

# High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity

Po-Ling Loh

UC Berkeley

NIPS 2011

December 13, 2011

Joint work with Martin Wainwright

- High-dimensional problems: # parameters  $p \gg$  # observations  $n$
- Numerous applications in science and engineering
  - DNA microarray analysis
  - Health studies, longitudinal analysis
  - Portfolio optimization
  - Compressed sensing, MRI/fMRI
  - Face recognition, spam filtering, astronomy, climatology ...
  - $p \approx 10,000, n \approx 100$

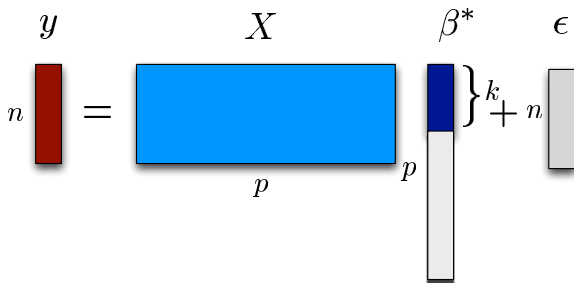
# Sparse linear regression

$$\begin{matrix} y \\ n \end{matrix} = \begin{matrix} X \\ p \end{matrix} \begin{matrix} \beta^* \\ p \end{matrix} + \begin{matrix} \epsilon \\ n \end{matrix}$$

- Linear model:

$$y_i = x_i^T \beta^* + \epsilon_i, \quad i = 1, \dots, n$$

# Sparse linear regression



- Linear model:

$$y_i = x_i^T \beta^* + \epsilon_i, \quad i = 1, \dots, n$$

- When  $p \gg n$ , assume sparsity:  $\|\beta^*\|_0 \leq k$

- Additional complications when  $Z$  observed in place of  $X$

# Corrupted variables

- Additional complications when  $Z$  observed in place of  $X$
- **Additive noise:**  $Z = X + W$ , where  $X \perp\!\!\!\perp W$  and  $\Sigma_w$  is known

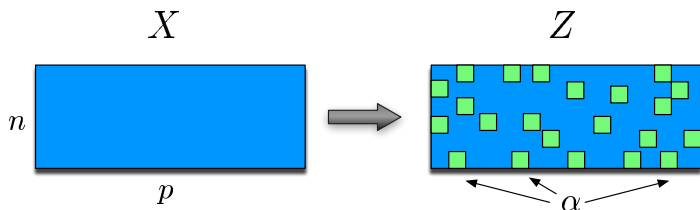
$$Z = \begin{matrix} & & X & & \\ & & & & W \\ n & \text{[Blue Box]} & + & \text{[Noisy Grid]} & \\ & & & & \\ & & & & p \end{matrix}$$

The diagram illustrates the equation  $Z = X + W$ . On the left, the variable  $Z$  is shown with a vertical dimension of  $n$ . The variable  $X$  is represented by a solid blue rectangle with a horizontal dimension of  $p$ . A plus sign follows. The variable  $W$  is represented by a 10x10 grid of squares in various shades of gray, representing noise. The grid is positioned to the right of the plus sign, and its horizontal dimension is labeled as  $p$  below it.

- **Ex:** Medical or experimental data, portfolio optimization

# Corrupted variables

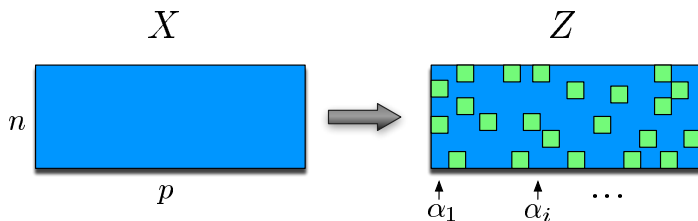
- Additional complications when  $Z$  observed in place of  $X$
- **Missing data:** entries of  $X$  missing independently with probability  $\alpha$



