

# Fused MCP with applications in signal processing

Xianshi Yu

*Joint work with Bing-yi Jing, Guangren Yang and Cun-Hui Zhang*

Department of Mathematics, HKUST

Jan 2015

# Outline

- 1 Review
- 2 Signal Denoising and FLSA
- 3 Method
- 4 Algorithms
- 5 Oracle Property
- 6 Simulations
- 7 Summary

# linear regression model

- Given  $(y_i, x_{i1}, \dots, x_{ip})$ ,  $i = 1, \dots, n$ , assume

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i, \quad i = 1, \dots, n.$$

- In matrix form,

$$\mathbf{Y} = \mathbf{1}\beta_0 + \mathbf{X}\beta + \epsilon$$

where

$$\mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ x_{21} & \cdots & x_{2p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{bmatrix}, \quad \epsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix},$$

# Least square estimate (LSE)

- LSE:

$$\hat{\beta}^{lse} = \arg \min_{\beta} \left\{ \sum_i (y_i - \sum_j x_{ij} \beta_j)^2 \right\} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \mathbf{Y}$$

- LSE is ill-posed if  $p > n$
- In some situations, we can impose some assumptions on  $\beta$

# Lasso

- Lasso (Tibshirani (1996)):

$$\hat{\beta}^{lasso} = \arg \min_{\beta} \frac{1}{2} \left\{ \sum_i (y_i - \sum_j x_{ij} \beta_j)^2 \right\} + \lambda \sum_{j=1}^p |\beta_j|$$

- When  $X^T X = I$ , then

$$\hat{\beta}^{lasso} = \text{sgn}(\hat{\beta}^{lse}) \left( |\hat{\beta}^{lse}| - \lambda \right)^+$$

- It works for  $p > n$  as well.
- $L_1$  penalty imposes sparsity on  $\beta$
- It does shrinkage and variable selection simultaneously.

# Fused Lasso

- In some situations, the features have an inherent order.
- Examples:
  - protein mass spectroscopy data
  - gene expression data
- Fused Lasso (Tibshirani (2005)) minimizes:

$$\frac{1}{2} \left\{ \sum_i (y_i - \sum_j x_{ij} \beta_j)^2 \right\} + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=2}^p |\beta_j - \beta_{j-1}|$$

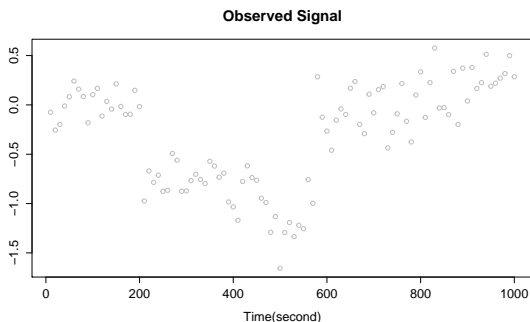
- It encourages both sparsity and local constancy in  $\beta$

# Outline

- 1 Review
- 2 Signal Denoising and FLSA**
- 3 Method
- 4 Algorithms
- 5 Oracle Property
- 6 Simulations
- 7 Summary

# Signal Denoising Problem

- In signal processing scenario,  $y_i = \beta_i + \epsilon_i$



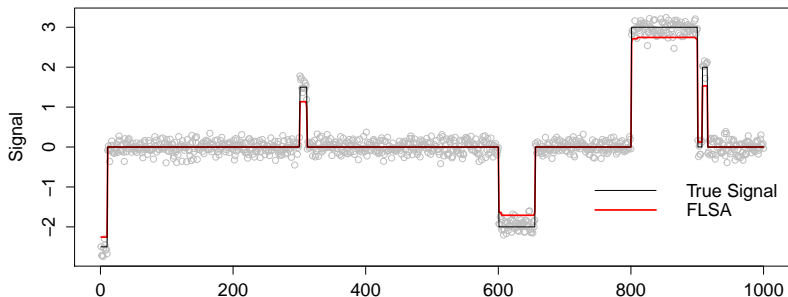
- Fusion penalty  $\sum_{j=2}^p |\beta_j - \beta_{j-1}|$  can be applied to denoise the corrupted signal



