilde{A} = Q'\hat{A}'\hat{A}Q = Q'Q$ Finally,  $||B^i|| = ||B^iQ'||$  and ||Y - XBA'|| = ||Y - XBQ'(AQ')'||, we know that  $(\hat{A}, \hat{B})$  is the solution. Then  $\hat{B} = \tilde{B}Q'$ 



### Definition 2.1

If the entire row j of B is zero, then the predictor variable  $X_i$  is called a nonactive variable, otherwise it is called an active variable.



### Definition 2.1

If the entire row j of B is zero, then the predictor variable  $X_i$  is called a nonactive variable, otherwise it is called an active variable.

### Lemma 2.2

The set of active variables obtained by the optimization problem is uniquely determined.





- 1 Stability approach to regularization selection
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- wavelet-based sparse reduce-rank regression

## iterative optimization



For fixed B, the optimization problem is equivalent with

$$\min_{A} ||Y - XBA'||$$
 such that  $A'A = I$ 

which is an orthogobal Procrustes Problem(Gower&Dijksterhuis 2004), and the solution is  $\hat{A} = UV$ . where U, V is the SVD of YXB, i.e. YXB = UDV

## iterative optimization



For fixed B, the optimization problem is equivalent with

$$\min_{A} ||Y - XBA'||$$
 such that  $A'A = I$ 

which is an orthogobal Procrustes Problem(Gower&Dijksterhuis 2004), and the solution is  $\hat{A} = UV$ , where U, V is the SVD of YXB, i.e. YXB = UDVThen with the fixed A, considering the column vector of A is orthogonal, we can let  $(A, A^{\perp})$  be an orthogonal matrix. Then

$$||Y - XBA'||^2 = ||(Y - XBA')(A, A^{\perp})||^2 = ||YA - XB||^2 + ||YA^{\perp}||^2$$

## iterative optimization



For fixed B, the optimization problem is equivalent with

$$\min_{A} ||Y - XBA'||$$
 such that  $A'A = I$ 

which is an orthogobal Procrustes Problem(Gower&Dijksterhuis 2004), and the solution is  $\hat{A} = UV$ , where U, V is the SVD of YXB, i.e. YXB = UDVThen with the fixed A, considering the column vector of A is orthogonal, we can let  $(A, A^{\perp})$  be an orthogonal matrix. Then

$$|| \mathbf{Y} - \mathbf{X} \mathbf{B} \mathbf{A}' ||^2 = || (\mathbf{Y} - \mathbf{X} \mathbf{B} \mathbf{A}') (\mathbf{A}, \mathbf{A}^{\perp}) ||^2 = || \mathbf{Y} \mathbf{A} - \mathbf{X} \mathbf{B} ||^2 + || \mathbf{Y} \mathbf{A}^{\perp} ||^2$$

The optimization problem is equivalent with, for A fixed,  $\min_{B} ||YA - XB||^2 + \sum_{i=1}^{p} \lambda_i ||B^i||$ 

### subgradient method



The subgradient quations about  $B^{\ell}$  is defined as follows

$$2X'_{\ell}(XB - YA) + \lambda_{\ell}s_{\ell} \quad \forall \ell = 1, 2, \cdots, p$$

where 
$$s_{\ell} = \frac{B^{\ell}}{||B^{\ell}||}, ||B^{\ell}|| \neq 0$$

### subgradient method



The subgradient quations about  $B^{\ell}$  is defined as follows

$$2X'_{\ell}(XB - YA) + \lambda_{\ell}s_{\ell} \quad \forall \ell = 1, 2, \cdots, p$$

where  $s_{\ell} = \frac{B^{\ell}}{||B^{\ell}||}, ||B^{\ell}|| \neq 0$ If  $||B^{\ell}|| = 0$ , the equation becomes

$$2X_{\ell}(\sum_{k\neq\ell}^{p}X_{k}B^{k}-YA)+\lambda_{\ell}s_{\ell}=0$$

, then 
$$s_\ell=-rac{2}{\lambda_\ell}X'_\ell(\sum_{k
eq\ell}^\rho X_kB^k-Y\!A):=rac{2}{\lambda_\ell}X'_\ell R_\ell$$

### subgradient method



resulting in

$$B^{\ell} = (X_{\ell}X_{\ell} + \frac{\lambda_{\ell}}{2||B^{\ell}||})^{-1}X_{\ell}R_{\ell}$$

. Noting that the RHS involves  $||B^{\ell}||$ , we let  $c = ||B^{\ell}||$  and solve the equation and plug in  $c=\frac{||X'_{\ell}R_{\ell}||-\frac{1}{2}\lambda_{\ell}}{||X_{\ell}||^2}$  we get the final solution

$$B^{\ell} = \frac{1}{X_{\ell}' X_{\ell}} (1 - \frac{\lambda_{\ell}}{2||X_{\ell}' R_{\ell}|})_{+} X_{\ell}' R_{\ell}$$

which is a vector version of the soft-thresholding rule.

#### variational method



Noting that the truth that  $\min_{C} \frac{1}{2}(cx^2 + \frac{1}{c}) = |x|$ , then the problem is equivalent

$$f = ||YA - XB||^2 + \sum_{i=1}^p \frac{\lambda_i}{2} (\mu_i ||B^i||^2 + \frac{1}{\mu_i}) \quad \textit{jointly} \quad \textit{over} \quad B \quad \textit{and} \quad \mu_i$$

#### variational method



Noting that the truth that  $\min_{C} \frac{1}{2}(cx^2 + \frac{1}{c}) = |x|$ , then the problem is equivalent

$$f = ||YA - XB||^2 + \sum_{i=1}^{p} \frac{\lambda_i}{2} (\mu_i ||B^i||^2 + \frac{1}{\mu_i})$$
 jointly over  $B$  and  $\mu_i$ 

For fixed B, the solution of  $\mu_i$  is that  $\mu_i = \frac{1}{||B^i||}, i = 1, \dots, p$ For fixed  $\mu_i$ , we have  $\frac{\partial f}{\partial B^i} = -2X_i(YA - XB) + \lambda_i\mu_iB^i = 0$ 

#### variational method



Noting that the truth that  $\min_{C} \frac{1}{2}(cx^2 + \frac{1}{c}) = |x|$ , then the problem is equivalent

$$f = ||YA - XB||^2 + \sum_{i=1}^p \frac{\lambda_i}{2} (\mu_i ||B^i||^2 + \frac{1}{\mu_i})$$
 jointly over  $B$  and  $\mu_i$ 

For fixed B, the solution of  $\mu_i$  is that  $\mu_i = \frac{1}{||B^i||}, i = 1, \dots, p$ For fixed  $\mu_i$ , we have  $\frac{\partial f}{\partial B^i} = -2X_i(YA - XB) + \lambda_i\mu_iB^i = 0$ As a result, the final solution is

$$B = \{X'X + \frac{1}{2}diag(\lambda_i\mu_i, \cdots, \lambda_p\mu_p)\}^{-1}X'YA$$



We give the following algorithm of iteration Input:  $X, Y, \lambda$ Output: A, B while (A, B) are not convergent for fixed B, we get the solution of A by SVD while(B is not convergent) for every  $\ell$  solve  $B^{\ell}$  and check whether B is convergent





- 1 Stability approach to regularization selection
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- wavelet-based sparse reduce-rank regression

