# Variable Selection in High Dimensional Time-varying Effect Model: based on Iterative Shrinkage Thresholding Algorithm (ISTA)

Yujie Liao

Department of Statistics

## Table of contents

- Motivation
- Two challenges:
  1. ISTA
  2. Variant of standard ISTA
- Backtracking rule
- Simulation

- Variable selection in high-dimensional time-varying effect model

$$Y_i = \sum_{k=1}^{K} \beta_k(t_i)X_{ik} + \epsilon_i$$

- iid $(Y_i, \boldsymbol{X}_i, t_i)$, $i = 1, \cdots, n$ (not longitudinal case for simplicity)
- General method: approximate $\beta_k(t) \approx \sum_{l=1}^{L} \beta_{kl}B_l(t)$
- Interests: $\boldsymbol{\beta}_k^* = (\beta_{k1}, \cdots, \beta_{kL})^T$ and $\boldsymbol{\beta}^* = (\boldsymbol{\beta}_1^{*T}, \cdots, \boldsymbol{\beta}_K^{*T})^T$
- Select the non-zero $\beta_k(t)$ by minimizing

$$Q(\beta^*) = \frac{1}{2n} \sum_{i=1}^{n} \{Y_i - \sum_{k-1}^{K} \sum_{l=1}^{L} \beta_{kl}B_l(t_i)X_{ik}\}^2 + g(\boldsymbol{\beta}^*)$$

$$= \frac{1}{2n} \sum_{i=1}^{n} (Y_i - \beta^{*T}\boldsymbol{X}_i^*)^2 + g(\boldsymbol{\beta}^*)$$

where $\boldsymbol{X}_i^* = (X_{i1}\boldsymbol{B}(t_i)^T, \cdots, X_{ik}\boldsymbol{B}(t_i)^T)$, $g(\boldsymbol{\beta}^*)$ is the penalty function (mentioned later).

- Two challenges: $Q(\beta^*) = f(\beta^*) + g(\beta^*)$

# Iterative Shrinkage Thresholding Algorithm (ISTA)

- Target: $Q(\beta) = f(\beta) + g(\beta)$, where $f(\cdot)$ is smooth (probably complicated) and $g(\cdot)$ is non-smooth

- Local isotropic quadratic approximation:

  1. $Q(\beta) \approx Q_A(\beta) = f_A(\beta|\beta_{r-1}) + g(\beta)$, where $f_A(\beta|\beta_{r-1})$ is

  $$f(\beta_{r-1}) + f^{'}(\beta_{r-1})^T(\beta - \beta_{r-1}) + \frac{1}{2}(\beta - \beta_{r-1})^T \frac{\partial^2 f(\beta_{r-1})}{\partial\beta\partial\beta^T}(\beta - \beta_{r-1})$$

  2. By Raleigh-Ritz Theorem: $\forall x \neq 0 \in \mathbb{R}^n$, $A : n \times n$ Hermitian matrix,

  $$\lambda_{min}(A)x^T x \leq x^T A x \leq \lambda_{max}(A)x^T x$$

  If $\{\frac{\partial^2 f(\beta_{r-1})}{\partial\beta\partial\beta^T}\}^{-1} \approx s_r I$, then

  $$f_A(\beta) = f(\beta_{r-1}) + f^{'}(\beta_{r-1})^T(\beta - \beta_{r-1}) + \frac{1}{2s_r}\|\beta - \beta_{r-1}\|^2$$

  3. Goal:
  $\beta_r = \arg\min_\beta Q_A(\beta|\beta_{r-1}, s_r) = \arg\min_\beta\{f_A(\beta_r|\beta_{r-1}, s_r) + g(\beta_r)\}$

