	Algorithm	Variable Splitting	

Differential Inclusion Method in High Dimensional Statistics

Yuan Yao

HKUST

February 12, 2018

Acknowledgements

- Theory
 - Stanley Osher, Wotao Yin (UCLA)
 - Feng Ruan (Stanford & PKU)
 - Jiechao Xiong, Chendi Huang (PKU)
- Applications:
 - Qianqian Xu, Jiechao Xiong, Chendi Huang, Xinwei Sun (PKU)
 - Lingjing Hu (BCMU)
 - Ming Yan, Zhimin Peng (UCLA)
- Grants:
 - National Basic Research Program of China (973 Program), NSFC

Outline	LASSO vs. Differential Inclusions	Algorithm 0000000000000	Variable Splitting 00000	

From LASSO to Differential Inclusions

- LASSO and Bias
- Differential Inclusions
- A Theory of Path Consistency

2 Large Scale Algorithm

- Linearized Bregman Iteration
- Generalizations
- Cran R package: Libra

3 Variable Splitting

A Weaker Irrepresentable/Incoherence Condition

4 Summary

	ŧ١		•	

LASSO vs. Differential Inclusions

Algorithm

Sparse Linear Regression

Assume that $\beta^* \in \mathbb{R}^{\rho}$ is sparse and unknown. Consider recovering β^* from *n* linear measurements

$$y = X\beta^* + \epsilon, \qquad y \in \mathbb{R}^n$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is **noise**.

- Basic Sparsity: $S := \operatorname{supp}(\beta^*)$ (s = |S|) and T be its complement.
 - $X_S(X_T)$ be the columns of X with indices restricted on S(T)
 - X is *n*-by-*p*, with $p \gg n \ge s$.
- Or Structural Sparsity: γ^{*} = Dβ^{*} is sparse, where D is a linear transform (wavelet, gradient, etc.), S = supp(γ^{*})
- How to recover β* (or γ*) sparsity pattern (sparsistency) and estimate values with variations (consistency)?

Our Best Possible in Basic Setting: The Oracle Estimator

Had God revealed S to us, the oracle estimator was the subset least square solution (MLE) with $\tilde{\beta}_{T}^{*}=0$ and

$$\tilde{\beta}_{S}^{*} = \beta_{S}^{*} + \frac{1}{n} \Sigma_{n}^{-1} X_{S}^{\mathsf{T}} \epsilon, \quad \text{where } \Sigma_{n} = \frac{1}{n} X_{S}^{\mathsf{T}} X_{S} \tag{1}$$

"Oracle properties"

- Model selection consistency: $\operatorname{supp}(\tilde{\beta}^*) = S$;
- Normality: $\tilde{\beta}_{S}^{*} \sim \mathcal{N}(\beta^{*}, \frac{\sigma^{2}}{n} \Sigma_{n}^{-1}).$

So $\tilde{\beta}^*$ is unbiased, i.e. $\mathbb{E}[\tilde{\beta}^*] = \beta^*$.

	LASSO vs. Differential Inclusions	Algorithm 0000000000000	Variable Splitting 00000	
LASSO and Bias				

Recall LASSO

LASSO:

$$\min_{\beta} \|\beta\|_1 + \frac{t}{2n} \|y - X\beta\|_2^2.$$

optimality condition:

$$\frac{\rho_t}{t} = \frac{1}{n} X^T (y - X\beta_t), \tag{2a}$$
$$\rho_t \in \partial \|\beta_t\|_1, \tag{2b}$$

where $\lambda = 1/t$ is often used in literature.

- Chen-Donoho-Saunders'1996 (BPDN)
- Tibshirani'1996 (LASSO)

	LASSO vs. Differential Inclusions		
	000000000000000000000000000000000000000		
LASSO and Rise			

The Bias of LASSO

LASSO is biased, i.e. $\mathbb{E}(\hat{\beta}) \neq \beta^*$

• e.g. X = Id, n = p = 1, LASSO is soft-thresholding

$$\hat{\beta}_{\tau} = \begin{cases} 0, & \text{if } \tau < 1/\tilde{\beta}^*, \\ \tilde{\beta}^* - \frac{1}{\tau}, & \text{otherwise} \end{cases}$$

• e.g. $n = 100, \ p = 256, \ X_{ij} \sim \mathcal{N}(0,1), \ \epsilon_i \sim \mathcal{N}(0,0.1)$



True vs LASSO (t hand-tuned)

	LASSO vs. Differential Inclusions	Algorithm	Variable Splitting	
	000000000000000000000000000000000000000			
LACCO L D				

LASSO Estimator is Biased at Path Consistency

Even when the following **path consistency** (conditions given by Zhao-Yu'06, Zou'06, Yuan-Lin'07, Wainwright'09, etc.) is reached at τ_n :

$$\exists \tau_n \in (0,\infty) \text{ s.t. } \operatorname{supp}(\hat{\beta}_{\tau_n}) = S,$$

LASSO estimate is biased away from the oracle estimator

$$(\hat{\beta}_{\tau_n})_S = \tilde{\beta}_S^* - \frac{1}{\tau_n} \Sigma_{n,S}^{-1} \operatorname{sign}(\beta_S^*), \quad \tau_n > 0.$$

How to remove the bias and return the Oracle Estimator?

