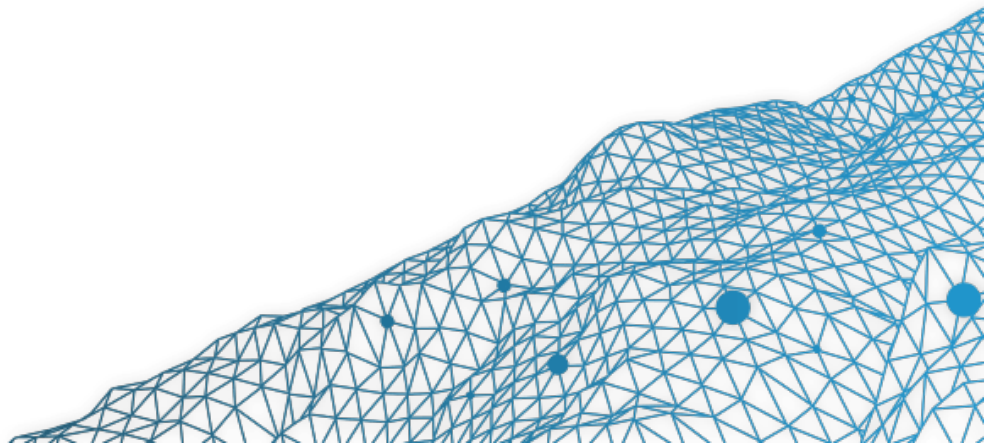


Advanced Econometrics #3: Model & Variable Selection

A. Charpentier (Université de Rennes 1)

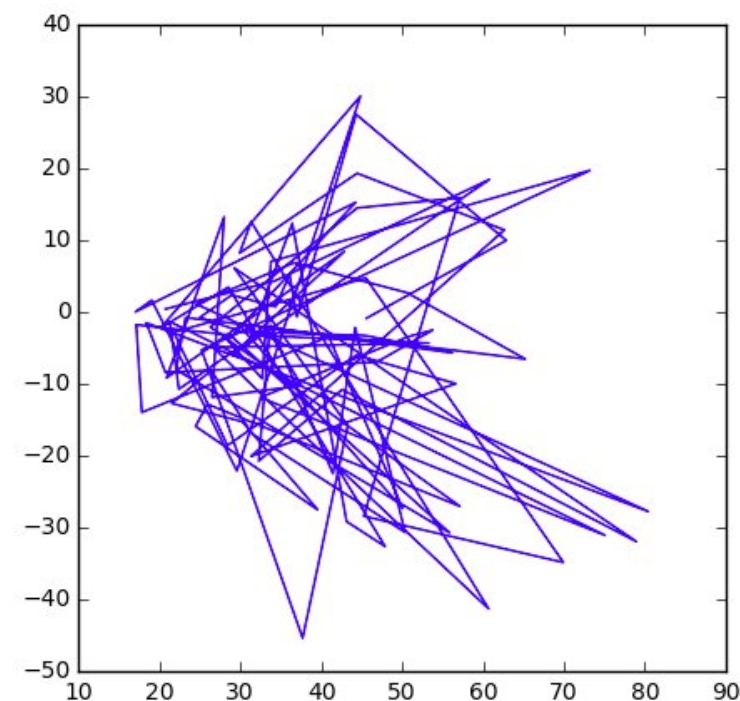
Université de Rennes 1,
Graduate Course, 2017.



“Great plot.

Now need to find the theory that explains it”

Deville (2017) <http://twitter.com>



Preliminary Results: Numerical Optimization

Problem : $\mathbf{x}^* \in \operatorname{argmin}\{f(\mathbf{x}); \mathbf{x} \in \mathbb{R}^d\}$

Gradient descent : $\mathbf{x}_{k+1} = \mathbf{x}_k - \eta \nabla f(\mathbf{x}_k)$ starting from some \mathbf{x}_0

Problem : $\mathbf{x}^* \in \operatorname{argmin}\{f(\mathbf{x}); \mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d\}$

Projected descent : $\mathbf{x}_{k+1} = \Pi_{\mathcal{X}}(\mathbf{x}_k - \eta \nabla f(\mathbf{x}_k))$ starting from some \mathbf{x}_0

A constrained problem is said to be **convex** if

$$\left\{ \begin{array}{ll} \min\{f(\mathbf{x})\} & \text{with } f \text{ convex} \\ \text{s.t. } g_i(\mathbf{x}) = 0, \forall i = 1, \dots, n & \text{with } g_i \text{ linear} \\ h_i(\mathbf{x}) \leq 0, \forall i = 1, \dots, m & \text{with } h_i \text{ convex} \end{array} \right.$$

Lagrangian : $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^n \lambda_i g_i(\mathbf{x}) + \sum_{i=1}^m \mu_i h_i(\mathbf{x})$ where \mathbf{x} are primal variables and $(\boldsymbol{\lambda}, \boldsymbol{\mu})$ are dual variables.

Remark \mathcal{L} is an affine function in $(\boldsymbol{\lambda}, \boldsymbol{\mu})$

Preliminary Results: Numerical Optimization

Karush–Kuhn–Tucker conditions : a convex problem has a solution \mathbf{x}^* if and only if there are $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ such that the following conditions hold

- stationarity : $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{0}$ at $(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$
- primal admissibility : $g_i(\mathbf{x}^*) = 0$ and $h_i(\mathbf{x}^*) \leq 0, \forall i$
- dual admissibility : $\boldsymbol{\mu}^* \geq 0$

Let L denote the associated dual function $L(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \min_{\mathbf{x}} \{\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})\}$

L is a convex function in $(\boldsymbol{\lambda}, \boldsymbol{\mu})$ and the **dual problem** is $\max_{\boldsymbol{\lambda}, \boldsymbol{\mu}} \{L(\boldsymbol{\lambda}, \boldsymbol{\mu})\}$.

References

Motivation

Banerjee, A., Chandrasekhar, A.G., Duflo, E. & Jackson, M.O. (2016). **Gossip: Identifying Central Individuals in a Social Networks.**

References

Belloni, A. & Chernozhukov, V. 2009. **Least squares after model selection in high-dimensional sparse models.**

Hastie, T., Tibshirani, R. & Wainwright, M. 2015 **Statistical Learning with Sparsity: The Lasso and Generalizations.** CRC Press.

Preamble

Assume that $y = m(\mathbf{x}) + \varepsilon$, where ε is some idiosyncatic unpredictable noise.

The error $\mathbb{E}[(y - m(\mathbf{x}))^2]$ is the sum of three terms

- variance of the estimator : $\mathbb{E}[(y - \hat{m}(\mathbf{x}))^2]$
- bias² of the estimator : $[m(\mathbf{x}) - \hat{m}(\mathbf{x})]^2$
- variance of the noise : $\mathbb{E}[(y - m(\mathbf{x}))^2]$

(the latter exists, even with a ‘perfect’ model).

Preamble

Consider a parametric model, with true (unkown) parameter θ , then

$$\text{mse}(\hat{\theta}) = \mathbb{E} \left[(\hat{\theta} - \theta)^2 \right] = \underbrace{\mathbb{E} \left[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2 \right]}_{\text{variance}} + \underbrace{\mathbb{E} \left[(\mathbb{E}[\hat{\theta}] - \theta)^2 \right]}_{\text{bias}^2}$$

Let $\tilde{\theta}$ denote an unbiased estimator of θ . Then

$$\hat{\theta} = \frac{\theta^2}{\theta^2 + \text{mse}(\tilde{\theta})} \cdot \tilde{\theta} = \tilde{\theta} - \underbrace{\frac{\text{mse}(\tilde{\theta})}{\theta^2 + \text{mse}(\tilde{\theta})}}_{\text{penalty}} \cdot \tilde{\theta}$$

satisfies $\text{mse}(\hat{\theta}) \leq \text{mse}(\tilde{\theta})$.

Occam's Razor

The “law of parsimony”, “*lex parsimoniae*”

CORE PRINCIPLES IN RESEARCH



OCCAM'S RAZOR

"WHEN FACED WITH TWO POSSIBLE EXPLANATIONS, THE SIMPLER OF THE TWO IS THE ONE MOST LIKELY TO BE TRUE."



OCCAM'S PROFESSOR

"WHEN FACED WITH TWO POSSIBLE WAYS OF DOING SOMETHING, THE MORE COMPLICATED ONE IS THE ONE YOUR PROFESSOR WILL MOST LIKELY ASK YOU TO DO."

WWW.PHDCOMICS.COM

JORGE CHAM © 2009

Penalize too complex models

James & Stein Estimator

Let $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbb{I})$. We want to estimate $\boldsymbol{\mu}$.

$$\hat{\boldsymbol{\mu}}_{\text{mle}} = \overline{\mathbf{X}}_n \sim \mathcal{N}\left(\boldsymbol{\mu}, \frac{\sigma^2}{n} \mathbb{I}\right).$$

From James & Stein (1961) **Estimation with quadratic loss**

$$\hat{\boldsymbol{\mu}}_{\text{JS}} = \left(1 - \frac{(d-2)\sigma^2}{n\|\overline{\mathbf{y}}\|^2}\right) \overline{\mathbf{y}}$$

where $\|\cdot\|$ is the Euclidean norm.

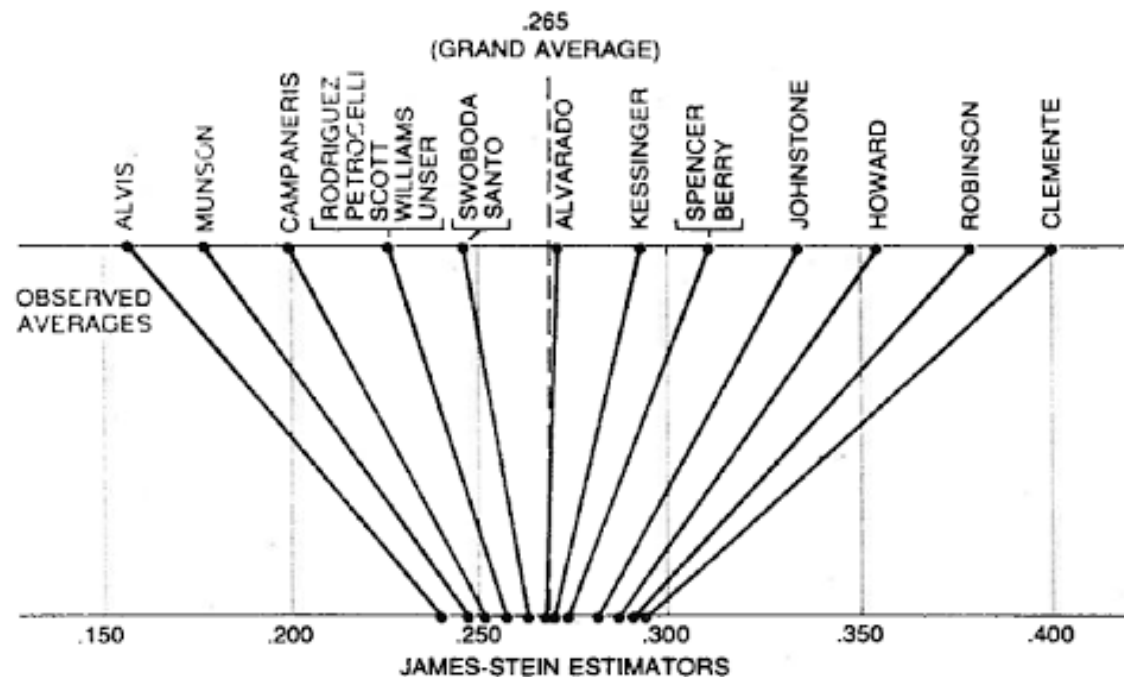
One can prove that if $d \geq 3$,

$$\mathbb{E}[(\hat{\boldsymbol{\mu}}_{\text{JS}} - \boldsymbol{\mu})^2] < \mathbb{E}[(\hat{\boldsymbol{\mu}}_{\text{mle}} - \boldsymbol{\mu})^2]$$

Samworth (2015) **Stein's paradox**, “*one should use the price of tea in China to obtain a better estimate of the chance of rain in Melbourne*”.