- **Ex:** Voting records, survey data, broken sensor arrays

# Corrupted variables

- Additional complications when  $Z$  observed in place of  $X$
- **Missing data:** entries of  $X$  missing independently with probability  $\alpha$

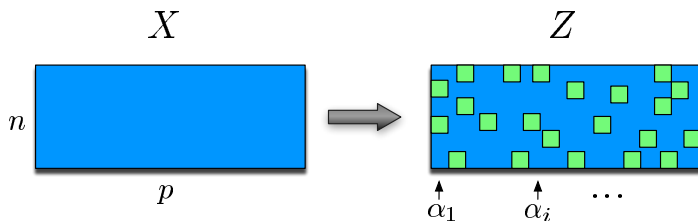


- Each column may have separate probability  $\alpha_i$  of missing entries



# Corrupted variables

- Additional complications when  $Z$  observed in place of  $X$
- **Missing data:** entries of  $X$  missing independently with probability  $\alpha$

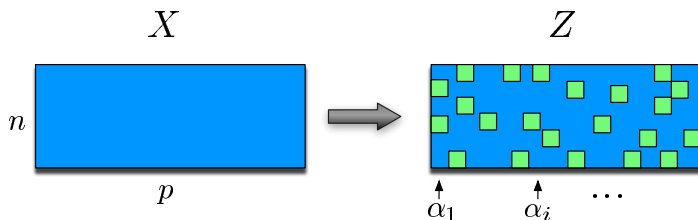


- **Model:**

$$y = X\beta^* + \epsilon$$

# Corrupted variables

- Additional complications when  $Z$  observed in place of  $X$
- **Missing data:** entries of  $X$  missing independently with probability  $\alpha$



- **Unlike EM methods, our method converges to a near-global optimum despite non-convexity**

- Note that

$$\beta^* \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \beta^T \Sigma_x \beta - \beta^{*T} \Sigma_x \beta \right\}$$

- Note that

$$\beta^* \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \beta^T \Sigma_x \beta - \beta^{*T} \Sigma_x \beta \right\}$$

- Compare to Lasso (Tibshirani '96):

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \left\{ \frac{1}{2n} \|y - X\beta\|_2^2 \right\}$$

- Note that

$$\beta^* \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \beta^T \Sigma_x \beta - \beta^{*T} \Sigma_x \beta \right\}$$

- Compare to Lasso (Tibshirani '96):

$$\begin{aligned} \hat{\beta} &\in \arg \min_{\|\beta\|_1 \leq R} \left\{ \frac{1}{2n} \|y - X\beta\|_2^2 \right\} \\ &= \arg \min_{\|\beta\|_1 \leq R} \left\{ \frac{1}{2} \beta^T \left( \frac{X^T X}{n} \right) \beta - \frac{y^T X}{n} \beta \right\} \end{aligned}$$

- Note that

$$\beta^* \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2} \beta^T \Sigma_x \beta - \beta^{*T} \Sigma_x \beta \right\}$$

- Idea:** form unbiased estimators  $(\hat{\Gamma}, \hat{\gamma})$  of  $(\Sigma_x, \text{Cov}(X, y))$  based on  $(y, Z)$ , solve constrained program

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \left\{ \frac{1}{2} \beta^T \hat{\Gamma} \beta - \hat{\gamma}^T \beta \right\}$$

## Example: Additive noise

- Since  $Z = X + W$  and  $X \perp\!\!\!\perp W$ , we have  $\Sigma_Z = \Sigma_X + \Sigma_W$  and  $\text{Cov}(y, X) = \text{Cov}(y, Z)$

## Example: Additive noise

- Since  $Z = X + W$  and  $X \perp\!\!\!\perp W$ , we have  $\Sigma_Z = \Sigma_X + \Sigma_W$  and  $\text{Cov}(y, X) = \text{Cov}(y, Z)$
- Use

$$\hat{\beta} = \frac{Z^T Z}{n} - \Sigma_W, \quad \hat{\gamma} = \frac{Z^T y}{n}$$