### Assumptions



### Assuming p, q are fixed with n going to infinity

- 1 :There's a positive definite matrix  $\Sigma$  such that  $\frac{XX'}{n} \to \Sigma(n \to \infty)$
- 2 :The first  $p_0$  variables are important and the rest are irrelevant i.e.  $||C_*^i|| > 0 (i \le p_0)$  and  $||C_*^i|| = 0 (i > p_0)$ where  $C_{*}^{i}$  is the ith row of  $C^{*}$ , which is the rank-r coefficient matrix used to generate data in the model



# Theorem 2.1 (consistency of parameter estimation) suppose $\frac{\lambda_i}{\sqrt{n}} = \frac{\lambda_{n,i}}{\sqrt{n}} \to 0, \forall i \leq p_0 Then$

- ▶ There is a local minimizer  $\hat{C}$  that is  $\sqrt{n}$ -consistent in estimating C<sub>\*</sub>
- $ightharpoonup \hat{C} = \hat{U}\hat{D}\hat{V}$  is SVD,  $\hat{C} = \hat{U}\hat{D}\hat{V}$  are  $\sqrt{n}$ -consistent in estimating  $U_*, D_*, V_*$



# Theorem 2.1 (consistency of parameter estimation) suppose $\frac{\lambda_i}{\sqrt{n}} = \frac{\lambda_{n,i}}{\sqrt{n}} \to 0, \forall i \leq p_0 Then$

- ▶ There is a local minimizer  $\hat{C}$  that is  $\sqrt{n}$ -consistent in estimating C<sub>\*</sub>
- $ightharpoonup \hat{C} = \hat{U}\hat{D}\hat{V}$  is SVD,  $\hat{C} = \hat{U}\hat{D}\hat{V}$  are  $\sqrt{n}$ -consistent in estimating  $U_*, D_*, V_*$

## Theorem 2.2 (concistency of variable selection)

if 
$$\frac{\lambda_{n,i}}{\sqrt{n}} \to 0$$
 for  $i \le p_0$  and  $\frac{\lambda_{n,i}}{\sqrt{n}} \to \infty (i > p_0)$ , then  $P(\hat{C}^i = 0) = P(\hat{U}^i = 0) \to 1$   $i > p_0$ 





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- 7 wavelet-based sparse reduce-rank regression



$$Y = XC + E$$
  $C \in \mathbb{R}^{p \times q}$   $rankC \leq r$ 

Using the polar decomposition, we have  $C = U\tilde{V}$ , resulting in Y = XUV + E



$$Y = XC + E$$
  $C \in \mathbb{R}^{p \times q}$   $rankC \leq r$ 

Using the polar decomposition, we have  $C = U\tilde{V}$ , resulting in  $Y = XU\tilde{V} + E$ 

This equation is related to a factor analysis model: XU can be regarded as a common factor matrix and  $\tilde{V}$  can be regarded as a loading matrix. Furthermore, if we assume

$$\mathbb{E}[x_i] = 0$$
,  $cov[x_i] = \Gamma_i$ , then  $cov[U'X] = U'\Gamma U = I_r$ 



$$Y = XC + E$$
  $C \in \mathbb{R}^{p \times q}$   $rankC < r$ 

Using the polar decomposition, we have  $C = U\tilde{V}$ , resulting in  $Y = XU\tilde{V} + E$ 

This equation is related to a factor analysis model: XU can be regarded as a common factor matrix and  $\tilde{V}$  can be regarded as a loading matrix. Furthermore, if we assume  $\mathbb{E}[x_i] = 0$ ,  $cov[x_i] = \Gamma_i$ , then  $cov[U'X] = U'\Gamma U = I_r$ 

 $\mathbb{E}[X_i] = 0$ ,  $cov[X_i] = 1$ , then  $cov[UX] = U1U = I_r$ We have the SFAR model(sequential co-sparse factor regression):

$$Y = XUDV$$
 such that  $U'\Gamma U = I_r$   $V'V = I_r$ 

with the coefficient matrix is C = UDV'



We have the optimization problem

$$\min_{U,D,V} \frac{1}{2} ||Y - XUDV||^2 + \lambda_1 \sum_{i=1}^p \sum_{j=1}^r w_{ij}^{(u)} |u_{ij}| + \lambda_2 \sum_{i=1}^q \sum_{j=1}^r w_{ij}^{(v)} |v_{ij}|$$

such that 
$$U'\frac{X'X}{n}U=I_r$$
  $VV=I_r$ 

where  $w_{ii}^{(u)}$ ,  $w_{ii}^{(v)}$  is called adaptive weights with positive values



The minimization problem for the kth latent factor is given by

$$\min_{dk, \mathbf{u_k}, \mathbf{v_k}} \frac{1}{2} || Y_k - d_k X \mathbf{u_k} \mathbf{v_k}^T ||^2 + \sum_{i=1}^p w_{ki}^{(u)} |u_{ki}| + \sum_{i=1}^q w_{ki}^{(w)} |v_{ki}|$$

such that 
$$d_k \geq 0$$
  $u_k^T X' X u_k = n$   $v_k v_k = 1$ 

where  $Y_k = Y - \sum_{i=1}^{k-1} d_i X u_i v_i^T (SeCURE algorithm)$ Mishra(2017))

# via manifold optimization



#### Definition 3.1

- ►  $St(r,q) := \{V \in \mathbb{R}^{q \times r} | VV = I_r\} (q \ge r)$  called Stiefel manifold
- ▶  $GSt(r, p) := \{U \in \mathbb{R}^{p \times r} | U'GU = I_r\} (p \ge r), G \in \mathbb{R}^{p \times p} \text{ is definite positive, called generalized Stiefel manifold.}$

In this paper, we use  $G = \frac{X'X}{n}$  The optimization is equivalent with

$$\min_{U \in GSt(r,p), D \in \mathbb{R}^{r \times r}, V \in St(r,q)} \frac{1}{2} ||Y - XUDV||^2 + n\lambda_1 \sum_{i=1}^{p} \sum_{j=1}^{r} w_{ij}^{(u)} |u_{ij}| +$$

$$n\lambda_2 \sum_{i=1}^{q} \sum_{j=1}^{r} w_{ij}^{(v)} |v_{ij}|$$

We propose the following minimization problem

$$\min_{\textit{U} \in \textit{GSt}(\textit{r},\textit{p}),\textit{D} \in \mathbb{R}^{\textit{r} \times \textit{r}},\textit{V} \in \textit{St}(\textit{r},\textit{q})} \frac{1}{2} ||\textit{Y} - \textit{XUDV}||^2 + \textit{n} \lambda_1 \sum_{i=1}^{\textit{p}} \sum_{j=1}^{\textit{r}} \textit{w}_{ij}^{(\textit{u})} |\textit{u}_{ij}| +$$

$$\textit{n}\alpha\lambda_2\sum_{i=1}^{q}\sum_{j=1}^{r}\textit{w}_{ij}^{(\textit{v})}|\textit{v}_{ij}|+\textit{n}\sqrt{\textit{q}}(1-\alpha)\lambda_2\sum_{i=1}^{r}\textit{w}_{i}^{(\textit{d})}\textit{I}(\textit{\textbf{v}}_{\textit{\textbf{i}}}\neq 0)$$

 $w_i^{(d)}$  is an adaptive weight with a positive value proposed by Zou(2006) and the group selection in the fourth term plays the role of the rank selection of the coefficient matrix C.





