

CorReg : Regression for correlated variables and application in steel industry

Clément THERY

February 19, 2014

Abstract. Linear regression generally suppose to have decorrelated covariates. This hypothesis is often irrealist with industrial datasets that contains many highly correlated covariates due to the process, physcial laws, *etc.* The proposed generative model consists in explicit modeling of the correlations with a family of linear regressions between the covariates permitting to obtain by marginalization a parsimonious correlation-free regression model, easily understandable and compatible with variable selection methods. The structure of correlations is found with an MCMC algorithm. An R package (CORREG) available on the CRAN implements this new method which will be illustrated on both simulated datasets and real-life datasets from steel industry.

Keywords. Regression, correlations, industry, variable selection, generative models

1 Introduction

When one wants to explain a phenomenon based on some covariates, the first statistical method tried frequently is the linear regression. It provides a predictive model with a good interpretability and is simple to learn for non-statistician. Therefore, linear regression is used in nearly all the fields where statistics are made, from industry (ballistic models to calibrate the process) to sociology (predicting some numerical properties of a population). Linear regression is a very classic situation that faces an also classical problem : the variance of the estimators. This variance increases based on two aspects :

- The dimension p (number of covariates) of the model : the more covariates you have the greater variance you get.
- The correlations within the covariates : strongly correlated covariates give bad-conditioning and increase variance of the estimators .

With the rise of informatic, datasets contains more and more covariates and thus more and more useless covariates. So dimension reduction becomes a necessity. Moreover, when you use more covariates, you increase the chance to have correlated ones. For example, this work takes place in an industrial context with a big set of covariates (many parameters of the whole process without any a priori) highly correlated (physical laws, process rules, etc). In such a context, variance of the estimators can lead to arbitrary results or even no results at all. Prediction and interpretation are both strongly needed, with a preference for interpretation in industrial context (better to improve the process when possible than to only predict defects).

When estimating the parameters of the regression we have to compute the inverse of a matrix[11] which will be ill-conditioned or even singular if some covariates depend linearly from each other. For a model defined by

$$Y = X\beta + \varepsilon \tag{1}$$

where X is the $n \times p$ matrix of the explicative variables, Y the response vector and $\varepsilon \sim \mathcal{N}(0, \sigma_Y^2)$ we have the following Ordinary Least Squares (OLS) estimators :

$$\hat{\beta} = (X'X)^{-1} X'Y \quad (2)$$

With variance

$$\text{Var}(\hat{\beta}) = \sigma_Y^2 (X'X)^{-1} \quad (3)$$

And when correlations between covariates are strong, the matrix to invert is ill-conditioned and the variance explodes.

Because it is the minimum-variance unbiased estimator, penalized methods try to reduce the variance introducing some bias to improve the bias-variance trade-off and get better prediction. Ridge regression[9] proposes a biased estimator :

$$\hat{\beta} = (X'X + kI)^{-1} X'Y \text{ with } k \geq 0 \quad (4)$$

But Ridge regression is not efficient to select covariates (it's an assumed choice) because coefficients tends to 0 but don't reach 0. So it gives difficult interpretations and is not adapted for our industrial context. We need to reduce the dimension of the model. Our goal is not just to predict but also to understand the response variable.

Real datasets implies many irrelevant variables (datasets based on the whole process without any a priori) so we have to use variable selection methods.

We note classical norms: $\|\beta\|_2^2 = \sum_{i=1}^p (\beta_i)^2$ and $\|\beta\|_1 = \sum_{i=1}^p |\beta_i|$

The Least Absolute Shrinkage and Selection Operator (LASSO)[12] consists in a shrinkage of the regression coefficients based on a λ parametric L_1 penalty.

$$\hat{\beta} = \text{argmin} \left\{ \|Y - X\beta\|_2^2 \right\} \text{ subject to } \|\beta\|_1 \leq \lambda \quad (5)$$

The Least Angle Regression[2] (LAR) Algorithm offers a very efficient way to obtain the whole LASSO path and is very attractive. It requires only the same order of magnitude of computational effort as OLS applied to the full set of covariates. And it really selects covariates with coefficients set exactly to 0. But LASSO also faces consistency problems[17] when confronted with correlated covariates. This point will be developed further (see 4.5) with numerical results. Another limitation of the LASSO is that it preserves at most n predictors (troublesome when in high dimension).

Elastic net[18] is a method developed to be a compromise between Ridge regression and the LASSO. Elastic net can be written:

$$\hat{\beta} = (1 + \lambda_2) \text{argmin} \left\{ \|Y - X\beta\|_2^2 \right\}, \text{ subject to } (1 - \alpha) \|\beta\|_1 + \alpha \|\beta\|_2^2 \leq t \text{ for some } t \quad (6)$$

where $\alpha = \frac{\lambda_2}{(\lambda_1 + \lambda_2)}$ It seems to give good predictions. But it is based on the grouping effect and if the dataset contains two identical variables they will obtain the same coefficient whereas LASSO will choose between one of them and will then obtain same predictions with a more parsimonious model.