## Example: Additive noise

- Since  $Z = X + W$  and  $X \perp\!\!\!\perp W$ , we have  $\Sigma_Z = \Sigma_X + \Sigma_W$  and  $\text{Cov}(y, X) = \text{Cov}(y, Z)$
- Use

$$\hat{\Gamma} = \frac{Z^T Z}{n} - \Sigma_W, \quad \hat{\gamma} = \frac{Z^T y}{n}$$

- Objective:

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \left\{ \frac{1}{2} \beta^T \left( \frac{Z^T Z}{n} - \Sigma_W \right) \beta - \frac{y^T Z}{n} \beta \right\}$$

## Example: Missing data

- $X \sim N(0, \Sigma_x)$ ,  $Z \in \mathbb{R}^{n \times p}$  is observed data with probability  $\alpha$  of missing values

## Example: Missing data

- $X \sim N(0, \Sigma_x)$ ,  $Z \in \mathbb{R}^{n \times p}$  is observed data with probability  $\alpha$  of missing values

- Let

$$\hat{Z}_{ij} = \begin{cases} \frac{Z_{ij}}{1-\alpha} & \text{if } Z_{ij} \text{ is observed} \\ 0 & \text{otherwise} \end{cases}$$

- Then

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

satisfies  $\mathbb{E}(\hat{\Gamma}) = \Sigma_x$  and  $\text{Cov}(\hat{Z}, y) = \text{Cov}(X, y)$

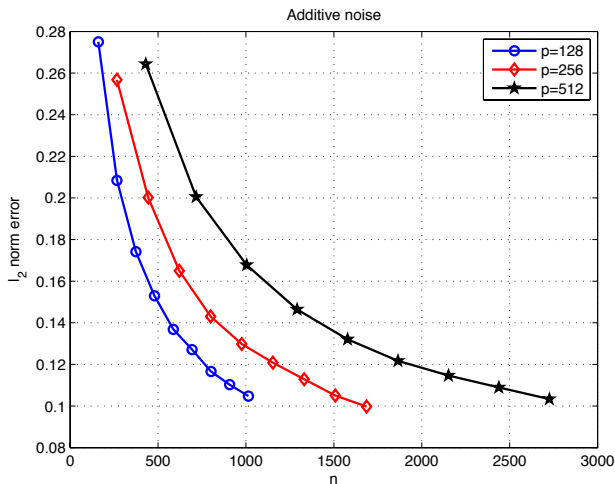
- Objective:

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \left\{ \frac{1}{2} \beta^T \hat{\Gamma} \beta - \frac{y^T \hat{Z}}{n} \beta \right\}$$

# High-dimensional consistency?

# High-dimensional consistency?

- Modified Lasso with additive noise,  $k \approx \sqrt{p}$
- Consistency:  $\|\hat{\beta} - \beta^*\|_2 \rightarrow 0$  as  $n \rightarrow \infty$



- Under restricted eigenvalue conditions on  $X$  (Bickel, Ritov & Tsybakov '08, van de Geer & Bühlmann '09),

$$\|\hat{\beta} - \beta^*\|_1 = \mathcal{O}\left(k\sqrt{\frac{\log p}{n}}\right), \quad \|\hat{\beta} - \beta^*\|_2 = \mathcal{O}\left(\sqrt{\frac{k \log p}{n}}\right)$$

- RE conditions hold w.h.p. when  $X$  is a random matrix with rows sampled i.i.d. from a (sub)-Gaussian distribution (Raskutti et al. '09)

## Theorem (Statistical error)

Under modified RE condition  $\hat{\Gamma}$  and deviation conditions on  $(\hat{\gamma}, \hat{\Gamma})$ , any global optimum  $\hat{\beta}$  satisfies

$$\|\hat{\beta} - \beta^*\|_1 \lesssim \varphi(\sigma_\epsilon) \left( k \sqrt{\frac{\log p}{n}} \right), \quad \|\hat{\beta} - \beta^*\|_2 \lesssim \varphi(\sigma_\epsilon) \left( \sqrt{\frac{k \log p}{n}} \right)$$

