

Statistics for high-dimensional data: Introduction, and the Lasso for linear models

Peter Bühlmann and Sara van de Geer

Seminar für Statistik, ETH Zürich

May 2012

High-dimensional data

Riboflavin production with *Bacillus Subtilis*

(in collaboration with DSM (Switzerland))

goal: improve riboflavin production rate of *Bacillus Subtilis*
using clever genetic engineering

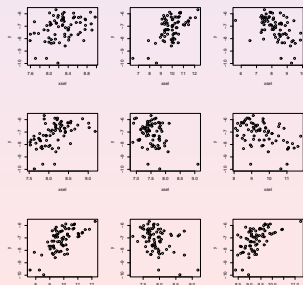
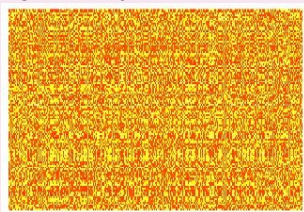
response variables $Y \in \mathbb{R}$: riboflavin (log-) production rate

covariates $X \in \mathbb{R}^p$: expressions from $p = 4088$ genes

sample size $n = 115$, $p \gg n$

Y versus 9 “reasonable” genes

gene expression data



general framework:

Z_1, \dots, Z_n (with some "i.i.d. components")

$\dim(Z_i) \gg n$

for example:

$Z_i = (X_i, Y_i)$, $X_i \in \mathcal{X} \subseteq \mathbb{R}^p$, $Y_i \in \mathcal{Y} \subseteq \mathbb{R}$: regression with $p \gg n$

$Z_i = (X_i, Y_i)$, $X_i \in \mathcal{X} \subseteq \mathbb{R}^p$, $Y_i \in \{0, 1\}$: classification for $p \gg n$

numerous applications:

biology, imaging, economy, environmental sciences, ...

High-dimensional linear models

$$Y_i = \sum_{j=1}^p \beta_j^0 X_i^{(j)} + \varepsilon_i, \quad i = 1, \dots, n$$

$$p \gg n$$

$$\text{in short: } \mathbf{Y} = \mathbf{X}\beta + \varepsilon$$

goals:

- ▶ prediction, e.g. w.r.t. squared prediction error
- ▶ estimation of β^0 , e.g. w.r.t. $\|\hat{\beta} - \beta^0\|_q$ ($q = 1, 2$)
- ▶ variable selection
i.e. estimating the active set with the effective variables
(having corresponding coefficient $\neq 0$)

we need to **regularize**...

and there are many proposals

- ▶ Bayesian methods for regularization
- ▶ greedy algorithms: aka forward selection or boosting
- ▶ preliminary dimension reduction
- ▶ ...

e.g. 3'190'000 entries on Google Scholar for
“high dimensional linear model” ...

we need to **regularize**...
and there are many proposals

- ▶ Bayesian methods for regularization
- ▶ greedy algorithms: aka forward selection or boosting
- ▶ preliminary dimension reduction
- ▶ ...

e.g. 3'190'000 entries on Google Scholar for
“high dimensional linear model” ...

we need to **regularize**...
and there are many proposals

- ▶ Bayesian methods for regularization
- ▶ greedy algorithms: aka forward selection or boosting
- ▶ preliminary dimension reduction
- ▶ ...

e.g. 3'190'000 entries on Google Scholar for
“high dimensional linear model” ...

Penalty-based methods

if true β^0 is sparse w.r.t.

- ▶ $\|\beta^0\|_0^0 =$ number of non-zero coefficients
 \leadsto regularize with the $\|\cdot\|_0$ -penalty:
 $\operatorname{argmin}_{\beta}(n^{-1}\|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda\|\beta\|_0^0)$, e.g. AIC, BIC
 \leadsto computationally infeasible if p is large (2^p sub-models)
- ▶ $\|\beta^0\|_1 = \sum_{j=1}^p |\beta_j^0|$
 \leadsto penalize with the $\|\cdot\|_1$ -norm, i.e. Lasso:
 $\operatorname{argmin}_{\beta}(n^{-1}\|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda\|\beta\|_1)$
 \leadsto convex optimization:
 computationally feasible and very fast for large p

The Lasso (Tibshirani, 1996)