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- joint variable and rank selection for parsimonious
   Introduction procedure and property numerical performance
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- wavelet-based sparse reduce-rank regression





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- Introduction procedure and property numerical performance
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- wavelet-based sparse reduce-rank regression



$$Y = XA + E$$
  $r = rankA$   $q = rankX$ 

Let J denotes the index set of the nonzero rows of A. Only r(n+|J|-r) free parameters need to be estimated by SVD, where |J| = #J

We can reduce X of rank q to an  $m \times q$  matrix with q independent columns and always assume that  $|J| \leq q$ . Using penalized least squares methods, removing predictor  $X_i$  from the model is equivalent with setting the jth row in A to zero.





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- Introduction procedure and property numerical performance
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- wavelet-based sparse reduce-rank regression



JRRS is short for the single-stage joint rank and row selection estimator. We can modify the rank selection criterion (RSC) as

$$\hat{B} = \arg\min_{B} \{||Y - XB||^2 + pen(B)\}$$
 (1)

where

$$pen(B)=c\sigma^2 r(B)(2n+\log(2e)|J(B)|+|J(B)|\log\frac{ep}{|J(B)|})$$
 ,and we call the equation (1) JRRS1

JRRS is short for the single-stage joint rank and row selection estimator. We can modify the rank selection criterion (RSC) as

$$\hat{B} = \arg\min_{B} \{||Y - XB||^2 + pen(B)\}$$
 (1)

where

$$pen(B) = c\sigma^2 r(B)(2n + log(2e)|J(B)| + |J(B)|log\frac{ep}{|J(B)|})$$
, and we call the equation (1) JRRS1

#### Theorem 4.1

The single-stage JRRS estimator  $\hat{B}$  using pen(B) with  $c = 12^3$ satisfies

$$\mathbb{E}[||\textit{XA} - \textit{X}\hat{\textit{B}}||^2] \leq 10||\textit{XA} - \textit{XB}||^2 + 8\textit{pen}(\textit{B}) + 768\textit{n}\sigma^2\textit{e}^{-\frac{\textit{n}}{2}} \quad \forall \textit{r}(\textit{B}) \geq 1$$

In particular if r(A) >> 1

$$\mathbb{E}[||XA - X\hat{B}||^2] \lesssim \sigma^2 r(A)(n + |J(A)|\log \frac{p}{|J(A)|})$$



#### Theorem 4.2

For any collection of (random) nonzero matrices  $B_1, B_2, \cdots$ , the single-stage JRRS estimator

$$\tilde{B} = arg \min_{B_j}(||Y - XB_j||^2 + pen(B_j))$$
 with  $c = 12$ , satisfies

$$\mathbb{E}[||XA - X\tilde{B}||^2] \le \inf\{10\mathbb{E}[||XA - XB_j||^2] + 8\mathbb{E}[pen(B_j)]\} + 768n\sigma^2 e^{-\frac{n}{2}}$$

### Rank-constrained predictor model



We propose a convex relaxation for the pen(B) with  $||B||_{2,1} = \sum_{i=1}^{p} ||b_i||_2$  rows  $b_i$  of B

$$\hat{B}_k = \arg\min_{r(B) \le k} \{ ||Y - XB||^2 + 2\lambda ||B||_{2,1} \}$$
 (2)

The expression (2) is called rank constrained group lasso(RCGL)(here we have two parameters nedded to be tuned k and  $\lambda$ )

## Rank-constrained predictor model



We propose a convex relaxation for the pen(B) with  $||B||_{2,1} = \sum_{i=1}^{p} ||b_i||_2$  rows  $b_i$  of B

$$\hat{B}_k = \arg\min_{r(B) \le k} \{ ||Y - XB||^2 + 2\lambda ||B||_{2,1} \}$$
 (2)

The expression (2) is called **rank constrained group** lasso(RCGL)(here we have two parameters nedded to be tuned k and  $\lambda$ )

**Assumption A**:We say  $\Sigma \in \mathbb{R}^{p \times p}$  satisfies condition  $A(I, \delta_I)$ for an index set  $I \subset \{1, \dots, p\}$  and  $\delta_I > 0 \iff tr(M'\Sigma M) \ge \delta_I \sum_{i \in I} ||m_i||_2^2 \text{ for all } p \times n \text{ matrices}$  $M(\text{with rows } m_i)$  satisfying  $\sum_{i \in I} ||m_i||_2 \le 2 \sum_{i \in I^c} ||m_i||_2 (\text{which } m_i)$ is  $\ell_2$  norm) In our paper,  $\Sigma = \frac{X'X}{m}$ , and we have the following remark.

### Rank-constrained predictor model



#### Note

- the constant 2 can be replaced by any constant>1
- A sufficient condition is: there exists a diagonal matrix D with  $D_{ii} = \delta_I$  for  $j \in I$  and  $D_{ii}$  otherwise such that  $\Sigma - D > 0$

Let  $\lambda_1(\Sigma)$  denotes the largest eigenvalue of  $\Sigma$  and set the parameter  $\lambda = c\sigma \sqrt{\lambda_1(\Sigma) kmlog(ep)} (c > 0)$ 



#### Theorem 4.3

Let  $B_k$  be the global minimizer corresponding  $\lambda$  above with c large enough. Then

$$\mathbb{E}[||X\hat{B}_{K}-XA||^{2}] \lesssim ||XB-XA||^{2} + k\sigma^{2}(n + (1 + \frac{\lambda_{1}(\Sigma)}{\delta_{J(B)}})|J(B)|\log(p))$$

 $\forall B \in \mathbb{R}^{p \times n}$  with  $1 \leq r(B) \leq k$  provided  $\Sigma$  satisfies Assumption  $A(J(B), \delta_{J(B)})$ 

## method 1(RSC $\rightarrow$ RCGL)



- ▶ Use RSC to select  $k = \hat{r}$  (review the last report) as the number of singular values of PY that exceed  $\sigma(\sqrt{2n} + \sqrt{2q}), P = X(X'X)^{-1}X'$
- $\triangleright$  Compute the rank constrained GLASSO estimator  $B_k$ with  $k = \hat{r}$  to obtain the final estimator  $\hat{B}^{(1)} = \hat{Br}$

## method 1(RSC $\rightarrow$ RCGL)



- ▶ Use RSC to select  $k = \hat{r}$  (review the last report) as the number of singular values of PY that exceed  $\sigma(\sqrt{2n} + \sqrt{2q}), P = X(X'X)^{-1}X'$
- $\triangleright$  Compute the rank constrained GLASSO estimator  $B_k$ with  $k = \hat{r}$  to obtain the final estimator  $\hat{B}^{(1)} = \hat{B}r$

This two-step estimator adapts to both rank and row sparsity under 2 additional mild restrictions

- Ass  $C_1$ :  $d_r(XA) > 2\sqrt{2}\sigma(\sqrt{n} + \sqrt{q})$  (RSC)
- Ass  $C_2$ :  $log(||XA||_F) \le (\sqrt{2}-1)^2 \frac{n+q}{4}$

## method 1(RSC $\rightarrow$ RCGL)



- ▶ Use RSC to select  $k = \hat{r}$  (review the last report) as the number of singular values of PY that exceed  $\sigma(\sqrt{2n} + \sqrt{2q}), P = X(X'X)^{-1}X'$
- $\triangleright$  Compute the rank constrained GLASSO estimator  $B_k$ with  $k = \hat{r}$  to obtain the final estimator  $\hat{B}^{(1)} = \hat{B}r$

This two-step estimator adapts to both rank and row sparsity under 2 additional mild restrictions

- Ass  $C_1$ :  $d_r(XA) > 2\sqrt{2}\sigma(\sqrt{n} + \sqrt{q})$  (RSC)
- Ass  $C_2$ :  $log(||XA||_F) < (\sqrt{2}-1)^2 \frac{n+q}{4}$

And we have the following property.