- Deviation conditions:

$$\|\hat{\gamma} - \text{Cov}(X, y)\|_\infty, \quad \|(\hat{\Gamma} - \Sigma_x)\beta^*\|_\infty \lesssim \varphi(\sigma_\epsilon) \left( \sqrt{\frac{\log p}{n}} \right)$$

## Theorem (Statistical error)

Under modified RE condition  $\hat{\Gamma}$  and deviation conditions on  $(\hat{\gamma}, \hat{\Gamma})$ , any global optimum  $\hat{\beta}$  satisfies

$$\|\hat{\beta} - \beta^*\|_1 \lesssim \varphi(\sigma_\epsilon) \left( k \sqrt{\frac{\log p}{n}} \right), \quad \|\hat{\beta} - \beta^*\|_2 \lesssim \varphi(\sigma_\epsilon) \left( \sqrt{\frac{k \log p}{n}} \right)$$

- Deviation conditions:

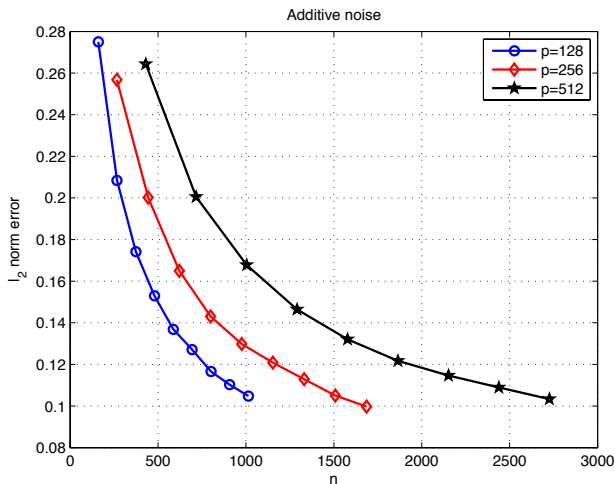
$$\|\hat{\gamma} - \text{Cov}(X, y)\|_\infty, \quad \|(\hat{\Gamma} - \Sigma_x)\beta^*\|_\infty \lesssim \varphi(\sigma_\epsilon) \left( \sqrt{\frac{\log p}{n}} \right)$$

- $\varphi(\sigma_\epsilon)$  is a function of corruption pattern and noise variance, decreases with SNR and increases with  $\alpha$



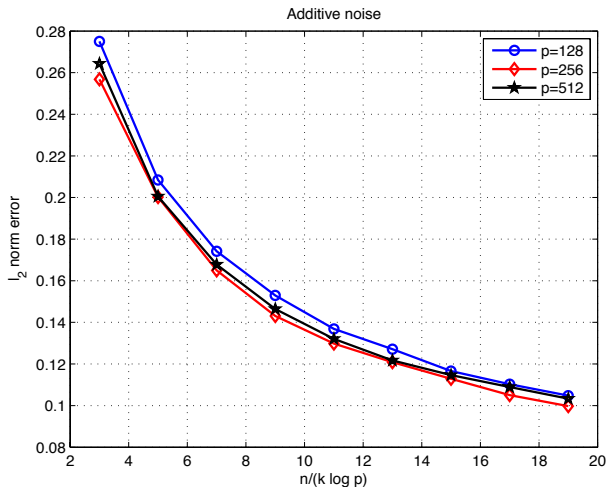
# High-dimensional consistency?