The CLusterwise Effect REgression[14] (CLERE) tries to reduce the dimension by considering the β_j no longer as fixed effect parameters but as unobserved independant random variables with β following a Gaussian Mixture distribution.

$$\beta_j | \mathbf{z}_j \sim \mathcal{N} \left(\sum_{k=1}^g b_k z_{jk}, \gamma^2 \right) \text{ with} \quad (7)$$

$$\mathbf{z}_j = (z_{j1}, \dots, z_{jg}) \sim \mathcal{M}(\pi_1, \dots, \pi_g) \text{ and} \quad (8)$$

$$\forall k = 1, \dots, g \quad \sum_{j=1}^p z_{jk} \geq 1 \quad (9)$$

The idea is to hope that this mixture will have few enough components to have a number of parameters to estimate significantly lower than p . In such a case, it improves interpretability and ability to yield reliable prediction with a smaller variance on $\hat{\beta}$. But we need to suppose having many

covariates with the same level of effect on the response variable and seems to be less efficient in prediction than elastic net.

When some try to reduce the dimension and then just hope to have small correlations in the remaining dimensions, we propose to focus on the correlations, giving a model with orthogonal covariates. In fact we search the greatest set of orthogonal covariates to keep the maximum but with an orthogonality constraint. This can be viewed as a pretreatment on the dataset allowing to use then dimension reduction tools without suffering from correlations. We only consider strong correlations (i.e. : problematic ones) thus we keep most of the information contained in the dataset. We will in a second time be able to use the remaining part of the information (sequential approach).

Our work is based on the assumption that if we know that correlations are a problem and if we precisely know the correlations, we could use this knowledge to avoid the problem. The idea is to suppose explicitly a linear structure between the covariates. It gives a recursive Simultaneous Equation Model (SEM)[1]. That can be viewed as a system of linear regression.

$$Y = X\beta_Y + \varepsilon_Y \quad (10)$$

$$X = XB + \varepsilon_X \quad (11)$$

Recursive sem don't really have a specific estimator because general system estimators (Seemingly Unrelated Regression (SUR)[15] and Two-Stage Least Squares (2SLS)) are equivalent to independent Ordinary Least Squares when applied to recursive SEM [13]. Our work can be viewed as a new way of estimating recursive SEM based on their own structure. In this work, we decide to distinguish the response variable from the other variables that are on the left of a regression. Thus we don't have a system of regressions but one regression on our response variable and a system of subregressions (without the response variable). The structure is supposed to be the source of the correlations and allows us to define a reduced set of independent covariates. Thus we reduce dimension and correlations in the same time. The structure justifies the eviction of the redundant covariates without significant information loss. It can be seen as a pretreatment on the dataset based on the hypothesis of a strong structure between the covariates (i.e. : small ε_X).

We can use any variable selection method on the reduced dataset with improved efficiency (reduced variance) due to dimension reduction and correlation suppression. So we obtain two kinds of zeros in our first model : coerced zeros due to correlations (redundant information) and estimated ones with classical variable selection methods applied on remaining variables. This two kinds of zero won't be interpreted in the same way and thus consistency issues don't mean interpretation issues any more. So we dodge the drawbacks of both grouping effect and variable selection.

But to work, we need an explicit structure between the covariates. SEM are often used in social sciences and economy where a structure is supposed "by hand" but here we want to be able to find a structure without any a priori (possibility to include some known structure remains). Graphical LASSO [4] offers a method to obtain a structure based on the precision matrix (inverse of the variance-covariance matrix). It consists in a selection in the precision matrix, setting some covariances to zero. But the resulting matrix is symmetric and we need an oriented graph for our SEM. So we developed an MCMC algorithm to find it (R package CorReg on CRAN). However, Graphical LASSO can be used in the initialization step of our MCMC. This structure is based on gaussian mixture models to fit better real datasets and to allow identifiability of the structure in terms of complexity (number of parameters).

This paper will first present the reduced model and its properties before describing in part 3 the algorithm used to find the structure. We will then look at some numerical results on simulated (part 4) and real industrial datasets (part 5) before concluding and giving some perspectives in the sixth part.

2 Model to decorrelate the covariates

2.1 The generative model

We define explicitly the correlations within a set $X \in \mathcal{R}^{n \times p}$ of p correlated covariates with a family of p_2 internal regressions between covariates with I_2 the set of indices of endogenous variables in X (explained ones) and $I_1 = \{I_1^1, \dots, I_1^{p_1}\}$ the set of the sets of indices of exogenous covariates (explaining ones) with $\forall j \notin I_2, I_1^j = \emptyset$. Then we have an explicit structure $S = (I_1, I_2, p_1, p_2)$ where $p_1 = (p_1^1, \dots, p_1^{p_2})$ is the vector of the number of covariates in each internal regression. Thus we have $p_2 = \#I_2$ and $p_1^j = \#I_1^j$.

