

Lecture 2 – Part 1

Optimization

(January 16, 2015)

Mu Zhu
University of Waterloo

Need for Optimization

- $\mathbb{E}(y|\mathbf{x})$, $\mathbb{P}(y|\mathbf{x})$ — want to “go after” them
- first, [model](#)
 - some examples last week
- then, [estimate](#)
 - didn't discuss last week
 - often requires solving an [optimization](#) problem

Ex I: Linear Regression

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^n (y_i - \mathbf{x}_i^{\top} \boldsymbol{\beta})^2 = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$$

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^{\top} \\ \vdots \\ \mathbf{x}_n^{\top} \end{bmatrix}$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{y}$$

Ex II: Logistic Regression

$$y_i \sim \text{Bernoulli}(p_i)$$

$$p_i \equiv \mathbb{P}(y_i = 1 | \mathbf{x}_i)$$

$$\log \frac{p_i}{1 - p_i} = \mathbf{x}_i^T \boldsymbol{\beta} \quad \text{or} \quad p_i = \frac{\exp(\mathbf{x}_i^T \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta})}$$

$$\mathcal{L}(p_i; y_i) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1 - y_i}$$

Ex II: Logistic Regression

- log-likelihood

$$\begin{aligned}\ell(p_i; y_i) &= \log \mathcal{L}(p_i; y_i) \\ &= \sum_{i=1}^n y_i \log(p_i) + (1 - y_i) \log(1 - p_i) \\ &= \sum_{i=1}^n y_i \log \frac{p_i}{1 - p_i} + \log(1 - p_i)\end{aligned}$$

$$\ell(\boldsymbol{\beta}; y_i, \mathbf{x}_i) = \sum_{i=1}^n y_i \mathbf{x}_i^T \boldsymbol{\beta} - \log(1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}})$$

- to estimate $\boldsymbol{\beta}$,

$$\max_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta}; y_i, \mathbf{x}_i)$$

by [Newton-Raphson](#)

Ex II: Logistic Regression

$$\begin{aligned}\ell'(\boldsymbol{\beta}) &= \sum_{i=1}^n y_i \mathbf{x}_i - \left[\frac{e^{\mathbf{x}_i^T \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}}} \right] \mathbf{x}_i = \sum_{i=1}^n (y_i - p_i) \mathbf{x}_i \\ &= \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n \end{bmatrix} \begin{bmatrix} y_1 - p_1 \\ \vdots \\ y_n - p_n \end{bmatrix} = \mathbf{X}^T (\mathbf{y} - \mathbf{p})\end{aligned}$$

Ex II: Logistic Regression

$$\begin{aligned} & \ell''(\boldsymbol{\beta}) \\ = & - \sum_{i=1}^n \mathbf{x}_i \left[\frac{(e^{\mathbf{x}_i^T \boldsymbol{\beta}})(1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}}) - (e^{\mathbf{x}_i^T \boldsymbol{\beta}})(e^{\mathbf{x}_i^T \boldsymbol{\beta}})}{(1 + e^{\mathbf{x}_i^T \boldsymbol{\beta}})^2} \right] \mathbf{x}_i^T \\ = & - \sum_{i=1}^n \mathbf{x}_i [p_i - p_i^2] \mathbf{x}_i^T \\ = & - \begin{bmatrix} \mathbf{x}_1 & \dots & \mathbf{x}_n \end{bmatrix} \begin{bmatrix} p_1(1 - p_1) & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & p_n(1 - p_n) \end{bmatrix} \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \\ = & -\mathbf{X}^T \mathbf{W} \mathbf{X} \end{aligned}$$

Ex II: Logistic Regression

$$\begin{aligned}\beta_{new} &= \beta_{old} - \left[\ell''(\beta_{old}) \right]^{-1} \left[\ell'(\beta_{old}) \right] \\ &= \beta_{old} + \left[\mathbf{X}^T \mathbf{W} \mathbf{X} \right]^{-1} \left[\mathbf{X}^T (\mathbf{y} - \mathbf{p}) \right] \\ &= \underbrace{\left[\mathbf{X}^T \mathbf{W} \mathbf{X} \right]^{-1} \mathbf{X}^T \mathbf{W}}_{\text{weighted least squares}} \left[\mathbf{X} \beta_{old} + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p}) \right]\end{aligned}$$

