

Statistica Sinica Preprint No: SS-2025-0211	
Title	High-dimensional Inference for Model Averaging Estimators
Manuscript ID	SS-2025-0211
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.202025.0211
Complete List of Authors	Lise Léonard, Eugen Pircalabelu and Rainer von Sachs
Corresponding Authors	Eugen Pircalabelu
E-mails	eugen.pircalabelu@uclouvain.be
Notice: Accepted author version.	

INFERENCE FOR HIGH-DIMENSIONAL MODEL AVERAGING ESTIMATORS

Lise Léonard, Eugen Pircalabelu and Rainer von Sachs

UCLouvain, Institute of Statistics, Biostatistics and Actuarial Sciences

Abstract: **Contents of the Abstract.**

Selection methods for high-dimensional models are well developed, but they do not take into account the choice of the model, which leads to an underestimation of the variability of the estimator. We propose a procedure for model averaging in high-dimensional regression models that allows inference even when the number of predictors is larger than the sample size. The proposed estimator is constructed from the debiased Lasso and the weights are chosen to reduce the prediction risk. We derive the asymptotic distribution of the estimator within a high-dimensional framework and offer guarantees for the minimal loss prediction obtained using our choice of the weights. In contrast to existing approaches, our proposed method combines the advantages of model averaging with the possibility of inference based on asymptotic normality. The estimator shows a smaller prediction risk than its competitors when applied to a real, high-dimensional dataset and along various simulation studies, confirming our theoretical results.

Keywords and phrases: Debiased Lasso, High-Dimensional Inference, Model Averaging, Prediction Risk.

1. Introduction

With the rapid evolution of technology, the amount of data to analyze is growing exponentially, leading to new challenges. In particular, samples with a size much smaller than the number of features are increasingly common in many research fields. For example, in biology, transcriptomic data have thousands or tens of thousands of gene expressions per sample unit. This situation does not occur only in biology, and in the last decades a great deal of work has already been done in the area of high-dimensional statistics. For an introduction to the subject and a survey of the existing literature, see Bühlmann and van de Geer (2011) and Giraud (2014).

In high-dimensional settings, model selection methods are very common and intensively used in practice to reduce the number of parameters. The most famous estimator for regression is the Lasso introduced by Tibshirani (1996), which selects the explanatory variables with the largest signal and shrinks to zero the coefficients related to the remaining variables. In a Bayesian framework, Park and Casella (2008) developed a similar shrinkage method. Other model selection methods are based on information criteria, such as the AIC (Akaike, 1979), BIC (Schwarz, 1978), or Mallows's C_p (Mallows, 1973). In particular, Shao (1997) inspected the asymptotic validity of several model selection procedures, concluding that Mallows's C_p is asymptotically loss efficient for low dimensional regression models. See Burnham and Anderson (2002) and Claeskens and Hjort (2008) for an overview of model selec-

tion in statistics. A disadvantage of such methods is the loss of information for the unselected variables. Especially when it comes to performing inference, estimation of the residual error variance does not take into account the selection step, failing to account for all the sources of variability (Burnham and Anderson, 2002).

An alternative to model selection is model averaging (MA) as it consists in aggregating several models instead of selecting just one. The aims are not only to propose an alternative to selection but also to reduce the risk associated with the estimator by incorporating information from different models. Such averaging strategies were first developed for Bayesian techniques; for a literature review see Hoeting et al. (1999) and Raftery and Zheng (2003). In frequentist theory, MA techniques have been developed more recently, namely over the last two decades. What is missing in the literature is an MA procedure which, in a high-dimensional setting, achieves both a low prediction loss and is asymptotically normal. This is precisely the gap that we will fill with our proposed procedure.

MA estimators for regression models are constructed by a weighted mean of different estimators obtained from different models, and most of the time, the method is based on least squares with nested models. For example, the first model is a simple regression with one explanatory variable, the second and further models add subsequently multiple variables until the full model is reached. However, in MA, the choice of the weights remains an important point of research. On this topic, Buckland et al.

(1997) first proposed weights based on the AIC, using a similar construction as used in the Bayesian framework. Cade (2015) argued that weights based on an information criterion, such as the AIC, are not an effective solution if there is multicollinearity in the data. An alternative is to construct the weights based on an optimization problem. This approach was adopted in Hansen (2007), where an estimator based on Mallows's C_p weights was proposed. This was shown further in Hansen (2014) to provide minimal prediction risk for homoscedastic models. This approach is supported in Le and Clarke (2022), where further asymptotic properties of the estimator are established. Leveraging the optimal performance of Mallows C_p , our proposal in Section 3 also uses this criterion to achieve minimal prediction loss. Finally, Hansen and Racine (2012) proposed a jackknife MA strategy for models with heteroscedastic errors, Chen et al. (2022) studied the consistency of MA weights based on BIC and Zhang et al. (2015) developed a Kullback-Leibler loss model averaging. A larger discussion on the choice of weights can be found in Fletcher (2018).

Model averaging techniques have also been developed for frameworks wider than linear regression models. Zhang et al. (2016) proposed an MA estimator for generalized linear models, and Zhang and Wang (2019) constructed an estimator for partial linear regression. For a larger overview of MA techniques and their applications, we refer to Moral-Benito (2015) and, more recently, Steel (2020), from both a frequentist and a Bayesian perspective.

In high-dimensional settings, MA is not well developed and most methods still rely on model selection first. For example, Ando and Li (2014) developed a two-step procedure that separates covariates into different models and uses leave-one-out cross-validation to construct the weights. Xie et al. (2021) extended the method to models with missing responses, and Lan et al. (2018) proposed a sequential model averaging procedure that works even in ultra-high dimensions. On the other hand, Schomaker (2012) proposed using the path of penalized methods such as the Ridge to define nested models. More recently, this idea was explored further in Feng and Liu (2020) without developing any theoretical results. Differently from what has been proposed in the MA literature for high-dimensional models, our proposal in Section 3 explores a sequence of debiased Lasso estimators, with tractable asymptotic distributions, allowing as well for valid inferential strategies (see van de Geer et al., 2014).

One is also interested in the asymptotic properties of the MA estimators. For low-dimensional models, Hjort and Claeskens (2003) studied the asymptotic properties of model averaging estimators constructed from likelihood-based models. The paper by Peng (2024) on asymptotic risk analysis explored the link between MA estimators and Stein shrinkage estimators, showing that Stein shrinkage estimators can be viewed as modified MA estimators. Regarding the distribution, only few results are available. Pötscher (2006) showed that is impossible to estimate the distribution

of a considered least squares model averaging, and Charkhi et al. (2016) derived a non-standard distribution for another averaged estimator. Indeed, when the weights are selected from the data, they are random, and hence, even when each estimator itself is Gaussian, there is no guarantee to obtain a standard distribution for the average estimator (see Liu, 2015).

The novel contribution of this paper is a model averaging estimator for high-dimensional regression based on the debiased Lasso estimator, exploiting its asymptotic normality in high-dimensional settings. Our proposed averaging is done on the entire Lasso path, and the weights are chosen based on an optimization problem that aims to reduce the prediction risk. We show that the proposed estimator is asymptotically Gaussian, and that the chosen weight vector leads to the best asymptotic prediction loss among all weight vectors. Thus, the proposed estimator combines asymptotic normality and loss optimality in a high-dimensional setting. Furthermore, to the best of our knowledge, it is the only MA estimator that has been shown to be asymptotically Gaussian and prediction efficient.

The paper is organized as follows: Section 2 introduces the high-dimensional model, the debiased Lasso estimator and it discusses the challenges associated with the choice of the regularization parameter. In Section 3, we present the proposed estimator and the weight construction. Section 4 shows the theoretical properties of the new estimator. Section 5 uses a simulation study to show the practical perfor-

mance of the estimator, and in Section 6 we show an application to a real dataset. Section 7 concludes with a discussion.