In the following, we note X^j the j^{th} column of a matrix X . For lighter notation we define $X_2 = X^{I_2}$ the matrix of the endogenous covariates and $X_1 = X \setminus X_2$ the matrix of the remaining exogenous covariates. We make the hypothesis of the uncrossing rule $I_1 \cap I_2 = \emptyset$ *i.e.* endogenous variables don't explain other covariates.

Let $Y \in \mathcal{R}^n$ be the response variable. We can now write the generative model:

- Main regression between Y and X :

$$Y_{|S} = XA + \varepsilon_Y = X_1A_1 + X_2A_2 + \varepsilon_Y \text{ with } \varepsilon_Y \sim \mathcal{N}(0, \sigma_Y^2 I_n); \quad (12)$$

where $A = (A_1, A_2) \in \mathcal{R}^p$ is the vector of the regression coefficients and I_n the identity matrix,

- Family of p_2 regressions within correlated covariates in X :

$$\forall j \in I_2 : X_{|X_1, S}^j = X_1 B_1^j + \varepsilon_j \text{ with } \varepsilon_j \sim (0, \sigma_j^2 I_n); \quad (13)$$

where $B_1^j \in \mathcal{R}^{(p-p_2)}$ are the vectors of the regression coefficients between the covariates (containing some zeros according to I_1^j).

- $p - p_2$ remaining independent exogenous Gaussian mixtures:

$$\forall j \notin I_2 : X_{|S}^j \sim f(\theta_j) = \mathcal{GM}(\pi_j; \mu_j; \sigma_j^2) \text{ with } \pi_j, \mu_j, \sigma_j^2 \text{ vectors of size } k_j; \quad (14)$$

This generative model is conditionnal to S , the discrete structure model that is identifiable because we can't permute some regressions in (13) and obtain the same joint distribution $P(X, Y)$, the residuals (ε_j) would not stay Gaussian. These residuals are in the following supposed independent but one can suppose dependencies between them and then use appropriate tools to estimate them and the B_1^j like SUR with Feasible Generalized Least Squares (FGLS) by Zellner [15].

If there are exact regressions in (13), the structure won't be identifiable. But it only means that several structure will have the same likelihood and they will have the same interpretation. So it's not really a problem. Moreover, when an exact subregression is found, we can delete one of the implied variables without any loss of information and the structure will define a list of variable from which to delete. CorReg (Our R package) prints a warning to point out exact regressions when found.

To better fit industrial variables (Figure 1), we suppose in (14) that variables in X_1 follow Gaussian mixtures. The great flexibility [10] of such models makes our model more robust. But one can use other laws if needed. Gaussian case is just a special case ($k_j = 1$) of Gaussian mixture so it is included in our hypothesis.

2.2 Estimator and properties

We note that (12) and (13) also give by simple integration on X_2 a regression model on Y *depending only on uncorrelated covariates* X_1 :

$$Y = X_1(A_1 + \sum_{j \in I_2} B_1^j A_j) + \sum_{j \in I_2} \varepsilon_j A_j + \varepsilon_Y = X_1 \alpha_1 + \varepsilon_\alpha \quad (15)$$

So we have the unbiased estimator:

$$\hat{\alpha}_1 = (X_1' X_1)^{-1} X_1' Y \text{ and } \hat{\alpha}_2 = 0 \quad (16)$$

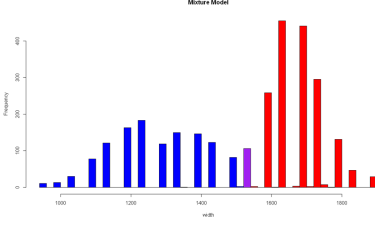


Figure 1: Exemple of non-Gaussian real variable easily modeled by a Gaussian mixture

with variance:

$$\text{Var}[\hat{\alpha}_1|X] = (\sigma_Y^2 + \sum_{j \in I_2} \sigma_j^2 A_j^2)(X_1' X_1)^{-1} \text{ and } \text{Var}[\hat{\alpha}_2|X] = 0 \quad (17)$$

We see that the variance is reduced (no correlations and smaller matrix give better conditioning) for small values of σ_j *i.e.* strong correlations. Moreover, the explicit structure is parsimonious and simply consists in linear regressions and thus is easily understood by non statistician, allowing them to have a better knowledge of the phenomenon inside the dataset. Expert knowledge can even be added to the structure. This new model is reduced even without variable selection and is just a linear regression so every method for variable selection in linear regression can be used. Hence we hope to obtain a parsimonious model with *two kinds of zeros*: those from decorrelating step and those from selection step, with different meanings.

