

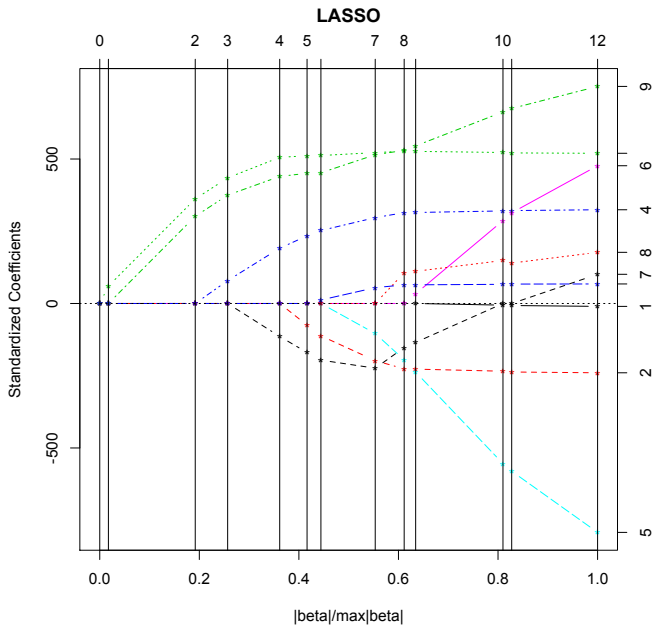
Interplay between Statistics and Optimization

Hui Zou

University of Minnesota

SAMSI

August 29, 2016



- Microarray data in early 2000s
 - Large scale multiple testing—false discovery rate (FDR) control (Benjamini and Hochberg), local FDR (Efron), Higher criticism (Donoho and Jin), SAM (Tibshirani)
 - Regression analysis—Least Angle Regression (Efron, Hastie, Johnstone and Tibshirani), Lasso
- Various penalization techniques: SCAD, MCP, Elastic Net, Adaptive Lasso, fused Lasso, group Lasso,...
- More sophisticated models/problems: GLM, GAM, precision matrix estimation, covariance matrix estimation, ...
- Compressed sensing, Matrix completion, Robust PCA
- Tensor regression, Tensor completion, Tensor decomposition

My personal view

- Optimization for Statistics: model fitting, model formulation, theoretical analysis
- Statistics for Optimization: new research thrusts

Today's talk:

- Majorization-Minimization (MM)
- Alternating Direction Method of Multipliers (ADMM)

Majorization-Minimization

Solve $\operatorname{argmin}_{\theta} C(\theta)$

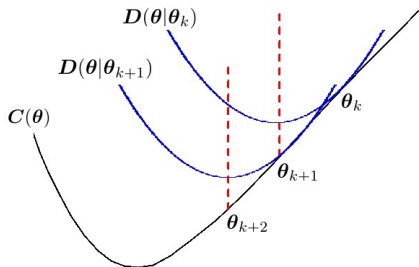
Majorization step:

$$C(\theta) < D(\theta|\theta_k) \text{ for any } \theta \neq \theta_k,$$

$$C(\theta_k) = D(\theta_k|\theta_k)$$

Minimization step:

$$\theta_k \leftarrow \theta_{k+1} = \operatorname{argmin}_{\theta} D(\theta|\theta_k)$$



Lange, Hunter and Yang (2000)[optimization transfer]; Hunter & Lange (2000) [MM]; Wu and Lange (2010) [EM and MM]

LLA

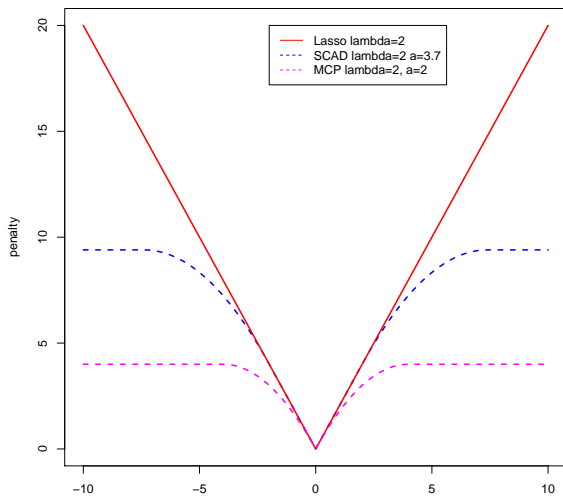
Zou and Li (2008)

Fan, Xue and Zou (2014)

Nonconvex penalized regression

$$\min_{\beta} \ell_n(\beta) + \sum_j P_\lambda(|\beta_j|)$$

- ℓ_n is convex and represents the statistical inference model
 - least squares loss
 - Huber 's M loss or least absolute loss
 - logistic regression: negative log-Bernoulli-likelihood
 - quantile regression: check loss
 - Ising model: composite conditional likelihood (Xue, Zou and Cai, 2012)
- $P_\lambda(t)$ is a non-decreasing concave function for $t \in (0, \infty)$
 - L_q norm penalty ($0 < q < 1$)
 - SCAD (Fan and Li, 2001)
 - MCP (Zhang, 2010)



LLA

$$\operatorname{argmin}_{\boldsymbol{\beta}} \left\{ \ell_n(\boldsymbol{\beta}) + \sum_{j=1}^p P_{\lambda}(\beta_j) \right\}$$

- 1 Start with some initial estimator $\boldsymbol{\beta}^{(0)}$.
- 2 At step k , define

$$Q_{\lambda}(\beta_j) = P_{\lambda}(|\beta_j^{(k)}|) + P'_{\lambda}(|\beta_j^{(k)}| +)(|\beta_j| - |\beta_j^{(k)}|)$$

- 3 Solve $\boldsymbol{\beta}^{(k+1)} = \operatorname{argmin}_{\boldsymbol{\beta}} \left\{ \ell_n(\boldsymbol{\beta}) + \sum_{j=1}^d Q_{\lambda}(\beta_j) \right\}$.

