# High-dimensional inference in misspecified linear models

### Peter Bühlmann and Sara van de Geer

March 24, 2015

#### Abstract

We consider high-dimensional inference when the assumed linear model is misspecified. We describe some correct interpretations and corresponding sufficient assumptions for valid asymptotic inference of the model parameters, which still have a useful meaning when the model is misspecified. We largely focus on the de-sparsified Lasso procedure but we also indicate some implications for (multiple) sample splitting techniques. In view of available methods and software, our results contribute to robustness considerations with respect to model misspecification.

## 1 Introduction

The construction of confidence intervals and statistical hypothesis tests is a primary goal for assessing uncertainty in high-dimensional inference. Most of the recent contributions for this task discuss some methods and approaches for high-dimensional linear models (Bühlmann, 2013; Zhang and Zhang, 2014; van de Geer et al., 2014; Javanmard and Montanari, 2014; Meinshausen, 2015; Foygel Barber and Candès, 2014), but generalized linear models (Meinshausen et al., 2009; Minnier et al., 2011; van de Geer et al., 2014), undirected graphical models (Ren et al., 2015; Jankova and van de Geer, 2014), instrumental variable models (Belloni et al., 2012) or very general models (Meinshausen and Bühlmann, 2010) have been considered as well, and all of these latter references cover linear models as special case. Another philosophy for inference in the high-dimensional setting is based on selective inference (Benjamini and Yekutieli, 2005; Lockhart et al., 2014; Taylor et al., 2014), but we do not consider this here. Our goal is to interpret and analyze the meaning of inference procedures when the linear model is misspecified. We address this issue in greater detail for the de-sparsified (or de-biased) Lasso (Zhang and Zhang, 2014), but we make a few more general comments in Section 6.1.

More concretely, we describe the correct interpretations and corresponding (sufficient) assumptions which guarantee valid asymptotic inference for the parameters in a high-dimensional, misspecified linear model. That is, we assume that the data is generated from an underlying true nonlinear model  $Y = f(X) + \xi$  but we fit the wrong linear model  $Y = X\beta^0 + \varepsilon$  to the data; see for example Wasserman (2014) who describes such settings as "weak modeling". Precise definitions of the models are given later. Some arising questions are: first, what is the interpretation of  $\beta^0$ ; and secondly, is the standard de-sparsified Lasso procedure valid for construction of statistical hypothesis tests and confidence intervals for the components  $\beta_j^0$  (j = 1, ..., p). Regarding the first issue, it is important to distinguish between random and fixed design scenarios. Regarding the second point, we do give sufficient conditions for asymptotic correctness of the de-sparsified Lasso procedure, although for the random design case, one has to estimate the asymptotic variance differently than for correctly specified models.

The novelty of this work is that we explicitly discuss the implications of linear model misspecification for construction of confidence intervals and hypothesis testing in high dimensions. We believe that this is a missing piece which should be addressed and which is informally often treated according to the folklore that the procedure leads to inference for the "best projected regression parameters": we make this precise and also show that some modifications are necessary for the random design case (see above). The latter are implemented in the statistical R-software package hdi (Meier et al., 2014) which includes various methods for frequentist high-dimensional inference (Dezeure et al., 2014).

# 2 The de-sparsified Lasso for potentially misspecified linear models

We consider n data points  $(Y^{(1)}, X^{(1)}), \ldots, (Y^{(n)}, X^{(n)})$  with univariate responses  $Y^{(i)}$  and p-dimensional covariables  $X^{(i)}$ . We denote by  $Y = (Y^{(1)}, \ldots, Y^{(n)})^T$  and  $X_j = (X_j^{(1)}, \ldots, X_j^{(n)})^T$   $(j = 1, \ldots, p)$  the  $n \times 1$  vectors, and by  $\mathbf{X} = (X_1, \ldots, X_p)$  the  $n \times p$  design matrix.

We fit a potentially misspecified linear model

$$Y = \mathbf{X}\beta^0 + \varepsilon,\tag{1}$$

where the model assumptions are as follows: i.i.d. distributed rows of  $\mathbf{X}$  (if random), and i.i.d. components of  $\varepsilon$  having mean zero, variance  $\sigma_{\varepsilon}^2$  and which are uncorrelated from  $\mathbf{X}$ . In a misspecified setting, the meaning of the parameter vector  $\beta^0$  and of the errors  $\varepsilon$  depends on the context, in particular whether the design is random or fixed. The different interpretations are presented in Sections 3 and 4 below.

For constructing confidence intervals and hypothesis tests for the individual parameters  $\beta_j^0$  (j = 1, ..., p), we consider the de-sparsified Lasso, originally proposed by Zhang and Zhang (2014). The procedure is as follows. First, do a Lasso (Tibshirani, 1996) or square root Lasso (Belloni et al., 2011) regression fit of  $X_j$  versus all other variables from  $\mathbf{X}_{-j}$ , the  $n \times (p-1)$  design matrix whose columns correspond to the variables  $\{X_k; k \neq j\}$ . That is, for the Lasso,

$$\hat{\gamma}_j = \operatorname{argmin}_{\gamma \in \mathbb{R}^{p-1}} \left( \|X_j - \mathbf{X}_{-j}\gamma\|_2^2 / n + \lambda_X \|\gamma\|_1 \right). \tag{2}$$

or using the square root Lasso,

$$\hat{\gamma}_j = \operatorname{argmin}_{\gamma \in \mathbb{R}^{p-1}} \left( \|X_j - \mathbf{X}_{-j}\gamma\|_2 / \sqrt{n} + \lambda_X \|\gamma\|_1 \right). \tag{3}$$

The residuals of such a regression are denoted by

$$Z_j = X_j - \mathbf{X}_{-j}\hat{\gamma}_j.$$

We then project the response Y onto this residual vector: if the model (1) were correct, we have

$$\frac{Z_j^T Y}{Z_j^T X_j} = \beta_j^0 + \sum_{k \neq j} \frac{Z_j^T \mathbf{X}_k}{Z_j^T X_j} \beta_k^0 + \frac{Z_j^T \varepsilon}{Z_j^T X_j}.$$

This suggests a bias correction as follows. Pursue a Lasso regression of Y versus X:

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}_{\boldsymbol{\beta}} \left( \| \boldsymbol{Y} - \mathbf{X} \boldsymbol{\beta} \|_2^2 / n + \lambda \| \boldsymbol{\beta} \|_1 \right),$$

plug it into the bias term and subtract the estimated bias. This leads to the de-sparsified Lasso estimator:

$$\hat{b}_{j} = \frac{Z_{j}^{T} Y}{Z_{j}^{T} X_{j}} - \sum_{k \neq j} \frac{Z_{j}^{T} \mathbf{X}_{k}}{Z_{j}^{T} X_{j}} \hat{\beta}_{k} \quad (j = 1, \dots, p).$$
(4)

From the construction and assuming that model (1) is correct, we heuristically obtain:

$$\frac{Z_j^T X_j}{\sqrt{n\omega_{p;jj}}} (\hat{b}_j - \beta_j^0) = \sum_{k \neq j} \frac{Z_j^T \mathbf{X}_k}{\sqrt{n\omega_{p;jj}}} (\hat{\beta}_k - \beta_k^0) + \frac{Z_j^T \varepsilon}{\sqrt{n\omega_{p;jj}}} \approx \frac{Z_j^T \varepsilon}{\sqrt{n\omega_{p;jj}}} \approx \mathcal{N}(0, 1),$$

where we assume for the first approximation that the error in estimating the bias is negligible, and where  $\omega_{p;jj}^2$  is the asymptotic variance of  $Z_j^T \varepsilon / \sqrt{n}$ . This reasoning has been made rigorous in earlier work, assuming some conditions (Zhang and Zhang, 2014; van de Geer et al., 2014). When the model (1) is wrong, however, the heuristics above needs to be justified anew. Also from a practical point of view, we need to characterize the meaning for  $\beta^0$  and we need to determine the correct specification of  $\omega_{p;jj}^2$  in order to construct asymptotically correct confidence intervals and tests. The details are described in the following Sections 3 and 4.

The procedure for the de-sparsified Lasso  $b_j$  in (4) remains (essentially) the same regardless whether the linear model is correct or not. Referring to the parenthesis in the previous sentence, what potentially changes relative to a correctly specified model is the proper asymptotic variance  $\omega_{p;jj}^2$ , see Section 3.1, and this new feature is now also implemented in the R-software package hdi (Meier et al., 2014).

Throughout the paper, the asymptotic statements are for the setting where the dimension  $p=p_n$  is allowed to depend on n (and hence also the random variables in the model), and we consider the behavior as  $n \to \infty$ , typically with  $p=p_n \to \infty$  at a much faster rate than n. We often suppress the index n in the notation.

# 3 Random design model

Consider the true model

$$Y^{(0)} = f^{0}(X^{(0)}) + \xi^{(0)}, \tag{5}$$

where  $\xi^{(0)}$  is independent of  $X^{(0)}$  with  $\mathbb{E}[\xi^{(0)}] = 0$ . For simplicity, we assume that  $\mathbb{E}[f^0(X^{(0)})] = 0$  as well as  $\mathbb{E}[X^{(0)}] = 0$ , and that furthermore the second moments of  $X^{(0)}$  and  $Y^{(0)}$  exist. We assume that the data are realizations of  $(Y^{(1)}, X^{(1)}), \dots, (Y^{(n)}, X^{(n)})$  of i.i.d. copies of  $(Y^{(0)}, X^{(0)})$  from model (5).

