

The not-so-rough landscape of nonconvex M -estimators

Po-Ling Loh

University of Wisconsin - Madison
Department of Statistics

KITP workshop on rough high-dimensional landscapes
UC Santa Barbara

January 7, 2019

- **Prediction/regression problem:** Observe $\{(x_i, y_i)\}_{i=1}^n$, estimate

$$\beta^* = \arg \min_{\beta \in \mathbb{R}^p} \mathbb{E}[\ell(\beta; x_i, y_i)], \quad x_i \in \mathbb{R}^p, \quad y_i \in \mathbb{R}$$

- **Prediction/regression problem:** Observe $\{(x_i, y_i)\}_{i=1}^n$, estimate

$$\beta^* = \arg \min_{\beta \in \mathbb{R}^p} \mathbb{E}[\ell(\beta; x_i, y_i)], \quad x_i \in \mathbb{R}^p, \quad y_i \in \mathbb{R}$$

- Statistical M -estimator:

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(\beta; x_i, y_i) \right\}$$

in high dimensions, may be ill-conditioned, large solution space

- **Prediction/regression problem:** Observe $\{(x_i, y_i)\}_{i=1}^n$, estimate

$$\beta^* = \arg \min_{\beta \in \mathbb{R}^p} \mathbb{E}[\ell(\beta; x_i, y_i)], \quad x_i \in \mathbb{R}^p, \quad y_i \in \mathbb{R}$$

- Regularized M -estimator:

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(\beta; x_i, y_i)}_{\mathcal{L}_n(\beta)} + \rho_{\lambda}(\beta) \right\}$$

Example: ℓ_1 -regularized OLS regression

- Linear model: $y_i = x_i^T \beta^* + \epsilon_i, \quad \|\beta^*\|_0 \leq k$

Example: ℓ_1 -regularized OLS regression

- Linear model: $y_i = x_i^T \beta^* + \epsilon_i, \quad \|\beta^*\|_0 \leq k$
- Low-dimensional M -estimator:

$$\hat{\beta}_{\text{OLS}} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 \right\}$$

Example: ℓ_1 -regularized OLS regression

- Linear model: $y_i = x_i^T \beta^* + \epsilon_i, \quad \|\beta^*\|_0 \leq k$
- Low-dimensional M -estimator:

$$\hat{\beta}_{\text{OLS}} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 \right\} = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n y_i x_i \right)$$

Example: ℓ_1 -regularized OLS regression

- Linear model: $y_i = x_i^T \beta^* + \epsilon_i, \quad \|\beta^*\|_0 \leq k$
- Low-dimensional M -estimator:

$$\hat{\beta}_{\text{OLS}} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 \right\} = \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n y_i x_i \right)$$

- High-dimensional **regularized** M -estimator:

$$\hat{\beta}_{\text{Lasso}} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_1 \right\}$$

Sources of nonconvexity

- May arise in **loss** or **regularizer**

Sources of nonconvexity

- May arise in **loss** or **regularizer**
- Nonconvex **loss** used to correct noise/errors, increase efficiency

Sources of nonconvexity

- May arise in **loss** or **regularizer**
- Nonconvex **loss** used to correct noise/errors, increase efficiency
- Nonconvex **regularizer** used to reduce bias

Example: Errors-in-variables regression

- Model:

$$y_i = x_i^T \beta^* + \epsilon_i$$

observe $\{(z_i, y_i)\}_{i=1}^n$, infer β^*

Example: Errors-in-variables regression

- Model:

$$y_i = x_i^T \beta^* + \epsilon_i$$

observe $\{(z_i, y_i)\}_{i=1}^n$, infer β^*

- OLS estimator

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (z_i^T \beta - y_i)^2 + \lambda \|\beta\|_1 \right\}$$

statistically inconsistent

- L. & Wainwright '12 propose natural method for correcting loss for linear regression:

$$\hat{\beta}_{\text{OLS}} \in \arg \min_{\beta} \left\{ \frac{1}{2} \beta^T \frac{X^T X}{n} \beta - \frac{y X^T}{n} \beta + \rho_{\lambda}(\beta) \right\}$$

$$\hat{\beta}_{\text{corr}} \in \arg \min_{\beta} \left\{ \frac{1}{2} \beta^T \hat{\Gamma} \beta - \hat{\gamma}^T \beta + \rho_{\lambda}(\beta) \right\}$$

$(\hat{\Gamma}, \hat{\gamma})$ estimators for $(\text{Cov}(x_i), \text{Cov}(x_i, y_i))$ based on $\{(z_i, y_i)\}_{i=1}^n$

Example: Additive noise

- Additive noise: $Z = X + W$, use

$$\hat{\Gamma} = \frac{Z^T Z}{n} - \Sigma_w, \quad \hat{\gamma} = \frac{Z^T y}{n}$$

