

Estimators based on non-convex programs: Statistical and computational guarantees

Martin Wainwright

UC Berkeley
Statistics and EECS

Based on joint work with:

Po-Ling Loh (UC Berkeley)

Introduction

Prediction/regression problems arise throughout statistics:

- vector of predictors/covariates $x \in \mathbb{R}^p$
- response variable $y \in \mathbb{R}$
- In “big data” setting:
 - ▶ both sample size n and ambient dimension p are large
 - ▶ many problems have $p \gg n$

Introduction

Prediction/regression problems arise throughout statistics:

- vector of predictors/covariates $x \in \mathbb{R}^p$
- response variable $y \in \mathbb{R}$
- In “big data” setting:
 - ▶ both sample size n and ambient dimension p are large
 - ▶ many problems have $p \gg n$

Regularization is essential:

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \underbrace{\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\theta; x_i, y_i)}_{\mathcal{L}_n(\theta; x_1^n, y_1^n)} + \underbrace{\mathcal{R}_\lambda(\theta)}_{\text{Regularizer}} \right\}.$$

Introduction

Prediction/regression problems arise throughout statistics:

- In “big data” setting:
 - ▶ both sample size n and ambient dimension p are large
 - ▶ many problems have $p \gg n$

Regularization is essential:

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \underbrace{\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\theta; x_i, y_i)}_{\mathcal{L}_n(\theta; x_1^n, y_1^n)} + \underbrace{\mathcal{R}_\lambda(\theta)}_{\text{Regularizer}} \right\}.$$

For non-convex problems: “Mind the gap!”

- any global optimum is “statistically good”
- but efficient algorithms only find local optima

Introduction

Prediction/regression problems arise throughout statistics:

- In “big data” setting:
 - ▶ both sample size n and ambient dimension p are large
 - ▶ many problems have $p \gg n$

Regularization is essential:

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \underbrace{\frac{1}{n} \sum_{i=1}^n \mathcal{L}(\theta; x_i, y_i)}_{\mathcal{L}_n(\theta; x_1^n, y_1^n)} + \underbrace{\mathcal{R}_\lambda(\theta)}_{\text{Regularizer}} \right\}.$$

For non-convex problems: “Mind the gap!”

- any global optimum is “statistically good”
- but efficient algorithms only find local optima

Question

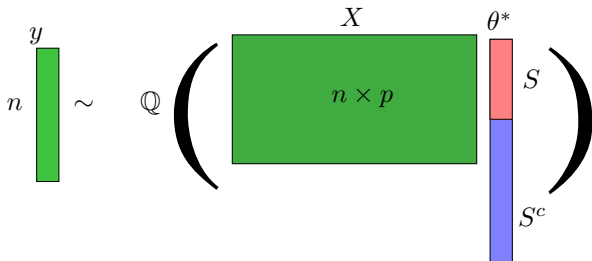
How to close this undesirable gap between statistics and computation?

Vignette A: Regression with non-convex penalties

Set-up: **Observe** (y_i, x_i) pairs for $i = 1, 2, \dots, n$, where

$$y_i \sim \mathbb{Q}(\cdot \mid \langle \theta^*, x_i \rangle),$$

where $\theta \in \mathbb{R}^p$ has “low-dimensional structure”

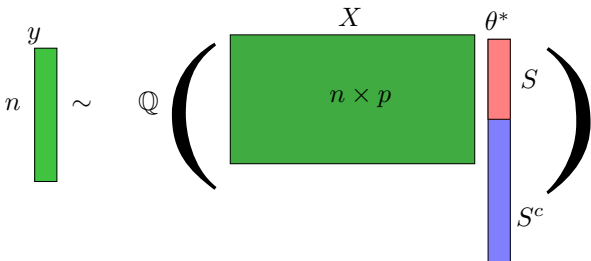


Vignette A: Regression with non-convex penalties

Set-up: **Observe** (y_i, x_i) pairs for $i = 1, 2, \dots, n$, where

$$y_i \sim \mathbb{Q}(\cdot \mid \langle \theta^*, x_i \rangle),$$

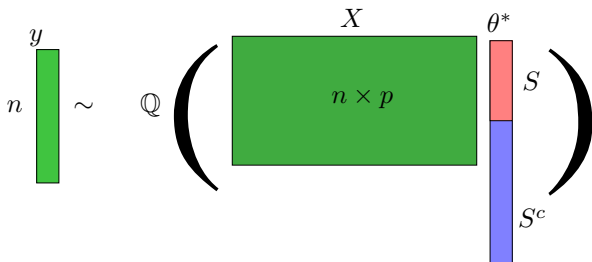
where $\theta \in \mathbb{R}^p$ has “low-dimensional structure”



Estimator: \mathcal{R}_λ -regularized likelihood

$$\hat{\theta} \in \arg \min_{\theta} \left\{ -\frac{1}{n} \sum_{i=1}^n \log \mathbb{Q}(y_i \mid \langle x_i, \theta \rangle) + \mathcal{R}_\lambda(\theta) \right\}.$$

Vignette A: Regression with non-convex penalties



Example: Logistic regression for binary responses $y_i \in \{0, 1\}$:

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \frac{1}{n} \sum_{i=1}^n \{ \log(1 + e^{\langle x_i, \theta \rangle}) - y_i \langle x_i, \theta \rangle \} + \mathcal{R}_{\lambda}(\theta) \right\}.$$