2. Framework

In this section, we briefly present the context in which we work. In particular, we succinctly introduce the high-dimensional model and the debiased Lasso. We then present the challenges associated with this estimator and illustrate the advantages of the proposed method with a simple toy example.

The classical high-dimensional linear regression model we use here is

$$Y = \mathbf{X}\beta_0 + \epsilon, \quad (2.1)$$

where \mathbf{X} is an $n \times p$ random design matrix, Y is an $n \times 1$ response vector, β_0 is the $p \times 1$ vector of coefficients and $\epsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ is the vector of errors with I_n the $n \times n$ identity matrix. Let $S_0 = \{\beta_{0,j} : \beta_{0,j} \neq 0\}$ be the true active set with cardinality denoted by s_0 . In this work, we allow for $p \geq n$ but we assume that $(\log p)/n = o(1)$.

The Lasso is defined as

$$\hat{\beta}_\lambda^{Lasso} := \arg \min_{\beta \in \mathbb{R}^p} \left(\frac{1}{n} \|Y - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 \right), \quad (2.2)$$

where $\|v\|_2^2 = v^T v$ denotes the squared Euclidean norm and $\|\beta\|_1 = \sum_{j=1}^p |\beta_j|$ is the ℓ_1 norm of the coefficient vector $\beta = (\beta_1, \dots, \beta_p)^T$. The parameter λ is a regularization parameter yet to be determined.

Even if the Lasso estimator is consistent in a high-dimensional framework, it is biased in finite sample setting. In particular, the signal for the active variables is underestimated due to the penalization term (Fan and Li, 2001). To reduce this bias and also to perform inference on all the coefficients, one possibility is to use the debiased Lasso estimator, a desparsified estimator constructed from the Lasso.

The debiased Lasso was introduced by Zhang and Zhang (2014) and further developed in van de Geer et al. (2014) and Javanmard and Montanari (2014). The estimator was derived from the Karush-Kuhn-Tucker conditions associated with Equation (2.2). The authors observed that the bias of the Lasso estimator could be well approximated, and proposed the following construction, in which the second term is designed to remove this bias:

$$\hat{\beta}_{\lambda}^d := \hat{\beta}_{\lambda}^{Lasso} + \frac{1}{n} \hat{\Omega} \mathbf{X}^T \left(Y - \mathbf{X} \hat{\beta}_{\lambda}^{Lasso} \right), \quad (2.3)$$

where $\hat{\beta}_{\lambda}^{Lasso}$ is an initial Lasso estimator with some chosen λ and $\hat{\Omega}$ is a $p \times p$ matrix that plays the role of the inverse of the sample variance-covariance matrix, $\hat{\Sigma} := (1/n) \mathbf{X}^T \mathbf{X}$. Indeed, in high-dimensional settings when $p > n$, $\text{rank}(\mathbf{X}^T \mathbf{X}) < p$ and so the matrix $\hat{\Sigma}$ is not invertible. In low-dimensional settings, when $\hat{\Omega} = \hat{\Sigma}^{-1}$, the estimator corresponds to the least squares estimator. The main motivation for the debiased Lasso estimator can be seen in the following decomposition obtained

from Equation (2.3) with direct algebraic decompositions,

$$\hat{\beta}_{\lambda}^d = \beta_0 + \frac{1}{n} \hat{\Omega} \mathbf{X}^T \epsilon + \left(I_p - \hat{\Omega} \hat{\Sigma} \right) \left(\hat{\beta}_{\lambda}^{Lasso} - \beta_0 \right). \quad (2.4)$$

The intuition is that the debiased Lasso estimator converges to the true parameter if $\hat{\Omega} \hat{\Sigma}$ is close to the identity matrix and if the Lasso estimator converges to the true vector in ℓ_1 norm. The estimator is then motivated by the second term on the right-hand side, as conditionally on the design, the term $(1/\sqrt{n}) \hat{\Omega} \mathbf{X}^T \epsilon$ is Gaussian. This allows showing the asymptotic normality of the estimator. Another important motivation for using this estimator is the significant bias reduction it achieves in high-dimensional settings compared to the Lasso, hence its name. The complete construction of the debiased Lasso estimator, as well as further details and motivations, can be found in van de Geer et al. (2014) and Javanmard and Montanari (2014) among others.

The debiased Lasso estimator depends on the Lasso estimator which itself depends on the regularization parameter λ . This parameter is essential since it determines the weight of the penalization term in the Lasso equation and thus the sparsity of the final Lasso estimator. A small value for λ will result in a Lasso estimator with a small bias but a high variability, which carries over to the debiased Lasso estimator. On the other hand, a high value for λ results in a smaller variability and a larger degree of sparsity for the estimator but with a larger bias. Thus, the choice of this parameter is a trade-off between bias and variability. From a theoretical

point of view, a parameter that is large enough and proportional to $\sigma\sqrt{(\log p)/n}$ is recommended to ensure that the mean squared prediction error converges to zero when n grows and p is a function of n (Lahiri, 2021). However, these theoretical recommendations can not help in the practical choice of the parameter, as the noise level σ is unknown, as well as the proportionality constant.

Many methods for selecting the regularization parameter have been proposed and a data-driven choice has prevailed in the literature. The most common methods are based on the minimization of a certain criterion, such as the AIC and the BIC with some modifications as in Chen and Chen (2008). These methods work well for selecting a useful model (Wang et al., 2009), but the theoretical results hold only when $p/n \rightarrow 0$ and not for high-dimensional designs (see Homrighausen and McDonald, 2018). Another range of methods is based on resampling techniques such as cross-validation which are risk-consistent in high-dimensional settings (see Homrighausen and McDonald, 2017). Nowadays, there is still no consensus on the method to use even if some methods are particularly popular. The choice of a regularization parameter can be seen as a model selection problem since it dictates which variables are retained in the model. In this paper, we propose an averaging procedure that avoids the choice of λ by incorporating all the fitted models associated with each possible regularization parameter into a single, global model. Moreover, the proposed method aims to minimize the prediction risk associated with the proposal

model averaging estimator.

To provide a motivating example, Table 1 shows the performance of the proposed Model Averaging debiased estimator (MA-d), presented thoroughly in Section 3, for a toy example. We compare MA-d with several commonly used methods for choosing λ and we inspect the out-of-sample squared prediction error loss

$$L_n(\hat{\beta}) := \frac{1}{n} \left\| \mathbf{X}\beta_0 - \mathbf{X}\hat{\beta} \right\|_2^2. \tag{2.5}$$

Table 1: Out-of-sample prediction error loss and its standard deviation for the proposed method, MA-d, and the debiased Lasso estimator with different values of the penalization parameter. CV-1se denotes the 'one standard error' rule (Chen and Yang, 2021), AIC and BIC are the λ values that minimize the AIC and the BIC criteria, respectively. CV is based on 5-fold cross-validation, and LOO stands for leave-one-out cross-validation. The simulation settings are detailed in the Supplementary Material.

	MA-d	Debiased Lasso with				
		$\lambda_{\text{CV-1se}}$	λ_{AIC}	λ_{BIC}	λ_{CV}	λ_{LOO}
L_n	1.88	4.51	3.26	3.77	3.03	2.82
std.	0.38	1.38	0.89	1.08	0.81	0.82

The table clearly shows that the different methods for λ selection can provide very

different results, and that the proposed procedure, MA-d, that optimally averages the estimates across the entire path of λ values, performs better than any competitor.

3. Proposed procedure

3.1 Model averaging

For low-dimensional models, the regularization parameter for the Lasso can take any value in \mathbb{R}^+ . However, there is a point, λ_{\max} , such that $\forall \lambda > \lambda_{\max}$, $\hat{\beta}_{\lambda}^{Lasso}$ is estimated at zero for each entry. Note that λ_{\max} depends on the design matrix, as well as on the response vector and it can be determined numerically. In the high-dimensional case, there is moreover a minimal positive value, λ_{\min} , such that the number of active coefficients estimated by $\hat{\beta}_{\lambda_{\min}}^{Lasso}$ is $\min(n, p)$ and problem (2.2) has no solution for any $\lambda < \lambda_{\min}$ (Tibshirani, 2013). In practical applications, a choice for the regularization parameter is always taken inside the interval $[\lambda_{\min}, \lambda_{\max}]$ and these bounds are obtained empirically.