The CORREG package proposes to make selection with both LASSO (with LAR), elasticnet, adaptative lasso, and others more classical methods. A last Ordinary Least Square step is added to improve estimation because penalisation is made to select variables but also shrinks remaining coefficients. This last step allows to keep the benefits of shrinkage (variable selection) without any impact on remaining coefficients (see [16]).

2.3 Why grouping effect is misleading

In industrial context, when a model explain why things go wrong, one will try to fix the problem. If $X_1 = X_2 + e$ and we have the grouping effect, we will obtain a model like $Y = aX_1 + aX_2$. Then when one will try to modify Y he will modify one of the covariates and both will change so he won't get expected results. Nothing constrains us to give only one equation. It is clearly better to give the user another equation (or system for more complex models) describing the correlations. So you get the following model : $Y = aX_1 + aX_2$ and $X_1 = X_2 + e$. So you have more information and are able to decide better actions. With such a model, grouping effect is no more useful because when saying $Y = 2aX_2$ and $X_1 = X_2 + e$ you clearly show that X_2 is correlated with both Y and X_1 . So it is possible to combine the advantages of grouping effect and selection just giving several equations. Each equation here is very simple so you don't really increase complexity of the model. Uncrossed model ($I_1 \cap I_2 = \emptyset$) guarantee to keep a simple structure easily interpretable.

3 Estimating subregressions with Markov chain

3.1 Structure comparison

All our work is based on S , the linear structure between the covariates. Our generative model allows us to compare structures with criterions like the Bayesian Information Criterion (BIC) which penalize the log-likelihood according to the complexity of the structure [8]. But BIC tends to give too complex structures because we test a great range of models. Thus CORREG allows to choose to penalise the complexity a bit more with a hierarchical uniform *a priori* law $P(S) = P(I_1|p_1, I_2, p_2)P(p_1|I_2, p_2)P(I_2|p_2)P(p_2)$ instead of a simple uniform law on S . Thus we have :

$$BIC^*(X|S) = BIC(X|S) + \ln(P(S)) \quad (18)$$

It increases penalty on complexity for $p_2 < \frac{p}{2}$. Hence when using *BIC** this hypothesis is made and will become a constraint in the MCMC. But we can imagine to use other criterions, like the *RIC* (Risk Inflation Criterion [3]) that choose a penalty in $\log p$ instead of $\log n$ and thus gives more parsimonious models when p is larger than n (high dimension) or any other criterion [5] thought to be better in a given context. Variables in X_1 are supposed to be independent. Thus if one have some hypothesis on the distribution of some variables (exponentially distributed for example) he can compute corresponding *BIC* separately, give it as an input of CORREG and then improve the efficiency of the algorithm (it will find a structure only if it is really relevant).

We will now use the following notation: $\psi(S) = \psi(X|S)$ where $\psi(X|S)$ is the chosen criterion that could be *BIC*, *BIC**, *RIC*, *etc.*

3.2 The neighbourhood

Let's define \mathcal{S} the ensemble of feasible structures (those with $I_1 \cap I_2 = \emptyset$).

For each step, starting from $S \in \mathcal{S}$ we define a neighbourhood $\mathcal{V}_{S,\phi}$ with $\phi \sim \mathcal{U}(\{1, \dots, p\})$ and for $\phi = j$:

$$\mathcal{V}_{S,j} = \{S^{(i,j)} | 1 \leq i \leq p\} \cup \{S\} \quad (19)$$

With $S^{(i,j)}$ defined by the following algorithm :

- if $i \notin I_1^j$ (add):
 - $I_1^j = I_1^j \cup \{i\}$
 - $I_1^i = \emptyset$ (explicative variables can't depend on others : column-wise relaxation)
 - $I_1 = I_1 \setminus \{j\}$ (dependent variables can't explain others : row-wise relaxation)
- else (remove): $I_1^j = I_1^j \setminus \{i\}$

At every moment, coherence between I_1 and others parts of S can be done by $I_2 = \{j | \#I_1^j > 0\}$, $p_2 = \#I_2$, $\forall 1 \leq j \leq p : p_1^j = \#I_1^j$.

We have here $\#\mathcal{V}_{S,j} = p$ but some other constraints can be added on the definition of \mathcal{S} and will consequently modify the size of the neighbourhood (for example a maximum complexity for the internal regressions or the whole structure, a maximum number of internal regressions, *etc.*). CORREG allows to modify this neighbourhood to better fit users constraints. Relaxation (column-wise and row-wise) is optional but gives more stability to the number of feasible candidates at each step and allows to modify several parts of I_1 in only one step when needed. Hence it improves efficiency by a significant reinforcement of the irreducibility of the Markov chain.