Lasso for linear models

$$\hat{\beta}(\lambda) = \underset{\beta}{\operatorname{argmin}} \left(n^{-1} \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \underbrace{\lambda}_{\geq 0} \underbrace{\|\beta\|_1}_{\sum_{j=1}^p |\beta_j|} \right)$$

↪ **convex** optimization problem

- ▶ Lasso **does variable selection**
some of the $\hat{\beta}_j(\lambda) = 0$
(because of “ ℓ_1 -geometry”)
- ▶ $\hat{\beta}(\lambda)$ is a **shrunk LS-estimate**

more about “ ℓ_1 -geometry”

equivalence to primal problem

$$\hat{\beta}_{\text{primal}}(R) = \operatorname{argmin}_{\beta; \|\beta\|_1 \leq R} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2/n,$$

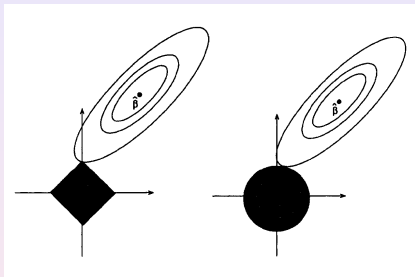
with a correspondence between λ and R which depends on the data $(X_1, Y_1), \dots, (X_n, Y_n)$

[such an equivalence holds since

- ▶ $\|\mathbf{Y} - \mathbf{X}\beta\|_2^2/n$ is convex in β
- ▶ convex constraint $\|\beta\|_1 \leq R$

see e.g. [Bertsekas \(1995\)](#)]

$p=2$



left: ℓ_1 -“world”

residual sum of squares reaches a minimal value (for certain constellations of the data) if its contour lines hit the ℓ_1 -ball in its corner

$$\leadsto \hat{\beta}_1 = 0$$

l_2 -“world” is different

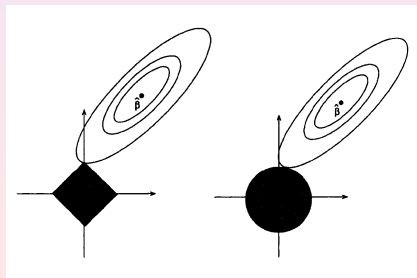
Ridge regression,

$$\hat{\beta}_{\text{Ridge}}(\lambda) = \operatorname{argmin}_{\beta} \left(\|\mathbf{Y} - \mathbf{X}\beta\|_2^2/n + \lambda\|\beta\|_2^2 \right),$$

equivalent primal equivalent solution

$$\hat{\beta}_{\text{Ridge};\text{primal}}(R) = \operatorname{argmin}_{\beta; \|\beta\|_2 \leq R} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2/n,$$

with a one-to-one correspondence between λ and R



A note on the Bayesian approach

model:

$$\beta_1, \dots, \beta_p \text{ i.i.d. } \sim p(\beta) d\beta,$$

$$\text{given } \beta : \mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\beta, \sigma^2 I_n) \text{ with density } f(\mathbf{y}|\sigma^2, \beta)$$

posterior density:

$$p(\beta|\mathbf{Y}, \sigma^2) = \frac{f(\mathbf{Y}|\beta, \sigma^2)p(\beta)}{\int f(\mathbf{Y}|\beta, \sigma^2)p(\beta) d\beta} \propto f(\mathbf{Y}|\beta, \sigma^2)p(\beta)$$

and hence for the MAP (Maximum A-Posteriori) estimator:

$$\begin{aligned} \hat{\beta}_{\text{MAP}} &= \operatorname{argmax}_{\beta} p(\beta|\mathbf{Y}, \sigma^2) = \operatorname{argmin}_{\beta} -\log \left(f(\mathbf{Y}|\beta, \sigma^2)p(\beta) \right) \\ &= \operatorname{argmin}_{\beta} \left(\frac{1}{2\sigma^2} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 - \sum_{j=1}^p \log(p(\beta_j)) \right) \end{aligned}$$

examples:

1. Double-Exponential prior $\text{DExp}(\xi)$:

$$p(\beta) = \frac{\tau}{2} \exp(-\tau\beta)$$

$\leadsto \hat{\beta}_{\text{MAP}}$ equals the Lasso with penalty parameter $\lambda = n^{-1}2\sigma^2\tau$

