

# Two Papers on Sparse Feature Selection in Multivariate Regression

Gregory Faletto

Department of Data Sciences and Operations  
Statistics Group  
University of Southern California Marshall School of Business

DSO 607, April 1 2019

- 1 Background on Multivariate Regression
- 2 Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)
  - Background and Problem Statement
  - Main Results
  - Selected Simulation Studies
- 3 Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)
  - Background and Problem Statement
  - Optimization Problem
  - Main Results
  - Selected Simulation Study

# Setup and Notation for Multivariate Linear Regression (single observation)

$$y_i^T = x_i^T B^* + w_i^T$$

# Setup and Notation for Multivariate Linear Regression (single observation)

Formulation for individual response  $\mathbf{y}_i$ :

$$\underbrace{\mathbf{y}_i^T}_{1 \times K} = \underbrace{\mathbf{x}_i^T}_{1 \times p} \underbrace{\mathbf{B}^*}_{p \times K} + \underbrace{\mathbf{w}_i^T}_{1 \times K} \quad (1)$$

where

- $\mathbf{y}_i \in \mathbb{R}^K$  is a  $K$ -vector of responses (rather than a scalar response as in univariate linear regression),
- $\mathbf{x}_i \in \mathbb{R}^p$  is a  $p$ -vector of predictors associated with  $\mathbf{y}_i$  (same as in univariate regression),
- $\mathbf{B}^* \in \mathbb{R}^{p \times K}$  is a  $p \times K$  matrix of coefficients (rather than a  $p$ -vector of coefficients as in univariate regression),
- $\mathbf{w}_i \in \mathbb{R}^K$  is a  $K$ -vector of random errors.

# Setup and Notation for Multivariate Linear Regression (full data set)

The diagram illustrates the matrix equation for multivariate linear regression. It consists of three main parts: a vector of target values, a vector of feature values, and a vector of weights, all connected by an equals sign and a plus sign.

- Left side (Green boxes):** A vertical stack of boxes representing the target vector  $y^T$ . The boxes are labeled  $y_1^T$ ,  $y_2^T$ ,  $y_3^T$ , followed by a vertical ellipsis  $\vdots$ , and finally  $y_n^T$ .
- Equality:** An equals sign  $=$  connects the target vector to the feature vector.
- Middle (Blue boxes):** A vertical stack of boxes representing the feature vector  $x^T$ . The boxes are labeled  $x_1^T$ ,  $x_2^T$ ,  $x_3^T$ , followed by a vertical ellipsis  $\vdots$ , and finally  $x_n^T$ .
- Matrix:** A purple square box labeled  $B^*$  represents the coefficient matrix, positioned between the feature vector and the weight vector.
- Addition:** A plus sign  $+$  connects the matrix  $B^*$  to the weight vector.
- Right side (Red boxes):** A vertical stack of boxes representing the weight vector  $w^T$ . The boxes are labeled  $w_1^T$ ,  $w_2^T$ ,  $w_3^T$ , followed by a vertical ellipsis  $\vdots$ , and finally  $w_n^T$ .

# Setup for Multivariate Linear Regression (full data set)

Formulation for full data set:

$$\underbrace{\mathbf{Y}}_{n \times K} = \underbrace{\mathbf{X}}_{n \times p} \underbrace{\mathbf{B}^*}_{p \times K} + \underbrace{\mathbf{W}}_{n \times K} \quad (2)$$

where

- $\mathbf{Y} \in \mathbb{R}^{n \times K} = (\mathbf{y}_1, \dots, \mathbf{y}_n)^T$  is an  $n \times K$  matrix of responses (row  $i$  contains response  $\mathbf{y}_i$ ),
- $\mathbf{X} \in \mathbb{R}^{n \times p}$  is an  $n \times p$  matrix of predictors (same as in univariate regression),
- $\mathbf{B}^* \in \mathbb{R}^{p \times K}$  is a  $p \times K$  matrix of coefficients,
- $\mathbf{W} \in \mathbb{R}^{n \times K} = (\mathbf{w}_1, \dots, \mathbf{w}_n)^T$  is an  $n \times K$  matrix of random errors. The error vectors are assumed to be IID with  $\mathbb{E}(\mathbf{w}_i) = \mathbf{0}$ .

# Some Motivating Examples

- Predicting both Math and Verbal SAT score as a function of parental income, extracurricular activity participation, etc.
- Responses of several measures of health to a drug in the presents of other covariates (age, weight, gender, etc.)
- Prediction of the prices of several stocks as a function of economic indicators
- Used to find sets of genes that are expressed in different forms of cancer in Chen, Chan, and Stenseth (2012)

# Two Papers on Sparse Feature Selection in Multivariate Regression

## └ Background on Multivariate Regression

## └ Some Motivating Examples

- Predicting both Math and Verbal SAT score as a function of parental income, extracurricular activity participation, etc.
- Responses of several measures of health to a drug in the presence of other covariates (age, weight, gender, etc.)
- Prediction of the prices of several stocks as a function of economic indicators
- Used to find sets of genes that are expressed in different forms of cancer in Chen, Chan, and Stenseth (2012)

- Sometimes called “multi-task learning” in machine learning
- Won’t go through example of gene expression because it’s a different paradigm from the usual case we’ll discuss



- 1 Background on Multivariate Regression
- 2 Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)
  - Background and Problem Statement
  - Main Results
  - Selected Simulation Studies
- 3 Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)
  - Background and Problem Statement
  - Optimization Problem
  - Main Results
  - Selected Simulation Study

# Meaning of “Support Union Recovery”

- Feature selection in multivariate setting
- Sparsity assumption: column  $k$  of coefficient matrix  $\mathbf{B}^*$  has nonzero entries in a subset

$$S_k := \{i \in \{1, \dots, p\} \mid \beta_{ik}^* \neq 0\} \quad (3)$$

of size  $s_k := |S_k|$ .

- We seek to *recover* the *union* of the *supports*: the set

$$S := \bigcup_{k=1}^K S_k$$

of size  $s := |S|$ .