3.2.1 The walk

We first make the approximation

$$P(S|X) \approx \exp(\psi(S)). \quad (20)$$

The algorithm follows a time-homogeneous markov chain whose transition matrix \mathcal{P} has $\#\mathcal{S}$ rows and columns (combinatory so we'll just compute the probabilities when we need them). At each step the markov chain moves with probabiliy:

$$\forall (S, \tilde{S}) \in \mathcal{S}^2 : \mathcal{P}(S, \tilde{S}) = \mathbf{1}_{\{\tilde{S} \in \mathcal{V}_{S,\phi}\}} \frac{\exp(-\frac{1}{2}\psi(\tilde{S}))}{\sum_{S_l \in \mathcal{V}_S} \exp(-\frac{1}{2}\psi(S_l))} \quad (21)$$

$$= \frac{1}{p} \quad (22)$$

And \mathcal{S} is a finite state space.

And then we can note $\forall (S, \tilde{S}) \in \mathcal{S}^2$:

$$\mathcal{P}(S, \tilde{S}) = \frac{1}{p} \sum_{j=1}^p q(\tilde{S}, \mathcal{V}_{S,j})$$

The output will be the best structure seen in terms of BIC. If we have some knowledge about some sub-regressions (physical models for example) we can add them in the found structure. So the model is really expert-friendly. Because the walk follows a regular and thus ergodic markov chain with a finite state space, it has exactly one stationary distribution [7] : π and every rows of $\lim_{k \rightarrow \infty} \mathcal{P}^k = W$ equals π .

With $\forall S \in \mathcal{S}$:

$$\begin{aligned} 0 &\leq \pi(S) \leq 1 \\ \sum_{S \in \mathcal{S}} \pi(S) &= 1 \\ \pi(S) &= \sum_{\tilde{S} \in \mathcal{S}} \pi(\tilde{S}) \mathcal{P}(\tilde{S}, S) \end{aligned} \tag{23}$$

$$\tag{24}$$

3.2.2 initialisation

The initial structure can be based on a first warming algorithm taking the correlations into account. coefficients are randomly placed into I_1 , weighted by the absolute value of the correlations. Then this structure can be reduced by the hadamard product with the binary matrix obtained by Graphical Lasso[4].

3.2.3 reduced neighbourhood

If the algorithm did not have time to converge, it can be finished with a few step for which the neighborhood would only contain smaller candidates. It is equivalent to ask for each element in the structure if the criterion would be better without it. Thus it can be seen as a final cleaning step. But in fact, it's just continuing the MCMC with a reduced neighborhood.

3.2.4 multiple try

One would rather test multiple short chains than lose time in initialisation or long chains [6]. It also helps to face local extrema.

4 Numerical results on simulated datasets

Here are some results on simulated datasets.

4.1 Finding the structure

4.1.1 How to evaluate found structure ?

The first criterion is BIC^* which is minimised in the MCMC. We can compare this criterion for both the true Structure and the found one. But if the structures differ, the BIC^* will also differ and comparison won't show how far the found structure is from the true one. So we need more precise criterions to compare the true structure S and the found one \tilde{S} . Global indicators :

- True left : the number of found dependent variables that really are dependent $TL = \#(I_2 \cap \tilde{I}_2)$
- Wrong left : the number of found dependent variables that are not dependent $WL = \#(\tilde{I}_2) - TL$
- Missing left : the number of really dependent variables not found $ML = \#I_2 - TL$
- Δp_2 : the gap between the number of sub-regression in both model : $\Delta p_2 = \#I_2 - \#(\tilde{I}_2)$. The sign defines if \tilde{S} is too complex or too simple
- $\Delta compl$: the difference in complexity between both model : $\Delta compl = \sum_{j \in p_2} p_1^j - \sum_{j \in \tilde{p}_2} \tilde{p}_1^j$

And we have aussi some criterion to compare between the structures :

- Mean true found R^2 (or σ^2) : the mean of R^2 (in the true model) for sub-regressions found by the algorithm and existing in the true model (comparing only the left-side variable). $MTF_{R^2} = \frac{1}{TL} \sum_{j \in I_2 \cap \tilde{I}_2} R^2(j)$.
- Mean Missing R^2 (or σ^2) : the mean of R^2 for sub-regressions only in the true structure (comparing only the left-side variable). $MM_{R^2} = \frac{1}{ML} \sum_{j \in I_2 \setminus (I_2 \cap \tilde{I}_2)} R^2(j)$.
- Mean found R^2 (or σ^2) : the mean of R^2 (in the found model) for sub-regressions found by the algorithm and existing in the true model (comparing only the left-side variable). $MF_{R^2} = \frac{1}{TL} \sum_{j \in I_2 \cap \tilde{I}_2} \tilde{R}^2(j)$.
- Mean Wrong R^2 (or σ^2) : the mean of R^2 for sub-regressions in only in the found structure (comparing only the left-side variable). $MW_{R^2} = \frac{1}{WL} \sum_{j \in \tilde{I}_2 \setminus (I_2 \cap \tilde{I}_2)} \tilde{R}^2(j)$.
- Mean R^2 : the mean of R^2 of all sub-regression for a given structure. We can then compare this value for both structures
- Mean σ^2 : the mean of σ^2 of all sub-regression for a given structure. We can then compare this value for both structures
- Complexity : global complexity of a model : $compl = \sum_{j \in p_2} p_1^j$. We can then compare both global complexities.