- However, corrected objective **nonconvex**:

$$\hat{\beta}_{\text{corr}} \in \arg \min_{\beta} \left\{ \frac{1}{2} \beta^T \left(\frac{Z^T Z}{n} - \Sigma_w \right) \beta - \frac{y^T Z}{n} \beta + \rho_{\lambda}(\beta) \right\}$$

Example: Additive noise

- Additive noise: $Z = X + W$, use

$$\hat{\Gamma} = \frac{Z^T Z}{n} - \Sigma_w, \quad \hat{\gamma} = \frac{Z^T y}{n}$$

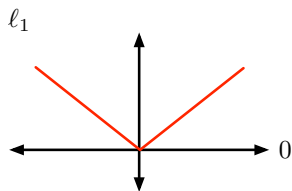
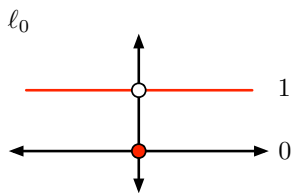
- However, corrected objective **nonconvex**:

$$\hat{\beta}_{\text{corr}} \in \arg \min_{\beta} \left\{ \frac{1}{2} \beta^T \left(\frac{Z^T Z}{n} - \Sigma_w \right) \beta - \frac{y^T Z}{n} \beta + \rho_{\lambda}(\beta) \right\}$$

- Fortunately, local optima have good properties

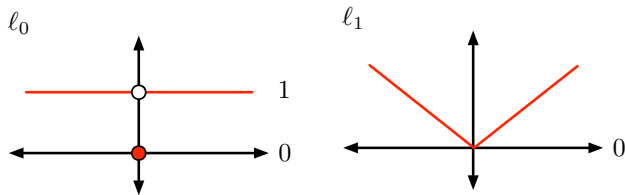
Nonconvex regularizers

- l_1 is “convexified” version of l_0



Nonconvex regularizers

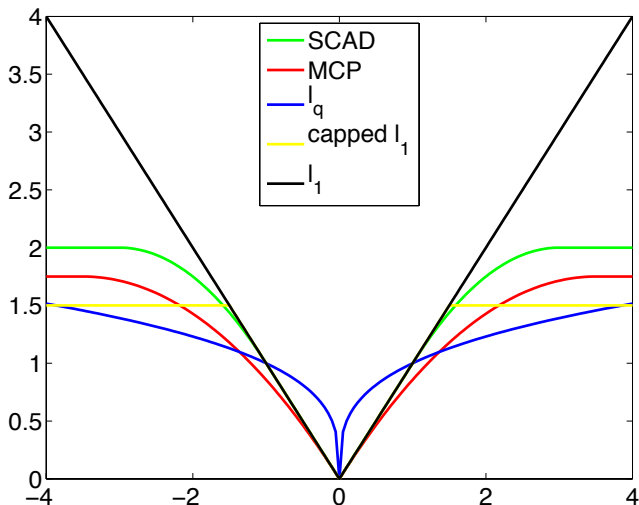
- l_1 is “convexified” version of l_0



- But** l_1 penalizes larger coefficients more, causes *solution bias*

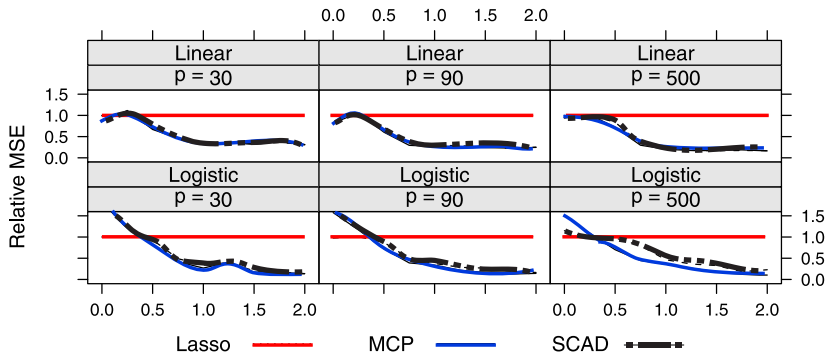
Alternative regularizers

- Various nonconvex regularizers in literature (Fan & Li '01, Zhang '10, etc.)



Empirical benefits

- Nonconvex regularizers show **significant improvement** (Breheny & Huang '11)

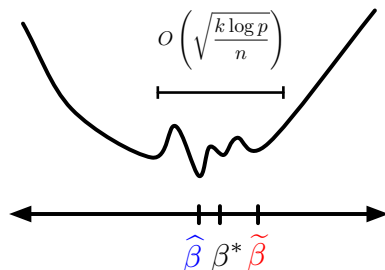


Local vs. global optima

- Optimization algorithms only guaranteed to find *local optima* (stationary points)
- Statistical theory only guarantees consistency of *global optima*

Local vs. global optima

- Optimization algorithms only guaranteed to find *local optima* (stationary points)
- Statistical theory only guarantees consistency of *global optima*



- **L. & Wainwright '13:** All stationary points of $\mathcal{L}_n(\beta) + \rho_\lambda(\beta)$ close when **nonconvexity** smaller than **curvature**