Consider the linear projection

$$Y^{(0)} = (X^{(0)})^T \beta^0 + \varepsilon^{(0)},$$
  

$$\beta^0 = \operatorname{argmin}_{\beta} \mathbb{E} |f^0(X^{(0)}) - (X^{(0)})^T \beta|^2,$$
(6)

where, due to the projection property,  $\mathbb{E}[\varepsilon^{(0)}X^{(0)}] = \operatorname{Cov}(\varepsilon^{(0)}, X^{(0)}) = 0$ . We denote the support of  $\beta_0$  by  $S_0 = \{j; \ \beta_j^0 \neq 0\}$ . While  $\mathbb{E}[\varepsilon^{(0)}] = 0$  we typically have that  $\mathbb{E}[\varepsilon^{(0)}|X^{(0)}] \neq 0$ , because  $\mathbb{E}[\varepsilon^{(0)}|X^{(0)}] = f^0(X^{(0)}) - X^{(0)}\beta^0$ . Thus, when conditioning on  $X^{(0)}$  the assumption of zero mean

for the error is not valid. However, when the inference for  $\beta^0$  is unconditional (not conditioning on  $X^{(0)}$ ), then we have zero mean for the error: therefore, due to model misspecification, the inference with random design should always be *unconditional* on  $X^{(0)}$ .

We note that  $\beta^0$  still has interesting model-free (and well known) interpretations such as: the jth component  $\beta_j^0 = L_j \cdot \operatorname{Parcorr}(Y^{(0)}, X_j^{(0)} | \{X_k^{(0)}; k \neq j\})$  equals the partial correlation between  $Y^{(0)}$  and  $X_j^{(0)}$  given all other variables, up to a constant  $L_j = \sqrt{K_{jj}/K_{YY}}$ , where  $K^{-1}$  is the  $(p+1) \times (p+1)$  covariance matrix of (Y,X); thus,  $\beta_j^0$  measures the linear effect of  $X_j$  on Y after adjusting for the linear effects of all other variables  $X_k$   $(k \neq j)$  on Y. In addition, for Gaussian design, we have the following important interpretation: if  $\beta_j^0 \neq 0$ , then the variable  $X_j^{(0)}$  is in the active set (i.e., relevant) of the nonlinear true function  $f^0$ , see Proposition 3.

We consider here a concrete set of assumptions for Theorem 1 below. Denote by

$$\gamma_j^0 = \operatorname{argmin}_{\gamma} \mathbb{E} |X_j^{(0)} - \sum_{k \neq j} \gamma_{j,k} X_k^{(0)}|^2,$$

$$Z_j^{(0)} = X_j^{(0)} - \sum_{k \neq j} \gamma_{j,k}^0 X_k^{(0)}$$

the population regression vector and residual variables when regressing the random variable  $X_j^{(0)}$  on all other variables  $\{X_k^{(0)}; k \neq j\}$ . It is well known that  $\gamma_j^0 = -(\Sigma^{-1})_{\bullet j}/(\Sigma^{-1})_{jj}$ , where  $(\Sigma^{-1})_{\bullet j}$  denotes the *j*th column vector of  $\Sigma^{-1}$  (assuming it exists, see (A1)).

## Assumptions.

The covariables are such that:

- (A1)  $Cov(X^{(0)}) = \Sigma$  has smallest eigenvalue  $\Lambda^2_{\min}(\Sigma) \geq C_1 > 0$ ;
- (A2)  $\max_{j} ||X_{j}^{(0)}||_{\infty} \le C_{2} < \infty;$
- (A3)  $||Z_j^{(0)}||_{\infty} \le C_3 < \infty;$
- (A4) We have either:
  - (a)  $\|\gamma_j^0\|_1 = o(\sqrt{n/\log(p)}), \|\gamma_j^0\|_r^r = o((n/\log(p))^{\frac{1-r}{2}}\log(p)^{-1/2})$  for 0 < r < 1, and the maximal eigenvalue of  $\mathbf{X}_{S_j}^T\mathbf{X}_{S_j}/n$  satisfies  $\hat{\Lambda}_{\max}^2(S_j) = O_P(1)$ , where  $\mathbf{X}_{S_j}$  denotes the submatrix of the design with columns corresponding to  $S_j = \{k; \gamma_{j,k}^0 \neq 0\}$ ;

or

(b) 
$$s_j = |S_j| = ||\gamma_j^0||_0^0 = \sum_{k \neq j} I((\Sigma^{-1})_{jk} \neq 0) = o(\sqrt{n}/\log(p)).$$

Regarding the structure of the regression:

- (A5) The sparsity satisfies either:
  - (a)  $\|\beta^0\|_1 = o(\sqrt{n/\log(p)})$ ,  $\|\beta^0\|_r^r \hat{\Lambda}_{\max}^r(S_0) = o_P\left((n/\log(p))^{\frac{1-r}{2}}\log(p)^{-1/2}\right)$  for 0 < r < 1, and the maximal eigenvalue of  $\mathbf{X}_{S_0}^T \mathbf{X}_{S_0}/n$  satisfies  $\hat{\Lambda}_{\max}^2(S_0) = O_P(1)$ , where  $\mathbf{X}_{S_0}$  denotes the submatrix of the design with columns corresponding to  $S_0 = \{j; \beta_j^0 \neq 0\}$ ;

or

(b) 
$$s_0 = |S_0| = \|\beta^0\|_0^0 = \sum_{i=1}^p I(\beta_i^0 \neq 0) = o(\sqrt{n}/\log(p)).$$

- (A6) For the second moment  $\omega_{p;jj}^2 := \mathbb{E}|\varepsilon^{(0)}Z_j^0|^2$ :  $\omega_{p;jj}^2 \geq C_4$  for some constant  $C_4 > 0$ . (The existence of  $\omega_{p;jj} < \infty$  is implied by (A3) and (A7)).
- (A7) The error satisfies one of the following conditions:
  - (a)  $|\varepsilon^{(0)}| \leq V$ , where V is a fixed random variable (not depending on p) with  $\mathbb{E}|V|^2 < \infty$ ;
  - (b)  $\mathbb{E}|\varepsilon^{(0)}|^{2+\delta} \leq C_5 < \infty$  for some  $\delta > 0$ .

Either of the conditions implies that for some constant  $C_6 < \infty$ ,  $\mathbb{E}|\varepsilon^{(0)}|^2 \le C_6 < \infty$ .

The assumptions (A2) and (A3) are somewhat restrictive (see also (B1) in van de Geer et al. (2014)). Assumption (A3) is implied by (A2) and assuming that  $\|\gamma_j^0\|_1$  is bounded. Examples where (A7) holds are discussed in Section 3.2. Regarding the assumptions (A4) and (A5) we first note that:

- (A4) can be replaced by (D2) in Section 7.1.1,
- (A5) can be replaced by (D3) in Section 7.1.1,

see Section 7.1 and Lemma 2. Furthermore, for  $\ell_r$  sparsity in (A4,a) and (A5,a), the condition on the maximal eigenvalue can be relaxed by requiring for e.g. (A5,a) that

$$S^* = \{j; |\beta_j^0| > C\sqrt{\log(p)/n}/\hat{\Lambda}_{\max}(S_0)\},$$

for some  $0 < C < \infty$ , has cardinality  $S^* = o(n/\log(p))$ ; and analogously for condition (A4,a). Requiring some sparsity for the design as in (A4) is due to our proof of Proposition 8: this is in contrast for fixed design, where no sparsity condition on the design is needed when using the nodewise square root Lasso in (3) (see Theorem 2). Finally, a sparsity assumption as in (A5) is typical for the de-sparsified Lasso (Zhang and Zhang, 2014; van de Geer et al., 2014; van de Geer, 2014).

**Theorem 1.** Consider the de-sparsified Lasso in (4) with (2) or (3), and the parameter  $\beta^0$  in (6) induced by the random design model (5). Assume (A1)-(A7). If  $\lambda = D_1 \sqrt{\log(p)/n}$  and  $\lambda_X = D_2 \sqrt{\log(p)/n}$  for  $D_1, D_2$  sufficiently large, then:

$$\sqrt{n} \frac{Z_j^T X_j / n}{\omega_{p;jj}} (\hat{b}_j - \beta_j^0) \Rightarrow \mathcal{N}(0,1) \ (n \to \infty),$$

where  $\omega_{p;jj}^2 = \mathbb{E}|\varepsilon^{(0)}Z_j^0|^2$ .

A proof is given in Section 7. The representation of the normalization factor should facilitate to recognize its order of magnitude  $\sqrt{n}$ . For construction of confidence intervals and hypothesis tests we need to consistently estimate the quantity  $\omega_{p;jj}$ : this is discussed in the following Section 3.1.