# Two Papers on Sparse Feature Selection in Multivariate Regression

└ Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)

- Feature selection in multivariate setting
- Sparsity assumption: column  $k$  of coefficient matrix  $\mathbf{B}^*$  has nonzero entries in a subset

$$S_k := \{i \in \{1, \dots, p\} \mid \beta_{ik}^* \neq 0\} \quad (3)$$

of size  $s_k := |S_k|$ .

- We seek to recover the union of the supports: the set

$$S := \bigcup_{k=1}^K S_k$$

of size  $s := |S|$ .

- Responses are related, so we expect supports  $S_k$  should have some overlap, and that the relationship between the supports and responses (i.e. the coefficients, or the vectors in the columns of  $\mathbf{B}$  are related).

# Block Regularization

- Used in many settings when goal is to choose whether or not to include groups of variables, rather than considering them individually.
- Example:  $\ell_1/\ell_q$  norm of a matrix takes the sum of the  $\ell_q$  norms of each row:

$$\|\mathbf{B}\|_{\ell_1/\ell_q} := \sum_{i=1}^p \left( \sum_{j=1}^K |\beta_{ij}|^q \right)^{1/q} = \sum_{i=1}^p \|\beta_i\|_q \quad (4)$$

where  $(\beta_{ik})_{1 \leq i \leq p, 1 \leq k \leq K}$  are the entries of  $\mathbf{B}$ .

- This paper makes use of the  $\ell_1/\ell_2$  norm

$$\|\mathbf{B}\|_{\ell_1/\ell_2} := \sum_{i=1}^p \sqrt{\sum_{j=1}^K \beta_{ij}^2} = \sum_{i=1}^p \|\beta_i\|_2 \quad (5)$$

# Multivariate Group Lasso

- We solve the following optimization problem to estimate  $\mathbf{B}^*$ :

$$\arg \min_{\mathbf{B} \in \mathbb{R}^{p \times K}} \left\{ \frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2 + \lambda_n \|\mathbf{B}\|_{\ell_1/\ell_q} \right\} \quad (6)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm:

$$\|A\|_F := \sqrt{\sum_{i,j} A_{i,j}^2}.$$

- Penalizing the  $\ell_1/\ell_1$  norm is equivalent to simply solving  $K$  separate lasso solutions.
- Penalizing the  $\ell_1/\ell_2$  norm groups rows together; the authors call this the **multivariate group lasso**.

# Two Papers on Sparse Feature Selection in Multivariate Regression

└ Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)

- We solve the following optimization problem to estimate  $\mathbf{B}^*$ :

$$\arg \min_{\mathbf{B} \in \mathbb{R}^{p \times m}} \left\{ \frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2 + \lambda_n \|\mathbf{B}\|_{\ell_1/\ell_2} \right\} \quad (6)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm:

$$\|\mathbf{A}\|_F := \sqrt{\sum_{i,j} A_{ij}^2}$$

- Penalizing the  $\ell_1/\ell_2$  norm is equivalent to simply solving  $K$  separate lasso solutions.
- Penalizing the  $\ell_1/\ell_2$  nor groups rows together; the authors call this the **multivariate group lasso**.

- Note that when  $K = 1$  this is exactly the lasso (regardless of  $q$ ):

$$\|\mathbf{B}\|_{\ell_1/\ell_1} := \sum_{i=1}^p \left( \sum_{j=1}^K |\beta_{ij}|^1 \right)^{1/1} = \sum_{i=1}^p |\beta_i| \quad (7)$$

- When  $q = 1$  the sum decouples (show on board)

# Multivariate Group Lasso (continued)

- The Frobenius norm is convex; therefore the optimization problem (6) is convex and can be solved efficiently. (In particular, it is a second-order cone program (SOCP) and can be solved efficiently with interior point methods.)

$$\arg \min_{\mathbf{B} \in \mathbb{R}^{p \times K}} \left\{ \frac{1}{2n} \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F^2 + \lambda_n \|\mathbf{B}\|_{\ell_1/\ell_q} \right\} \quad (6)$$

# Main Assumptions

- (A1) *Bounded eigenspectrum*: There exist fixed constants  $C_{\min} > 0$  and  $C_{\max} < +\infty$  such that all eigenvalues of the  $s \times s$  matrix  $\Sigma_{SS}$  (the covariance matrix of the true support) are contained in the interval  $[C_{\min}, C_{\max}]$ .
- (A2) *Irrepresentable Condition*: There exists a fixed parameter  $\gamma \in (0, 1]$  such that

$$\|\Sigma_{S^c S}(\Sigma_{SS})^{-1}\|_{\infty} \leq 1 - \gamma.$$

- (A3) *Self-incoherence*: There exists some  $D_{\max} < +\infty$  such that

$$\|(\Sigma_{SS})^{-1}\|_{\infty} \leq D_{\max}.$$



# Two Papers on Sparse Feature Selection in Multivariate Regression

└ Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)

## Main Assumptions

- (A1) *Bounded eigenspectrum*: There exist fixed constants  $C_{\min} > 0$  and  $C_{\max} < +\infty$  such that all eigenvalues of the  $s \times s$  matrix  $\Sigma_{SS}$  (the covariance matrix of the true support) are contained in the interval  $[C_{\min}, C_{\max}]$ .
- (A2) *Irrepresentable Condition*: There exists a fixed parameter  $\gamma \in (0, 1]$  such that

$$\|\Sigma_{S^c S}(\Sigma_{SS})^{-1}\|_{\infty} \leq 1 - \gamma.$$

- (A3) *Self-incoherence*: There exists some  $D_{\max} < +\infty$  such that

$$\|(\Sigma_{SS})^{-1}\|_{\infty} \leq D_{\max}.$$

- Intuition of bounded eigenspectrum condition: prevents excess dependence among elements of the design matrix associated with the support  $S$ .)
- Intuition of Irrepresentable condition: correlation between features in true support and noise features is not too high.