James & Stein Estimator

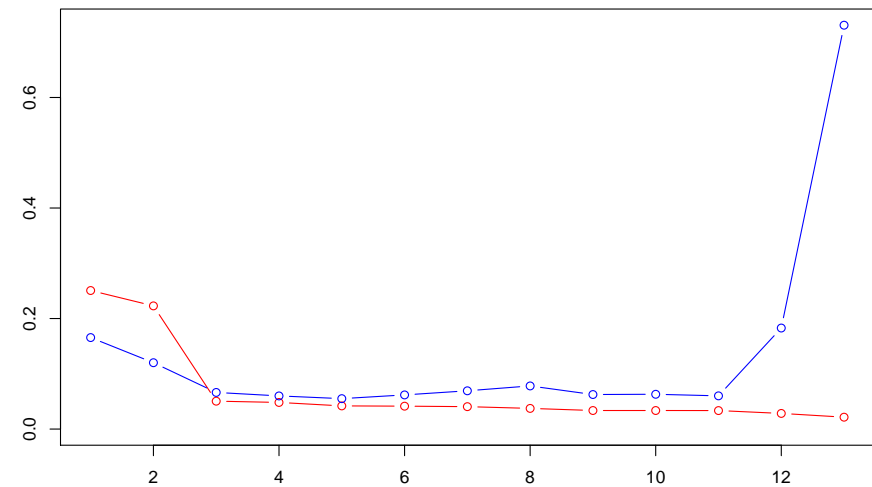
Heuristics : consider a biased estimator, to decrease the variance.



See Efron (2010) [Large-Scale Inference](#)

Motivation: Avoiding Overfit

Generalization : the model should perform well on **new** data (and not only on the training ones).



Reducing Dimension with PCA

Use **principal components** to reduce dimension (on centered and scaled variables):
we want d vectors $\mathbf{z}_1, \dots, \mathbf{z}_d$ such that

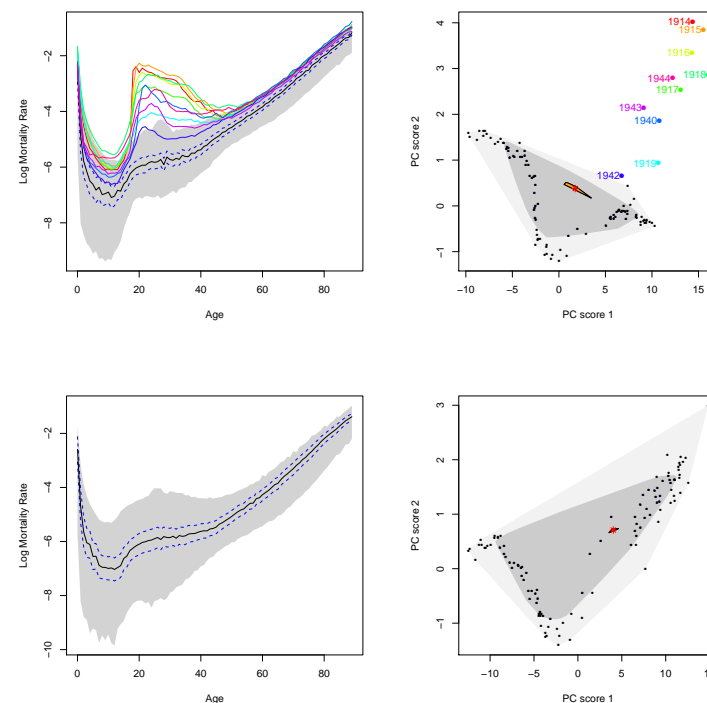
First Component is $\mathbf{z}_1 = \mathbf{X}\boldsymbol{\omega}_1$ where

$$\boldsymbol{\omega}_1 = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \left\{ \|\mathbf{X} \cdot \boldsymbol{\omega}\|^2 \right\} = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \left\{ \boldsymbol{\omega}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\omega} \right\}$$

Second Component is $\mathbf{z}_2 = \mathbf{X}\boldsymbol{\omega}_2$ where

$$\boldsymbol{\omega}_2 = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \left\{ \|\widetilde{\mathbf{X}}^{(1)} \cdot \boldsymbol{\omega}\|^2 \right\}$$

$$\text{with } \widetilde{\mathbf{X}}^{(1)} = \mathbf{X} - \underbrace{\mathbf{X}\boldsymbol{\omega}_1\boldsymbol{\omega}_1^\top}_{\mathbf{z}_1}.$$



Reducing Dimension with PCA

A regression on (the d) principal components, $y = \mathbf{z}^\top \boldsymbol{\beta} + \boldsymbol{\eta}$ could be an interesting idea, unfortunately, principal components have no reason to be correlated with y . First component was $\mathbf{z}_1 = \mathbf{X}\boldsymbol{\omega}_1$ where

$$\boldsymbol{\omega}_1 = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \left\{ \|\mathbf{X} \cdot \boldsymbol{\omega}\|^2 \right\} = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \left\{ \boldsymbol{\omega}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\omega} \right\}$$

It is a non-supervised technique.

Instead, use **partial least squares**, introduced in Wold (1966) **Estimation of Principal Components and Related Models by Iterative Least squares**. First component is $\mathbf{z}_1 = \mathbf{X}\boldsymbol{\omega}_1$ where

$$\boldsymbol{\omega}_1 = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \left\{ \langle \mathbf{y}, \mathbf{X} \cdot \boldsymbol{\omega} \rangle \right\} = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \left\{ \boldsymbol{\omega}^\top \mathbf{X}^\top \mathbf{y} \mathbf{y}^\top \mathbf{X} \boldsymbol{\omega} \right\}$$

Terminology

Consider a dataset $\{y_i, \mathbf{x}_i\}$, assumed to be generated from Y, \mathbf{X} , from an unknown distribution \mathbb{P} .

Let $m_0(\cdot)$ be the “true” model. Assume that $y_i = m_0(\mathbf{x}_i) + \varepsilon_i$.

In a regression context (**quadratic loss function** function), the **risk** associated to m is

$$\mathcal{R}(m) = \mathbb{E}_{\mathbb{P}}[(Y - m(\mathbf{X}))^2]$$

An optimal model m^* within a class \mathcal{M} satisfies

$$\mathcal{R}(m^*) = \inf_{m \in \mathcal{M}} \{\mathcal{R}(m)\}$$

Such a model m^* is usually called **oracle**.

Observe that $m^*(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$ is the solution of

$$\mathcal{R}(m^*) = \inf_{m \in \mathcal{M}} \{\mathcal{R}(m)\} \text{ where } \mathcal{M} \text{ is the set of measurable functions}$$

The **empirical risk** is

$$\mathcal{R}_n(m) = \frac{1}{n} \sum_{i=1}^n (y_i - m(\mathbf{x}_i))^2$$

For instance, m can be a linear predictor, $m(\mathbf{x}) = \beta_0 + \mathbf{x}^\top \boldsymbol{\beta}$, where $\boldsymbol{\theta} = (\beta_0, \boldsymbol{\beta})$ should be estimated (trained).

$\mathbb{E}[R_n(\hat{m})] = \mathbb{E}[(\hat{m}(\mathbf{X}) - Y)^2]$ can be expressed as

$$\begin{aligned} & \mathbb{E}[(\hat{m}(\mathbf{X}) - \mathbb{E}[\hat{m}(\mathbf{X})|\mathbf{X}])^2] \quad \text{variance of } \hat{m} \\ + & \mathbb{E}[(\mathbb{E}[\hat{m}(\mathbf{X})|\mathbf{X}] - \underbrace{\mathbb{E}[Y|\mathbf{X}]}_{m_0(\mathbf{X})})^2] \quad \text{bias of } \hat{m} \\ + & \mathbb{E}[(Y - \underbrace{\mathbb{E}[Y|\mathbf{X}]}_{m_0(\mathbf{X})})^2] \quad \text{variance of the noise} \end{aligned}$$

The third term is the risk of the “optimal” estimator m , that cannot be decreased.

Mallows Penalty and Model Complexity

Consider a linear predictor (see #1), i.e. $\hat{\mathbf{y}} = \hat{m}(\mathbf{x}) = \mathbf{A}\mathbf{y}$.

Assume that $\mathbf{y} = m_0(\mathbf{x}) + \boldsymbol{\varepsilon}$, with $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$ and $\text{Var}[\boldsymbol{\varepsilon}] = \sigma^2 \mathbb{I}$.

Let $\|\cdot\|$ denote the Euclidean norm

Empirical risk : $\hat{\mathcal{R}}_n(m) = \frac{1}{n} \|\mathbf{y} - m(\mathbf{x})\|^2$

Vapnik's risk : $\mathbb{E}[\hat{\mathcal{R}}_n(m)] = \frac{1}{n} \|m_0(\mathbf{x}) - m(\mathbf{x})\|^2 + \frac{1}{n} \mathbb{E}(\|\mathbf{y} - m_0(\mathbf{x})\|^2)$ with $m_0(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$.

Observe that

$$n\mathbb{E}[\hat{\mathcal{R}}_n(\hat{m})] = \mathbb{E}(\|\mathbf{y} - \hat{m}(\mathbf{x})\|^2) = \|(\mathbb{I} - \mathbf{A})m_0\|^2 + \sigma^2 \|\mathbb{I} - \mathbf{A}\|^2$$

while

$$= \mathbb{E}(\|m_0(\mathbf{x}) - \hat{m}(\mathbf{x})\|^2) = \underbrace{\|(\mathbb{I} - \mathbf{A})m_0\|^2}_{\text{bias}} + \underbrace{\sigma^2 \|\mathbf{A}\|^2}_{\text{variance}}$$

Mallows Penalty and Model Complexity

One can obtain

$$\mathbb{E}[\mathcal{R}_n(\hat{m})] = \mathbb{E}[\hat{\mathcal{R}}_n(\hat{m})] + 2\frac{\sigma^2}{n}\text{trace}(\mathbf{A}).$$

If $\text{trace}(\mathbf{A}) \geq 0$ the empirical risk underestimate the true risk of the estimator.