**Remark 1.** If the assumptions in (A3), (A4) and (A6) hold uniformly in j, we can rephrase the statement of Theorem 1 as follows:

$$\frac{Z_j^T X_j}{\sqrt{n\omega_{p;jj}}} (\hat{b}_j - \beta_j^0) = \Delta_j + W_j,$$

$$\max_{j=1,\dots,p} |\Delta_j| = o_P(1), \ W_j \Rightarrow \mathcal{N}(0,1)$$

#### 3.1 Estimation of the variance

We can estimate  $\omega_{p;jj}^2 = \mathbb{E}|\varepsilon^{(0)}Z_j^0|^2$  by the empirical variance of  $\hat{\varepsilon}_i Z_{j;i}$ ,

$$n^{-1} \sum_{i=1}^{n} (\hat{\varepsilon}_i Z_{j;i} - n^{-1} \sum_{r=1}^{n} \hat{\varepsilon}_r Z_{j;r})^2, \quad \hat{\varepsilon} = Y - \mathbf{X} \hat{\beta}.$$

**Proposition 1.** Consider the random design model (5) with the projected parameter  $\beta^0$  in (6). Assume (A1), (A2), (A3),  $\|\beta^0\|_1 = o(\sqrt{n/\log(p)})$  (which is part of assumption (A5)), (A6), (A7) and (D2) from Section 7 (the latter is implied by the additional assumption (A4)). Then,

$$\hat{\omega}_{p;jj}^2/\omega_{p;jj}^2 = 1 + o_P(1).$$

A proof is given in Section 7. We have as an estimate of the normalizing factor in Theorem 1 the following expression:

$$\frac{Z_j^T X_j}{\sqrt{n\hat{\omega}_{n:jj}}},\tag{7}$$

corresponding to the "sandwich formula" in the case with p < n (Eicker, 1967; Huber, 1967; White, 1980; Freedman et al., 1981).

In particular the formula in (7) is different than the usual expression for correctly specified high-dimensional linear models, used in van de Geer et al. (2014),

$$\frac{Z_j^T X_j}{\|Z_j\|_2 \hat{\sigma_\varepsilon}},\tag{8}$$

where  $\hat{\sigma}_{\varepsilon}^2$  is an estimate of the error variance  $\sigma_{\varepsilon}^2$ , e.g.,  $\hat{\sigma}_{\varepsilon}^2 = n^{-1} \sum_{i=1}^n (\hat{\varepsilon}_i - n^{-1} \sum_{r=1}^n \hat{\varepsilon}_r)^2$  with  $\hat{\varepsilon} = Y - \mathbf{X}\hat{\beta}$ . While the formula in (8) is asymptotically valid for correctly specified models, the analogue in (7) is robust and valid irrespective whether the model is correct or not. The expression in (7) is now also implemented in the R-software package hdi.

# 3.2 Sparsity of the projection and implications on the error $\varepsilon^{(0)}$

The statement in Theorem 1 depends, among other conditions, on assumptions (A5)-(A7) which are depending on the projection of the nonlinear to a linear model. In particular, (A5) requires sparsity of the projected parameter vector: even if the underlying true nonlinear regression function depends only on a few covariables, the projected parameter  $\beta^0$  in (6) is not necessarily sparse. We provide here some sufficient conditions ensuring a sparse  $\beta^0$ .

Throughout this subsection,  $\beta^0$  is as in (6). We know that

$$\beta^{0} = \Sigma^{-1}\Gamma,$$

$$\Sigma = \text{Cov}(X^{(0)}), \ \Gamma = (\text{Cov}(f^{0}(X^{(0)}), X_{1}^{(0)}), \dots, \text{Cov}(f^{0}(X^{(0)}), X_{p}^{(0)})^{T}.$$

Therefore,

$$\beta_j^0 = \sum_{\ell=1}^p (\Sigma^{-1})_{j\ell} \Gamma_\ell. \tag{9}$$

Denote by  $\|\Sigma^{-1}\|_{\infty} = \max_{jk} |(\Sigma^{-1})_{jk}|$  and by  $(\Sigma^{-1})_{\bullet\ell}$  the  $\ell$ th column of  $\Sigma^{-1}$ , and generally by  $\|u\|_0^0 = \sum_{r=1}^d I(u_r \neq 0)$  the  $\ell_0$ -sparsity of a d-dimensional vector u.

**Proposition 2.** Consider the random design model (5) with the projected parameter  $\beta^0$  in (6). Assume that  $\Sigma$  is positive definite (but not requiring bounds on its eigenvalues). The following holds:

1.  $\ell_r$ -sparsity for  $0 < r \le 1$ :

$$\|\beta^0\|_r \le \max_{\ell} \|(\Sigma^{-1})_{\bullet\ell}\|_r \|\Gamma\|_r,$$

which implies, for  $s_{\ell} = \|\gamma_{\ell}^0\|_0^0 = \sum_{k \neq \ell} I((\Sigma^{-1})_{k\ell} \neq 0)$ ,

$$\|\beta^0\|_r \le (\max_{\ell} s_{\ell} + 1)^{1/r} \|\Sigma^{-1}\|_{\infty} \|\Gamma\|_r.$$

2.  $\ell_0$ -sparsity:

$$\|\beta^0\|_0^0 \le \sum_{\ell \in S_{\Gamma}} (s_{\ell} + 1), \ S_{\Gamma} = \{j; \ \Gamma_j \ne 0\},$$

which implies

$$\|\beta^0\|_0^0 \le (\max_{\ell} s_{\ell} + 1) \|\Gamma\|_0^0.$$

A proof is given in Section 7. As an example, consider the case where  $\Sigma$  is block-diagonal with maximal block-size equal to  $b_{\text{max}}$ . We then have that  $\max_{\ell} s_{\ell} + 1 = b_{\text{max}}$  and hence by Proposition 2:

$$\|\beta^0\|_r \le b_{\max}^{1/r} \|\Sigma^{-1}\|_{\infty} \|\Gamma\|_r \quad (0 < r \le 1),$$
$$\|\beta^0\|_0^0 \le b_{\max} \|\Gamma\|_0^0.$$

Block dependence. Assume now that the predictor variables exhibit block dependence with blocks corresponding to the associated block-diagonal covariance matrix  $\Sigma$ . That is, there are blocks of variables, where the variables from different blocks are (jointly) independent, and these blocks induce a block-diagonal covariance matrix. Denote by  $S_{f^0} \subseteq \{1, \ldots, p\}$  the support of  $f^0(\cdot)$  which contains all the variables which have an influence in  $f^0(\cdot)$ .

Corollary 1. Assume the conditions of Proposition 2. In addition, assume block dependence with maximal block-size equal to  $b_{max}$ . We have that

$$\|\Gamma\|_0^0 \le b_{\max} |S_{f^0}|,$$

and, due to Proposition 2,

$$\|\beta^0\|_0^0 \le b_{\max}^2 |S_{f^0}|.$$

A proof is given in Section 7.

Proposition 2 and Corollary 1 obviously lead to justifications of the assumption on the sparsity  $s_0$  in (A5), but also for the conditions in (A7). Regarding the latter: if  $\|\beta^0\|_1 \le C_9 < \infty$  (which is implied by  $\|\beta^0\|_0^0$  bounded and  $\max_j |\beta_j^0|$  bounded) and assuming (A2) we have that

$$|\varepsilon^{(0)}| = |Y^{(0)} - (X^{(0)})^T \beta^0| \le |Y^{(0)}| + C_9 C_2.$$

Thus, assuming either  $|Y^{(0)}| \leq V$  for some fixed random variable V with  $\mathbb{E}|V|^2 < \infty$  or  $\mathbb{E}|Y^{(0)}|^{2+\delta} \leq M_3 < \infty$  (which are both rather weak assumptions) implies either (A7,a) or (A7,b), respectively.

#### 3.3 Gaussian design

The bound in Proposition 8 and Corollary 1 for  $\ell_0$ -sparsity can be much improved when assuming that  $X^{(0)}$  has a joint Gaussian distribution. This is in conflict with assumption (A2). However, for the case with Gaussian design, thereby dropping (A2) and (A3), it would be easier to derive the statements from Theorem 1 and Proposition 1.

**Proposition 3.** Consider the random design model (5) with the projected parameter  $\beta^0$  in (6). Assume that  $X^{(0)}$  has a joint Gaussian distribution with positive definite covariance matrix  $\Sigma$  (but not requiring bounds on its eigenvalues). Then,

$$S_0 \subseteq S_{f^0}$$
.

A proof is given in Section 7. This is an important result saying that if we infer a variable as an active variable (significantly different from zero) in the misspecified linear model, it must be an active variable in the nonlinear true model.