# Sparsity Overlap function

- *Sparsity-overlap function*  $\phi(\mathbf{B}^*)$  measures the sparsity of the matrix  $\mathbf{B}^*$  as well as the overlap between the different regressions (the different columns of  $\mathbf{B}^*$ ).
- In case of univariate regression  $K = 1$ , entries of design matrix are i.i.d. standard Gaussian; then  $\phi(\mathbf{B}^*) = s$ .
- For suitably correlated designs: sparsity overlap is  $\phi(\mathbf{B}^*) = s/K$ .
- Lemma 1(a): Under assumption (A1), the sparsity overlap  $\phi(\mathbf{B}^*)$  obeys the bounds

$$\frac{s}{c_{\max} K} \leq \phi(\mathbf{B}^*) \leq \frac{s}{c_{\min}}.$$

# Two Papers on Sparse Feature Selection in Multivariate Regression

└ Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)

- Sparsity-overlap function  $\phi(\mathbf{B}^*)$  measures the sparsity of the matrix  $\mathbf{B}^*$  as well as the overlap between the different regressions (the different columns of  $\mathbf{B}^*$ ).
- In case of univariate regression  $K = 1$ , entries of design matrix are i.i.d. standard Gaussian; then  $\phi(\mathbf{B}^*) = s$ .
- For suitably correlated designs: sparsity overlap is  $\phi(\mathbf{B}^*) = s/K$ .
- Lemma 1(a): Under assumption (A1), the sparsity overlap  $\phi(\mathbf{B}^*)$  obeys the bounds

$$\frac{s}{C_{\max} K} \leq \phi(\mathbf{B}^*) \leq \frac{s}{C_{\min}}$$

- We will see in Theorem 1 (next slide) that when the sparsity overlap is small, multivariate group lasso performs more favorably. Therefore we want  $C_{\min}$  to be large (not too close to 0, hopefully close to 1).

- 1 Background on Multivariate Regression
- 2 Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)
  - Background and Problem Statement
  - **Main Results**
  - Selected Simulation Studies
- 3 Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)
  - Background and Problem Statement
  - Optimization Problem
  - Main Results
  - Selected Simulation Study

# Main Results

(under assumptions A1 - A3 and some other mild conditions)

- (Summary of Theorem 1.) The *sample complexity*

$$\theta(n, p, s) := \frac{n}{2\psi(\mathbf{B}^*) \log(p - s)},$$

determines whether the multivariate group lasso recovers the exact row pattern (with high probability).

- (Summary of Corollary 1.) In particular, in the special case of the standard Gaussian ensemble, if  $\theta(n, p, s)$  exceeds a critical level  $\theta$ , the multivariate group lasso has a unique solution  $\hat{\mathbf{B}}$  with row support  $S(\mathbf{B})$  that is contained within the true row support  $S(\mathbf{B}^*)$ . Additionally,  $\|\hat{\mathbf{B}} - \mathbf{B}^*\|_{\ell_\infty/\ell_2}$  is bounded. If  $\theta(n, p, s)$  is below the critical level  $\theta$ , the multivariate group lasso fails.
- More generally, the multivariate group lasso succeeds for problem sequences  $(n, p, s)$  such that  $\theta(n, p, s)$  exceeds a critical level  $\theta_u$  and fails for sequences such that  $\theta(n, p, s)$  lies below a critical level  $\theta_\ell$ .

# Main Results (continued)

(under assumptions A1 - A3 and some other mild conditions)

- Recall Lemma 1(a): under assumption (A1) we have bounds on  $\phi(\mathbf{B}^*)$ , so  $\theta(n, p, s)$  ranges between

$$\frac{nC_{\min}}{2s \log(p-s)} \leq \theta(n, p, s) \leq \frac{nC_{\max}K}{2s \log(p-s)}. \quad (8)$$

This generalizes a previous result from Wanwright (2009): the lasso succeeds in performing exact support recovery if the ratio  $n/[s \log(p-s)]$  exceeds a certain threshold.

- (Summary of Corollary 3 and part of Corollary 1.) If  $\mathbf{X}$  is uncorrelated for the variables corresponding to the active rows of  $\mathbf{B}^*$ ,  $\ell_1/\ell_2$  group regularization never harms performance relative to an ordinary lasso approach.

# Two Papers on Sparse Feature Selection in Multivariate Regression

## Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)

- Per (8),  $\ell_1/\ell_2$  regularization for multivariate regression can yield substantial improvements in sample complexity (up to a factor of  $K$ ) when the coefficient vectors are suitably orthogonal.

- Recall Lemma 1(a): under assumption (A1) we have bounds on  $\phi(\mathbf{B}^*)$ , so  $\theta(n, p, s)$  ranges between

$$\frac{nC_{\min}}{2s \log(p-s)} \leq \theta(n, p, s) \leq \frac{nC_{\max}K}{2s \log(p-s)} \quad (8)$$

This generalizes a previous result from Wainwright (2009): the lasso succeeds in performing exact support recovery if the ratio  $n/[s \log(p-s)]$  exceeds a certain threshold.

- (Summary of Corollary 3 and part of Corollary 1.) If  $\mathbf{X}$  is uncorrelated for the variables corresponding to the active rows of  $\mathbf{B}^*$ ,  $\ell_1/\ell_2$  group regularization never harms performance relative to an ordinary lasso approach.

# Main Results (continued)

(under assumptions A1 - A3 and some other mild conditions)

(Summary of Corollary 2.) Given a method that returns the row supports  $S$  (with high probability) with  $|S| \ll p$ , recovering the individual supports  $S_k$  with high probability is easy using the following procedure:

- ① Compute the (restricted) multivariate ordinary least squares estimate

$$\tilde{\mathbf{B}}_{\hat{S}} = \arg \min_{\mathbf{B}_{\hat{S}}} \{ \|\mathbf{Y} - \mathbf{X}_{\hat{S}} \mathbf{B}_{\hat{S}}\|_F \} \quad (9)$$

where  $\mathbf{B}_{\hat{S}}$  is the coefficient matrix  $\mathbf{B}$  but only with rows in the estimated support union  $\hat{S}$  and  $\mathbf{X}_{\hat{S}}$  contains only the columns of  $\mathbf{X}$  corresponding to those features.

- ② Compute  $T(\tilde{\mathbf{B}}_{\hat{S}})$  by setting every entry in  $\tilde{\mathbf{B}}_{\hat{S}}$  with absolute value less than  $2\sqrt{2 \log(K|\hat{S}|)/(C_{\min} n)}$  equal to 0. Then the support is the nonzero entries of  $T(\tilde{\mathbf{B}}_{\hat{S}})$ .