# FLSA

Friedman(2007) introduced FLSA:

$$\text{Minimize } \frac{1}{2} \|\mathbf{Y} - \beta\|_2^2 + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=2}^p |\beta_j - \beta_{j-1}|.$$



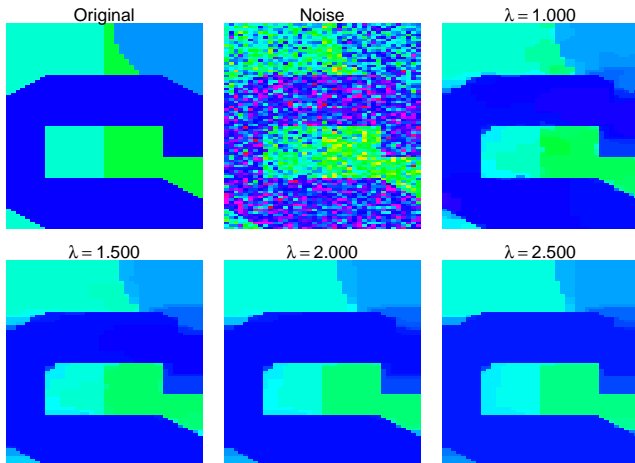
# FLSA

FLSA is the discrete analogue of *total variation denoising* established by Rudin *et al.* (1992).

$$\min_u \int_{\Omega} |\nabla u| \, du \quad \text{subject to} \quad \|u - y\|^2 = \sigma^2$$

# 2-d FLSA

$$\text{Minimize } \frac{1}{2} \|\mathbf{Y} - \beta\|_F^2 + \lambda \sum |\beta_{i,j} - \beta_{i,j-1}| + \lambda \sum |\beta_{i,j} - \beta_{i-1,j}|$$



# Comments on FLSA

- FLSA captures the profile of the signals
- Contrast in signal is shrunked
- Jump points in 1-d recovered signal is not clear
- Edges in recovered image is not clear

# Outline

- 1 Review
- 2 Signal Denoising and FLSA
- 3 Method**
- 4 Algorithms
- 5 Oracle Property
- 6 Simulations
- 7 Summary

# Mcp penalty

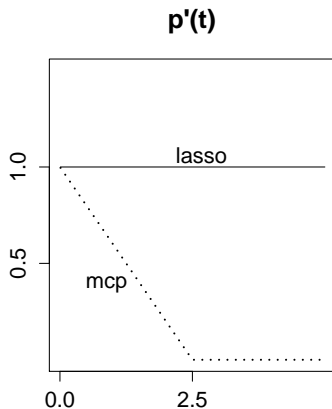
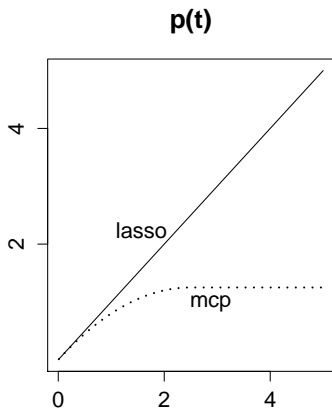
Mcp (Zhang (2010)):

$$\hat{\beta}^{mcp} = \arg \min_{\beta} \frac{1}{2} \left\{ \sum_i (y_i - \sum_j x_{ij} \beta_j)^2 \right\} + \sum_{j=1}^p \rho(|\beta_j|; \lambda)$$

where  $\rho(t; \lambda) = \lambda \int_0^t (1 - x/(\gamma\lambda))_+ dx$

# Mcp penalty

## Comparison of Mcp and Lasso



# Mcp penalty

- $\rho(t; \lambda)$  is nonconvex
- Mcp introduces sparsity
- Mcp dose not shrink  $\beta$  when it is large
- Mcp is nearly unbiased
- Computation: PLUS—a path algorithm



# Fused Mcp

- 1-d Fused Mcp minimizes

$$\frac{1}{2} \sum_i (y_i - \beta_j)^2 + \sum_{j=1}^p \rho(|\beta_j|; \lambda_1) + \sum_{j=2}^p \rho(|\beta_j - \beta_{j-1}|; \lambda_2)$$

