

Degrees of freedom for piecewise Lipschitz estimators

Frederik Riis Mikkelsen and Niels Richard Hansen

Department of Mathematical Sciences, University of Copenhagen, Universitetsparken 5, 2100 Copenhagen Ø, Denmark.

E-mail: frm@math.ku.dk; Niels.R.Hansen@math.ku.dk

Received 10 June 2016; revised 10 February 2017; accepted 22 February 2017

Abstract. A representation of the degrees of freedom akin to Stein's lemma is given for a class of estimators of a mean value parameter in \mathbb{R}^n . Contrary to previous results our representation holds for a range of discontinuous estimators. It shows that even though the discontinuities form a Lebesgue null set, they cannot be ignored when computing degrees of freedom. Estimators with discontinuities arise naturally in regression if data driven variable selection is used. Two such examples, namely best subset selection and lasso-OLS, are considered in detail in this paper. For lasso-OLS the general representation leads to an estimate of the degrees of freedom based on the lasso solution path, which in turn can be used for estimating the risk of lasso-OLS. A similar estimate is proposed for best subset selection. The usefulness of the risk estimates for selecting the number of variables is demonstrated via simulations with a particular focus on lasso-OLS.

Résumé. Une représentation des degrés de liberté comparable au lemme de Stein est donnée pour une classe d'estimateurs du paramètre de la moyenne dans \mathbb{R}^n . Contrairement aux résultats précédents, notre représentation est valable pour une famille d'estimateurs discontinus. Cela montre que même si les discontinuités sont de mesure de Lebesgue zéro, elles ne peuvent pas être ignorées lors du calcul des degrés de liberté. Les estimateurs avec discontinuités apparaissent naturellement dans les modèles de régression si de la sélection de variables par les données est utilisée. Deux tels exemples, la sélection du meilleur sous-ensemble et le lasso-OLS, sont considérés en détail dans l'article. Pour le Lasso-OLS, la représentation générale mène à une estimation des degrés de liberté basée sur le chemin de la solution Lasso, qui à son tour peut être utilisée pour estimer le risque du lasso-OLS. Une estimation similaire est proposée pour la sélection du meilleur sous-ensemble. L'utilité des estimées de risque pour le choix du nombre de variables est démontrée par des simulations qui se concentrent en particulier sur lasso-OLS.

MSC: 62J05; 62J07

Keywords: Best subset selection; Lasso-OLS; Degrees of freedom; Stein's lemma

1. Introduction

Representations of the effective dimension of a statistical model have been studied extensively in many different frameworks. For classical model selection criteria such as AIC and Mallows's C_p the dimension of the parameter space is used to adjust the empirical risk for its optimism so as to provide a fair model score across different dimensions. A number of extensions to models or methods without a well defined dimension exist, such as the trace of the smoother matrix for scatter plot smoothers, see e.g. [13], and the use of the divergence of a sufficiently differentiable estimator based on Stein's lemma as described in [5]. Stein's lemma was used by Zou et al. [28] and Tibshirani and Taylor [25] to demonstrate that for the lasso estimator in a linear regression model with Gaussian errors, the number of estimated non-zero parameters is an appropriate estimate of the effective dimension.

It is well known that neither Mallows's C_p nor AIC or related information criteria correctly adjust for the optimism that results from selecting one model among a number of models of equal dimension. The usage of such methods for model selection without adequate adjustments was called "a quiet scandal in the statistical community" by Breiman

[1], who proposed a bootstrap based method for risk estimation as an alternative. Ye [27] defined the notion of generalized degrees of freedom for an estimator of the mean in a Gaussian model and showed how to use this number for risk estimation. The results by Ye apply to discontinuous estimators that involve model selection, but his proposal for computing the degrees of freedom was similarly to Breiman's based on refitting models to perturbed data.

If the estimator satisfies the differentiability requirements for Stein's lemma, Lemma 2 in [21], the divergence of the estimator w.r.t. the data is an unbiased estimate of the degrees of freedom in the generalized sense of [27]. This was used by Donoho and Johnstone [4], Meyer and Woodroffe [18], Zou et al. [28], Kato [14] and Tibshirani and Taylor [25] among others to derive formulas for the degrees of freedom of estimators that are Lipschitz continuous.

For estimators with discontinuities Stein's lemma generally breaks down and the divergence will not be an unbiased estimate of the degrees of freedom. Note that an estimator can be continuous or even differentiable almost everywhere – it can be a projection locally – and still be defined globally in such a way that it has non-ignorable discontinuities. This is, in particular, the case in regression when data adaptive variable selection is used to select among a number of projection estimators. Best subset selection is one central example, but variable selection procedures lead in general to non-ignorable discontinuities. A variable selection procedure effectively divides the sample space into a finite number of disjoint regions, with the estimator being a projection, say, on each region. The resulting estimator consisting of a selection step and a projection step will generally be discontinuous on the boundary between two regions.

