

Computation & optimization for Lasso - part 2

Luyang Han & Janosch Ott

ETH Zürich

22 October 2018

2018-10-02

Computation & optimization

Overview

1. Coordinate Descent
2. A Simulation Study
3. Least Angle Regression
4. ADMM
5. Minor-Max Algorithms
6. Alternating Minimizations
7. Screening Rules

2018-10-02

Computation & optimization

└ Overview

1. Coordinate Descent
2. A Simulation Study
3. Least Angle Regression
4. ADMM
5. Minor-Max Algorithms
6. Alternating Minimizations
7. Screening Rules

Alternating Direction Method of Multipliers (ADMM)

Problem

$$\underset{\beta \in \mathbb{R}^m, \theta \in \mathbb{R}^n}{\text{minimize}} \quad f(\beta) + g(\theta) \quad \text{subject to} \quad \mathbf{A}\beta + \mathbf{B}\theta = c$$

Lagrangian

$$f(\beta) + g(\theta) + \rho \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

Augmented Lagrangian

$$L_\rho(\beta, \theta, \mu) := f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle + \frac{\rho}{2} \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

2018-10-02

Computation & optimization

└ ADMM

└ Alternating Direction Method of Multipliers
(ADMM)Augmented: scalar product with μ gets added

Problem

$$\underset{\beta \in \mathbb{R}^m, \theta \in \mathbb{R}^n}{\text{minimize}} \quad f(\beta) + g(\theta) \quad \text{subject to} \quad \mathbf{A}\beta + \mathbf{B}\theta = c$$

Lagrangian

$$f(\beta) + g(\theta) + \rho \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

Augmented Lagrangian

$$L_\rho(\beta, \theta, \mu) := f(\beta) + g(\theta) + \langle \mu, \mathbf{A}\beta + \mathbf{B}\theta - c \rangle + \frac{\rho}{2} \|\mathbf{A}\beta + \mathbf{B}\theta - c\|_2^2$$

Dual variable update

$$\beta^{t+1} = \arg \min_{\beta \in \mathbb{R}^m} L_{\rho}(\beta, \theta^t, \mu^t)$$

$$\theta^{t+1} = \arg \min_{\theta \in \mathbb{R}^m} L_{\rho}(\beta^{t+1}, \theta, \mu^t)$$

$$\mu^{t+1} = \mu^t + \rho(\mathbf{A}\beta^{t+1} + \mathbf{B}\theta^{t+1} - c)$$

2018-10-02

Computation & optimization

└ ADMM

└ Dual variable update

Dual variable update

$$\begin{aligned}\beta^{t+1} &= \arg \min_{\beta \in \mathbb{R}^m} L_{\rho}(\beta, \theta^t, \mu^t) \\ \theta^{t+1} &= \arg \min_{\theta \in \mathbb{R}^m} L_{\rho}(\beta^{t+1}, \theta, \mu^t) \\ \mu^{t+1} &= \mu^t + \rho(\mathbf{A}\beta^{t+1} + \mathbf{B}\theta^{t+1} - c)\end{aligned}$$

ADMM - Why?

- convex problems with nondifferentiable constraints
- blockwise computation
 - sample blocks
 - feature blocks

2018-10-02

Computation & optimization

└ ADMM

└ ADMM - Why?

- convex problems with nondifferentiable constraints
- blockwise computation
 - sample blocks
 - feature blocks

ADMM for the Lasso

Problem in Lagrangian form

$$\underset{\beta \in \mathbb{R}^p, \theta \in \mathbb{R}^p}{\text{minimize}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\theta\|_1 \right\} \quad \text{such that } \beta - \theta = 0$$