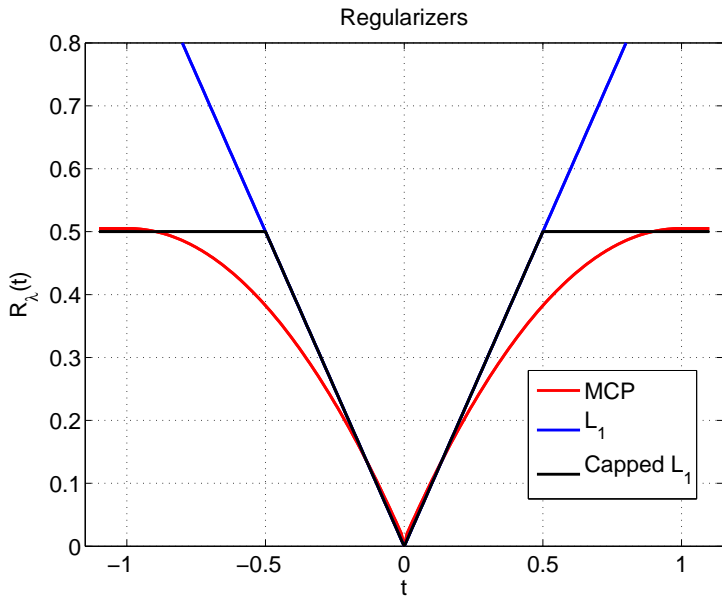
Many **non-convex penalties** are possible:

- capped ℓ_1 -penalty
- SCAD penalty
- MCP penalty

(Fan & Li, 2001)

(Zhang, 2006)

Convex and non-convex regularizers



Statistical error versus optimization error

Algorithm generating sequence of iterates $\{\theta^t\}_{t=0}^{\infty}$ to solve

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\theta; x_i, y_i) + \mathcal{R}_{\lambda}(\theta) \right\}.$$

Global minimizer of population risk

$$\theta^* := \arg \min_{\theta} \underbrace{\mathbb{E}_{X,Y} [\mathcal{L}(\theta; X, Y)]}_{\bar{\mathcal{L}}(\theta)}$$

Statistical error versus optimization error

Algorithm generating sequence of iterates $\{\theta^t\}_{t=0}^{\infty}$ to solve

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\theta; x_i, y_i) + \mathcal{R}_{\lambda}(\theta) \right\}.$$

Global minimizer of population risk

$$\theta^* := \arg \min_{\theta} \underbrace{\mathbb{E}_{X,Y} [\mathcal{L}(\theta; X, Y)]}_{\bar{\mathcal{L}}(\theta)}$$

Goal of statistician

Provide bounds on **Statistical error:** $\|\theta^t - \theta^*\|$ or $\|\hat{\theta} - \theta^*\|$

Statistical error versus optimization error

Algorithm generating sequence of iterates $\{\theta^t\}_{t=0}^{\infty}$ to solve

$$\hat{\theta} \in \arg \min_{\theta} \left\{ \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\theta; x_i, y_i) + \mathcal{R}_{\lambda}(\theta) \right\}.$$

Global minimizer of population risk

$$\theta^* := \arg \min_{\theta} \underbrace{\mathbb{E}_{X,Y} [\mathcal{L}(\theta; X, Y)]}_{\bar{\mathcal{L}}(\theta)}$$

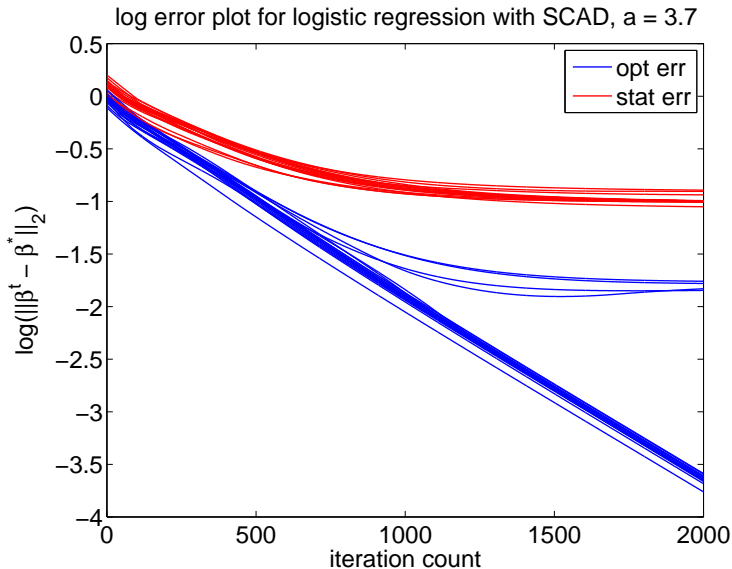
Goal of statistician

Provide bounds on **Statistical error:** $\|\theta^t - \theta^*\|$ or $\|\hat{\theta} - \theta^*\|$

Goal of optimization-theorist

Provide bounds on **Optimization error:** $\|\theta^t - \hat{\theta}\|$

Logistic regression with non-convex regularizer



What phenomena need to be explained?

Empirical observation #1:

From a statistical perspective, **all local optima** are essentially as good as a global optimum.

What phenomena need to be explained?

Empirical observation #1:

From a statistical perspective, **all local optima** are essentially as good as a global optimum.

Some past work:

- for least-squares loss, certain local optima are good (Zhang & Zhang, 2012)
- if initialized at Lasso solution with ℓ_∞ -guarantees, local algorithm has good behavior (Fan et al., 2012)

What phenomena need to be explained?

Empirical observation #1:

From a statistical perspective, **all local optima** are essentially as good as a global optimum.

Some past work:

- for least-squares loss, certain local optima are good (Zhang & Zhang, 2012)
- if initialized at Lasso solution with ℓ_∞ -guarantees, local algorithm has good behavior (Fan et al., 2012)

Empirical observation #2:

First-order methods converge **as fast as possible** up to **statistical precision**.

Vignette B: Error-in-variables regression

Begin with high-dimensional sparse regression:

$$y = X\theta^* + \varepsilon$$

The diagram illustrates the equation $y = X\theta^* + \varepsilon$ with the following components:

- y : A green vertical vector of size n .
- X : A gray rectangular matrix of size $n \times p$.
- θ^* : A vertical vector composed of two parts: a red top part labeled S and a blue bottom part labeled S^c .
- ε : A purple vertical vector representing the error term.

