Context
Proposed Models
Structure estimation
Results
Missing values
Tools

# CorReg

Clément Théry, Christophe Biernacki, Gaétan Loridant

ArcelorMittal Dunkerque, Université de Lille 1,équipe MΘdal Inria

February 8, 2015

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Context

Proposed Models

Structure estimation

Results

Missing values

Tools

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Industrial context
Statistical context

1. Steel industry databases.
2. Goal: To understand and prevent quality problems on finished product, knowing the whole process, <u>without a priori</u>.

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Industrial context
Statistical context

## Regression

$$\boldsymbol{Y} = \boldsymbol{X}\beta + \varepsilon \qquad (1)$$

where $\varepsilon \sim \mathcal{N}(0, \sigma_Y^2 \boldsymbol{I}_n)$

Context
Proposed Models
Structure estimation
Results
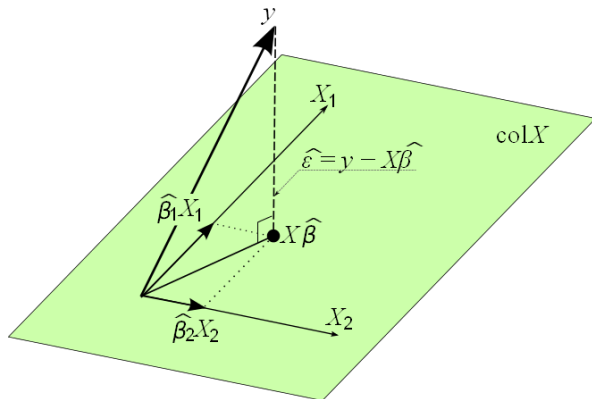Missing values
Tools

Industrial context
Statistical context

# OLS



Figure: Multiple linear regression with Ordinary Least Squares seen as a projection on the $d-$dimensional hyperplane spanned by the regressors $X$. Public domain image.

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Industrial context
Statistical context

## OLS

$\beta$ can be estimated by $\hat{\beta}$ with Ordinary Least Squares (OLS), that is the unbiased maximum likelihood estimator [Saporta, 2006, Dodge and Rousson, 2004]:

$$\hat{\beta}_{OLS} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{Y} \qquad (2)$$

with variance matrix

$$\text{Var}(\hat{\beta}_{OLS}) = \sigma_Y^2\,(\boldsymbol{X}'\boldsymbol{X})^{-1}. \qquad (3)$$

In fact it is the Best Linear Unbiased Estimator (BLUE). The theoretical MSE is given by

$$\text{MSE}(\hat{\beta}_{OLS}) = \sigma_Y^2\,\text{Tr}((\boldsymbol{X}'\boldsymbol{X})^{-1}).$$

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Industrial context
Statistical context

## Running example

$\boldsymbol{X}^1, \boldsymbol{X}^2, \boldsymbol{X}^4, \boldsymbol{X}^5 \sim \mathcal{N}(0, 1)$ and $\boldsymbol{X}^3 = \boldsymbol{X}^1 + \boldsymbol{X}^2 + \varepsilon_1$ where
$\varepsilon_1 \sim \mathcal{N}(\boldsymbol{0}, \sigma_1^2 \boldsymbol{I}_n)$.
Two *scenarii* for $\boldsymbol{Y}$:
$\beta = (1, 1, 1, 1, 1)'$ and $\sigma_Y \in \{10, 20\}$.
It is clear that $\boldsymbol{X}'\boldsymbol{X}$ will become more ill-conditioned as $\sigma_1$ gets
smaller. $R^2$ stands for the coefficient of determination which is
here:

$$R^2 = 1 - \frac{\mathsf{Var}(\varepsilon_1)}{\mathsf{Var}(\boldsymbol{X}^3)} \qquad (4)$$

Context
Proposed Models
Structure estimation
Results
Missing values
Tools
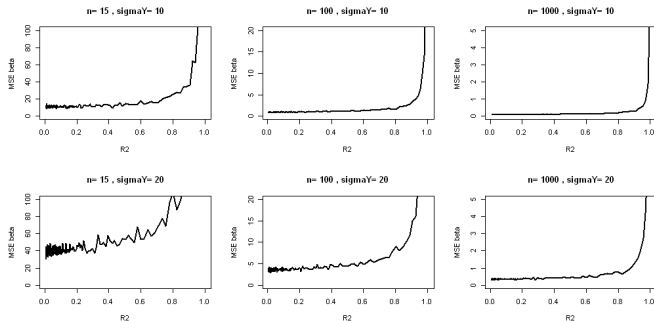
Industrial context
Statistical context

Figure: Evolution of observed Mean Squared error on $\hat{\beta}_{OLS}$ with the strength of the correlations for various sample sizes and strength of regression. $d = 5$ covariates (running example).