(both \mathbf{W} and \mathbf{p} depend on β_{old})

$w_{ii} \equiv p_i(1 - p_i)$ **max** at $p_i = 1/2$ and **min** at $p_i = 0$ or $p_i = 1$



estimates influenced mostly by points near decision boundary

Enter “Big Data”

- if all this sounds easy, think about $\boldsymbol{x} \in \mathbb{R}^d$ and d is large relative to n
- end up **estimating too much with too little**
- resulting estimate cannot be very good
- in statistically terms, $\text{Var}(\hat{\boldsymbol{\beta}})$ inflated
- to reduce **variance**, introduce **bias**

Penalized Regression

- let $\mathbf{x}_1, \dots, \mathbf{x}_d$ denote **columns** of \mathbf{X}
- suppose $\mathbf{y}, \mathbf{x}_1, \dots, \mathbf{x}_d$ all standardized ($\mathbf{1}^\top \mathbf{y} = 0$, $\|\mathbf{y}\| = 1$, etc)
- bias each β_j by introducing **penalty** $J(\beta_j)$

$$\min_{\boldsymbol{\beta}} \frac{1}{2} \|\mathbf{y} - (\beta_1 \mathbf{x}_1 + \dots + \beta_d \mathbf{x}_d)\|^2 + \sum_{j=1}^d J(\beta_j)$$

A Class of Penalty Functions

$$J(\beta_j) = \lambda |\beta_j|^\alpha \quad \Rightarrow \quad \sum_{j=1}^d J(\beta_j) = \lambda \sum_{j=1}^d |\beta_j|^\alpha \equiv \lambda \|\boldsymbol{\beta}\|_\alpha$$

	α	$\ \cdot\ _\alpha$	difficulty	algorithm
ridge regression	2	ℓ_2	convex	analytic (exercise)
LASSO	1	ℓ_1	convex	coordinate descent
subset selection	0	ℓ_0	NP-hard	heuristic

$$|\beta_j|^0 = \begin{cases} 0, & \beta_j = 0 \\ 1, & \beta_j \neq 0 \end{cases} \quad \Rightarrow \quad \sum |\beta_j|^0 = \sum I(\beta_j \neq 0)$$

Bias-Variance Trade-off

Exercise Show that, under typical model assumptions,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \mathbb{E}(\boldsymbol{\varepsilon}) = \mathbf{0}, \quad \text{and} \quad \text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I},$$

(i) the “usual” estimate, $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$, is unbiased;

(ii) the ridge regression estimate — call it $\hat{\boldsymbol{\beta}}_\lambda$ — is biased, but \exists orthonormal matrix \mathbf{V} such that

$$\text{Var}(\hat{\boldsymbol{\beta}}_\lambda) = \sigma^2 \mathbf{V} \mathbf{D}_\lambda \mathbf{V}^T, \quad \text{Var}(\hat{\boldsymbol{\beta}}) = \sigma^2 \mathbf{V} \mathbf{D} \mathbf{V}^T,$$

and

$$\mathbf{D}_\lambda(j, j) \leq \mathbf{D}(j, j) \quad \forall \quad j.$$

(Hint: Use the [singular value decomposition](#) of \mathbf{X} .)

Ex III: LASSO

- coordinate descent — at each iteration, solve a univariate problem,

$$\min_{\beta_j} L(\beta_j) = \frac{1}{2} \|\mathbf{z} - \beta_j \mathbf{x}_j\|_2^2 + \lambda |\beta_j| + c_j,$$

where

$$\mathbf{z} \equiv \mathbf{y} - \sum_{k \neq j} \beta_k \mathbf{x}_k, \quad c_j = \sum_{k \neq j} \lambda |\beta_k|,$$

while fixing all β_k , $k \neq j$.