The number of degrees of freedom of the (linear) predictor is related to $\text{trace}(\mathbf{A})$

$2\frac{\sigma^2}{n}\text{trace}(\mathbf{A})$ is called Mallows's penalty C_L .

If \mathbf{A} is a projection matrix, $\text{trace}(\mathbf{A})$ is the dimension of the projection space, p , then we obtain Mallows's C_P , $2\frac{\sigma^2}{n}p$.

Remark : Mallows (1973) **Some Comments on C_p** introduced this penalty while focusing on the R^2 .

Penalty and Likelihood

C_P is associated to a quadratic risk

an alternative is to use a distance on the (conditional) distribution of Y , namely

Kullback-Leibler distance

discrete case:
$$D_{\text{KL}}(P\|Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

continuous case :

$$D_{\text{KL}}(P\|Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx \quad D_{\text{KL}}(P\|Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

Let f denote the true (unknown) density, and f_θ some parametric distribution,

$$D_{\text{KL}}(f\|f_\theta) = \int_{-\infty}^{\infty} f(x) \log \frac{f(x)}{f_\theta(x)} dx = \int f(x) \log[f(x)] dx - \underbrace{\int f(x) \log[f_\theta(x)] dx}_{\text{relative information}}$$

Hence

$$\text{minimize } \{D_{\text{KL}}(f\|f_\theta)\} \iff \text{maximize } \{\mathbb{E}[\log[f_\theta(X)]]\}$$

Penalty and Likelihood

Akaike (1974) **A new look at the statistical model identification** observe that for n large enough

$$\mathbb{E}[\log[f_{\theta}(X)]] \sim \log[\mathcal{L}(\hat{\theta})] - \dim(\theta)$$

Thus

$$AIC = -2 \log \mathcal{L}(\hat{\theta}) + 2\dim(\theta)$$

Example : in a (Gaussian) linear model, $y_i = \beta_0 + \mathbf{x}_i^{\top} \boldsymbol{\beta} + \varepsilon_i$

$$AIC = n \log \left(\frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i \right) + 2[\dim(\boldsymbol{\beta}) + 2]$$

Penalty and Likelihood

Remark : this is valid for large sample (rule of thumb $n/\dim(\theta) > 40$), otherwise use a corrected AIC

$$AIC_c = AIC + \underbrace{\frac{2k(k+1)}{n-k-1}}_{\text{bias correction}} \quad \text{where } k = \dim(\theta)$$

see Sugiura (1978) **Further analysis of the data by Akaike's information criterion and the finite corrections** second order AIC.

Using a Bayesian interpretation, Schwarz (1978) **Estimating the dimension of a model** obtained

$$BIC = -2 \log \mathcal{L}(\hat{\theta}) + \log(n) \dim(\theta).$$

Observe that the criteria considered is

$$\text{criteria} = -\text{function}(\mathcal{L}(\hat{\theta})) + \text{penalty}(\text{complexity})$$

Estimation of the Risk

Consider a naive bootstrap procedure, based on a bootstrap sample

$$\mathcal{S}_b = \{(y_i^{(b)}, \mathbf{x}_i^{(b)})\}.$$

The **plug-in** estimator of the empirical risk is

$$\hat{\mathcal{R}}_n(\hat{m}^{(b)}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{m}^{(b)}(\mathbf{x}_i))^2$$

and then

$$\hat{\mathcal{R}}_n = \frac{1}{B} \sum_{b=1}^B \hat{\mathcal{R}}_n(\hat{m}^{(b)}) = \frac{1}{B} \sum_{b=1}^B \frac{1}{n} \sum_{i=1}^n (y_i - \hat{m}^{(b)}(\mathbf{x}_i))^2$$

Estimation of the Risk

One might improve this estimate using a **out-of-bag** procedure

$$\hat{\mathcal{R}}_n = \frac{1}{n} \sum_{i=1}^n \frac{1}{\#\mathcal{B}_i} \sum_{b \in \mathcal{B}_i} (y_i - \hat{m}^{(b)}(\mathbf{x}_i))^2$$

where \mathcal{B}_i is the set of all bootstrap sample that contain (y_i, \mathbf{x}_i) .

Remark: $\mathbb{P}((y_i, \mathbf{x}_i) \notin \mathcal{S}_b) = \left(1 - \frac{1}{n}\right)^n \sim e^{-1} = 36,78\%$.

Linear Regression Shortcoming

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

Unbiased Estimator $\mathbb{E}[\hat{\beta}] = \beta$

Variance $\text{Var}[\hat{\beta}] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$

which can be (extremely) large when $\det[(\mathbf{X}^\top \mathbf{X})] \sim 0$.

$$\mathbf{X} = \begin{bmatrix} 1 & -1 & 2 \\ 1 & 0 & 1 \\ 1 & 2 & -1 \\ 1 & 1 & 0 \end{bmatrix} \quad \text{then } \mathbf{X}^\top \mathbf{X} = \begin{bmatrix} 4 & 2 & 2 \\ 2 & 6 & -4 \\ 2 & -4 & 6 \end{bmatrix} \quad \text{while } \mathbf{X}^\top \mathbf{X} + \mathbb{I} = \begin{bmatrix} 5 & 2 & 2 \\ 2 & 7 & -4 \\ 2 & -4 & 7 \end{bmatrix}$$

eigenvalues : $\{10, 6, 0\}$

$\{11, 7, 1\}$

Ad-hoc strategy: use $\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}$

Linear Regression Shortcoming

Evolution of $(\beta_1, \beta_2) \mapsto \sum_{i=1}^n [y_i - (\beta_1 x_{1,i} + \beta_2 x_{2,i})]^2$

when $\text{cor}(X_1, X_2) = r \in [0, 1]$, on top.

Below, **Ridge regression**

$(\beta_1, \beta_2) \mapsto \sum_{i=1}^n [y_i - (\beta_1 x_{1,i} + \beta_2 x_{2,i})]^2 + \lambda(\beta_1^2 + \beta_2^2)$

where $\lambda \in [0, \infty)$, below,

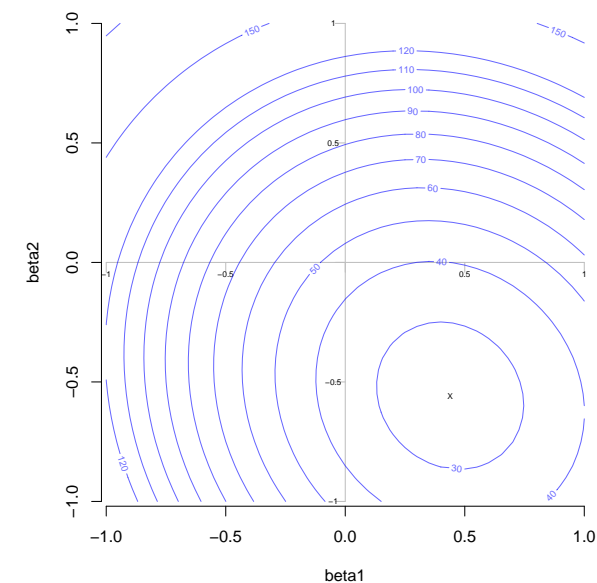
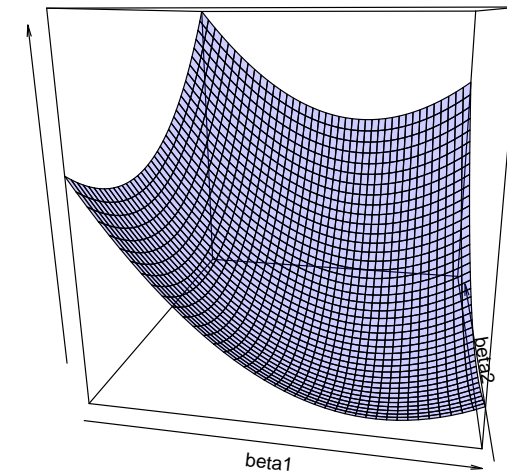
when $\text{cor}(X_1, X_2) \sim 1$ (colinearity).

Normalization : Euclidean ℓ_2 vs. Mahalanobis

We want to penalize complicated models :
if β_k is “too small”, we prefer to have $\beta_k = 0$.



Instead of $d(\mathbf{x}, \mathbf{y}) = (\mathbf{x} - \mathbf{y})^\top (\mathbf{x} - \mathbf{y})$
use $d_{\Sigma}(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})^\top \Sigma^{-1} (\mathbf{x} - \mathbf{y})}$



Ridge Regression

... like the least square, but it shrinks estimated coefficients towards 0.

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \operatorname{argmin} \left\{ \sum_{i=1}^n (y_i - \mathbf{x}_i^{\top} \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \operatorname{argmin} \left\{ \underbrace{\|\mathbf{y} - \mathbf{X}\beta\|_{\ell_2}^2}_{=\text{criteria}} + \underbrace{\lambda \|\beta\|_{\ell_2}^2}_{=\text{penalty}} \right\}$$

$\lambda \geq 0$ is a tuning parameter.