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Industrial context
Statistical context

## Ridge Regression

[Hoerl and Kennard, 1970, Marquardt and Snee, 1975] proposes a possibly biased estimator for $\beta$ that can be written in terms of a parametric $L_2$ penalty:

$$\hat{\beta} = \text{argmin}_\beta \left\{ \| \boldsymbol{Y} - \boldsymbol{X}\beta \|_2^2 \right\} \text{ subject to } \| \beta \|_2^2 \leq \eta \text{ with } \eta > 0 \tag{5}$$

But this penalty is not guided by correlations. The solution of the ridge regression is given by

$$\hat{\beta} = \left( \boldsymbol{X}'\boldsymbol{X} - \lambda \boldsymbol{I}_n \right)^{-1} \boldsymbol{X}'\boldsymbol{Y} \tag{6}$$

Methods do exist to automatically choose a good value for $\lambda$ [Cule and De Iorio, 2013, Er et al., 2013] and a R package called `ridge` is on CRAN [Cule, 2014].

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

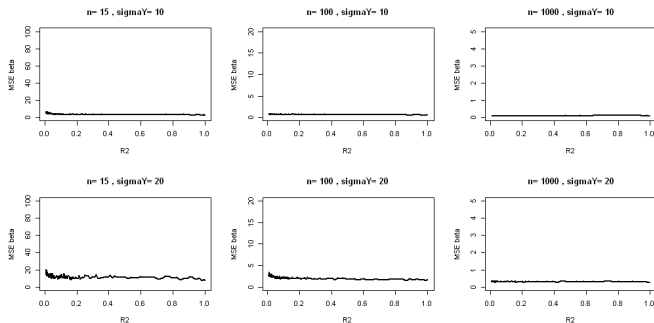Industrial context
Statistical context

# Ridge Regression



Figure: Evolution of observed Mean Squared error on $\hat{\beta}_{ridge}$ with the strength of the correlations for various sample sizes and strength of regression. $d = 5$ covariates.

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Industrial context
**Statistical context**

# LASSO

The Least Absolute Shrinkage and Selection Operator (LASSO, [Tibshirani, 1996] and [Tibshirani et al., ]) consists in a shrinkage of the regression coefficients based on a $\lambda$ parametric $L_1$ penalty to obtain zeros in $\hat{\beta}$ instead of the $L_2$ penalty of the ridge regression:

$$\hat{\boldsymbol{\beta}} = \text{argmin}\left\{\| \boldsymbol{Y} - \boldsymbol{X}\beta \|_2^2\right\} \text{ subject to } \| \beta \|_1 \leq \lambda \text{ with } \lambda > 0.$$

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Industrial context
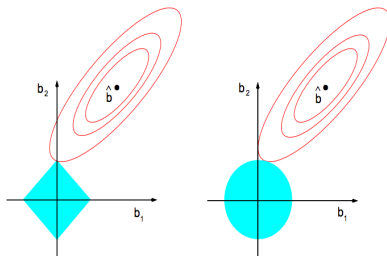Statistical context

# LASSO



Figure: Geometric view of the Penalty for the LASSO (left) compared to ridge regression (right) as shown in the book from Hastie [Hastie et al., 2009]

Figure shows the contour of error (red) and constraint function (blue). The axis stands for the regression coefficients.

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Industrial context
**Statistical context**

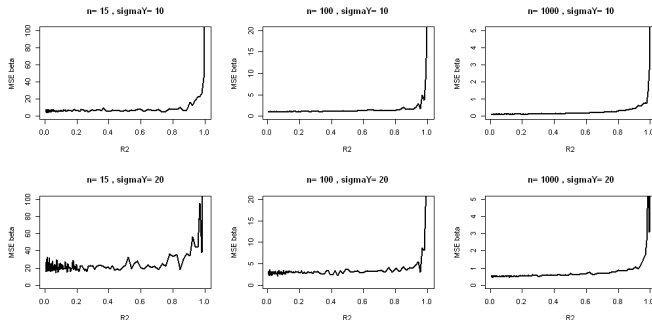Figure: Evolution of observed Mean Squared error on $\hat{\beta}_{lar}$ with the strength of the correlations for various sample sizes and strength of regression. $d = 5$ covariates.

lars package on CRAN ([Hastie and Efron, 2013]).