	LASSO vs. Differential Inclusions	Algorithm 0000000000000	Variable Splitting 00000	
LASSO and Bias				

Nonconvex Regularization?

• To reduce bias, non-convex regularization was proposed (Fan-Li's SCAD, Zhang's MPLUS, Zou's Adaptive LASSO, l_q (q < 1), etc.)

$$\min_{\beta}\sum_{i}p(|\beta_{i}|)+\frac{t}{2n}\|y-X\beta\|_{2}^{2}.$$



- Yet it is generally hard to locate the global optimizer
- Any other simple scheme?

	LASSO vs. Differential Inclusions	Algorithm 00000000000000	Variable Splitting 00000		
Differential Inclusior	s				
New Idea					

• LASSO:

$$\min_{\beta} \|\beta\|_1 + \frac{t}{2n} \|y - X\beta\|_2^2.$$

	LASSO vs. Differential Inclusions		
	0000 0000 000000		
Differential Inclusions			

• LASSO:

$$\min_{\beta} \|\beta\|_1 + \frac{t}{2n} \|y - X\beta\|_2^2.$$

• KKT optimality condition:

$$\Rightarrow \rho_t = \frac{1}{n} X^{\mathsf{T}} (y - X \beta_t) t$$

	LASSO vs. Differential Inclusions	Algorithm	Variable Splitting	
	0000 0000 000000			
Differential Inclusions				

• LASSO:

$$\min_{\beta} \|\beta\|_1 + \frac{t}{2n} \|y - X\beta\|_2^2.$$

• KKT optimality condition:

$$\Rightarrow \rho_t = \frac{1}{n} X^{\mathsf{T}} (y - X \beta_t) t$$

• Taking derivative (assuming differentiability) w.r.t. t

$$\Rightarrow \dot{\rho}_t = \frac{1}{n} X^T (y - X(\dot{\beta}_t t + \beta_t)), \quad \rho_t \in \partial \|\beta_t\|_1$$

	LASSO vs. Differential Inclusions	Algorithm	Variable Splitting	
	0000 0000 000000			
Differential Inclusions				

• LASSO:

$$\min_{\beta} \|\beta\|_1 + \frac{t}{2n} \|y - X\beta\|_2^2$$

• KKT optimality condition:

$$\Rightarrow \rho_t = \frac{1}{n} X^T (y - X \beta_t) t$$

• Taking derivative (assuming differentiability) w.r.t. t

$$\Rightarrow \dot{\rho}_t = \frac{1}{n} X^T (y - X(\dot{\beta}_t t + \beta_t)), \quad \rho_t \in \partial \|\beta_t\|_1$$

• Assuming sign-consistency in a neighborhood of τ_n ,

for
$$i \in S$$
, $\rho_{\tau_n}(i) = \operatorname{sign}(\beta^*(i)) \in \pm 1 \Rightarrow \dot{\rho}_{\tau_n}(i) = 0$,
 $\Rightarrow \dot{\beta}_{\tau_n} \tau_n + \beta_{\tau_n} = \tilde{\beta}^*$

	LASSO vs. Differential Inclusions	Algorithm	Variable Splitting	
	0000 0000 000000			
Differential Inclusions				

• LASSO:

$$\min_{\beta} \|\beta\|_1 + \frac{t}{2n} \|y - X\beta\|_2^2.$$

• KKT optimality condition:

$$\Rightarrow \rho_t = \frac{1}{n} X^{\mathsf{T}} (y - X \beta_t) t$$

• Taking derivative (assuming differentiability) w.r.t. t

$$\Rightarrow \dot{\rho}_t = \frac{1}{n} X^T (y - X(\dot{\beta}_t t + \beta_t)), \quad \rho_t \in \partial \|\beta_t\|_1$$

• Assuming sign-consistency in a neighborhood of τ_n ,

for
$$i \in S$$
, $\rho_{\tau_n}(i) = \operatorname{sign}(\beta^*(i)) \in \pm 1 \Rightarrow \dot{\rho}_{\tau_n}(i) = 0$,
 $\Rightarrow \dot{\beta}_{\tau_n}\tau_n + \beta_{\tau_n} = \tilde{\beta}^*$