- Various measures of statistical consistency

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(\beta; x_i, y_i) + \rho_{\lambda}(\beta) \right\}$$

- Various measures of statistical consistency

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(\beta; \mathbf{x}_i, y_i) + \rho_{\lambda}(\beta) \right\}$$

- **Estimation:** $\|\hat{\beta} - \beta^*\| \rightarrow 0$
- **Prediction:** $\frac{1}{n} \sum_{i=1}^n \ell(\hat{\beta}; \mathbf{x}_i, y_i) \rightarrow 0$
- **Variable selection:** $\text{supp}(\hat{\beta}) \rightarrow \text{supp}(\beta^*)$

- Various measures of statistical consistency

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(\beta; \mathbf{x}_i, y_i) + \rho_{\lambda}(\beta) \right\}$$

- **Estimation:** $\|\hat{\beta} - \beta^*\| \rightarrow 0$
- **Prediction:** $\frac{1}{n} \sum_{i=1}^n \ell(\hat{\beta}; \mathbf{x}_i, y_i) \rightarrow 0$
- **Variable selection:** $\text{supp}(\hat{\beta}) \rightarrow \text{supp}(\beta^*)$

- Interested in cases where ℓ and ρ_{λ} possibly *nonconvex*

- Composite objective function

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \left\{ \mathcal{L}_n(\beta) + \sum_{j=1}^p \rho_\lambda(\beta_j) \right\}$$

- Composite objective function

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \left\{ \mathcal{L}_n(\beta) + \sum_{j=1}^p \rho_\lambda(\beta_j) \right\}$$

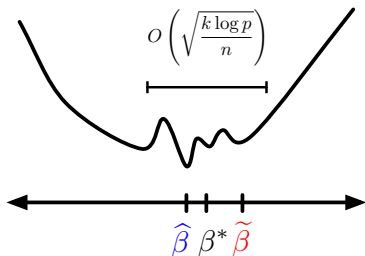
- \mathcal{L}_n satisfies **restricted strong convexity** with curvature α
- ρ_λ has bounded subgradient at 0, and $\rho_\lambda(t) + \mu t^2$ convex

- Composite objective function

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \left\{ \mathcal{L}_n(\beta) + \sum_{j=1}^p \rho_\lambda(\beta_j) \right\}$$

- \mathcal{L}_n satisfies **restricted strong convexity** with curvature α
- ρ_λ has bounded subgradient at 0, and $\rho_\lambda(t) + \mu t^2$ convex
- **L. & Wainwright '13**: All stationary points of $\mathcal{L}_n(\beta) + \rho_\lambda(\beta)$ close when $\alpha > \mu$

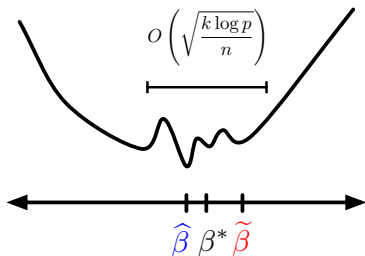
More formally



- **Stationary points** statistically indistinguishable from **global optima**

$$\langle \nabla \mathcal{L}_n(\tilde{\beta}) + \nabla \rho_\lambda(\tilde{\beta}), \beta - \tilde{\beta} \rangle \geq 0, \quad \forall \beta \text{ feasible}$$

More formally



- **Stationary points** statistically indistinguishable from **global optima**

$$\langle \nabla \mathcal{L}_n(\tilde{\beta}) + \nabla \rho_\lambda(\tilde{\beta}), \beta - \tilde{\beta} \rangle \geq 0, \quad \forall \beta \text{ feasible}$$

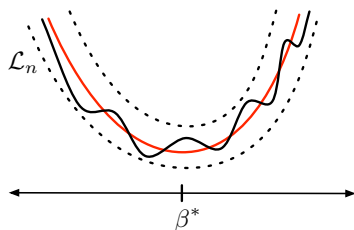
- **Nonasymptotic rates:** For $\lambda \asymp \sqrt{\frac{\log p}{n}}$ and $R \asymp \frac{1}{\lambda}$,

$$\|\tilde{\beta} - \beta^*\|_2 \leq c \sqrt{\frac{k \log p}{n}} \approx \text{statistical error}$$