- 1 Background on Multivariate Regression
- 2 Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)
  - Background and Problem Statement
  - Main Results
  - Selected Simulation Studies
- 3 Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)
  - Background and Problem Statement
  - Optimization Problem
  - Main Results
  - Selected Simulation Study

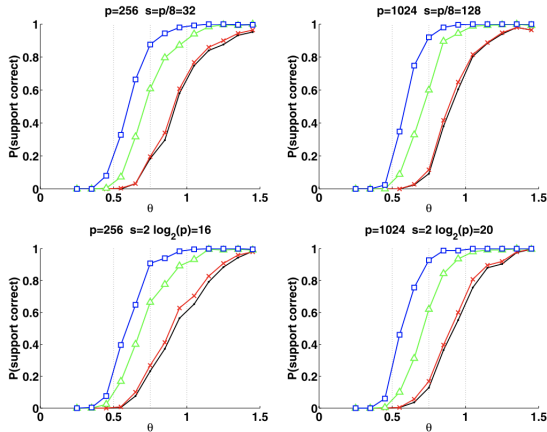


FIG. 3. Plots of support recovery probability,  $\mathbb{P}[\widehat{S} = S]$ , versus the control parameter  $\theta = n/[2s \log(p - s)]$  for two different type of sparsity: logarithmic sparsity on top ( $s = \mathcal{O}(\log(p))$ ) and linear sparsity on bottom ( $s = \alpha p$ ), and for increasing values of  $p$  from left to right. The noise level is set at  $\sigma = 0.1$ . Each graph shows four curves (black, red, green, blue) corresponding to the case of independent  $\ell_1$  regularization, and, for  $\ell_1/\ell_2$  regularization, the cases of identical regression, intermediate angles and orthonormal regressions. Note how curves corresponding to the same case across different problem sizes  $p$  all coincide, as predicted by Theorems 1 and 2. Moreover, consistent with the theory, the curves for the identical regression group reach  $\mathbb{P}[\widehat{S} = S] \approx 0.50$  at  $\theta \approx 1$ , whereas the orthonormal regression group reaches 50% success substantially earlier.



# Two Papers on Sparse Feature Selection in Multivariate Regression

## Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)

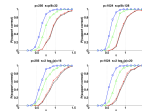


FIG. 1. Plot of support recovery probability,  $\Pr(\text{support union})$ , versus the number of features,  $p$ , for different regression groups. The curves are labeled as in the legend. The top row shows results for  $p=10$  and the bottom row for  $p=100$ . The left column shows results for  $\theta=0.5$  and the right column for  $\theta=1$ . The curves are labeled as in the legend. The blue curves show the highest recovery probability, followed by green, red, and black. The curves are labeled as in the legend. The blue curves show the highest recovery probability, followed by green, red, and black.

- Black curve is independent  $\ell_1$  regularization (ordinary lasso), which we see does the worst in pretty much all settings. Red is  $\ell_1/\ell_2$  regularization for identical regression (the columns of  $\mathbf{B}^*$  are identical). For orthonormal regressions (blue lines), the columns of  $\mathbf{B}^*$  are orthonormal, which is the most favorable setting for this method. For intermediate angles (green lines), the columns of  $\mathbf{B}^*$  are at a 60 degree angle.
- Note that the improvement of this method over lasso is largest when the regressions are orthonormal; in the other extreme when the columns are identical, it does about the same as lasso.
- Note that varying  $p$  doesn't really change anything if the control parameter  $\theta$  is equal.
- In identical regression group, Pr support recovery reaches 0.5 at  $\theta = 1$ ; improvement in other regimes.

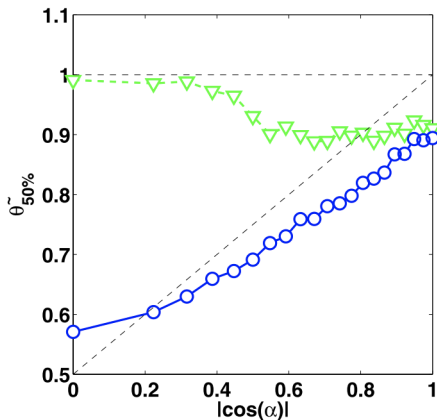


FIG. 5. Plots of the Lasso sample complexity  $\theta = n/[2s \log(p - s)]$  for which the probability of support union recovery exceeds 50% empirically as a function of  $|\cos(\alpha)|$  for  $\ell_1$ -based recovery and  $\ell_1/\ell_2$ -based recovery, where  $\alpha$  is the angle between  $Z^{(1)*}$  and  $Z^{(2)*}$  for the family  $\mathcal{B}_1$ . We consider the two following methods for performing row selection: Ordinary Lasso ( $\ell_1$ , green triangles) and multivariate group Lasso (blue circles).

## Two Papers on Sparse Feature Selection in Multivariate Regression

- Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)

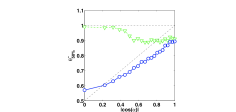


FIG. 5. Plots of the Lasso sample complexity  $\theta = n(2s \log(p) - c)$  for which the probability of support union recovery exceeds 50% empirically as a function of  $\alpha$  (the angle between columns of  $\mathbf{B}^*$ ) for  $\ell_1/\ell_2$ -based recovery and  $\ell_1/\ell_2$ -based recovery, where  $\alpha$  is the angle between  $\mathbf{Z}^{(1)}$  and  $\mathbf{Z}^{(2)}$  for the family  $\mathbf{B}_1$ . We consider the two following methods for performing row selection: Ordinary Lasso ( $\ell_1$ , green triangles) and multivariate group Lasso (blue circles).

- We have fixed  $p = 2048$  and sparsity  $s = \log_2(p) = 22$ .
- $\alpha$  is the angle between columns of  $\mathbf{B}^*$ , so as its cosine increases, the columns come closer together.
- We see that lasso needs a higher value of  $\theta$  to get a 50% chance of support recovery, but the gap decreases as the columns of  $\mathbf{B}^*$  come closer and closer together, until there is no gap when they are equal.
- Theory from this paper predicts that the circles and triangles should lie at or below dotted lines.
- Intuition of why this method does better when columns are orthonormal: in that case independent estimates of the support are more likely to include (by union) spurious covariates in the row support.