Iterate between Steps 2 and 3.

LLA and EM

Condition on P_λ : if there is a positive function $H(t)$ such that

$$\exp(-nP_\lambda(|\beta|)) = \int_0^\infty H(t)e^{-t|\beta|} dt. \quad (*)$$

Let $\pi(t) = \frac{2}{t}H(\frac{1}{t})$ and $p(\beta_j|\tau_j) = \frac{1}{2\tau_j}e^{-\frac{|\beta_j|}{\tau_j}}$. Then $(*)$ yields

$$\exp(-nP_\lambda(|\beta_j|)) = \int_0^\infty p(\beta_j|\tau_j)\pi(\tau_j)d\tau_j. \quad (**)$$

$(**)$ represents a hierarchical Bayesian model and suggests an EM algorithm for maximizing the penalized likelihood by treating τ_j s as “missing values”.

Under condition $(*)$ EM=LLA.

The issue of multiple local solutions

- The folded concave penalization problem usually has multiple local solutions, but the theory (namely, the oracle property) is established only for one of the unknown local solutions (Fan and Li, 2001; Fan and Peng, 2004; Lv and Fan, 2008; Fan and Lv, 2011; ...).
- Over a decade, the challenging fundamental issue still remains that it is not clear whether the local optimal solution computed by a given optimization algorithm possesses those nice theoretical properties.

Numeric demonstration

Simulation model: $y \sim \text{Bernoulli}\left(\frac{\exp(X\beta^*)}{1+\exp(X\beta^*)}\right)$, where $X \sim N_p(0, \Sigma)$
with $\Sigma_{ij} = 0.5^{|i-j|}$ and $\beta^* = (3, 1.5, 0, 0, 2, 0_{p-5})$.

		$n = 200 \ \& \ p = 1000$			
		ℓ_1 loss	ℓ_2 loss	# FP	# FN
		Sparse logistic regression			
Lasso		5.67	2.37	24.02	0.04
		(0.05)	(0.02)	(0.44)	(0.01)
SCAD-CD		4.50	2.13	13.99	0.08
		(0.06)	(0.02)	(0.31)	(0.01)
SCAD-LLA-zero		2.16	1.32	0.31	0.22
		(0.11)	(0.06)	(0.05)	(0.02)
SCAD-LLA-Lasso		2.08	1.28	0.26	0.19
		(0.10)	(0.06)	(0.04)	(0.02)

LLA closes the theoretical gap

In Fan, Xue and Zou (2014) it is shown that

Theorem

- ★ *If the initial estimator is Lasso, then the two-step LLA procedure finds the oracle solution with high probability.*
- ★ *If the initial estimator is zero, then the three-step LLA procedure finds the oracle solution with high probability.*

As illustration, the theory is verified for penalized least squares, penalized logistic regression, penalized quantile regression and penalized graphical model estimation.

The philosophical root of our theory

- In the classical MLE theory, when the log-likelihood function is not concave, one of the local maximizers of the log-likelihood function is shown to be asymptotic efficient, but how to compute that estimator is very challenging and often unclear.
- Le Cam (1956) (and later Bickel 1975) overcame this technical difficulty by focusing on a specially designed one-step Newton-Raphson estimator initialized by a root-n estimator.
- Le Cam did not try to get the global maximizer nor the theoretical local maximizer of the likelihood.

The search for the global minimizer

Mixed Integer Programming has been used to get the global minimizer of L_0 penalized and SCAD penalized least squares.

- Dimitris Bertsimas, Angela King and Rahul Mazumder (2016, AoS). Best Subset Selection via a Modern Optimization Lens.
- Hongcheng Liu, Tao Yao, Runze Li (2016). Global solutions to folded concave penalized nonconvex learning. AoS, 44(2), 629-659.

Extension to more general models?

BMD

Yang and Zou (2013)

Coordinate descent for lasso

$$\operatorname{argmin}_{\beta_1, \dots, \beta_p} f(\beta_1, \dots, \beta_p) + \sum_{j=1}^p \lambda |\beta_j|$$

- 1 Initialization of $\tilde{\beta}$
- 2 Cyclic coordinate descent: for $j = 1, 2, \dots, p, 1, 2, \dots$, update β_j by minimizing the objective function

$$\tilde{\beta}_j^{\text{update}} \leftarrow \operatorname{argmin}_{\beta_j} f(\tilde{\beta}_1, \dots, \tilde{\beta}_{j-1}, \beta_j, \tilde{\beta}_{j+1}, \tilde{\beta}_p) + \lambda |\beta_j|$$

- 3 Repeat (2) till convergence.

Lasso regression: $f(\beta_1, \dots, \beta_p) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$

$$\tilde{\beta}_j^{update} \leftarrow \underset{\beta_j}{\operatorname{argmin}} \left\| \mathbf{y} - \sum_{k \neq j} x_k \tilde{\beta}_k - x_j \beta_j \right\|_2^2 + \lambda |\beta_j|$$

reduces to soft-thresholding.

- Fu (1998) proposed the algorithm named “shooting”.
- Friedman, Hastie and Tibshirani (2008) glmnet, the same CD but with **clever implementation tricks such as active set, warm start and later strong rule**.
- For lasso logistic regression, Friedman, Hastie and Tibshirani (2008) did CD within a Newton-Raphson loop. Genkin, Lewis and Madigan (2007) did the standard CD by solving the one-dimensional optimization repeatedly.

