

# Penalised regression

Ridge, LASSO and elastic net regression

COWIDUR

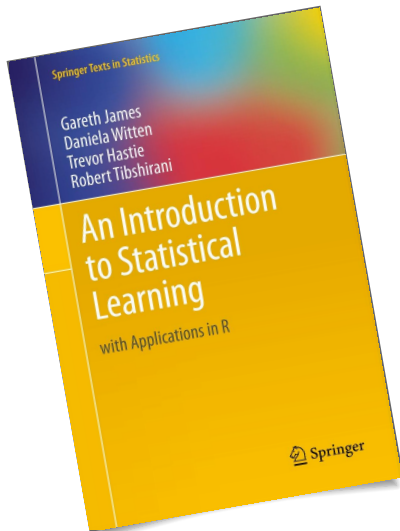
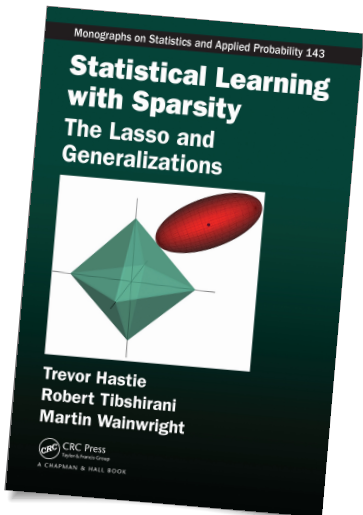
Torben Tvedebrink  
tvede@math.aau.dk

Department of Mathematical Sciences



**AALBORG UNIVERSITY**  
DENMARK

*Version: 08/05/2019 11:05*



Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

1 Regularised regression

Ridge regression

LASSO regression

Extensions

# Bet on sparsity principle



*Use a procedure that does well in sparse problems,  
since no procedure does well in dense problems.*

When  $p \gg n$  (the “short, fat data problem”), two things go wrong:

- ▶ The Curse of Dimensionality is acute.
- ▶ There are insufficient degrees of freedom to estimate the full model.

However, there is a substantial body of practical experience which indicates that, in some circumstances, one can actually make good statistical inferences and predictions.

Penalised  
regression

Torben Tvedebrink  
tvede@math.aau.dk

2 Regularised  
regression

Ridge regression

LASSO regression

Extensions

# Our point of departure

In linear regression we assume that the  $i$ th response,  $y_i$ , can be modelled using a linear relationship between some covariates and the response with an additive error term with constant variance

$$y_i = \beta_0 + \sum_{j=1}^p x_{ij}\beta_j + \varepsilon_i$$



Penalised  
regression

Torben Tvedebrink  
tvede@math.aau.dk

3 Regularised  
regression

Ridge regression

LASSO regression

Extensions

# Our point of departure



In linear regression we assume that the  $i$ th response,  $y_i$ , can be modelled using a linear relationship between some covariates and the response with an additive error term with constant variance

$$y_i = \beta_0 + \sum_{j=1}^p x_{ij}\beta_j + \varepsilon_i$$

If we have observations,  $i = 1, \dots, n > p$ , we have that the least squares estimator for  $\beta_0$  and  $\beta = (\beta_1, \dots, \beta_p)$  is given by

$$(\hat{\beta}_0, \hat{\beta}) = \arg \min_{\beta_0, \beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2$$

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

3 Regularised regression

Ridge regression

LASSO regression

Extensions

# Least squares

On a *budget*



Penalised  
regression

Torben Tvedebrink  
tvede@math.aau.dk

4 Regularised  
regression

Ridge regression

LASSO regression

Extensions

Imagine that we only had a limited *budget* of regression coefficients,  $t$ , such that the sum  $\sum_{j=1}^p h(\beta_j)$  was restricted by  $t$ , then the solution should obey this constraint

$$\min_{\beta_0, \beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 \quad \text{such that} \quad \sum_{j=1}^p h(\beta_j) \leq t$$

# Least squares

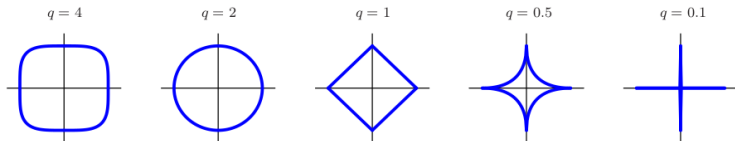
On a *budget*



Imagine that we only had a limited *budget* of regression coefficients,  $t$ , such that the sum  $\sum_{j=1}^p h(\beta_j)$  was restricted by  $t$ , then the solution should obey this constraint

$$\min_{\beta_0, \beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 \quad \text{such that} \quad \sum_{j=1}^p h(\beta_j) \leq t$$

Constraint regions for  $\sum_{j=1}^p h(\beta_j) = |\beta_j|^q \leq 1$ .