- 1 Background on Multivariate Regression
- 2 Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)
  - Background and Problem Statement
  - Main Results
  - Selected Simulation Studies
- 3 Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)
  - Background and Problem Statement
  - Optimization Problem
  - Main Results
  - Selected Simulation Study

# Singular Value Decomposition of the coefficient matrix $\mathbf{B}^*$

- Reduced rank regression model:  $\mathbf{Y} = \mathbf{X}\mathbf{B}^* + \mathbf{W}$  as before. We will be concerned with  $\text{rank}(\mathbf{B}^*) = r^* \leq \min\{p, K\}$ .
- Singular value decomposition (SVD):

$$\mathbf{B}^* = \mathbf{U}\mathbf{D}\mathbf{V}^T = \sum_{k=1}^{r^*} d_k \mathbf{u}_k \mathbf{v}_k^T = \sum_{k=1}^{r^*} \mathbf{B}_k \quad (10)$$

- $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_{r^*}) \in \mathbb{R}^{p \times r^*}$  consists of orthonormal left singular vectors
- $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_{r^*}) \in \mathbb{R}^{K \times r^*}$  consists of orthonormal right singular vectors
- $\mathbf{D} \in \mathbb{R}^{r^* \times r^*}$  is a diagonal matrix with positive singular values  $d_1 > \dots > d_{r^*}$  on its diagonal.



# Two Papers on Sparse Feature Selection in Multivariate Regression

## Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)

Singular Value Decomposition of the coefficient matrix  $B^*$

- Reduced rank regression model:  $\mathbf{Y} = \mathbf{X}\mathbf{B}^* + \mathbf{W}$  as before. We will be concerned with  $\text{rank}(\mathbf{B}^*) = r^* \leq \min\{p, K\}$ .
- Singular value decomposition (SVD):

$$\mathbf{B}^* = \mathbf{U}\mathbf{D}\mathbf{V}^T = \sum_{k=1}^{r^*} d_k \mathbf{u}_k \mathbf{v}_k^T = \sum_{k=1}^{r^*} \mathbf{B}_k \quad (10)$$

- $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_{r^*}) \in \mathbb{R}^{p \times r^*}$  consists of orthonormal left singular vectors
- $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_{r^*}) \in \mathbb{R}^{K \times r^*}$  consists of orthonormal right singular vectors
- $\mathbf{D} \in \mathbb{R}^{r^* \times r^*}$  is a diagonal matrix with positive singular values  $d_1 > \dots > d_{r^*}$  on its diagonal.

- Using singular value decomposition,  $\mathbf{B}^*$  can be expressed as a sum of  $r^*$  unit rank matrices  $\mathbf{B}_k$  that are proportional to the outer product of the left and right singular vectors.
- (The singular values are assumed to be distinct, so the decomposition is unique up to the signs of the singular vectors. In practice this is usually the case.)

# Singular Value Decomposition of $B^*$ (continued)

- **Key insight:**  $B^*$  is the sum of  $r^*$  orthogonal layers of decreasing importance. For each layer  $k$ ,  $\mathbf{u}_k$  are the predictor effects,  $\mathbf{v}_k$  are the response effects, and  $d_k$  indicates the relative importance of the association.
- If we believe that each pathway of association involves only a subset of the responses and predictors, then the left and right singular vectors should be sparse.

- 1 Background on Multivariate Regression
- 2 Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)
  - Background and Problem Statement
  - Main Results
  - Selected Simulation Studies
- 3 Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)
  - Background and Problem Statement
  - **Optimization Problem**
  - Main Results
  - Selected Simulation Study

# Optimization Problem

We plug the expression for  $\mathbf{B}^*$  from (10) into the typical penalized least squares objective function (similar to (6)) to obtain the objective function to optimize:

$$\arg \min_{d_k, \mathbf{u}_k, \mathbf{v}_k, k \in 1, \dots, r^*} \left\{ \frac{1}{2} \left\| \mathbf{Y} - \mathbf{X} \sum_{k=1}^{r^*} d_k \mathbf{u}_k, \mathbf{v}_k^T \right\|_F^2 + \sum_{k=1}^{r^*} p_\lambda(\lambda_k, (d_k, \mathbf{u}_k, \mathbf{v}_k)) \right\} \quad (11)$$

where  $p_\lambda$  is a penalty function (to be specified),  $\|\mathbf{u}_k\|_2 = 1$  and  $\|\mathbf{v}_k\|_2 = 1$ .

# Optimization Problem

## Penalty Function

- The penalty function used to encourage sparsity is

$$\begin{aligned} p_{\lambda}(\lambda_k, (d_k, \mathbf{u}_k, \mathbf{v}_k)) &= \lambda_k \sum_{i=1}^p \sum_{j=1}^K w_{ijk} |d_k u_{ik} v_{jk}| \\ &= \lambda_k (w_k^{(d)} d_k) \left( \sum_{i=1}^p w_{ik}^{(u)} |u_{ik}| \right) \left( \sum_{j=1}^K w_{jk}^{(v)} |v_{jk}| \right) \end{aligned} \quad (12)$$

- Penalizes each singular vector in the SVD layer, creating automatic adjustment of sparsity between  $\mathbf{u}_k$  and  $\mathbf{v}_k$ .
- $w_{ijk}$  is a weighting term (similar to the adaptive lasso (Zou 2006))

# Two Papers on Sparse Feature Selection in Multivariate Regression

└ Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)

- Weight term allows for flexibility in how much we penalize individual terms in each SVD layer.