2. Gaussian prior $\mathcal{N}(0, \tau^2)$:

$$p(\beta) = \frac{1}{\sqrt{2\pi\tau}} \exp(-\beta^2/(2\tau^2))$$

$\leadsto \hat{\beta}_{\text{MAP}}$ equals the Ridge estimator with penalty parameter $\lambda = n^{-1}\sigma^2/\tau^2$

but we will argue that Lasso is also good if the truth is sparse with respect to $\|\beta^0\|_0$, e.g. if prior is (much) more spiky around zero than Double-Exponential distribution

Orthonormal design

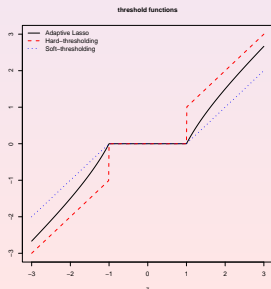
$$\mathbf{Y} = \mathbf{X}\beta + \varepsilon, \quad n^{-1}\mathbf{X}^T\mathbf{X} = \mathbf{I}$$

Lasso = soft-thresholding estimator

$$\hat{\beta}_j(\lambda) = \text{sign}(Z_j)(|Z_j| - \lambda/2)_+, \quad Z_j = \underbrace{(n^{-1}\mathbf{X}^T\mathbf{Y})_j}_{=\text{OLS}}$$

$$\hat{\beta}_j(\lambda) = g_{\text{soft}}(Z_j),$$

[this is Exercise 2.1]



Prediction

goal: predict a new observation Y_{new}

consider expected (w.r.t. new data; and random X) squared error loss:

$$\begin{aligned}\mathbb{E}_{X_{\text{new}}, Y_{\text{new}}} [(Y_{\text{new}} - X_{\text{new}}\hat{\beta})^2] &= \sigma^2 + \mathbb{E}_{X_{\text{new}}} [(X_{\text{new}}(\beta^0 - \hat{\beta}))^2] \\ &= \sigma^2 + (\hat{\beta} - \beta^0)^T \underbrace{\Sigma}_{\text{Cov}(X)} (\hat{\beta} - \beta^0)\end{aligned}$$

→ terminology “prediction error”:

for random design \mathbf{X} : $(\hat{\beta} - \beta^0)^T \Sigma (\hat{\beta} - \beta^0) = \mathbb{E}_{X_{\text{new}}} [(X_{\text{new}}(\hat{\beta} - \beta^0))^2]$

for fixed design \mathbf{X} : $(\hat{\beta} - \beta^0)^T \hat{\Sigma} (\hat{\beta} - \beta^0) = \|\mathbf{X}(\hat{\beta} - \beta^0)\|_2^2/n$

binary lymph node classification using gene expressions:

a high noise problem

$n = 49$ samples, $p = 7130$ gene expressions

despite that it is classification: $\mathbb{P}[Y = 1|X = x] = \mathbb{E}[Y|X = x]$

$\leadsto \hat{p}(x)$ via linear model; can then do classification

cross-validated misclassification error (2/3 training; 1/3 test)

Lasso	L_2 Boosting	FPLR	Pelora	1-NN	DLDA	SVM
21.1%	17.7%	35.25%	27.8%	43.25%	36.12%	36.88%

with variable selection

best 200 genes (Wilcoxon test)
no additional variable selection

Lasso selected on CV-average **13.12 out of $p = 7130$** genes

from a practical perspective:

if you trust in cross-validation: can validate how good we are
i.e. prediction may be a black box, but we can evaluate it!

binary lymph node classification using gene expressions:

a high noise problem

$n = 49$ samples, $p = 7130$ gene expressions

despite that it is classification: $\mathbb{P}[Y = 1|X = x] = \mathbb{E}[Y|X = x]$

$\leadsto \hat{p}(x)$ via linear model; can then do classification

cross-validated misclassification error (2/3 training; 1/3 test)

Lasso	L_2 Boosting	FPLR	Pelora	1-NN	DLDA	SVM
21.1%	17.7%	35.25%	27.8%	43.25%	36.12%	36.88%

with variable selection

best 200 genes (Wilcoxon test)
no additional variable selection

Lasso selected on CV-average **13.12 out of $p = 7130$** genes

from a practical perspective:

if you trust in cross-validation: can validate how good we are
i.e. prediction may be a black box, but we can evaluate it!