Group lasso regression

$$\min_{(\beta_0, \boldsymbol{\beta})} \left[\frac{1}{2} \left\| \mathbf{y} - \beta_0 - \sum_k \mathbf{X}_{(k)} \boldsymbol{\beta}^{(k)} \right\|_2^2 + \lambda \sum_{k=1}^K \sqrt{p_k} \|\boldsymbol{\beta}^{(k)}\|_2 \right].$$

- Group Lasso penalty was introduced in Turlach, Vanebles and Wright (2004) and Yuan and Lin (2006).
- A blockwise descent algorithm under **a groupwise orthonormal condition: $\mathbf{X}_{(k)}$ columns are orthonormal.**
- **Orthonormal condition is incompatible with cross-validation, bootstrap, sub-sampling.**

A general group lasso problem

$$\arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n \tau_i \Phi(y_i, \beta^T \mathbf{x}_i) + \lambda \sum_{k=1}^K w_k \|\beta^{(k)}\|_2$$

where $\tau_i \geq 0$ and $w_k \geq 0$ for all i, k .

- The observation weights τ_i s are introduced in order to cover methods such as weighted regression and weighted large margin classification (biased sampling, unequal cost classification).
- The penalty weights w_k s make a more flexible model. The default choice for w_k is $\sqrt{p_k}$. If we do not want to penalize a group of predictors, simply let the corresponding weight be zero.

Loss functions

- Least squares: $\Phi(y, f) = \frac{1}{2}(y - f)^2$
- Logistic regression: $\Phi(y, f) = \log(1 + e^{-yf})$, $y = \pm 1$
- Squared hinge loss: $\Phi(y, f) = [(1 - yf)_+]^2$, $y = \pm 1$
- Huberized SVM loss: $\Phi(y, f) = \text{hsvm}(yf)$, $y = \pm 1$ where

$$\text{hsvm}(t) = \begin{cases} 0, & t > 1 \\ (1 - t)^2 / 2\delta, & 1 - \delta < t \leq 1 \\ 1 - t - \delta/2, & t \leq 1 - \delta. \end{cases}$$

Let \mathbf{D} denote the data $\{\mathbf{y}, \mathbf{X}\}$ and define

$$L(\boldsymbol{\beta} \mid \mathbf{D}) = \frac{1}{n} \sum_{i=1}^n \tau_i \Phi(y_i, \boldsymbol{\beta}^\top \mathbf{x}_i).$$

Definition

The loss function Φ is said to satisfy the QM condition, if

- (i). $\nabla L(\boldsymbol{\beta} \mid \mathbf{D})$ exists everywhere.
- (ii). There exists a $p \times p$ matrix \mathbf{H} , which may only depend on the data \mathbf{D} , such that for all $\boldsymbol{\beta}, \boldsymbol{\beta}^*$,

$$\begin{aligned} L(\boldsymbol{\beta} \mid \mathbf{D}) &\leq L(\boldsymbol{\beta}^* \mid \mathbf{D}) + (\boldsymbol{\beta} - \boldsymbol{\beta}^*)^\top \nabla L(\boldsymbol{\beta}^* \mid \mathbf{D}) \\ &\quad + \frac{1}{2} (\boldsymbol{\beta} - \boldsymbol{\beta}^*)^\top \mathbf{H} (\boldsymbol{\beta} - \boldsymbol{\beta}^*). \end{aligned}$$

Loss	$-\nabla L(\boldsymbol{\beta} \mid \mathbf{D})$	\mathbf{H}
Least squares	$\frac{1}{n} \sum_{i=1}^n \tau_i (y_i - \mathbf{x}_i^\top \boldsymbol{\beta}) \mathbf{x}_i$	$\mathbf{X}^\top \Gamma \mathbf{X} / n$
Logistic regression	$\frac{1}{n} \sum_{i=1}^n \tau_i y_i \mathbf{x}_i \frac{1}{1 + \exp(y_i \mathbf{x}_i^\top \boldsymbol{\beta})}$	$\frac{1}{4} \mathbf{X}^\top \Gamma \mathbf{X} / n$
Squared hinge loss	$\frac{1}{n} \sum_{i=1}^n 2\tau_i y_i \mathbf{x}_i (1 - y_i \mathbf{x}_i^\top \boldsymbol{\beta})_+$	$4\mathbf{X}^\top \Gamma \mathbf{X} / n$
Huberized SVM loss	$\frac{1}{n} \sum_{i=1}^n \tau_i y_i \mathbf{x}_i \text{hsvm}'(y_i \mathbf{x}_i^\top \boldsymbol{\beta})$	$\frac{2}{\delta} \mathbf{X}^\top \Gamma \mathbf{X} / n$

$$\Gamma = \text{diag}(\tau_1, \dots, \tau_n)$$

Write β such that $\beta^{(k')} = \tilde{\beta}^{(k')}$ for $k' \neq k$.

Given $\beta^{(k')} = \tilde{\beta}^{(k')}$ for $k' \neq k$, the optimal $\beta^{(k)}$ is defined as

$$\underset{\beta^{(k)}}{\operatorname{argmin}} L(\beta \mid \mathbf{D}) + \lambda w_k \|\beta^{(k)}\|_2.$$

By QM condition,

$$L(\beta \mid \mathbf{D}) \leq L(\tilde{\beta} \mid \mathbf{D}) + (\beta - \tilde{\beta})^\top \nabla L(\tilde{\beta} \mid \mathbf{D}) + \frac{1}{2} (\beta - \tilde{\beta})^\top \mathbf{H} (\beta - \tilde{\beta}).$$

Write $U(\tilde{\beta}) = -\nabla L(\tilde{\beta} \mid \mathbf{D})$.

$$\begin{aligned} L(\beta \mid \mathbf{D}) &\leq L(\tilde{\beta} \mid \mathbf{D}) - (\beta^{(k)} - \tilde{\beta}^{(k)})^\top U^{(k)} \\ &\quad + \frac{1}{2} (\beta^{(k)} - \tilde{\beta}^{(k)})^\top \mathbf{H}^{(k)} (\beta^{(k)} - \tilde{\beta}^{(k)}). \end{aligned}$$