Tibshirani [24] recently made headway with the computation of the degrees of freedom for some discontinuous estimators. Specifically, he considered a linear regression model with an orthogonal design and showed how to compute the degrees of freedom for hard thresholding, which for orthogonal designs is equivalent to the Lagrangian formulation of best subset selection. He also gave an extension of Stein's lemma to some discontinuous estimators, though it was not shown if this extension applies to subset selection estimators. Hansen and Sokol [12] gave a different generalization of Stein's lemma for all estimators that are metric projections onto a closed set. This generalization applies to subset selection and other estimators with non-convex constraints, but did not lead to a readily computable representation of the contribution to the degrees of freedom that are due to the discontinuities of the metric projection.

The first main contribution of this paper is the general Theorem 2.4, which is a version of Stein's lemma for estimators that are locally Lipschitz continuous on each of a finite number of open sets, whose union makes up Lebesgue almost all of \mathbb{R}^n . This is a broad class of estimators containing a number of regression estimators that include variable selection. Compared to existing results, Theorem 2.4 holds under verifiable conditions without putting restrictions on the design matrix such as orthogonality.

As a main example the lasso-OLS estimator in a linear regression setup is investigated in detail in Section 3. The lasso-OLS estimator consists of two steps: variable selection using lasso followed by ordinary least squares estimation using the selected variables. This estimator was referred to as the LARS-OLS hybrid in [6], and it is a limit case of the relaxed lasso as considered in [17]. We follow the terminology of [2], p. 34, and call it the lasso-OLS estimator.

The second main contribution of this paper is a derivation of a computable estimate of the degrees of freedom – and thus the risk – for lasso-OLS, which only involves the computation of a single lasso solution path and corresponding OLS estimators along the path. Simulation studies reported in Section 4 demonstrated that the resulting risk estimate leads to reliable model selection across a range of different designs and parameter settings, and that the risk estimate itself has smaller mean squared error than the computationally more demanding cross-validation estimate.

For the Lagrangian formulation of best subset selection it is also demonstrated that Theorem 2.4 holds, but the situation is more complicated than for lasso-OLS. However, it is possible to derive an approximation, which is exact for orthogonal designs, as shown in Section 5.

The proof of Theorem 2.4 and some auxiliary technical results are in the [appendix](#).

2. A general representation of degrees of freedom

Throughout the paper we consider the multivariate Gaussian model $\mathcal{N}(\mu, \sigma^2 I)$ on \mathbb{R}^n with μ the unknown parameter, and we let $\hat{\mu} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ denote an estimator of μ . A typical application is to linear regression estimators of the form $X\hat{\beta}$ where X denotes an $n \times p$ matrix and $\hat{\beta}$ denotes an estimator of the parameters in the linear regression model. When the estimator $\hat{\beta}$ sets some of the parameters to exactly zero we say that the estimator does variable selection. The lasso, [22], is an example of a globally Lipschitz continuous estimator that does variable selection, while best subset selection is a discontinuous estimator that does variable selection. The lasso-OLS – as studied

intensively in Section 3 – is another example of a discontinuous regression estimator that does variable selection. Though discontinuous regression estimators that do variable selection constitute the main motivation for the present paper, the general results are more conveniently formulated in terms of estimators of the mean μ without reference to the regression setup.

Letting $Y \sim \mathcal{N}(\mu, \sigma^2 I)$ the risk of the estimator is defined as

$$\text{Risk}(\hat{\mu}) := E \|\mu - \hat{\mu}(Y)\|_2^2,$$

provided that $\hat{\mu}(Y)$ has finite second moment, which will thus be assumed throughout. The risk is a quantification of the error of $\hat{\mu}$, and tuning parameters are often chosen by minimising an estimate of the risk. Our main interest is to estimate the risk under the Gaussian model. The following definition introduces two notions of degrees of freedom that are useful when we want to estimate the risk. In the definition, $\psi(y; \mu, \sigma^2)$ denotes the density for the $\mathcal{N}(\mu, \sigma^2 I)$ distribution and $\langle \cdot, \cdot \rangle$ denotes the standard inner product on \mathbb{R}^n . The divergence operator is also needed. It is the differential operator defined as

$$\text{div}(f) = \sum_{i=1}^n \partial_i f_i$$

for $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ Lebesgue almost everywhere differentiable and with ∂_i denoting the partial derivative w.r.t. the i th coordinate.