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Industrial context
Statistical context

# SEM

Modélisation de la structure mais à la main et aucun impact sur l'estimation

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Industrial context
**Statistical context**

## Selvarclust

Semble très bien mais n'aboutit pas vers la régression donc on le prolonge en CorReg

Context
**Proposed Models**
Structure estimation
Results
Missing values
Tools

Marginal model
Plug-in model

## Hypothesis 1

There are $d_r \geq 0$ "sub-regressions", each sub-regression $j = 1, \ldots, d_r$ having the covariate $\boldsymbol{X}^{J_r^j}$ as *response* variable ($J_r^j \in \{1, \ldots, p\}$ and $J_r^j \neq J_r^{j'}$ if $j \neq j'$) and having the $d_p^j > 0$ covariates $\boldsymbol{X}^{J_p^j}$ as *predictor* variables ($J_p^j \subset \{1, \ldots, d\} \backslash J_r^j$ and $d_p^j = |J_p^j|$ the cardinal of $J_p^j$):

$$\boldsymbol{X}^{J_r^j} = \boldsymbol{X}^{J_p^j} \boldsymbol{\alpha}_j + \varepsilon_j, \tag{7}$$

where $\boldsymbol{\alpha}_j \in \mathbb{R}^{d_p^j}$ ($\alpha_j^h \neq 0$ for all $j = 1, \ldots, d_r$ and $h = 1, \ldots, d_p^j$) and $\varepsilon_j \sim \mathcal{N}_n(\boldsymbol{0}, \sigma_j^2 \boldsymbol{I})$.

Context
**Proposed Models**
Structure estimation
Results
Missing values
Tools

Marginal model
Plug-in model

## Hypothesis 2

the response covariates and the predictor covariates are totally disjoint: for any sub-regression $j = 1, \ldots, d_r$, $J_p^j \subset J_f$ where $J_r = \{J_r^1, \ldots, J_r^{d_r}\}$ is set of all response covariates and $J_f = \{1, \ldots, d\} \backslash J_r$ is the set of all *non* response covariates of cardinal $d_f = d - d_r = |J_f|$. We call this hypothesis the underline{uncrossing rule}. Then:

$$\boldsymbol{Y} = \boldsymbol{X}_f \beta_f + \boldsymbol{X}_r \beta_r + \varepsilon_Y. \tag{8}$$

Context
**Proposed Models**
Structure estimation
Results
Missing values
Tools

Marginal model
Plug-in model

## Hypotheses 3

We assume that all errors $\varepsilon_Y$ and $\varepsilon_j$ ($j = 1, \ldots, d_r$) are *mutually independent*. It implies in particular that conditional response covariates $\{\boldsymbol{X}^{j_r^j} | \boldsymbol{X}^{j_p^j}, \boldsymbol{S}; \boldsymbol{\alpha}_j, \sigma_j^2\}$ are *mutually independent*:

$$\mathbb{P}(\boldsymbol{X}_r | \boldsymbol{X}_f, \boldsymbol{S}; \boldsymbol{\alpha}, \boldsymbol{\sigma}^2) = \prod_{j=1}^{d_r} \mathbb{P}(\boldsymbol{X}^{j^j} | \boldsymbol{X}^{j_p^j}, \boldsymbol{S}; \boldsymbol{\alpha}_j, \sigma_j^2). \qquad (9)$$

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Marginal model
Plug-in model

## Marginal model

We obtain for the distribution of $\{\boldsymbol{Y}|\boldsymbol{X}_f, \boldsymbol{S}; \beta, \boldsymbol{\alpha}, \sigma_Y^2, \boldsymbol{\sigma}^2\}$:

$$
\begin{aligned}
\boldsymbol{Y} &= \boldsymbol{X}_f(\beta_f + \sum_{j=1}^{d_r} \beta_{j_r^j}\boldsymbol{\alpha}_j^*) + \sum_{j=1}^{d_r} \beta_{j_r^j}\boldsymbol{\varepsilon}_j + \boldsymbol{\varepsilon}_Y \qquad (10) \\
&= \boldsymbol{X}_f\beta_f^* + \boldsymbol{\varepsilon}_Y^*, \qquad (11)
\end{aligned}
$$