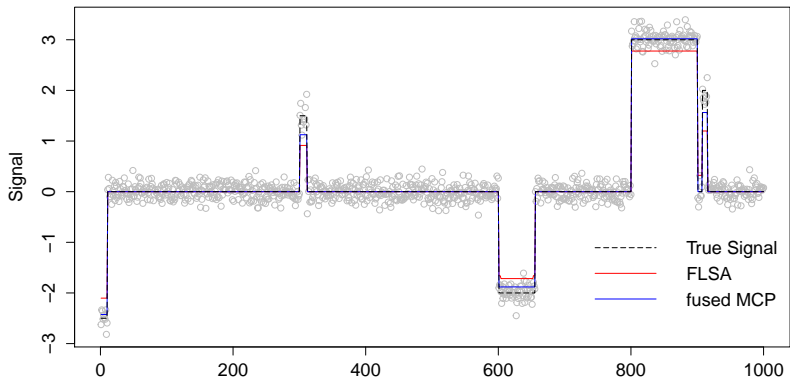
- Specifically, when  $\lambda_1 = 0$ , it is called fusion MCP.

$$\text{Minimize } \frac{1}{2} \sum_i (y_i - \beta_j)^2 + \sum_{j=2}^p \rho(|\beta_j - \beta_{j-1}|; \lambda)$$

- 2-d Fused Mcp minimizes

$$\frac{1}{2} \sum_i (y_i - \beta_j)^2 + \sum \rho(|\beta_{i,j} - \beta_{i,j-1}|; \lambda) + \sum \rho(|\beta_{i,j} - \beta_{i-1,j}|; \lambda)$$

# Fused MCP vs. FLSA

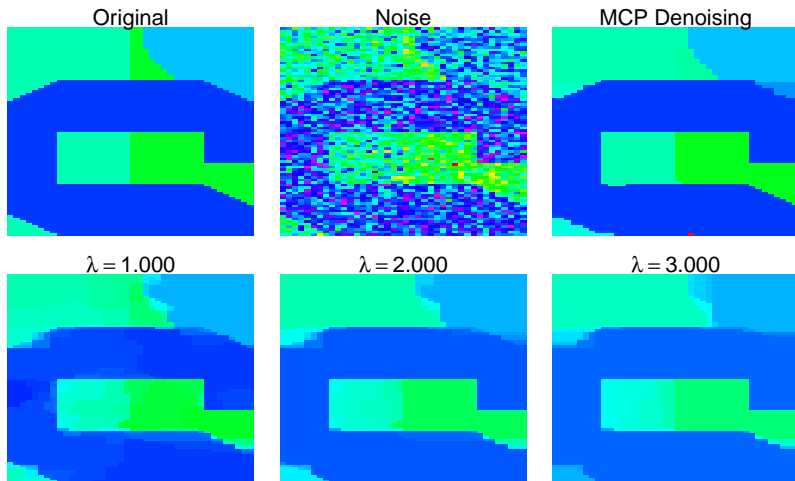


# Fused MCP vs. FLSA

Fused MCP has better performance in jump point and jump size detection than FLSA.

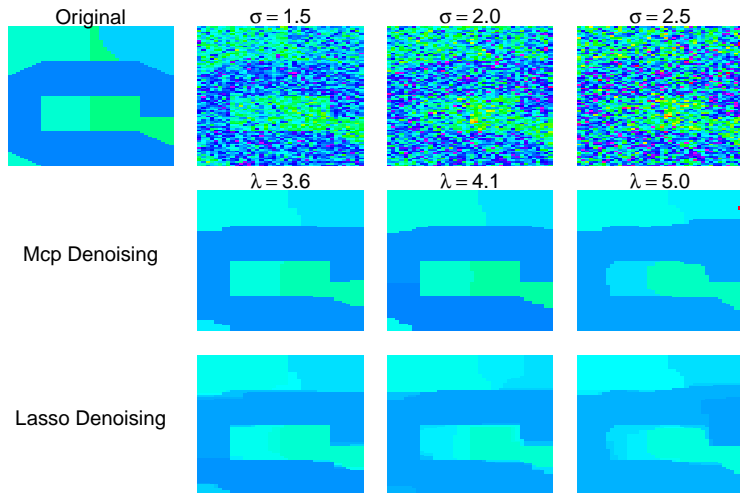
# 2-d fused MCP vs. 2-d FLSA

- When the noise is small



# 2-d fused MCP vs. 2-d FLSA

- When the magnitude of noise increases



## 2-d fused MCP vs. 2-d FLSA

2-d fused Mcp keeps the contrast of the colors and finds the edges in the images effectively. It largely improves 2-d FLSA.