The constant is usually unpenalized. The *true* equation is

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \operatorname{argmin} \left\{ \underbrace{\|\mathbf{y} - (\beta_0 + \mathbf{X}\beta)\|_{\ell_2}^2}_{=\text{criteria}} + \underbrace{\lambda \|\beta\|_{\ell_2}^2}_{=\text{penalty}} \right\}$$

Ridge Regression

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \operatorname{argmin} \left\{ \|\mathbf{y} - (\beta_0 + \mathbf{X}\beta)\|_{\ell_2}^2 + \lambda \|\beta\|_{\ell_2}^2 \right\}$$

can be seen as a constrained optimization problem

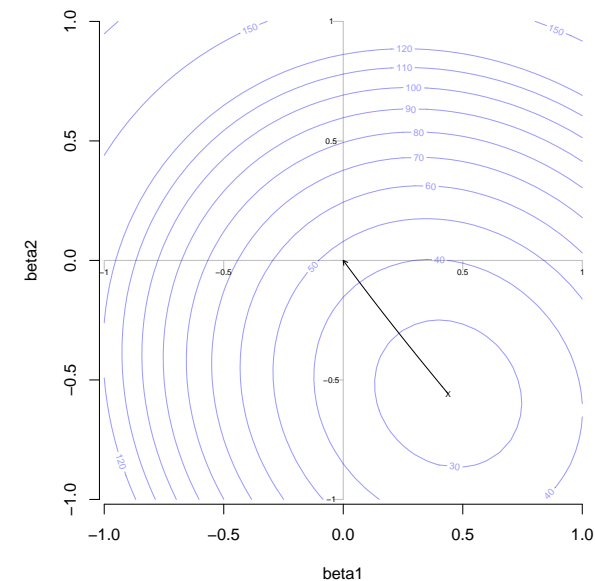
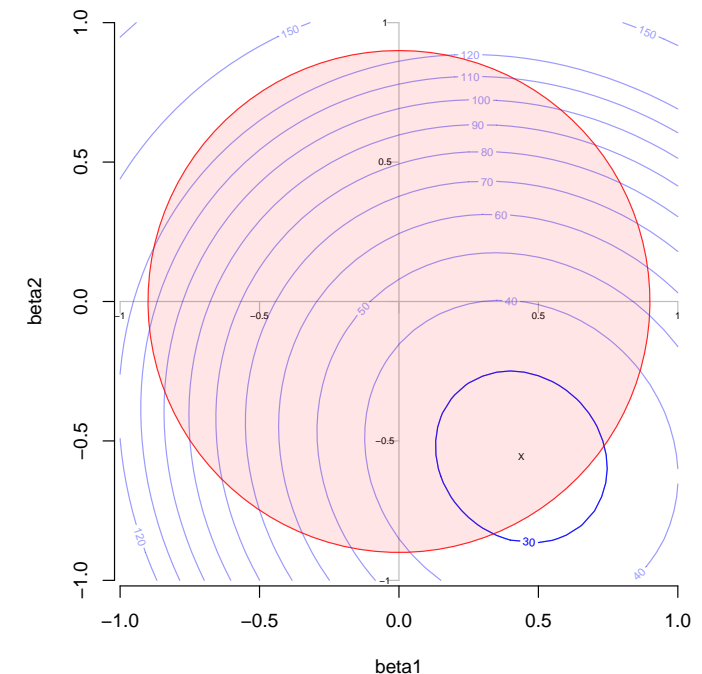
$$\hat{\beta}_{\lambda}^{\text{ridge}} = \operatorname{argmin}_{\|\beta\|_{\ell_2}^2 \leq h_{\lambda}} \left\{ \|\mathbf{y} - (\beta_0 + \mathbf{X}\beta)\|_{\ell_2}^2 \right\}$$

Explicit solution

$$\hat{\beta}_{\lambda} = (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbb{I})^{-1} \mathbf{X}^{\top} \mathbf{y}$$

If $\lambda \rightarrow 0$, $\hat{\beta}_0^{\text{ridge}} = \hat{\beta}^{\text{ols}}$

If $\lambda \rightarrow \infty$, $\hat{\beta}_{\infty}^{\text{ridge}} = \mathbf{0}$.



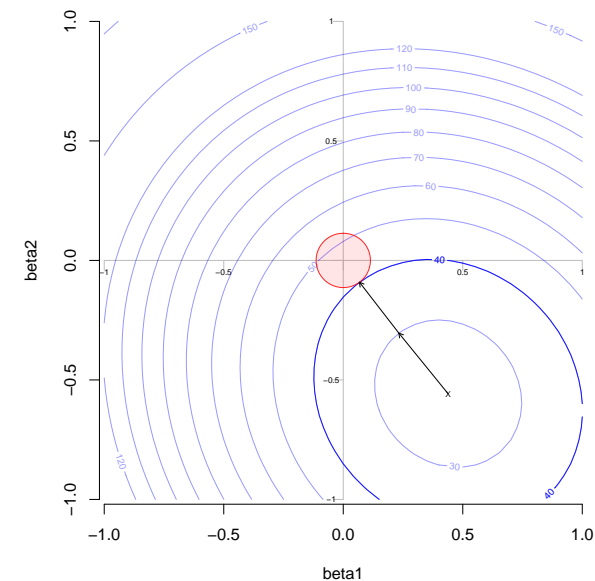
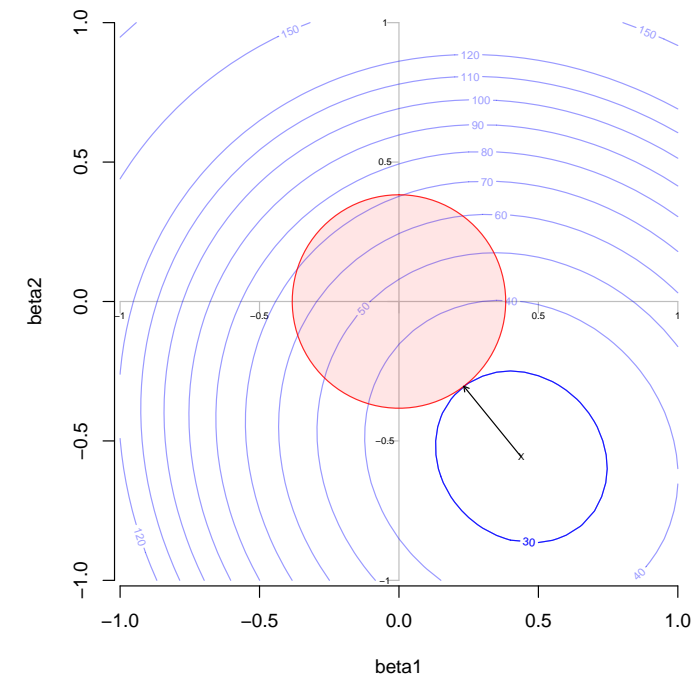
Ridge Regression

This penalty can be seen as rather unfair if components of \mathbf{x} are not expressed on the same scale

- center: $\bar{\mathbf{x}}_j = 0$, then $\hat{\beta}_0 = \bar{y}$
- scale: $\mathbf{x}_j^\top \mathbf{x}_j = 1$

Then compute

$$\hat{\beta}_\lambda^{\text{ridge}} = \operatorname{argmin} \left\{ \underbrace{\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_{\ell_2}^2}_{=\text{loss}} + \underbrace{\lambda \|\boldsymbol{\beta}\|_{\ell_2}^2}_{=\text{penalty}} \right\}$$



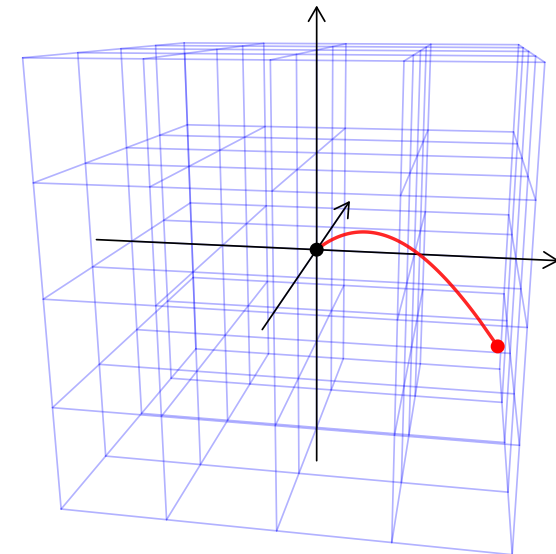
Ridge Regression

Observe that if $\mathbf{x}_{j_1} \perp \mathbf{x}_{j_2}$, then

$$\hat{\beta}_{\lambda}^{\text{ridge}} = [1 + \lambda]^{-1} \hat{\beta}_{\lambda}^{\text{ols}}$$

which explain relationship with shrinkage.

But generally, it is not the case...



Theorem There exists λ such that $\text{mse}[\hat{\beta}_{\lambda}^{\text{ridge}}] \leq \text{mse}[\hat{\beta}_{\lambda}^{\text{ols}}]$

Ridge Regression

$$\mathcal{L}_\lambda(\beta) = \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

$$\frac{\partial \mathcal{L}_\lambda(\beta)}{\partial \beta} = -2\mathbf{X}^\top \mathbf{y} + 2(\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I})\beta$$

$$\frac{\partial^2 \mathcal{L}_\lambda(\beta)}{\partial \beta \partial \beta^\top} = 2(\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I})$$

where $\mathbf{X}^\top \mathbf{X}$ is a semi-positive definite matrix, and $\lambda \mathbb{I}$ is a positive definite matrix, and

$$\hat{\beta}_\lambda = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

The Bayesian Interpretation

From a Bayesian perspective,

$$\underbrace{\mathbb{P}[\boldsymbol{\theta}|\mathbf{y}]}_{\text{posterior}} \propto \underbrace{\mathbb{P}[\mathbf{y}|\boldsymbol{\theta}]}_{\text{likelihood}} \cdot \underbrace{\mathbb{P}[\boldsymbol{\theta}]}_{\text{prior}} \quad \text{i.e.} \quad \log \mathbb{P}[\boldsymbol{\theta}|\mathbf{y}] = \underbrace{\log \mathbb{P}[\mathbf{y}|\boldsymbol{\theta}]}_{\text{log likelihood}} + \underbrace{\log \mathbb{P}[\boldsymbol{\theta}]}_{\text{penalty}}$$

If $\boldsymbol{\beta}$ has a prior $\mathcal{N}(\mathbf{0}, \tau^2 \mathbb{I})$ distribution, then its posterior distribution has mean

$$\mathbb{E}[\boldsymbol{\beta}|\mathbf{y}, \mathbf{X}] = \left(\mathbf{X}^\top \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbb{I} \right)^{-1} \mathbf{X}^\top \mathbf{y}.$$

Properties of the Ridge Estimator

$$\hat{\beta}_{\lambda} = (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbb{I})^{-1} \mathbf{X}^{\top} \mathbf{y}$$

$$\mathbb{E}[\hat{\beta}_{\lambda}] = \mathbf{X}^{\top} \mathbf{X} (\lambda \mathbb{I} + \mathbf{X}^{\top} \mathbf{X})^{-1} \beta.$$

i.e. $\mathbb{E}[\hat{\beta}_{\lambda}] \neq \beta$.

Observe that $\mathbb{E}[\hat{\beta}_{\lambda}] \rightarrow \mathbf{0}$ as $\lambda \rightarrow \infty$.