Update

$$\beta^{t+1} = (\mathbf{X}^T \mathbf{X} + \rho \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{y} + \rho \theta^t - \mu^t)$$

$$\theta^{t+1} = \mathcal{S}_{\lambda/\rho}(\beta^{t+1} + \mu^t/\rho)$$

$$\mu^{t+1} = \mu^t + \rho(\beta^{t+1} - \theta^{t+1})$$

where $\mathcal{S}_{\lambda/\rho}(z) = \text{sign}(z)(|z| - \frac{\lambda}{\rho})_+$.

$$\underset{\beta \in \mathbb{R}^p, \theta \in \mathbb{R}^p}{\text{minimize}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\theta\|_1 \right\} \quad \text{such that } \beta - \theta = 0$$

$$\beta^{t+1} = (\mathbf{X}^T \mathbf{X} + \rho \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{y} + \rho \theta^t - \mu^t)$$

$$\theta^{t+1} = \mathcal{S}_{\lambda/\rho}(\beta^{t+1} + \mu^t/\rho)$$

$$\mu^{t+1} = \mu^t + \rho(\beta^{t+1} - \theta^{t+1})$$

$$\text{where } \mathcal{S}_{\lambda/\rho}(z) = \text{sign}(z)(|z| - \frac{\lambda}{\rho})_+.$$

Computational cost: Initially $\mathcal{O}(p^3)$, which is a lot, for the SVD(singular value decomposition of \mathbf{X}), after that comparable to coordinate descent or composite gradient from earlier

Minorization-Maximization Algorithms (MMA)

- Problem: minimize $f(\beta)$ over $\beta \in \mathbb{R}^p$
for f possibly non-convex
- Introduce additional variable θ
- Use θ to majorize (bound from above) the objective
function to be minimized

Majorization-Minimization Algorithms work analogously.

2018-10-02

Computation & optimization

└ Minor-Max Algorithms

└ Minorization-Maximization Algorithms (MMA)

- Problem: minimize $f(\beta)$ over $\beta \in \mathbb{R}^p$
for f possibly non-convex
- Introduce additional variable θ
- Use θ to majorize (bound from above) the objective
function to be minimized

Majorization-Minimization Algorithms work analogously.

MMA visually

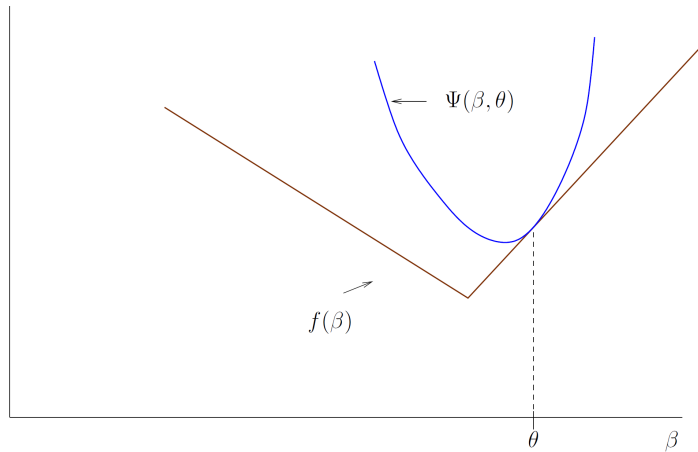


Figure: Figure 5.10 from [Hastie et al., 2015]

2018-10-02

Computation & optimization

└ Minor-Max Algorithms

└ MMA visually

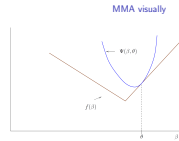


Figure: Figure 5.10 from [Hastie et al., 2015]

Computation
&
optimization

Luyang Han &
Janosch Ott

Coordinate
Descent

A Simulation
Study

Least Angle
Regression

ADMM

Minor-Max
Algorithms

Alternating
Minimizations

Screening
Rules

MMA analytically

2018-10-02

Computation & optimization

└ Minor-Max Algorithms

└ MMA analytically

MMA analytically

Dual Polytope Projection (DPP)