In this work, we use this idea by discretizing the interval $[\lambda_{\min}, \lambda_{\max}]$ into M values $\lambda_1, \dots, \lambda_M$, where M is a positive integer independent of p and n . Therefore, for each λ_m with $m = 1 \dots, M$, one can obtain a solution for optimization problem (2.2). This defines a finite sequence of estimated parameters $\hat{\beta}_1^{Lasso}, \dots, \hat{\beta}_M^{Lasso}$. Using further Equation (2.3), we construct a sequence of M debiased Lasso estimators

3.2 Choice of the weight vector

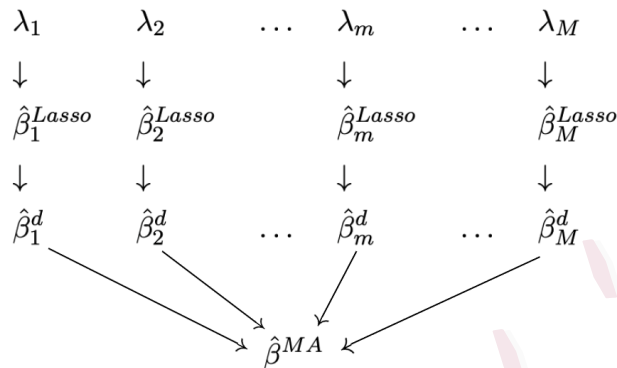


Figure 1: Illustration of the construction of the model averaging estimator: from a grid of penalization parameter values, we estimate first the Lasso, then construct the debiased Lasso and finally the model averaging estimator.

$\hat{\beta}_1^d, \dots, \hat{\beta}_M^d$. Finally, the model averaging estimator is constructed as

$$\hat{\beta}^{MA} := \sum_{m=1}^M w_m \hat{\beta}_m^d,$$

where $\hat{\beta}_m^d := \hat{\beta}_{\lambda_m}^d$ is the debiased Lasso estimator for model m and $\mathbf{w} = (w_1, \dots, w_M)^T$ is a vector of weights satisfying that $\mathbf{w} \in [0, 1]^M$ and $\sum_{m=1}^M w_m = 1$. Figure 1 summarizes these steps.

3.2 Choice of the weight vector

The choice of the weight vector \mathbf{w} is crucial to ensure good properties of the estimator. Hansen (2014) has shown that for low-dimensional model averaging based on the OLS estimator, the weight that minimizes Mallows' Cp leads to the best prediction risk

3.2 Choice of the weight vector

for the model averaging estimator in the case of homoscedastic errors. The Mallows' Cp criterion is not directly defined for high-dimensional regression, but following the same idea of minimization of an empirical loss with a penalization term, we use the following quantity to derive the weights:

$$\begin{aligned} \hat{\mathbf{w}} &:= \arg \min_{\mathbf{w} \in \mathcal{H}} C_n(\mathbf{w}), \quad \mathcal{H} := \left\{ \mathbf{w} : \mathbf{w} \in [0, 1]^M; \sum_{m=1}^M w_m = 1 \right\}, \\ C_n(\mathbf{w}) &:= \frac{1}{n} \left\| Y - \mathbf{X} \left(\sum_{m=1}^M w_m \hat{\beta}_m^d \right) \right\|_2^2 + \frac{2\hat{\sigma}^2}{n} \sum_{m=1}^M w_m \hat{s}_m, \end{aligned} \quad (3.6)$$

where \hat{s}_m is the number of active coefficients estimated by the Lasso for model m , i.e. $\hat{s}_m = |\hat{S}_m| = |\{\hat{\beta}_{m,j}^{Lasso} : \hat{\beta}_{m,j}^{Lasso} \neq 0\}|$ and $\hat{\sigma}^2$ is an estimator for σ^2 . The set \mathcal{H} is the convex set of all weight vectors for which the components are between 0 and 1 and sum to 1; this set is fixed and does not depend on n , nor on p .

The penalization term in the function $C_n(\cdot)$ is included to avoid overfitting and to favor models that are not too complex. Indeed, for two models with the same squared prediction error, the one with the smaller number of estimated active variables by the Lasso is privileged. The number of estimated active coefficients is an unbiased estimator for the degrees of freedom of the Lasso regression (Zou et al., 2007), so this criterion can be seen as an extension of the Mallows' Cp criterion for high-dimensional models, and so optimization problem (3.6) is very similar in spirit to the one used in Hansen (2014). Moreover, the $C_n(\cdot)$ function has already shown good empirical performance in Feng and Liu (2020) for model averaging based on the Lasso only,

3.2 Choice of the weight vector

but no statistical guarantees, nor distributional results were provided.

Optimization problem (3.6) can be re-written as

$$\min_{\mathbf{w} \in \mathcal{H}} \frac{1}{2n} \mathbf{w}^T \bar{E}^T \bar{E} \mathbf{w} + \frac{\hat{\sigma}^2}{n} K^T \mathbf{w}, \quad (3.7)$$

where $\mathbf{w} = (w_1, \dots, w_M)^T$ is the $M \times 1$ vector of weights, $K = (\hat{s}_1, \dots, \hat{s}_M)^T$ is an $M \times 1$ vector and $\bar{E} = (e_1, \dots, e_M)$ is an $n \times M$ matrix containing the residual vectors of each model, $e_m = Y - \mathbf{X} \hat{\beta}_m^d$, for all $m \in \{1, \dots, M\}$.

Remark 1. The optimization program in (3.7) is a quadratic program, so the solution is unique if and only if the matrix $(1/n) \bar{E}^T \bar{E}$ is positive definite. This implies that the matrix \bar{E} needs to be full rank which is not to be expected in practice. In fact, the residuals of model k may be highly correlated with the residuals of model $(k+1)$. However, even if the weights are not necessarily unique, the fitted values *are* unique as shown in Proposition 1. This result implies that a sufficient condition for the uniqueness of the MA-d estimator is to have an invertible design matrix, and this even if the weight vector is not unique.

Proposition 1. *For any fixed sequence $\lambda_1, \lambda_2, \dots, \lambda_M$, any Y , any matrix \mathbf{X} and for any estimator $\hat{\sigma}^2$ of σ^2 , the optimization problem (3.6) is such that*

1. *there exists either one solution $\hat{\mathbf{w}}$ or an infinite number of solutions,*
2. *the fitted values $\mathbf{X} \hat{\beta}^{MA}(\hat{\mathbf{w}})$ associated to the solutions are unique,*

3. if $\hat{\mathbf{w}}_1$ and $\hat{\mathbf{w}}_2$ are two different solutions, then $\sum_{m=1}^M \hat{s}_m \hat{w}_{m,1} = \sum_{m=1}^M \hat{s}_m \hat{w}_{m,2}$.

The proof of this proposition can be found in the Supplementary Material.

4. Theoretical properties

We present here novel theoretical results on the MA-d estimator. The two main results focus on the asymptotic normality of the estimator and the loss reduction. All the rates presented in this paper are valid for an increasing n and for p being a function of n . We also rely on the asymptotic normality of the debiased Lasso, hence the first two assumptions below are from van de Geer et al. (2014). The regression model is defined in Equation (2.1) and the rows of the design matrix are assumed to be drawn randomly from a $\mathcal{N}(0, \Sigma)$ distribution where 0 is the $p \times 1$ zero vector.