For all  $q < 1$  the constraint region is non-convex.

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

4 Regularised regression

Ridge regression

LASSO regression

Extensions

# Least squares

On a *budget*

Imagine that we only had a limited *budget* of regression coefficients,  $t$ , such that the sum  $\sum_{j=1}^p h(\beta_j)$  was restricted by  $t$ , then the solution should obey this constraint

$$\min_{\beta_0, \beta} \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j)^2 \quad \text{such that} \quad \sum_{j=1}^p h(\beta_j) \leq t$$

For

- ▶  $h(\beta_j) = |\beta_j|$  we term the regression problem the *LASSO*, and
- ▶  $h(\beta_j) = \beta_j^2$  we refer to the problem as *ridge regression*.



Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

4

Regularised regression

Ridge regression

LASSO regression

Extensions



# Reasons for abandoning least squares



- ▶ The *prediction accuracy* can sometimes be improved because even though least squares has zero bias, its high variance may cause bad prediction ability. Hence, shrinking some coefficients, or setting the *noisy terms* to zero, may improve the accuracy.

## Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

5

## Regularised regression

Ridge regression

LASSO regression

Extensions

# Reasons for abandoning least squares



- ▶ The *prediction accuracy* can sometimes be improved because even though least squares has zero bias, its high variance may cause bad prediction ability. Hence, shrinking some coefficients, or setting the *noisy terms* to zero, may improve the accuracy.
- ▶ The second reason is *interpretation*. The fewer terms to interpret the easier it gets.

## Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

## 5 Regularised regression

Ridge regression

LASSO regression

Extensions

# Reasons for abandoning least squares



- ▶ The *prediction accuracy* can sometimes be improved because even though least squares has zero bias, its high variance may cause bad prediction ability. Hence, shrinking some coefficients, or setting the *noisy terms* to zero, may improve the accuracy.
- ▶ The second reason is *interpretation*. The fewer terms to interpret the easier it gets.
- ▶ The third reason being that it fails for *wide* data, i.e. data for which  $p \gg n$

## Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

### 5 Regularised regression

Ridge regression

LASSO regression

Extensions

# Standardisation of $\mathbf{X}$



As the *numerical value* of coefficients is sensitive to the scale of the covariates, it is typically preferred to standardise the  $\mathbf{X}$  matrix before estimating the coefficients. That is,

$$\sum_{i=1}^n x_{ij} = 0 \quad \text{and} \quad \sum_{i=1}^n x_{ij}^2 = n$$

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

6 Regularised regression

Ridge regression

LASSO regression

Extensions

# Standardisation of $\mathbf{X}$ and centering of $y$



As the *numerical value* of coefficients is sensitive to the scale of the covariates, it is typically preferred to standardise the  $\mathbf{X}$  matrix before estimating the coefficients. That is,

$$\sum_{i=1}^n x_{ij} = 0 \quad \text{and} \quad \sum_{i=1}^n x_{ij}^2 = n$$

And in order to discard the intercept,  $\beta_0$ , from the regularisation in the case of linear regression we center the response

$$\sum_{i=1}^n y_i = 0$$

Penalised  
regression

Torben Tvedebrink  
tvede@math.aau.dk

6 Regularised  
regression

Ridge regression

LASSO regression

Extensions

# The *wide* data problem

In the case where  $p \gg n$ , the least squares estimator is undefined as  $(\mathbf{X}^T \mathbf{X})$  isn't invertible because  $\mathbf{X}$  is not of full rank. Hence,  $\hat{\beta}^{\text{ols}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  cannot be evaluated.



Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

7 Ridge regression

LASSO regression

Extensions

# The *wide* data problem



In the case where  $p \gg n$ , the least squares estimator is undefined as  $(\mathbf{X}^\top \mathbf{X})$  isn't invertible because  $\mathbf{X}$  is not of full rank. Hence,  $\hat{\beta}^{\text{ols}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$  cannot be evaluated.