and in fact: we will hear that

- ▶ Lasso is consistent for prediction assuming “essentially nothing”
- ▶ Lasso is optimal for prediction assuming the “compatibility condition” for \mathbf{X}

Estimation of regression coefficients

$$\mathbf{Y} = \mathbf{X}\beta^0 + \varepsilon, \quad p \gg n$$

with fixed (deterministic) design \mathbf{X}

problem of identifiability:

for $p > n$: $\mathbf{X}\beta^0 = \mathbf{X}\theta$

for any $\theta = \beta^0 + \xi$, ξ in the null-space of \mathbf{X}

\leadsto cannot say anything about $\|\hat{\beta} - \beta^0\|$ without further assumptions!

\leadsto we will work with the compatibility assumption (see later by Sara)

and Sara will explain: under compatibility condition

$$\|\hat{\beta} - \beta^0\|_1 \leq C \frac{s_0}{\phi_0^2} \sqrt{\log(p)/n},$$

$$s_0 = |\text{supp}(\beta^0)| = |\{j; \beta_j^0 \neq 0\}|$$

Variable selection

Example: Motif regression

for finding HIF1 α transcription factor binding sites in DNA seq.

Müller, Meier, PB & Ricci



$Y_i \in \mathbb{R}$: univariate response measuring binding intensity of HIF1 α on coarse DNA segment i (from CHIP-chip experiments)

$X_i = (X_i^{(1)}, \dots, X_i^{(p)}) \in \mathbb{R}^p$:

$X_i^{(j)}$ = abundance score of candidate motif j in DNA segment i (using sequence data and computational biology algorithms, e.g. MDSCAN)

question: relation between the binding intensity Y and the abundance of short candidate motifs?

~> linear model is often reasonable

“motif regression” (Conlon, X.S. Liu, Lieb & J.S. Liu, 2003)

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon, \quad n = 287, \quad p = 195$$

goal: variable selection

~> find the relevant motifs among the $p = 195$ candidates

Lasso for variable selection

$$\hat{S}(\lambda) = \{j; \hat{\beta}_j(\lambda) \neq 0\}$$

for

$$S_0 = \{j; \beta_j^0 \neq 0\}$$

no significance testing involved
it's convex optimization only!

(and that can be a problem... see later)

Motif regression

for finding HIF1 α transcription factor binding sites in DNA seq.

$Y_i \in \mathbb{R}$: univariate response measuring binding intensity on coarse DNA segment i (from CHIP-chip experiments)

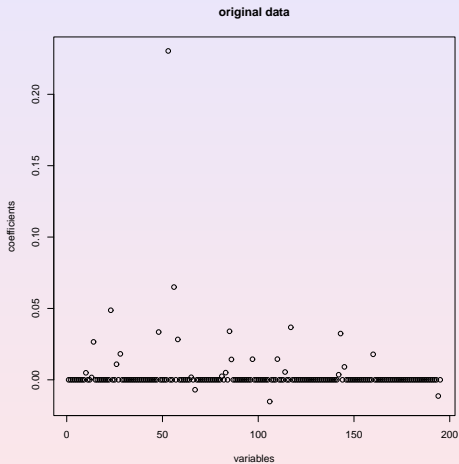
$X_i^{(j)}$ = abundance score of candidate motif j in DNA segment i

variable selection in linear model $Y_i = \beta_0 + \sum_{j=1}^p \beta_j X_i^{(j)} + \varepsilon_i$,

$i = 1, \dots, n = 287, p = 195$

\leadsto Lasso selects 26 covariates and $R^2 \approx 50\%$
i.e. 26 interesting candidate motifs

motif regression: estimated coefficients $\hat{\beta}(\hat{\lambda}_{CV})$



“Theory” for variable selection with Lasso

for (fixed design) linear model $\mathbf{Y} = \mathbf{X}\beta^0 + \varepsilon$ with
active set $S_0 = \{j; \beta_j^0 \neq 0\}$
two key assumptions

1. neighborhood stability condition for design \mathbf{X}
 \Leftrightarrow irrepresentable condition for design \mathbf{X}
2. beta-min condition

$$\min_{j \in S_0} |\beta_j^0| \geq C \sqrt{s_0 \log(p)/n}, \quad C \text{ suitably large}$$

both conditions are **sufficient and “essentially” necessary** for

$$\hat{S}(\lambda) = S_0 \text{ with high probability, } \lambda \gg \underbrace{\sqrt{\log(p)/n}}_{\text{larger than for pred.}}$$

already proved in **Meinshausen & PB, 2004 (publ: 2006)**
and both assumptions are restrictive!