Let η_k be the largest eigenvalue of $\mathbf{H}^{(k)}$. We set $\gamma_k = (1 + 10^{-4})\eta_k$

$$L(\boldsymbol{\beta} \mid \mathbf{D}) \leq L(\tilde{\boldsymbol{\beta}} \mid \mathbf{D}) - (\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)})^\top \mathbf{U}^{(k)} + \frac{1}{2}\gamma_k \|(\boldsymbol{\beta}^{(k)} - \tilde{\boldsymbol{\beta}}^{(k)})\|_2^2 \quad (*)$$

"=" holds if only if $\boldsymbol{\beta}^{(k)} = \tilde{\boldsymbol{\beta}}^{(k)}$

The minimizer of the right hand side of (*) is

$$\tilde{\boldsymbol{\beta}}^{(k)}(\text{new}) = \frac{1}{\gamma_k} \left(\mathbf{U}^{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)} \right) \left(1 - \frac{\lambda w_k}{\|\mathbf{U}^{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)}\|_2} \right)_+.$$

The whole process drives the objective **strictly downhill** unless the optimal solution is reached (i.e., KKT conditions are satisfied).

BMD for group lasso

- For $k = 1, \dots, K$, compute γ_k , the largest eigenvalue of $\mathbf{H}^{(k)}$
- $\gamma_k = (1 + 10^{-4})\gamma_k$ (for nontrivial groups with size ≥ 2)
- Initialize $\tilde{\boldsymbol{\beta}}$.
- Repeat the following cyclic blockwise updates until convergence:
 - ★ for $k = 1, \dots, K$, do (1)–(3)
 - ★ (1) Compute $U(\tilde{\boldsymbol{\beta}}) = -\nabla L(\tilde{\boldsymbol{\beta}}|\mathbf{D})$.
 - ★ (2) Compute

$$\tilde{\boldsymbol{\beta}}^{(k)}(\text{new}) = \frac{1}{\gamma_k} \left(U^{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)} \right) \left(1 - \frac{\lambda w_k}{\|U^{(k)} + \gamma_k \tilde{\boldsymbol{\beta}}^{(k)}\|_2} \right)_+.$$
 - ★ (3) Set $\tilde{\boldsymbol{\beta}}^{(k)} = \tilde{\boldsymbol{\beta}}^{(k)}(\text{new})$.

gglasso package: also uses active set, strong rule and warm start.

Competitors

- block coordinate gradient descent `grplasso`: Meier et al. (2008) for group-lasso logistic regression.
- **ISTA-BC** algorithm: Qin et al. (2010), an extension of the ISTA/FISTA (Beck & Teboulle 2009) based on variable step-lengths.
- **SLEP** implemented Nesterov's method: Liu et al. (2009)

Dataset	Type	n	q	p	Data Source
<i>Autompg</i>	R	392	7	31	(Quinlan 1993)
<i>Bardet</i>	R	120	200	1000	(Scheetz et al. 2006)
<i>Cardiomyopathy</i>	R	30	6319	31595	(Segal et al. 2003)
<i>Spectroscopy</i>	R	103	100	500	(Sabo et al. 2008)
<i>Breast</i>	C	42	22283	111415	(Graham et al. 2010)
<i>Colon</i>	C	62	2000	10000	(Alon et al. 1999)
<i>Prostate</i>	C	102	6033	30165	(Singh et al. 2002)
<i>Sonar</i>	C	208	60	300	(Gorman et al. 1988)

Some real datasets. n is the number of instances. q is the number of original variables. p is the number of predictors after expansion. "R" means regression and "C" means classification.

Group-lasso GAM regression, timing performance

Dataset	<i>Autompg</i>	<i>Bardet</i>	<i>Cardiomyopathy</i>	<i>Spectroscopy</i>
SLEP	3.14	9.96	78.23	9.37
ISTA-BC	5.66	1.55	2.43	1.31
gglasso	2.51	0.77	2.48	0.76

All experiments were carried out on an Intel Xeon X5560 (Quad-core 2.8 GHz) processor.

 Group-lasso GAM classification, timing performance

Dataset	<i>Colon</i>	<i>Prostate</i>	<i>Sonar</i>	<i>Breast</i>
grplasso (Logit)	60.42	111.75	24.55	439.76
SLEP (Logit)	75.31	166.91	5.49	358.75
gglasso (Logit)	1.13	3.877	1.54	9.62
gglasso (HSVM)	1.15	3.53	0.66	9.15

All experiments were carried out on an Intel Xeon X5560 (Quad-core 2.8 GHz) processor.

A Small Trick

Yang and Zou (2012)

A counterintuitive phenomenon

Consider the glmnet for fitting elastic net penalized regression.

W.L.O.G. assume $\sum_{i=1}^N x_{ij} = 0$, $\frac{1}{N} \sum_{i=1}^N x_{ij}^2 = 1$, for $j = 1, \dots, p$.

$$R(\beta_0, \beta) = \frac{1}{2N} \sum_{i=1}^N (y_i - \beta_0 - x_i^\top \beta)^2 + P_{\lambda, \alpha}(\beta),$$

where $P_{\lambda, \alpha}(\beta)$ is the elastic net penalty

$$P_{\lambda, \alpha}(\beta) = \lambda \sum_{j=1}^p p_{\alpha}(\beta_j) = \lambda \sum_{j=1}^p \left[\frac{1}{2} (1 - \alpha) \beta_j^2 + \alpha |\beta_j| \right].$$

`glmnet` implements the standard CD algorithm in which we iteratively solve a univariate elastic net problem

$$\hat{\beta}_j = \arg \min_{\beta_j} R(\beta_j | \tilde{\beta}_0, \tilde{\beta}),$$

where

$$R(\beta_j | \tilde{\beta}_0, \tilde{\beta}) = \frac{1}{2} (\beta_j - \tilde{\beta}_j)^2 - \frac{1}{N} \sum_{i=1}^N r_i x_{ij} (\beta_j - \tilde{\beta}_j) + \lambda p_\alpha(\beta_j).$$

$$\hat{\beta}_j = \frac{S\left(\frac{1}{N} \sum_{i=1}^N x_{ij} r_i + \tilde{\beta}_j, \lambda \alpha\right)}{1 + \lambda(1 - \alpha)},$$

where $S(z, t) = (|z| - t)_+ \text{sgn}(z)$.