To make further statements, we represent the function  $f^0$  as follows:

$$f^{0}(x) = \sum_{k=1}^{d} f_{k}^{0}(x_{S_{k}}),$$
  
$$\{S_{1}, \dots, S_{d}\} \text{ a partition: } S_{f^{0}} = \bigcup_{k=1}^{d} S_{k}, \ S_{k} \cap S_{\ell} = \emptyset \ (k \neq \ell),$$

where  $x_A$  denotes the subvector of x with components in  $A \subseteq \{1, \ldots, p\}$  and  $\mathbb{E}[f_k^0(X_{S_k})] = 0$ ; and the partition is finest in the sense that the representation of  $f^0$  is given with the  $S_k$ 's of smallest possible cardinality. For example, for the function considered in Section 5

$$f^{0}(x) = -5 + 5\sin(\pi x_{1}x_{2}) + 4(x_{3} - 0.5)^{2} + 2x_{5} + x_{6},$$
(10)

we have the partition  $S_1 = \{1, 2\}, S_2 = \{3\}, S_3 = \{5\}, S_4 = \{6\}.$ 

**Proposition 4.** Consider the random design model (5) with the projected parameter  $\beta^0$  in (6). Assume that  $X^{(0)}$  has a joint Gaussian distribution with positive definite covariance matrix  $\Sigma$  (but not requiring bounds on its eigenvalues). Consider the projected parameter in the submodel with variables from  $S_k$  ( $k \in \{1, ..., d\}$ ):

$$\tilde{\beta}(S_k) = \operatorname{argmin}_{\beta \in \mathbb{R}^{|S_k|}} \mathbb{E} |f_k^0(X_{S_k}^{(0)}) - (X_{S_k}^{(0)})^T \beta|^2.$$

For  $j \in S_k$  we denote by c(j) the index of the component in  $\tilde{\beta}(S_k)$  which corresponds to variable  $X_j^{(0)}$ . Then,

$$\beta_j^0 = \tilde{\beta}_{c(j)}(S_k),$$

saying that we can infer  $\beta_j^0$  with  $j \in S_k$  from the submodel with variables  $X_{S_k}^{(0)}$ .

A proof is given in the Appendix. As an example, we consider again  $f^0$  from (10). Proposition 4 then implies:

$$\begin{split} &(\beta_1^0,\beta_2^0)^T = \mathrm{argmin}_{\beta \in \mathbb{R}^2} \mathbb{E} |5\sin(\pi X_1^{(0)} X_2^{(0)}) - (X_1^{(0)},X_2^{(0)})\beta|^2 = (0,0)^T, \\ &\beta_3^0 = \mathrm{argmin}_{\beta \in \mathbb{R}} \mathbb{E} |4(X_3^{(0)} - 0.5)^2 - 5 - X_3^{(0)}\beta|^2 = -4, \\ &\beta_5^0 = \mathrm{argmin}_{\beta \in \mathbb{R}} \mathbb{E} |2X_5^{(0)} - X_5^{(0)}\beta|^2 = 2, \\ &\beta_5^0 = \mathrm{argmin}_{\beta \in \mathbb{R}} \mathbb{E} |X_5^{(0)} - X_6^{(0)}\beta|^2 = 1, \end{split}$$

and all  $\beta_j^0 = 0$  for  $j \notin S_{f^0}$ . For the numerical values of  $\beta_1^0, \beta_2^0$  and  $\beta_3^0$ , we used that  $X^{(0)}$  has mean zero.

# 4 Fixed design model

Consider the model as in (5) but now with fixed design:

$$Y^{(i)} = f^{0}(X^{(i)}) + \xi^{(i)}, \ i = 1, \dots, n,$$
(11)

where  $\xi^{(1)}, \dots, \xi^{(n)}$  are i.i.d. with  $\mathbb{E}[\xi^{(i)}] = 0$  and  $\mathbb{E}|\xi^{(i)}|^2 = \sigma^2$ . As before, we denote the  $n \times p$  design matrix by  $\mathbf{X}$  and the  $n \times 1$  response vector by  $Y = (Y^{(1)}, \dots, Y^{(n)})^T$ . We assume that  $\mathrm{rank}(\mathbf{X}) = n \leq p$  and thus, we can always represent the vector  $\mathbf{f}^0 = (f^0(X^{(1)}), \dots, f(X^{(n)}))^T$  as  $\mathbf{X}\beta^{\dagger}$ . The vector  $\beta^{\dagger}$  is not unique, but we can look for some sparsest solution. We consider the basis pursuit solution (Chen et al., 1998), known also as the solution from compressed sensing (Candès and Tao, 2006; Donoho, 2006):

$$\beta^0 = \operatorname{argmin}_{\beta} \{ \|\beta\|_1; \mathbf{X}\beta = \mathbf{f}^0 \}. \tag{12}$$

Thus, the model in (11) is correctly specified as a linear model

$$Y = \mathbf{X}\beta^0 + \varepsilon \text{ with } \beta^0 \text{ as in (12)}, \tag{13}$$

where  $\varepsilon = (\xi_1, \dots, \xi_n)^T$ . In particular, due to correct specification, the interpretation of  $\beta^0$  is standard.

We refer to this  $\beta^0$  in (12) throughout this section (unless stated otherwise). We assume the following:

**(B1)** 
$$\lambda_X \simeq \sqrt{\log(p)/n} \text{ and } ||Z_j||_2^2/n \ge C > 0;$$

**(B2)** 
$$\|\hat{\beta}(\lambda) - \beta^0\|_1 = o_P(1/\sqrt{\log(p)}).$$

We justify these assumptions below.

**Theorem 2.** Consider the de-sparsified Lasso in (4) with (2) or (3), and the fixed design model (11) with rank( $\mathbf{X}$ ) = n and linear representation as in (13) with  $\beta^0$  as in (12). Assume either Gaussian errors or condition (A7) and assume that  $\sigma^2 \geq L > 0$ . Suppose that (B1) and (B2) hold when using the nodewise Lasso (2), or only (B2) when using the nodewise square root Lasso (3). Then

$$\frac{Z_j^T X_j}{\sigma \|Z_i\|_2} (\hat{b}_j - \beta_j^0) \Rightarrow \mathcal{N}(0, 1).$$

Proof: This follows from van de Geer et al. (2014, Th.2.1) for Gaussian errors. For non-Gaussian errors, we invoke the Lindeberg condition and proceed as for the proof of Theorem 1 (Proposition 7).

We argue first that (B1) holds with high probability. Assume the following.

Consider the setting where the rows of **X** arise as fixed i.i.d. realizations of a p-dimensional random variable X with covariance matrix  $\Sigma$ .

- (C1) (i)  $0 < C_7 \le 1/(\Sigma^{-1})_{jj} = \mathbb{E}|Z_j^{(0)}|^2 \ge C_8 < \infty$  (the upper bound is implied by (A3); the lower bound is the analogue of (A6));
  - (ii)  $\max_{j} ||X_{j}||_{\infty} \le C_{2} < \infty$  (which is assumption (A2));
  - (iii)  $\|\gamma_i^0\|_1 = o(\sqrt{n/\log(p)})$  (which is part of the assumption (A4a)).
- (C2) (A1), (A2), (A5) and (A7).

**Proposition 5.** (for nodewise Lasso only) Assume that (C1) holds. Then, for  $\lambda_X = D_2 \sqrt{\log(p)/n}$  with  $D_2$  sufficiently large, assumption (B1) holds with probability tending to one.

A proof is given in Section 7.

**Proposition 6.** Consider the fixed design model (11) having a linear representation as in (13) with  $\beta^0$  as in (12). Assume that (C2) holds. Then, for  $\lambda = D_1 \sqrt{\log(p)/n}$  with  $D_1$  sufficiently large, assumption (B2) holds with probability tending to one.

Proof. The statement can be derived as in the proof of statement 2 in Lemma 2 in Section 7.□

Sparse solutions and misspecification. We note that for a fixed design linear model, misspecification with respect to the linearity in the unknown parameters cannot happen. The same is true when conditioning on the covariables X. In this scenario, we do not need to employ the "sandwich" variance formula in (7) but we can use the more standard expression from (8). What is important though is the interpretation of the parameter  $\beta^0$  and of the output of the de-sparsified Lasso: the inferential statements are valid for a sparse approximation. We focused here on the choice of the basis pursuit solution in (12) which is perhaps among the simplest and which can be computed. But in fact, any solution of  $\mathbf{X}\beta = \mathbf{f}^0$  satisfying assumption (B2) is good enough: or in view of Proposition 6, any solution which is weak  $\ell_r$ - (0 < r < 1) or  $\ell_0$ -sparse, see (A5), is fine. A confidence interval then means that it covers any sufficiently  $\ell_r$ - and  $\ell_0$ -sparse solution  $\beta^0$  of  $\mathbf{X}\beta = \mathbf{f}^0$ . This itself is a nice and "strong" interpretation of a confidence interval, namely that despite non-uniqueness, it covers all sparse solutions.