Assumption 0. *The sample size, n , and the number of parameters, p , are allowed to grow simultaneously, satisfying $(\log p)/n = o(1)$.*

Assumption 1. *The rows of \mathbf{X} are i.i.d. realizations of a Gaussian distribution whose p -dimensional inner product matrix Σ has a strictly positive smallest eigenvalue $\Lambda_{\min}(\Sigma)$ satisfying $1/\{\Lambda_{\min}(\Sigma)\} = O(1)$. Furthermore for all p , $\max_{j=1,\dots,p} \Sigma_{j,j} = O(1)$.*

Assumption 2. *The sparsity index for the coefficient vector from model (2.1), s_0 , satisfies $s_0 = o(\sqrt{n}/\log p)$ and $\max_{j=1,\dots,p} s_j = o(n/\log p)$, where $s_j := |\{\Omega_{j,k} : j \neq k \text{ and } \Omega_{j,k} \neq 0\}|$ denotes the sparsity in the rows of the precision matrix, $\Omega := \Sigma^{-1}$.*

Assumption 3. *The design matrix \mathbf{X} satisfies $\|\mathbf{X}\|_\infty = O_{\mathbb{P}}(\sqrt{n})$, where $\|\cdot\|_\infty$ denotes the entry-wise sup norm of the matrix.*

Assumption 1 concerns the variance-covariance population matrix Σ . In particular, it implies the compatibility condition required for Lasso convergence (Bühlmann and van de Geer, 2011, Chap. 6). Assumption 2 concerns the sparsity of the coefficient vector β_0 and of the Ω matrix. As studied in Javanmard and Montanari (2018), the sparsity assumption for the debiased Lasso is stricter than the usual one for Lasso, which is $o\left(\sqrt{n/\log p}\right)$. Assumption 3 controls the boundedness of the entries in the design matrix and it is needed to ensure the convergence of the prediction error in Theorem 2. A sufficient condition for this assumption is that $p = O(n^2)$. Indeed, by Theorem 4.4.5 in Vershynin (2018), we have $\|\mathbf{X}\|_\infty \leq \|\mathbf{X}\|_{spect} = \sqrt{\Lambda_{\max}(\mathbf{X}^T \mathbf{X})} = O_{\mathbb{P}}(p^{1/4})$, where $\|\mathbf{X}\|_{spect}$ is the spectral norm.

Our first result focuses on the asymptotic normality of the proposal estimator. Due to the randomness of the weight vector, the result is not trivial, as usually model averaging estimators have an unknown and non-standard distribution even if all the estimators are Gaussian.

Theorem 1. *Under Assumptions 0 – 2, for the linear model with Gaussian error term $\epsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ where $\sigma^2 = O(1)$ and whose $\hat{\mathbf{w}} = (\hat{w}_1, \dots, \hat{w}_M)^T$ denotes the weights selected by (3.6) with $\hat{\sigma}^2$ as estimator for the noise level, if for every model*

the penalization parameter satisfies $\lambda_m = O\left\{\sqrt{(\log p)/n}\right\}$, we have:

$$\sqrt{n}(\hat{\beta}^{MA} - \beta_0) = W + \sum_{m=1}^M \hat{w}_m \Delta_m,$$

$$W|\mathbf{X} \sim \mathcal{N}(0, \sigma^2 \hat{\Omega} \hat{\Sigma} \hat{\Omega}^T),$$

$$\left\| \sum_{m=1}^M \hat{w}_m \Delta_m \right\|_{\infty} = o_{\mathbb{P}}(1),$$

where $W := \hat{\Omega} \mathbf{X}^T \epsilon / \sqrt{n}$ and $\Delta_m := \sqrt{n}(I_p - \hat{\Omega} \hat{\Sigma})(\hat{\beta}_m^{Lasso} - \beta_0)$ is a bias term.

Proof of Theorem 1. By construction of the debiased estimator, we have

$$\sqrt{n}(\hat{\beta}^{MA} - \beta_0) = \sum_{m=1}^M \hat{w}_m \sqrt{n}(\hat{\beta}_m^d - \beta_0) = \sum_{m=1}^M \hat{w}_m (W + \Delta_m) = W + \sum_{m=1}^M \hat{w}_m \Delta_m,$$

where, as in the proof of Theorem 2.2 of van de Geer et al. (2014), $W = \hat{\Omega} \mathbf{X}^T \epsilon / \sqrt{n}$ is Gaussian when we condition on the design matrix and does not depend on the penalization parameter.

Under Assumption 1, for the bias term $\Delta_m := \sqrt{n}(I_p - \hat{\Omega} \hat{\Sigma})(\hat{\beta}_m^{Lasso} - \beta_0)$, we have

$$\|\Delta_m\|_{\infty} \leq \sqrt{n} \|I_p - \hat{\Omega} \hat{\Sigma}\|_{\infty} \|\hat{\beta}_m^{Lasso} - \beta_0\|_1 = O_{\mathbb{P}}\left(s_0 \frac{\log p}{\sqrt{n}}\right),$$

where $\|\Delta_m\|_{\infty} = \max_{j=1, \dots, p} |\Delta_{m,j}|$ denotes the sup norm of the vector and where $\|I_p - \hat{\Omega} \hat{\Sigma}\|_{\infty} = O_{\mathbb{P}}\{\sqrt{(\log p)/n}\}$ is shown in Theorem 2.2 from van de Geer et al. (2014).

Since $\hat{\mathbf{w}} \in [0, 1]^M$ for all n and p , and as the number of considered models M is fixed,

$$\left\| \sum_{m=1}^M \hat{w}_m \Delta_m \right\|_{\infty} \leq \sum_{m=1}^M \|\hat{w}_m \Delta_m\|_{\infty} \leq \sum_{m=1}^M \|\Delta_m\|_{\infty} = O_{\mathbb{P}}\left(s_0 \frac{\log p}{\sqrt{n}}\right).$$

Moreover, under Assumption 2, which controls the sparsity index s_0 , we have

$$\left\| \sum_{m=1}^M \hat{w}_m \Delta_m \right\|_{\infty} = O_{\mathbb{P}} \left(s_0 \frac{\log p}{\sqrt{n}} \right) = o_{\mathbb{P}}(1).$$

□

Theorem 1 shows that the MA-d estimator can be decomposed in two parts. The first one does not depend on the weights and is asymptotically Gaussian and the second part is a bias term whose sup norm converges in probability to zero.

Due to its construction, the decomposition of the MA-d estimator in Theorem 1 shows that crucially, the weights play a role *only* in the second term. This gives us the result of Proposition 2 below which determines the rate of convergence for the estimated coefficients.

For the next results, one needs to control the noise level of the statistical problem and more precisely, the term $\mathbf{X}^T \epsilon / n$. For this purpose, we work on the event $\xi_{\infty} := \{ \|\mathbf{X}^T \epsilon / n\|_{\infty} \leq \lambda_0 \}$, with $\lambda_0 = O \left\{ \sqrt{(\log p)/n} \right\}$ and by Bühlmann and van de Geer (2011, Lemma 6.2), directly applicable to our setting, the probability of the event tends to one for increasing samples sizes.

Proposition 2. *Under Assumptions 0 – 2, for the linear model with Gaussian error term $\epsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ where $\sigma^2 = O(1)$, if for every model $m \in \{1, \dots, M\}$, the penalization parameter satisfies $\lambda_m = O \left\{ \sqrt{(\log p)/n} \right\}$, we have*

$$\left\| \hat{\beta}^{MA}(\hat{\mathbf{w}}) - \beta_0 \right\|_{\infty} = O_{\mathbb{P}} \left(\max_{j=1, \dots, p} \sqrt{s_j} \sqrt{\frac{\log p}{n}} + s_0 \frac{\log p}{n} \right) = o_{\mathbb{P}}(1),$$

where $\hat{\beta}^{MA}(\hat{\mathbf{w}})$ denotes the model averaging estimator with weights selected by (3.6) with $\hat{\sigma}^2$ as estimator for the noise level.

The proof is provided in the Supplementary Material. The main idea is to decompose the MA-d estimator as in Theorem (1) and derive a bound for each term.

These two results show that, unlike many model averaging estimators, the proposed estimator is asymptotically Gaussian and maintains the same rate of convergence in sup-norm *even after* the averaging step. In addition to these interesting properties, the following theorem shows that the weights obtained by (3.6) lead to the asymptotically best loss compared to *any* other choice of weights. In particular, the MA-d is no worse than any individual debiased model.