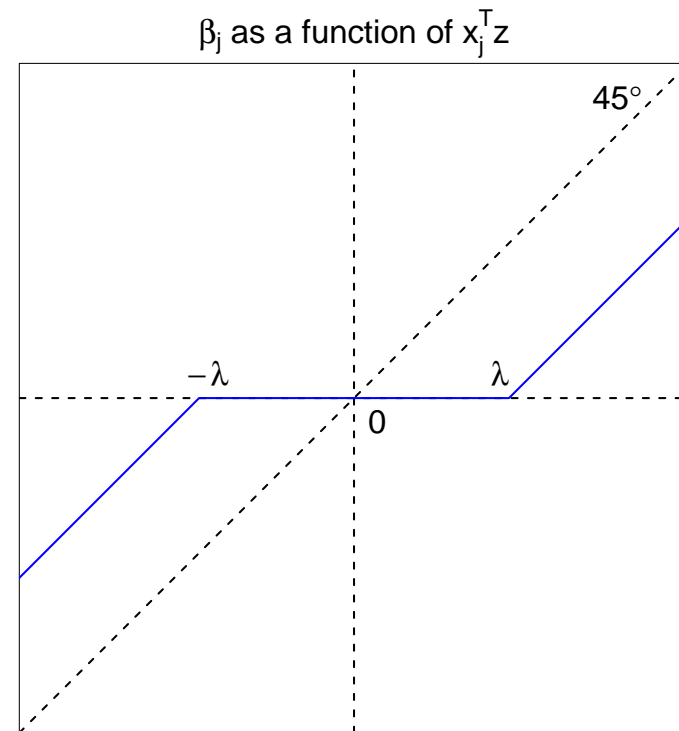
- cycle through $j = 1, 2, \dots, d, 1, 2, \dots, d, \dots$

Ex III: LASSO

$$\frac{d}{d\beta_j} L(\beta_j) = 0$$

$$\beta_j \underbrace{(\mathbf{x}_j^T \mathbf{x}_j)}_{\|\mathbf{x}_j\|^2=1} + \lambda \cdot \text{sgn}(\beta_j) = \mathbf{x}_j^T \mathbf{z}$$

$$\beta_j = \begin{cases} \mathbf{x}_j^T \mathbf{z} - \lambda, & \beta_j > 0; \\ \mathbf{x}_j^T \mathbf{z} + \lambda, & \beta_j < 0 \end{cases}$$



solution will typically contain many zeros, i.e., be [sparse](#)
selection effect of the LASSO

The ℓ_1 Penalty

D. L. Donoho (2006), “For most large underdetermined systems of linear equations the minimal ℓ_1 -norm solution is also the sparsest solution”, *Communications on Pure and Applied Mathematics* **59**, pp. 797–829.

- ℓ_1 -problem as a [convex relaxation](#) of ℓ_0 -problem

Ex IV: Graphical LASSO

- from $\mathbf{x}_i \stackrel{iid}{\sim} \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, estimate $\boldsymbol{\Omega} \equiv \boldsymbol{\Sigma}^{-1}$ (hard for d large)
- recall **linear discriminant analysis** requires $\boldsymbol{\Sigma}^{-1}$ (last week)
- for $\mathbf{x} = (x_1, \dots, x_d)^T \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, variables x_j and x_k **conditionally independent** given **all other** variables if and only if $\boldsymbol{\Omega}_{jk} = 0$ (Dempster, 1972; *Biometrics*)
- for **Gaussian graphical models**, \exists edge between nodes j and k if and only if $\boldsymbol{\Omega}_{jk} \neq 0$

Ex IV: Graphical LASSO

$$\hat{\boldsymbol{\mu}} = \bar{\boldsymbol{x}} = \frac{1}{n} \sum_{i=1}^n \boldsymbol{x}_i, \quad \boldsymbol{S} = \frac{1}{n} \sum_{i=1}^n (\boldsymbol{x}_i - \hat{\boldsymbol{\mu}})(\boldsymbol{x}_i - \hat{\boldsymbol{\mu}})^\top$$

$$\mathcal{L}(\boldsymbol{\Sigma}) = \prod_{i=1}^n \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left[-\frac{1}{2} (\boldsymbol{x}_i - \hat{\boldsymbol{\mu}})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_i - \hat{\boldsymbol{\mu}}) \right]$$

$$\begin{aligned} \ell(\boldsymbol{\Omega}) &= \text{const} + \frac{n}{2} \log |\boldsymbol{\Sigma}^{-1}| - \frac{1}{2} \sum_{i=1}^n \text{tr} [(\boldsymbol{x}_i - \hat{\boldsymbol{\mu}})^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_i - \hat{\boldsymbol{\mu}})] \\ &= \text{const} + \frac{n}{2} \log |\boldsymbol{\Omega}| - \frac{n}{2} \text{tr} (\boldsymbol{\Omega} \boldsymbol{S}) \end{aligned}$$