A solution to this is to add an invertible matrix to  $\mathbf{X}^\top \mathbf{X}$  to obtain an invertible matrix. The simplest such candidate is  $\lambda \mathbf{I}_p$ , for some positive  $\lambda \in \mathbb{R}$ :

$$\hat{\beta}^{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^\top \mathbf{y},$$

which is what is referred to as the ridge regression estimator.

Penalised  
regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised  
regression

7 Ridge regression

LASSO regression

Extensions

# Ridge regression



For the least squares regression problem with a budget on the squared entries of  $\beta$  we have

$$\min_{\beta} \sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 \quad \text{such that} \quad \sum_{j=1}^p \beta_j^2 \leq t.$$

This can also be stated as

$$\min_{\beta} \sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p \beta_j^2.$$

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

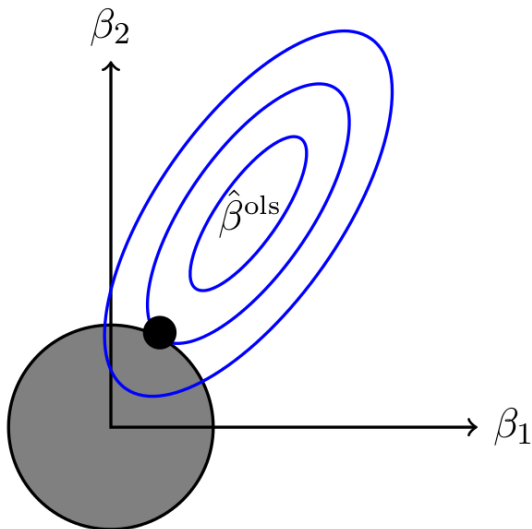
8 Ridge regression

LASSO regression

Extensions



# Visual representation of $\hat{\beta}^{\text{ridge}}$ Compared to $\hat{\beta}^{\text{ols}}$ (in two dimensions)



Penalised  
regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised  
regression

9 Ridge regression

LASSO regression

Extensions

# LASSO regression



Now, what happens if we instead of using a squared penalty,  $\beta_j^2$ , uses the absolute penalty,  $|\beta_j|$ ?

Well – we obtain the LASSO

$$\min_{\beta} \sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 \quad \text{such that} \quad \sum_{j=1}^p |\beta_j| \leq t.$$

and again an equivalent form

$$\min_{\beta} \sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j|.$$

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

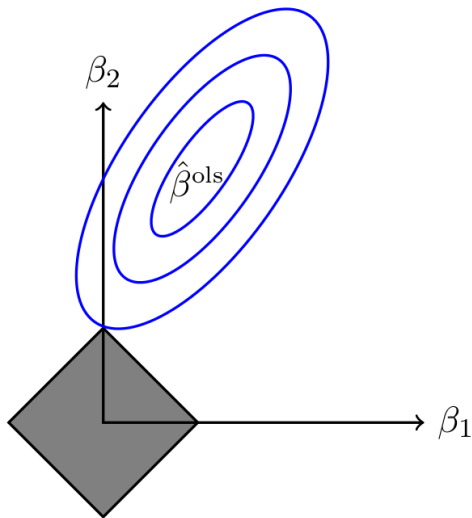
Ridge regression

10 LASSO regression

Extensions

# Visual representation of $\hat{\beta}^{\text{lasso}}$

Compared to  $\hat{\beta}^{\text{ols}}$  (in two dimensions)



Penalised  
regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised  
regression

Ridge regression

11 LASSO regression

Extensions

# LASSO solution

## Comparison to Least Squares solution

With a standardized predictor, the LASSO solution is a soft-thresholded version of the ordinary least-squares (OLS) estimate  $\hat{\beta}$

$$\hat{\beta}_j = \begin{cases} \hat{\beta}_j^{(\text{OLS})} + \lambda, & \hat{\beta}_j^{(\text{OLS})} < -\lambda \\ 0, & -\lambda \leq \hat{\beta}_j^{(\text{OLS})} \leq \lambda \\ \hat{\beta}_j^{(\text{OLS})} - \lambda, & \hat{\beta}_j^{(\text{OLS})} > \lambda. \end{cases}$$



Penalised  
regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised  
regression

Ridge regression

12 LASSO regression

Extensions

# LASSO solution

## Comparison to Least Squares solution

With a standardized predictor, the LASSO solution is a soft-thresholded version of the ordinary least-squares (OLS) estimate  $\hat{\beta}$

$$\hat{\beta}_j = \begin{cases} \hat{\beta}_j^{(\text{OLS})} + \lambda, & \hat{\beta}_j^{(\text{OLS})} < -\lambda \\ 0, & -\lambda \leq \hat{\beta}_j^{(\text{OLS})} \leq \lambda \\ \hat{\beta}_j^{(\text{OLS})} - \lambda, & \hat{\beta}_j^{(\text{OLS})} > \lambda. \end{cases}$$

This relationship also holds (in a slightly modified way) in case where the  $\hat{\beta}^{(\text{OLS})}$  do not exist.