Theorem 2. *Under Assumptions 0 – 3, for the linear model with Gaussian error term $\epsilon \sim \mathcal{N}(0, \sigma^2 I_n)$, let $\hat{\mathbf{w}} = (\hat{w}_1, \dots, \hat{w}_M)^T$ denote the weights selected by (3.6) with $\hat{\sigma}^2$ as estimator for the noise level. If for every model the penalization parameter satisfies $\lambda_m = O\left\{\sqrt{(\log p)/n}\right\}$ and if $\hat{\sigma}^2$ satisfies $|(\hat{\sigma}^2/\sigma^2) - 1| = o_{\mathbb{P}}(1)$ we have:*

$$\frac{L_n\left\{\hat{\beta}^{MA}(\hat{\mathbf{w}})\right\}}{\inf_{\mathbf{w} \in \mathcal{H}} L_n\left\{\hat{\beta}^{MA}(\mathbf{w})\right\}} \rightarrow 1 \text{ in probability,}$$

where $L_n(\cdot)$ denotes the quadratic loss function defined in Equation (2.5).

The proof of Theorem 2 can be found in the Supplementary Material, here we just explain briefly the main idea. We begin by showing that if the weights minimize

$C_n(\mathbf{w})$ by definition of optimization problem (3.6), they also minimize a function $\tilde{C}_n(\mathbf{w})$ defined as $\tilde{C}_n(\mathbf{w}) = L_n\{\hat{\beta}^{MA}(\mathbf{w})\} + R\{\hat{\beta}^{MA}(\mathbf{w})\}$, where $R\{\hat{\beta}^{MA}(\mathbf{w})\}$ is a remainder term to be defined in the proof. Afterward, we explore under which conditions these two terms converge to zero, and we show that the convergence rate of $L_n\{\hat{\beta}^{MA}(\mathbf{w})\}$ is slower than the one of $R\{\hat{\beta}^{MA}(\mathbf{w})\}$. The closing argument is then: asymptotically choosing the weights by minimizing $C_n(\mathbf{w})$ is equivalent to choosing the weights by minimizing $L_n(\cdot)$.

The next proposition extends the result of Theorem 1 to the case where the number of models, M , increases with n . For this, one needs to control all the regularization parameters in the grid for each n and M . Interestingly, our proposition shows that the growth rate of M does not affect the normality result. The proof can be found in the Supplementary Material.

Proposition 3. *Under Assumption 0 – 2, for the linear model with Gaussian error term $\epsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ where $\sigma^2 = O(1)$, let $M_n = M(n)$ be the number of models and $\hat{\mathbf{w}} = (\hat{w}_1, \dots, \hat{w}_{M_n})^T$ the weights selected by (3.6) with $\hat{\sigma}^2$ as estimator for the noise level. If for every model and every sample size n , we have $\lambda_m \propto \sqrt{(\log p)/n}$ and*

$\sup_{m \leq M_n} \lambda_m < K \sqrt{(\log p)/n}$ with $K < \infty$ for all M_n , then

$$\sqrt{n}(\hat{\beta}^{MA} - \beta_0) = W + \sum_{m=1}^{M_n} \hat{w}_m \Delta_m,$$

$$W|\mathbf{X} \sim \mathcal{N}(0, \sigma^2 \hat{\Omega} \hat{\Sigma} \hat{\Omega}^T),$$

$$\left\| \sum_{m=1}^{M_n} \hat{w}_m \Delta_m \right\|_{\infty} = o_{\mathbb{P}}(1),$$

where $W := \hat{\Omega} \mathbf{X}^T \epsilon / \sqrt{n}$ and $\Delta_m := \sqrt{n}(I_p - \hat{\Omega} \hat{\Sigma})(\hat{\beta}_m^{Lasso} - \beta_0)$ is a bias term.

5. Simulation

We illustrate the performance of the proposed method with a numerical study. There are two main results we expect to verify with this study: (i) the optimality of the loss as shown in Theorem 2 and (ii) the asymptotic normality of the estimator as shown in Theorem 1.

We work under model (2.1) where the rows of the design matrix \mathbf{X} are i.i.d. Gaussian, having a correlation matrix Σ that follows a Toeplitz structure, $\Sigma_{ab} = \rho^{|a-b|}$ with $\rho = 0.5$. The vector of coefficients is set to $\beta_0 = (2, \dots, 2, 0, \dots, 0)^T$ with the number of non-null entries, the sparsity index s_0 , set at $10n^{1/4}/\log p$ rounded to the nearest integer. The sample size is $n \in \{100, 200, 300, 500, 750, 1000\}$ and the number of covariates is $p = 2n$, so that each design corresponds to a high-dimensional setting. For the noise level σ_0^2 , we allow it to vary in the set $\{0.75, 1.5, 3\}$.

The MA-d estimator is constructed with an estimated value for the noise level in the weights optimization problem (3.6). Here the estimator $\hat{\sigma}_{\text{scaled}}^2$ from Sun and Zhang (2012) is used. It is a consistent estimator in high-dimensional settings and

it is obtained as

$$\left(\hat{\beta}_{\text{scaled}}, \hat{\sigma}_{\text{scaled}}\right) := \arg \min_{\beta, \sigma} \left(\frac{\|Y - \mathbf{X}\beta\|_2^2}{2n\sigma} + \frac{\sigma}{2} + \lambda_0 \|\beta\|_1 \right), \quad (5.8)$$

where $\lambda_0 \approx \sqrt{(2/n) \log p}$ is chosen using the quantile method defined and studied in Sun and Zhang (2013). In practice, we use the **scalreg** R package with the default parameters to obtain the estimates.

The number of models is set to $M = 200$. The competitors are the debiased Lasso and the Lasso with λ obtained by 10-fold cross-validation. We also included MA-Lasso, a model averaging estimator on the Lasso path developed by Feng and Liu (2020) for which there are no theoretical results. This estimator is a weighted mean of the Lasso estimates, where the weight vector is obtained via an optimization similar to the one we proposed in Equation (3.6).

$$\begin{aligned} \hat{\beta}^{\text{MA-Lasso}} &:= \sum_{m=1}^M \hat{w}_m \hat{\beta}_m^{\text{Lasso}}, \\ \hat{\mathbf{w}} &:= \arg \min_{\mathbf{w} \in \mathcal{H}_M} \frac{1}{n} \left\| Y - \mathbf{X} \sum_{m=1}^M w_m \hat{\beta}_m^{\text{Lasso}} \right\|_2^2 + \frac{2\sigma^2}{n} \sum_{m=1}^M \hat{s}_m w_m, \end{aligned}$$

with $\mathcal{H}_M := \{\mathbf{w} \in [0, 1]^M; \sum_{m=1}^M w_m = 1\}$, the set of weight vectors in \mathbb{R}^M and \hat{s}_m is the number of non zero estimated coefficients in the solution $\hat{\beta}_m^{\text{Lasso}}$.

For all lasso-based methods, the regularization parameter is chosen from the same path of M values. For the construction of the inference metrics, $\hat{\sigma}_{\text{scaled}}^2$ is used for all competitors that require an estimate of the noise level. Additional simulations

comparing different estimators for σ^2 are shown in the Supplementary Material, but they all point roughly to similar conclusions so we skip presenting the results here.

Table 2 shows, averaged over $R = 1000$ repetitions, the following loss ratio

$$\text{LR} := \frac{1}{R} \sum_{r=1}^R \left[\frac{L_n \{ \hat{\beta}^{(r)} \}}{\min_{\mathbf{w} \in \mathcal{H}} L_n \{ \hat{\beta}^{MA, (r)}(\mathbf{w}) \}} \right]. \quad (5.9)$$

The denominator is the minimum loss function for model averaging with oracle weights. These weights are the ones that directly minimize $L_n \{ \hat{\beta}^{MA}(\mathbf{w}) \}$, which is known in this simulation setting. The numerator is the loss function for the estimator at the r th iteration. This ratio allows us to evaluate if the weights obtained by minimizing (3.6) give a loss similar to the one of the oracle.