where $\boldsymbol{\alpha}_j^* \in \mathbb{R}^{d_f}$ with $(\boldsymbol{\alpha}_j^*)_{J_p^j} = \boldsymbol{\alpha}_j$ and $(\boldsymbol{\alpha}_j^*)_{J_f \setminus J_p^j} = \boldsymbol{0}$. We define $\boldsymbol{\alpha}^* \in \mathbb{R}^{(d_f \times d_r)}$ to use more compact notations:

$$
\begin{aligned}
\boldsymbol{X}_r &= \boldsymbol{X}_f\boldsymbol{\alpha}^* + \boldsymbol{\varepsilon} \\
\boldsymbol{Y} &= \boldsymbol{X}_f(\beta_f + \boldsymbol{\alpha}^*\beta_r) + \boldsymbol{\varepsilon}\beta_r + \boldsymbol{\varepsilon}_Y \qquad (12)
\end{aligned}
$$

Where $\boldsymbol{\varepsilon}$ is the $n \times d_r$ matrix whose columns are the $\boldsymbol{\varepsilon}_j$, the noises of the sub-regressions.

Context
**Proposed Models**
Structure estimation
Results
Missing values
Tools

Marginal model
**Plug-in model**

## Plug-in model

$$\varepsilon_Y^* = \varepsilon \beta_r + \varepsilon_Y. \tag{13}$$

Then the Best Linear Unbiased Estimator ($\mathrm{BLUE}$) for $\beta_r$ is given ($\mathrm{MLE}$ estimator) by:

$$\hat{\beta}_r = (\varepsilon'\varepsilon)^{-1}\varepsilon'\varepsilon_Y^*. \tag{14}$$

And we have the following estimators:

$$
\begin{aligned}
\hat{\varepsilon} &= \boldsymbol{X}_r - \boldsymbol{X}_f\hat{\boldsymbol{\alpha}}^* \text{ and} \\
\hat{\varepsilon}_Y^* &= \boldsymbol{Y} - \boldsymbol{X}_f\hat{\beta}_f^*
\end{aligned}
$$

that we can use by plug-in.

Context
**Proposed Models**
Structure estimation
Results
Missing values
Tools

Marginal model
**Plug-in model**

## Plug-in model

$$\hat{\beta}_r^{\varepsilon} = (\hat{\varepsilon}'\hat{\varepsilon})^{-1}\hat{\varepsilon}'\hat{\varepsilon}_Y^*$$

that depends on all covariates in $\boldsymbol{X}$ and relies on the estimated coefficients of sub-regressions $\hat{\boldsymbol{\alpha}}^*$ and on the estimate $\hat{\boldsymbol{\beta}}_f^*$ of the coefficients in the marginal model. Then we can estimate $\boldsymbol{Y}$ by:

$$\hat{\boldsymbol{Y}}_{plug-in} = \boldsymbol{X}_f\hat{\boldsymbol{\beta}}_f^* + \hat{\varepsilon}\hat{\beta}_r^{\varepsilon}. \tag{15}$$

We can improve estimation of $\boldsymbol{\beta}_f$ (in terms of bias) by doing an additional identification step. We know that $\boldsymbol{\beta}_f^* = \boldsymbol{\beta}_f + \boldsymbol{\alpha}^*\boldsymbol{\beta}_r$ so we naturally define the following estimator:

$$\hat{\boldsymbol{\beta}}_f^{\varepsilon} = \hat{\boldsymbol{\beta}}_f^* - \hat{\boldsymbol{\alpha}}^*\hat{\boldsymbol{\beta}}_r^{\varepsilon}.$$

# Marginal properties

biased

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Marginal model
Plug-in model