A tiny change to glmnet

We change the univariate update formula to

$$\hat{\beta}_j^B = \frac{S\left(\frac{1}{N} \sum_{i=1}^N x_{ij} r_i + f \cdot \tilde{\beta}_j, \lambda \alpha\right)}{f \cdot 1 + \lambda(1 - \alpha)} \quad (f \geq 1)$$

- Yang made a code error by using $f = 2$ in glmnet, but still got good/even better results.
- As long as $f \geq 1$ the iterative process converges to the desired solution.
- A bigger f means a smaller step size along each coordinate direction. For an orthogonal design, $f = 1$ is the best choice.

Simulation

FHT model:

We simulated data with N observations and p predictors where each pair of predictors X_j and $X_{j'}$ have the same population correlation ρ , with ρ ranges from zero to 0.95.

The response variable was generated by

$$Y = \sum_{j=1}^p X_j \beta_j + k \cdot N(0, 1),$$

where $\beta_j = (-1)^j \exp(-(2j - 1)/20)$ and k is set to make the signal-to-noise ratio equal 3.

We compared $f = 1$ (glmnet) and $f = 2$ (glmnet2).

		Correlation					
		0	0.1	0.2	0.5	0.8	0.95
		$\alpha = 1$					
		$N = 100, p = 5000$					
glmnet		0.2222	0.2339	0.2979	0.4606	0.7919	1.9016
glmnet2		0.2533	0.2519	0.2886	0.3758	0.5450	1.0735
		$\alpha = 0.5$					
		$N = 100, p = 5000$					
glmnet		0.2107	0.2189	0.2356	0.3669	0.7765	2.1528
glmnet2		0.2225	0.2285	0.2414	0.2861	0.4876	1.3335

A simple explanation

$$v_j = (0, \dots, \frac{x_j^\top y}{fN}, \dots, 0) \quad \mathbf{u}_j = (u_{kj})_{p \times 1} = \begin{cases} -\frac{1}{f} & k = j \\ -\frac{\rho}{f} & k \neq j \end{cases}$$

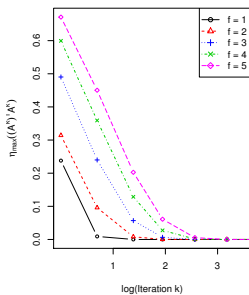
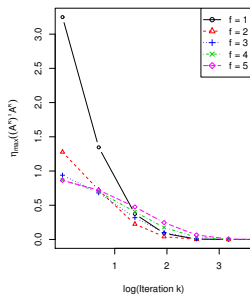
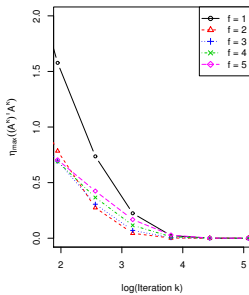
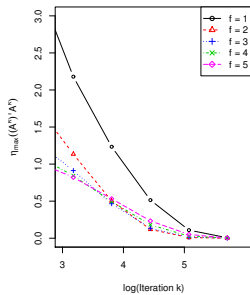
$$\mathbf{W}_j = \mathbf{I}_{p \times p} + \begin{bmatrix} \mathbf{0}_{p \times (j-1)} & \mathbf{u}_j & \mathbf{0}_{p \times (N-j)} \end{bmatrix}$$

$$\mathbf{A} = \prod_{j=1}^p \mathbf{W}_j \quad \mu = \sum_{s=1}^{p-1} \left(v_s \prod_{j=s+1}^p \mathbf{W}_j \right) + v_p$$

If apply the CD and CMD to the LS problem, after a complete cycle from $j = 1$ to $j = p$, we get

$$\tilde{\beta}^{(k)} = \tilde{\beta}^{(k-1)} \mathbf{A} + \mu$$

The convergence rate is basically the maximum eigenvalue of $(\mathbf{A}^k)^\top \mathbf{A}^k$, which is affected by both f and ρ .

(a) $\rho = 0.1$ (b) $\rho = 0.5$ (c) $\rho = 0.8$ (d) $\rho = 0.95$ 

	Colon	Prostate	WBCD	Ionosphere	Sonar
N	62	102	569	351	208
p	2000	6033	495 (30)	560 (32)	1890 (60)
α_{CV}	0.6	0.5	0.6	0.4	0.4
Test Error	8.3%	5%	1.77%	2.86%	24.39%
glmnet	0.1166	0.3283	9.4039	0.5158	2.0828
glmnet2	0.0910	0.2938	4.9593	0.3667	1.0945
Improv. %	+28%	+11.7%	+89.6%	+40.6%	+90.3%

ADMM

- Douglas & Rachford (1956); Lions & Mercier (1979); Eckstein & Bertsekas (1992)
- Goldstein & Osher (2009); Yin, Osher, Goldfard, and Darbon (2008); Goldfarb & Ma (2012)
- many applications in signal processing, statistics, machine learning

Improving MPT

Xue, Ma and Zou (2012)

An investor has p assets. Asset j makes up ω_j proportion of the investor's portfolio.

$$\omega_j \geq 0 \quad \sum_{j=1}^p \omega_j = 1.$$

Asset j delivers return R_j which has mean μ_j and variance σ_j^2 .