Vignette B: Error-in-variables regression

Begin with high-dimensional sparse regression:

$$y = X\theta^* + \varepsilon$$

The diagram illustrates the equation $y = X\theta^* + \varepsilon$. On the left, a green vertical bar labeled y has a dimension of n . This is equal to a gray rectangle labeled X with dimensions $n \times p$. This is followed by a vertical bar for θ^* , which is split into a red top section labeled S and a blue bottom section labeled S^c . Finally, a purple vertical bar labeled ε is added to the right.

Observe $y \in \mathbb{R}^n$ and $Z \in \mathbb{R}^{n \times p}$

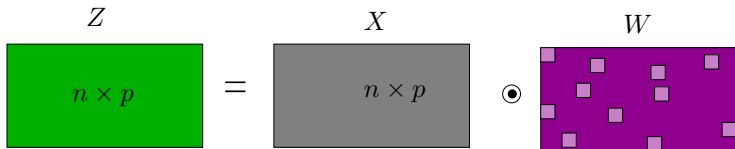
$$Z = \mathbf{F}\left(X, W\right)$$

The diagram illustrates the equation $Z = \mathbf{F}(X, W)$. On the left, a green rectangle labeled Z has dimensions $n \times p$. This is equal to the function \mathbf{F} applied to two inputs: a gray rectangle labeled X with dimensions $n \times p$, and a purple rectangle labeled W with dimensions $n \times p$.

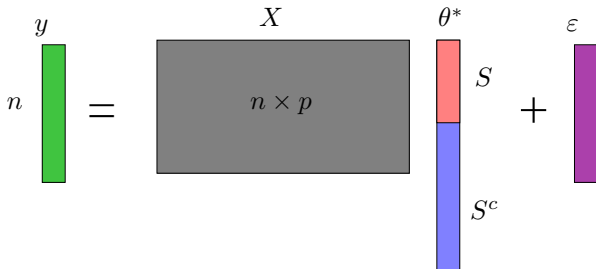
Here $W \in \mathbb{R}^{n \times p}$ is a stochastic perturbation.

Example: Missing data

Missing data:



Here $W \in \mathbb{R}^{n \times p}$ is multiplicative perturbation (e.g., $W_{ij} \sim \text{Ber}(\alpha)$.)



Example: Additive perturbations

Additive noise in covariates:

$$\begin{array}{c} Z \\ n \times p \end{array} = \begin{array}{c} X \\ n \times p \end{array} + \begin{array}{c} W \\ n \times p \end{array}$$

Here $W \in \mathbb{R}^{n \times p}$ is an additive perturbation (e.g., $W_{ij} \sim N(0, \sigma^2)$).

$$\begin{array}{c} y \\ n \end{array} = \begin{array}{c} X \\ n \times p \end{array} \begin{array}{c} \theta^* \\ S \\ S^c \end{array} + \begin{array}{c} \varepsilon \end{array}$$

A second look at regularized least-squares

Equivalent formulation:

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \left\{ \frac{1}{2} \theta^T \left(\frac{X^T X}{n} \right) \theta - \left\langle \theta, \frac{X^T y}{n} \right\rangle + \mathcal{R}_\lambda(\theta) \right\}.$$

A second look at regularized least-squares

Equivalent formulation:

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \left\{ \frac{1}{2} \theta^T \left(\frac{X^T X}{n} \right) \theta - \left\langle \theta, \frac{X^T y}{n} \right\rangle + \mathcal{R}_\lambda(\theta) \right\}.$$

Population view: unbiased estimators

$$\text{cov}(x_1) = \mathbb{E} \left[\frac{X^T X}{n} \right], \quad \text{and} \quad \text{cov}(x_1, y_1) = \mathbb{E} \left[\frac{X^T y}{n} \right].$$

Corrected estimators

Equivalent formulation:

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \left\{ \frac{1}{2} \theta^T \left(\frac{X^T X}{n} \right) \theta - \left\langle \theta, \frac{X^T y}{n} \right\rangle + \mathcal{R}_\lambda(\theta) \right\}.$$

Population view: unbiased estimators

$$\text{cov}(x_1) = \mathbb{E} \left[\frac{X^T X}{n} \right], \quad \text{and} \quad \text{cov}(x_1, y_1) = \mathbb{E} \left[\frac{X^T y}{n} \right].$$

A general family of estimators

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \left\{ \frac{1}{2} \theta^T \hat{\Gamma} \theta - \theta^T \hat{\gamma} + \lambda_n^2 \|\theta\|_1^2 \right\},$$

where $(\hat{\Gamma}, \hat{\gamma})$ are unbiased estimators of $\text{cov}(x_1)$ and $\text{cov}(x_1, y_1)$.

Example: Estimator for missing data

- observe corrupted version $Z \in \mathbb{R}^{n \times p}$

$$Z_{ij} = \begin{cases} X_{ij} & \text{with probability } 1 - \alpha \\ \star & \text{with probability } \alpha. \end{cases}$$

Example: Estimator for missing data

- observe corrupted version $Z \in \mathbb{R}^{n \times p}$

$$Z_{ij} = \begin{cases} X_{ij} & \text{with probability } 1 - \alpha \\ \star & \text{with probability } \alpha. \end{cases}$$

- Natural unbiased estimates: set $\star \equiv 0$ and $\hat{Z} := \frac{Z}{(1-\alpha)}$:

$$\hat{\Gamma} = \frac{\hat{Z}^T \hat{Z}}{n} - \alpha \text{diag} \left(\frac{\hat{Z}^T \hat{Z}}{n} \right), \quad \text{and} \quad \hat{\gamma} = \frac{\hat{Z}^T y}{n},$$

Example: Estimator for missing data

- observe corrupted version $Z \in \mathbb{R}^{n \times p}$

$$Z_{ij} = \begin{cases} X_{ij} & \text{with probability } 1 - \alpha \\ \star & \text{with probability } \alpha. \end{cases}$$