## 5 Some empirical results

We consider two non-linear models as in (5) (or versions thereof for fixed design, see Section 5.2). The first one uses a nonlinear regression function from Friedman's (1991) MARS paper but with smaller signal to noise ratio:

(M1)

$$X^{(0)} \sim \mathcal{N}_p(0, \Sigma), \ \Sigma_{j,j} = 1 \ \forall j, \ \Sigma_{3,4} = \Sigma_{4,3} = 0.8, \ \Sigma_{j,k} = 0 \ (j \neq k; \ j,k \notin \{3,4\}),$$
  
 $f^0(x) = -5 + 2\sin(\pi x_1 x_2) + 4(x_3 - 0.5)^2 + 2x_5 + x_6,$   
 $\xi^{(0)} \sim \mathcal{N}(0,1).$ 

(M2)

$$X^{(0)}$$
 as in (M1),  
 $f^{0}(x) = \sin(\pi/2x_1)x_2 + x_3^3/5 + x_5 + x_6/2,$   
 $\xi^{(0)} \sim \mathcal{N}(0, 1).$ 

(M3)

$$X^{(0)} \sim \mathcal{N}_p(0, \Sigma), \ \Sigma_{j,k} = 0.8^{|j-k|}, \ f^0 \text{ as in (M1)}.$$

(M4)

$$X^{(0)} \sim \text{ as in (M3)}, f^0 \text{ as in (M2)}.$$

The intercept -5 in the function  $f^0$  in (M1) and (M3) ensures that  $\mathbb{E}[f^0(X^{(0)})] = 0$ .

## 5.1 Simulations for random design

For random design, the corresponding parameters  $\beta^0$  in (6) are as follows:

for model (M1),(M3): 
$$\beta^0 = (0, 0, -4, 0, 2, 1, 0, \dots, 0)^T$$
  
for model (M2),(M4):  $\beta^0 = (0, 0, 0.6, 0, 1, 0.5, 0, \dots, 0)^T$ .

The values are in accordance with Proposition 4, because of Gaussianity of the design: the active set  $S_0 = \{3, 5, 6\} \subset S_{f^0} = \{1, 2, 3, 5, 6\}$ . Figure 1 displays  $\|\beta\|_r^r$  as a function of r for  $0 \le r \le 1$ . The log-sparsity is approximately a linear function in r, once increasing (for (M1),(M3)) and once

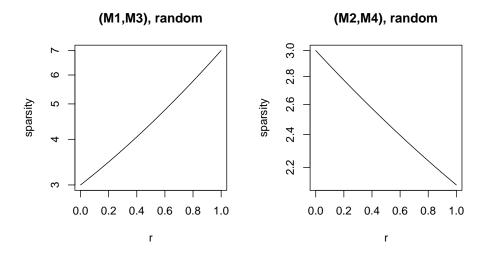


Figure 1: Random design for models (M1),(M3) and (M2),(M4) with p = 1000. Plot of  $\|\beta^0\|_r^r$  (on log-scale) as a function of  $r \in [0,1]$  (r = 0 corresponds to the  $\ell_0$ -sparsity), where  $\beta^0$  is as in (6).

decaying (for (M2),(M4)). Our theory requires either weak  $\ell_r$ -sparsity or  $\ell_0$ -sparsity of  $\beta^0$  (see (A5,a) or (A5,b)) and hence a possibly more realistic assumption than  $\ell_0$ -sparsity alone.

For simulations with random design, we generate n independent data points according to the models (M1)-(M4) where for each realization, we generate the X and  $\xi$  variables anew. We consider the case with sample size n=200 and dimension p=1000. We use the de-sparsified Lasso procedure as described in (4) with the nodewise Lasso (2) and tuning parameters  $\lambda$  and  $\lambda_X$  (the same for all j) from the default in the R-software package hdi (Meier et al., 2014). For estimation of the asymptotic variance we use (7).

Table 1 and Figure 2 report empirical results based on 100 independent simulations. Denoting by  $CI_j$  a confidence interval for  $\beta_j^0$ , the average coverage is

$$\operatorname{avgcov}(S_0) = |S_0|^{-1} \sum_{j \in S_0} \mathbb{P}[\beta_j^0 \in \operatorname{CI}_j],$$

$$\operatorname{avgcov}(S_0^c) = |S_0^c|^{-1} \sum_{j \in S_0^c} \mathbb{P}[\beta_j^0 \in \operatorname{CI}_j],$$
(14)

and the empirical analogue by replacing the probability " $\mathbb{P}$ " by an empirical average over the 100 simulations. We consider the average expected length of the confidence intervals

$$\operatorname{avglen}(S_0) = |S_0|^{-1} \sum_{j \in S_0} \mathbb{E}[\operatorname{length}(\operatorname{CI}_j)],$$

$$\operatorname{avglen}(S_0^c) = |S_0^c|^{-1} \sum_{j \in S_0^c} \mathbb{E}[\operatorname{length}(\operatorname{CI}_j)],$$
(15)

and the empirical analogue by replacing the expectation "E" with an empirical average. The actual

model	avg. coverage $S_0$	avg. coverage $S_0^c$	avg. length $S_0$	avg. length $S_0^c$
(M1)	0.98	0.99	3.01	2.19
(M2)	0.91	0.95	0.48	0.41
(M3)	0.98	0.99	4.18	3.56
(M4)	0.95	0.95	0.70	0.65

Table 1: Random design. Average coverage and average length of confidence intervals (empirical versions of (14) and (15)), for  $S_0$  and  $S_0^c$  separately (note that  $S_0^c = \emptyset$  for (M3) and (M4)). Nominal level equal to 0.95. Sample size n = 200 and dimension p = 1000.

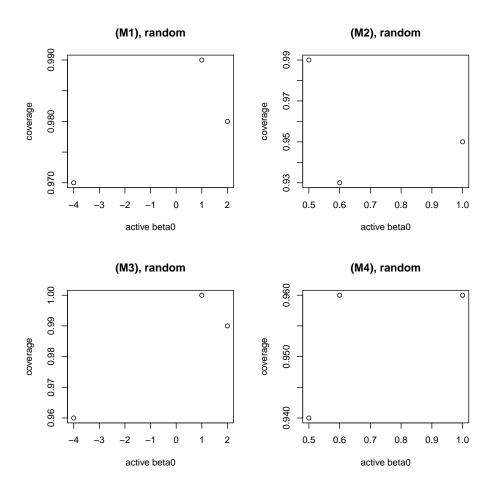


Figure 2: Random design. Coverage as a function of the coefficients  $\beta_j^0$  of the active variables with  $j \in S_0$ . Nominal level equal to 0.95. Sample size n = 200 and dimension p = 1000.

coverage results in Table 1 and the more detailed view given in Figure 2 are very satisfactory. We note that the lengths of the confidence intervals are not constant for the same covariance model for

X. The reason is that at least asymptotically (see Theorem 1), the length depends, among other things, on  $\mathbb{E}|Z_j^{(0)}\varepsilon^{(0)}|^2$ , and the error term  $\varepsilon^{(0)}$  itself depends on the true function  $f^0$ . This is in contrast to fixed design, where the asymptotic length of the confidence intervals is a function of  $\mathbb{E}|Z_j^{(0)}|^2 = 1/(\Sigma^{-1})_{jj}$  and  $\sigma^2 = \mathbb{E}|\xi_i|^2 = \mathbb{E}|\varepsilon_i|^2$  only (see Theorem 2 and formula (17) and (19)).

#### 5.2 Simulations for fixed design

We consider the same models (M1)-(M4) but now with fixed design with n = 200 and p = 1000, where we use a fixed realization of the X variables in the corresponding model. We generate n independent data points according to the models (M1)-(M4) where for each realization, we generate only the  $\xi$  error variables anew.

We note that for all the four models with fixed design we have that  $|S_0| = n = 200$ . Figure 3 displays  $\|\beta^0\|_r^r$  as a function of r for  $0 \le r \le 1$ , where  $\beta^0$  is the basis pursuit solution from (12) and the parameter of interest, for 100 different independent simulation runs. The log-sparsity is approximately a linear decreasing function in r. Even more pronounced here for fixed than random

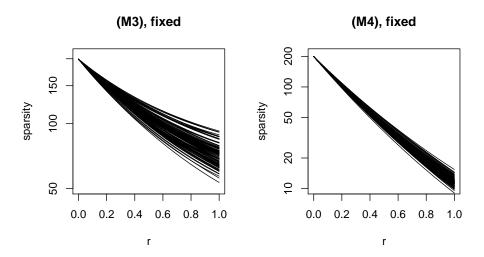


Figure 3: Fixed design for models (M3) and (M4) with n=200 and p=1000. 100 independent realizations and corresponding basis pursuit solutions  $\beta^0$  as in (12): the lines correspond to the 100 different values of  $\|\beta^0\|_r^r$  (on log-scale) as a function of  $r \in [0,1]$  (r=0 corresponds to the  $\ell_0$ -sparsity).

design, we conclude that weak  $\ell_r$ -sparsity, as required by our theory, seems to be a much more realistic assumption than  $\ell_0$ -sparsity which is always equal to n = 200. However, we also see that for model (M3), the parameter  $\beta^0$  is not very  $\ell_r$ -sparse. Thus, it might be difficult that a confidence interval would achieve good coverage, see also Figure 4 and the last paragraph of this section.

We use the de-sparsified Lasso procedure as described in (4) with the nodewise Lasso (2) and tuning parameters  $\lambda$  and  $\lambda_X$  (the same for all j) from the default in the R-software package hdi (Meier et al., 2014). For estimation of the asymptotic variance we use (8). Table 2 and Figure 4 report empirical results for the basis pursuit solution  $\beta^0$  in (12), based on 100 independent

simulations where the design is a fixed realization from the models (M1)-(M4). The actual

model	avg. coverage $S_0$	avg. coverage $S_0^c$	avg. length $S_0$	avg. length $S_0^c$
$\overline{\mathrm{(M1)}}$	0.97	0.98	1.68	1.69
(M2)	0.95	0.97	0.41	0.41
(M3)	0.96	0.97	3.26	3.27
(M4)	0.96	0.96	0.95	0.95

Table 2: Fixed design. Average coverage and average length of confidence intervals (empirical versions of (14) and (15)) for the basis pursuit solution  $\beta^0$  in (12), for  $S_0$  and  $S_0^c$  separately. Nominal level equal to 0.95. Sample size n = 200 and dimension p = 1000.