3

# Why interesting?

$$\beta_r = \arg\min_{\beta}\{f^{'}(\beta_{r-1})^T\beta + \frac{1}{2s_r}\left\|\beta - \beta_{r-1}\right\|^2 + g(\beta_r)\}$$

$$= \arg\min_{\beta}\left[\frac{1}{2s_r}\left\|\beta - \{\beta_{r-1} - s_r f'(\beta_{r-1})\}\right\|^2 + g(\beta_r)\right]$$

$$= \arg\min_{\beta}\{\frac{1}{2s_r}\left\|\beta - \tilde{\beta}_r\right\|^2 + g(\beta_r)\} \quad \text{where} \quad \tilde{\beta}_r = \beta_{r-1} - s_r f'(\beta_{r-1})$$

- Pros:
    1. Standard variable selection subject function
    2. Not depend on $f(\cdot)$ at all, only on $g(\cdot)$
    3. $f(\cdot)$ can be complicated, we only need to consider its gradient
    4. closed form for many important functions (e.g. lasso $\rightarrow$ component-wise $\rightarrow$ soft-thresholding rule)

## Also called Proximal Gradient Descent

- Define $prox_{g,s}(\boldsymbol{z}) = \arg\min_{\boldsymbol{x}} \frac{1}{2s} \|\boldsymbol{x} - \boldsymbol{z}\|^2 + g(\boldsymbol{x})$

- Choose initialize $\boldsymbol{x}_0$, repeat

$$\boldsymbol{x}_r = prox_{g,s_r}\{\boldsymbol{x}_{r-1} - s_r f'(\boldsymbol{x}_{r-1})\}, \quad r = 1, 2, \cdots$$

- Then $\boldsymbol{x}_r = \boldsymbol{x}_{r-1} - s_r \cdot G_{s_r}(\boldsymbol{x}_{r-1})$,
  where $G_s$ is the generalized gradient of $f$,

$$G_s(\boldsymbol{x}) = \frac{\boldsymbol{x} - prox_{g,s}(\boldsymbol{x} - s \cdot g'(\boldsymbol{x}))}{s}.$$

## Another Challenge

- $Q(\beta^*) = \frac{1}{2n} \|\boldsymbol{Y} - \boldsymbol{X}^*\beta^*\|^2 + g(\beta^*)$, $\beta^* = (\beta_1^{*T}, \cdots, \beta_K^{*T})^T \in \mathbb{R}^{K \times L}$

- $\beta_k^* = (\beta_{k1}, \cdots, \beta_{kL})^T \in \mathbb{R}^L$ is associated with group k, since $\beta_k(t) \approx \sum_{l=1}^L \beta_{kl} B_l(t)$

- Time-varying model: not penalty on each component of $\beta^* \in \mathbb{R}^{K \times L}$

- Group variable selection: need coefficients in a group to be in or out of the model at the same time $\rightarrow$ sparsity between groups, not within groups

- Group LASSO, group SCAD, other constraints: $\sum_{k=1}^K p_\lambda(\|\beta_k^*\|_2)$ not component-wise, but group-wise

- Solution:
  $\min_{\beta^*} f(\beta^*) = \frac{1}{2n} \|\boldsymbol{Y} - \boldsymbol{X}^*\beta^*\|^2$ s.t. $\tau(\{k : \|\beta_k^*\|_2 > 0\}) \leq m$

$$\beta_{k,r}^* = \arg\min_{\beta_k^*} f_A(\beta_k^*|\beta_{k,r-1}^*)$$

$$= \arg\min_{\beta_k^*}\{f(\beta_{k,r-1}^*) + f^{'}(\beta_{k,r-1}^*)^T(\beta_k^* - \beta_{k,r-1}^*) + \frac{1}{2s_r}\left\|\beta_k^* - \beta_{k,r-1}^*\right\|^2\}$$

$$= \arg\min_{\beta^*}\{\frac{1}{2s_r}\left\|\beta_k^* - \tilde{\beta}_{k,r}^*\right\|^2\} = h_A(\beta_k^*),$$

where $\tilde{\beta}_r^* = \beta_{r-1}^* - s_r f'(\beta_{r-1}^*)$, for $\forall k = 1, \cdots, K$,

s.t. $\tau(\{k : \left\|\beta_k^*\right\|_2 > 0\}) \le m$.