- Modified Lasso with additive noise,  $k \approx \sqrt{p}$
- Consistency:  $\|\hat{\beta} - \beta^*\|_2 \rightarrow 0$  as  $n \rightarrow \infty$



# High-dimensional consistency?

- $\ell_2$ -error vs. rescaled sample size  $n/(k \log p)$
- Curves stack up, verifying theoretical results



- Corrected objective is **not** convex

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \left\{ \frac{1}{2} \beta^T \left( \frac{Z^T Z}{n} - \Sigma_w \right) \beta - \frac{y^T Z}{n} \beta \right\}$$

- Hessian has at least  $p - n$  negative eigenvalues

- Corrected objective is **not** convex

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \left\{ \frac{1}{2} \beta^T \left( \frac{Z^T Z}{n} - \Sigma_w \right) \beta - \frac{y^T Z}{n} \beta \right\}$$

- Hessian has at least  $p - n$  negative eigenvalues
- **Pretend objective is convex, apply projected gradient descent algorithm**

# Projected gradient descent

- Solve constrained optimization problem

$$\min_{\beta} \underbrace{\frac{1}{2n} \|y - X\beta\|_2^2}_{\mathcal{L}(\beta)} \quad \text{s.t.} \quad \|\beta\|_1 \leq R$$

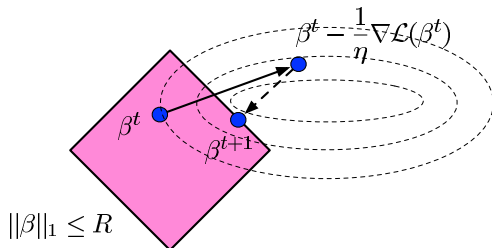
# Projected gradient descent

- Solve constrained optimization problem

$$\min_{\beta} \underbrace{\frac{1}{2n} \|y - X\beta\|_2^2}_{\mathcal{L}(\beta)} \quad \text{s.t.} \quad \|\beta\|_1 \leq R$$

- Produces iterates  $\beta^t$  with

$$\beta^{t+1} = \Pi \left( \beta^t - \frac{1}{\eta} \nabla \mathcal{L}(\beta^t) \right), \quad \text{stepsize } \eta > 0$$



- Linear convergence when  $\mathcal{L}$  is smooth and strongly convex (Bertsekas '95):

$$\|\beta^t - \hat{\beta}\|_2 \leq \gamma^t \|\beta^0 - \hat{\beta}\|_2$$

- Linear convergence when  $\mathcal{L}$  is smooth and strongly convex (Bertsekas '95):

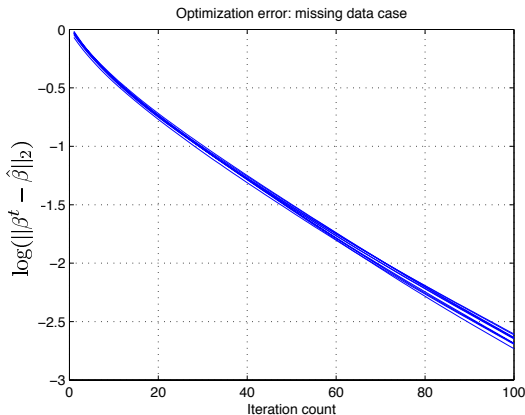
$$\|\beta^t - \hat{\beta}\|_2 \leq \gamma^t \|\beta^0 - \hat{\beta}\|_2$$

- **When  $\mathcal{L}$  non-convex, projected gradient descent may not converge, or converge at slower rates**



# Global linear convergence observed in practice

- For fixed problem instance, 10 runs of projected gradient descent, plotted optimization error  $\|\beta^t - \hat{\beta}\|_2$
- $p = 512$ ,  $k \approx \sqrt{p}$ ,  $n \approx 5k \log p$



## Theorem (Optimization error)

*For the modified Lasso,*

$$\|\beta^t - \hat{\beta}\|_2 \leq \gamma^t \|\beta^0 - \hat{\beta}\|_2 + o\left(\sqrt{\frac{k \log p}{n}}\right)$$