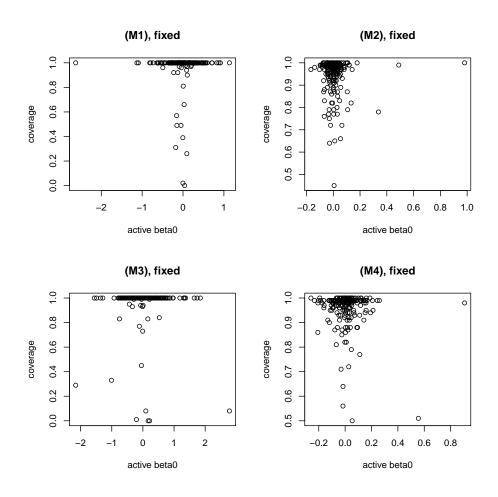


Figure 4: Fixed design. Coverage as a function of the coefficients  $\beta_j^0$  (from basis pursuit in (12)) of the active variables with  $j \in S_0$ . Nominal level equal to 0.95. Sample size n = 200 and dimension p = 1000.

average coverage results in Table 2 are very fine. However, with the more detailed view in Figure 4,

the coverage can be quite poor for a few coefficients although this should be interpreted cautiously, as explained below. The poor coverage is particularly visible for the models (M1) and (M3): a reason might be that the degree of weak  $\ell_r$ -sparsity of the basis pursuit solution  $\beta^0$  in (12) is not as high as for (M2) and (M4) ((shown for (M3), (M4) in Figure 3). Regarding the lengths of the confidence intervals: we cannot confirm the asymptotic behavior saying that they are equal for the same covariance model for the realized  $\mathbf{X}$  and the same error variances (e.g. (M1) and (M2)), regardless of the true underlying nonlinear regression function.

It is important to interpret the obtained confidence intervals as described in the last paragraph of Section 4: any solution of  $\mathbf{X}\beta = \mathbf{f}^0$  which is weak  $\ell_r$ -sparse (0 < r < 1) or  $\ell_0$ -sparse is fine and should be covered by the confidence interval. Our findings in Figure 4 are for the basis pursuit solution only, and the latter is not very sparse (see Figure 3). This doesn't imply though that there isn't another solution  $\beta^0$  which is  $\ell_r$ - or  $\ell_0$ -sparse and whose components would be covered well by the obtained confidence intervals. Unfortunately, the latter statement is uncheckable due to the involved computational complexity; in contrast to the findings for the basis pursuit solution which can be easily computed with a linear program. Therefore, the somewhat negative findings indicated in Figure 4 should be down-weighted.

### 6 Discussion

The current work offers a precise description of interpretation and (sufficient) assumptions for inference in a misspecified high-dimensional linear model. The following Table 3 summarizes the main points with respect to interpretation and modification of the de-sparsified Lasso procedure. A modification of the variance as in (7) is needed for the case of a random design misspecified model. Such a modification seems always advisable for the random design case, as it is consistent irrespective whether the model is correct or not and hence offers some robustness against model misspecification; see for example Huber (1967). The conceptual parts, as indicated in Table 3, will not change for generalized linear models as one can link them to weighted linear regression. One should decide beforehand, whether the inference should be performed with fixed  $\mathbf{X}$  (or conditional on  $\mathbf{X}$ ) or whether  $\mathbf{X}$  is considered as random. The interpretation of the parameter  $\beta^0$  (see Table 3) changes when the true underlying regression function is non-linear, perhaps more dramatically than expected. For the special case of Gaussian random design we have the interesting property that  $S_0 \subseteq S_{f^0}$  (Proposition 3), saying that if a variable is significant in the misspecified linear model, it must be relevant in the true nonlinear model.

#### 6.1 Sample splitting methods

Regarding other methods for construction of p-values and confidence intervals, we briefly discuss sample splitting techniques. Such procedures, including the preferred multiple sample instead of single sample splitting (Meinshausen et al., 2009), can be used for the random design misspecified case. The reason is that the sample splitting device implicitly assumes the same probability distribution in split samples, and this holds for random  $\mathbf{X}$  (but typically not for fixed  $\mathbf{X}$ ) and implies the same projected parameter  $\beta^0$  in (6) in split samples. If the linear model is correct with the same sparse true  $\beta^0$  for every sample point, sample splitting can also be used for fixed design cases (because both split samples are from a fixed design linear model with parameter vector  $\beta^0$ ). However, for the fixed design model as in (13), the issue is different since e.g. the basis pursuit

design	interpretation of $\beta^0$	modification
random design	via projection in (6);	modified variance in (7)
	with model-free interp. described after (6);	
	for Gaussian des.: active set property (Prop. 3)	
fixed design	any sparse solution of $\mathbf{X}\beta = \mathbf{f}^0$	no modification
	(e.g. basis pursuit solution in (12));	
	with standard interp. (since no misspecif.)	

Table 3: Conceptual summary of interpretation and required modification of the de-sparsified Lasso procedure for misspecified high-dimensional linear model. The required assumptions for asymptotic validity of the method are described in Theorems 1 and 2. In case of fixed design where the true underlying regression function is linear with a corresponding sparsest "true" parameter vector  $\beta^0$ , the basis pursuit solution typically coincides with  $\beta^0$  (see compressed sensing literature (Candès and Tao, 2007, cf.)).

solution  $\beta^0$  in (12) would be different for every split sample.

A modification is necessary though for the misspecified random design case: even for low-dimensional inference, which is what is used after screening for variables in the first half of the sample, one has to use a modified estimator for the variance, analogously to the estimator in (7) which is robust against model misspecification.

## 7 Proofs

### 7.1 Proof of Theorem 1 for random design

We prove here the statement of Theorem 1 under slightly weaker assumptions than in condition (A). In this section, **X** is always random and the parameter  $\beta^0$  as in (6).

#### 7.1.1 Preliminary results

We show here that the following conditions hold:

**(D1)** 
$$\max_{k \neq j} |\varepsilon^T \mathbf{X}_k / n| = O_P(\sqrt{\log(p)/n}).$$

(**D2**) For either the nodewise Lasso in (2) or the square root Lasso in (3):  $\|\hat{\gamma}_j(\lambda_X) - \gamma_j^0\|_1 = o_P(1/\sqrt{\log(p)})$ .

**(D3)** 
$$\|\hat{\beta}(\lambda) - \beta^0\|_1 = o_P(1/\sqrt{\log(p)}).$$

**Lemma 1.** For random **X**, assume (A2) and  $\mathbb{E}|\varepsilon^{(0)}|^2 \leq C < \infty$  for some constant C > 0 (the latter is implied by (A7)). Then, (D1) holds, that is:

$$\max_{k \neq j} |\varepsilon^T \mathbf{X}_k / n| = O_P(\sqrt{\log(p)/n}).$$

Proof: Using Nemirovski's inequality (Bühlmann and van de Geer, 2011, Lemma 14.24) we obtain:

$$\mathbb{E}[\max_{1 \le j \le p} |n^{-1} \varepsilon^T X_j|^2] \le 8 \log(2p) C_2^2 C_6 / n = O(\log(p) / n).$$

Thus, since  $\mathbb{E}[\varepsilon^T X_i] = 0$  and using Markov's inequality:

$$\mathbb{P}[\max_{j=1,\dots,p}|n^{-1}\varepsilon^TX_j|>c] \leq \mathbb{E}[\max_{j=1,\dots,p}|n^{-1}\varepsilon^TX_j|]/c \leq \sqrt{\mathbb{E}[\max_{j=1,\dots,p}|n^{-1}\varepsilon^TX_j|^2]}/c = O(\sqrt{\log(p)/n})/c.$$

This completes the proof.

**Lemma 2.** For random X, assume (A1) and (A2).

- 1. Then, for  $\lambda_X = D_2 \sqrt{\log(p)/n}$  with  $D_2$  sufficiently large, (A3) and (A4) imply (D2).
- 2. If  $\mathbb{E}|\varepsilon^{(0)}|^2 \leq C < \infty$  for some constant C > 0 (the latter is implied by (A7)), then for  $\lambda = D_1 \sqrt{\log(p)/n}$  with  $D_1$  sufficiently large, (A5) implies (D3).

Proof: The first and second statement can be proved analogously. For the first one, due to (A3), the error when regressing  $X_j$  versus  $X_{-j} = \{X_k; k \neq j\}$  is bounded.