Definition 2.1. For a measurable map $\hat{\mu} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ such that $\hat{\mu}(Y)$ has finite second moment the degrees of freedom of $\hat{\mu}$ is

$$\text{df}(\hat{\mu}) := \sum_{i=1}^n \frac{\text{cov}(Y_i, \hat{\mu}(Y)_i)}{\sigma^2} = \int \frac{\langle y - \mu, \hat{\mu}(y) \rangle}{\sigma^2} \psi(y; \mu, \sigma^2) dy. \tag{1}$$

If $\hat{\mu}$ is differentiable in Lebesgue almost all points and $\text{div}(\hat{\mu})$ has finite first moment Stein’s degrees of freedom of $\hat{\mu}$ is

$$\text{df}_S(\hat{\mu}) := E(\text{div}(\hat{\mu})(Y)). \tag{2}$$

A simple expansion of the risk yields

$$\text{Risk} = E \|Y - \hat{\mu}(Y)\|_2^2 - n\sigma^2 + 2\sigma^2 \text{df}(\hat{\mu}). \tag{3}$$

Hence $\|Y - \hat{\mu}(Y)\|_2^2 - n\sigma^2 + 2\sigma^2 \widehat{\text{df}}$ is an unbiased risk estimate if $\widehat{\text{df}}$ is an unbiased estimate of $\text{df}(\hat{\mu})$. In practice, σ^2 must be estimated as well and a bias of $\widehat{\text{df}}$ can also be preferable if it reduces the variance. Hence exact unbiasedness of a risk estimate based on (3) is of secondary interest, but it is of interest to find adequate corrections of the squared error $\|Y - \hat{\mu}(Y)\|_2^2$ that can be used for model assessment and comparison.

If $\hat{\mu}$ is *almost differentiable* then $\text{df}(\hat{\mu}) = \text{df}_S(\hat{\mu})$ due to Stein’s lemma (Lemma 2 in [21]), in which case $\text{div}(\hat{\mu})(Y)$ is an unbiased estimate of $\text{df}(\hat{\mu})$. However, most estimators with discontinuities are not almost differentiable, and for such estimators it is not clear if $\text{div}(\hat{\mu})(Y)$ is a useful estimate of the degrees of freedom. Indeed, our main result, Theorem 2.4, provides a representation of $\text{df}(\hat{\mu}) - \text{df}_S(\hat{\mu})$, which is nonzero for a range of estimators. The result provides the theoretical basis for establishing more adequate estimates of the degrees of freedom and thus the risk. Furthermore, Theorem 3.2 provides a quite remarkable connection between $\text{df}(\hat{\mu})$ and $\text{df}_S(\hat{\mu})$ for the lasso-OLS estimator, which can be used to derive an estimate of $\text{df}(\hat{\mu})$. This result is directly applicable in practice and provides fast and accurate risk estimation without the need for cross-validation, say.

Our main result is derived under the assumptions on the estimator as stated below. To fix notation we let $B(x, r)$ denote the closed ball in \mathbb{R}^n of radius r and center x . Additionally, we let \mathcal{H}^{n-1} denote the $n - 1$ dimensional Hausdorff measure – a generalisation of the surface measure of $n - 1$ dimensional hypersurfaces in \mathbb{R}^n (see e.g. [7] for details).

Assumption 2.2. The estimator $\hat{\mu}$ can be written as $\hat{\mu} = \sum_{i=1}^N 1_{U_i} \hat{\mu}_i$ for a collection of open and disjoint sets $\{U_i\}_{i=1}^N$ with $\bigcup_{i=1}^N \overline{U_i} = \mathbb{R}^n$. Additionally, for each $i = 1, \dots, N$:

- (a) The map $\hat{\mu}_i : \overline{U_i} \rightarrow \mathbb{R}^n$ is locally Lipschitz.
- (b) The random variable $1_{U_i} \operatorname{div}(\hat{\mu}_i)(Y)$ has finite first moment and $\|\hat{\mu}_i\|$ is polynomially bounded on U_i .
- (c) The function $r \mapsto \mathcal{H}^{n-1}(\partial U_i \cap B(0, r))$ is polynomially bounded.

Remark 2.3. The following points are worth noting:

- (a) *Boundary values of the estimator.* Assumption 2.2(c) implies that the boundaries of the sets U_i are Lebesgue null sets, and thus that $\mathbb{R}^n \setminus \bigcup_i U_i$ has Lebesgue measure zero. The estimator $\hat{\mu}$ is here defined to be zero on this null set, but with Y having an absolutely continuous distribution its value on a null set is irrelevant. Note, however, that Assumption 2.2(a) ensures that $\hat{\mu}_i$ is uniquely defined on ∂U_i . In a concrete case there may be a natural way to define $\hat{\mu}$ on the common boundary between U_i and U_j , say, but we make no abstract attempt to select between μ_i and μ_j on the boundary.
- (b) *Degrees of freedom.* Assumption 2.2(a) implies by Rademacher's theorem (Theorems 3.1.6 and 3.1.7 in [9]) that $\operatorname{div}(\hat{\mu}_i)$ is defined Lebesgue a.e. Combining this with Assumption 2.2(b) we conclude that under Assumption 2.2 both $\operatorname{df}(\hat{\mu})$ and $\operatorname{df}_S(\hat{\mu})$ are well defined.
- (c) *Existence of normal vectors.* Assumption 2.2(c) implies that the sets U_i have locally finite perimeter (see Theorem 5.11.1 in [7]), thus a measure theoretic outer unit normal η_i is defined on a subset of ∂U_i . In fact, by Lemma A.2 Assumption 2.2(c) only needs to hold for the reduced boundary $\partial^* U_i$ (see Definition 5.7 and Lemma 5.8.1 in [7]). Whenever ∂U_i is smooth the measure theoretic unit normal coincides with the usual point-wise unit normal.