The mean of the return of the entire portfolio is $\sum_{j=1}^p \omega_j \mu_j$ and the variance of the portfolio's return is

$$\sum_i^p \sum_j^p \omega_i \omega_j \sigma_i \sigma_j \rho_{ij}$$

where ρ_{ij} is the correlation between R_i and R_j .

$$\mathbf{w} = (\omega_1, \dots, \omega_p)^\top, \boldsymbol{\mu} = (\mu_1, \dots, \mu_p)^\top,$$

Σ is the covariance matrix of return vector $(R_1, \dots, R_p)^\top$.

MPT

$$\underset{\mathbf{w}}{\operatorname{argmin}} \mathbf{w}^\top \Sigma \mathbf{w}$$

s.t.

$$\mathbf{w}^\top \boldsymbol{\mu} = \mu_p, \mathbf{w}^\top \vec{1} = 1,$$

$$\omega_j \geq 0, j = 1, \dots, p.$$

MPT (1952, J. of Finance) won 1990 Nobel Prize in Economics.

The usual implementation of MPT

$\hat{\boldsymbol{\mu}} = (\mu_1, \dots, \mu_p)^\top$ is the sample mean vector

$\hat{\boldsymbol{\Sigma}}_n$ is the sample covariance matrix

Empirical MPT

$$\hat{\boldsymbol{w}} = \arg \min_{\boldsymbol{w}} \boldsymbol{w}^\top \hat{\boldsymbol{\Sigma}}_n \boldsymbol{w}$$

s.t.

$$\boldsymbol{w}^\top \hat{\boldsymbol{\mu}} = \mu_p, \boldsymbol{w}^\top \vec{\mathbf{1}} = 1,$$

$$\omega_j \geq 0, j = 1, \dots, p.$$

When p is relatively large

- The sample cov. matrix performs poorly (Johnstone, 2001). It leads to bias and undesirable risk issues in the empirical MPT (El Karoui, 2010, Brodie et al., 2009; DeMiguel et al., 2009; Fan et al., 2012).
- Under some suitable “**sparsity**” assumption on Σ , an optimal estimator can be obtained by **Thresholding** (Bickel and Levina 2008a; El Karoui, 2008, Cai and Zhou, 2011)
- Let $\hat{\sigma}_{ij}$ be the ij entry of the sample covariance matrix.

$$\hat{\Sigma}_{\text{thresholding}} = \{s_{\lambda}(\hat{\sigma}_{ij})\}_{1 \leq i, j \leq p}$$

The difficulty is how to preserve both P.D. and Sparsity simultaneously.

Notation: $|\Sigma|_1 = \sum_{i \neq j} |\sigma_{ij}|$, $\|\Sigma\|_F^2 = \sum_{i,j} \sigma_{ij}^2$.

The soft-thresholding estimator is the global solution of

$$\operatorname{argmin}_{\Sigma} \frac{1}{2} \|\Sigma - \hat{\Sigma}_n\|_F^2 + \lambda |\Sigma|_1.$$

PSD sparse covariance estimator

$$\hat{\Sigma}^+ = \operatorname{argmin}_{\Sigma \succeq \epsilon \mathbf{I}} \frac{1}{2} \|\Sigma - \hat{\Sigma}_n\|_F^2 + \lambda |\Sigma|_1.$$

$\epsilon = 10^{-6}$. ϵ can be other positive constant depending on the application.

Algorithm

The augmented Lagrangian function for some given parameter μ ,

$$L(\Theta, \Sigma; \Lambda) = \frac{1}{2} \|\Sigma - \hat{\Sigma}_n\|_F^2 + \lambda \|\Sigma\|_1 - \langle \Lambda, \Theta - \Sigma \rangle + \frac{1}{2\mu} \|\Theta - \Sigma\|_F^2,$$

where Λ is the Lagrange multiplier.

For $i = 0, 1, 2, \dots$,

$$\Theta \text{ step : } \Theta^{i+1} = \arg \min_{\Theta \succeq \epsilon \mathbf{I}} L(\Theta, \Sigma^i; \Lambda^i)$$

$$\Sigma \text{ step : } \Sigma^{i+1} = \arg \min_{\Sigma} L(\Theta^{i+1}, \Sigma; \Lambda^i)$$

$$\Lambda \text{ step : } \Lambda^{i+1} = \Lambda^i - \frac{1}{\mu} (\Theta^{i+1} - \Sigma^{i+1}).$$

⊖ step

$$L(\Theta, \Sigma; \Lambda) = \frac{1}{2} \|\Sigma - \hat{\Sigma}_n\|_F^2 + \lambda |\Sigma|_1 - \langle \Lambda, \Theta - \Sigma \rangle + \frac{1}{2\mu} \|\Theta - \Sigma\|_F^2$$

$$\begin{aligned} \Theta^{i+1} &= \arg \min_{\Theta \succeq \epsilon \mathbf{I}} L(\Theta, \Sigma^i; \Lambda^i) \\ &= \arg \min_{\Theta \succeq \epsilon \mathbf{I}} -\langle \Lambda^i, \Theta \rangle + \frac{1}{2\mu} \|\Theta - \Sigma^i\|_F^2 \\ &= \arg \min_{\Theta \succeq \epsilon \mathbf{I}} \|\Theta - (\Sigma^i + \mu \Lambda^i)\|_F^2 \\ &= (\Sigma^i + \mu \Lambda^i)_+. \end{aligned}$$

Let \mathbf{Z} 's eigen-decomposition be $\sum_{j=1}^p \lambda_j \mathbf{v}_j^T \mathbf{v}_j$, then define

$$(\mathbf{Z})_+ = \sum_{j=1}^p \max(\lambda_j, \epsilon) \mathbf{v}_j^T \mathbf{v}_j.$$