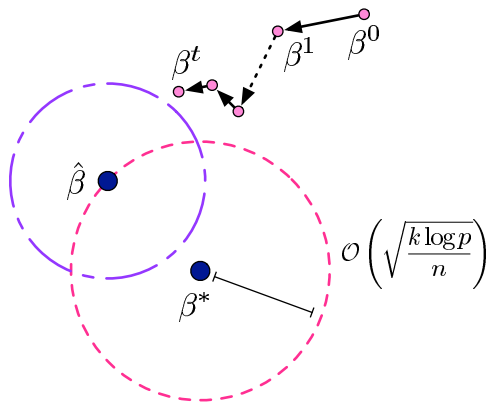
## Theorem (Optimization error)

For the modified Lasso,

$$\|\beta^t - \hat{\beta}\|_2 \leq \gamma^t \|\beta^0 - \hat{\beta}\|_2 + o\left(\sqrt{\frac{k \log p}{n}}\right)$$

- Use results from Agarwal, Negahban & Wainwright (NIPS '10), applied to non-convex objective
- Requires restricted strong convexity (RSC) and restricted smoothness (RSM), holding w.h.p. in settings of interest

# Illustration of statistical and optimization error

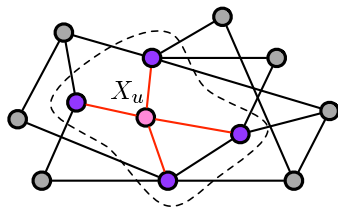
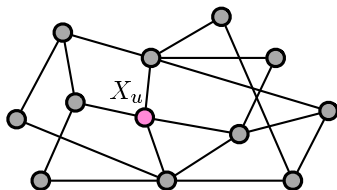


- Statistical error:  $\|\hat{\beta} - \beta^*\|_2 = \mathcal{O}\left(\sqrt{\frac{k \log p}{n}}\right)$
- Optimization error:  $\|\beta^t - \hat{\beta}\|_2 = \gamma^t \|\beta^0 - \hat{\beta}\|_2 + o\left(\sqrt{\frac{k \log p}{n}}\right)$

# Application: Gaussian graphical models

- Conditional independence property for graphical model:

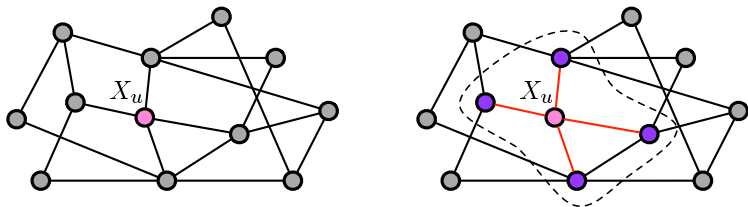
$$X_u \mid X_{V \setminus \{u\}} \stackrel{d}{=} X_u \mid X_{N(u)}$$



# Application: Gaussian graphical models

- Conditional independence property for graphical model:

$$X_u \mid X_{V \setminus \{u\}} \stackrel{d}{=} X_u \mid X_{N(u)}$$



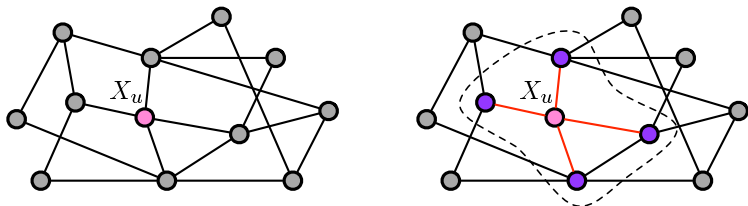
- When  $X \sim N(0, \Sigma)$ , entries of  $\Theta = \Sigma^{-1}$  may be recovered via nodewise linear regression (Meinshausen and Bühlmann '06, Yuan '10)

sparsity  $k \iff$  max degree of vertex

# Application: Gaussian graphical models

- Conditional independence property for graphical model:

$$X_u \mid X_{V \setminus \{u\}} \stackrel{d}{=} X_u \mid X_{N(u)}$$



- When  $X \sim N(0, \Sigma)$ , entries of  $\Theta = \Sigma^{-1}$  may be recovered via nodewise linear regression (Meinshausen and Bühlmann '06, Yuan '10)

sparsity  $k \iff$  max degree of vertex

- For corrupted observations, use noisy regression to recover  $\Theta$

## Theorem (Spectral norm consistency)

For estimate  $\hat{\Theta}$  based on corrupted observations of a Gaussian graphical model,

$$\|\hat{\Theta} - \Theta\|_{op} = \mathcal{O}\left(k\sqrt{\frac{\log p}{n}}\right)$$