Estimators that involve data driven variable selection will generally fulfil Assumption 2.2 with each U_i corresponding to a set of selected variables. Example 2.5 provides a thorough characterization of U_i in the lasso-OLS setup. Moreover, a similar characterization of U_i is given in Example 3.4 for a class of estimators defined via minimization of a penalized loss function.

The conditions in Assumption 2.2 are typically easy to verify, except perhaps the third condition, as it involves bounding Hausdorff measures. Section A.1 provides some results that can be helpful for verifying the third condition. For estimators satisfying Assumption 2.2 we have the following representation of the degrees of freedom.

Theorem 2.4. *If $\hat{\mu}$ satisfies Assumption 2.2 then*

$$\operatorname{df}(\hat{\mu}) = \operatorname{df}_S(\hat{\mu}) + \frac{1}{2} \sum_{i \neq j} \int_{\overline{U_i} \cap \overline{U_j}} \langle \hat{\mu}_j - \hat{\mu}_i, \eta_i \rangle \psi(\cdot; \mu, \sigma^2) d\mathcal{H}^{n-1}, \quad (4)$$

where η_i denotes the measure theoretic outer unit normal to ∂U_i .

The proof is in Section A.2. The essential part is an application of a generalized version of Gauss–Green's formula combined with a dominated convergence argument. Note that though $\overline{U_i} \cap \overline{U_j}$ is a Lebesgue null set – on which $\hat{\mu}$ is defined to be zero – $\hat{\mu}_j$ and $\hat{\mu}_i$ are uniquely defined by Assumption 2.2(a) and generally non-zero and different, cf. also Remark 2.3(a).

If $\hat{\mu}$ satisfies Assumption 2.2 and is continuous then (4) reduces to $\operatorname{df}(\hat{\mu}) = \operatorname{df}_S(\hat{\mu})$, which is Stein's lemma for a class of locally Lipschitz continuous estimators. The boundary integrals therefore account for potential jumps of $\hat{\mu}$ across the boundary of any two adjacent regions U_i and U_j . For two-step procedures consisting of a model selection step followed by a parameter estimation step, df_S generally only accounts for the contribution to the degrees of freedom by the estimation step, and the boundary integrals account for the contribution from the selection step.

The following example illustrates how to verify Assumption 2.2 for the lasso-OLS estimator, which is the estimator that will also be the main focus of the subsequent section.

Example 2.5 (The lasso-OLS estimator). Let X be an $n \times p$ -matrix. For any subset $A \subseteq \{1, \dots, p\}$, X_A denotes the matrix whose columns are those of X indexed by A , and similarly, $\beta_A \in \mathbb{R}^{|A|}$ denotes $(\beta_i)_{i \in A}$ for $\beta \in \mathbb{R}^p$. We let

$$\mathcal{S} := \{S = \operatorname{col}(X_A) \mid A \subseteq \{1, \dots, p\}\}$$

denote the set of subspaces spanned by columns of X . The orthogonal projection onto a subspace $S \in \mathcal{S}$ is denoted by Π_S .

A *lasso estimator* $\hat{\mu}_{\text{lasso}}^\lambda(y)$ with tuning parameter $\lambda > 0$ is defined as $\hat{\mu}_{\text{lasso}}^\lambda(y) = X\hat{\beta}^\lambda$ where

$$\hat{\beta}^\lambda \in \arg \min_{\beta} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1.$$

We do not make any assumptions on X , and therefore it may happen that multiple $\hat{\beta}^\lambda$ -solutions exist. For a solution $\hat{\beta}^\lambda$, the support, $\text{supp}(\hat{\beta}^\lambda) \subseteq \{1, \dots, p\}$, is called an *active set*. The lasso estimator $\hat{\mu}_{\text{lasso}}^\lambda(y) = X\hat{\beta}^\lambda$ belongs to the space $\text{col}(X_A)$ for $A = \text{supp}(\hat{\beta}^\lambda)$, and it follows by Lemma 7 in [25] that there exists a Lebesgue null set N , such that $\text{col}(X_A)$ is invariant with respect to the choice of the active set of solutions for $y \notin N$. The map $\widehat{S}^\lambda : \mathbb{R}^n \setminus N \rightarrow \mathcal{S}$ returning $\text{col}(X_A)$ when there is a solution $\hat{\beta}^\lambda$ with active set $A = \text{supp}(\hat{\beta}^\lambda)$ is therefore well defined. The *lasso-OLS estimator* $\hat{\mu}_{\text{1-OLS}}^\lambda := \Pi_{\widehat{S}^\lambda}$ is defined as the projection onto the space selected by the lasso, and is thus well-defined Lebesgue almost everywhere.