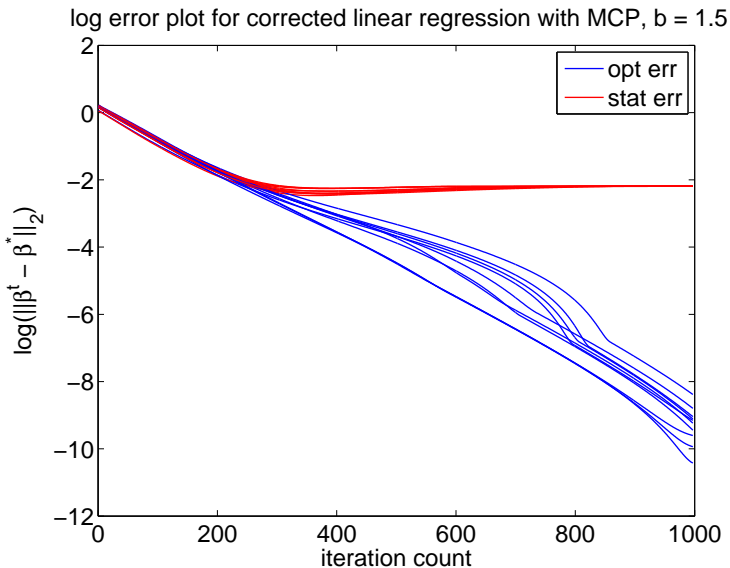
- Natural unbiased estimates: set $\star \equiv 0$ and $\hat{Z} := \frac{Z}{(1-\alpha)}$:

$$\hat{\Gamma} = \frac{\hat{Z}^T \hat{Z}}{n} - \alpha \text{diag} \left(\frac{\hat{Z}^T \hat{Z}}{n} \right), \quad \text{and} \quad \hat{\gamma} = \frac{\hat{Z}^T y}{n},$$

- solve (doubly non-convex) optimization problem: (Loh & W., 2012)

$$\hat{\theta} \in \arg \min_{\theta \in \Omega} \left\{ \frac{1}{2} \theta^T \hat{\Gamma} \theta - \langle \hat{\gamma}, \theta \rangle + \mathcal{R}_\lambda(\theta) \right\}.$$

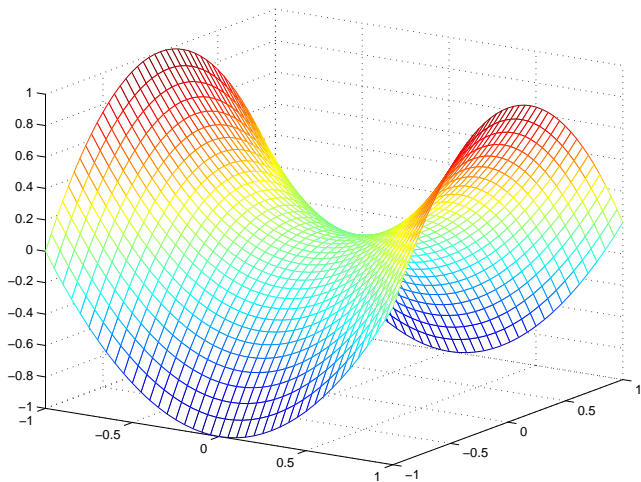
Non-convex quadratic and non-convex regularizer



Remainder of talk

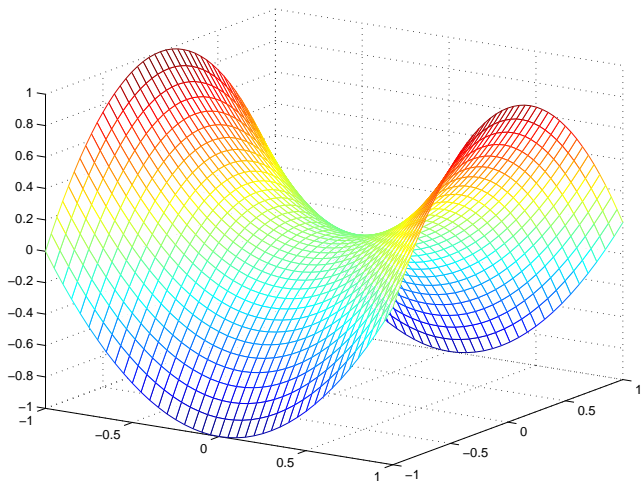
- ① Why are all local optima “statistically good”?
 - ▶ Restricted strong convexity
 - ▶ A general theorem
 - ▶ Various examples
- ② Why do first-order gradient methods converge quickly?
 - ▶ Composite gradient methods
 - ▶ Statistical versus optimization error
 - ▶ Fast convergence for non-convex problems

Geometry of a non-convex quadratic loss



- Loss function has directions of both positive and negative curvature.

Geometry of a non-convex quadratic loss



- Loss function has directions of both positive and negative curvature.
- Negative directions **must be forbidden** by regularizer.

Restricted strong convexity

Here defined with respect to the ℓ_1 -norm:

Definition

The loss function \mathcal{L}_n satisfies RSC with parameters (α_j, τ_j) , $j = 1, 2$ if

$$\underbrace{\langle \nabla \mathcal{L}_n(\theta^* + \Delta) - \nabla \mathcal{L}_n(\theta^*), \Delta \rangle}_{\text{Measure of curvature}} \geq \begin{cases} \alpha_1 \|\Delta\|_2^2 - \tau_1 \frac{\log p}{n} \|\Delta\|_1^2 & \text{if } \|\Delta\|_2 \leq 1 \\ \alpha_2 \|\Delta\|_2 - \tau_2 \sqrt{\frac{\log p}{n}} \|\Delta\|_1 & \text{if } \|\Delta\|_2 > 1. \end{cases}$$

Restricted strong convexity

Here defined with respect to the ℓ_1 -norm:

Definition

The loss function \mathcal{L}_n satisfies RSC with parameters (α_j, τ_j) , $j = 1, 2$ if

$$\underbrace{\langle \nabla \mathcal{L}_n(\theta^* + \Delta) - \nabla \mathcal{L}_n(\theta^*), \Delta \rangle}_{\text{Measure of curvature}} \geq \begin{cases} \alpha_1 \|\Delta\|_2^2 - \tau_1 \frac{\log p}{n} \|\Delta\|_1^2 & \text{if } \|\Delta\|_2 \leq 1 \\ \alpha_2 \|\Delta\|_2 - \tau_2 \sqrt{\frac{\log p}{n}} \|\Delta\|_1 & \text{if } \|\Delta\|_2 > 1. \end{cases}$$

- holds with $\tau_1 = \tau_2 = 0$ for any function that is locally strongly convex around θ^*

Restricted strong convexity

Here defined with respect to the ℓ_1 -norm:

Definition

The loss function \mathcal{L}_n satisfies RSC with parameters (α_j, τ_j) , $j = 1, 2$ if

$$\underbrace{\langle \nabla \mathcal{L}_n(\theta^* + \Delta) - \nabla \mathcal{L}_n(\theta^*), \Delta \rangle}_{\text{Measure of curvature}} \geq \begin{cases} \alpha_1 \|\Delta\|_2^2 - \tau_1 \frac{\log p}{n} \|\Delta\|_1^2 & \text{if } \|\Delta\|_2 \leq 1 \\ \alpha_2 \|\Delta\|_2 - \tau_2 \sqrt{\frac{\log p}{n}} \|\Delta\|_1 & \text{if } \|\Delta\|_2 > 1. \end{cases}$$

- holds with $\tau_1 = \tau_2 = 0$ for any function that is locally strongly convex around θ^*
- holds for a variety of loss functions (convex and non-convex):
 - ▶ ordinary least-squares (Raskutti, W. & Yu, 2010)
 - ▶ likelihoods for generalized linear models (Negahban et al., 2012)
 - ▶ certain non-convex quadratic functions (Loh & W, 2012)

Well-behaved regularizers

Properties defined at the univariate level $\mathcal{R}_\lambda : \mathbb{R} \rightarrow [0, \infty]$.

- Satisfies $\mathcal{R}_\lambda(0) = 0$, and is symmetric around zero ($\mathcal{R}_\lambda(t) = \mathcal{R}_\lambda(-t)$.)
- Non-decreasing and subadditive $\mathcal{R}_\lambda(s + t) \leq \mathcal{R}_\lambda(s) + \mathcal{R}_\lambda(t)$.
- Function $t \mapsto \frac{\mathcal{R}_\lambda(t)}{t}$ is nonincreasing for $t > 0$
- Differentiable for all $t \neq 0$, subdifferentiable at $t = 0$ with subgradients bounded in absolute value by λL .
- For some $\mu > 0$, the function $\tilde{\mathcal{R}}_\lambda(t) = \mathcal{R}_\lambda(t) + \mu t^2$ is convex.

Well-behaved regularizers

Properties defined at the univariate level $\mathcal{R}_\lambda : \mathbb{R} \rightarrow [0, \infty]$.

- Satisfies $\mathcal{R}_\lambda(0) = 0$, and is symmetric around zero ($\mathcal{R}_\lambda(t) = \mathcal{R}_\lambda(-t)$.)
- Non-decreasing and subadditive $\mathcal{R}_\lambda(s + t) \leq \mathcal{R}_\lambda(s) + \mathcal{R}_\lambda(t)$.
- Function $t \mapsto \frac{\mathcal{R}_\lambda(t)}{t}$ is nonincreasing for $t > 0$
- Differentiable for all $t \neq 0$, subdifferentiable at $t = 0$ with subgradients bounded in absolute value by λL .
- For some $\mu > 0$, the function $\tilde{\mathcal{R}}_\lambda(t) = \mathcal{R}_\lambda(t) + \mu t^2$ is convex.

Includes (among others):

- rescaled ℓ_1 loss: $\mathcal{R}_\lambda(t) = \lambda|t|$.
- MCP penalty and SCAD penalties (Fan et al., 2001; Zhang, 2006)
- does **not include** capped ℓ_1 -penalty

Main statistical guarantee

- regularized M -estimator

$$\hat{\theta} \in \arg \min_{\|\theta\|_1 \leq M} \left\{ \mathcal{L}_n(\theta) + \mathcal{R}_\lambda(\theta) \right\}.$$

- loss function satisfies (α, τ) RSC, and regularizer is regular (with parameters (μ, L))

Main statistical guarantee

- regularized M -estimator

$$\hat{\theta} \in \arg \min_{\|\theta\|_1 \leq M} \left\{ \mathcal{L}_n(\theta) + \mathcal{R}_\lambda(\theta) \right\}.$$

- loss function satisfies (α, τ) RSC, and regularizer is regular (with parameters (μ, L))
- local optimum** $\hat{\theta}$ defined by conditions

$$\langle \nabla \mathcal{L}_n(\hat{\theta}) + \nabla \mathcal{R}_\lambda(\hat{\theta}), \theta - \hat{\theta} \rangle \geq 0 \quad \text{for all feasible } \theta.$$

Main statistical guarantee

- regularized M -estimator

$$\hat{\theta} \in \arg \min_{\|\theta\|_1 \leq M} \left\{ \mathcal{L}_n(\theta) + \mathcal{R}_\lambda(\theta) \right\}.$$

- loss function satisfies (α, τ) RSC, and regularizer is regular (with parameters (μ, L))
- local optimum** $\hat{\theta}$ defined by conditions

$$\langle \nabla \mathcal{L}_n(\hat{\theta}) + \nabla \mathcal{R}_\lambda(\hat{\theta}), \theta - \hat{\theta} \rangle \geq 0 \quad \text{for all feasible } \theta.$$

Theorem (Loh & W., 2013)