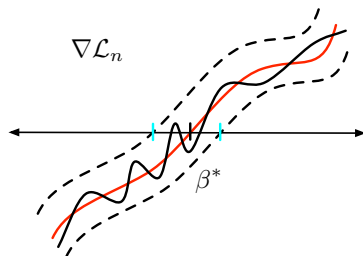
Geometric intuition

- **Population-level** convexity, **finite-sample** nonconvexity

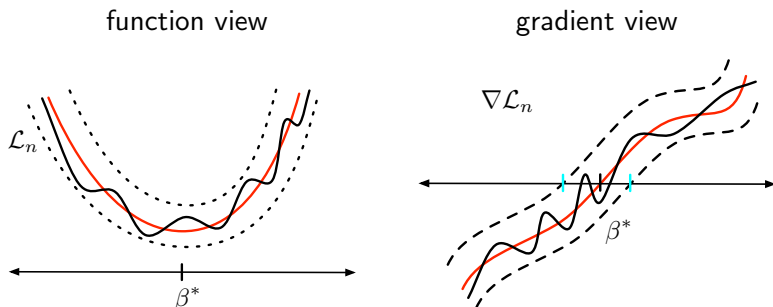
function view



gradient view



- Population-level convexity, **finite-sample** nonconvexity



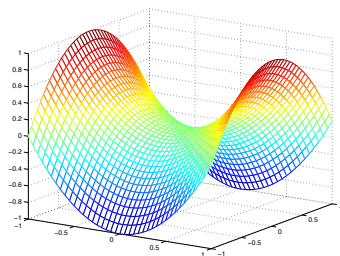
- Population-level objective \mathcal{L} strongly convex, $\alpha > \mu$
- RSC quantifies convergence rate of $\nabla \mathcal{L}_n \rightarrow \nabla \mathcal{L}$

- Requirements on loss and regularizer to ensure consistency of stationary points

- Requirements on loss and regularizer to ensure consistency of stationary points
 - Restricted strong convexity of \mathcal{L}_n
 - Bound on nonconvexity of ρ_λ

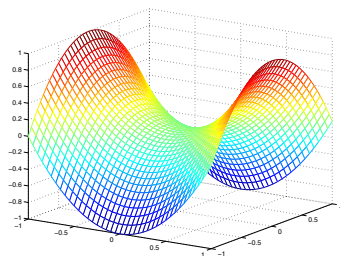
- **Restricted strong convexity** (Negahban et al. '12):

$$\langle \nabla \mathcal{L}_n(\beta^* + \Delta) - \nabla \mathcal{L}_n(\beta^*), \Delta \rangle \geq \begin{cases} \alpha \|\Delta\|_2^2 - \tau \frac{\log p}{n} \|\Delta\|_1^2, & \forall \|\Delta\|_2 \leq r \\ \alpha \|\Delta\|_2 - \tau \sqrt{\frac{\log p}{n}} \|\Delta\|_1, & \text{o.w.} \end{cases}$$



- **Restricted strong convexity** (Negahban et al. '12):

$$\langle \nabla \mathcal{L}_n(\beta^* + \Delta) - \nabla \mathcal{L}_n(\beta^*), \Delta \rangle \geq \begin{cases} \alpha \|\Delta\|_2^2 - \tau \frac{\log p}{n} \|\Delta\|_1^2, & \forall \|\Delta\|_2 \leq r \\ \alpha \|\Delta\|_2 - \tau \sqrt{\frac{\log p}{n}} \|\Delta\|_1, & \text{o.w.} \end{cases}$$



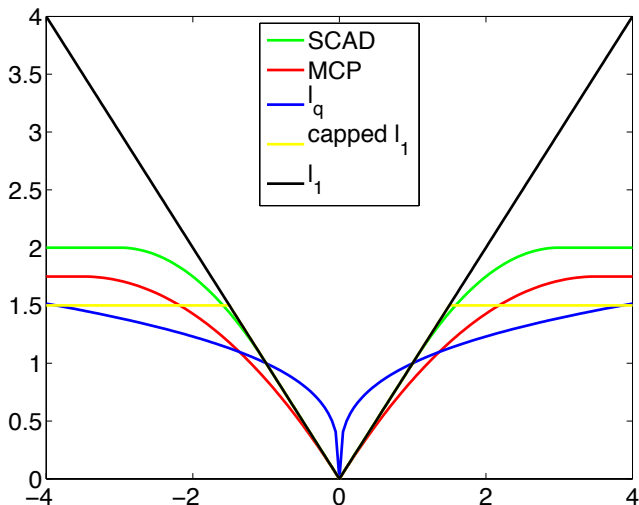
- Holds for various convex/nonconvex losses:
 - OLS & corrected OLS for linear regression, log likelihood for GLMs
 - Huber loss for robust regression

- Focus on *amenable* regularizers $\rho_\lambda(\beta) = \sum_{j=1}^p \rho_\lambda(\beta_j)$ satisfying:

- Focus on *amenable* regularizers $\rho_\lambda(\beta) = \sum_{j=1}^p \rho_\lambda(\beta_j)$ satisfying:
 - $\rho_\lambda(0) = 0$, symmetric around 0
 - Nondecreasing on \mathbb{R}^+
 - $t \mapsto \frac{\rho_\lambda(t)}{t}$ nonincreasing on \mathbb{R}^+
 - $q_\lambda(t) := \lambda|t| - \rho_\lambda(t)$ differentiable everywhere
 - $\rho_\lambda(t) + \mu t^2$ convex for some $\mu > 0$

Alternative regularizers

- Various nonconvex regularizers in literature (Fan & Li '01, Zhang '10, etc.)