By defining the disjoint selection events

$$U_S^\lambda := (\widehat{S}^\lambda = S)$$

for each $S \in \mathcal{S}$, we immediately see from Lemma 6 in [25] that each selection event is open and that $\mathbb{R}^n = \bigcup_{S \in \mathcal{S}} \overline{U_S^\lambda}$. We can safely ignore any empty U_S^λ . From the proof of Lemma 6 in [25] we see that $\partial U_S^\lambda \subseteq (\bigcup_{T \in \mathcal{S}} U_T^\lambda)^c$ is a finite union of affine subspaces of dimensions $\leq n - 1$, and $r \mapsto \mathcal{H}^{n-1}(\partial U_S^\lambda \cap B(0, r))$ is thus polynomially bounded. This follows by elementary considerations, but it is also a consequence of Lemma A.1. Consequently,

$$\hat{\mu}_{\text{1-OLS}}^\lambda = \sum_{S \in \mathcal{S}} 1_{U_S^\lambda} \Pi_S \quad \text{almost everywhere,}$$

and it satisfies all conditions in Assumption 2.2. Figure 1 provides an illustration of the partition of \mathbb{R}^n for $n = p = 2$ for different choices of angles between the columns in X .

Note that since $\hat{\mu}_{\text{1-OLS}}^\lambda = \Pi_S$ on the open set U_S^λ , its divergence equals $\dim(S)$, hence Stein’s degrees of freedom is

$$\text{df}_S(\hat{\mu}_{\text{1-OLS}}^\lambda) = E(\dim(\widehat{S}^\lambda)).$$

From Lemma 3 in [23] it follows that $\dim(\widehat{S}^\lambda) = |\text{supp}(\hat{\beta}^\lambda)|$ whenever the columns of X are in general position, which is useful for practical computations.

The arguments above are based on results in [25], but see also [15] for related characterizations of the selection events for lasso.

3. Risk estimation for lasso-OLS

It is not obvious how the general formula in Theorem 2.4 for $\text{df}(\hat{\mu})$ can be used for computing or estimating the degrees of freedom. The first term of (4), $\text{df}_S(\hat{\mu})$, may be estimated by $\text{div}(\hat{\mu})(Y)$, but the second term is more difficult. In this section we show how this second term can be related to the derivative of $\lambda \mapsto \text{df}_S(\hat{\mu}_{\text{1-OLS}}^\lambda)$ for lasso-OLS. First we recapitulate the computations in [24] of the degrees of freedom for lasso-OLS with X orthogonal, which will reveal the general formula shown below.

Example 3.1 (Continuation of Example 2.5). Assume that $n = p$ and $X = I$. In this case it is well known that the lasso and the lasso-OLS estimators become the soft and hard thresholding estimators, respectively. That is,

$$\hat{\mu}_{\text{lasso},i}^\lambda = \begin{cases} Y_i - \lambda \text{sign}(Y_i) & \text{if } |Y_i| > \lambda, \\ 0 & \text{otherwise,} \end{cases} \quad \text{and} \quad \hat{\mu}_{\text{1-OLS},i}^\lambda = \begin{cases} Y_i & \text{if } |Y_i| > \lambda, \\ 0 & \text{otherwise.} \end{cases}$$