4.1.2 Results

Tableau de la forme :

n	time	trueBIC	BICempty	BIC_opt	True1	False1	missing1	$\Delta p2$	True_left	False_left
40	??	??	??	??	??	??	??	??	??	??
60	??	??	??	??	??	??	??	??	??	??
100	??	??	??	??	??	??	??	??	??	??

Table 1: p variables. Markov chain was XX seconds long for $n = 100$ (mean observed).

Ordre des critères de comparaison : MSE sur X, Vraigauche, fauxgauche, bics (les 3), vrais1, faux1, missing1, $\Delta p2$

L'idée serait de n'avoir qu'une seule configuration (vu qu'ici on ne dépend pas de Y) et garder la même pour tous les tableaux suivants (pour pouvoir s'appuyer sur celui-ci dans l'interprétation). Tous les tableaux seraient générés en même temps. Pour chaque base générée, on génèrerait plusieurs Y de plusieurs manières pour avoir tous les cas sur les mêmes données. La parallélisation des expériences se ferait alors sur le nombre de répliques. Les résultats seraient toujours basés sur Zchapeau (et donc Bchapeau).

On devrait constater que quand l'algo a le temps de converger, on trouve pour n petit des BICs meilleurs que le vrai modèle (d'où un bruit sur la structure). quand n augmente, ce surapprentissage devrait disparaître et on devrait donc converger vers le vrai Z (et le vrai BIC).

4.2 Y depends only on some covariates in X^{I_1}

4.2.1 without selection

On doit constater qu'on est meilleurs que OLS, que l'explicatif gagne (vrai modèle possible) mais que le prédictif reste bon. On doit aussi voir que quand n grandit OLS commence à redevenir correct.

n	OLS	(sd)	explicative	(sd)	predictive	(sd)
40	NA	NA	??	??	??	??
60	??	??	??	??	??	??
100	??	??	??	??	??	??

Table 2: Y only depends on X^{I_1} . $p = 50$ and $\text{Var}(Y) \simeq 3.10^8$

4.2.2 with selection

Avec le même Y que pour le cas sans sélection, (et les mêmes données) on teste simplement d'autres modèles :

- package lm
- modèle complet avec lasso (et LAR)
- modèle complet avec elastic net
- modèle complet ridge
- modèle explicatif elastic net
- modèle prédictif elastic net
- modèle explicatif LASSO
- modèle prédictif LASSO

Il y aurait un tableau par valeur de n pour pouvoir donner en plus des MSE des valeurs de sparsité et de validité du modèle (comparaison des positions des 0). Je n'ai pas pour l'instant de quoi utiliser CLERE mais l'article CLERE montre qu'elastic net est meilleur en prédiction donc pour les MSE ce n'est pas trop un problème.

4.3 Y depends only on some covariates in X^{I_2}

même chose qu'avant mais on part avec un handicap. les notions d'explicatif et prédictif finaux devraient alors prendre tout leur sens.

4.4 global case

Y dépend un peu de tout le monde... me semble trop compliqué car beaucoup trop de cas possibles. la conclusion étant de toute manières qu'on sera quelque part entre les deux cas précédents. Je mettrais bien des exemples simples et poussés (3 variables explicatives comme dans l'article sur la consistance du lasso) pour que les gens puissent facilement refaire le test chez eux, même sans notre package (hypothèse du vrai Z). la simplicité de l'exemple permettrait aussi de voir ce qui se passe si Z_{chapeau} est une version permutée du vrai Z . On testerait là aussi tous les modèles concurrents abordés plus haut.

Attention : on a une variabilité due à la validation croisée. Sur les mêmes données, quand on lance plusieurs fois la sélection on ne trouve pas toujours exactement les mêmes 0 (tout de même relativement stable, peut s'arranger en choisissant un meilleur K pour la validation croisée).