Statistical consistency

- Regularized M -estimator

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \{ \mathcal{L}_n(\beta) + \rho_\lambda(\beta) \},$$

loss function satisfies (α, τ) -RSC and regularizer is μ -amenable

Statistical consistency

- Regularized M -estimator

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \{ \mathcal{L}_n(\beta) + \rho_\lambda(\beta) \},$$

loss function satisfies (α, τ) -RSC and regularizer is μ -amenable

Theorem (L. & Wainwright '13)

Suppose R is chosen s.t. β^* is feasible, and λ satisfies

$$\max \left\{ \|\nabla \mathcal{L}_n(\beta^*)\|_\infty, \alpha \sqrt{\frac{\log p}{n}} \right\} \lesssim \lambda \lesssim \frac{\alpha}{R}.$$

Statistical consistency

- Regularized M -estimator

$$\hat{\beta} \in \arg \min_{\|\beta\|_1 \leq R} \{ \mathcal{L}_n(\beta) + \rho_\lambda(\beta) \},$$

loss function satisfies (α, τ) -RSC and regularizer is μ -amenable

Theorem (L. & Wainwright '13)

Suppose R is chosen s.t. β^* is feasible, and λ satisfies

$$\max \left\{ \|\nabla \mathcal{L}_n(\beta^*)\|_\infty, \alpha \sqrt{\frac{\log p}{n}} \right\} \lesssim \lambda \lesssim \frac{\alpha}{R}.$$

For $n \geq \frac{C\tau^2}{\alpha^2} R^2 \log p$, any stationary point $\tilde{\beta}$ satisfies

$$\|\tilde{\beta} - \beta^*\|_2 \lesssim \frac{\lambda \sqrt{k}}{\alpha - \mu}, \quad \text{where } k = \|\beta^*\|_0.$$

- 1 Convexity of population-level objective \implies tractable landscape of empirical loss

Robust regression

- Robust statistics introduced in 1960s (Huber, Tukey, Hampel, et al.)

Robust regression

- Robust statistics introduced in 1960s (Huber, Tukey, Hampel, et al.)
- **Goals:**
 - 1 Develop estimators $T(\cdot)$ that are reliable under deviations from model assumptions
 - 2 Quantify performance with respect to deviations

- Robust statistics introduced in 1960s (Huber, Tukey, Hampel, et al.)
- **Goals:**
 - ① Develop estimators $T(\cdot)$ that are reliable under deviations from model assumptions
 - ② Quantify performance with respect to deviations
- Local stability captured by *influence function*

$$IF(x; T, F) = \lim_{t \rightarrow 0} \frac{T((1-t)F + t\delta_x) - T(F)}{t}$$

- Robust statistics introduced in 1960s (Huber, Tukey, Hampel, et al.)
- **Goals:**
 - ① Develop estimators $T(\cdot)$ that are reliable under deviations from model assumptions
 - ② Quantify performance with respect to deviations
- Local stability captured by *influence function*

$$IF(x; T, F) = \lim_{t \rightarrow 0} \frac{T((1-t)F + t\delta_x) - T(F)}{t}$$

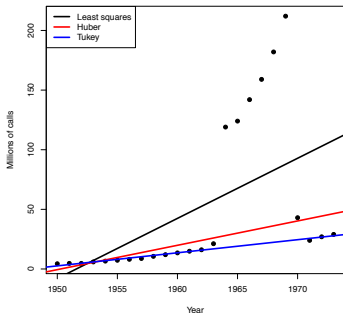
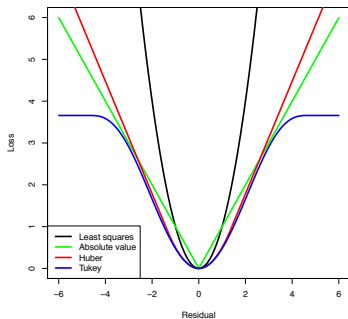
- Global stability captured by *breakdown point*

$$\epsilon^*(T; X_1, \dots, X_n) = \min \left\{ \frac{m}{n} : \sup_{X^m} \|T(X^m) - T(X)\| = \infty \right\}$$

“Robust” M -estimators

- Generalization of OLS suitable for heavy-tailed/contaminated errors:

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(x_i^T \beta - y_i) \right\}$$

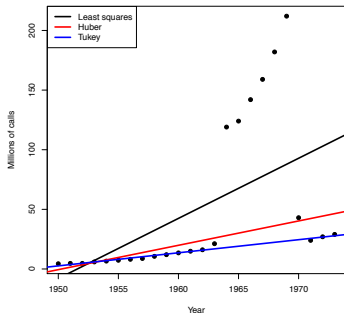
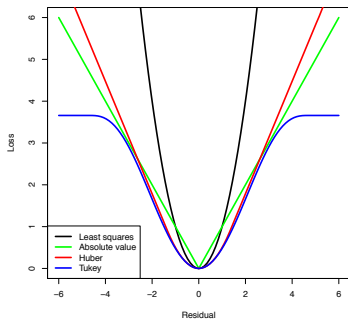