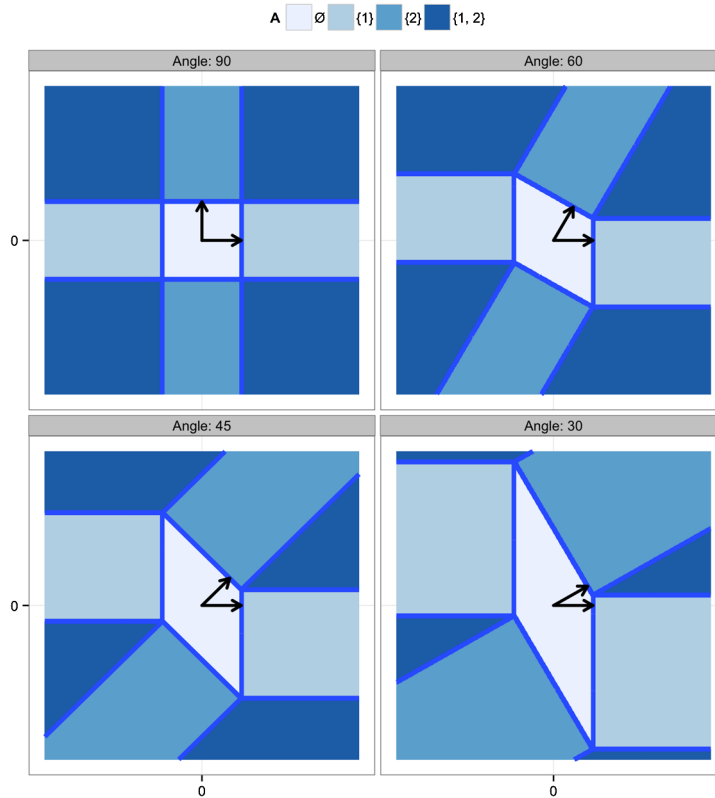


Fig. 1. Illustrations of the decomposition of \mathbb{R}^2 into the four sets U_{\emptyset}^1 , $U_{\{1\}}^1$, $U_{\{2\}}^1$ and $U_{\{1,2\}}^1$ according to the lasso estimator with $\lambda = 1$. The set U_{\emptyset}^1 consists of the points shrunk to zero, the sets $U_{\{1\}}^1$ and $U_{\{2\}}^1$ to the points where either the second or the first coordinate, respectively, is shrunk to zero and $U_{\{1,2\}}^1$ to the set where none of the coordinates are shrunk to zero. The decomposition depends on the angle between the two columns in X .

We can write up closed form expressions for $\text{df}(\hat{\mu}_{\text{I-OLS}}^\lambda)$ and $\text{df}_S(\hat{\mu}_{\text{I-OLS}}^\lambda)$:

$$\begin{aligned} \text{df}_S(\hat{\mu}_{\text{I-OLS}}^\lambda) &= \int \psi(y; \mu, \sigma^2) \sum_i 1_{(|y_i| > \lambda)} dy = \sum_i \int_{(|y_i| > \lambda)} \psi(y_i; \mu_i, \sigma^2) dy_i \\ &= \sum_i \Phi\left(\frac{-\lambda - \mu_i}{\sigma}\right) + \left(1 - \Phi\left(\frac{\lambda - \mu_i}{\sigma}\right)\right), \end{aligned}$$

and as in [24]

$$\begin{aligned} \text{df}(\hat{\mu}_{\text{I-OLS}}^\lambda) &= \sum_i \int_{\lambda}^{\infty} \psi(y_i; \mu_i, \sigma^2) \frac{y_i(y_i - \mu_i)}{\sigma^2} dy_i + \int_{-\infty}^{-\lambda} \psi(y_i; \mu_i, \sigma^2) \frac{y_i(y_i - \mu_i)}{\sigma^2} dy_i \\ &= \sum_i [-\psi(y_i; \mu_i, \sigma^2) y_i]_{\lambda}^{\infty} + \int_{\lambda}^{\infty} \psi(y_i; \mu_i, \sigma^2) dy_i \\ &\quad + [-\psi(y_i; \mu_i, \sigma^2) y_i]_{-\infty}^{-\lambda} + \int_{-\infty}^{-\lambda} \psi(y_i; \mu_i, \sigma^2) dy_i \\ &= \lambda \sum_i (\psi(\lambda; \mu_i, \sigma^2) + \psi(-\lambda; \mu_i, \sigma^2)) + \text{df}_S(\hat{\mu}_{\text{I-OLS}}^\lambda). \end{aligned}$$

Letting ∂_λ denote the differential operator with respect to λ we observe that

$$\text{df}(\hat{\mu}_{1\text{-OLS}}^\lambda) = \text{dfs}(\hat{\mu}_{1\text{-OLS}}^\lambda) - \lambda \partial_\lambda \text{dfs}(\hat{\mu}_{1\text{-OLS}}^\lambda), \tag{5}$$

which is a striking identity. This is because the formula for $\text{df}(\hat{\mu}_{1\text{-OLS}}^\lambda)$, though explicit, involves the unknown parameter μ and is not readily estimable. But we have the divergence estimator, $\sum_i 1_{(|y_i| > \lambda)}$, of $\text{dfs}(\hat{\mu}_{1\text{-OLS}}^\lambda)$, and if we from this can estimate its derivative as well, the formula above suggests how to estimate $\text{df}(\hat{\mu}_{1\text{-OLS}}^\lambda)$.

The remarkable fact that we will show is that (5) holds without the orthogonality assumption on X .

Theorem 3.2. *For the lasso-OLS estimator defined in Example 2.5 it holds that*

$$\text{df}(\hat{\mu}_{1\text{-OLS}}^\lambda) = \text{dfs}(\hat{\mu}_{1\text{-OLS}}^\lambda) - \lambda \partial_\lambda \text{dfs}(\hat{\mu}_{1\text{-OLS}}^\lambda), \tag{6}$$

where ∂_λ denotes differentiation w.r.t. λ .