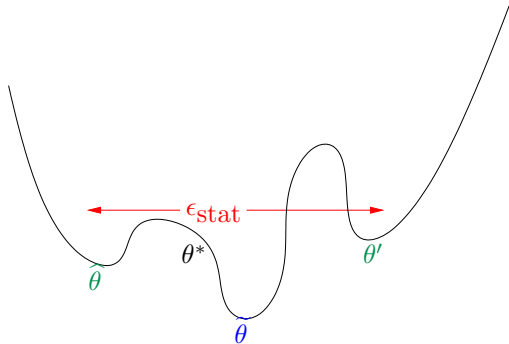
Suppose M is chosen such that θ^* is feasible, and λ satisfies the bounds

$$\max \left\{ \|\nabla \mathcal{L}_n(\theta^*)\|_\infty, \alpha_2 \sqrt{\frac{\log p}{n}} \right\} \leq \lambda \leq \frac{\alpha_2}{6LM}$$

Then **any local optimum** $\hat{\theta}$ satisfies the bound

$$\|\hat{\theta} - \theta^*\|_2 \leq \frac{6 \lambda_n \sqrt{s}}{4(\alpha - \mu)} \quad \text{where } s = \|\beta^*\|_0.$$

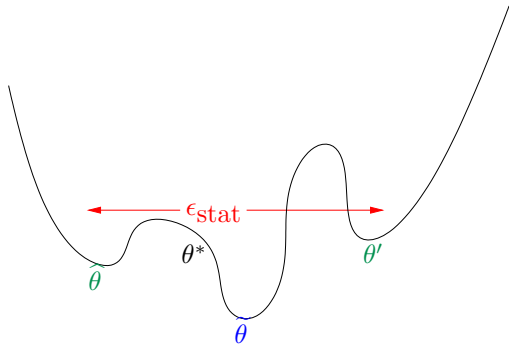
Geometry of local/global optima



Consequence:

All { local, global } optima are within distance ϵ_{stat} of the target θ^* .

Geometry of local/global optima



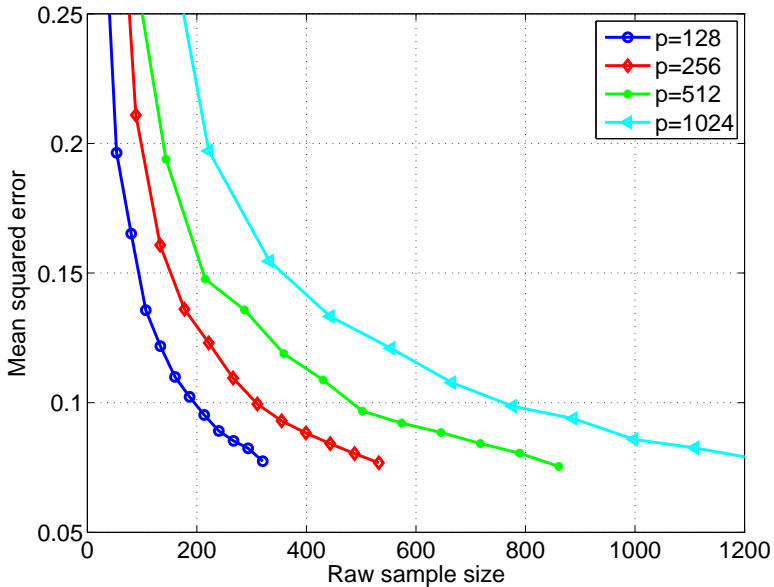
Consequence:

All { local, global } optima are within distance ϵ_{stat} of the target θ^* .

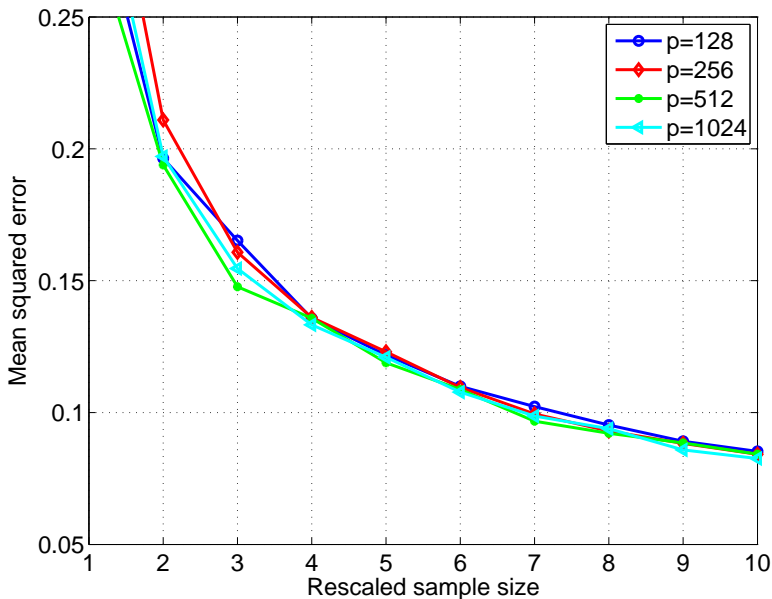
With $\lambda = c \sqrt{\frac{\log p}{n}}$, statistical error scales as

$$\epsilon_{\text{stat}} \asymp \sqrt{\frac{s \log p}{n}}, \quad \text{which is minimax optimal.}$$

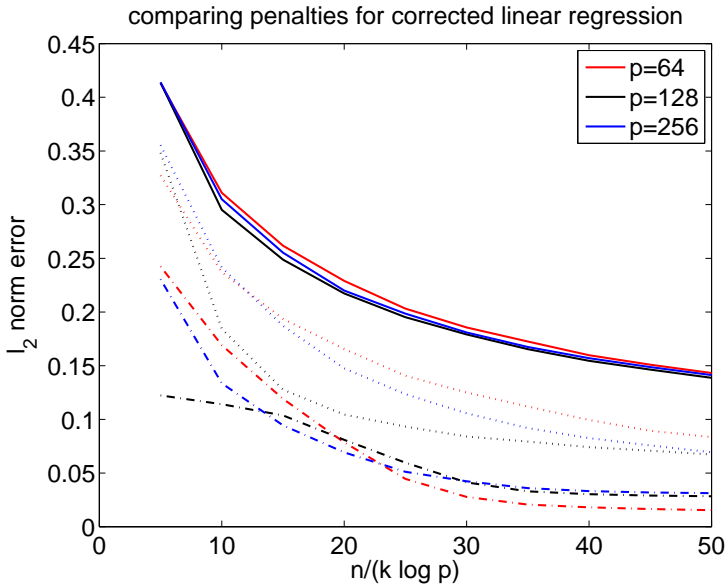
Empirical results (unrescaled)



Empirical results (rescaled)



Comparisons between different penalties



First-order algorithms and fast convergence

Thus far....