Ex IV: Graphical LASSO

- maximize ℓ_1 -penalized likelihood (Friedman, Hastie & Tibshirani, 2008; *Biostatistics*):

$$\max_{\Omega \succeq 0} \underbrace{\log |\Omega| - \text{tr}(\Omega S)}_{\ell(\Omega)} - \lambda \|\Omega\|_1$$

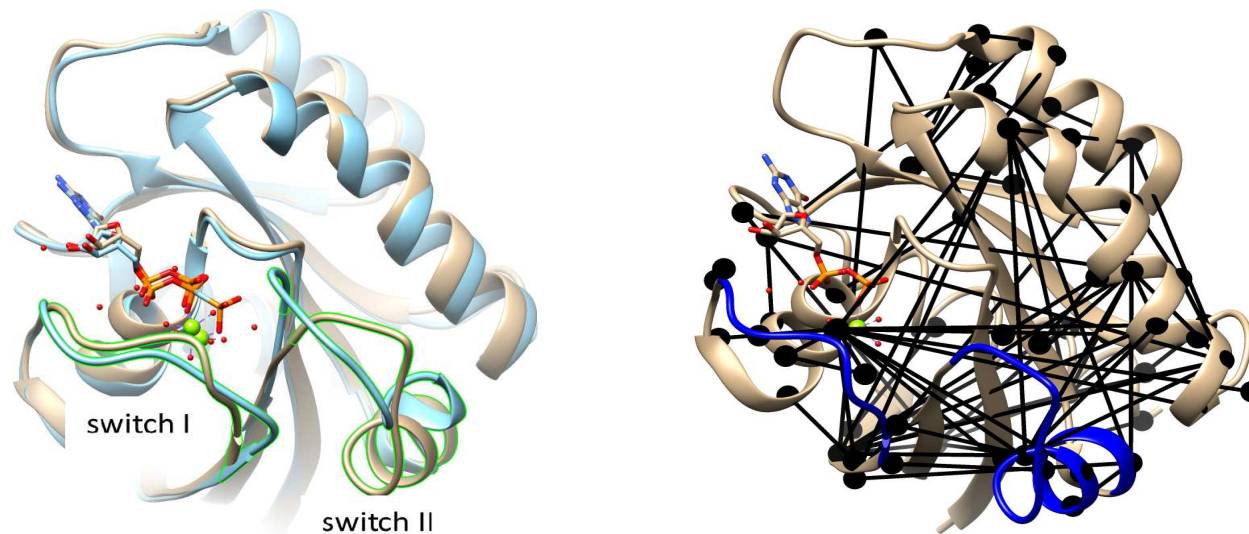
where

$$\|\Omega\|_1 = \sum_{j,k} |\Omega_{jk}|$$

by [coordinate descent](#) (one row/column at a time)

- resulting graphical model is [sparse](#) in the number of edges

Ex IV: Graphical LASSO



H-Ras. *Left*: Inactive form (4q21) & active form (6q21). Regions labelled “Switch I” & “Switch II” are known to undergo major conformational changes between 4q21 and 6q21. *Right*: Estimated graphical model showing conditional dependence structures, obtained by analysing 4q21 alone. (L. Soltan-Ghoraie, F. Burkowski & M. Zhu)

Ex V: The Netflix Problem (Part 1)

- high-profile, million-dollar Netflix contest (2006-09)
- rating matrix \mathbf{R} where

$$r_{ui} = \text{rating of item } i \text{ by user } u$$

- set $\mathbb{T} = \{(u, i) : r_{ui} \text{ observed}\}$
- want to predict the missing entries of \mathbf{R}