- The penalty function used to encourage sparsity is

$$\begin{aligned} \rho_1(\lambda_1, (d_k, u_k, v_k)) &= \lambda_1 \sum_{i=1}^p \sum_{j=1}^K w_{ja} |d_k u_{ia} v_{ja}| \\ &= \lambda_1 (w_k^{(d)} d_k) \left( \sum_{i=1}^p w_{ia}^{(u)} |u_{ia}| \right) \left( \sum_{j=1}^K w_{ja}^{(v)} |v_{ja}| \right) \end{aligned} \quad (12)$$

- Penalizes each singular vector in the SVD layer, creating automatic adjustment of sparsity between  $u_k$  and  $v_k$
- $w_{ja}$  is a weighting term (similar to the adaptive lasso (Zou 2006))

# Adaptive Lasso (review)

- Recall from Zou (2006): adaptive lasso estimates  $\hat{\beta}^{*(n)}$  are given by

$$\hat{\beta}^{*(n)} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_n \sum_{j=1}^p \hat{w}_j |\beta_j| \right\} \quad (13)$$

where  $\hat{\mathbf{w}} = (|\hat{\beta}_1|^{-2}, \dots, |\hat{\beta}_p|^{-2})$  where  $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_p)$  is the OLS estimate for  $\beta$ .

- Zou shows the adaptive lasso is consistent in variable selection and that the asymptotic distribution of  $\hat{\beta}^{*(n)}$  is normal, unbiased, and efficient.

# Two Papers on Sparse Feature Selection in Multivariate Regression

- Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)

- Recall from Zou (2006): adaptive lasso estimates  $\tilde{\beta}^{(n)}$  are given by

$$\tilde{\beta}^{(n)} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_n \sum_{j=1}^p \hat{w}_j |\beta_j| \right\} \quad (13)$$

where  $\hat{\mathbf{w}} = (|\hat{\beta}_1|^{-2}, \dots, |\hat{\beta}_p|^{-2})$  where  $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_p)$  is the OLS estimate for  $\beta$ .

- Zou shows the adaptive lasso is consistent in variable selection and that the asymptotic distribution of  $\tilde{\beta}^{(n)}$  is normal, unbiased, and efficient.

- in the (univariate) linear model  $\mathbf{y} = \mathbf{X}\beta + \epsilon$
- (We will see that CCS prove similar results for their estimator.)



# Optimization Problem

## Penalty Function (continued)

- In the rank one case ( $r^* = 1$ , so we don't index over  $k$ ), Chen, Chan and Stenseth (CCS) choose the penalty weightings  $w^{(d)}$ ,  $\mathbf{w}^{(u)}$ ,  $\mathbf{w}^{(v)}$  in (12) as follows:

- 1 Estimate  $\mathbf{B}^*$  (call the estimate  $\tilde{\mathbf{B}}$ )
- 2 Find the SVD of  $\tilde{\mathbf{B}}$ :  $\tilde{\mathbf{B}} = \tilde{\mathbf{d}}\tilde{\mathbf{u}}\tilde{\mathbf{v}}^T$
- 3 Set

$$\begin{cases} w^{(d)} = |\tilde{d}|^{-2} \\ \mathbf{w}^{(u)} = (|\tilde{u}_1|^{-2}, \dots, |\tilde{u}_p|^{-2}) \\ \mathbf{w}^{(v)} = (|\tilde{v}_1|^{-2}, \dots, |\tilde{v}_K|^{-2}) \end{cases} \quad (14)$$

- CCS provide a (more complicated) generalization of this rule for the general case ( $r^* \in \mathbb{N}$ ).

# Two Papers on Sparse Feature Selection in Multivariate Regression

## Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)

- Note the similarity to the adaptive lasso weightings.
- Other choices besides -2 for the exponent are possible; CCS choose -2 based on simulation studies as well as the suggestion of Zou (2006).

- In the rank one case ( $r^* = 1$ , so we don't index over  $k$ ), Chen, Chan and Stenseth (CCS) choose the penalty weightings  $w^{(1)}$ ,  $w^{(2)}$ ,  $w^{(3)}$  in (12) as follows:

- Estimate  $\hat{B}^T$  (call the estimate  $\hat{B}$ )

- Find the SVD of  $\hat{B}$ :  $\hat{B} = U \Sigma V^T$

- Set

$$\begin{cases} w^{(1)} = |\hat{\lambda}|^{-2} \\ w^{(2)} = (|\hat{\lambda}_1|^{-2}, \dots, |\hat{\lambda}_p|^{-2}) \\ w^{(3)} = (|\hat{\lambda}_1|^{-2}, \dots, |\hat{\lambda}_p|^{-2}) \end{cases} \quad (14)$$

- CCS provide a (more complicated) generalization of this rule for the general case ( $r^* \in \mathbb{N}$ ).

# Exclusive Extraction Algorithm (EEA) for sparse SVD estimation of $B^*$

Basic idea: Start with an initial estimator  $\tilde{B} = \sum_{k=1}^{r^*} \tilde{d}_k \tilde{u}_k \tilde{v}_k^T$  for  $B^*$ , then compute  $r^*$  parallel sparse unit rank regression problems.

(a) For each  $k \in \{1, \dots, r^*\}$ :

(i) Construct the adaptive weights

$$w_k^{(d)} = |\tilde{d}_k|^{-2}, \mathbf{w}_k^{(u)} = (|\tilde{u}_{1k}|^{-2}, \dots, |\tilde{u}_{pk}|^{-2}), \text{ and} \\ \mathbf{w}_k^{(v)} = (|\tilde{v}_{1k}|^{-2}, \dots, |\tilde{v}_{Kk}|^{-2}).$$

(ii) Let  $\tilde{B}_k = \tilde{d}_k \tilde{u}_k \tilde{v}_k^T$ , and construct the exclusive layer  $\mathbf{Y}_k = \mathbf{Y} - \mathbf{X}(\tilde{B} - \tilde{B}_k)$ .

(iii) Find  $(\hat{d}_k, \hat{u}_k, \hat{v}_k)$  by performing the sparse unit rank regression of  $\mathbf{Y}_k$  on  $\mathbf{X}$ .

(b) The final estimator of  $B^*$  is given by  $\hat{B} = \sum_{k=1}^{r^*} \hat{d}_k \hat{u}_k \hat{v}_k^T$ .

# Two Papers on Sparse Feature Selection in Multivariate Regression

## Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)

Basic idea: Start with an initial estimator  $\tilde{B} = \sum_{i=1}^{r^*} \tilde{d}_i \tilde{u}_i \tilde{v}_i^T$  for  $B^*$ , then compute  $r^*$  parallel sparse unit rank regression problems.