Suppose we want to calculate a lasso solution at $\lambda < \lambda_{\max}$.
The DPP rule discards the j^{th} variable if

$$\left| \mathbf{x}_j^T \mathbf{y} \right| < \lambda_{\max} - \|\mathbf{x}_j\|_2 \|\mathbf{y}\|_2 \frac{\lambda_{\max} - \lambda}{\lambda}$$

Sequential DPP rule

Suppose we have the lasso solution $\hat{\beta}(\lambda')$ at λ' and want to
screen variables for solutions at $\lambda < \lambda'$. We discard the j^{th}
variable if

$$\left| \mathbf{x}_j^T (\mathbf{y} - \mathbf{X} \hat{\beta}(\lambda')) \right| < \lambda' - \|\mathbf{x}_j\|_2 \|\mathbf{y}\|_2 \frac{\lambda_{\max} - \lambda}{\lambda}$$

$$\left| \mathbf{x}_j^T \mathbf{y} \right| < \lambda_{\max} - \|\mathbf{x}_j\|_2 \|\mathbf{y}\|_2 \frac{\lambda_{\max} - \lambda}{\lambda}$$

$$\left| \mathbf{x}_j^T (\mathbf{y} - \mathbf{X} \hat{\beta}(\lambda')) \right| < \lambda' - \|\mathbf{x}_j\|_2 \|\mathbf{y}\|_2 \frac{\lambda_{\max} - \lambda}{\lambda}$$

Global Strong Rule

Suppose we want to calculate a lasso solution at $\lambda < \lambda_{\max}$.

The global strong rule discards the j^{th} variable if

$$\left| \mathbf{x}_j^T \mathbf{y} \right| < \lambda - (\lambda_{\max} - \lambda) = 2\lambda - \lambda_{\max}$$

Sequential Strong Rule

Suppose we have the lasso solution $\hat{\beta}(\lambda')$ at λ' and want to screen variables for solutions at $\lambda < \lambda'$. We discard the j^{th} variable if

$$\left| \mathbf{x}_j^T (\mathbf{y} - \mathbf{X} \hat{\beta}(\lambda')) \right| < 2\lambda - \lambda'$$

$$\left| \mathbf{x}_j^T \mathbf{y} \right| < \lambda - (\lambda_{\max} - \lambda) = 2\lambda - \lambda_{\max}$$

$$\left| \mathbf{x}_j^T (\mathbf{y} - \mathbf{X} \hat{\beta}(\lambda')) \right| < 2\lambda - \lambda'$$

Sed diam enim, sagittis nec condimentum sit amet, ullamcorper sit amet libero. Aliquam vel dui orci, a porta odio. Nullam id suscipit ipsum. Aenean lobortis commodo sem, ut commodo leo gravida vitae. Pellentesque vehicula ante iaculis arcu pretium rutrum eget sit amet purus. Integer ornare nulla quis neque ultrices lobortis. Vestibulum ultrices tincidunt libero, quis commodo erat ullamcorper id.

Bullet Points

- Lorem ipsum dolor sit amet, consectetur adipiscing elit
- Aliquam blandit faucibus nisi, sit amet dapibus enim tempus eu
- Nulla commodo, erat quis gravida posuere, elit lacus lobortis est, quis porttitor odio mauris at libero
- Nam cursus est eget velit posuere pellentesque
- Vestibulum faucibus velit a augue condimentum quis convallis nulla gravida

2018-10-02

Computation & optimization

└ Screening Rules

└ Bullet Points

Bullet Points

- Lorem ipsum dolor sit amet, consectetur adipiscing elit
- Aliquam blandit faucibus nisi, sit amet dapibus enim tempus eu
- Nulla commodo, erat quis gravida posuere, elit lacus lobortis est, quis porttitor odio mauris at libero
- Nam cursus est eget velit posuere pellentesque
- Vestibulum faucibus velit a augue condimentum quis convallis nulla gravida