Penalised  
regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised  
regression

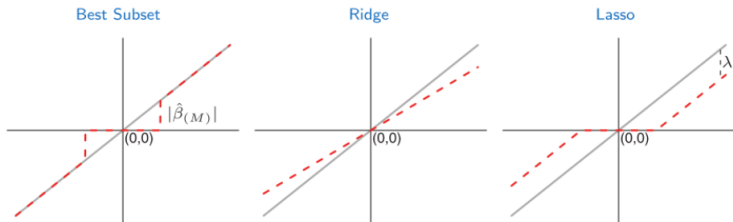
Ridge regression

12 LASSO regression

Extensions

# Soft thresholding

Modifications of the OLS estimates (if they exists)



Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

13 LASSO regression

Extensions

# Elastic Net

The best from two worlds?

A downside with the Lasso is that it may have difficulties when several variables are collinear, such that linear combinations of them are hard to distinguish.

In such a case the Ridge Regression is better as it will typically form an average of the variables. Hence, for stable selection of variables in this case Ridge Regression may be preferred.

However, Ridge Regression seldom sets any parameters to zero, i.e. no variable selection which is what we would like in the end...



Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

14

Elastic Net  
Estimation  
Group LASSO  
Bayesian perspective  
Bootstrap

25

Department of  
Mathematical Sciences

# Elastic Net

The best from two worlds?

The solution to the problem is Elastic Net, which incorporates both the Lasso and Ridge penalties in a convex way:

$$\min_{\beta} \sum_{i=1}^2 (y_i - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p \{\alpha |\beta_j| + (1 - \alpha) \beta_j^2\},$$

where  $\alpha$  is yet another tuning parameter deciding the amount of Lasso ( $\alpha = 1$ ) and Ridge ( $\alpha = 0$ ) penalty that goes into the solution.

Both  $\alpha$  and  $\lambda$  are selected based on cross-validation.



Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

15

Elastic Net  
Estimation  
Group LASSO  
Bayesian perspective  
Bootstrap

25

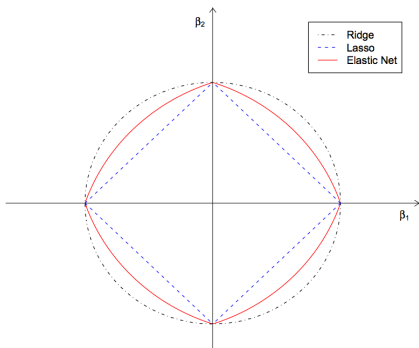
Department of  
Mathematical Sciences



# Elastic Net

The best from two worlds?

In the Figure below we see the three types of regularisation discussed above. The shape of the Elastic Net solution area depends on  $\alpha$  - the closer to 1 the more square it is, and the closer to 0 the more spherical.



Penalised  
regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised  
regression

Ridge regression

LASSO regression

Extensions

16

Elastic Net  
Estimation  
Group LASSO  
Bayesian perspective  
Bootstrap

25

Department of  
Mathematical Sciences

# A brief history of LASSO algorithms

And practical limits (in terms of number of covariates,  $p$ )



As mentioned earlier, the lasso penalty lacks a closed form solution in general.

As a result, optimisation algorithms must be employed to find the minimising solution

The historical efficiency of algorithms to fit lasso models can be summarized as follows:

Year	Algorithm	Operations	Practical limit
1996	QP <sup>†</sup>	$O(n^2p)$	$\sim 100$
2003	LARS <sup>‡</sup>	$O(np^2)$	$\sim 10,000$
2008	Coordinate descent	$O(np)$	$\sim 1,000,000$

†: Quadratic Programming

‡: Least Angle Regression

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

Elastic Net

Group LASSO

Bayesian perspective

Bootstrap

17

25

# Group LASSO

Setting groups of coefficients to zero

The LASSO penalises each  $\beta_j$  coefficient individually by assessing the correlation between the partial residuals and the explanatory variable.