Assume that \mathbf{X} is an orthogonal design matrix, i.e. $\mathbf{X}^{\top} \mathbf{X} = \mathbb{I}$, then

$$\hat{\beta}_{\lambda} = (1 + \lambda)^{-1} \hat{\beta}^{\text{ols}}.$$

Properties of the Ridge Estimator

Set $\mathbf{W}_\lambda = (\mathbb{I} + \lambda[\mathbf{X}^\top \mathbf{X}]^{-1})^{-1}$. One can prove that

$$\mathbf{W}_\lambda \hat{\boldsymbol{\beta}}^{\text{ols}} = \hat{\boldsymbol{\beta}}_\lambda.$$

Thus,

$$\text{Var}[\hat{\boldsymbol{\beta}}_\lambda] = \mathbf{W}_\lambda \text{Var}[\hat{\boldsymbol{\beta}}^{\text{ols}}] \mathbf{W}_\lambda^\top$$

and

$$\text{Var}[\hat{\boldsymbol{\beta}}_\lambda] = \sigma^2 (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I})^{-1} \mathbf{X}^\top \mathbf{X} [(\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I})^{-1}]^\top.$$

Observe that

$$\text{Var}[\hat{\boldsymbol{\beta}}^{\text{ols}}] - \text{Var}[\hat{\boldsymbol{\beta}}_\lambda] = \sigma^2 \mathbf{W}_\lambda [2\lambda(\mathbf{X}^\top \mathbf{X})^{-2} + \lambda^2(\mathbf{X}^\top \mathbf{X})^{-3}] \mathbf{W}_\lambda^\top \geq \mathbf{0}.$$

Properties of the Ridge Estimator

Hence, the confidence ellipsoid of ridge estimator is indeed smaller than the OLS,

If \mathbf{X} is an orthogonal design matrix,

$$\text{Var}[\hat{\boldsymbol{\beta}}_{\lambda}] = \sigma^2(1 + \lambda)^{-2}\mathbb{I}.$$

$$\text{mse}[\hat{\boldsymbol{\beta}}_{\lambda}] = \sigma^2 \text{trace}(\mathbf{W}_{\lambda}(\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{W}_{\lambda}^{\top}) + \boldsymbol{\beta}^{\top}(\mathbf{W}_{\lambda} - \mathbb{I})^{\top}(\mathbf{W}_{\lambda} - \mathbb{I})\boldsymbol{\beta}.$$

If \mathbf{X} is an orthogonal design matrix,

$$\text{mse}[\hat{\boldsymbol{\beta}}_{\lambda}] = \frac{p\sigma^2}{(1 + \lambda)^2} + \frac{\lambda^2}{(1 + \lambda)^2}\boldsymbol{\beta}^{\top}\boldsymbol{\beta}$$

Properties of the Ridge Estimator

$$\text{mse}[\hat{\beta}_{\lambda}] = \frac{p\sigma^2}{(1+\lambda)^2} + \frac{\lambda^2}{(1+\lambda)^2} \beta^{\top} \beta$$

is minimal for

$$\lambda^* = \frac{p\sigma^2}{\beta^{\top} \beta}$$

Note that there exists $\lambda > 0$ such that $\text{mse}[\hat{\beta}_{\lambda}] < \text{mse}[\hat{\beta}_0] = \text{mse}[\hat{\beta}^{\text{ols}}]$.

SVD decomposition

Consider the singular value decomposition $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$. Then

$$\hat{\beta}^{\text{ols}} = \mathbf{V} \underbrace{\mathbf{D}^{-2}\mathbf{D}} \mathbf{U}^\top \mathbf{y}$$

$$\hat{\beta}_\lambda = \mathbf{V} \underbrace{(\mathbf{D}^2 + \lambda \mathbb{I})^{-1}\mathbf{D}} \mathbf{U}^\top \mathbf{y}$$

Observe that

$$D_{i,i}^{-1} \geq \frac{D_{i,i}}{D_{i,i}^2 + \lambda}$$

hence, the ridge penalty shrinks singular values.

Set now $\mathbf{R} = \mathbf{U}\mathbf{D}$ ($n \times n$ matrix), so that $\mathbf{X} = \mathbf{R}\mathbf{V}^\top$,

$$\hat{\beta}_\lambda = \mathbf{V}(\mathbf{R}^\top \mathbf{R} + \lambda \mathbb{I})^{-1} \mathbf{R}^\top \mathbf{y}$$

Hat matrix and Degrees of Freedom

Recall that $\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$ with

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

Similarly

$$\mathbf{H}_\lambda = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I})^{-1} \mathbf{X}^\top$$

$$\text{trace}[\mathbf{H}_\lambda] = \sum_{j=1}^p \frac{d_{j,j}^2}{d_{j,j}^2 + \lambda} \rightarrow 0, \text{ as } \lambda \rightarrow \infty.$$

Sparsity Issues

In several applications, k can be (very) large, but a lot of features are just noise: $\beta_j = 0$ for many j 's. Let s denote the number of relevant features, with $s \ll k$, cf Hastie, Tibshirani & Wainwright (2015) [Statistical Learning with Sparsity](#),

$$s = \text{card}\{\mathcal{S}\} \text{ where } \mathcal{S} = \{j; \beta_j \neq 0\}$$

The model is now $y = \mathbf{X}_{\mathcal{S}}^{\top} \boldsymbol{\beta}_{\mathcal{S}} + \varepsilon$, where $\mathbf{X}_{\mathcal{S}}^{\top} \mathbf{X}_{\mathcal{S}}$ is a full rank matrix.

Going further on sparsity issues

The Ridge regression problem was to solve

$$\hat{\beta} = \underset{\beta \in \{\|\beta\|_{\ell_2} \leq s\}}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^\top \beta\|_{\ell_2}^2 \}$$

Define $\|\mathbf{a}\|_{\ell_0} = \sum \mathbf{1}(|a_i| > 0)$.

Here $\dim(\beta) = k$ but $\|\beta\|_{\ell_0} = s$.

We wish we could solve

$$\hat{\beta} = \underset{\beta \in \{\|\beta\|_{\ell_0} = s\}}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^\top \beta\|_{\ell_2}^2 \}$$

Problem: it is usually not possible to describe all possible constraints, since $\binom{s}{k}$ coefficients should be chosen here (with k (very) large).

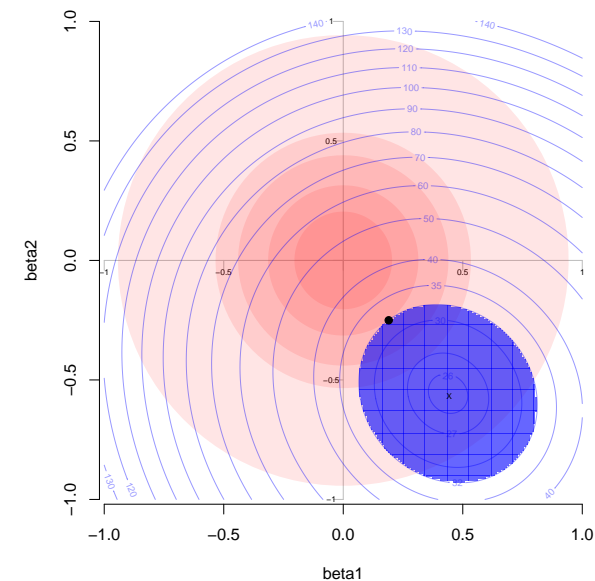
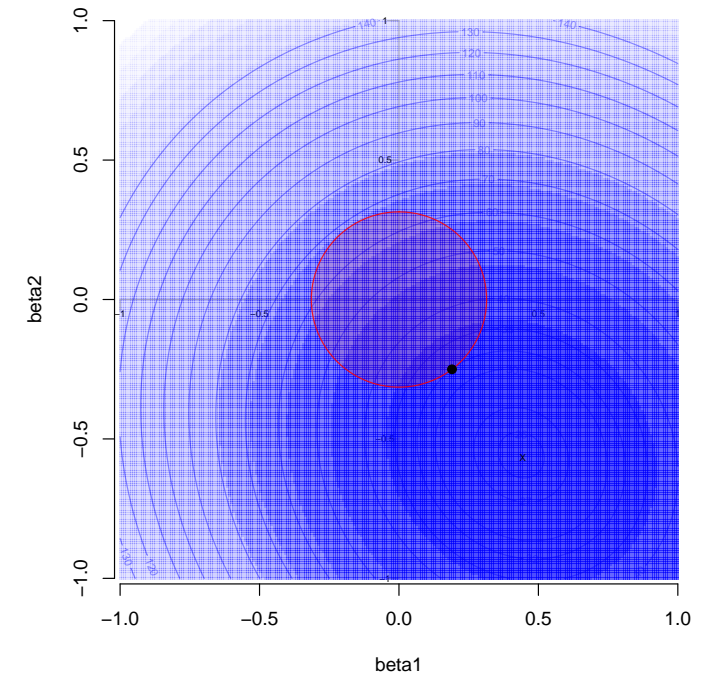
Going further on sparcity issues

In a convex problem, solve the **dual problem**,
e.g. in the Ridge regression : primal problem

$$\min_{\beta \in \{\|\beta\|_{\ell_2} \leq s\}} \{\|Y - X^T \beta\|_{\ell_2}^2\}$$

and the dual problem

$$\min_{\beta \in \{\|Y - X^T \beta\|_{\ell_2} \leq t\}} \{\|\beta\|_{\ell_2}^2\}$$

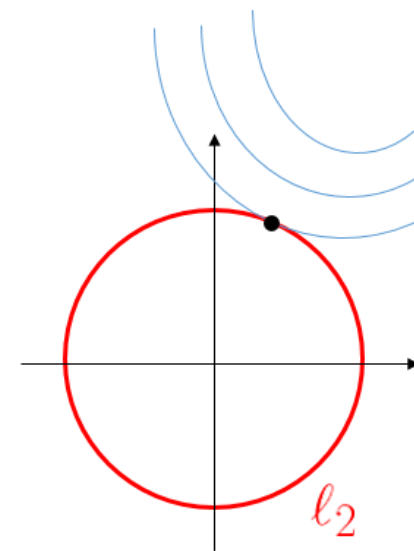
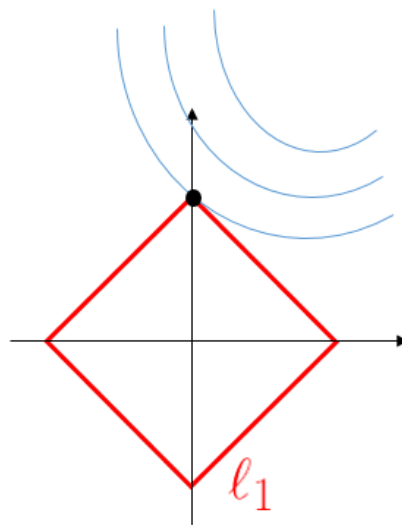
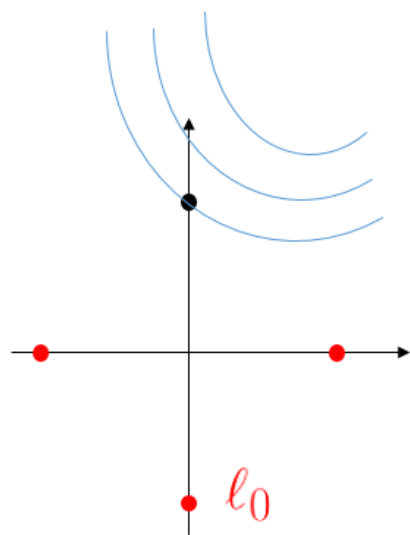


Going further on sparsity issues

Idea: solve the dual problem

$$\hat{\beta} = \underset{\beta \in \{\|\mathbf{Y} - \mathbf{X}^\top \beta\|_{\ell_2} \leq h\}}{\operatorname{argmin}} \{ \|\beta\|_{\ell_0} \}$$

where we might **convexify the ℓ_0 norm**, $\|\cdot\|_{\ell_0}$.