Theorem 3.2 suggests that $\text{df}(\hat{\mu}_{1\text{-OLS}}^\lambda)$ can be estimated by differentiation of an estimate of $\text{dfs}(\hat{\mu}_{1\text{-OLS}}^\lambda)$. The divergence estimate of Stein’s degrees of freedom is, however, not differentiable as a function of λ , and we need to somehow smooth it. To this end it is convenient to reparametrize the penalization in terms of $\delta = \log(\lambda)$, so that with

$$h(\delta) := \text{dfs}(\hat{\mu}_{1\text{-OLS}}^{\text{exp}(\delta)}),$$

then

$$\text{df}(\hat{\mu}_{1\text{-OLS}}^{\text{exp}(\delta)}) = h(\delta) - h'(\delta).$$

In simulations h was found to be monotonically decreasing, and thus h' to be negative, but we cannot prove that this is generally the case. The integral representation of h' from Theorem 2.4 is not particularly helpful as the integrand can, in fact, be negative. Based on our computational observations – and to reduce variance of the resulting estimate – our proposal is based on the assumption that h' is negative. It is effectively a kernel smoother that estimates the intensity of jumps for a monotone jump process.

We note that $\text{dim}(\hat{S}^{\text{exp}(\delta)})$ is an unbiased estimate of $h(\delta)$ and that the function $\delta \mapsto \text{dim}(\hat{S}^{\text{exp}(\delta)})$ is a step function. The problem of estimating the derivative, h' , of its mean is thus analogous to estimating the intensity for a jump process with one main difference; the step function can have jumps of negative as well as positive sign, though most jumps will be negative. Our proposed estimate ignores the positive excursions of the step function and is computed as follows:

- Compute the jump points, λ_i and jump sizes, $\Delta_i := \inf_{\lambda < \lambda_i} \text{dim}(\hat{S}^\lambda) - \text{dim}(\hat{S}^{\lambda+})$, of the decreasing function $\lambda \mapsto \inf_{\lambda' < \lambda} \text{dim}(\hat{S}^{\lambda'})$ for $i = 1, \dots, M$.
- Apply a kernel density smoother to the points $\delta_i = \log(\lambda_i)$ for $i = 1, \dots, M$ counted with the multiplicities Δ_i . In the simulations presented in this paper an adaptive Gaussian kernel density smoother was used (see Section 10.4.3.2 in [11]).
- Rescale the density estimate by the total number of jumps, that is, by $\sum_{i=1}^M \Delta_i$.

As mentioned above, we can think of the proposed estimate of h' as a non-parametric estimate of the intensity of the jumps for a monotonically decreasing jump process. Alternatively, we can think of it as smoothing the jumps by a sigmoidal function (the anti-derivative of the kernel) to obtain a smooth estimate of Stein’s degrees of freedom, which can then be differentiated. Note that even if Δ_i may always be 1 in theory, the jumps are in practice computed on a grid and may thus be larger than 1, which the procedure accounts for. The estimate of $-\lambda \partial_\lambda \text{dfs}(\hat{\mu}_{1\text{-OLS}}^\lambda)$ resulting from the procedure above is denoted by $\hat{\partial}$.

Using $\text{dim}(\hat{S}^\lambda) + \hat{\partial}$ as an estimate of degrees of freedom leads to the risk estimate

$$\widehat{\text{Risk}}_{\text{df}} := \|Y - \hat{\mu}_{1\text{-OLS}}^\lambda\|_2^2 - n\sigma^2 + 2\sigma^2(\text{dim}(\hat{S}^\lambda) + \hat{\partial}). \tag{7}$$