“Robust” M -estimators

- Generalization of OLS suitable for heavy-tailed/contaminated errors:

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(x_i^T \beta - y_i) \right\}$$

- Extensive theory (consistency, asymptotic normality) for p fixed, $n \rightarrow \infty$



Classes of loss functions

- **Bounded** ℓ' limits influence of outliers:

$$IF((x, y); T, F) = \lim_{t \rightarrow 0^+} \frac{T((1-t)F + t\delta_{(x,y)}) - T(F)}{t}$$
$$\propto \ell'(x^T \beta - y)x$$

where $F \sim F_\beta$ and T minimizes M -estimator

Classes of loss functions

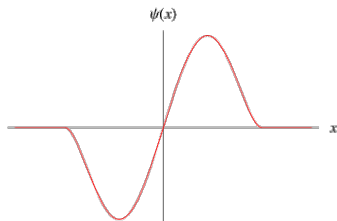
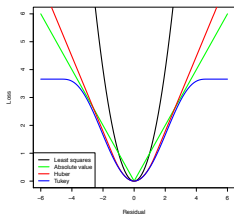
- **Bounded** ℓ' limits influence of outliers:

$$IF((x, y); T, F) = \lim_{t \rightarrow 0^+} \frac{T((1-t)F + t\delta_{(x,y)}) - T(F)}{t} \\ \propto \ell'(x^T \beta - y)x$$

where $F \sim F_\beta$ and T minimizes M -estimator

- **Redescending** M -estimators have *finite rejection point*:

$$\ell'(u) = 0, \quad \text{for } |u| \geq c$$



Classes of loss functions

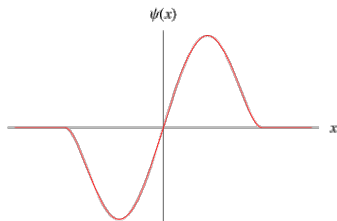
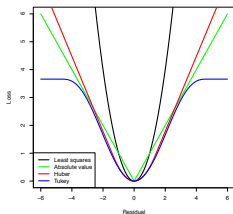
- **Bounded** ℓ' limits influence of outliers:

$$IF((x, y); T, F) = \lim_{t \rightarrow 0^+} \frac{T((1-t)F + t\delta_{(x,y)}) - T(F)}{t} \\ \propto \ell'(x^T \beta - y)x$$

where $F \sim F_\beta$ and T minimizes M -estimator

- **Redescending** M -estimators have *finite rejection point*:

$$\ell'(u) = 0, \quad \text{for } |u| \geq c$$



- **But bad for optimization!!**

- **Natural idea:** For $p > n$, use regularized version:

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(x_i^T \beta - y_i) + \lambda \|\beta\|_1 \right\}$$

- **Natural idea:** For $p > n$, use regularized version:

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(x_i^T \beta - y_i) + \lambda \|\beta\|_1 \right\}$$

Complications:

- Optimization for nonconvex ℓ ?
- Statistical theory? Are certain losses provably better than others?

- **Natural idea:** For $p > n$, use regularized version:

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(x_i^T \beta - y_i) + \lambda \|\beta\|_1 \right\}$$

Complications:

- Optimization for nonconvex ℓ ?
- Statistical theory? Are certain losses provably better than others?
- Population-level convexity *no longer satisfied*

- When $\|\ell'\|_\infty < C$, **global optima** of high-dimensional M -estimator satisfy

$$\|\widehat{\beta} - \beta^*\|_2 \leq C \sqrt{\frac{k \log p}{n}},$$

regardless of distribution of ϵ_i

- When $\|\ell'\|_\infty < C$, **global optima** of high-dimensional M -estimator satisfy

$$\|\hat{\beta} - \beta^*\|_2 \leq C \sqrt{\frac{k \log p}{n}},$$

regardless of distribution of ϵ_i

- **Compare to Lasso theory:** Requires sub-Gaussian ϵ_i 's

- When $\|\ell'\|_\infty < C$, **global optima** of high-dimensional M -estimator satisfy

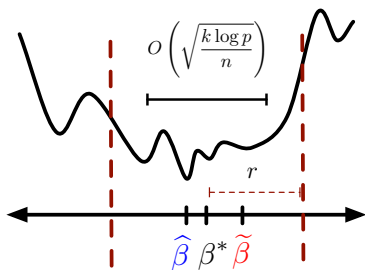
$$\|\hat{\beta} - \beta^*\|_2 \leq C \sqrt{\frac{k \log p}{n}},$$

regardless of distribution of ϵ_i

- **Compare to Lasso theory:** Requires sub-Gaussian ϵ_i 's
- If $\ell(u)$ is *locally convex/smooth* for $|u| \leq r$, any **local optima** within radius cr of β^* satisfy