Computation & optimization	2018-10-02	Computation & optimization	Blocks of Highlighted Text
Luyang Han & Janosch Ott		Screening Rules	Block 1 Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer lectus nisl, ultricies in feugiat rutrum, porttitor sit amet augue. Aliquam ut tortor mauris. Sed volutpat ante purus, quis accumsan dolor.
Coordinate Descent		Blocks of Highlighted Text	Block 2 Pellentesque sed tellus purus. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Vestibulum quis magna at risus dictum tempor eu vitae velit.
A Simulation Study			Block 3 Suspendisse tincidunt sagittis gravida. Curabitur condimentum, enim sed venenatis rutrum, ipsum neque consectetur orci, sed blandit justo nisi ac lacus.
Least Angle Regression			
ADMM			
Minor-Max Algorithms			
Alternating Minimizations			
Screening Rules			

<div>Computation & optimization</div> <div>Luyang Han & Janosch Ott</div> <div>Coordinate Descent</div> <div>A Simulation Study</div> <div>Least Angle Regression</div> <div>ADMM</div> <div>Minor-Max Algorithms</div> <div>Alternating Minimizations</div> <div>Screening Rules</div>	<div>Multiple Columns</div> <div>Heading</div> <div><div>1 Statement</div><div>2 Explanation</div><div>3 Example</div></div> <div>Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer lectus nisl, ultricies in feugiat rutrum, porttitor sit amet augue. Aliquam ut tortor mauris. Sed volutpat ante purus, quis accumsan dolor.</div>	<div>2018-10-02</div> <div>Computation & optimization</div> <div>└ Screening Rules</div> <div>└ Multiple Columns</div>	<div>Multiple Columns</div> <div>Heading</div> <div><div>1 Statement</div><div>2 Explanation</div><div>3 Example</div></div> <div>Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer lectus nisl, ultricies in feugiat rutrum, porttitor sit amet augue. Aliquam ut tortor mauris. Sed volutpat ante purus, quis accumsan dolor.</div>
---	--	--	--

2018-10-02	Computation & optimization	Table
	└ Screening Rules	
	└ Table	
Table		
Table: Table caption		

Treatments	Response 1	Response 2
Treatment 1	0.0003262	0.562
Treatment 2	0.0015681	0.910
Treatment 3	0.0009271	0.296

Theorem

Theorem (Mass–energy equivalence)

$$E = mc^2$$

2018-10-02

Computation & optimization

└ Screening Rules

└ Theorem

Theorem

Theorem (Mass–energy equivalence)
 $E = mc^2$

Verbatim

Example (Theorem Slide Code)

```
\begin{frame}  
\frametitle{Theorem}  
\begin{theorem}[Mass--energy equivalence]  
$E = mc^2$  
\end{theorem}  
\end{frame}
```

2018-10-02

Example (Theorem Slide Code)

```
\begin{frame}  
\frametitle{Theorem}  
\begin{theorem}[Mass--energy equivalence]  
$E = mc^2$  
\end{theorem}  
\end{frame}
```

Figure

Uncomment the code on this slide to include your own image
from the same directory as the template .TeX file.

2018-10-02

Computation & optimization

└ Screening Rules

└ Figure

Citation

An example of the `\cite` command to cite within the presentation:

This statement requires citation [Hastie et al., 2015].

2018-10-02

Computation & optimization

└ Screening Rules

└ Citation

References



Trevor Hastie, Robert Tibshirani, and Martin Wainwright (2015)
Statistical learning with sparsity: the Lasso and generalizations
CRC Press; Boca Raton, FL

2018-10-02

Computation & optimization

└ Screening Rules

└ References

Computation
&
optimization

Luyang Han &
Janosch Ott

Coordinate
Descent

A Simulation
Study

Least Angle
Regression

ADMM

Minor-Max
Algorithms

Alternating
Minimizations

Screening
Rules

The End

2018-10-02

Computation & optimization
└ Screening Rules

The End