- have shown that all local optima are “statistically good”
- how to obtain a local optimum quickly?

First-order algorithms and fast convergence

Thus far....

- have shown that all local optima are “statistically good”
- how to obtain a local optimum quickly?

Composite gradient descent for regularized objectives:

(Nesterov, 2007)

$$\min_{\theta \in \Omega} \left\{ f(\theta) + g(\theta) \right\}$$

where f is differentiable, and g is convex, sub-differentiable.

First-order algorithms and fast convergence

Thus far....

- have shown that all local optima are “statistically good”
- how to obtain a local optimum quickly?

Composite gradient descent for regularized objectives:

(Nesterov, 2007)

$$\min_{\theta \in \Omega} \left\{ f(\theta) + g(\theta) \right\}$$

where f is differentiable, and g is convex, sub-differentiable.

Simple updates:

$$\theta^{t+1} = \arg \min_{\theta \in \Omega} \left\{ \|\theta - \alpha^t \nabla f(\theta^t)\|_2^2 + g(\theta) \right\}.$$

First-order algorithms and fast convergence

Thus far....

- have shown that all local optima are “statistically good”
- how to obtain a local optimum quickly?

Composite gradient descent for regularized objectives:

(Nesterov, 2007)

$$\min_{\theta \in \Omega} \left\{ f(\theta) + g(\theta) \right\}$$

where f is differentiable, and g is convex, sub-differentiable.

Simple updates:

$$\theta^{t+1} = \arg \min_{\theta \in \Omega} \left\{ \|\theta - \alpha^t \nabla f(\theta^t)\|_2^2 + g(\theta) \right\}.$$

Not directly applicable with $f = \mathcal{L}_n$ and $g = \mathcal{R}_\lambda$ (since \mathcal{R}_λ can be non-convex).

Composite gradient on a convenient splitting

- Define modified loss functions and regularizers:

$$\tilde{\mathcal{L}}_n(\theta) := \underbrace{\mathcal{L}_n(\theta) - \mu\|\theta\|_2^2}_{\text{non-convex}}, \quad \text{and} \quad \tilde{\mathcal{R}}_\lambda(\theta) := \underbrace{\mathcal{R}_\lambda(\theta) + \mu\|\theta\|_2^2}_{\text{convex}}.$$

Composite gradient on a convenient splitting

- Define modified loss functions and regularizers:

$$\tilde{\mathcal{L}}_n(\theta) := \underbrace{\mathcal{L}_n(\theta) - \mu\|\theta\|_2^2}_{\text{non-convex}}, \quad \text{and} \quad \tilde{\mathcal{R}}_\lambda(\theta) := \underbrace{\mathcal{R}_\lambda(\theta) + \mu\|\theta\|_2^2}_{\text{convex}}.$$

- Apply composite gradient descent to the objective

$$\min_{\theta \in \Omega} \left\{ \tilde{\mathcal{L}}_n(\theta) + \tilde{\mathcal{R}}_\lambda(\theta) \right\}.$$

Composite gradient on a convenient splitting

- Define modified loss functions and regularizers:

$$\tilde{\mathcal{L}}_n(\theta) := \underbrace{\mathcal{L}_n(\theta) - \mu\|\theta\|_2^2}_{\text{non-convex}}, \quad \text{and} \quad \tilde{\mathcal{R}}_\lambda(\theta) := \underbrace{\mathcal{R}_\lambda(\theta) + \mu\|\theta\|_2^2}_{\text{convex}}.$$

- Apply composite gradient descent to the objective

$$\min_{\theta \in \Omega} \left\{ \tilde{\mathcal{L}}_n(\theta) + \tilde{\mathcal{R}}_\lambda(\theta) \right\}.$$

- converges to local optimum $\hat{\theta}$ (Nesterov, 2007)

$$\langle \nabla \tilde{\mathcal{L}}_n(\theta) + \nabla \tilde{\mathcal{R}}_\lambda(\hat{\theta}), \theta - \hat{\theta} \rangle \geq 0 \quad \text{for all feasible } \theta.$$

Composite gradient on a convenient splitting

- Define modified loss functions and regularizers:

$$\tilde{\mathcal{L}}_n(\theta) := \underbrace{\mathcal{L}_n(\theta) - \mu \|\theta\|_2^2}_{\text{non-convex}}, \quad \text{and} \quad \tilde{\mathcal{R}}_\lambda(\theta) := \underbrace{\mathcal{R}_\lambda(\theta) + \mu \|\theta\|_2^2}_{\text{convex}}.$$

- Apply composite gradient descent to the objective

$$\min_{\theta \in \Omega} \left\{ \tilde{\mathcal{L}}_n(\theta) + \tilde{\mathcal{R}}_\lambda(\theta) \right\}.$$

- converges to local optimum $\hat{\theta}$ (Nesterov, 2007)

$$\langle \nabla \tilde{\mathcal{L}}_n(\theta) + \nabla \tilde{\mathcal{R}}_\lambda(\hat{\theta}), \theta - \hat{\theta} \rangle \geq 0 \quad \text{for all feasible } \theta.$$