When invoking the  $\ell_0$ -sparsity assumptions (A4,b) or (A5,b), respectively, we know that the compatibility condition holds with probability tending to one: because of (A1), (A2) and the  $\ell_0$ -sparsity assumption (Bühlmann and van de Geer, 2011, cf. Ch. 6.12)). Therefore, and using Lemma 1, we obtain the statements invoking some oracle inequality for the Lasso (Bühlmann and van de Geer, 2011, cf. Th.6.1) or the square root Lasso (van de Geer, 2014, Th.1.4.2).

When invoking the  $\ell_r$ -sparsity (0 < r < 1) assumptions (A4,a) or (A5,a), respectively, we can use the results from van de Geer (2015, Sec.5) which apply not only for the square root Lasso but also for the Lasso (van de Geer, 2014, cf.Th.1.3.2). We need to argue that the compatibility condition holds with probability tending to one for, e.g. when proving the second statement, the set:

$$S^* = \{j; |\beta_i^0| > C\sqrt{\log(p)/n}/\hat{\Lambda}_{\max}(S_0)\}$$

Due to the assumption on  $\ell_1$ -sparsity and due to the assumption that  $\hat{\Lambda}(S_0)$  is bounded, we have that  $|S^*| = o(n/\log(p))$ . Therefore, due to (A1) and (A2), the compatibility condition holds for  $S^*$  with probability tending to one (Bühlmann and van de Geer, 2011, cf. Ch. 6.12)).

#### 7.1.2 Proof

Denote by  $Z_j^0 = X_j - \mathbf{X}_{-j}\gamma_j^0$ , analogously as in Section 3 but now for  $n \times 1$  vectors. We first analyze the behavior of the part  $Z_j^T \varepsilon/n$ . We have that

$$\mathbb{E}[\varepsilon_i X_{k;i}] = 0 \ \forall k,$$

and hence  $\mathbb{E}[(Z_i^0)^T \varepsilon] = 0$ .

**Proposition 7.** Assume (A1), (A3), (A6) (only that  $\omega_{p;jj} > 0$ ) and (A7). Denote by  $\omega_{p;jj}^2 =$  $\mathbb{E}|\varepsilon^{(0)}Z_i^{(0)}|^2$ . Then:

$$\sqrt{n} \frac{\varepsilon^T Z_j^0/n}{\omega_{\nu;jj}} \Rightarrow \mathcal{N}(0,1) \ (n \to \infty).$$

Note that  $p = p_n$  is allowed to depend on n.

Proof. Denote by  $W_{p;i} = \varepsilon_i Z_{j;i}^0$ . Since  $Cov(\varepsilon_i, X_{k;i}) = 0 \ \forall k$ , we have that  $\mathbb{E}[W_{p;i}] = 0$ . Furthermore,  $W_{p;1}, \ldots, W_{p;n}$  are independent. We verify the Lindeberg condition. For  $\kappa > 0$ ,

$$\lim_{n \to \infty} \frac{1}{\omega_{p;jj}^2} \int_{|W_p| > \kappa \sqrt{n} \omega_{p;jj}} W_p^2 dP = 0.$$

Assuming (A7,a), we invoke the dominated convergence theorem:

$$|W_p|^2 I_{|W_p| > \kappa \sqrt{n\omega_{p;jj}}} \le |W_p|^2 \le |\varepsilon^{(0)}|^2 |Z_j^{(0)}|^2 \le V^2 C_3^2.$$

Because  $I(|W_p| > \kappa \sqrt{n\omega_{p;jj}}) = 0 \ (n \to \infty)$  in probability, and hence

$$|W_p|^2 I_{|W_p| > \kappa \sqrt{n}\omega_{p;jj}} = o_P(1),$$

and because of the dominated convergence theorem we conclude that the Lindeberg condition holds. Assuming (A7,b), we have that  $\mathbb{E}|W_{p,i}|^{2+\delta} \leq \mathbb{E}|\varepsilon_i|^{2+\delta}C_3^{2+\delta} \leq C_5C_3^{2+\delta}$ . The Lindeberg condition is then implied by the Lyapunov theorem.

**Proposition 8.** (with  $Z_j$  instead of  $Z_j^0$ ) Assume (A1), (A3), (A6), (A7), (D1) and (D2). Then:

$$\sqrt{n} \frac{\varepsilon^T Z_j/n}{\omega_{p;jj}} \Rightarrow \mathcal{N}(0,1) \ (n \to \infty).$$

Proof. We only need to control the difference  $\varepsilon^T(Z_j - Z_j^0)/n$ . We have that

$$|\varepsilon^T (Z_j^0 - Z_j)/n| \le \max_{k \ne j} |\varepsilon^T \mathbf{X}_k/n| \|\hat{\gamma}_j - \gamma_j^0\|_1.$$

The statement then follows from Proposition 7 and invoking (D1) and (D2).

**Proposition 9.** Assume (A2), (A3), (A6), (A7), (D1), (D2) and (D3). Then:

$$\sqrt{n} \frac{Z_j^T X_j / n}{\omega_{p;jj}} (\hat{b}_j - \beta_j^0) \Rightarrow \mathcal{N}(0,1) \ (n \to \infty).$$

Proof. The statement follows by standard arguments as in van de Geer et al. (2014), requiring (D3), and using Proposition 8. For the case with the square root Lasso in (3), the proof is analogous. One can easily show that  $||Z_j||_2/\sqrt{n} = \sqrt{\mathbb{E}|Z_j^{(0)}|^2} + o_P(1)$ , due to (A2), (A3), and (D2), and  $\mathbb{E}|Z_j^{(0)}|^2$ is upper bounded by (A3).

Using the results from Section 7.1.1 and Proposition 9 establish the result from Theorem 1.  $\Box$ 

#### 7.2 Proof of Proposition 1

We write

$$n^{-1} \sum_{i=1}^{n} (\hat{\varepsilon}_i Z_{j,i})^2 = n^{-1} \sum_{i=1}^{n} (\varepsilon_i + (\hat{\varepsilon}_i - \varepsilon_i))^2 (Z_{j,i}^0 + (Z_{j,i} - Z_{j,i}^0))^2.$$

We then get

$$n^{-1} \sum_{i=1}^{n} (\hat{\varepsilon}_i Z_{j;i})^2 = n^{-1} \sum_{i=1}^{n} (\varepsilon_i Z_{j;i}^0)^2 + \Delta.$$

One can easily show that  $\Delta = o_P(1)$  by using Hölder's inequality (for  $\ell_1 - \ell_\infty$ ; and Cauchy-Schwarz for  $\ell_2 - \ell_2$ ) and invoking the following:

$$\begin{aligned} & \max_{i} |Z_{j;i}^{0}| \leq C_{3} < \infty \text{ due to (A3)} \\ & \max_{i} |Z_{j;i} - Z_{j;i}^{0}| \leq \max_{i} |\mathbf{X}_{i,j}| \|\hat{\gamma}_{j} - \gamma_{j}^{0}\|_{1} = o_{P}(1) \text{ due to (A2) and (D2),} \\ & \|\hat{\varepsilon} - \varepsilon\|_{2}^{2}/n = \|\mathbf{X}(\hat{\beta} - \beta^{0})\|_{2}^{2}/n = o_{P}(1) \text{ due to (A2), } \|\beta^{0}\|_{1} = o(\sqrt{n/\log(p)}) \text{ and (A7),} \end{aligned}$$

where the last bound follows from e.g. Bühlmann and van de Geer (2011, Cor.6.1). Therefore,

$$n^{-1} \sum_{i=1}^{n} (\hat{\varepsilon}_i Z_{j,i})^2 = \mathbb{E} |\varepsilon^{(0)} Z_j^{(0)}|^2 + o_P(1).$$

Furthermore, and simpler to obtain:

$$n^{-1} \sum_{i=1}^{n} \hat{\varepsilon}_i Z_{j,i} = \mathbb{E}[\varepsilon^{(0)} Z_j^{(0)}] + o_P(1) = o_P(1).$$

Due to (A6), the latter two displayed formulae complete the proof.

#### 7.3 Proof of Proposition 2

For statement 1, consider:

$$\begin{split} & \sum_{j=1}^{p} |\beta_{j}^{0}|^{r} \\ \leq & \sum_{j=1}^{p} (\sum_{\ell=1}^{p} |(\Sigma^{-1})_{j\ell}| |\Gamma_{\ell}|)^{r} \leq \sum_{j=1}^{p} \sum_{\ell=1}^{p} |(\Sigma^{-1})_{j\ell}|^{r} |\Gamma_{\ell}|^{r} = \sum_{\ell=1}^{p} \|(\Sigma^{-1})_{\bullet\ell}\|_{r}^{r} |\Gamma_{\ell}|^{r} \leq \max_{\ell} \|(\Sigma^{-1})_{\bullet\ell}\|_{r}^{r} \|\Gamma\|_{r}^{r}. \end{split}$$

Furthermore, we have that  $\max_{\ell} \|(\Sigma^{-1})_{\bullet \ell}\|_r^r \leq (\max_{\ell} s_{\ell} + 1) \|\Sigma^{-1}\|_{\infty}^r$  and therefore statement 1 is complete.

Regarding statement 2, we use the following argument. Every point  $\ell \in S_{\Gamma}$  can lead to at most  $s_{\ell} + 1$  non-zero values of the components of  $\beta^0$ , due to formula (9). Hence we obtain both bounds for  $\|\beta^0\|_0^0$ .