• Equivalently, the blue part removes bias of LASSO automatically

$$\beta_{\tau_n}^{\text{lasso}} = \tilde{\beta}^* - \frac{1}{\tau_n} \Sigma_n^{-1} \text{sign}(\beta^*) \Rightarrow \dot{\beta}_{\tau_n}^{\text{lasso}} \tau_n + \beta_{\tau_n}^{\text{lasso}} = \tilde{\beta}^*(\text{oracle})!$$

()	tı		ρ	

Differential Inclusions

Differential Inclusion: Inverse Scaled Spaces (ISS)

Differential inclusion replacing $\dot{\beta}_{\tau_n}^{\textit{lasso}}\tau_n+\beta_{\tau_n}^{\textit{lasso}}$ by β_t

$$\dot{\rho}_t = \frac{1}{n} X^T (y - X \beta_t), \tag{3a}$$

$$\rho_t \in \partial \|\beta_t\|_1. \tag{3b}$$

starting at t = 0 and $\rho(0) = \beta(0) = \mathbf{0}$.

• Replace ho/t in LASSO KKT by ${
m d}
ho/{
m d}t$

$$\frac{\rho_t}{t} = \frac{1}{n} X^T (y - X \beta_t)$$

 Burger-Gilboa-Osher-Xu'06 (in image recovery it recovers the objects in an inverse-scale order as t increases (larger objects appear in β_t first))

	LASSO vs. Differential Inclusions		
	0000 00000 000000		
Differential Inclusions			

Examples

• e.g. X = Id, n = p = 1, hard-thresholding

$$eta_{ au} = \left\{ egin{array}{cc} 0, & ext{if } au < 1/(ilde{eta}^*); \ ilde{eta}^*, & ext{otherwise}, \end{array}
ight.$$

• the same example shown before





	LASSO vs. Differential Inclusions	Algorithm	Variable Splitting	
	0000 00000 000000			
D'6 CLUB C				

Differential Inclusions

Solution Path: Sequential Restricted Maximum Likelihood Estimate

• ρ_t is piece-wise linear in t,

$$\rho_t = \rho_{t_k} + \frac{t - t_k}{n} X^T (y - X \beta_{t_k}), \quad t \in [t_k, t_{k+1})$$

where $t_{k+1} = \sup\{t > t_k : \rho_{t_k} + \frac{t - t_k}{n} X^T (y - X \beta_{t_k}) \in \partial \|\beta_{t_k}\|_1 \}$

β_t is piece-wise constant in t: β_t = β_{tk} for t ∈ [t_k, t_{k+1}) and β_{tk+1} is the sequential restricted Maximum Likelihood Estimate by solving nonnegative least square (Burger et al.'13; Osher et al.'16)

$$\beta_{t_{k+1}} = \arg \min_{\beta} \qquad \|y - X\beta\|_2^2$$

subject to $(\rho_{t_{k+1}})_i \beta_i \ge 0 \qquad \forall \ i \in S_{k+1}, \qquad (4)$
 $\beta_j = 0 \qquad \forall \ j \in T_{k+1}.$

• Note: Sign consistency $\rho_t = sign(\beta^*) \Rightarrow \beta_t = \tilde{\beta}^*$ the oracle estimator

Algorithm 00000000000000 Variable Splitting

Summary

Differential Inclusions

Example: Regularization Paths of LASSO vs. ISS



Figure: Diabetes data (Efron et al.'04) and regularization paths are different, yet bearing similarities on the order of parameters being nonzero

A Theory of Path Consistency

How does it work? A Path Consistency Theory

Our aim is to show that under nearly the same conditions for sign-consistency of LASSO, there exists points on their paths $(\beta(t), \rho(t))_{t \ge 0}$, which are

- sparse
- sign-consistent (the same sparsity pattern of nonzeros as true signal)
- the oracle estimator which is unbiased, better than the LASSO estimate.
- Early stopping regularization is necessary to prevent overfitting noise!

():	111	

LASSO vs. Differential Inclusions

Variable Splitting

Summary

A Theory of Path Consistency

Intuition



	ŧ١		•	

Variable Splitting

Summary

A Theory of Path Consistency

History: two traditions of regularizations

- Penalty functions
 - ℓ_2 : Ridge regression/Tikhonov regularization: $\frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i^T \beta) + \lambda \|\beta\|_2^2$
 - ℓ_1 (sparse): Basis Pursuit/LASSO (ISTA): $\frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i^T \beta) + \lambda \|\beta\|_1^2$
- Early stopping of dynamic regularization paths
 - ℓ_2 -equivalent: Landweber iterations/gradient descent/ ℓ_2 -Boost

$$\frac{d\beta_t}{dt} = -\frac{1}{n} \sum_{i=1}^n \nabla_\beta \ell(y_i, x_i^T \beta), \quad \beta_t = \nabla \left\{ \frac{1}{2} \|\beta_t\|^2 \right\}$$