“Theory” for variable selection with Lasso

for (fixed design) linear model $\mathbf{Y} = \mathbf{X}\beta^0 + \varepsilon$ with
active set $S_0 = \{j; \beta_j^0 \neq 0\}$
two key assumptions

1. neighborhood stability condition for design \mathbf{X}
 \Leftrightarrow irrepresentable condition for design \mathbf{X}
2. beta-min condition

$$\min_{j \in S_0} |\beta_j^0| \geq C \sqrt{s_0 \log(p)/n}, \quad C \text{ suitably large}$$

both conditions are **sufficient and “essentially” necessary** for

$$\hat{S}(\lambda) = S_0 \text{ with high probability, } \lambda \gg \underbrace{\sqrt{\log(p)/n}}_{\text{larger than for pred.}}$$

already proved in Meinshausen & PB, 2004 (publ: 2006)
and **both assumptions are restrictive!**

neighborhood stability condition \Leftrightarrow irrepresentable condition

(Zhao & Yu, 2006)

$$n^{-1} \mathbf{X}^T \mathbf{X} = \hat{\Sigma}$$

active set $S_0 = \{j; \beta_j \neq 0\} = \{1, \dots, s_0\}$ consists of the first s_0 variables; partition

$$\hat{\Sigma} = \begin{pmatrix} \hat{\Sigma}_{S_0, S_0} & \hat{\Sigma}_{S_0, S_0^c} \\ \hat{\Sigma}_{S_0^c, S_0} & \hat{\Sigma}_{S_0^c, S_0^c} \end{pmatrix}$$

irrep. condition : $\|\hat{\Sigma}_{S_0^c, S_0} \hat{\Sigma}_{S_0, S_0}^{-1} \text{sign}(\beta_1^0, \dots, \beta_{s_0}^0)^T\|_\infty < 1$

not very realistic assumptions... what can we expect?

recall: under compatibility condition

$$\|\hat{\beta} - \beta^0\|_1 \leq C \frac{s_0}{\phi_0^2} \sqrt{\log(p)/n}$$

consider the relevant active variables

$$S_{\text{relev}} = \{j; |\beta_j^0| > C \frac{s_0}{\phi_0^2} \sqrt{\log(p)/n}\}$$

then, clearly,

$$\hat{S} \supseteq S_{\text{relev}} \text{ with high probability}$$

screening for detecting the relevant variables is possible!

without beta-min condition and assuming compatibility condition only

in addition: assuming beta-min condition

$$\min_{j \in S_0} |\beta_j^0| > C \frac{s_0}{\phi_0^2} \sqrt{\log(p)/n}$$

$\hat{S} \supseteq S_0$ with high probability

screening for detecting the true variables

Tibshirani (1996):

LASSO = Least Absolute Shrinkage and Selection Operator

new translation:

LASSO = Least Absolute Shrinkage and **S**creening Operator

Practical perspective

choice of λ : $\hat{\lambda}_{CV}$ from cross-validation
empirical and theoretical indications (Meinshausen & PB, 2006)
that

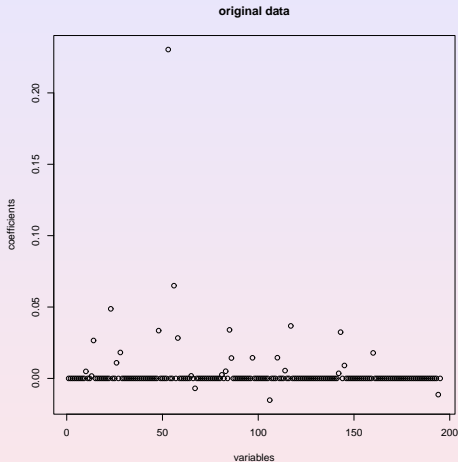
$$\hat{S}(\hat{\lambda}_{CV}) \supseteq S_0 \quad (\text{or } S_{\text{relev}})$$

moreover

$$|\hat{S}(\hat{\lambda}_{CV})| \leq \min(n, p) (= n \text{ if } p \gg n)$$

\leadsto **huge dimensionality reduction** (in the original covariates)

motif regression: estimated coefficients $\hat{\beta}(\hat{\lambda}_{CV})$



which variables in \hat{S} are false positives?
(p-values would be very useful!)