## 7.4 Proof of Corollary 1

The bound above for  $\|\Gamma\|_0^0$  follows by a similar argument as for statement 2. in Proposition 2: every support point in  $S_{f^0}$  exhibits a dependence with at most  $b_{\text{max}}$  X-variables: therefore there are at most  $b_{\text{max}}|S_{f^0}|$  non-zero covariances between  $f^0(X)$  and the X-variables.

#### 7.5 Proof of Proposition 3

It is well known that

$$\beta_j^0 = \mathbb{E}[Z_j^0 f^0(X)] = \mathbb{E}[Z_j^0 f(X_{S_{f^0}})].$$

Furthermore, since  $Z_j^{(0)}$  is the residual when projecting  $X_j^{(0)}$  onto  $X_{-j}^{(0)} = \{X_k^{(0)}; k \neq j\}$  and due to the Gaussian assumption:  $Z_j^{(0)}$  is independent of  $\{X_k^{(0)}; k \neq j\}$ .

Therefore, if  $j \notin S_{f^0}$ ,  $Z_j^{(0)}$  is independent also of  $X_{S_{f^0}}^{(0)}$  and therefore, using the representation for  $\beta_j^0$  above:  $\beta_j^0 = \mathbb{E}[Z_j^{(0)}]\mathbb{E}[f^0(X_{S_{f^0}})] = 0$ , saying that  $j \notin S^0$ . This proves the claim.

#### 7.6 Proof of Proposition 4

As mentioned already in the proof of Proposition 3 we know that  $Z_j^{(0)}$  is independent of  $\{X_k^{(0)}; k \neq j\}$ . Therefore, for  $j \in S_k$ :

$$\beta_j^0 = \mathbb{E}[Z_j^{(0)} f^0(X^{(0)})] = \mathbb{E}[Z_j^{(0)} \left( f_1^0(X_{S_1}^{(0)}) + \ldots + f_d^0(X_{S_d}^{(0)}) \right)] = \mathbb{E}[Z_j^{(0)} f_k^0(X_{S_k}^{(0)})].$$

This means that we can obtain  $\beta_j^0$  from projecting  $f^0(X_{S_k}^{(0)})$  onto  $\{X_i^{(0)}; j=1,\ldots,p\}$ :

$$\gamma = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \mathbb{E} |f^0(X_{S_k}^{(0)}) - (X^{(0)})^T \beta|^2, \tag{16}$$

and  $\beta_j^0 = \gamma_j$ . But we know from Proposition 3 that for the support of  $\gamma$ :

$$S(\gamma) = \{j; \ \gamma_j \neq 0\} \subseteq S_k.$$

Therefore, we can restrict the projection in (16) to the variables from  $S_k$ :

$$\tilde{\gamma} = \operatorname{argmin}_{\beta \in \mathbb{R}^{|S_k|}} \mathbb{E} |f^0(X_{S_k}^{(0)}) - (X_{S_k}^{(0)})^T \beta|^2,$$

and  $\beta_j^0 = \tilde{\gamma}_{c(j)}$ , where c(j) the index of the component in  $\tilde{\gamma}$  which corresponds to variable  $X_j^{(0)}$ . This completes the proof.

## 7.7 Proof of Proposition 5

We write

$$||Z_{j}||_{2}^{2}/n = ||Z_{j}^{0}||_{2}^{2}/n + ||\mathbf{X}_{-j}(\hat{\gamma}_{j} - \gamma_{j}^{0})||_{2}^{2}/n + \Xi,$$
  

$$|\Xi| \leq 2||Z_{j}^{0}||_{2}/\sqrt{n}||\mathbf{X}_{-j}(\hat{\gamma}_{j} - \gamma_{j}^{0})||_{2}/\sqrt{n}.$$
(17)

Due to (C1,i) we have that

$$||Z_j^0||_2^2/n \ge C_7/2$$
 with probability tending to one. (18)

We can also establish, analogous to Bühlmann and van de Geer (2011, Cor.6.1) invoking (C1,iii), but now controlling  $\max_{k\neq j} |(Z_j^0)^T \mathbf{X}_k|/n = O_P(\sqrt{\log(p)/n})$  (see Lemma 1 and using (C1,i) and (C1,ii)):

$$\|\mathbf{X}_{-i}(\hat{\gamma}_i - \gamma_i^0)\|_2^2 / n = o_P(1). \tag{19}$$

By (17), (18) and (19) we complete the proof.

## References

- Belloni, A., Chen, D., Chernozhukov, V., and Hansen, C. (2012). Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica*, 80:2369–2429.
- Belloni, A., Chernozhukov, V., and Wang, L. (2011). Square-root Lasso: pivotal recovery of sparse signals via conic programming. *Biometrika*, 98:791–806.
- Benjamini, Y. and Yekutieli, D. (2005). False discovery rate-adjusted multiple confidence intervals for selected parameters. *Journal of the American Statistical Association*, 100:71–81.
- Bühlmann, P. (2013). Statistical significance in high-dimensional linear models. *Bernoulli*, 19:1212–1242.
- Bühlmann, P. and van de Geer, S. (2011). Statistics for High-Dimensional Data: Methods, Theory and Applications. Springer.
- Candès, E. and Tao, T. (2007). The Dantzig selector: statistical estimation when p is much larger than n (with discussion). *Annals of Statistics*, 35:2313–2404.
- Candès, E. J. and Tao, T. (2006). Near-optimal signal recovery from random projections: Universal encoding strategies? *IEEE Transactions on Information Theory*, 52:5406–5425.
- Chen, S. S., Donoho, D. L., and Saunders, M. A. (1998). Atomic decomposition by basis pursuit. SIAM Journal on Scientific Computing, 20:33–61.
- Dezeure, R., Bühlmann, P., Meier, L., and Meinshausen, N. (2014). High-dimensional inference: confidence intervals, p-values and R-software hdi. Preprint arXiv:1408.4026.
- Donoho, D. L. (2006). Compressed sensing. *IEEE Transactions on Information Theory*, 52:1289–1306.
- Eicker, F. (1967). Limit theorems for regressions with unequal and dependent errors. In *Proceedings* of the fifth Berkeley symposium on mathematical statistics and probability, volume 1, pages 59–82.
- Foygel Barber, R. and Candès, E. (2014). Controlling the false discovery rate via knockoffs. arXiv:1404.5609.
- Freedman, D. A. et al. (1981). Bootstrapping regression models. Annals of Statistics, 9:1218–1228.

- Friedman, J. H. (1991). Multivariate adaptive regression splines (with discussion). *Annals of Statistics*, 19:1–67.
- Huber, P. J. (1967). The behavior of maximum likelihood estimates under nonstandard conditions. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 221–233.
- Jankova, J. and van de Geer, S. (2014). Confidence intervals for high-dimensional inverse covariance estimation. Preprint arXiv:1403.6752.
- Javanmard, A. and Montanari, A. (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *Journal of Machine Learning Research*, 15:2869–2909.
- Lockhart, R., Taylor, J., Tibshirani, R. J., and Tibshirani, R. (2014). A significance test for the lasso (with discussion). *Annals of Statistics*, 42:413–468.
- Meier, L., Meinshausen, N., and Dezeure, R. (2014). hdi: High-Dimensional Inference. R package version 0.1-2.
- Meinshausen, N. (2015). Group-bound: confidence intervals for groups of variables in sparse high-dimensional regression without assumptions on the design. *Journal of the Royal Statistical Society, Series B (to appear); Preprint arXiv:1309.3489.*
- Meinshausen, N. and Bühlmann, P. (2010). Stability Selection (with discussion). *Journal of the Royal Statistical Society, Series B*, 72:417–473.
- Meinshausen, N., Meier, L., and Bühlmann, P. (2009). P-values for high-dimensional regression. Journal of the American Statistical Association, 104:1671–1681.
- Minnier, J., Tian, L., and Cai, T. (2011). A perturbation method for inference on regularized regression estimates. *Journal of the American Statistical Association*, 106:1371–1382.
- Ren, Z., Sun, T., Zhang, C.-H., and Zhou, H. (2015). Asymptotic normality and optimalities in estimation of large Gaussian graphical model. To appear in the Annals of Statistics; Preprint arXiv:1309.6024.
- Taylor, J., Lockhart, R., Tibshirani, R. J., and Tibshirani, R. (2014). Exact post-selection inference for forward stepwise and least angle regression. Preprint arXiv:1401.3889.
- Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288.
- van de Geer, S. (2014). Statistical theory for high-dimensional models. Preprint arXiv:1409.8557.
- van de Geer, S. (2015).  $\chi^2$ -confidence sets in high-dimensional regression. Preprint arXiv:1502.07131.
- van de Geer, S., Bühlmann, P., Ritov, Y., and Dezeure, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *Annals of Statistics*, 42:1166–1202.
- Wasserman, L. (2014). Discussion: a significance test for the Lasso. Annals of Statistics, 42:501–508.

- White, H. (1980). A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity. *Econometrica: Journal of the Econometric Society*, 48:817–838.
- Zhang, C.-H. and Zhang, S. S. (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society, Series B*, 76:217–242.