4.5 Consistency Issues

Consistency issues of the LASSO are well known and Zhao [17] gives a very simple example to illustrate it. We have taken the same example to show how our method is more consistent. Here $p = 3$ and $n = 1000$. We define $X_1, X_2, \varepsilon_Y, \varepsilon_X i.i.d. \sim \mathcal{N}(0, 1)$ and then $X_3 = \frac{2}{3}X_1 + \frac{2}{3}X_2 + \frac{1}{3}\varepsilon_X$ and $Y = 2X_1 + 3X_2 + \varepsilon_Y$. We compare consistencies of complete, explicative and predictive model with LASSO (and LAR) for selection. It happens that the algorithm don't find the true strucure but a permuted

one so we look at the results obtained with the true Z (but \hat{B} is used) and with the structure found by the Markov chain after a few seconds.

True Z is found 340 times on 1000 tries.

	Classical LASSO	Explicative	Predictive
True Z	1.006479	1.005468	1.006093
\hat{Z}	1.006479	1.884175	1.006517

Table 3: MSE observer on a validation sample (1000 individuals)

We observe as we hoped that explicative model is better when using true Z (coercing real zeros) and that explicative with \hat{Z} is penalized (coercing wrong coefficients to be zeros). But the main point is that the predictive model stay better than the classical one with the true Z and corrects enough the explicative model to follow the classical LASSO closely when using \hat{Z} . And when we look at the consistency :

	Classical LASSO	Explicative	Predictive
True Z	0	1000	830
\hat{Z}	0	340	621

Table 4: number of consistent model found (Y depending on X_1, X_2 and only them) on 1000 tries

299 times on 1000 tries, the predictive model using \hat{Z} is better than classical LASSO in terms of MSE and consistent (classical LASSO is never consistent).

We also made the same experiment but with X_1, X_2 (and consequently X_3) following gaussian mixtures (to improve identifiability) randomly generated by our R package. True Z is now found 714 times on 1000 tries . So it confirms that non-gaussian models are easier to identify.

	Classical LASSO	Explicative	Predictive
True Z	1.571029	1.569559	1.570801
\hat{Z}	1.005402	1.465768	1.005066

Table 5: MSE observed on a validation sample (1000 individuals)

And when we look at the consistency :

	Classical LASSO	Explicative	Predictive
True Z	0	1000	789
\hat{Z}	0	714	608

Table 6: number of consistent model found (Y depending on X_1, X_2 and only them) on 1000 tries

299 times on 1000 tries, the predictive model using \hat{Z} is better than classical LASSO in terms of MSE and consistent (classical LASSO is never consistent).

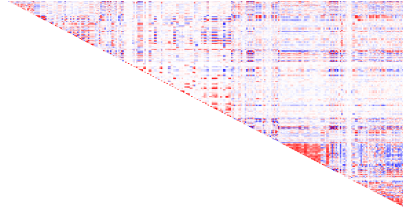
5 Numerical results on real datasets

This work takes place in steel industry context, with quality oriented objective : to understand and prevent quality problems on finished product, knowing the whole process. In particular, we focus on regression problems.

5.1 The dataset

- variables from the whole process
- The stakes : hundreds euros per ton (for information: Dunkerque site produces up to 7.5 millions tons a year)

Figure 2: Correlation matrix of the dataset ($p = 293$, $n = 3000$)



Some observed correlations with physical meaning :

- Width and Weight : 0.905
- Temperature before and after a tool : 0.983
- Roughness of both faces : 0.919
- Mean & Max of a curve : 0.911

The method was tested on 205 variables without missing values.

5.2 Results

The algorithm gives a structure with 82 subregressions with a mean of 5 regressors. Some found subregressions with physical meaning :

- $\text{Mean.weight} = f(\text{Min.weight}, \text{Max.weight}, \text{Sigma.weight})$ and other same-shaped subregressions.
- $\text{Width} = f(\text{Mean.flow}, \text{Mean.speed.CC})$
True Physical model (not linear) :
 - $\text{Width} = \text{flow} / (\text{speed} * \text{thickness})$ (thickness is constant)

Some of the other subregression represent physical models used to regulate the process and that were forgotten by the metallurgist we worked with. Found model has selected relevant variables (verified with metallurgist).

We used Elastic Net[18] on this dataset for selection (get better results than LASSO). Here are the observed MSE on a $n = 847$ validation sample. Predictive model (sequential elastic net base on estimated structure and using all the variables) is 5,82% better (Figure 5.2) than elastic net computed on the whole dataset.