#### Theorem 4.4

Let  $\Sigma$  satisfy  $A(J, \delta_J)$  with  $J = J(A) \neq \Phi$ , let  $\frac{\lambda_1(\Sigma)}{\delta_J}$  be bounded and let  $C_1$ ,  $C_2$  hold.

$$\mathbb{E}[||X\hat{B}^{(1)} - XA||^2] \lesssim nr + |J|r \cdot \log(p)$$

The practical choice of the threshold  $2\sigma(\sqrt{n}+\sqrt{q})$  can be done by CV

# method 2(RCGL $\rightarrow$ JRRS1)



- pre-specify a grid  $\Lambda$  of values for  $\lambda$  and use (2) to construct  $\mathcal{B} = \{\hat{B_{k,\lambda}} : \lambda \in \Lambda\}$

# method 2(RCGL→JRRS1)



- $\triangleright$  pre-specify a grid  $\Lambda$  of values for  $\lambda$  and use (2) to construct  $\mathcal{B} = \{\hat{B_{k,\lambda}} : \lambda \in \Lambda\}$
- ► Compute  $B^{(2)} = \arg\min_{B \in \mathcal{B}} (||Y XB||^2 + pen(B))$

#### Theorem 4.5

Provided  $\Sigma$  satisfies condition  $A(J, \delta_J)$  with  $J = J(A) \neq \Phi$ ,  $\lambda_1(\Sigma)/\delta_I$  is bounded, and  $\Lambda$  contains  $\lambda$  for c >> 1

$$\mathbb{E}[||XB^{(2)} - XA||^2] \lesssim nr + |J|\log(p)r$$

and  $\hat{B}^{(2)}$  has the same rate as  $\hat{B}^{(1)}$ 

# method 3(GLASSO $\rightarrow$ RSC)



- Select the predictors via GLASSO
- Based only on the selected predictors, use RSC to construct an adaptive estimator of reduced rank





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- Joint variable and rank selection for parsimonious

   Introduction procedure and property numerical

  performance
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- wavelet-based sparse reduce-rank regression



For minimizing  $F(B; \lambda) = \frac{1}{2}||Y - XB||^2 + \lambda||B||_{2,1}$  over all  $p \times n$  matrices B of rank less than or equal to k. By using the polar decomposition B = SV, where V is orthogonal and S is semi-positive definite.

$$(\hat{S}, \hat{V}) := F(S, V; \lambda) = \arg \min_{S \in \mathbb{R}^{p \times k}, V \in O^{n \times k}} \frac{1}{2} ||Y - XCV||^2 + \lambda ||S||_{2.1}$$
(3)

We propse the following iterative optimization precedure.

# algorithm



Given  $1 \le k \le m \land p \land n, \lambda \ge 0, V_{k,\lambda}^{(0)} \in O^{n \times k}$  (first k columns of  $I_{n \times n}$ )

 $i \leftarrow 0$ , converged  $\leftarrow$  FALSE while not converged do:

(a). 
$$S_{k,\lambda}^{(j+1)} \leftarrow arg\min_{S \in \mathbb{R}^{p \times k}} \frac{1}{2} ||YV_{k,\lambda}^{(j)} - XS||^2 + \lambda ||S||_{2,1}$$

(b). 
$$W \leftarrow YXS_{k,\lambda}^{(j+1)}, W \in \mathbb{R}^{n \times k}$$
, Using SVD  $W = U_wD_wVw$ 

(c). 
$$V_{k,\lambda}^{(j+1)} \leftarrow U_w V_w$$

(d). 
$$B_{k,\lambda}^{(j+1)} \leftarrow S_{k,\lambda}^{(j+1)}(V_{k,\lambda}^{(j+1)})'$$

(e). converged 
$$\leftarrow |F(B_{k,\lambda}^{(j+1)};\lambda) - F(B_{k,\lambda}^{(j)};\lambda)| < \epsilon$$

(f). 
$$j \leftarrow j + 1$$

end while and diliver  $\hat{B_{k,\lambda}} = B_{k,\lambda}^{(j+1)}, \hat{S_{k,\lambda}} = S_{k,\lambda}^{(j+1)}, \hat{V_{k,\lambda}} = V_{k,\lambda}^{(j+1)}$ 

#### some remarks



#### Note

We run the algorithm to obtain a solution path, for each  $(k, \lambda)$ in a 2-dimensional grid or a grid of  $\lambda$  with k determined by RSC

From the solution path, we get a series of candidate estimates. Then the single stage JRRS or other tuning criterion can be used to select the optimal estimate.



#### Note

We run the algorithm to obtain a solution path, for each  $(k, \lambda)$ in a 2-dimensional grid or a grid of  $\lambda$  with k determined by RSC

From the solution path, we get a series of candidate estimates. Then the single stage JRRS or other tuning criterion can be used to select the optimal estimate.

Step (a) needs to solve a GLASSO optimization problem.





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- Dimension reduction and coefficient estimation
   model setup property Tuning
- 6 Reduced rank stochastic regression
- 7 wavelet-based sparse reduce-rank regression





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 6 Reduced rank stochastic regression
- 7 wavelet-based sparse reduce-rank regression



$$Y = XB + E \rightarrow Y = F\Omega + E$$

where  $B = \Gamma\Omega$ ,  $F = X\Gamma$ ,  $\Gamma \in \mathbb{R}^{p \times r}$  for r < min(p, q). The  $\Omega \in \mathbb{R}^{r \times q}$  is called factor loading matrix. The columns of  $F, F_i (j = 1, \dots, r)$  represent the so-called factors

$$Y_j = XB_j + E_j$$



$$Y = XB + E \rightarrow Y = F\Omega + E$$

where  $B = \Gamma \Omega$ ,  $F = X\Gamma$ ,  $\Gamma \in \mathbb{R}^{p \times r}$  for r < min(p, q). The  $\Omega \in \mathbb{R}^{r \times q}$  is called factor loading matrix. The columns of  $F, F_i (j = 1, \dots, r)$  represent the so-called factors

$$Y_j = XB_j + E_j$$

(i represents the columns of a matrix) The basic idea of dimension reduction is that the regression coefficient  $B_1, \cdots, B_q$  actually come from a linear space  $\mathcal{B}$  of dimension lower than p.



$$Y = XB + E \rightarrow Y = F\Omega + E$$

where  $B = \Gamma \Omega$ ,  $F = X\Gamma$ ,  $\Gamma \in \mathbb{R}^{p \times r}$  for r < min(p, q). The  $\Omega \in \mathbb{R}^{r \times q}$  is called factor loading matrix. The columns of  $F, F_i (j = 1, \dots, r)$  represent the so-called factors

$$Y_j = XB_j + E_j$$

(i represents the columns of a matrix) The basic idea of dimension reduction is that the regression coefficient  $B_1, \dots, B_q$  actually come from a linear space  $\mathcal{B}$  of dimension lower than p.

As a result, we have a set of basis elements  $\{\eta_1, \dots, \eta_p\}$  for  $\mathbb{R}^p$  and a subset  $\mathcal{A} \subset \{1, \cdots, p\}$  such that  $\mathcal{B} \subset span\{\eta_i : i \in \mathcal{A}\}$ 



Now we have the model as follows:

$$Y = F\Omega + E = X\Gamma\Omega + E = XB + E$$

where 
$$F = F_1, \dots, F_p$$
,  $F_i = X\eta_i, \Gamma = (\eta_1, \dots, \eta_p), B = \Gamma\Omega = (\eta_1, \dots, \eta_p)\Omega$ .