- For each  $k \in \{1, \dots, r^*\}$ :
  - Construct the adaptive weights  $w_k^{(0)} = |\tilde{d}_k|^{-2}$ ,  $\mathbf{w}_k^{(0)} = (|d_{k1}|^{-2}, \dots, |d_{km}|^{-2})$ , and  $\mathbf{w}_k^{(0)} = (|d_{k1}|^{-2}, \dots, |d_{km}|^{-2})$ .
  - Let  $\tilde{B}_k = \tilde{d}_k \tilde{u}_k \tilde{v}_k^T$ , and construct the exclusive layer  $\mathbf{Y}_k = \mathbf{Y} - \mathbf{X}(\tilde{B} - \tilde{B}_k)$ .
  - Find  $(\tilde{d}_k, \tilde{u}_k, \tilde{v}_k)$  by performing the sparse unit rank regression of  $\mathbf{Y}_k$  on  $\mathbf{X}$ .
- The final estimator of  $B^*$  is given by  $\tilde{B} = \sum_{i=1}^{r^*} \tilde{d}_i \tilde{u}_i \tilde{v}_i^T$ .

- Note that this requires you to know  $r^*$ ; CCS show that this method has some robustness if  $r^*$  is misspecified.
- Typically the initial estimator is the reduced rank least squares estimator (ordinary regression but allowing for possibility that coefficient matrix is rank-deficient (e.g. two columns are identical))
- If instead of using  $\mathbf{Y}_k = \mathbf{Y} - \mathbf{X}(\tilde{B} - \tilde{B}_k)$  you use  $\mathbf{Y}_k = \mathbf{Y} - \mathbf{X}\tilde{B}$  the you get the sequential extraction algorithm (SEA) method which will come up later in simulations.
- (The optimal  $\lambda_k$  is usually chosen by a BIC specified by CCS, because cross-validation becomes infeasible due to the number of parameters to optimize over. CCS did simulations and found that results were similar with five-fold CV versus BIC.)
- Computational cost is linear in  $r^*$ ; estimation for different layers can be done in parallel

# Iterative Exclusive Extraction Algorithm (IEEA) and Orthogonality Relaxation

- Iterative exclusive extraction algorithm (IEEA): perform EEA once, then uses this as the initial estimate for another EEA iteration, repeating until the difference between estimates is sufficiently small.
- Recall that  $\mathbf{U}$  and  $\mathbf{V}$  are typically orthonormal in SVD. Relax the orthogonality condition to obtain sparsity in SVD.
- The estimators of different layers are consistent (under Theorem 2), so the estimators are “asymptotically orthogonal,” even though the solutions yielded are in general not exactly orthogonal.

# Two Papers on Sparse Feature Selection in Multivariate Regression

└ Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)

Iterative Exclusive Extraction Algorithm (IEEA) and Orthogonality Relaxation

- Iterative exclusive extraction algorithm (IEEA): perform EEA once, then uses this as the initial estimate for another EEA iteration, repeating until the difference between estimates is sufficiently small.
- Recall that  $U$  and  $V$  are typically orthonormal in SVD. Relax the orthogonality condition to obtain sparsity in SVD.
- The estimators of different layers are consistent (under Theorem 2), so the estimators are "asymptotically orthogonal," even though the solutions yielded are in general not exactly orthogonal.

- The optimization (11) is carried out locally near an initial consistent estimator of  $\mathbf{B}^*$  without enforcing exact orthogonality.
- The relaxation of the orthogonality requirement also improves local search efficiency.
- Usually only need a small number of additional iterations (less than 10) to get convergence (according to simulations even 1 or 2 gets you very close to optimal solution)

# Optimization Problem

## Optimization Algorithm

- The optimization problem (11) is nonconvex, but CCS developed an efficient parallelized coordinate descent optimization algorithm.
- Typically convergence occurs within only a few iterations.

- 1 Background on Multivariate Regression
- 2 Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)
  - Background and Problem Statement
  - Main Results
  - Selected Simulation Studies
- 3 Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)
  - Background and Problem Statement
  - Optimization Problem
  - **Main Results**
  - Selected Simulation Study



- (Summary of Theorem 1) Under very mild assumptions, the optimization problem (11) solved using the IEEA has a local minimum with  $\sqrt{n}$  convergence.
- (Summary of Theorem 2) Under very mild assumptions, the asymptotic distributions of the IEEA estimators of the nonzero elements of  $\mathbf{U}$  and  $\mathbf{V}$  are normal with zero mean. Further, the estimators are consistent.
- (Summary of Theorem 3) Under very mild assumptions, the IEEA estimators are selection consistent; that is, the elements of  $\mathbf{U}$  and  $\mathbf{V}$  that equal 0 will be set to 0 by IEEA with probability 1 as the number of iterations tends to infinity.

# Two Papers on Sparse Feature Selection in Multivariate Regression

## └ Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)

### Theorems

- (Summary of Theorem 1) Under very mild assumptions, the optimization problem (11) solved using the IEEA has a local minimum with  $\sqrt{n}$  convergence.
- (Summary of Theorem 2) Under very mild assumptions, the asymptotic distributions of the IEEA estimators of the nonzero elements of  $\mathbf{U}$  and  $\mathbf{V}$  are normal with zero mean. Further, the estimators are consistent.
- (Summary of Theorem 3) Under very mild assumptions, the IEEA estimators are selection consistent; that is, the elements of  $\mathbf{U}$  and  $\mathbf{V}$  that equal 0 will be set to 0 by IEEA with probability 1 as the number of iterations tends to infinity.