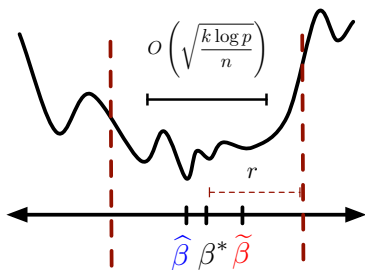
$$\|\tilde{\beta} - \beta^*\|_2 \leq C' \sqrt{\frac{k \log p}{n}}$$

Some optimization theory (L. '17)



- Local optima may be obtained via **two-step algorithm** (L. '17)

Some optimization theory (L. '17)



- Local optima may be obtained via **two-step algorithm** (L. '17)

Algorithm

- Run composite gradient descent on **convex**, robust loss + ℓ_1 -penalty until convergence, output $\hat{\beta}_H$
- Run composite gradient descent on **nonconvex**, robust loss + μ -amenable penalty, input $\beta^0 = \hat{\beta}_H$

Motivating calculation

- Lasso analysis (e.g., van de Geer '07, Bickel et al. '08):

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \underbrace{\frac{1}{n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1}_{F(\beta)} \right\}$$

Motivating calculation

- Lasso analysis (e.g., van de Geer '07, Bickel et al. '08):

$$\hat{\beta} \in \arg \min_{\beta} \left\{ \underbrace{\frac{1}{n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1}_{F(\beta)} \right\}$$

- Rearranging *basic inequality* $F(\hat{\beta}) \leq F(\beta^*)$ and assuming $\lambda \geq 2 \left\| \frac{X^T \epsilon}{n} \right\|_{\infty}$, obtain

$$\|\hat{\beta} - \beta^*\|_2 \leq c\lambda\sqrt{k}$$

Motivating calculation

- Lasso analysis (e.g., van de Geer '07, Bickel et al. '08):

$$\hat{\beta} \in \arg \min_{\beta} \underbrace{\left\{ \frac{1}{n} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}}_{F(\beta)}$$

- Rearranging *basic inequality* $F(\hat{\beta}) \leq F(\beta^*)$ and assuming $\lambda \geq 2 \left\| \frac{X^T \epsilon}{n} \right\|_{\infty}$, obtain

$$\|\hat{\beta} - \beta^*\|_2 \leq c\lambda\sqrt{k}$$

- Sub-Gaussian assumptions on x_i 's and ϵ_i 's provide $\mathcal{O}\left(\sqrt{\frac{k \log p}{n}}\right)$ bounds, minimax optimal

Motivating calculation

- **Key observation:** For general loss function, if $\lambda \geq 2 \left\| \frac{X^T \ell'(\epsilon)}{n} \right\|_\infty$, obtain

$$\|\hat{\beta} - \beta^*\|_2 \leq c\lambda\sqrt{k}$$

Motivating calculation

- **Key observation:** For general loss function, if $\lambda \geq 2 \left\| \frac{X^T \ell'(\epsilon)}{n} \right\|_\infty$, obtain

$$\|\hat{\beta} - \beta^*\|_2 \leq c\lambda\sqrt{k}$$

- $\ell'(\epsilon)$ sub-Gaussian whenever ℓ' bounded

Motivating calculation

- **Key observation:** For general loss function, if $\lambda \geq 2 \left\| \frac{X^T \ell'(\epsilon)}{n} \right\|_{\infty}$, obtain

$$\|\hat{\beta} - \beta^*\|_2 \leq c\lambda\sqrt{k}$$

- $\ell'(\epsilon)$ sub-Gaussian whenever ℓ' bounded
 \implies can achieve estimation error

$$\|\hat{\beta} - \beta^*\|_2 \leq c\sqrt{\frac{k \log p}{n}},$$

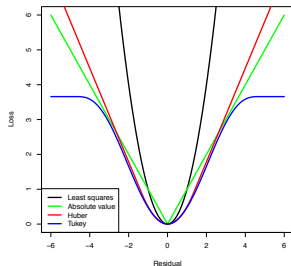
without assuming ϵ_j is sub-Gaussian

Technical challenges

- Lasso analysis also requires verifying restricted eigenvalue (RE) condition on design matrix, more complicated for general ℓ

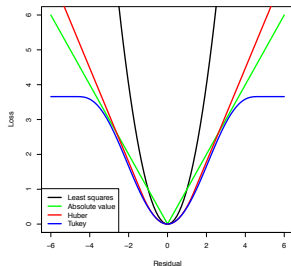
Technical challenges

- Lasso analysis also requires verifying restricted eigenvalue (RE) condition on design matrix, more complicated for general ℓ
- Addressed by local curvature of robust losses around origin



Technical challenges

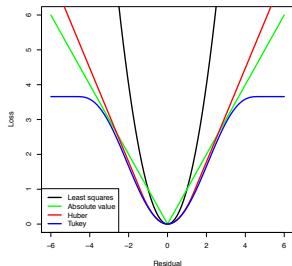
- Lasso analysis also requires verifying restricted eigenvalue (RE) condition on design matrix, more complicated for general ℓ
- Addressed by local curvature of robust losses around origin