• ℓ_1 (sparse)-equiv.: Orthogonal Matching Pursuit, Linearized Bregman Iteration (sparse Mirror Descent) (not ISTA! –later)

$$\frac{d\rho_t}{dt} = -\frac{1}{n} \sum_{i=1}^n \nabla_\beta \ell(y_i, x_i^T \beta), \quad \rho_t \in \partial \|\beta_t\|_1$$

	LASSO vs. Differential Inclusions		
	000000000000000		
A Theory of Path Cor	isistency		

Assumptions

(A1) Restricted Strongly Convex: $\exists \gamma \in (0, 1]$,

$$\frac{1}{n}X_S^T X_S \ge \gamma I$$

(A2) Incoherence/Irrepresentable Condition: $\exists \eta \in (0,1)$,

$$\left\|\frac{1}{n}X_T^T X_S^{\dagger}\right\|_{\infty} = \left\|\frac{1}{n}X_T^T X_S \left(\frac{1}{n}X_S^T X_S\right)^{-1}\right\|_{\infty} \le 1 - \eta$$

- "Irrepresentable" means that one can not represent (regress) column vectors in X_T by covariates in X_S.
- The incoherence/irrepresentable condition is used independently in Tropp'04, Yuan-Lin'05, Zhao-Yu'06, and Zou'06, Wainwright'09, etc.

)			

Algorithm

Variable Splitting 00000 Summary

A Theory of Path Consistency

Understanding the Dynamics

ISS as restricted gradient descent:

$$\dot{
ho}_t = -
abla L(eta_t) = rac{1}{n} X^T (y - X eta_t), \ \
ho_t \in \partial \|eta_t\|_1$$

such that

- incoherence condition and strong signals ensure it firstly evolves on index set *S* to reduce the loss
- strongly convex in subspace restricted on index set $S \Rightarrow$ fast decay in loss
- early stopping after all strong signals are detected, before picking up the noise



():	111	

LASSO vs. Differential Inclusions

Variable Splitting 00000 Summary

A Theory of Path Consistency

Path Consistency

Theorem (Osher-Ruan-Xiong-Y.-Yin'2016)

Assume (A1) and (A2). Define an early stopping time

$$\overline{\tau} := \frac{\eta}{2\sigma} \sqrt{\frac{n}{\log p}} \left(\max_{j \in T} \|X_j\| \right)^{-1},$$

and the smallest magnitude $\beta^*_{\min} = \min(|\beta^*_i| : i \in S)$. Then

- No-false-positive: for all t ≤ τ, the path has no-false-positive with high probability, supp(β(t)) ⊆ S;
- Consistency: moreover if the signal is strong enough such that

$$\beta_{\min}^* \geq \left(\frac{4\sigma}{\gamma^{1/2}} \vee \frac{8\sigma(2 + \log s) \left(\max_{j \in T} \|X_j\|\right)}{\gamma \eta}\right) \sqrt{\frac{\log p}{n}},$$

there is $\tau \leq \overline{\tau}$ such that solution path $\beta(t) = \tilde{\beta}^*$ for every $t \in [\tau, \overline{\tau}]$.

Note: equivalent to LASSO with $\lambda^* = 1/\overline{\tau}$ (Wainwright'09) up to log s.

	Algorithm	Variable Splitting	
	••••••		

Large scale algorithm: Linearized Bregman Iteration

Damped Dynamics: continuous solution path

$$\dot{\rho}_t + \frac{1}{\kappa} \dot{\beta}_t = \frac{1}{n} X^T (y - X \beta_t), \quad \rho_t \in \partial \|\beta_t\|_1.$$
(5)

Linearized Bregman Iteration as forward Euler discretization proposed even earlier than ISS dynamics (Osher-Burger-Goldfarb-Xu-Yin'05, Yin-Osher-Goldfarb-Darbon'08): for $\rho_k \in \partial \|\beta_k\|_1$,

$$\rho_{k+1} + \frac{1}{\kappa}\beta_{k+1} = \rho_k + \frac{1}{\kappa}\beta_k + \frac{\alpha_k}{n}X^{\mathsf{T}}(y - X\beta_k), \tag{6}$$

where

- Damping factor: $\kappa > 0$
- Step size: $\alpha_k > 0$ s.t. $\alpha_k \kappa \|\Sigma_n\| \le 2$
- Moreau Decomposition: $z_k := \rho_k + \frac{1}{\kappa}\beta_k \Leftrightarrow \beta_k = \kappa \cdot Shrink(z_k, 1)$

T)	- 1	۰.	