However, in the case of regression involving factors, the usual dummy variable encoding implies that the different derived dummy variables are penalised individually.

This causes some problems as we prefer that *all* dummy variables are set to zero, i.e. *all* levels of the factor are insignificant.

– This why we in ordinary regression use `anova(lm(...))` to test for significance of factors and not the individual *t*-tests reported in `summary(lm(...))`.



Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

Elastic Net  
Estimation

18

**Group LASSO**  
Bayesian perspective  
Bootstrap

25

Department of  
Mathematical Sciences

# Group LASSO

## Adjusting the penalty



SLS use  $\theta$  for the group LASSO in order to avoid confusion between the LASSO with penalty on the individual  $\beta$  parameters. Hence, we may reformulate the minimisation problem as

$$\min_{\theta_0, \theta} \left\{ \frac{1}{2} \sum_{i=1}^n (y_i - \theta_0 - \sum_{j=1}^J z_{ij}^T \theta_j)^2 + \lambda \sum_{j=1}^J \|\theta_j\|_2 \right\},$$

where  $\|\theta_j\|_2 = \sqrt{\sum_{k=1}^{p_j} \theta_{jk}^2}$  is the  $\ell_2$ -norm.

For  $p_j = 1$  we have that  $\|\theta_j\|_2 = \sqrt{\theta_{j1}^2} = |\theta_{j1}|$ , which is just the LASSO penalty.

For  $p_j > 1$ , the  $\ell_2$ -penalty will imply that either  $\theta_j = 0$  or non-zero.

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

Elastic Net

Estimation

19

Group LASSO

Bayesian perspective

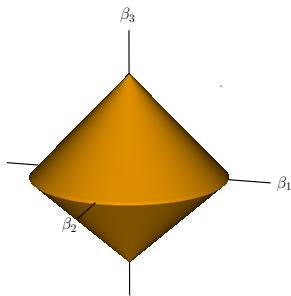
Bootstrap

25

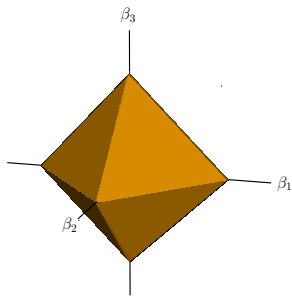
Department of  
Mathematical Sciences

# The group LASSO ball

In  $\mathbb{R}^3$



*Left:* Group LASSO;  
 $\theta_1 = (\beta_1, \beta_2)$  and  $\theta_2 = \beta_3$



*Right:* LASSO;  $\theta_j = \beta_j$

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

Elastic Net

Estimation

Group LASSO

Bayesian perspective

Bootstrap

20

25

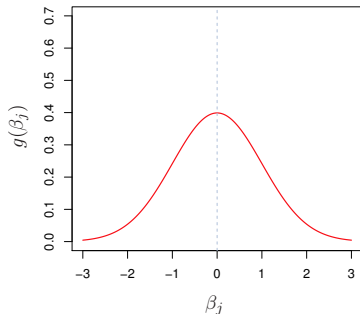
# Bayesian perspective

The Bridge and BLASSO priors – Marginal priors for  $\beta_j$



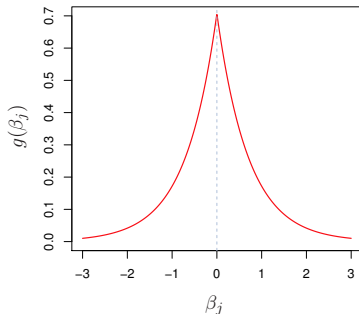
Normal prior

Ridge regression



Laplacian prior

LASSO



Interpretation:

For the LASSO we *a priori* expect more parameters (due to the peaked nature of the Laplace prior) to be zero than for the normal ridge regression prior.

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

Elastic Net

Estimation

Group LASSO

Bayesian perspective

Bootstrap

21

25

# Bootstrap

Yet another handy application of the resampling technique



The bootstrap procedure may be one of the most important contributions to modern statistics – simple and yet powerful.

We may use the bootstrap as a non-parametric alternative to the Bayesian LASSO in order to assess the coefficient variability.

By *permuting the data* the parameter estimates may differ substantially. Hence, in order to capture this *repetition* of the experiment over and over again, resampling the data with replacement reflects this uncertainty.