- When ℓ is nonconvex, local optima $\tilde{\beta}$ may exist that are not global optima

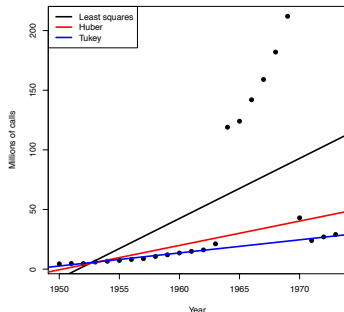
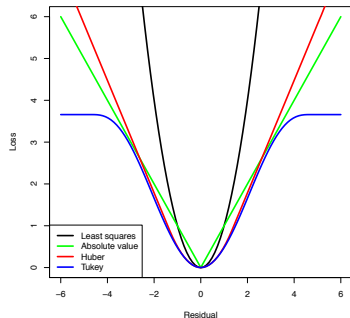
Technical challenges

- Lasso analysis also requires verifying restricted eigenvalue (RE) condition on design matrix, more complicated for general ℓ
- Addressed by local curvature of robust losses around origin



- When ℓ is nonconvex, local optima $\tilde{\beta}$ may exist that are not global optima
- Addressed by theoretical analysis of $\|\tilde{\beta} - \beta^*\|_2$ and derivation of suitable optimization algorithms

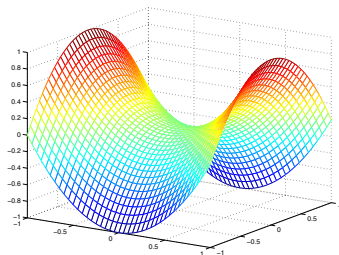
Local statistical consistency



- **Challenge in robust statistics:** Population-level nonconvexity of loss \implies need for *local* optimization theory

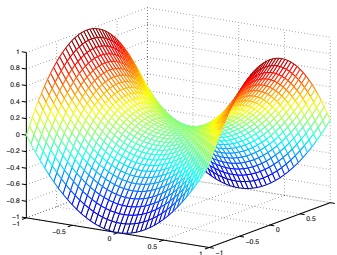
- **Local** restricted strong convexity: For $\Delta := \beta_1 - \beta_2$,

$$\langle \nabla \mathcal{L}_n(\beta_1) - \nabla \mathcal{L}_n(\beta_2), \Delta \rangle \geq \alpha \|\Delta\|_2^2 - \tau \frac{\log p}{n} \|\Delta\|_1^2, \quad \forall \|\beta_j - \beta^*\|_2 \leq r$$



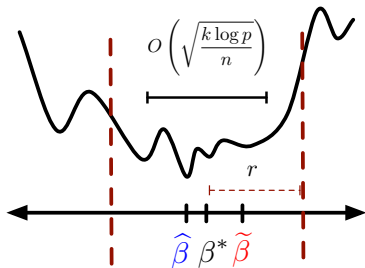
- **Local** restricted strong convexity: For $\Delta := \beta_1 - \beta_2$,

$$\langle \nabla \mathcal{L}_n(\beta_1) - \nabla \mathcal{L}_n(\beta_2), \Delta \rangle \geq \alpha \|\Delta\|_2^2 - \tau \frac{\log p}{n} \|\Delta\|_1^2, \quad \forall \|\beta_j - \beta^*\|_2 \leq r$$



- Only requires restricted curvature within constant-radius region around β^*

Consistency of local stationary points



Theorem (L. '17)

Suppose \mathcal{L}_n satisfies α -local RSC and ρ_λ is μ -amenable, with $\alpha > \mu$.

Suppose $\|\ell'\|_\infty \leq C$ and $\lambda \asymp \sqrt{\frac{\log p}{n}}$. For $n \gtrsim \frac{\tau}{\alpha - \mu} k \log p$, any stationary point $\tilde{\beta}$ s.t. $\|\tilde{\beta} - \beta^*\|_2 \leq r$ satisfies

$$\|\tilde{\beta} - \beta^*\|_2 \lesssim \frac{\lambda \sqrt{k}}{\alpha - \mu}.$$

- 1 Convexity of population-level objective \implies tractable landscape of empirical loss
- 2 Local convexity of population-level objective \implies empirical loss landscape **locally** well-behaved

- 1 Convexity of population-level objective \implies tractable landscape of empirical loss
- 2 Local convexity of population-level objective \implies empirical loss landscape **locally** well-behaved
- 3 Global optimum of **convex surrogate** may provide appropriate initial point

- **P. Loh** and M. J. Wainwright (2015). Regularized M -estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Journal of Machine Learning Research*.
- **P. Loh** (2018). Statistical consistency and asymptotic normality for high-dimensional robust M -estimators. *Annals of Statistics*.

Thank you!