Algorithm

Variable Splitting

Summary

Linearized Bregman Iteration

Easy for Parallel Implementation



Figure: Linear speed-ups on a 16-core machine with synchronized parallel computation of matrix-vector products.

	Algorithm	Variable Splitting	
	00000 00000000		

Comparison with ISTA

Linearized Bregman (LB) iteration:

$$z_{t+1} = z_t - \alpha_t X^{\mathsf{T}}(\kappa X \operatorname{Shrink}(z_t, 1) - y)$$

which is not ISTA:

$$z_{t+1} = \frac{Shrink(z_t - \alpha_t X^T(Xz_t - y), \lambda)}{\lambda}.$$

Comparison:

• ISTA:

• as $t \to \infty$ solves LASSO: $\frac{1}{n} ||y - X\beta||_2^2 + \lambda ||\beta||_1$

- parallel run ISTA with $\{\lambda_k\}$ for LASSO regularization paths
- LB: a single run generates the whole regularization path at same cost of ISTA-LASSO estimator for a fixed regularization

	Algorithm
	000000000000

LB generates regularization paths





LB κ=64



Figure: As $\kappa \to \infty$, LB paths have a limit as piecewise-constant ISS path

	Algorithm	
0000000000000	000000000000000000000000000000000000000	

Accuracy: LB may be less biased than LASSO



- Left shows (the magnitudes of) nonzero entries of β^* .
- Middle shows the regularization path of LB.
- Right shows the regularization path of LASSO vs. $t = 1/\lambda$.

	Algorithm	
	000000000000000000000000000000000000000	

Path Consistency in Discrete Setting

Theorem (Osher-Ruan-Xiong-Y.-Yin'2016)

Assume that κ is large enough and α is small enough, with $\kappa \alpha \|X_S^*X_S\| < 2$,

$$\overline{\tau} := \frac{(1 - B/\kappa\eta)\eta}{2\sigma} \sqrt{\frac{n}{\log p}} \left(\max_{j \in T} \|X_j\| \right)^{-1}$$
$$\beta_{max}^* + 2\sigma \sqrt{\frac{\log p}{\gamma n}} + \frac{\|X\beta^*\|_2 + 2s\sqrt{\log n}}{n\sqrt{\gamma}} \triangleq B \le \kappa\eta,$$

then all the results for ISS can be extended to the discrete algorithm.

Note: it recovers the previous theorem as $\kappa \to \infty$ and $\alpha \to 0$, so LB can be less biased than LASSO.

	Algorithm	Variable Splitting	
	000000000000000000000000000000000000000		
Generalizations			

General Loss and Regularizer

$$\dot{\eta}_t = -\frac{\kappa_0}{n} \sum_{i=1}^n \nabla_{\eta} \ell(\mathbf{x}_i, \theta_t, \eta_t)$$
(7a)

$$\dot{\rho}_t + \frac{\dot{\theta}_t}{\kappa_1} = -\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \ell(\mathbf{x}_i, \theta_t, \eta_t)$$
(7b)

$$\rho_t \in \partial \|\theta_t\|_* \tag{7c}$$

where

- ℓ(x_i, θ) is a loss function: negative logarithmic likelihood, non-convex loss (neural networks), etc.
- $\|\theta_t\|_*$ is the Minkowski-functional (gauge) of dictionary convex hulls:

 $\|\theta\|_* := \inf\{\lambda \ge 0 : \theta \in \lambda K\}, \quad K \text{ is a symmetric convex hull of } \{a_i\}$

• it can be generalized to non-convex regularizers

	Algorithm	Variable Splitting	
	00000 0 000000		
Companylingstinger			

Linearized Bregman Iteration Algorithms

Differential inclusion (7) admits the following Euler Forward discretization

$$\eta_{t+1} = \eta_t - \frac{\alpha_k \kappa_0}{n} \sum_{i=1}^n \nabla_\eta \ell(x_i, \theta_t, \eta_t)$$
(8a)

$$z_{t+1} = z_t - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla_{\theta} \ell(x_i, \theta_t, \eta_t)$$
(8b)

$$\theta_{t+1} = \kappa_1 \cdot \operatorname{prox}_{\|\cdot\|_*}(z_{t+1}) \tag{8c}$$

where (8c) is given by Moreau Decomposition with

$$\operatorname{prox}_{\|\cdot\|_{*}}(z_{t}) = \arg\min_{x} \frac{1}{2} \|x - z_{t}\|^{2} + \|x\|_{*},$$

and

- $\alpha_k > 0$ is step-size while $\alpha_k \kappa_i \|\nabla^2_{\theta} \hat{\mathbb{E}} \ell(x, \theta)\| < 2$
- as simple as ISTA, easy to parallel implementation