- When $\hat{\beta}_k^* \neq 0$, then $\hat{\beta}_k^* = \tilde{\beta}_k^*$ $\Rightarrow$ $h_{A1}(\beta_k^*) = 0$
  When $\hat{\beta}_k^* = 0$ $\Rightarrow$ $h_{A2}(\beta_k^*) = \frac{1}{2s_r}\left\|\tilde{\beta}_k^*\right\|^2$
- Set $\beta_k^* = 0$ if $h_{A2}(\beta_k^*) - h_{A1}(\beta_k^*) = \frac{1}{2s_r}\left\|\tilde{\beta}_k^*\right\|^2$ is small,
  i.e. $\left\|\tilde{\beta}_k^*\right\|^2$ is small
- Let $g_k = \tilde{\beta}_k^{*T}\tilde{\beta}_k^*$, sort $g_i$ so that $g_{(1)} \ge g_{(2)} \ge \cdots g_{(K)}$.
  By hard-thresholding rule, $\hat{\beta}_k^* = \tilde{\beta}_k^* I\{g_k > g_{(m+1)}\}$

- $f(\beta^*) \to f_A(\beta^*|\beta^*_{r-1}) =$
  $f(\beta^*_{r-1}) + f'(\beta^*_{r-1})^T(\beta^* - \beta^*_{r-1}) + \frac{1}{2s_r} \left\| \beta^* - \beta^*_{r-1} \right\|^2$
- Check two conditions of MM algorithm:
  1. $f_A(\beta^*_{r-1}|\beta^*_{r-1}) = f(\beta^*_{r-1})$
  2. Majorizing-minimization: choose $s_r$ by backtracking rule
     2.1 Set $s_0 > 0$, $0 < \delta < 1$, $\beta^*_0$.
     2.2 Find the smallest non-negative integer $i_r$ s.t. with $s = \delta^{i_r} s_{r-1}$,

     $$f(\beta^*_{r,s}) \leq f_A(\beta^*_{r,s}|\beta^*_{r-1}, s),$$

     i.e. $\quad f(\beta^*_{r,s}) \leq f(\beta^*_{r-1}) + f'(\beta^*_{r-1})^T(\beta^* - \beta^*_{r-1}) + \frac{1}{2s} \left\| \beta^* - \beta^*_{r-1} \right\|^2,$

     where $\beta^*_{k,r,s} \leftarrow \hat{\beta}^*_{k,r,s} = \tilde{\beta}^*_{k,r,s} I\{g_k > g_{(m+1)}\}$ is a function of $s$.
     2.3 Set $s_r = \delta^{i_r} s_{r-1} \quad \Rightarrow \quad \hat{\beta}^*_{k,r,s_r}$

- $(t_i^{'}, \boldsymbol{X}_i^{'}) \sim N_{K+1}(\boldsymbol{0}, \boldsymbol{\Sigma})$, where $\boldsymbol{\Sigma} = (\sigma_{ij})$
  $\sigma_{ij} = 1$ if $i = j$; $\sigma_{ij} = \rho$ if $i \neq j$; $\rho = 0.5$

- $t_i = \Phi(t_i^{'})$, where $\Phi(\cdot)$ is the CDF of $N(0, 1)$

- $L = 5$, $Rep = 1000$, $K = 800$, $n = 200$, $m = [\frac{n^{4/5}}{log(n^{4/5})}]$

- True coefficient functions:
  $\beta_1(t) = 0.95cos(\frac{\pi t}{2}) + 3.36, \quad \beta_2(t) = 1.5sin(2\pi t) + 3,$
  $\beta_3(t) = -2t^2 - 4, \quad \beta_4(t) = 0.52t + t^3 + 2.9$

- Criteria: $P_k$: the proportion of submodels $\hat{\mathcal{M}}$ with size m that contain $X_k$ among $Rep$ repetitions

- Results:

| $P_1$ | $P_2$ | $P_3$ | $P_4$ |
|-------|-------|-------|-------|
| 1 | 0.94 | 0.88 | 1 |