Going further on sparsity issues

On $[-1, +1]^k$, the convex hull of $\|\beta\|_{\ell_0}$ is $\|\beta\|_{\ell_1}$

On $[-a, +a]^k$, the convex hull of $\|\beta\|_{\ell_0}$ is $a^{-1}\|\beta\|_{\ell_1}$

Hence, why not solve

$$\hat{\beta} = \underset{\beta; \|\beta\|_{\ell_1} \leq \tilde{s}}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^\top \beta\|_{\ell_2} \}$$

which is equivalent (Kuhn-Tucker theorem) to the Lagrangian optimization problem

$$\hat{\beta} = \operatorname{argmin} \{ \|\mathbf{Y} - \mathbf{X}^\top \beta\|_{\ell_2}^2 + \lambda \|\beta\|_{\ell_1} \}$$

LASSO *Least Absolute Shrinkage and Selection Operator*

$$\hat{\beta} \in \operatorname{argmin}\{\|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2}^2 + \lambda \|\beta\|_{\ell_1}\}$$

is a convex problem (several algorithms[★]), but not strictly convex (no unicity of the minimum). Nevertheless, predictions $\hat{\mathbf{y}} = \mathbf{x}^T \hat{\beta}$ are unique.

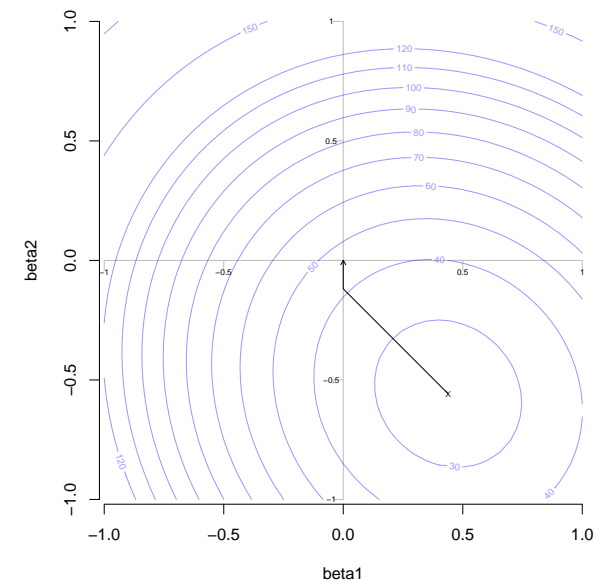
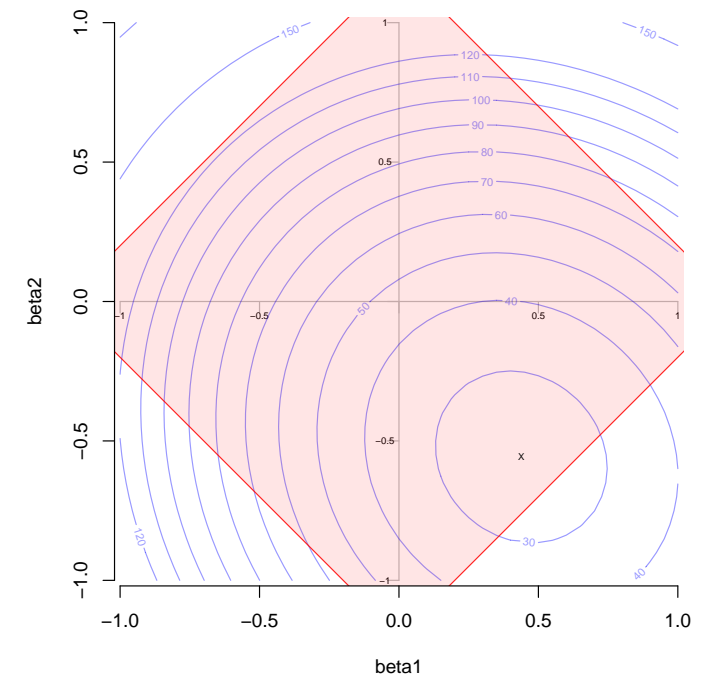
★ MM, minimize majorization, coordinate descent Hunter & Lange (2003) [A Tutorial on MM Algorithms](#).

LASSO Regression

No explicit solution...

If $\lambda \rightarrow 0$, $\hat{\beta}_0^{\text{lasso}} = \hat{\beta}^{\text{ols}}$

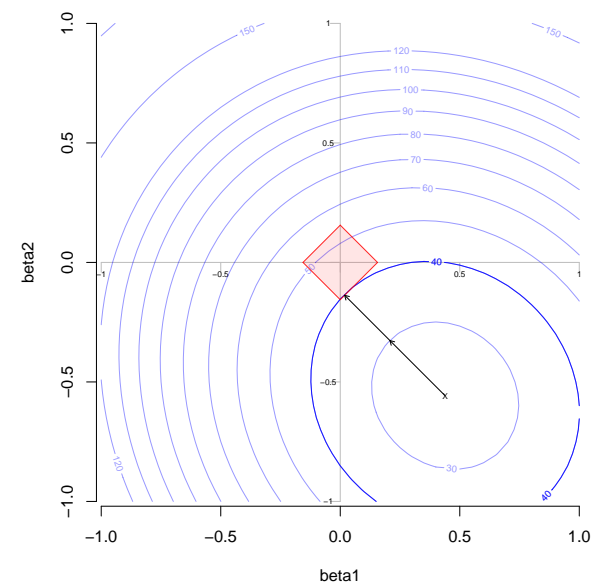
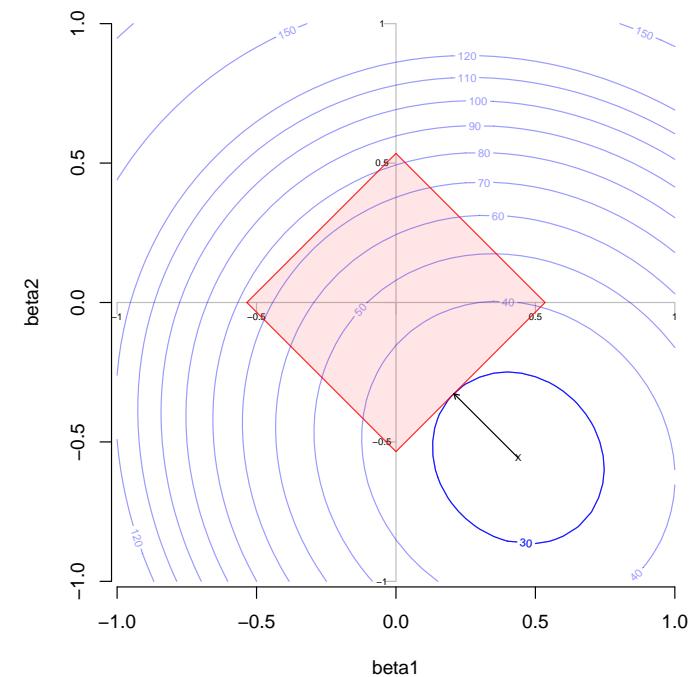
If $\lambda \rightarrow \infty$, $\hat{\beta}_\infty^{\text{lasso}} = \mathbf{0}$.



LASSO Regression

For some λ , there are k 's such that $\hat{\beta}_{k,\lambda}^{\text{lasso}} = 0$.

Further, $\lambda \mapsto \hat{\beta}_{k,\lambda}^{\text{lasso}}$ is **piecewise linear**

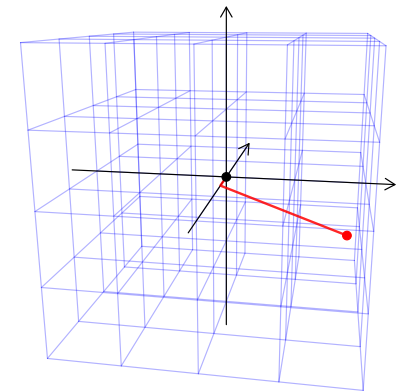


LASSO Regression

In the orthogonal case, $\mathbf{X}^\top \mathbf{X} = \mathbb{I}$,

$$\hat{\beta}_{k,\lambda}^{\text{lasso}} = \text{sign}(\hat{\beta}_k^{\text{ols}}) \left(|\hat{\beta}_k^{\text{ols}}| - \frac{\lambda}{2} \right)$$

i.e. the LASSO estimate is related to the soft threshold function...



Optimal LASSO Penalty

Use cross validation, e.g. K -fold,

$$\hat{\beta}_{(-k)}(\lambda) = \operatorname{argmin} \left\{ \sum_{i \notin \mathcal{I}_k} [y_i - \mathbf{x}_i^\top \beta]^2 + \lambda \|\beta\|_{\ell_1} \right\}$$

then compute the sum of the squared errors,

$$Q_k(\lambda) = \sum_{i \in \mathcal{I}_k} [y_i - \mathbf{x}_i^\top \hat{\beta}_{(-k)}(\lambda)]^2$$

and finally solve

$$\lambda^* = \operatorname{argmin} \left\{ \overline{Q}(\lambda) = \frac{1}{K} \sum_k Q_k(\lambda) \right\}$$

Note that this might overfit, so Hastie, Tibshiriani & Friedman (2009) [Elements of Statistical Learning](#) suggest the largest λ such that

$$\overline{Q}(\lambda) \leq \overline{Q}(\lambda^*) + \operatorname{se}[\lambda^*] \quad \text{with} \quad \operatorname{se}[\lambda]^2 = \frac{1}{K^2} \sum_{k=1}^K [Q_k(\lambda) - \overline{Q}(\lambda)]^2$$

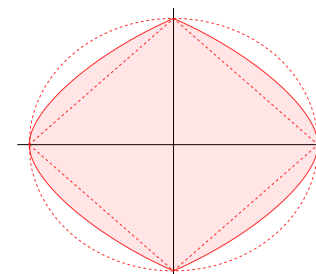
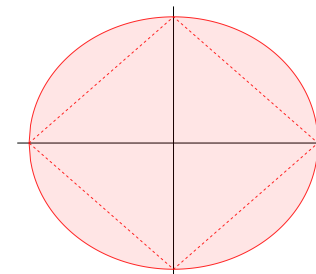
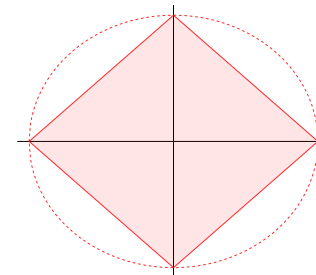
LASSO and Ridge, with R

```

1 > library(glmnet)
2 > chicago=read.table("http://freakonometrics.free.fr/
    chicago.txt",header=TRUE,sep=";")
3 > standardize <- function(x) {(x-mean(x))/sd(x)}
4 > z0 <- standardize(chicago[, 1])
5 > z1 <- standardize(chicago[, 3])
6 > z2 <- standardize(chicago[, 4])
7 > ridge <-glmnet(cbind(z1, z2), z0, alpha=0, intercept=
    FALSE, lambda=1)
8 > lasso <-glmnet(cbind(z1, z2), z0, alpha=1, intercept=
    FALSE, lambda=1)
9 > elastic <-glmnet(cbind(z1, z2), z0, alpha=.5,
    intercept=FALSE, lambda=1)