From Table 2, left-hand-side panel, we conclude that the Lasso and the MA-Lasso are the only methods with a ratio below 1, which was expected since they are sparse methods, the true model is sparse and the denominator in (5.9) is the minimum loss function for model averaging with oracle weights. Among the non-sparse methods, we observe that for each sample size and noise level, the MA-d has the smallest ratio. As the noise level increases, the ratio increases as well, and as the sample size increases, the ratio decreases to unity as expected. The ratio of the competitor also converges to unity, but at a much slower rate. The fast decay for MA-d shows that in practice the weights obtained by minimizing (3.6) are very close to the optimal ones, pointing to substantial gains for finite sample analyzes.

Table 2: Prediction loss ratio (5.9) averaged over $R = 1000$ repetitions and average coverage for active and non-active coefficients.

		Prediction loss ratio				Inference metrics (<i>nominal level: .95</i>)			
σ_0^2	n	MA-d	Debiased Lasso	Lasso	MA-Lasso	MA-d		Debiased Lasso	
						Active	Non-active	Active	Non-active
0.75	100	1.08	2.81	0.11	0.10	.87	.98	.85	.95
	200	1.01	2.19	0.05	0.04	.89	.98	.87	.95
	300	1.00	1.78	0.03	0.03	.90	.98	.88	.95
	500	1.00	1.43	0.01	0.01	.92	.97	.90	.95
	750	1.00	1.24	0.01	0.01	.93	.97	.91	.95
	1000	1.00	1.13	0.01	0.01	.93	.96	.91	.95
1.5	100	1.23	3.18	0.14	0.13	.86	.98	.86	.96
	200	1.08	2.78	0.06	0.06	.88	.98	.89	.96
	300	1.02	2.34	0.04	0.04	.89	.98	.89	.96
	500	1.00	1.91	0.02	0.02	.91	.98	.91	.96
	750	1.00	1.63	0.01	0.01	.92	.98	.92	.95
	1000	1.00	1.45	0.01	0.01	.92	.97	.92	.95
3	100	1.41	3.56	0.17	0.16	.85	.98	.87	.96
	200	1.27	3.43	0.08	0.08	.88	.98	.90	.96
	300	1.13	3.03	0.05	0.05	.88	.98	.90	.96
	500	1.02	2.57	0.03	0.03	.89	.98	.91	.96
	750	1.00	2.23	0.02	0.02	.91	.98	.92	.96
	1000	1.00	1.98	0.01	0.01	.91	.98	.92	.95

The prediction loss ratio of the MA-Lasso is slightly smaller than that one of the Lasso, confirming the results obtained in Feng and Liu (2020). However, the improvement due to the model averaging procedure is modest, compared to the difference between the MA-d estimator and the debiased Lasso, for which the model averaging procedure reduces the prediction loss ratio by a factor two.

Another point of interest is the asymptotic distribution of the estimator. The right-hand side panel in Table 2 shows the average coverage for active and non-active variables. The target coverage is 95%. The only competitor is the debiased Lasso estimator, since it is the only other considered competitor with a known distribution in the high-dimensional setting. The empirical coverage is calculated as

$$\begin{aligned} \text{Cvg}(\hat{\beta}_j) &:= \frac{1}{R} \sum_{r=1}^R \mathbf{1} \left[\beta_{0,j} \in \widehat{CI}_{1-\alpha} \left\{ \hat{\beta}_j^{(r)} \right\} \right], \\ \text{avgCvg}_a(\hat{\beta}) &:= \frac{1}{s} \sum_{j=1}^s \text{Cvg}(\hat{\beta}_j), \quad \text{avgCvg}_{na}(\hat{\beta}) := \frac{1}{p-s} \sum_{j=s+1}^p \text{Cvg}(\hat{\beta}_j), \end{aligned}$$

where $\mathbf{1}[A] = 1$ if A takes place, is the indicator function and $\widehat{CI}_{1-\alpha} \left\{ \hat{\beta}_j^{(r)} \right\}$ is the confidence interval at level $1 - \alpha$ for the j th coefficient at iteration r . The confidence interval is computed with the asymptotic variance presented in Theorem 1, which is, for $\hat{\beta}_j$, $\hat{\sigma}^2(\hat{\Omega}\hat{\Sigma}\hat{\Omega})_{j,j}$. For both estimators, the noise level is estimated by $\hat{\sigma}_{\text{scaled}}^2$.

Table 2 shows that the performance of the MA-d and that of the debiased Lasso estimator, with respect to the empirical coverage, are comparable. The coverage for both methods is close and gets closer to the target as the sample size increases as

expected. The coverage is slightly lower for the active variables and slightly higher for the non-active variables and this phenomenon occurs for both methods. When the noise level is high, the coverage performance of MA-d is slightly worse. Since the standard error of the model averaging estimator does not depend on the weights, the lengths of the confidence intervals of the two methods are the same and as such, their values are not presented here.

We show next in Figure 2 the distribution of the MA-d estimator, for a randomly chosen active and non-active standardized coefficient for different values of n . The standardization was obtained by

$$\frac{\hat{\beta}_j^{MA} - \beta_{0,j}}{\text{sd}(\hat{\beta}_j^{MA})} = \frac{\hat{\beta}_j^{MA} - \beta_{0,j}}{\sqrt{\hat{\sigma}_{\text{scaled}}^2(\hat{\Omega}\hat{\Sigma}\hat{\Omega})_{j,j}}}.$$

Figure 2 complements well the results for coverage. It illustrates that the distribution for the non-active coefficients has slightly larger variance for small sample size, but it approaches the Gaussian limit as n increases.

Additional accuracy metrics such as bias and MSE for active and non-active coefficients can be found in the Supplementary Material, but the results show empirically that the MA-d estimator retains the bias reduction property of the debiased Lasso.

Remark 2. In Equation (2.3), we use the node-wise method from van de Geer et al. (2014) to obtain the matrix $\hat{\Omega}$ as it is the most used method and as there is a dedicated R package for it (Dezeure et al., 2015). With this method, the matrix

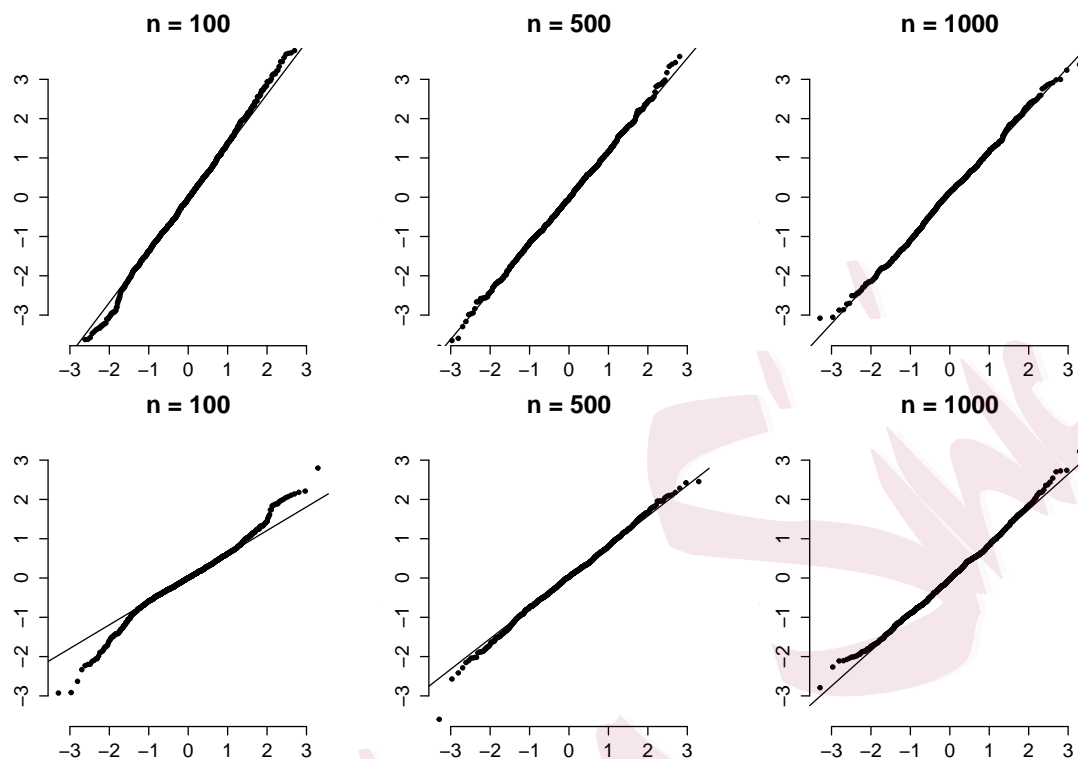


Figure 2: QQplot for an active (top row) and a non-active (bottom row) standardized estimated coefficient by MA-d over 1000 replications.