Ex V: The Netflix Problem (Part 1)













	Item 1	Item 2	Item 3	Item 4
User A				
User B				
User C				
User D				

Illustration of the Rating Matrix R

Ex V: The Netflix Problem (Part 1)

- want to solve

$$\min_{\hat{\mathbf{R}}} \text{rank}(\hat{\mathbf{R}}) \quad \text{s.t.} \quad \hat{r}_{ui} = r_{ui} \text{ for } (u, i) \in \mathbb{T}$$

- philosophy: only a few factors affect user preferences, so **rank** of rating matrix must be low
- but problem is **NP-hard**

Ex V: The Netflix Problem (Part 1)

- instead, solve **convex relaxation**,

$$\min_{\hat{\mathbf{R}}} \|\hat{\mathbf{R}}\|_* \quad \text{s.t.} \quad \hat{r}_{ui} = r_{ui} \text{ for } (u, i) \in \mathbb{T}$$

where $\|\cdot\|_*$ denotes the **nuclear norm** of a matrix

- let $\sigma_1, \sigma_2, \dots$ be **singular values** of \mathbf{A} ; then,

$$\|\mathbf{A}\|_* = \sum |\sigma_j|^1$$

whereas

$$\text{rank}(\mathbf{A}) = \sum I(\sigma_j \neq 0) = \sum |\sigma_j|^0$$

- a matter of ℓ_1 vs ℓ_0

Ex V: The Netflix Problem (Part 1)

E. J. Candès & B. Recht (2009), “Exact matrix completion via convex optimization”, *Found Comput Math* **9**, pp. 717–772.

B. Recht, M. Fazel, & P. A. Parrilo (2010), “Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization”, *SIAM Review* **52**, pp. 471–501.

- under certain conditions, $\min \| \cdot \|_* \Leftrightarrow \min \text{rank}(\cdot)$
- $\min \| \cdot \|_* \Leftrightarrow$ a semi-definite program

Ex VI: The Netflix Problem (Part 2)

- explicit parameterization

$$\mathbf{R} \approx \mathbf{P}\mathbf{Q}^T = \begin{bmatrix} \mathbf{p}_1^T \\ \vdots \\ \mathbf{p}_n^T \end{bmatrix} \begin{bmatrix} \mathbf{q}_1 & \cdots & \mathbf{q}_m \end{bmatrix}$$

- $\mathbf{p}_u, \mathbf{q}_i \in \mathbb{R}^K$ ($K \ll n, m$) are **latent coordinates**
- just **estimate** $\mathbf{p}_u, \mathbf{q}_i$ and **predict** missing entries with $\hat{r}_{ui} = \mathbf{p}_u^T \mathbf{q}_i$
- in reality, user- and item-effects are removed prior to doing this

Ex VI: The Netflix Problem (Part 2)

$$\min_{\mathbf{p}_u, \mathbf{q}_i} \sum_{(u,i) \in \mathbb{T}} (r_{ui} - \mathbf{p}_u^T \mathbf{q}_i)^2 + \lambda \left[\sum_u \|\mathbf{p}_u\|^2 + \sum_i \|\mathbf{q}_i\|^2 \right]$$