Σ step

$$L(\Theta, \Sigma; \Lambda) = \frac{1}{2} \|\Sigma - \hat{\Sigma}_n\|_F^2 + \lambda |\Sigma|_1 - \langle \Lambda, \Theta - \Sigma \rangle + \frac{1}{2\mu} \|\Theta - \Sigma\|_F^2$$

$$\begin{aligned} \Sigma^{i+1} &= \arg \min_{\Sigma} L(\Theta^{i+1}, \Sigma; \Lambda^i) \\ &= \arg \min_{\Sigma} \frac{1}{2} \|\Sigma - \hat{\Sigma}_n\|_F^2 + \lambda |\Sigma|_1 + \langle \Lambda^i, \Sigma \rangle + \frac{1}{2\mu} \|\Sigma - \Theta^{i+1}\|_F^2 \\ &= \arg \min_{\Sigma} \frac{1}{2} \left\| \Sigma - \frac{\mu(\hat{\Sigma}_n - \Lambda^i) + \Theta^{i+1}}{1 + \mu} \right\|_F^2 + \frac{\lambda\mu}{1 + \mu} |\Sigma|_1 \\ &= \frac{1}{1 + \mu} \mathbf{S}(\mu(\hat{\Sigma}_n - \Lambda^i) + \Theta^{i+1}, \lambda\mu). \end{aligned}$$

Define $\mathbf{S}(\mathbf{Z}, \tau) = \{s(z_{j\ell}, \tau)\}_{1 \leq j, \ell \leq p}$ with

$$s(z_{j\ell}, \tau) = \text{sign}(z_{j\ell}) \max(|z_{j\ell}| - \tau, 0) I_{\{j \neq \ell\}} + z_{j\ell} I_{\{j = \ell\}}.$$

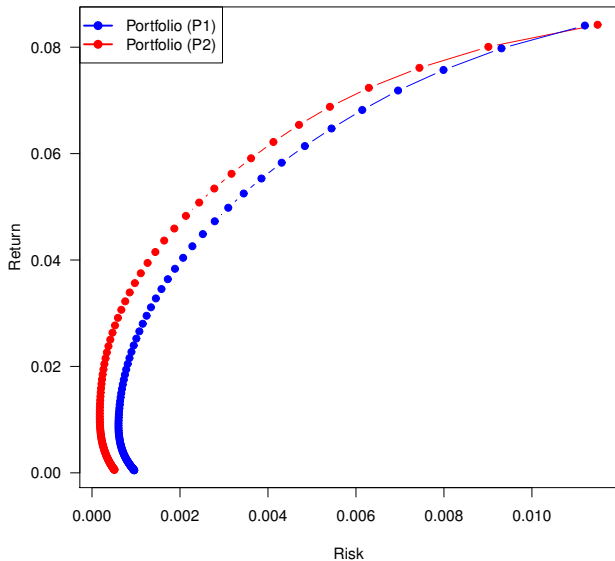
Improved Empirical MPT

$$\hat{w} = \arg \min_w w^T \hat{\Sigma}^+ w$$

s.t.

$$w^T \hat{\mu} = \mu_P, w^T \vec{1} = 1,$$

$$w_j \geq 0, j = 1, \dots, p.$$



Two MPT frontiers based on S&P 100 from Jan. 1990—Jan. 1993. Red: new; blue: traditional.

Latent Variable glasso

Ma, Xue and Zou (2013)

Latent variable Gaussian graphical model

- **observed \mathbf{X}** (p -dim.) and **unobserved \mathbf{Y}** (q -dim.) are jointly Gaussian

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \sim N_{p+q} \left(\begin{bmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_Y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_X & \boldsymbol{\Sigma}_{XY} \\ \boldsymbol{\Sigma}_{YX} & \boldsymbol{\Sigma}_Y \end{bmatrix} \right)$$

- **sparsity**: $\mathbf{X}|\mathbf{Y}$ has a sparse Gaussian graphical model representation.

-

$$(\boldsymbol{\Sigma})^{-1} = \boldsymbol{\Theta} = \begin{bmatrix} \boldsymbol{\Theta}_X & \boldsymbol{\Theta}_{XY} \\ \boldsymbol{\Theta}_{YX} & \boldsymbol{\Theta}_Y \end{bmatrix}$$

$\mathbf{X}|\mathbf{Y}$ is normal with precision matrix $\boldsymbol{\Theta}_X$.

- How to estimate $\boldsymbol{\Theta}_X$ just based on \mathbf{X} ?

A convex formulation

Chandrasekaran, Parrilo & Willsky (2012)

A key observation: $\Sigma_X^{-1} = \Theta_X - \Theta_{XY}\Theta_Y^{-1}\Theta_{YX}$

- Θ_X is sparse (assumption),
- Θ_{XY} has rank at most q , $\Theta_{XY}\Theta_Y^{-1}\Theta_{YX}$ ' rank is at most q .

If assume q is small (very reasonable in applications), we have a “sparse” -“low-rank” decomposition of Σ_X^{-1} —the marginal precision matrix of \mathbf{X} .

LVGM estimator

Write

$$\Sigma_X^{-1} = \mathbf{S} - \mathbf{L},$$

\mathbf{S} is a sparse PD matrix and \mathbf{L} is a low rank SPD matrix.

$$\begin{aligned} \min_{(\mathbf{S}, \mathbf{L})} \quad & \langle \hat{\Sigma}_X, \mathbf{S} - \mathbf{L} \rangle - \log \det(\mathbf{S} - \mathbf{L}) + \alpha \|\mathbf{S}\|_1 + \beta \text{Tr}(\mathbf{L}) \\ \text{subject to} \quad & \mathbf{S} - \mathbf{L} \succ 0, \mathbf{L} \succeq 0 \end{aligned}$$

$\|\mathbf{S}\|_1$ is a convex relaxation of the sparsity of \mathbf{S} . $\text{Tr}(\mathbf{L})$ is a convex relaxation of the rank of \mathbf{L} .

Chandrasekaran, Parrilo & Willsky (2012) viewed the above as a log-determinant semidefinite programming problem.