$\hat{\Omega}$ is constructed only from \mathbf{X} and is independent of the choice of the penalization parameter. Thus, it can be constructed once and be reused to estimate the debiased Lasso for different values of the penalization parameter.

Remark 3. In practice, a log-scale is used for the construction of the grid of candidate values for λ as the coefficients vary more for small values. On the other hand, the coarseness of the grid empirically seems not to affect the results too much. We

explore this point with a simulation study in the Supplementary Material, where we recommend a moderate value such as $M = 10\sqrt{n}$.

Remark 4. The list of competitors presented here is not exhaustive. Averaging competitors based on the Lasso and the debiased Lasso are presented in the Supplementary Material, where their performance is compared with that of the MA-d. Among other things, we observe that the MA-d estimator outperforms all other MA methods based on the debiased Lasso that we considered. Additionally, we present a sparse scenario in which the MA-d estimator provides better prediction loss performance relative to the Lasso and better coverage performance relative to the debiased Lasso.

6. Real data analysis

The Bardet-Biedl Syndrome (BBS) is a genetic disease that affects several organs. The main symptoms are obesity and retinal failure, but polydactia and learning disabilities can also occur. The `eyedata` database available in the R package `flare` (Li et al., 2015) contains the expression of 200 genes for 120 rats affected by BBS, and is derived from Scheetz et al. (2006). In that study, the TRIM32 gene was identified as a disease-causing gene by looking at the correlation with other known BBS genes.

The aim of this section is to identify whether the predictor genes for the TRIM32 gene are BBS genes that could play a role in the disease. To identify genes linked

to TRIM32, we perform a regression analysis on the dataset with TRIM32 as response and since there are more covariates than sampling units, the usage of high-dimensional methods is mandatory. On the whole dataset, the Lasso selects 57 coefficients when the regularization parameter is chosen by 10-fold cross-validation and 19 if it is chosen by the ‘1se’ rule. The debiased estimator based on this cross-validation Lasso estimates that 3 variables are significant for a confidence level of 95% and with a Bonferroni correction for multiplicity. The significant genes are genes labelled 10540, 16984 and 17599 if we refer to the notation of the dataset. The MA-d estimator identifies that only 2 coefficients are significant for the same confidence level and multiplicity correction; these are genes 10540 and 17599.

To test the performance of the estimators, we compute the leave-one-out prediction risk (LOO) and compare it with the debiased Lasso estimator, the Lasso and the Ridge estimator where the penalization parameter was chosen by 10-fold cross-validation. The noise level is estimated with $\hat{\sigma}_{\text{scaled}}$ defined in Equation (5.8) when needed. We define

$$\text{LOO}(\hat{\beta}) := \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \hat{\beta}_{-i})^2,$$

where $\hat{\beta}_{-i}$ is the estimator computed from all observations except the i th one.

Table 3: Leave-one-out error and its standard deviation on the *eyedata* for the MA-d, Lasso, the debiased Lasso and the Ridge. The penalization parameter was selected by 10-fold cross-validation for the methods that require it.

	MA-d	Debiased Lasso	Ridge	Lasso
LOO	0.48	3.87	0.40	0.45
std.	0.73	21.52	0.95	1.38

To obtain a more stable calculation, both the predictors and the response vector have been standardized, and the prediction risk is expressed in terms of these standardized data. Table 3 shows the results of this analysis, and clearly the MA-d estimator performs similarly to the Lasso or the Ridge with a smaller standard deviation. However, the debiased estimator has a very high prediction risk for this example.

7. Conclusion

In this paper, we proposed a new model averaging estimator for high-dimensional regression. The estimator is based on the debiased Lasso estimator, with averaging taken along the Lasso solution path, and the weights obtained through optimization of a sample loss. We showed the asymptotic normality of this novel averaging es-

timator and the optimality of the weights in terms of prediction risk. In practice, the proposed estimator outperformed competitors in terms of prediction loss and the asymptotic normality continues to hold even for the high-dimensional design. Moreover, since the procedure is based on the debiased estimator, our empirical results show that the MA-d estimator still maintains a bias advantage, even after averaging on the λ path.

While in this work we have focused on model averaging applied to the debiased Lasso estimator, it would be interesting to explore whether for other high-dimensional estimators, model averaging can lead to substantial gains in terms of prediction or accuracy. Another open question is the control of the False Discovery Rate (FDR) under model averaging. Linking to the previous idea, it would be interesting to evaluate if, under high-dimensional settings, a model averaging procedure using the SLOPE estimator (Bogdan et al., 2015) maintains FDR control.

Supplementary Material

The Supplementary Material contains the proofs of Theorem 2, Propositions 1, 2 and 3, bias and MSE metrics for the setting presented in Section 5, and additional simulations. Among these, we show that the MA-d estimator maintains good performance using different estimators for the noise level, and that the influence of the grid coarseness is negligible when the number of considered models is sufficiently

REFERENCES

large. Furthermore, we show that, in some sparse cases, the MA-d estimator can provide better prediction loss performance relative to the Lasso and better coverage performance relative to the debiased Lasso. We also constructed MA competitors based on the Lasso and the debiased Lasso and compared their performance with that of the proposed method. Additionally, we also explored the performance of the MA-d estimator when the precision matrix is no longer sparse and when the signal decreases as the sample size increases.

Acknowledgement

Computational resources have been provided by the supercomputing facilities of the Université catholique de Louvain (CISM/UCL) and the Consortium des Équipements de Calcul Intensif en Fédération Wallonie Bruxelles (CÉCI) funded by the Fond de la Recherche Scientifique de Belgique (F.R.S.-FNRS) under convention 2.5020.11 and by the Walloon Region.

References

- Akaike, H. (1979). A Bayesian Extension of the Minimum AIC Procedure of Autoregressive Model Fitting. *Biometrika* 66(2), 237–242.
- Ando, T. and K.-C. Li (2014). A Model-Averaging Approach for High-Dimensional Regression. *Journal of the American Statistical Association* 109(505), 254–265.

REFERENCES

- Bogdan, M., E. van den Berg, C. Sabatti, W. Su, and E. J. Candès (2015). SLOPE—Adaptive Variable Selection via Convex Optimization. *The Annals of Applied Statistics* 9(3), 1103–1140.
- Buckland, S. T., K. P. Burnham, and N. H. Augustin (1997). Model Selection: An Integral Part of Inference. *Biometrics* 53(2), 603–618.
- Bühlmann, P. and S. van de Geer (2011). *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer Series in Statistics. Berlin, Heidelberg: Springer.
- Burnham, K. P. and D. R. Anderson (2002). *Model Selection and Multimodel Inference*. New York, NY: Springer.
- Cade, B. S. (2015). Model averaging and muddled multimodel inferences. *Ecology* 96(9), 2370–2382.
- Charkhi, A., G. Claeskens, and B. E. Hansen (2016). Minimum Mean Squared Error Model Averaging in Likelihood Models. *Statistica Sinica* 26(2), 809–840.
- Chen, J. and Z. Chen (2008). Extended Bayesian Information Criteria for Model Selection with Large Model Spaces. *Biometrika* 95(3), 759–771.
- Chen, Y. and Y. Yang (2021). The One Standard Error Rule for Model Selection: Does It Work? *Stats* 4(4), 868–892.
- Chen, Z., J. Zhang, W. Xu, and Y. Yang (2022). Consistency of BIC Model Averaging. *Statistica Sinica* 32, 635–640.
- Claeskens, G. and N. L. Hjort (2008). *Model Selection and Model Averaging*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press.