```

Elastic net, $\lambda_1 \|\beta\|_{\ell_1} + \lambda_2 \|\beta\|_{\ell_2}^2$

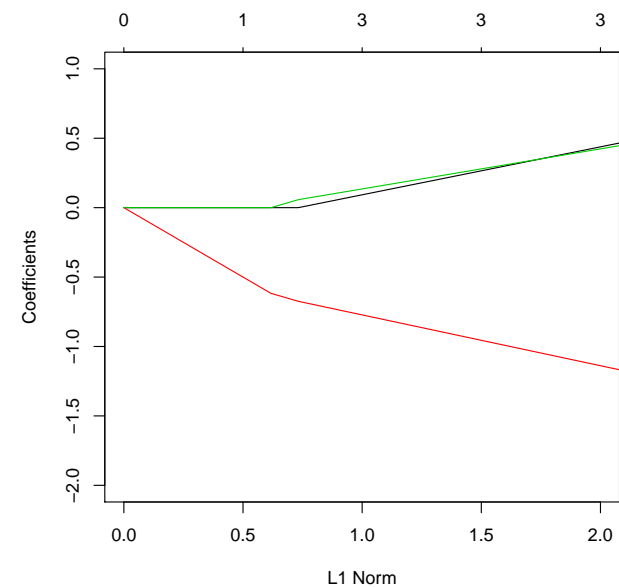
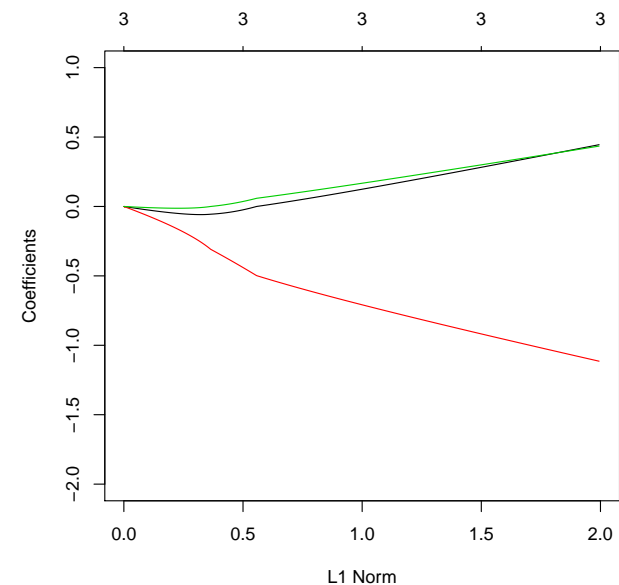
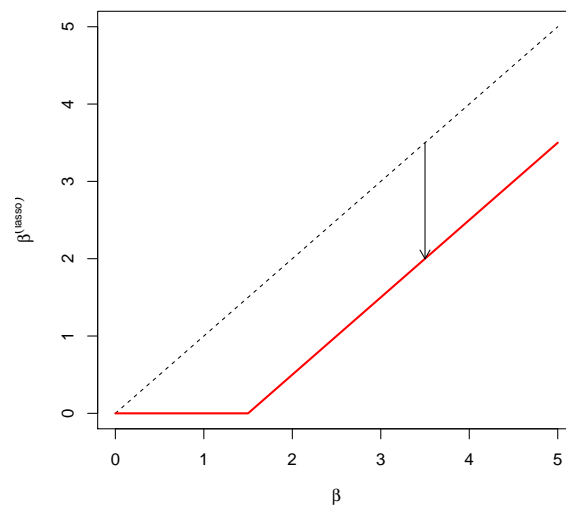
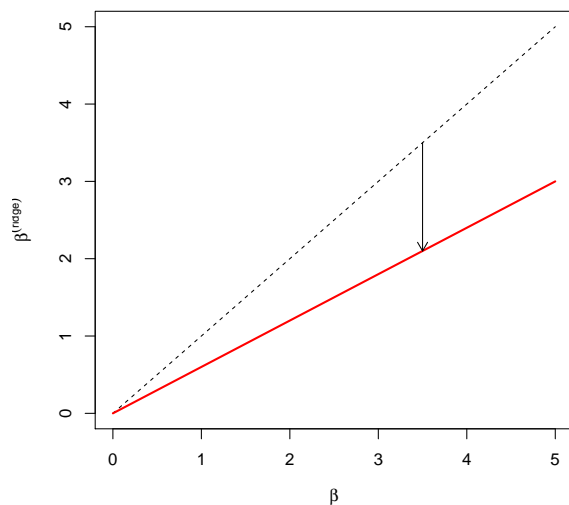


LASSO Regression, Smoothing and Overfit

LASSO can be used to avoid overfit.

Ridge vs. LASSO

Consider simulated data (output on the right).
With orthogonal variables, shrinkage operators are



Optimization Heuristics

First idea: given some initial guess $\beta_{(0)}$, $|\beta| \sim |\beta_{(0)}| + \frac{1}{2|\beta_{(0)}|}(\beta^2 - \beta_{(0)}^2)$

LASSO estimate can probably be derived from iterated Ridge estimates

$$\|\mathbf{y} - \mathbf{X}\beta_{(k+1)}\|_{\ell_2}^2 + \lambda\|\beta_{(k+1)}\|_{\ell_1} \sim \|\mathbf{X}\beta_{(k+1)}\|_{\ell_2}^2 + \frac{\lambda}{2} \sum_j \frac{1}{|\beta_{j,(k)}|} [\beta_{j,(k+1)}]^2$$

which is a **weighted ridge penalty function**

Thus,

$$\beta_{(k+1)} = (\mathbf{X}^\top \mathbf{X} + \lambda \Delta_{(k)})^{-1} \mathbf{X}^\top \mathbf{y}$$

where $\Delta_{(k)} = \text{diag}[|\beta_{j,(k)}|^{-1}]$. Then $\beta_{(k)} \rightarrow \hat{\beta}^{\text{lasso}}$, as $k \rightarrow \infty$.

Properties of LASSO Estimate

From this iterative technique

$$\hat{\beta}_{\lambda}^{\text{lasso}} \sim (\mathbf{X}^{\top} \mathbf{X} + \lambda \Delta)^{-1} \mathbf{X}^{\top} \mathbf{y}$$

where $\Delta = \text{diag}[|\hat{\beta}_{j,\lambda}^{\text{lasso}}|^{-1}]$ if $\hat{\beta}_{j,\lambda}^{\text{lasso}} \neq 0$, 0 otherwise.

Thus,

$$\mathbb{E}[\hat{\beta}_{\lambda}^{\text{lasso}}] \sim (\mathbf{X}^{\top} \mathbf{X} + \lambda \Delta)^{-1} \mathbf{X}^{\top} \mathbf{X} \beta$$

and

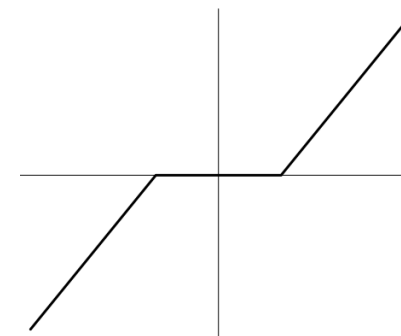
$$\text{Var}[\hat{\beta}_{\lambda}^{\text{lasso}}] \sim \sigma^2 (\mathbf{X}^{\top} \mathbf{X} + \lambda \Delta)^{-1} \mathbf{X}^{\top} \mathbf{X}^{\top} \mathbf{X} (\mathbf{X}^{\top} \mathbf{X} + \lambda \Delta)^{-1} \mathbf{X}^{\top}$$

Optimization Heuristics

Consider here a simplified problem, $\min_{a \in \mathbb{R}} \underbrace{\left\{ \frac{1}{2}(a - b)^2 + \lambda|a| \right\}}_{g(a)}$ with $\lambda > 0$.

Observe that $g'(0) = -b \pm \lambda$. Then

- if $|b| \leq \lambda$, then $a^* = 0$
- if $b \geq \lambda$, then $a^* = b - \lambda$
- if $b \leq -\lambda$, then $a^* = b + \lambda$



$$a^* = \operatorname{argmin}_{a \in \mathbb{R}} \left\{ \frac{1}{2}(a - b)^2 + \lambda|a| \right\} = S_\lambda(b) = \operatorname{sign}(b) \cdot (|b| - \lambda)_+,$$

also called **soft-thresholding** operator.

Optimization Heuristics

Definition for any convex function h , define the **proximal operator** operator of h ,

$$\text{proximal}_h(\mathbf{y}) = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_{\ell_2}^2 + h(\mathbf{x}) \right\}$$

Note that

$$\text{proximal}_{\lambda \|\cdot\|_{\ell_2}^2}(\mathbf{y}) = \frac{1}{1 + \lambda} \mathbf{x} \quad \text{shrinkage operator}$$

$$\text{proximal}_{\lambda \|\cdot\|_{\ell_1}}(\mathbf{y}) = S_\lambda(\mathbf{y}) = \operatorname{sign}(\mathbf{y}) \cdot (|\mathbf{y}| - \lambda)_+$$

Optimization Heuristics

We want to solve here

$$\hat{\boldsymbol{\theta}} \in \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^d} \left\{ \underbrace{\frac{1}{n} \|\mathbf{y} - m_{\boldsymbol{\theta}}(\mathbf{x})\|_{\ell_2}^2}_{f(\boldsymbol{\theta})} + \underbrace{\lambda \operatorname{penalty}(\boldsymbol{\theta})}_{g(\boldsymbol{\theta})} \right\}.$$

where f is convex and smooth, and g is convex, but not smooth...