recall:

$$\hat{S}(\hat{\lambda}_{CV}) \supseteq S_0 \quad (\text{or } S_{\text{relev}})$$

and we would then use a second-stage to reduce the number of false positive selections

→ re-estimation on much smaller model with variables from \hat{S}

- ▶ OLS on \hat{S} with e.g. BIC variable selection
- ▶ thresholding coefficients and OLS re-estimation
- ▶ adaptive Lasso (Zou, 2006)
- ▶ ...

recall:

$$\hat{S}(\hat{\lambda}_{CV}) \supseteq S_0 \quad (\text{or } S_{\text{relev}})$$

and we would then use a second-stage to reduce the number of false positive selections

- ↪ re-estimation on much smaller model with variables from \hat{S}
- ▶ OLS on \hat{S} with e.g. BIC variable selection
 - ▶ thresholding coefficients and OLS re-estimation
 - ▶ adaptive Lasso (Zou, 2006)
 - ▶ ...

Adaptive Lasso (Zou, 2006)

re-weighting the penalty function

$$\hat{\beta} = \operatorname{argmin}_{\beta} \left(\|\mathbf{Y} - \mathbf{X}\beta\|_2^2/n + \lambda \sum_{j=1}^p \frac{|\beta_j|}{|\hat{\beta}_{init,j}|} \right),$$

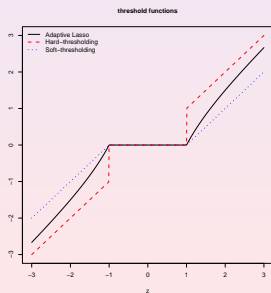
$\hat{\beta}_{init,j}$ from Lasso in first stage (or OLS if $p < n$)
Zou (2006)

for orthogonal design,

if $\hat{\beta}_{init} = \text{OLS}$:

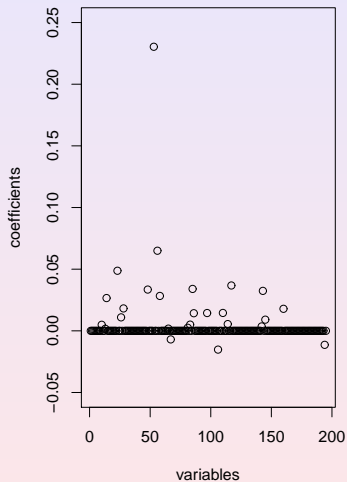
Adaptive Lasso = NN-garrote

\rightsquigarrow less bias than Lasso

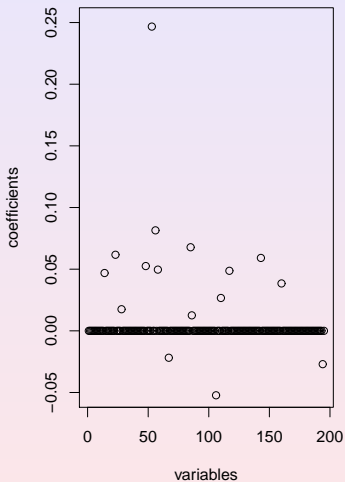


motif regression

Lasso



Adaptive Lasso



Lasso selects 26 variables Adaptive Lasso selects 16 variables

KKT conditions and Computation

characterization of solution(s) $\hat{\beta}$ as minimizer of the criterion function

$$Q_\lambda(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|_2^2/n + \lambda\|\beta\|_1$$

since $Q_\lambda(\cdot)$ is a convex function:

necessary and sufficient that subdifferential of $\partial Q_\lambda(\beta)/\partial\beta$ at $\hat{\beta}$ contains the zero element

Lemma 2.1 first part (in the book)

denote by $G(\beta) = -2\mathbf{X}^T(\mathbf{Y} - \mathbf{X}\beta)/n$ the gradient vector of $\|\mathbf{Y} - \mathbf{X}\beta\|_2^2/n$