REFERENCES

- Dezeure, R., P. Bühlmann, L. Meier, and N. Meinshausen (2015). High-Dimensional Inference: Confidence Intervals, p-Values and R-Software hdi. *Statistical Science* 30(4), 533–558.
- Fan, J. and R. Li (2001). Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties. *Journal of the American Statistical Association* 96(456), 1348–1360.
- Feng, Y. and Q. Liu (2020). Nested model averaging on solution path for high-dimensional linear regression. *Stat* 9(1), e317.
- Fletcher, D. (2018). *Model Averaging*. SpringerBriefs in Statistics. Berlin, Heidelberg: Springer.
- Giraud, C. (2014). *Introduction to High-Dimensional Statistics*. Boca Raton: Chapman and Hall/CRC.
- Hansen, B. E. (2007). Least Squares Model Averaging. *Econometrica* 75(4), 1175–1189.
- Hansen, B. E. (2014). Model averaging, asymptotic risk, and regressor groups. *Quantitative Economics* 5(3), 495–530.
- Hansen, B. E. and J. S. Racine (2012). Jackknife model averaging. *Journal of Econometrics* 167(1), 38–46.
- Hjort, N. L. and G. Claeskens (2003). Frequentist Model Average Estimators. *Journal of the American Statistical Association* 98(464), 879–899.
- Hoeting, J. A., D. Madigan, A. E. Raftery, and C. T. Volinsky (1999). Bayesian Model Averaging: A Tutorial. *Statistical Science* 14(4), 382–401.
- Homrighausen, D. and D. J. McDonald (2017). Risk Consistency of Cross-Validation with Lasso-Type Procedures. *Statistica Sinica* 27(3), 1017–1036.
- Homrighausen, D. and D. J. McDonald (2018). A study on tuning parameter selection for the high-

REFERENCES

- dimensional lasso. *Journal of Statistical Computation and Simulation* 88(15), 2865–2892.
- Javanmard, A. and A. Montanari (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *The Journal of Machine Learning Research* 15(1), 2869–2909.
- Javanmard, A. and A. Montanari (2018). Debiasing the lasso: Optimal sample size for Gaussian designs. *The Annals of Statistics* 46(6A), 2593–2622.
- Lahiri, S. N. (2021). Necessary and sufficient conditions for variable selection consistency of the LASSO in high dimensions. *The Annals of Statistics* 49(2), 820–844.
- Lan, W., Y. Ma, J. Zhao, H. Wang, and C.-L. Tsai (2018). Sequential Model Averaging for High Dimensional Linear Regression Models. *Statistica Sinica* 28(1), 449–469.
- Le, T. M. and B. S. Clarke (2022). Model Averaging Is Asymptotically Better Than Model Selection For Prediction. *Journal of Machine Learning Research* 23(33), 1–53.
- Li, X., T. Zhao, X. Yuan, and H. Liu (2015). The flare Package for High Dimensional Linear Regression and Precision Matrix Estimation in R. *The Journal of Machine Learning Research* 16(18), 553–557.
- Liu, C.-A. (2015). Distribution theory of the least squares averaging estimator. *Journal of Econometrics* 186(1), 142–159.
- Mallows, C. L. (1973). Some Comments on CP. *Technometrics* 15(4), 661–675.
- Moral-Benito, E. (2015). Model Averaging in Economics: An Overview. *Journal of Economic Surveys* 29(1), 46–75.
- Park, T. and G. Casella (2008). The Bayesian Lasso. *Journal of the American Statistical Association*

REFERENCES

- tion 103(482), 681–686.
- Peng, J. (2024). Model averaging: A shrinkage perspective. *Electronic Journal of Statistics* 18(2), 3535–3572.
- Pötscher, B. M. (2006). The Distribution of Model Averaging Estimators and an Impossibility Result regarding Its Estimation. *Lecture Notes-Monograph Series* 52, 113–129.
- Raftery, A. E. and Y. Zheng (2003). Discussion: Performance of Bayesian Model Averaging. *Journal of the American Statistical Association* 98(464), 931–938.
- Scheetz, T. E., K.-Y. A. Kim, R. E. Swiderski, A. R. Philp, T. A. Braun, K. L. Knudtson, A. M. Dorrance, G. F. DiBona, J. Huang, T. L. Casavant, V. C. Sheffield, and E. M. Stone (2006). Regulation of gene expression in the mammalian eye and its relevance to eye disease. *Proceedings of the National Academy of Sciences of the United States of America* 103(39), 14429–14434.
- Schomaker, M. (2012). Shrinkage averaging estimation. *Statistical Papers* 53(4), 1015–1034.
- Schwarz, G. (1978). Estimating the Dimension of a Model. *The Annals of Statistics* 6(2), 461–464.
- Shao, J. (1997). An Asymptotic Theory for Linear Model Selection. *Statistica Sinica* 7(2), 221–242.
- Steel, M. F. J. (2020). Model Averaging and Its Use in Economics. *Journal of Economic Literature* 58(3), 644–719.
- Sun, T. and C.-H. Zhang (2012). Scaled sparse linear regression. *Biometrika* 99(4), 879–898.
- Sun, T. and C.-H. Zhang (2013). Sparse Matrix Inversion with Scaled Lasso. *Journal of Machine Learning Research* 14(106), 3385–3418.

REFERENCES

- Tibshirani, R. (1996). Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* 58(1), 267–288.
- Tibshirani, R. J. (2013). The lasso problem and uniqueness. *Electronic Journal of Statistics* 7, 1456–1490.
- van de Geer, S., P. Bühlmann, Y. Ritov, and R. Dezeure (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics* 42(3), 1166–1202.
- Vershynin, R. (2018). *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press.
- Wang, H., B. Li, and C. Leng (2009). Shrinkage Tuning Parameter Selection with a Diverging Number of Parameters. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 71(3), 671–683.
- Xie, J., X. Yan, and N. Tang (2021). A Model-averaging method for high-dimensional regression with missing responses at random. *Statistica Sinica* 31(2), 1005–1026.
- Zhang, C.-H. and S. S. Zhang (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 76(1), 217–242.
- Zhang, X. and W. Wang (2019). Optimal Model Averaging Estimation for Partially Linear Models. *Statistica Sinica* 29(2), 693–718.
- Zhang, X., D. Yu, G. Zou, and H. Liang (2016). Optimal Model Averaging Estimation for Generalized Linear Models and Generalized Linear Mixed-Effects Models. *Journal of the American Statistical*

REFERENCES

Association 111(516), 1775–1790.

Zhang, X., G. Zou, and R. J. Carroll (2015). Model Averaging Based on Kullback-Leibler Distance. *Statistica Sinica* 25, 1583–1598.

Zou, H., T. Hastie, and R. Tibshirani (2007). On the “degrees of freedom” of the lasso. *The Annals of Statistics* 35(5), 2173–2192.

UCLouvain, Institute of Statistics, Biostatistics and Actuarial Sciences

E-mail: (lise.leonard@uclouvain.be)

UCLouvain, Institute of Statistics, Biostatistics and Actuarial Sciences

E-mail: (eugen.pircalabelu@uclouvain.be)

UCLouvain, Institute of Statistics, Biostatistics and Actuarial Sciences

E-mail: (rainer.vonsachs@uclouvain.be)