()		•

LASSO vs. Differential Inclusions

Algorithm

Variable Splitting

Summary

Cran R package: Libra

Cran R package: Libra

http://cran.r-project.org/web/packages/Libra/

CRAN - Package Libra ×
 https://cran.r-project.org/web/packages/Libra/index.html

Libra: Linearized Bregman Algorithms for Generalized Linear Models

Efficient procedures for fitting the regularization path for linear, hinomial, multinomial, Ising and Potts models with lasso, group lasso er ochum lassofoully for multinomial pleamly. The package uses Linearized Bregman Algorithm to solve the regularization path through literations. Bregman Inverse Scale Space Differential Inclusion solver is also provided for linear model with lasso penalty.

Version:	1.5
Depends:	R (≥ 3.0), nnls
Suggests:	lars, MASS, igraph
Published:	2016-02-17
Author:	Feng Ruan, Jiechao Xiong and Yuan Yao
Maintainer:	Jiechao Xiong <xiongjiechao at="" pku.edu.cn=""></xiongjiechao>
License:	GPL-2
URL:	http://arxiv.org/abs/1406.7728
NeedsCompilation:	yes
SystemRequirements:	GNU Scientific Library (GSL)
CRAN checks:	Libra results







Reference manual:	Libra.pdf
Package source:	Libra 1.5.tar.gz
Windows binaries:	r-devel: Libra 1.5.zip, r-release: Libra 1.5.zip, r-oldrel: Libra 1.5.zip
OS X Snow Leopard binaries:	r-release: Libra 1.5.tgz, r-oldrel: not available
OS X Mavericks binaries:	r-release: Libra 1.5.tgz
Old sources:	Libra archive



	Algorithm	
	000000000000000000000000000000000000000	

Cran R package: Libra

Libra (1.5) currently includes

Sparse statistical models:

- linear regression: ISS (differential inclusion), LB
- logistic regression (binomial, multinomial): LB
- graphical models (Gaussian, Ising, Potts): LB

Two types of regularization:

- LASSO: *I*₁-norm penalty
- Group LASSO: $I_2 I_1$ penalty

	Algorithm	Variable Splitting	
	000000000000000		

Cran R package: Libra

A logistic regression with early stopping regularization



Figure: Peter Hall vs. other COPSS award winners in sparse logistic regression [papers from AoS/JASA/Biometrika/JRSSB, 2003-2012]: true coauthors are merely Tony Cai, R.J. Carroll, and J. Fan

	Algorithm	
	00000000000000	

Cran R package: Libra

Early stopping against overfitting in sparse Ising model learning



a true Ising model of 2-D grid

a movie of LB path

Algorithm

Variable Splitting

Summary

Cran R package: Libra

Example: Dream of the Red Mansion (Xueqin Cao vs. E. Gao)



Figure: Left: main characters net in the first 80 chapters at sparsity 10%; Right: the remaining 40 chapters.

		Algorithm ○○○○○○○○○○○○○○	Variable Splitting 00000	
Cran R package:	Libra			
	М	ore reference		

- Logistic Regression: loss conditional likelihood, regularizer l₁ (Shi-Yin-Osher-Saijda'10,Huang-Yao'18)
- Graphical Models (Gaussian/Ising/Potts Model): loss likelihood, composite conditional likelihood, regularizer – l₁ and group l₁ (Huang-Yao'18)
- Fused LASSO/TV: split Bregman with composite *l*₂ loss and *l*₁ gauge (Osher-Burger-Goldfarb-Xu-Yin'06, Burger-Gilboa-Osher-Xu'06, Yin-Osher-Goldfarb-Darbon'08, Huang-Sun-Xiong-Yao'16)
- Matrix Completion/Regression: gauge the matrix nuclear norm (Cai-Candès-Shen'10)

Split LB vs. Generalized LASSO

Structural Sparse Regression:

$$y = X\beta^{\star} + \epsilon, \ \gamma^{\star} = D\beta^{\star} \ (S = \operatorname{supp}(\gamma^{\star}), \ s = |S| \ll p),$$
(9)

Loss that splits prediction vs. sparsity control

$$\ell(\beta,\gamma) := \frac{1}{2n} \|y - X\beta\|_2^2 + \frac{1}{2\nu} \|\gamma - D\beta\|_2^2 \quad (\nu > 0).$$
 (10)

Split LBI:

$$\beta_{k+1} = \beta_k - \kappa \alpha \nabla_\beta \ell(\beta_k, \gamma_k), \tag{11a}$$

$$z_{k+1} = z_k - \alpha \nabla_{\gamma} \ell(\beta_k, \gamma_k), \tag{11b}$$

$$\gamma_{k+1} = \kappa \cdot \operatorname{prox}_{\|\cdot\|_1}(z_{k+1}), \tag{11c}$$