Now we have the model as follows:

$$Y = F\Omega + E = X\Gamma\Omega + E = XB + E$$

where  $F = F_1, \dots, F_n$ ,  $F_i = X\eta_i, \Gamma = (\eta_1, \dots, \eta_n), B = \Gamma\Omega = \Gamma$  $(\eta_1, \cdots, \eta_p)\Omega$ . A family of estimates for this can be obtained by

$$min\{tr(Y-F\Omega)W(Y-F\Omega)'\}$$
 subject to  $\sum_{i=1}^{p}||\omega_i||_{\alpha}\leq t$ 

where  $\omega_i$  is the ith row of  $\Omega$ , W is a weight matrix with common choices  $\Sigma^{-1}$  or I(which is corresponding to Frobenius)norm).



Now we have the model as follows:

$$Y = F\Omega + E = X\Gamma\Omega + E = XB + E$$

where  $F = F_1, \dots, F_n$ ,  $F_i = X\eta_i, \Gamma = (\eta_1, \dots, \eta_n), B = \Gamma\Omega = \Gamma$  $(\eta_1, \cdots, \eta_p)\Omega$ . A family of estimates for this can be obtained by

$$min\{tr(Y-F\Omega)W(Y-F\Omega)'\}$$
 subject to  $\sum_{i=1}^{p}||\omega_i||_{\alpha}\leq t$ 

where  $\omega_i$  is the ith row of  $\Omega$ , W is a weight matrix with common choices  $\Sigma^{-1}$  or I(which is corresponding to Frobenius norm). Here we assume W = I.



$$Y = X\eta_1\omega_1 + \cdots + X\eta_p\omega_p$$

The ith factor will be included if and only if  $\omega_i$  is non-zero.



$$Y = X\eta_1\omega_1 + \cdots + X\eta_p\omega_p$$

The ith factor will be included if and only if  $\omega_i$  is non-zero. We choose  $\alpha = 2$  and we need to obtain  $\eta s$  first.



$$Y = X\eta_1\omega_1 + \cdots + X\eta_p\omega_p$$

The ith factor will be included if and only if  $\omega_i$  is non-zero. We choose  $\alpha=2$  and we need to obtain  $\eta s$  first. We choose  $\eta s$  to be the eigenvectors of BB', because this set of basis contains the basis of  $\mathcal{B}$ .



$$Y = X\eta_1\omega_1 + \cdots + X\eta_p\omega_p$$

The ith factor will be included if and only if  $\omega_i$  is non-zero. We choose  $\alpha=2$  and we need to obtain  $\eta s$  first. We choose  $\eta s$  to be the eigenvectors of BB', because this set of basis contains the basis of  $\mathcal{B}$ . We can understand this by the following truth: B = UDV' is the SVD, then  $BB' = UD^2U'$  we choose  $(\eta_1, \cdots, \eta_n) = U$  span the column space of B Then  $\Omega = DV$ ,  $\omega = \sigma_i v_i$ , where  $v_i$  is the ith column of V, and  $||\omega_i|| = \sigma_i$ .



$$Y = X\eta_1\omega_1 + \cdots + X\eta_p\omega_p$$

The ith factor will be included if and only if  $\omega_i$  is non-zero. We choose  $\alpha=2$  and we need to obtain  $\eta s$  first. We choose  $\eta s$  to be the eigenvectors of BB', because this set of basis contains the basis of  $\mathcal{B}$ . We can understand this by the following truth: B = UDV' is the SVD, then  $BB' = UD^2U'$  we choose  $(\eta_1, \cdots, \eta_n) = U$  span the column space of B Then  $\Omega = DV$ ,  $\omega = \sigma_i v_i$ , where  $v_i$  is the ith column of V.and  $||\omega_i|| = \sigma_i$ . Finally we have the optimization problem:

$$min(tr(Y-XB)(Y-XB)')$$
 subject to  $\sum_{i=1}^{min(p,q)} \sigma_i \leq t$  (4)



The last term is known as Ky Fan norm for B. There's is no restriction of B because once the estimation B is available, the basis  $\eta s$  can be obtained as its left singular vectors U. Therefore, we can also compute the factors  $F_i = X\eta_i$  and loading  $\Omega = DV$ .



The last term is known as Ky Fan norm for B. There's is no restriction of B because once the estimation B is available, the basis  $\eta s$  can be obtained as its left singular vectors U. Therefore, we can also compute the factors  $F_i = X\eta_i$  and loading  $\Omega = DV$ .

For the tuning  $\alpha$  cases, we get the optimization problem tr(Y-XB)(Y-XB)' subject to  $(\sum_i \sigma_i^{\alpha})^{\frac{1}{\alpha}} < t$  and RRR is another special case of expression with  $\alpha = 0^+$ 





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- Dimension reduction and coefficient estimation

  model setup property Tuning
- 6 Reduced rank stochastic regression
- 7 wavelet-based sparse reduce-rank regression

# orthogonal design



### Lemma 5.1

Let  $U^{\hat{L}S}D^{\hat{L}S}V^{\hat{L}S}$  be the SVD of the least squares estimate  $B^{\hat{L}S}$ . Then, under the orthogonal design where X'X = nI, the minimizer of (4) is  $\hat{B} = \hat{U}^{LS}\hat{D}(\hat{V}^{LS})', \hat{D}_{ii} = \max(\hat{D}^{LS}_{ii}, 0)$  (singular values are shrunk), and  $\lambda > 0$  is a constant such that  $\sum_{i} \hat{D}_{ii} = min(t, \sum_{i} \hat{D}_{ii}^{\hat{L}S})$ 

# orthogonal design



### Lemma 5.1

Let  $\hat{U^{LS}}\hat{D^{LS}}\hat{V^{LS}}$  be the SVD of the least squares estimate  $\hat{B^{LS}}$ . Then, under the orthogonal design where X'X = nI, the minimizer of (4) is  $\hat{B} = \hat{U}^{LS}\hat{D}(\hat{V}^{LS})', \hat{D}_{ii} = \max(\hat{D}^{LS}_{ii}, 0)$  (singular values are shrunk), and  $\lambda \geq 0$  is a constant such that  $\sum_{i} \hat{D}_{ii} = min(t, \sum_{i} \hat{D}_{ii}^{\hat{L}S})$ 

### Note

In fact, the  $\lambda$  arisen is the result of Lagrange multipilication

# orthogonal design



### Lemma 5.1

Let  $U^{\hat{L}S}D^{\hat{L}S}V^{\hat{L}S}$  be the SVD of the least squares estimate  $B^{\hat{L}S}$ . Then, under the orthogonal design where X'X = nI, the minimizer of (4) is  $\hat{B} = \hat{U}^{LS}\hat{D}(\hat{V}^{LS})', \hat{D}_{ii} = \max(\hat{D}^{LS}_{ii}, 0)$  (singular values are shrunk), and  $\lambda > 0$  is a constant such that  $\sum_{i} \hat{D}_{ii} = min(t, \sum_{i} \hat{D}_{ii}^{\hat{L}S})$ 