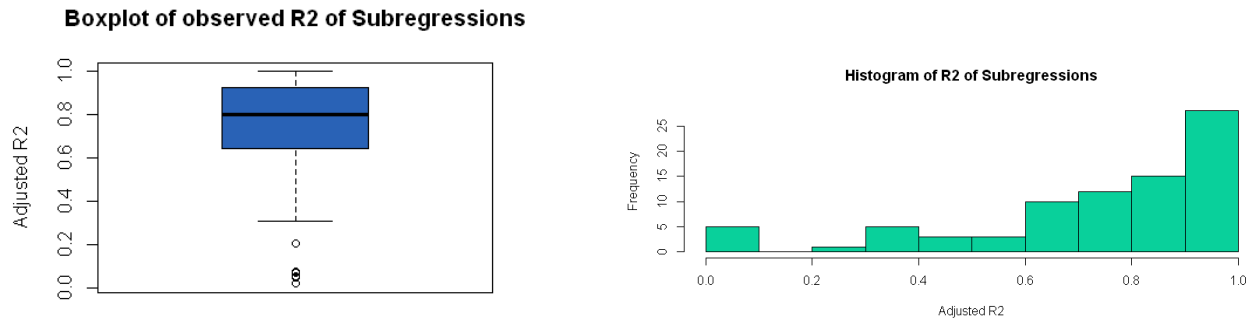


Figure 3: Adjusted R2 of found subregressions (industrial dataset)

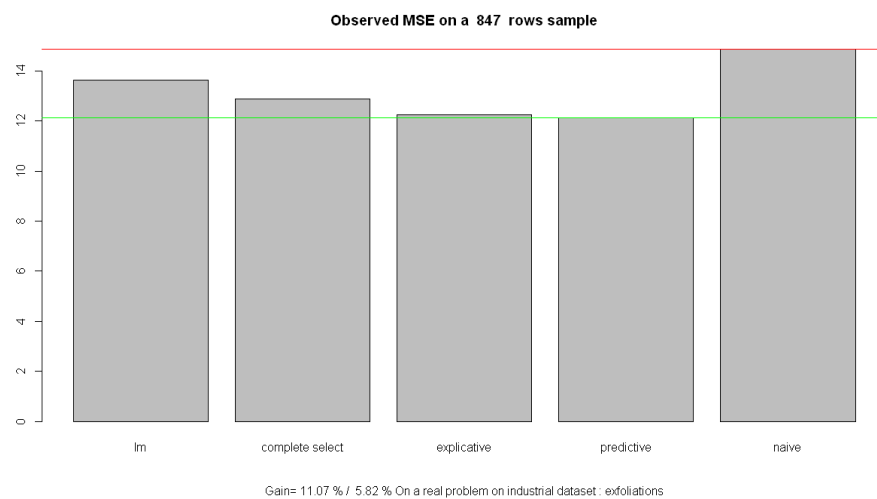


Figure 4: MSE comparison on industrial dataset

6 Conclusion

We have seen that correlations can lead to serious estimation and variable selection problems in linear regression and that in such a context, it can be useful to explicitly model the structure between the covariates and to use this structure (even sequentially) to avoid correlations issues. We also show that real industrial context faces this kind of situations so our model can help to interpret and predict physical phenomenon efficiently and to help to manage missing values. But for now we still need a full dataset to learn the structure between the covariates and the method only works with numerical values. Further work is needed to face these two challenges.

References

- [1] R. Davidson and J.G. MacKinnon. Estimation and inference in econometrics. *Oxford University Press Catalogue*, 1993.
- [2] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004.
- [3] Dean P Foster and Edward I George. The risk inflation criterion for multiple regression. *The Annals of Statistics*, pages 1947–1975, 1994.
- [4] J. Friedman, T. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.
- [5] E.I. George and R.E. McCulloch. Variable selection via gibbs sampling. *Journal of the American Statistical Association*, pages 881–889, 1993.
- [6] Walter R Gilks, Sylvia Richardson, and David J Spiegelhalter. *Markov chain Monte Carlo in practice*, volume 2. CRC press, 1996.
- [7] Charles Miller Grinstead and James Laurie Snell. *Introduction to probability*. American Mathematical Soc., 1997.
- [8] É. Lebarbier and T. Mary-Huard. Une introduction au critère bic: fondements théoriques et interprétation. *Journal de la SFdS*, 147(1):39–57, 2006.
- [9] D.W. Marquardt and R.D. Snee. Ridge regression in practice. *American Statistician*, pages 3–20, 1975.
- [10] Geoffrey McLachlan and David Peel. *Finite mixture models*. John Wiley & Sons, 2004.
- [11] G. Saporta. *Probabilités, analyse des données et statistique*. Editions Technip, 2006.
- [12] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- [13] N.H. Timm. *Applied multivariate analysis*. Springer Verlag, 2002.
- [14] Loic Yengo, Julien Jacques, Christophe Biernacki, et al. Variable clustering in high dimensional linear regression models. 2012.
- [15] A. Zellner. An efficient method of estimating seemingly unrelated regressions and tests for aggregation bias. *Journal of the American statistical Association*, pages 348–368, 1962.
- [16] Yongli Zhang and Xiaotong Shen. Model selection procedure for high-dimensional data. *Statistical Analysis and Data Mining*, 3(5):350–358, 2010.
- [17] Peng Zhao and Bin Yu. On model selection consistency of lasso. *J. Mach. Learn. Res.*, 7:2541–2563, December 2006.

- [18] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.