# Outline

- 1 Review
- 2 Signal Denoising and FLSA
- 3 Method
- 4 Algorithms**
- 5 Oracle Property
- 6 Simulations
- 7 Summary

# Algorithms

|                                  |                       |
|----------------------------------|-----------------------|
| fused MCP( $\lambda_1 > 0$ )     | adjusted MM Algorithm |
| fusion MCP( $\lambda_1 = 0, b$ ) | PLUS Algorithm        |
| 2-d fused MCP                    | adjusted MM Algorithm |



# fusion MCP

$$\min_{\beta} \frac{1}{2} \|\mathbf{Y} - \beta\|_2^2 + \sum_{j=2}^p \rho(|\beta_j - \beta_{j-1}|; \lambda)$$

can be transformed to the objective function of an MCP penalized regression problem:

$$\min_{\eta} \frac{1}{2} \|\mathbf{A} - \mathbf{B}\eta\|_2^2 + \sum_{j=1}^{p-1} \rho(|\eta_j|; \lambda).$$

Thus, PLUS algorithm can be applied.

# fused MCP

The objective function becomes more complex, PLUS algorithm is not applicable any more.

We design an adjusted majorization-minimization (MM) algorithm to solve this problem.

# MM Algorithm for fused MCP

## MM Algorithm

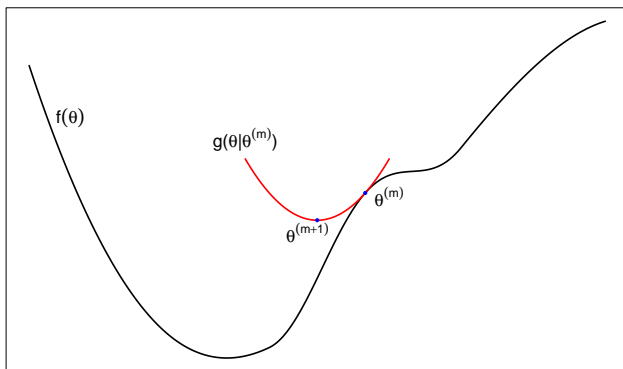


Figure: Illustration of one step in the MM algorithm

# MM Algorithm for fused MCP

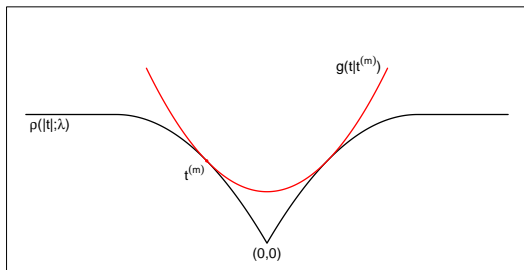
- $g(\theta | \theta^{(m)}) \geq f(\theta)$  for all  $\theta$
- $g(\theta^{(m)} | \theta^{(m)}) = f(\theta^{(m)})$
- $g(\theta | \theta^{(m)})$  is devised to be easy to solve
- The optimum  $\theta$  of  $f(\theta)$  is found by minimizing  $g(\theta | \theta^{(m)})$  iteratively

# MM Algorithm for fused MCP

The objective function is

$$\frac{1}{2} \sum_i (y_i - \beta_j)^2 + \sum_{j=1}^p \rho(|\beta_j|; \lambda_1) + \sum_{j=2}^p \rho(|\beta_j - \beta_{j-1}|; \lambda_2).$$

The term  $\rho(|t|; \lambda_1)$  is majorized by  $g(t | t^{(m)})$  and  $g$  is quadratic.

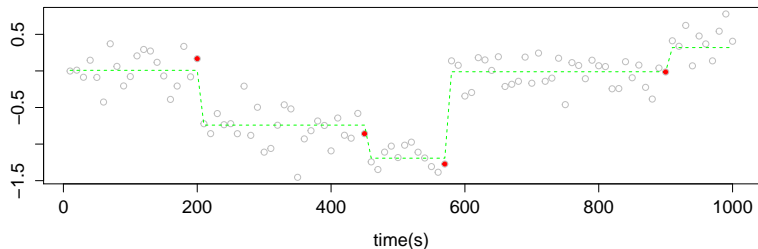


# Outline

- 1 Review
- 2 Signal Denoising and FLSA
- 3 Method
- 4 Algorithms
- 5 Oracle Property**
- 6 Simulations
- 7 Summary

# Oracle Estimation

If the jump points are known, they will partition the signal into several segments. A reasonable way to denoise the signal, knowing this partition, would be to average the observed signal in every segment.