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

Elastic Net

Estimation

Group LASSO

Bayesian perspective

Bootstrap

22

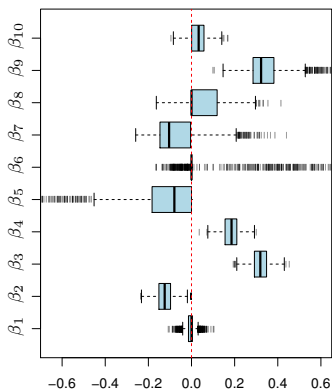
25

# Example

## Diabetes – non-parametric Bootstrapped parameters



Bootstrap Samples



Bootstrap Probability of 0

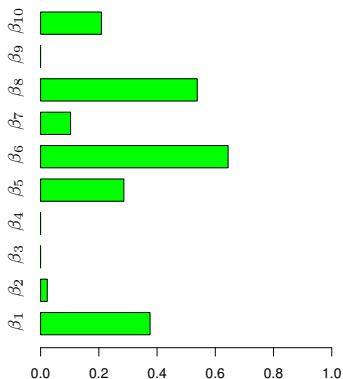


Figure 6.4 in SLS based on 1,000 bootstrap samples for fixed  $\hat{\lambda}_{CV}$ . *Left:* Parameter estimates  $\hat{\beta}^*$ ; *Right:* “Significance” of each covariate.

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

Elastic Net

Estimation

Group LASSO

Bayesian perspective

Bootstrap

23

25

Department of  
Mathematical Sciences



# Example

## Diabetes – Bayesian posteriors



Bayesian Posterior Samples

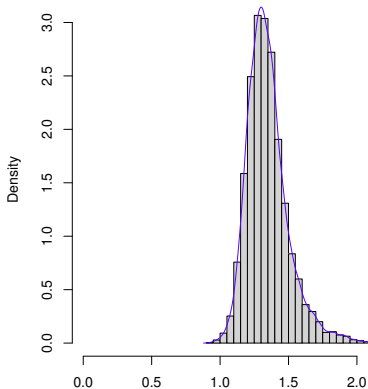
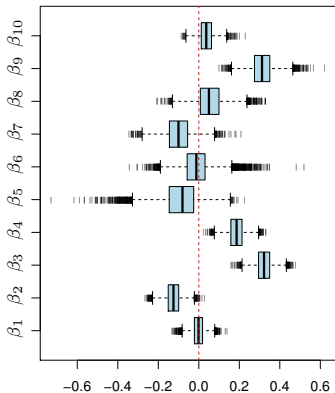


Figure 6.3 in SLS based on 10,000 samples from the posterior distributions (*left*) and  $\|\beta\|_1$  (*right*).

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

Elastic Net

Estimation

Group LASSO

Bayesian perspective

Bootstrap

24

25

# Example

## Diabetes – parametric Bootstrapped parameters

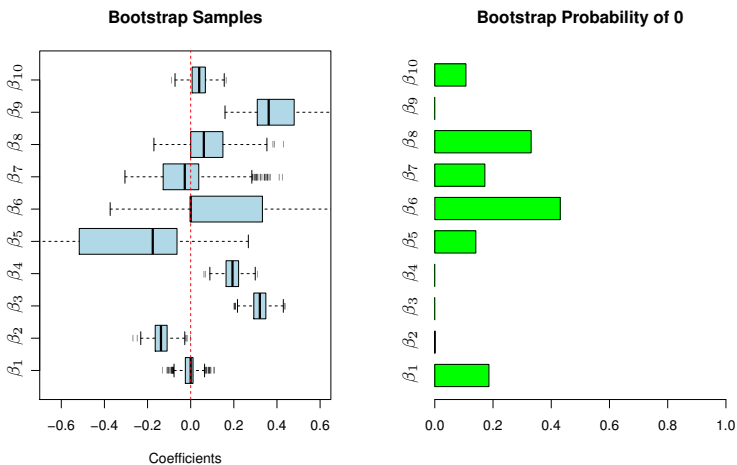


Figure 6.7 in SLS.  $\mathbf{y}^*$  is sampled from the estimated model with  $(\hat{\beta}, \hat{\sigma}^2)$  from the full data, and estimated for  $(\mathbf{y}^*, \mathbf{X})$ .

Penalised regression

Torben Tvedebrink  
tvede@math.aau.dk

Regularised regression

Ridge regression

LASSO regression

Extensions

Elastic Net

Estimation

Group LASSO

Bayesian perspective

Bootstrap

25

25

Department of  
Mathematical Sciences