Algorithm

■ $\mathbf{R} = \mathbf{S} - \mathbf{L}$

$$\min_{(\mathbf{R}, \mathbf{S}, \mathbf{L})} \quad \langle \hat{\Sigma}_X, \mathbf{R} \rangle - \log \det(\mathbf{R}) + \alpha \|\mathbf{S}\|_1 + \beta \text{Tr}(\mathbf{L})$$

$$\text{subject to} \quad \mathbf{R} \succ 0, \mathbf{L} \succeq 0$$

■ augmented Lagrangian

$$\begin{aligned} \mathcal{L}(\mathbf{R}, \mathbf{S}, \mathbf{L}; \Lambda) &= \langle \hat{\Sigma}_X, \mathbf{R} \rangle - \log \det(\mathbf{R}) + \alpha \|\mathbf{S}\|_1 + \beta \text{Tr}(\mathbf{L}) \\ &\quad - \langle \Lambda, \mathbf{R} - \mathbf{S} + \mathbf{L} \rangle + \frac{1}{2\mu} \|\mathbf{R} - \mathbf{S} + \mathbf{L}\|_F^2. \end{aligned}$$

■ alternating minimization

$$\left\{ \begin{array}{l} \mathbf{R}^{k+1} = \arg \min_{\mathbf{R}} \mathcal{L}(\mathbf{R}, \mathbf{S}^k, \mathbf{L}^k; \Lambda^k) \\ \mathbf{S}^{k+1} = \arg \min_{\mathbf{S}} \mathcal{L}(\mathbf{R}^{k+1}, \mathbf{S}, \mathbf{L}^k; \Lambda^k) \\ \mathbf{L}^{k+1} = \arg \min_{\mathbf{L} \succeq 0} \mathcal{L}(\mathbf{R}^{k+1}, \mathbf{S}^{k+1}, \mathbf{L}; \Lambda^k) \\ \Lambda^{k+1} = \Lambda^k - \frac{1}{\mu} (\mathbf{R}^{k+1} - \mathbf{S}^{k+1} + \mathbf{L}^{k+1}) \end{array} \right.$$

R step

$$\arg \min_{\mathbf{R}} \langle \hat{\Sigma}_X, \mathbf{R} \rangle - \log \det(\mathbf{R}) - \langle \Lambda^k, \mathbf{R} - \mathbf{S}^k + \mathbf{L}^k \rangle + \frac{1}{2\mu} \|\mathbf{R} - \mathbf{S}^k + \mathbf{L}^k\|_F^2$$

$$\arg \min_{\mathbf{R}} - \log \det(\mathbf{R}) + \frac{1}{2\mu} \|\mathbf{R} - \mathbf{G}\|_F^2$$

$$\mathbf{G} = \mathbf{S}^k - \mathbf{L}^k - \mu(\hat{\Sigma}_X - \Lambda^k)$$

$$\mathbf{R} - \mathbf{G} - \mu\mathbf{R}^{-1} = 0$$

Let $\mathbf{G} = U^T \sigma U$ (eigen-decomposition of \mathbf{G})

$$\mathbf{R} = U^T \gamma U$$

with

$$\gamma_i = \frac{\sigma_i + \sqrt{\sigma_i^2 + 4\mu}}{2}$$

S step

$$\arg \min_{\mathbf{S}} \alpha \|\mathbf{S}\|_1 - \langle \mathbf{\Lambda}^k, \mathbf{R}^{k+1} - \mathbf{S} + \mathbf{L}^k \rangle + \frac{1}{2\mu} \|\mathbf{R}^{k+1} - \mathbf{S} + \mathbf{L}^k\|_F^2$$

$$\mathbf{S}^{k+1} = \arg \min_{\mathbf{S}} \mu\alpha \|\mathbf{S}\|_1 + \frac{1}{2} \|\mathbf{Z} - \mathbf{S}\|_F^2$$

$$\mathbf{Z} = (\mathbf{R}^{k+1} + \mathbf{L}^k - \mu\mathbf{\Lambda}^k)$$

$$\tau = \mu\alpha$$

$$\mathbf{S}_{ij}^{k+1} = [\text{Shrink}(\mathbf{Z}, \tau)]_{ij} := \begin{cases} Z_{ii}, & \text{if } i = j \\ Z_{ij} - \tau, & \text{if } i \neq j \text{ and } Z_{ij} > \tau \\ Z_{ij} + \tau, & \text{if } i \neq j \text{ and } Z_{ij} < -\tau \\ 0, & \text{if } i \neq j \text{ and } -\tau \leq Z_{ij} \leq \tau. \end{cases}$$

L step

The above is equivalent to

$$\arg \min_{\mathbf{L} \succeq 0} \beta \text{Tr}(\mathbf{L}) - \langle \boldsymbol{\Lambda}^k, \mathbf{R}^{k+1} - \mathbf{S}^{k+1} + \mathbf{L} \rangle + \frac{1}{2\mu} \|\mathbf{R}^{k+1} - \mathbf{S}^{k+1} + \mathbf{L}\|_F^2$$

$$\mathbf{L}^{k+1} = \arg \min_{\mathbf{L} \succeq 0} (\mu\beta) \text{Tr}(\mathbf{L}) + \frac{1}{2} \|\mathbf{M} - \mathbf{L}\|_F^2$$

where

$$\mathbf{M} = (\mathbf{S}^{k+1} - \mathbf{R}^{k+1} + \mu\boldsymbol{\Lambda}^k)$$

$\mathbf{M} = U^T \boldsymbol{\sigma} U$ (eigen-decomposition of \mathbf{M}) then

$$\mathbf{L}^{k+1} = SVT(\mathbf{M}, \mu\beta) = U^T \boldsymbol{\gamma} U$$

with

$$\gamma_i = \max(\sigma_i - \mu\beta, 0)$$

Concluding remark

- ★ Tailoring optimization algorithms to the specific statistics problem
- ★ Efforts to polish the solver

Thank You