The oracle estimation is denoted as  $\beta^{\mathcal{O}}$ .

# Oracle Property of fusion MCP

## Theorem

Suppose  $\mathbf{Y}_{p \times 1} = \beta^0 + \varepsilon$ ,  $A^0 = \{j | \beta_j^0 \neq \beta_{j+1}^0\}$  and  $\tilde{\mathbf{H}}$  is a matrix that depends only on  $p$ .  $c_{\min}(\mathbf{M})$  denotes the smallest eigenvalue of matrix  $\mathbf{M}$ . let

$$\hat{A} \equiv \{j | \hat{\beta}_j \neq \hat{\beta}_{j+1}\},$$

$$\tilde{\mathbf{H}}_{\mathcal{O}} \equiv (\tilde{\mathbf{H}}_j, j \in A^0)_{p \times |A^0|},$$

$$(\omega_j^0, j \in A^0) \equiv \text{the diagonal elements of } (\tilde{\mathbf{H}}_{\mathcal{O}}^T \tilde{\mathbf{H}}_{\mathcal{O}})^{-1}.$$

...



## Theorem

If

$$\sup_{\|u\|_2=1} P(u^T \varepsilon > \sigma t) \leq e^{-t^2/2} \quad \forall t > 0, \quad (1)$$

$$\gamma > \frac{1}{c_{\min}(\tilde{\mathbf{H}}^T \tilde{\mathbf{H}})}, \quad (2)$$

$$P(\lambda_l \leq \lambda \leq \lambda_u) = 1 \text{ and } \lambda_u \gamma \leq \eta_* = \min_{j \in A^0} \left\{ \left| \beta_j^0 - \beta_{j+1}^0 \right| \right\}, \quad (3)$$

then

$$P(\hat{\mathbf{A}} \neq A^0) \leq P(\hat{\beta} \neq \beta^0) \leq \pi_1(\lambda_l) + \pi_2(\lambda_u), \quad (4)$$

where  $\pi_1(\lambda) = 2 \sum_{j \notin A^0} \exp \left\{ -\frac{\lambda^2}{2\sigma^2 \|\tilde{\mathbf{H}}_j\|_2^2} \right\}$ ,

$$\pi_2(\lambda) = \sum_{j \in A^0} \exp \left\{ -\frac{(\lambda\gamma - |\beta_j^0 - \beta_{j+1}^0|)^2}{2\omega_j^0 \sigma^2} \right\}.$$

# Outline

- 1 Review
- 2 Signal Denoising and FLSA
- 3 Method
- 4 Algorithms
- 5 Oracle Property
- 6 Simulations**
- 7 Summary

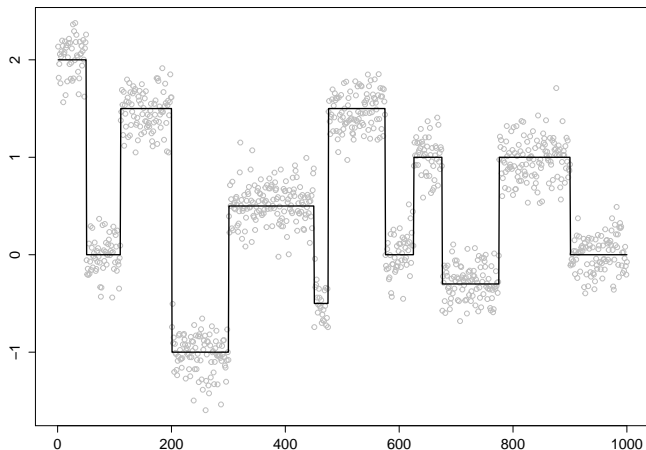
# Simulations

We compare the estimation and selection accuracy of fusion MCP and FLSA with  $\lambda_1 = 0$ . Two types of shapes of the signals are considered.