# Plug-in properties

asymptotically unbiased

Context
**Proposed Models**
Structure estimation
Results
Missing values
Tools

Marginal model
**Plug-in model**

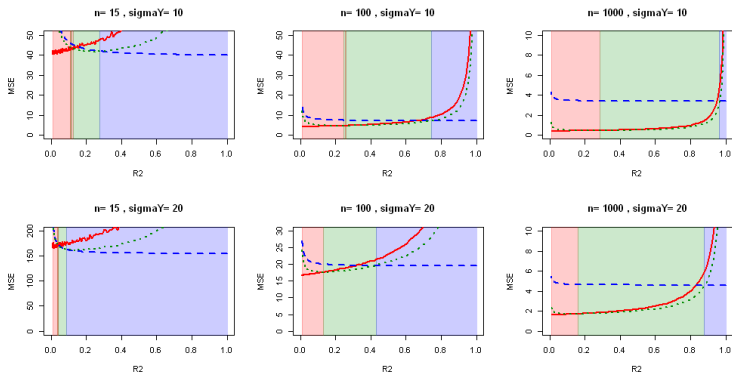Figure: MSE on $\hat{\boldsymbol{\beta}}$ of OLS (plain red) and CorReg marginal (blue dashed) and CorReg plug-in (green dotted) estimators for varying $R^2$ of the sub-regression, $n$ and $\sigma_Y$. Results obtained on the running example with $d = 5$ covariates.

Context
Proposed Models
Structure estimation
Results
Missing values
Tools

Marginal model
Plug-in model

## Lasso Consistency

Consistency issues of the LASSO are well known and Zhao
[Zhao and Yu, 2006] gives a very simple example to illustrate it.
We have taken the same example to show how our method is
better to find the true relevant covariates. Here $d = 3$ and
$n = 1\,000$.
We define $\boldsymbol{X}^1, \boldsymbol{X}^2, \varepsilon_Y, \varepsilon_1 \quad i.i.d. \sim \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_n)$ and then

$$\begin{aligned} \boldsymbol{X}^3 &= \frac{2}{3}\boldsymbol{X}^1 + \frac{2}{3}\boldsymbol{X}^2 + \frac{1}{3}\varepsilon_1 \text{ and} \\ \boldsymbol{Y} &= 2\boldsymbol{X}^1 + 3\boldsymbol{X}^2 + \varepsilon_Y. \end{aligned}$$

Context
**Proposed Models**
Structure estimation
Results
Missing values
Tools

Marginal model
Plug-in model

## Lasso Consistency

True **S** was found 991 times on 1 000 tries.

|        | Classical LASSO | CorReg marginal + LASSO | CorReg full plug-in + LASSO |
|--------|-----------------|--------------------------|------------------------------|
| True **S** | 1.003303 (0.046) | **1.002273** (0.046) | **1.002812** (0.046) |
| **Ŝ**  | 1.003303 (0.046) | 1.017622 (0.17) | **1.002812** (0.046) |

Table: MSE observed on a validation sample (1 000 individuals) and their standard deviation (between brackets).

We look at the consistency that is the real stake:

|        | Classical LASSO | CorReg marginal + LASSO | CorReg full plug-in + LASSO |
|--------|-----------------|--------------------------|------------------------------|
| True **S** | 0 | **1000** | 835 |
| **Ŝ**  | 0 | **991** | 829 |

Table: Number of consistent models found on 1 000 tries.

Context
Proposed Models
**Structure estimation**
Results
Missing values
Tools

•

Context
Proposed Models
Structure estimation
**Results**
Missing values
Tools

Simulation results
Industrial results

•

-

-

Modèle génératif complet avec dépendances

Context
Proposed Models
Structure estimation
Results
**Missing values**
Tools

•

Explosion des mélanges

Context
Proposed Models
Structure estimation
Results
**Missing values**
Tools

- 

SEM avec Gibbs

Context
Proposed Models
Structure estimation
Results
**Missing values**
Tools

- 

Bic pondéré

Context
Proposed Models
Structure estimation
Results
**Missing values**
Tools

- 

Résultats pourris

Context
Proposed Models
Structure estimation
Results
Missing values
**Tools**

- 

Excel, fonctions graphiques, arbres de décision

Context
Proposed Models
Structure estimation
Results
Missing values
**Tools**

📄 Cule, E. (2014).
*ridge: Ridge Regression with automatic selection of the penalty parameter*.
R package version 2.1-3.

📄 Cule, E. and De Iorio, M. (2013).
Ridge regression in prediction problems: automatic choice of the ridge parameter.
*Genetic epidemiology*, 37(7):704–714.

📄 Dodge, Y. and Rousson, V. (2004).
Analyse de régression appliquée: manuel et exercices corrigés (coll. eco sup, ).
*Recherche*, 67:02.

📄 Er, M. J., Shao, Z., and Wang, N. (2013).
A systematic method to guide the choice of ridge parameter in