- Matches rates for fully-observed case



- Provided a Lasso variant based on noisy observations  $(y, Z)$ , such that

$$\|\hat{\beta} - \beta^*\|_2 = \mathcal{O}\left(\sqrt{\frac{k \log p}{n}}\right)$$

- Derived an estimator for the inverse covariance matrix of a (noisy) Gaussian graphical model, such that

$$\|\hat{\Theta} - \Theta\|_{\text{op}} = \mathcal{O}\left(k\sqrt{\frac{\log p}{n}}\right)$$

- Demonstrated that global minimizers  $\hat{\beta}$  for the **non-convex** objective can be obtained via projected gradient descent

- Support recovery for corrupted observations
- Minimax lower bounds
- Additive noise model with unknown  $\Sigma_w$
- Other corruption patterns: multiplicative noise, censored data

- Additive noise:  $X_i \sim N(0, \sigma_x^2 I)$ ,  $W_i \sim N(0, \sigma_w^2 I)$ :

$$\varphi = \sqrt{1 + \frac{\sigma_w^2}{\sigma_x^2}} \sqrt{\frac{\sigma_w^2}{\sigma_x^2} + \frac{\sigma_\epsilon^2}{\sigma_x^2}}$$

- Missing data:  $X_i \sim N(0, \sigma_x^2 I)$ ,

$$\varphi = \frac{\sigma_\epsilon}{\sigma_x(1 - \alpha)} + \frac{1}{(1 - \alpha)^2}$$

## Algorithm:

- Perform  $p$  linear regressions of the variables  $Z^i$  upon the remaining variables  $Z^{-i}$ , using the modified Lasso program with estimators  $(\hat{\Gamma}^{(i)}, \hat{\gamma}^{(i)})$
- Estimate scalars  $a_i$  using plug-in estimator  $\hat{a}_i = -(\hat{\Gamma}_{ii} - \hat{\Gamma}_{i,-i}\hat{\theta}^i)^{-1}$
- Form the matrix  $\tilde{\Theta}$  with  $\tilde{\Theta}_{i,-i} = \hat{a}_i\hat{\theta}^i$  and  $\tilde{\Theta}_{ii} = -\hat{a}_i$
- Symmetrize:  $\hat{\Theta} \in \arg \min_{\Theta \in S^p} \left\| \Theta - \tilde{\Theta} \right\|_{\ell_1 \rightarrow \ell_1}$
  
- Last step is an LP, can be optimized with standard techniques

## Some references

- A. Agarwal, S. Negahban, and M.J. Wainwright (2011). Fast global convergence of gradient methods for high-dimensional statistical recovery. ArXiv paper.
- P. Loh and M.J. Wainwright (2011). High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity. ArXiv paper.
- N. Meinshausen and P. Bühlmann (2006). High-dimensional graphs and variable selection with the Lasso. *Annals of Statistics*.
- G. Raskutti, M.J. Wainwright, and B. Yu (2010). Restricted eigenvalue properties for correlated Gaussian random designs. *Journal of Machine Learning Research*.
- M. Rosenbaum and A.B. Tsybakov (2010). Sparse recovery under matrix uncertainty. *Annals of Statistics*.
- R. Tibshirani (1996). Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*.
- M. Yuan (2010). High-dimensional inverse covariance matrix estimation via linear programming. *Journal of Machine Learning Research*.