Generalized LASSO (genlasso):

$$\arg\min_{\beta} \left(\frac{1}{2n} \left\| y - X\beta \right\|_{2}^{2} + \lambda \left\| D\beta \right\|_{1} \right).$$
(12)

Algorithm

Variable Splitting

Summary

Split LBI vs. Generalized LASSO paths



Split LB may beat Generalized LASSO in Model Selection

genlasso		Split LBI		genlasso		Split LBI	
	u = 1	$\nu = 5$	u = 10	_	u = 1	u = 5	u = 10
.9426 (.0390)	.9845 (.0185)	.9969 (.0065)	.9982 (.0043)	.9705 (.0212)	.9955 (.0056)	.9996 (.0014)	.9998 (.0009)

• Example: n = p = 50, $X \in \mathbb{R}^{n \times p}$ with $X_j \sim N(0, I_p)$, $\epsilon \sim N(0, I_n)$

- (Right) 1-D fused (generalized) LASSO vs. Split LB (next page).
- In terms of Area Under the ROC Curve (AUC), LB has less false discoveries than genlasso
- *Why*? Split LB may need weaker irrepresentable conditions than generalized LASSO...

Outline LASSC	vs. Differential Inclusions	Variable Splitting	

Structural Sparsity Assumptions

- Define $\Sigma(\nu) := (I D(\nu X^*X + D^TD)^{\dagger}D^T)/\nu$.
- Assumption 1: Restricted Strong Convexity (RSC).

$$\Sigma_{S,S}(\nu) \succeq \lambda \cdot I. \tag{13}$$

• Assumption 2: Irrepresentable Condition (IRR).

$$\operatorname{IRR}(\nu) := \|\Sigma_{\mathcal{S}^{c},\mathcal{S}}(\nu) \cdot \Sigma_{\mathcal{S},\mathcal{S}}^{-1}(\nu)\|_{\infty} \le 1 - \eta.$$
(14)

- ν → 0: RSC and IRR above reduce to the RSC and IRR neccessary and sufficient for consistency of genlasso (Vaiter'13,LeeSunTay'13).
- $\nu \neq 0$: by allowing variable splitting in proximity, IRR above can be weaker than literature, bringing better variable selection consistency than genlasso (observed before)!

	Algorithm	Variable Splitting	

Identifiable Condition (IC) and Irrepresentable Condition (IRR)

• Let the columns of W form an orthogonal basis of $\ker(D_{S^c})$.

$$\Omega^{S} := \left(D_{S^{c}}^{\dagger} \right)^{T} \left(X^{*} X W \left(W^{T} X^{*} X W \right)^{\dagger} W^{T} - I \right) D_{S}^{T}, \qquad (15)$$

$$\operatorname{IC}_{0} := \left\| \Omega^{S} \right\|_{\infty}, \ \operatorname{IC}_{1} := \min_{u \in \operatorname{ker}(D_{S^{c}})} \left\| \Omega^{S} \operatorname{sign}\left(D_{S} \beta^{\star} \right) - u \right\|_{\infty}.$$
(16)

- The sign consistency of genlasso has been proved, under $\rm IC_1 < 1$ (Vaiter et al. 2013).
- We will show the sign consistency of Split LBI, under $IRR(\nu) < 1$.
- If $IRR(\nu) < IC_1$, then our IRR is easier to be met?

Outline

LASSO vs. Differential Inclusio

Algorithm 0000000000000 Variable Splitting

Summary

A Weaker Irrepresentable/Incoherence Condition

Split LB improves Irrepresentable Condition (Huang-Sun-Xiong-Y.'16)



Theorem (Huang-Sun-Xiong-Y.'2016)

- $IC_0 \ge IC_1$.
- $\operatorname{IRR}(\nu) \to \operatorname{IC}_0 (\nu \to 0).$
- $\operatorname{IRR}(\nu) \to C \ (\nu \to \infty). \ C = 0 \iff \ker(X) \subseteq \ker(D_S).$

		Algorithm 00000000000000	Variable Splitting ○●○○○			
A Weaker Irrepr	A Weaker Irrepresentable/Incoherence Condition					
		Consistency				

Theorem (Huang-Sun-Xiong-Y.'2016)

Under RSC and IRR, with large κ and small δ , there exists K such that with high probability, the following properties hold.