1. Focus on f : descent lemma, $\forall \boldsymbol{\theta}, \boldsymbol{\theta}'$

$$f(\boldsymbol{\theta}) \leq f(\boldsymbol{\theta}') + \langle \nabla f(\boldsymbol{\theta}'), \boldsymbol{\theta} - \boldsymbol{\theta}' \rangle + \frac{t}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_{\ell_2}^2$$

Consider a gradient descent sequence $\boldsymbol{\theta}_k$, i.e. $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - t^{-1} \nabla f(\boldsymbol{\theta}_k)$, then

$$f(\boldsymbol{\theta}) \leq \overbrace{f(\boldsymbol{\theta}_k) + \langle \nabla f(\boldsymbol{\theta}_k), \boldsymbol{\theta} - \boldsymbol{\theta}_k \rangle + \frac{t}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_k\|_{\ell_2}^2}^{\varphi(\boldsymbol{\theta}): \boldsymbol{\theta}_{k+1} = \operatorname{argmin}\{\varphi(\boldsymbol{\theta})\}}$$

Optimization Heuristics

2. Add function g

$$\overbrace{f(\boldsymbol{\theta}) + g(\boldsymbol{\theta})}^{\psi(\boldsymbol{\theta})} \leq f(\boldsymbol{\theta}_k) + \langle \nabla f(\boldsymbol{\theta}_k), \boldsymbol{\theta} - \boldsymbol{\theta}_k \rangle + \frac{t}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_k\|_{\ell_2}^2 + g(\boldsymbol{\theta})$$

And one can proof that

$$\boldsymbol{\theta}_{k+1} = \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^d} \left\{ \psi(\boldsymbol{\theta}) \right\} = \operatorname{proximal}_{g/t} \left(\boldsymbol{\theta}_k - t^{-1} \nabla f(\boldsymbol{\theta}_k) \right)$$

so called **proximal gradient descent algorithm**, since

$$\operatorname{argmin} \{ \psi(\boldsymbol{\theta}) \} = \operatorname{argmin} \left\{ \frac{t}{2} \left\| \boldsymbol{\theta} - \left(\boldsymbol{\theta}_k - t^{-1} \nabla f(\boldsymbol{\theta}_k) \right) \right\|_{\ell_2}^2 + g(\boldsymbol{\theta}) \right\}$$

Coordinate-wise minimization

Consider some convex differentiable $f : \mathbb{R}^k \rightarrow \mathbb{R}$ function.

Consider $\mathbf{x}^* \in \mathbb{R}^k$ obtained by minimizing along each coordinate axis, i.e.

$$f(x_1^*, x_{i-1}^*, \mathbf{x}_i, x_{i+1}^*, \dots, x_k^*) \geq f(x_1^*, x_{i-1}^*, \mathbf{x}_i^*, x_{i+1}^*, \dots, x_k^*)$$

for all i . Is \mathbf{x}^* a global minimizer? i.e.

$$f(\mathbf{x}) \geq f(\mathbf{x}^*), \quad \forall \mathbf{x} \in \mathbb{R}^k.$$

Yes. If f is convex and differentiable.

$$\nabla f(\mathbf{x})|_{\mathbf{x}=\mathbf{x}^*} = \left(\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_k} \right) = \mathbf{0}$$

There might be problem if f is not differentiable (except in each axis direction).

If $f(\mathbf{x}) = g(\mathbf{x}) + \sum_{i=1}^k h_i(x_i)$ with g convex and differentiable, yes, since

$$f(\mathbf{x}) - f(\mathbf{x}^*) \geq \nabla g(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) + \sum_i [h_i(x_i) - h_i(x_i^*)]$$

Coordinate-wise minimization

$$f(\mathbf{x}) - f(\mathbf{x}^*) \geq \sum_i \underbrace{[\nabla_i g(\mathbf{x}^*)^\top (x_i - x_i^*) h_i(x_i) - h_i(x_i^*)]}_{\geq 0} \geq 0$$

Thus, for functions $f(\mathbf{x}) = g(\mathbf{x}) + \sum_{i=1}^k h_i(x_i)$ we can use coordinate descent to find a minimizer, i.e. at step j

$$x_1^{(j)} \in \operatorname{argmin}_{x_1} f(x_1, x_2^{(j-1)}, x_3^{(j-1)}, \dots, x_k^{(j-1)})$$

$$x_2^{(j)} \in \operatorname{argmin}_{x_2} f(x_1^{(j)}, x_2, x_3^{(j-1)}, \dots, x_k^{(j-1)})$$

$$x_3^{(j)} \in \operatorname{argmin}_{x_3} f(x_1^{(j)}, x_2^{(j)}, x_3, \dots, x_k^{(j-1)})$$

Tseng (2001) **Convergence of Block Coordinate Descent Method**: if f is continuous, then \mathbf{x}^∞ is a minimizer of f .

Application in Linear Regression

Let $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$, with $\mathbf{y} \in \mathbb{R}^n$ and $\mathbf{A} \in \mathcal{M}_{n \times k}$. Let $\mathbf{A} = [\mathbf{A}_1, \dots, \mathbf{A}_k]$.

Let us minimize in direction i . Let \mathbf{x}_{-i} denote the vector in \mathbb{R}^{k-1} without x_i .

Here

$$0 = \frac{\partial f(\mathbf{x})}{\partial x_i} = \mathbf{A}_i^\top [\mathbf{A}\mathbf{x} - \mathbf{y}] = \mathbf{A}_i^\top [\mathbf{A}_i x_i + \mathbf{A}_{-i} \mathbf{x}_{-i} - \mathbf{y}]$$

thus, the optimal value is here

$$x_i^* = \frac{\mathbf{A}_i^\top [\mathbf{A}_{-i} \mathbf{x}_{-i} - \mathbf{y}]}{\mathbf{A}_i^\top \mathbf{A}_i}$$

Application to LASSO

Let $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \lambda \|\mathbf{x}\|_{\ell_1}$, so that the non-differentiable part is separable, since $\|\mathbf{x}\|_{\ell_1} = \sum_{i=1}^k |x_i|$.

Let us minimize in direction i . Let \mathbf{x}_{-i} denote the vector in \mathbb{R}^{k-1} without x_i . Here

$$0 = \frac{\partial f(\mathbf{x})}{\partial x_i} = \mathbf{A}_i^\top [\mathbf{A}_i x_i + \mathbf{A}_{-i} \mathbf{x}_{-i} - \mathbf{y}] + \lambda s_i$$

where $s_i \in \partial |x_i|$. Thus, solution is obtained by soft-thresholding

$$x_i^* = S_{\lambda / \|\mathbf{A}_i\|^2} \left(\frac{\mathbf{A}_i^\top [\mathbf{A}_{-i} \mathbf{x}_{-i} - \mathbf{y}]}{\mathbf{A}_i^\top \mathbf{A}_i} \right)$$

Convergence rate for LASSO

Let $f(\mathbf{x}) = g(\mathbf{x}) + \lambda \|\mathbf{x}\|_{\ell_1}$ with

- g convex, ∇g Lipschitz with constant $L > 0$, and $Id - \nabla g/L$ monotone inscreasing in each component
- there exists \mathbf{z} such that, componentwise, either $\mathbf{z} \geq S_\lambda(\mathbf{z} - \nabla g(\mathbf{z}))$ or $\mathbf{z} \leq S_\lambda(\mathbf{z} - \nabla g(\mathbf{z}))$

Saka & Tewari (2010), **On the finite time convergence of cyclic coordinate descent methods** proved that a coordinate descent starting from \mathbf{z} satisfies

$$f(\mathbf{x}^{(j)}) - f(\mathbf{x}^\star) \leq \frac{L \|\mathbf{z} - \mathbf{x}^\star\|^2}{2j}$$

Graphical Lasso and Covariance Estimation

We want to estimate an (unknown) covariance matrix Σ , or Σ^{-1} .

An estimate for Σ^{-1} is Θ^* solution of

$$\Theta \in \underset{\Theta \in \mathcal{M}_{k \times k}}{\operatorname{argmin}} \{ -\log[\det(\Theta)] + \operatorname{trace}[S\Theta] + \lambda \|\Theta\|_{\ell_1} \} \quad \text{where } S = \frac{X^\top X}{n}$$

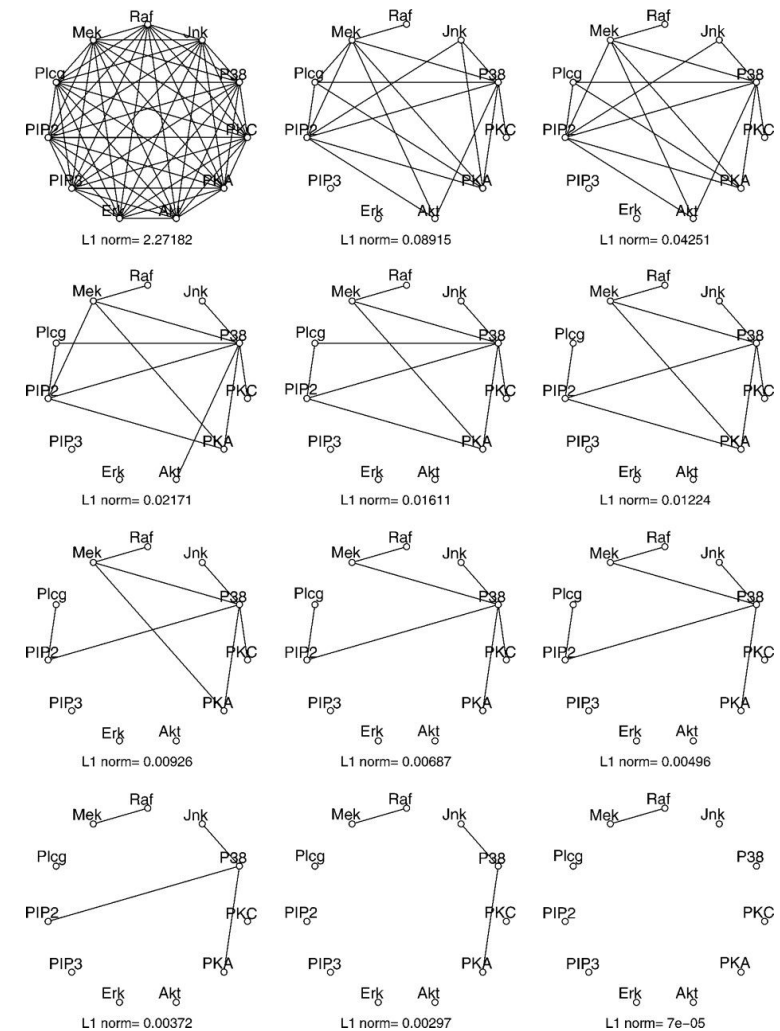
and where $\|\Theta\|_{\ell_1} = \sum |\Theta_{i,j}|$.

See van Wieringen (2016) [Undirected network reconstruction from high-dimensional data](#) and <https://github.com/kaizhang/glasso>

Application to Network Simplification

Can be applied on networks, to spot ‘significant’ connexions...

Source: <http://khughitt.github.io/graphical-lasso/>



Extention of Penalization Techniques

In a more general context, we want to solve

$$\hat{\boldsymbol{\theta}} \in \operatorname{argmin}_{\boldsymbol{\theta} \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(y_i, m_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \lambda \cdot \text{penalty}(\boldsymbol{\theta}) \right\}.$$