- Note the similarity to the adaptive lasso results discussed earlier

- 1 Background on Multivariate Regression
- 2 Support Union Recovery in High-Dimensional Multivariate Regression (Obozinski, Wainwright, and Jordan 2011)
  - Background and Problem Statement
  - Main Results
  - Selected Simulation Studies
- 3 Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)
  - Background and Problem Statement
  - Optimization Problem
  - Main Results
  - Selected Simulation Study

**Table 3.** Estimation and prediction accuracy of the IEEA, SEA, OLS, RRR and NNP estimators

<i>Model</i>	<i>SNR</i>	<i>Error</i>	<i>Results for the following methods:</i>				
			<i>IEEA</i>	<i>SEA</i>	<i>OLS</i>	<i>RRR</i>	<i>NNP</i>
I	0.125	Er-Est	3.29	4.28	55.32	9.90	10.99
		Er-Pred	52.10	67.37	416.34	106.40	139.10
	0.25	Er-Est	1.10	1.52	27.21	4.56	6.65
		Er-Pred	17.68	23.20	197.12	46.16	74.55
	0.5	Er-Est	0.57	0.99	15.12	2.44	3.75
		Er-Pred	9.48	15.03	110.45	25.89	43.74
II	1	Er-Est	0.23	0.52	7.41	1.17	2.19
		Er-Pred	4.01	7.59	52.95	11.87	23.23
	0.125	Er-Est	0.52	0.51	51.50	4.87	5.01
		Er-Pred	12.86	14.84	342.59	60.12	72.00
	0.25	Er-Est	0.15	0.29	28.09	3.89	4.22
		Er-Pred	4.47	7.85	176.22	23.64	45.44
	0.5	Er-Est	0.06	0.17	14.81	3.23	3.67
		Er-Pred	1.80	4.27	84.94	9.67	26.87
	1	Er-Est	0.03	0.10	8.40	2.75	3.13
		Er-Pred	0.84	2.35	42.57	4.55	15.43

# Two Papers on Sparse Feature Selection in Multivariate Regression

Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)

Table 2. Estimation and prediction accuracy of the IEEA, SEA, OLS, RRR and NNP estimators

Model	SNR	Error	Results for the following methods:				
			IEEA	SEA	OLS	RRR	NNP
I	0.125	Ex-Est	3.29	4.28	55.32	9.90	10.99
		Ex-Pred	52.10	67.37	416.34	186.48	139.10
		Ex-Inv	1.00	1.52	27.21	4.56	6.65
	0.25	Ex-Est	17.68	23.28	197.12	46.34	74.55
		Ex-Pred	0.57	0.89	15.12	2.44	3.75
		Ex-Inv	9.40	15.83	110.45	25.89	43.74
II	1	Ex-Est	0.23	0.52	7.41	1.37	2.18
		Ex-Pred	6.01	7.59	52.95	11.87	25.35
		Ex-Inv	0.52	0.51	51.50	4.87	5.01
	0.125	Ex-Est	12.86	14.84	362.59	60.12	77.00
		Ex-Pred	0.15	0.28	28.09	3.89	4.22
		Ex-Inv	6.47	7.15	176.22	23.64	46.46
	0.5	Ex-Est	0.86	0.17	14.81	3.23	3.67
		Ex-Pred	1.80	4.27	84.94	9.67	20.87
		Ex-Inv	0.83	0.19	8.40	2.75	3.15
	1	Ex-Est	0.84	2.35	42.57	4.35	15.45
		Ex-Pred					
		Ex-Inv					

- IEEA: proposed method; regularization parameter chosen using BIC.
- SEA: simplification of EEA. In SEA: sequentially perform sparse unit rank regression, each time with data matrix  $Y$  replaced by residual matrix  $Y - X B$ . Has been used in many penalized matrix decomposition problems. Correspond to a different decomposition of the coefficient matrix other than SVD, and need not produce SVD layers of  $B$ , so it is not suitable for recovering the desired sparse SVD structure in  $B$ . Regularization parameter chosen using BIC.
- OLS: least squares
- RRR: reduced rank regression (ordinary regression but allowing for possibility that coefficient matrix is rank-deficient (e.g. two columns are identical))
- NNP: nuclear-norm penalized regression (nuclear norm is sum of absolute values of eigenvalues)

# Two Papers on Sparse Feature Selection in Multivariate Regression

Reduced Rank Stochastic Regression with a Sparse Singular Value Decomposition (Chen, Chan, and Stenseth 2012)

Table 3. Estimation and prediction accuracy of the IEEA, SEA, OLS, RRR and NNP estimators

Model	SNR	Error	Results for the following methods:				
			IEEA	SEA	OLS	RRR	NNP
I	0.125	Er-Est	3.29	4.28	55.32	9.98	10.99
		Er-Pred	52.09	67.37	416.34	186.48	139.10
	0.25	Er-Est	1.90	1.52	27.21	4.56	6.65
		Er-Pred	17.68	23.28	197.12	46.16	74.55
	0.5	Er-Est	0.57	0.89	15.12	2.44	3.75
		Er-Pred	9.40	15.83	110.45	25.89	43.74
II	1	Er-Est	0.21	0.52	7.40	1.37	2.18
		Er-Pred	6.01	7.94	52.95	11.87	25.35
	0.125	Er-Est	0.52	0.51	51.90	4.87	5.01
		Er-Pred	12.86	14.84	362.59	60.12	77.09
	0.25	Er-Est	0.15	0.28	28.09	1.89	4.22
		Er-Pred	6.47	7.15	178.22	23.64	46.46
I	0.5	Er-Est	0.06	0.17	14.81	3.23	3.67
		Er-Pred	1.80	4.27	84.94	9.67	26.87
		Er-Pred	0.83	0.39	6.40	2.79	3.15
I	1	Er-Est	0.84	2.35	42.57	4.35	15.45

- Model I:  $p = K = 25, n = 50, r^* = 3$ .
- Model II:  $p = K = 60, n = 50, r^* = 3$  (same but more noise features, and more responses even though their columns of  $B$  don't add to rank of matrix).
- Then generate data via  $Y = XB + W$  with  $\sigma$  chosen according to desired SNR.
- Vary SNR as well as estimation and prediction error. Er-Est:  $\|B^* - \hat{B}\|_F^2 / (pq)$ . Er-Pred:  $\|XB^* - X\hat{B}\|_F^2 / (pq)$ .
- We see that authors are correct that IEEA pretty much uniformly outperforms SEA, although this is the second best method. (Sparsity assumption benefits method when model dimension is high and number of irrelevant responses or predictors is large.)
- Next RRR, then NNP, finally OLS.