- No-false-positive property: γ_k (k ≤ K) has no false-positive, i.e. supp(γ_k) ⊆ S = supp(γ^{*}).
- Sign consistency of γ_k: If γ^{*}_{min} := min(|γ^{*}_j| : j ∈ S) (the minimal signal) is not weak, then supp(γ_K) = supp(γ^{*}).
- ℓ_2 consistency of γ_k : $\|\gamma_K \gamma^*\|_2 \le C_1 \sqrt{s \log m/n}$.
- ℓ_2 "consistency" of β_k : $\|\beta_K \beta^\star\|_2 \leq C_2 \sqrt{s \log m/n} + C_3 \nu$.
- Issues due to variable splitting (despite benefit on IRR):
 - $D\beta_K$ does not follow the sparsity pattern of $\gamma^{\star} = D\beta^{\star}$.
 - β_K incurs an additional loss $C_3\nu$ ($\nu \sim \sqrt{s \log m/n}$ minimax optimal).

		Variable Splitting	
		00000	
A Weaker Irrepres	sentable/Incoherence Condition		

Consistency

Theorem (Huang-Sun-Xiong-Y.'2016)

Define

$$\tilde{\beta}_{k} := \operatorname{Proj}_{\ker(D_{S_{k}^{c}})}(\beta_{k}) \ (S_{k} = \operatorname{supp}(\gamma_{k}))$$
(17)

Under RSC and IRR, with large κ and small δ , there exists K such that with high probability, the following properties hold, if γ_{\min}^{\star} is not weak.

- Sign consistency of $D\tilde{\beta}_{\kappa}$: supp $(D\tilde{\beta}_{\kappa}) = \operatorname{supp}(D\beta^*)$.
- ℓ_2 consistency of $\tilde{\beta}_{\kappa}$: $\left\|\tilde{\beta}_{\kappa} \beta^*\right\|_2 \leq C_4 \sqrt{s \log m/n}.$

LASSO vs. Differential Inclusion

Algorithm

Variable Splitting

Summary

A Weaker Irrepresentable/Incoherence Condition

Application: Alzheimer's Disease Detection



Figure: [Sun-Hu-Y.-Wang'17] A split of prediction (β) vs. interpretability ($\tilde{\beta}$): $\tilde{\beta}$ corresponds to the degenerate voxels interpretable for AD, while β additionally leverages the procedure bias to improve the prediction

LASSO vs. Differential Inclusio

Algorithm

Variable Splitting

Summary

A Weaker Irrepresentable/Incoherence Condition

Application: Partial Order of Basketball Teams



Figure: Partial order ranking for basketball teams. Top left shows $\{\beta_{\lambda}\}$ $(t = 1/\lambda)$ by genlasso and $\tilde{\beta}_k$ $(t = k\alpha)$ by Split LBI. Top right shows the same grouping result just passing t_5 . Bottom is the FIBA ranking of all teams.

	Algorithm 00000000000000	Variable Splitting 00000	Summary

Summary

We have seen:

- The limit of Linearized Bregman iterations follows a restricted gradient flow: differential inclusions dynamics
- It passes the unbiased Oracle Estimator under sign-consistency
- Sign consistency under nearly the same condition as LASSO
 - Restricted Strongly Convex + Irrepresentable Condition
- Split extension: sign consistency under a weaker condition than generalized LASSO
 - under a provably weaker Irrepresentable Condition
- Early stopping regularization is exploited against overfitting under noise

A Renaissance of Boosting as restricted gradient descent ...

000000000000

ns Algorithm 000000 Variable Splittin

Some Reference

- Osher, Ruan, Xiong, Yao, and Yin, "Sparse Recovery via Differential Equations", Applied and Computational Harmonic Analysis, 2016
- Xiong, Ruan, and Yao, "A Tutorial on Libra: R package for Linearized Bregman Algorithms in High Dimensional Statistics", Handbook of Big Data Analytics, Eds. by Wolfgang Karl Härdle, Henry Horng-Shing Lu, and Xiaotong Shen, Springer, 2017
- Xu, Xiong, Cao, and Yao, "False Discovery Rate Control and Statistical Quality Assessment of Annotators in Crowdsourced Ranking", ICML 2016, arXiv:1604.05910
- Huang, Sun, Xiong, and Yao, "Split LBI: an iterative regularization path with structural sparsity", NIPS 2016, https://github.com/yuany-pku/split-lbi
- Sun, Hu, Wang, and Yao, "GSplit LBI: taming the procedure bias in neuroimaging for disease prediction", MICCAI 2017
- Huang and Yao, "A Unified Dynamic Approach to Sparse Model Selection", AISTATS 2018
- Huang, Sun, Xiong, and Yao, "Boosting with Structural Sparsity: A Differential Inclusion Approach", Applied and Computational Harmonic Analysis, 2018, arXiv: 1704.04833
- R package:
 - http://cran.r-project.org/web/packages/Libra/index.html