Then: $\hat{\beta}$ is a solution if and only if

$$\begin{aligned} G_j(\hat{\beta}) &= -\text{sign}(\hat{\beta}_j)\lambda \text{ if } \hat{\beta}_j \neq 0, \\ |G_j(\hat{\beta})| &\leq \lambda \text{ if } \hat{\beta}_j = 0 \end{aligned}$$

Lemma 2.1 second part (in the book)

If the solution of $\operatorname{argmin}_{\beta} Q_{\lambda}(\beta)$ is not unique (e.g. if $p > n$), and if $G_j(\hat{\beta}) < \lambda$ for some solution $\hat{\beta}$, then $\hat{\beta}_j = 0$ for all (other) solutions $\hat{\beta}$ in $\operatorname{argmin}_{\beta} Q_{\lambda}(\beta)$.

The zeroes are “essentially” unique

(“essentially” refers to the situation: $\hat{\beta}_j = 0$ and $G_j(\hat{\beta}) = \lambda$)

Proof: Exercise (optional), or see in the book

Coordinate descent algorithm for computation

general idea is to compute a solution $\hat{\beta}(\lambda_{\text{grid},k})$ and use it as a starting value for the computation of $\hat{\beta}(\lambda_{\text{grid},k-1})$

$\underbrace{\hspace{10em}}_{< \lambda_{\text{grid},k}}$

$\beta^{(0)} \in \mathbb{R}^p$ an initial parameter vector. Set $m = 0$.

REPEAT:

Increase m by one: $m \leftarrow m + 1$.

For $j = 1, \dots, p$:

if $|G_j(\beta_{-j}^{(m-1)})| \leq \lambda$: set $\beta_j^{(m)} = 0$,

otherwise: $\beta_j^{(m)} = \operatorname{argmin}_{\beta_j} Q_\lambda(\beta_{+j}^{(m-1)})$,

β_{-j} : parameter vector setting j th component to zero

$\beta_{+j}^{(m-1)}$: parameter vector which equals $\beta^{(m-1)}$ except for j th component equalling β_j

UNTIL numerical convergence

for squared error loss: explicit up-dating formulae (Exercise 2.7)

$$\begin{aligned}G_j(\beta) &= -2\mathbf{X}_j^T(\mathbf{Y} - \mathbf{X}\beta) \\ \beta_j^{(m)} &= \frac{\text{sign}(Z_j)(|Z_j| - \lambda/2)_+}{\hat{\Sigma}_{jj}}, \\ Z_j &= \mathbf{X}_j^T(\mathbf{Y} - \mathbf{X}\beta_{-j})/n, \quad \hat{\Sigma} = n^{-1}\mathbf{X}^T\mathbf{X}.\end{aligned}$$

→ componentwise soft-thresholding

this is very fast if true problem is sparse

active set strategy: can do non-systematic cycling, visiting mainly the active (non-zero) components

riboflavin example, $n=71$, $p=4088$

0.33 secs. CPU using `glmnet`-package in R

(Friedman, Hastie & Tibshirani, 2008)

coordinate descent algorithm converges to a stationary point
(Paul Tseng \approx 2000)

\leadsto convergence to a global optimum, due to convexity of the problem

main assumption:

objective function = smooth function + penalty
separable

here: “separable” means “additive”, i.e., $\text{pen}(\beta) = \sum_{j=1}^p p_j(\beta_j)$

failure of coordinate descent algorithm:

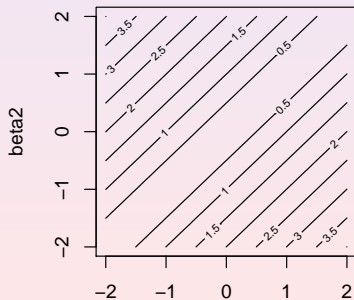
Fused Lasso

$$\hat{\beta} = \operatorname{argmin}_{\beta} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2/n + \lambda \sum_{j=2}^p |\beta_j - \beta_{j-1}| + \lambda_2 \|\beta\|_1$$

but $\sum_{j=2}^p |\beta_j - \beta_{j-1}|$ is non-separable

contour lines of penalties for $p = 2$

$|\beta_1 - \beta_2|$



$|\beta_1| + |\beta_2|$