### Note

In fact, the  $\lambda$  arisen is the result of Lagrange multipilication

### Lemma 5.2

Suppose that max(p, q) = o(n), under the orthogonal design, if  $\lambda \to 0$  in such a fashion that  $\frac{\max(p,q)}{n} = o(\lambda^2)$ . Then  $|\sigma_i(\hat{B}) - \sigma_i(B)| \to 0$  with probability if  $\sigma(B) > 0$  and  $P(\sigma(\hat{B}) = 0) \rightarrow 1 \text{ if } \sigma(B) = 0$ 





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- Dimension reduction and coefficient estimation
   model setup property Tuning
- 6 Reduced rank stochastic regression
- 7 wavelet-based sparse reduce-rank regression



We develop a GCV type of statistic for determining t. We give a lagrange form:

$$Q_n(B) = \frac{1}{2} tr(Y - XB)(Y - XB)' + n\lambda \sum_{i=1}^{p \wedge q} \sigma_i(B)$$
 (5)

The following lemma explicitly describes the relationship between t and  $\lambda$ 



We develop a GCV type of statistic for determining t. We give a lagrange form:

$$Q_n(B) = \frac{1}{2} tr(Y - XB)(Y - XB)' + n\lambda \sum_{i=1}^{p \wedge q} \sigma_i(B)$$
 (5)

The following lemma explicitly describes the relationship between t and  $\lambda$ 

### Lemma 5.3

write  $\hat{d}_i = \hat{D}_{ii}$  for  $i = 1, \dots, p \land q$ . For any  $t \leq \sum_i \hat{d}_i$ , the minimizer of equation (5) coincides with the minimizer of (4) if

$$n\lambda = \frac{1}{\#(\hat{d}_i > 0)} \sum_{\hat{d}_i > 0} (\tilde{X}_i' \tilde{Y}_i - \tilde{X}_i' \tilde{X}_i \hat{d}_i) \tag{6}$$



 $Y_i$  is the ith column of Y = YU and  $X_i$  is the ith column of  $\tilde{X} = X\hat{V}$ 



 $Y_i$  is the ith column of Y = YU and  $X_i$  is the ith column of  $\tilde{X} = X\hat{V}$ 

#### Note

We can transform  $\sum_{i=1}^{p \wedge q} \sigma_i(B)$  as follows:

$$\sum_{i=1}^{p\wedge q}\sigma_i(B)=\sum_{i=1}^{p\wedge q}\hat{D_{ii}}=\sum_{i=1}^p\sigma_i(\hat{B}K\hat{B}')=\operatorname{tr}(\hat{B}K\hat{B}')$$

where  $K = \sum_{\hat{D}_{ii} > 0} \frac{1}{\hat{D}_{ii}} \hat{v}_i \hat{v}_i$ .



 $Y_i$  is the ith column of Y = YU and  $X_i$  is the ith column of  $\tilde{X} = X\hat{V}$ 

#### Note

We can transform  $\sum_{i=1}^{p \wedge q} \sigma_i(B)$  as follows:

$$\sum_{i=1}^{p\wedge q}\sigma_i(B)=\sum_{i=1}^{p\wedge q}\hat{D_{ii}}=\sum_{i=1}^p\sigma_i(\hat{B}K\hat{B}')=\operatorname{tr}(\hat{B}K\hat{B}')$$

where  $K = \sum_{\hat{D}_{ii} > 0} \frac{1}{\hat{D}_{ii}} \hat{v}_i \hat{V}_i$ . Considering that  $\hat{B}K\hat{B}'=\sum_{\hat{D}_{ii}>0}\frac{1}{D_{ii}}\hat{B}\hat{v}_i(\hat{B}\hat{v}_i)'=\sum_i\frac{1}{D_{ii}}\sigma_iu_i\sigma_iu_i'=\sum_i\sigma_iu_iu_i'$  Using  $u_i \perp u_i$ , we get the eigenvalue of  $\hat{B}K\hat{B}'$  is  $\sigma_i$ . However,  $(\hat{B}K\hat{B}')(\hat{B}K\hat{B}')' = \sum_i \sigma_i^2 u_i u_i'$ . As a result, the singular value of  $\hat{B}K\hat{B}'$  is  $\sigma_i$ 



Then we can transform the Lagrange form into

$$\frac{1}{2}tr(Y-XB)(Y-XB)'+n\lambda tr(BKB')$$
 (7)

Since  $\hat{B}$  is the minimizer of (7), it can be expressed as  $\hat{B} = (X'X + 2n\lambda K)^{-1}X'Y$ 



Then we can transform the Lagrange form into

$$\frac{1}{2}tr(Y-XB)(Y-XB)'+n\lambda tr(BKB') \tag{7}$$

Since  $\hat{B}$  is the minimizer of (7), it can be expressed as  $\hat{B} = (X'X + 2n\lambda K)^{-1}X'Y$ We can define the hat matrix for expression (6) as  $H = X(X'X + 2n\lambda K)^{-1}X'$  and the degree of freedom as df(t) = qtrH



Then we can transform the Lagrange form into

$$\frac{1}{2}tr(Y-XB)(Y-XB)'+n\lambda tr(BKB') \tag{7}$$

Since  $\hat{B}$  is the minimizer of (7), it can be expressed as  $\hat{B} = (X'X + 2n\lambda K)^{-1}X'Y$ 

We can define the hat matrix for expression (6) as  $H = X(X'X + 2n\lambda K)^{-1}X'$  and the degree of freedom as df(t) = qtrH

The GCV score is given by  $GCV(t) = \frac{tr(Y-X\hat{B})(Y-X\hat{B})'}{qp-df(t)}$ . We choose a tuning parameter by minimizing GCV(t).



To sum up:

Step 1: for each candidate t-value

- (a). compute the minimizer of (4) (denote the solution B(t))
- (b). evaluate  $\lambda$  by using (6)
- (c). compute the GCV score

Step 2: denote  $t^*$  the minimizer of the GCV score. Return  $B(t^*)$  as the estimator of B





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression ■ intro ■ sparse unit regression ■ selection and higher rank
- wavelet-based sparse reduce-rank regression





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- intro sparse unit regression selection and higher rank
- wavelet-based sparse reduce-rank regression



$$Y = XC + E \quad rankC = r^*$$

We have SVD:

$$C = UDV = \sum_{k=1}^{r^*} d_k \mathbf{u_k} \mathbf{v_k} := \sum_{k=1}^{r^*} C_k$$

where  $U = (u_1, \dots, u_{r^*}), V = (v_1, \dots, v_{r^*}), C_k = d_k u_k v_k$  $C_k$  is the layer k unit rank matrix of C.



$$Y = XC + E \quad rankC = r^*$$

We have SVD.

$$C = UDV = \sum_{k=1}^{r^*} d_k \mathbf{u_k} \mathbf{v_k} := \sum_{k=1}^{r^*} C_k$$

where  $U = (u_1, \dots, u_{r^*}), V = (v_1, \dots, v_{r^*}), C_k = d_k u_k v_k$  $C_k$  is the layer k unit rank matrix of C.Here all the singular values are assumed to be distinct so that this SVD is unique because in practice, the singular values rarely conincide.