For each case, we repeat fusion MCP and FLSA for 500 times.

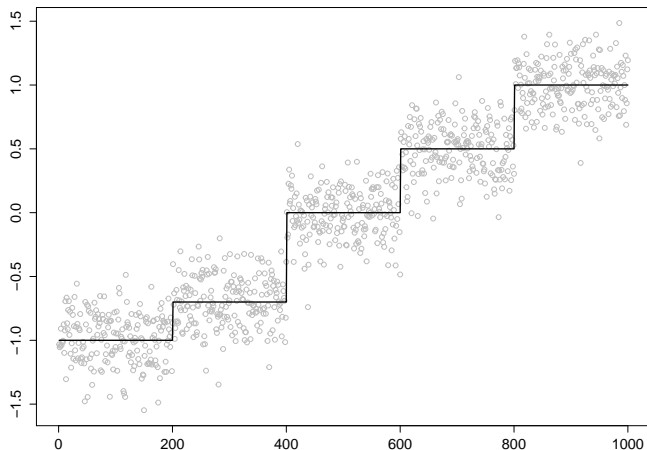
# Two Cases

- Case 1



# Two Cases

- Case 2



# Estimation and Selection Accuracy

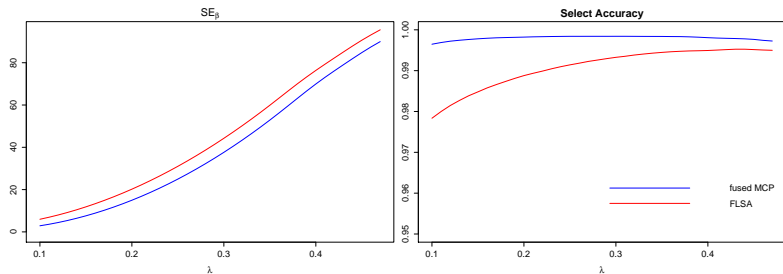


Figure: Case 1.  $SE_{\beta} = \|\hat{\beta} - \beta^0\|^2$ .

# Estimation and Selection Accuracy

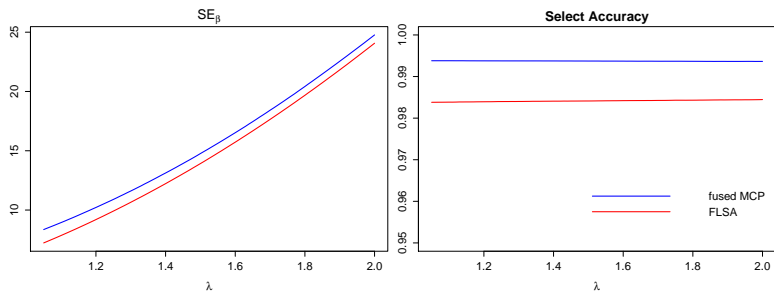


Figure: Case 2

# Outline


- 1 Review
- 2 Signal Denoising and FLSA
- 3 Method
- 4 Algorithms
- 5 Oracle Property
- 6 Simulations
- 7 Summary**



# Summary

- Fused Mcp dose not penalize large jumps
- It has a better performance in jump point(edge) detection
- It keeps the jump size(contrast) of the signal
- In signal processing scenario, both 1-d and 2-d fused Mcp problem can be solved efficiently

# References

-  Friedman, J., Hastie, T., Hofling, H. and Tibshirani, R. (2007). Pathwise coordinate optimization. *The Annals of Applied Statistics*, **1(2)**, 302-332.
-  Hunter, D. R., and Li, R. (2005). Variable selection using MM algorithms. *Annals of statistics*, **33(4)**, 1617.
-  Rudin, L. I., Osher, S. and Fatemi, E. (1992). Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, **60(1)**, 259-268.
-  Tibshirani, R. (1996). Regression shrinkage and selection via the lasso, *J. Royal. Statist. Soc. B.* **58**, 267-288.
-  Tibshirani, R., Saunders, M., Rosset, S., Zhu, J. and Knight, K. (2005). Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **67(1)**, 91-108.
-  Tibshirani, R. J. and Taylor, J. (2011). The solution path of the generalized lasso. *The Annals of Statistics*, **39(3)**, 1335-1371.
-  Zhang, C. H. (2010). Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, 894-942.