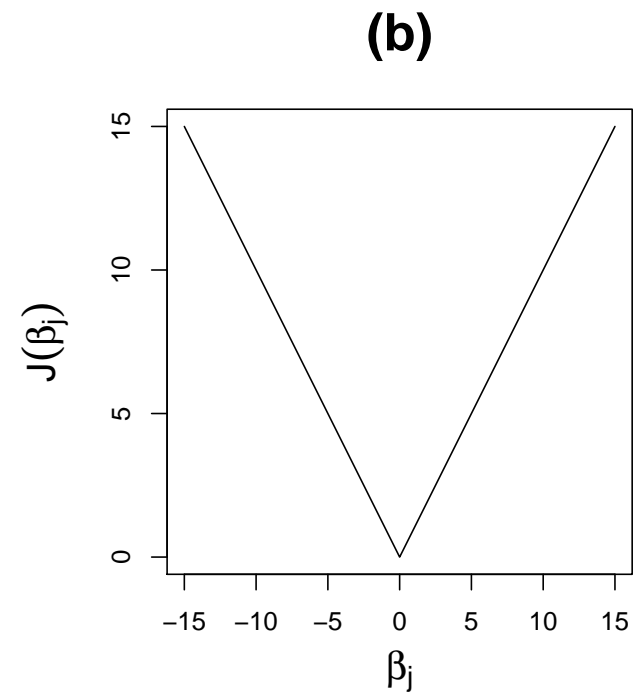
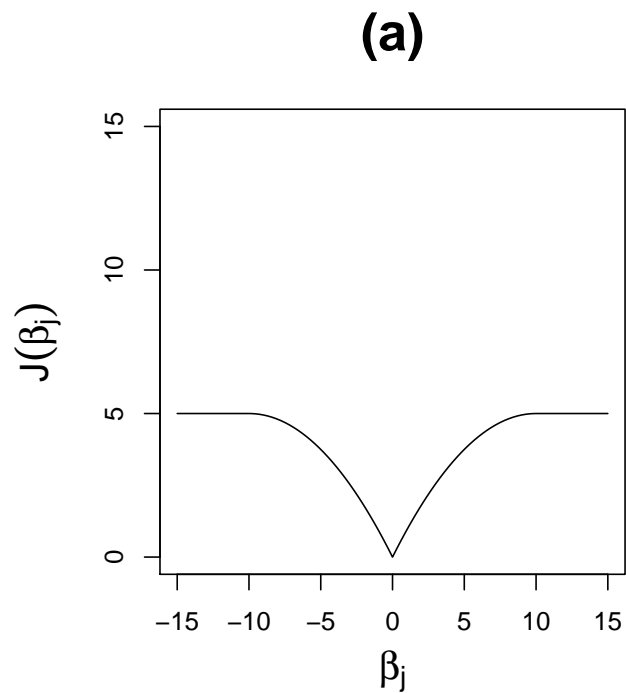
- **coordinate descent** still applies (over $\mathbf{p}_1, \dots, \mathbf{p}_n, \mathbf{q}_1, \dots, \mathbf{q}_m, \dots$)
- strictly speaking, blockwise coordinate descent
- each step **convex**, but overall **nonconvex** (both $\mathbf{p}_u, \mathbf{q}_i$ unknown)
- many traps (e.g., **local solutions**, **saddle points**)

Ex VII: MCP

For the penalized regression problem, Zhang (2010; *Ann. Stat.*) proposed the so-called **minimax concave penalty** (MCP):

$$J(\beta_j) = \lambda \int_0^{|\beta_j|} \left(1 - \frac{x}{\gamma\lambda}\right)_+ dx = \begin{cases} \lambda|\beta_j| - \frac{\beta_j^2}{2\gamma}, & |\beta_j| \leq \gamma\lambda; \\ \frac{1}{2}\gamma\lambda^2, & |\beta_j| > \gamma\lambda, \end{cases}$$

instead of $J(\beta_j) = \lambda|\beta_j|^\alpha$.



(a) MCP ($\lambda = 1, \gamma = 10$); (b) LASSO ($\lambda = 1$).

Ex VII: MCP

- for $|\beta_j|$ large,
 - $J(\cdot)$ constant “beyond a certain point” — reduces **bias**
- for $|\beta_j|$ small,
 - $J(\cdot)$ still “like the LASSO” — keeps **sparsity**
- in between,
 - smooth interpolation that minimizes maximal concavity
- coordinate descent applies,
 - but many traps (e.g., **local solutions**, **saddle points**)

Summary

- key ideas:
 - introduce bias to reduce variance; penalty functions
 - ℓ_1 norm; nuclear norm; nonconvex penalties
- specific methods:
 - ridge; LASSO; MCP
 - graphical LASSO
 - matrix completion; matrix factorization
 - Newton-Raphson; coordinate descent
- application areas:
 - proteins
 - recommender systems

Next ...

- lecture @ 2 pm by Professor S. Vavasis on [optimization](#)
- a short, 10-minute break
- research with my students (mostly, the Netflix Problem)