- will show that convergence is **geometrically fast** with constant stepsize

Theoretical guarantees on computational error

- implement Nesterov's composite method with constant stepsize to $(\tilde{\mathcal{L}}_n, \tilde{\mathcal{R}}_\lambda)$ split.
- fixed global optimum $\hat{\beta}$ defines the **statistical error** $\epsilon_{\text{stat}}^2 = \|\hat{\beta} - \beta^*\|_2$.
- population minimizer β^* is s -sparse

Theoretical guarantees on computational error

- implement Nesterov's composite method with constant stepsize to $(\tilde{\mathcal{L}}_n, \tilde{\mathcal{R}}_\lambda)$ split.
- fixed global optimum $\hat{\beta}$ defines the **statistical error** $\epsilon_{\text{stat}}^2 = \|\hat{\beta} - \beta^*\|_2$.
- population minimizer β^* is s -sparse
- loss function satisfies (α, τ) **RSC** and smoothness conditions, and regularizer is (μ, L) -good

Theoretical guarantees on computational error

- implement Nesterov's composite method with constant stepsize to $(\tilde{\mathcal{L}}_n, \tilde{\mathcal{R}}_\lambda)$ split.
- fixed global optimum $\hat{\beta}$ defines the **statistical error** $\epsilon_{\text{stat}}^2 = \|\hat{\beta} - \beta^*\|_2$.
- population minimizer β^* is s -sparse
- loss function satisfies (α, τ) RSC and smoothness conditions, and regularizer is (μ, L) -good

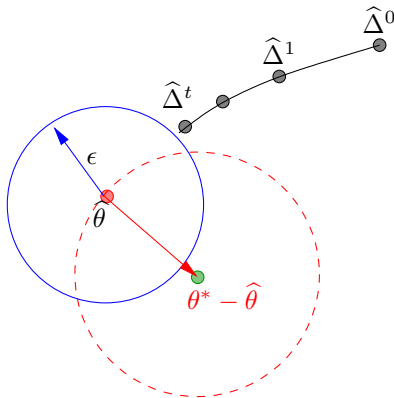
Theorem (Loh & W., 2013)

If $n \gtrsim s \log p$, there is a *contraction factor* $\kappa \in (0, 1)$ such that for any $\delta \geq \epsilon_{\text{stat}}$, we have

$$\|\theta^t - \hat{\theta}\|_2^2 \leq \frac{2}{\alpha - \mu} \left(\delta^2 + 128\tau \frac{s \log p}{n} \epsilon_{\text{stat}}^2 \right) \quad \text{for all } t \geq T(\delta) \text{ iterations,}$$

where $T(\delta) \asymp \frac{\log(1/\delta)}{\log(1/\kappa)}$.

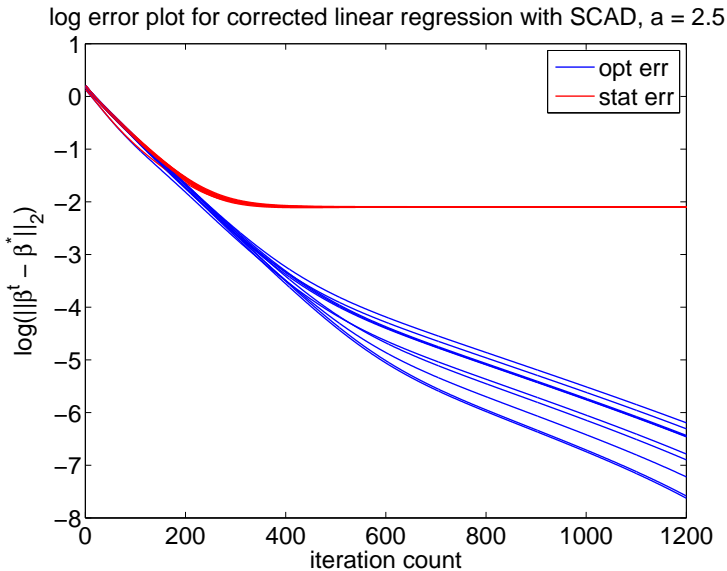
Geometry of result



Optimization error $\hat{\Delta}^t := \theta^t - \hat{\theta}$ decreases geometrically up to statistical tolerance:

$$\|\theta^{t+1} - \hat{\theta}\|^2 \leq \kappa^t \|\theta^0 - \hat{\theta}\|^2 + o\left(\underbrace{\|\theta^* - \hat{\theta}\|^2}_{\text{Stat. error } \epsilon_{\text{stat}}^2}\right) \quad \text{for all } t = 0, 1, 2, \dots$$

Non-convex linear regression with SCAD



Non-convex linear regression with SCAD



Summary

- M -estimators based on non-convex programs arise frequently
- under suitable regularity conditions, we showed that:
 - ▶ all local optima are “well-behaved” from the statistical point of view
 - ▶ simple first-order methods converge as fast as possible

Summary

- M -estimators based on non-convex programs arise frequently
- under suitable regularity conditions, we showed that:
 - ▶ all local optima are “well-behaved” from the statistical point of view
 - ▶ simple first-order methods converge as fast as possible
- many open questions
 - ▶ similar guarantees for more general problems?
 - ▶ geometry of non-convex problems in statistics?

Summary

- M -estimators based on non-convex programs arise frequently
 - under suitable regularity conditions, we showed that:
 - ▶ all local optima are “well-behaved” from the statistical point of view
 - ▶ simple first-order methods converge as fast as possible

 - many open questions
 - ▶ similar guarantees for more general problems?
 - ▶ geometry of non-convex problems in statistics?
-

Papers and pre-prints:

- Loh & W. (2013). Regularized M -estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Pre-print arXiv:1305.2436*
- Loh & W. (2012). High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity. *Annals of Statistics*, 40:1637–1664.
- Negahban et al. (2012). A unified framework for high-dimensional analysis of M -estimators. *Statistical Science*, 27(4): 538–557.