#### Introduction



We propse to estimate C by minimizing the following objective function with respect to  $(d_k, \mathbf{u_k}, \mathbf{v_k})$  for  $k = 1, \dots, r^*$ 

$$\frac{1}{2}||Y - X\sum_{k=1}^{r^*} d_k u_k v_k'||^2 + \sum_{k=1}^{r^*} Pe(\lambda_k, (d_k, u_k, v_k'))$$

#### Introduction



We propse to estimate C by minimizing the following objective function with respect to  $(d_k, \mathbf{u_k}, \mathbf{v_k})$  for  $k = 1, \dots, r^*$ 

$$\frac{1}{2}||Y - X\sum_{k=1}^{r^*} d_k u_k v_k'||^2 + \sum_{k=1}^{r^*} Pe(\lambda_k, (d_k, u_k, v_k'))$$

We consider

$$Pe = \lambda_k \sum_{i=1}^{p} \sum_{j=1}^{q} w_{ijk} |d_k u_{ik} v_{jk}| = \lambda_k (w_k^{(d)} d_k) (\sum_{i=1}^{p} w_{ik}^{(u)} |u_{ik}) (\sum_{j=1}^{q} w_{jk}^{(v)} |v_{jk}|)$$
(8)

### Introduction



We propse to estimate C by minimizing the following objective function with respect to  $(d_k, \mathbf{u_k}, \mathbf{v_k})$  for  $k = 1, \dots, r^*$ 

$$\frac{1}{2}||Y - X\sum_{k=1}^{r^*} d_k \mathbf{u_k} \mathbf{v_k}||^2 + \sum_{k=1}^{r^*} Pe(\lambda_k, (d_k, \mathbf{u_k}, \mathbf{v_k}))$$

We consider

$$Pe = \lambda_k \sum_{i=1}^{p} \sum_{j=1}^{q} w_{ijk} |d_k u_{ik} v_{jk}| = \lambda_k (w_k^{(d)} d_k) (\sum_{i=1}^{p} w_{ik}^{(u)} |u_{ik}) (\sum_{j=1}^{q} w_{jk}^{(v)} |v_{jk}|)$$
(8)

where  $w_{iik} = w_k^{(d)} w_{ik}^{(u)} w_{ik}^{(v)}$  are data-driven weights to be done below. It can be viewed as penalizing each of the singular vectors comprising the SVD layer.





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- intro sparse unit regression selection and higher rank
- wavelet-based sparse reduce-rank regression

### unit regression



The weights can be chosen as

$$w^{(d)} = |\tilde{d}|^{-\gamma}$$

$$w^{(u)} = (w_1^{(u)}, \dots, w_P^{(u)})' = |\tilde{u}|^{-\gamma}$$

$$w^{(v)} = (w_1^{(v)}, \dots, w_q^{(q)})' = |\tilde{v}|^{-\gamma}$$

where  $\gamma$  is a prespecified non-negative parameter and  $|\cdot|^{(-\gamma)}$  is defined componentwise for the enclosed vector ( and we use  $\gamma = 2$ ).

# unit regression



The weights can be chosen as

$$w^{(d)} = |\tilde{d}|^{-\gamma}$$

$$w^{(u)} = (w_1^{(u)}, \dots, w_P^{(u)})' = |\tilde{u}|^{-\gamma}$$

$$w^{(v)} = (w_1^{(v)}, \dots, w_q^{(q)})' = |\tilde{v}|^{-\gamma}$$

where  $\gamma$  is a prespecified non-negative parameter and  $|\cdot|^{(-\gamma)}$  is defined componentwise for the enclosed vector ( and we use  $\gamma = 2$ ). When  $r^* = 1$ , the problem is with respect to  $(d, \mathbf{u}, \mathbf{v})$ :

$$\frac{1}{2}||Y - dXuv'||^2 + \lambda \sum_{i=1}^{q} \sum_{j=1}^{q} w_{ij}|du_iv_j|$$



For fixed  $\boldsymbol{u}$  the problem with respect to  $(d, \boldsymbol{v})$  becomes:

$$\frac{1}{2}||y - X^{(v)}\check{v}||^2 + \lambda^{(v)} \sum_{j=1}^{q} |\check{v}_j|$$
 (9)

where

$$\check{v} = diag(dw^{(v)})v, y = vec(Y), X^{(v)} = diag(w^{(v)})^{-1} \otimes (Xu)$$
 and  $\lambda^{(v)} = \lambda w^{(d)}(\sum_{i=1}^p w_i^{(u)}|u_i|)$ 



For fixed  $\boldsymbol{u}$  the problem with respect to  $(d, \boldsymbol{v})$  becomes:

$$\frac{1}{2}||y - X^{(v)}\check{\mathbf{v}}||^2 + \lambda^{(v)} \sum_{j=1}^{q} |\check{\mathbf{v}}_j|$$
 (9)

where

$$\check{v} = diag(dw^{(v)})v, y = vec(Y), X^{(v)} = diag(w^{(v)})^{-1} \otimes (Xu)$$
 and  $\lambda^{(v)} = \lambda w^{(d)}(\sum_{i=1}^p w_i^{(u)}|u_i|)$ 

This model can be recognized as a lasso regression with respect to  $\check{v}$ 



In contrast, for fixed v, the problem with respect to  $(d, \mathbf{u})$ becomes

$$\frac{1}{2}||y - X^{(u)}\check{u}||^2 + \lambda^{(u)} \sum_{i=1}^{p} |\check{u}_i|$$
 (10)

where  $\check{u} = diag(dw^{(u)})u, X^{(u)} = v \otimes Xdiag(w^{(u)})^{-1}, \lambda^{(u)} = v \otimes Xdiag(w^{(u)})^{-1}$  $\lambda w^{(d)}(\sum_{j=1}^q w_j^{(v)}|v_j|).$  Again this is a lasso regression problem with respect to  $\check{u}$ 





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression ■ intro ■ sparse unit regression ■ selection and higher rank
- wavelet-based sparse reduce-rank regression

### parameter selection



Denote  $(\hat{d}^{(\lambda)}, \hat{u}^{(\lambda)}, \hat{v}^{(\lambda)})$  as the fitted value of  $(d, \mathbf{u}, \mathbf{v})$  with the regularization parameter being  $\lambda$ . Define BIC as:

$$BIC(\lambda) = log(SSE(\lambda)) + \frac{log(nq)}{nq}df(\lambda)$$

where 
$$SSE(\lambda) = ||Y - \hat{d}^{(\lambda)}X\hat{u}^{(\lambda)}\hat{v}^{(\lambda)}||^2$$
,  $\hat{d}f(\lambda) = \sum_{i=1}^p I(u_i^{(\lambda)} \neq 0) + \sum_{j=1}^q I(v_j^{(\lambda)} \neq 0) - 1$ 

# Extension to the higher rank case



Exclusive extraction algorithm(EEA)

This idea is to seek a  $\hat{C}$  with sparse SVD structure near some initial consistent estimator, e.g. the least squares reduced rank regression estimator C whose SVD is given by

$$\sum_{k=1}^{r^*} \tilde{d}_k \tilde{u}_k \tilde{v}_k = \sum_{k=1}^{r^*} \tilde{C}_k$$

# Extension to the higher rank case



Exclusive extraction algorithm(EEA)

This idea is to seek a  $\hat{C}$  with sparse SVD structure near some initial consistent estimator, e.g. the least squares reduced rank regression estimator C whose SVD is given by

$$\sum_{k=1}^{r^*} \tilde{d}_k \tilde{u}_k \tilde{v}_k = \sum_{k=1}^{r^*} \tilde{C}_k$$
The EEA is as follows:

- (a). for each  $k \in \{1, \dots, r^*\}$ 
  - (1). construct the adaptive weights  $w_k^{(d)}=|\tilde{d}_k|^{-\gamma}, w_k^{(u)}=|\tilde{u_k}|^{-\gamma}$  and  $w_k^{(v)}=|\tilde{v_k}|^{-\gamma}$
  - (2). construct the exclusive layer  $Y_k = Y X(\tilde{C} \tilde{C}_k)$
  - (3). find  $(\hat{d}_k, \hat{u}_k, \hat{v}_k)$  by performing the sparse unit rank regression of  $Y_k$  on X with  $\lambda_k$  chosen by BIC
- (b). The final estimator C is given by  $\hat{C} = \sum_{k=1}^{r^*} \hat{d}_k \hat{u}_k \hat{V}_{\nu}$



### Note

We can also have the method done iteratively called the iterative exclusive extraction algorithm e.g.

$$\frac{||\textit{C}^{(\hat{\textit{i}}+1)} - \hat{\textit{C}^{(\textit{i})}}||}{||\hat{\textit{C}^{(\textit{i})}}|} < \epsilon$$

for example  $\epsilon = 10^{-6}$ 





- 1 Stability approach to regularization selection
- 2 sparse reduced-rank regression for simultaneous
- 3 sparse RRR for simultaneous rank and variable
- 4 joint variable and rank selection for parsimonious
- 5 Dimension reduction and coefficient estimation
- 6 Reduced rank stochastic regression
- 7 wavelet-based sparse reduce-rank regression

### model establishment



$$\tilde{Y} = DWV + \tilde{N} \tag{11}$$

where D is an  $n \times n$  orthogonal matrix(using Schmidt orthogonalization), V is  $p \times r$  unknown orthogonal matrix.  $W \in \mathbb{R}^{n \times r}$ 



$$\tilde{Y} = DWV + \tilde{N}$$
 (11)

where D is an  $n \times n$  orthogonal matrix(using Schmidt orthogonalization), V is  $p \times r$  unknown orthogonal matrix.  $W \in \mathbb{R}^{n \times r}$ 

The optimization problem is as follows:

$$J(V, W) = \frac{1}{2} ||\tilde{Y} - DWV||^2 + \sum_{i} \lambda_{i} ||w_{(i)}||_1$$
 (12)

$$(\hat{V}, \hat{W}) = arg \min_{V,W} J(V, W)$$
 such that  $VV = I_r$  (13)

### cyclic descent algorithm



(1). W-step: Given a fixed V, the problem can be written as

$$arg \min_{W} \frac{1}{2} ||B - W||^2 + \sum_{i} \lambda_{i} ||w_{(i)}||_1$$
 (14)

where  $B = D'\tilde{Y}V$ , and the solution is

$$\hat{w}_{jj} = \max(0, |b_{ji}| - \lambda_i) \frac{b_{ij}}{|b_{ij}|}$$
 (15)

### cyclic descent algorithm



(1). W-step: Given a fixed V, the problem can be written as

$$\arg\min_{W} \frac{1}{2} ||B - W||^2 + \sum_{i} \lambda_{i} ||w_{(i)}||_1$$
 (14)

where  $B = D'\tilde{Y}V$ , and the solution is

$$\hat{w}_{jj} = \max(0, |b_{ji}| - \lambda_i) \frac{b_{ij}}{|b_{ij}|}$$
 (15)

(2). V-step: Given a fixed W

$$arg \min_{V} ||\tilde{Y} - DWV||^2$$
 such that  $VV = I_r$ 

which has a solution given by  $\hat{V} = QG'$  where Q and Gare computed using  $M = \hat{W}D\hat{Y} = Q\Sigma G$ 

#### initialization and selection



Using  $\tilde{Y}=U_{\tilde{Y}}S_{\tilde{Y}}V_{\tilde{V}}$ , and setting  $V_0=V_{\tilde{Y}}$ , we can proceed the algorithm above.

Noting that (14) is separable and can be written as

$$arg\min_{w_{(i)}} \frac{1}{2} ||b_{(i)} - w_{(i)}||^2 + \lambda_i ||w_{(i)}||_1$$

### initialization and selection



Using  $\tilde{Y} = U_{\tilde{Y}} S_{\tilde{Y}} V_{\tilde{Y}}$ , and setting  $V_0 = V_{\tilde{Y}}$ , we can proceed the algorithm above.

Noting that (14) is separable and can be written as

$$arg \min_{w_{(i)}} \frac{1}{2} ||b_{(i)} - w_{(i)}||^2 + \lambda_i ||w_{(i)}||_1$$

The **SURE criterion** is given by

$$SURE(\lambda_i) = ||\hat{w}_{(i)} - b_{(i)}||^2 - n + 2n_i$$

where  $n_i = \#\{j : |\hat{w}_{ii}| > \lambda_i\}$ 



The algorithm is from the journal: **Discovering genetic** associations with high-dimensionality neuroimaging phenotypes a sparse reduced-rank regression approach



The algorithm is from the journal: **Discovering genetic** associations with high-dimensionality neuroimaging phenotypes a sparse reduced-rank regression approach The full rank coefficient matrix  $C_{(R)}$  and the estimated residual covariance matrix  $\hat{S_{\epsilon\epsilon(Y)}} = (Y - X\hat{C_{(r)}})'(Y - X\hat{C_{(r)}}).$ 



The algorithm is from the journal: **Discovering genetic** associations with high-dimensionality neuroimaging phenotypes a sparse reduced-rank regression approach The full rank coefficient matrix  $C_{(R)}$  and the estimated residual covariance matrix  $\hat{S}_{\epsilon\epsilon(Y)} = (Y - X\hat{C}_{(r)})'(Y - X\hat{C}_{(r)})$ . The rank trace plotting is obtained by plotting, for all values of r in a range from 0 to R, the following 2 quantities:

$$\Delta \hat{C_{(r)}} = \frac{||\hat{C_{(R)}} - \hat{C_{(r)}}||}{||\hat{C_{(R)}} - \hat{C_{(0)}}||} \qquad \Delta \hat{S_{\epsilon\epsilon(r)}} = \frac{||\hat{S_{\epsilon\epsilon(R)}} - \hat{S_{\epsilon\epsilon(r)}}||}{||\hat{S_{\epsilon\epsilon(R)}} - \hat{S_{\epsilon\epsilon(0)}}||}$$



The algorithm is from the journal: **Discovering genetic** associations with high-dimensionality neuroimaging phenotypes a sparse reduced-rank regression approach The full rank coefficient matrix  $C_{(R)}$  and the estimated residual covariance matrix  $\hat{S}_{\epsilon\epsilon(Y)} = (Y - X\hat{C}_{(r)})'(Y - X\hat{C}_{(r)})$ . The rank trace plotting is obtained by plotting, for all values of r in a range from 0 to R, the following 2 quantities:

$$\Delta \hat{C_{(r)}} = \frac{||\hat{C_{(R)}} - \hat{C_{(r)}}||}{||\hat{C_{(R)}} - \hat{C_{(0)}}||} \qquad \Delta \hat{S_{\epsilon\epsilon(r)}} = \frac{||\hat{S_{\epsilon\epsilon(R)}} - \hat{S_{\epsilon\epsilon(r)}}||}{||\hat{S_{\epsilon\epsilon(R)}} - \hat{S_{\epsilon\epsilon(0)}}||}$$

As r varies from 0 to R in both X and Y axes, both coefficients take values in [0,1]. As more ranks are added, starting at the top-right corner with r = 0, the curve moves towards the origin of the plot. When a further rank addition doesn't produce a significant reduction the plot indicates an "optimal" rank  $R^*$  which has been found.



# 谢谢!