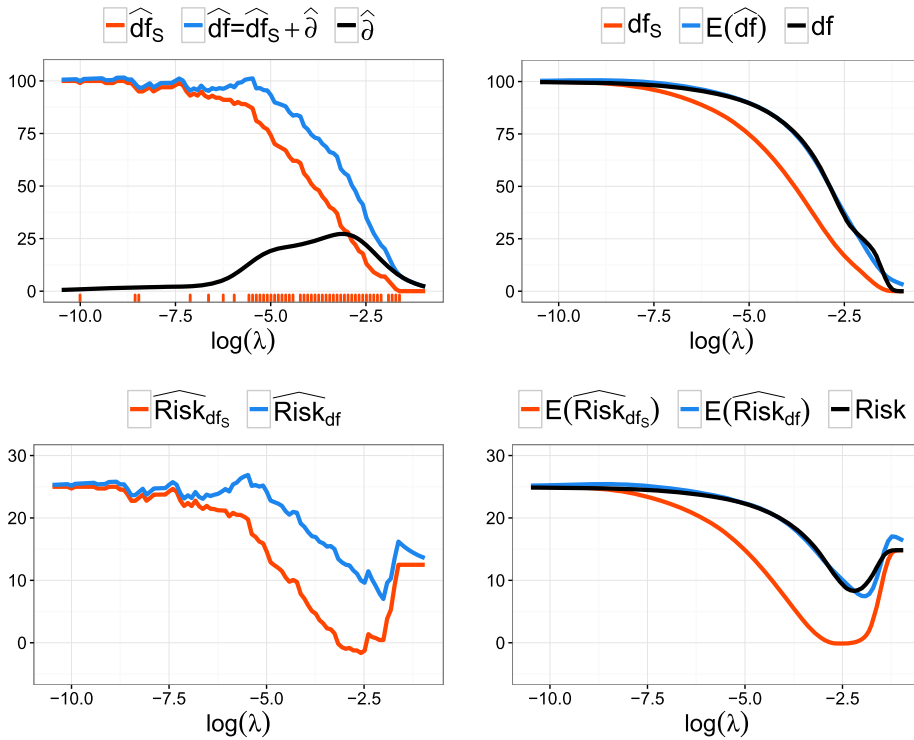


Fig. 2. Left: Realization of the estimates of degrees of freedom $\hat{df}_S = \dim(\hat{S}^\lambda)$ and $\hat{df} = \dim(\hat{S}^\lambda) + \hat{\delta}$ as well as the correction term $\hat{\delta}$ as a function of $\log(\lambda)$ (top) and corresponding estimates of the risk (bottom). Right: Similar to the left but mean values of the estimates obtained by averaging over 1000 samples along with the degrees of freedom $df = df(\hat{\mu}_{1-OLS}^\lambda)$ obtained from the 1000 samples using the covariance definition (1). The design parameters were: $\sigma = 0.5$, $n = p = 100$, $\gamma = 1$, $\alpha = 0.1$ and the design type was (S) with constant correlation of $\rho = 0.1$ (see Section 4).

For an example of the above estimate see Figure 2, where $\hat{\delta}$ and \widehat{Risk}_{df} are applied to a single realization of Y along with an average over 1000 replications.

To prove Theorem 3.2 we prove a more general intermediate result for estimators that are parametrized in a similar way by a tuning parameter. We use in the following D to denote the differential operator w.r.t. y .

Proposition 3.3. *Let $q > 0$ and suppose that $\hat{\mu}^\lambda = \sum_i 1_{U_i^\lambda} \hat{\mu}_i$ where*

$$U_i^\lambda = \lambda^q U_i^1, \quad \text{for all } i = 1, \dots, N. \tag{8}$$

Assume that $\text{div}(\hat{\mu}_i)$ is locally Lipschitz and both $\text{div}(\hat{\mu}_i)$ and $D(\text{div}(\hat{\mu}_i))$ are polynomially bounded for each $i = 1, \dots, N$. If $\hat{\mu}^1$ satisfies Assumption 2.2 then

$$-\frac{\lambda}{q} \partial_\lambda df_S(\hat{\mu}^\lambda) = \frac{1}{2} \sum_{i \neq j} \int_{U_i^\lambda \cap \bar{U}_j^\lambda} (\text{div}(\hat{\mu}_j)(y) - \text{div}(\hat{\mu}_i)(y)) \langle y, \eta_i \rangle \psi(y; \mu, \sigma^2) d\mathcal{H}^{n-1}(y). \tag{9}$$

Proof. First observe that $\partial U_i^\lambda \cap B(0, r) = \lambda^q (\partial U_i^1 \cap B(0, r/\lambda^q))$, hence if $\hat{\mu}^1$ satisfies Assumption 2.2 so does $\hat{\mu}^\lambda$ for all λ . Next, the change of variable formula yields

$$\begin{aligned} df_S(\hat{\mu}^\lambda) &= \int \psi(y) \text{div}(\hat{\mu}^\lambda)(y) dy = \sum_i \int_{U_i^\lambda} \psi(y) \text{div}(\hat{\mu}_i)(y) dy \\ &= \sum_i \int_{U_i^1} \lambda^{qn} (\psi \text{div}(\hat{\mu}_i))(\lambda^q z) dz. \end{aligned}$$

Here $\psi = \psi(\cdot; \mu, \sigma^2)$ to ease notation.

The last integrand is differentiable w.r.t. λ (for Lebesgue a.a. z) and its derivative is

$$\begin{aligned} & qn\lambda^{qn-1}(\psi \operatorname{div}(\hat{\mu}_i))(\lambda^q z) + \lambda^{qn} \langle D(\psi \operatorname{div}(\hat{\mu}_i))(\lambda^q z), q\lambda^{q-1} z \rangle \\ &= \frac{q}{\lambda} \lambda^{qn} (n(\psi \operatorname{div}(\hat{\mu}_i))(\lambda^q z) + \langle D(\psi \operatorname{div}(\hat{\mu}_i))(\lambda^q z), \lambda^q z \rangle), \end{aligned}$$

which is dominated in a neighbourhood of λ by an integrable function due to the polynomial bounds. Hence, by the change of variable formula

$$\begin{aligned} \frac{\lambda}{q} \partial_\lambda \operatorname{df}_S(\hat{\mu}^\lambda) &= \sum_i \int_{U_i^\lambda} \lambda^{qn} (n(\psi \operatorname{div}(\hat{\mu}_i))(\lambda^q z) + \langle D(\psi \operatorname{div}(\hat{\mu}_i))(\lambda^q z), \lambda^q z \rangle) dz \\ &= \sum_i \int_{U_i^\lambda} n(\psi \operatorname{div}(\hat{\mu}_i))(y) + \langle D(\psi \operatorname{div}(\hat{\mu}_i))(y), y \rangle dy \\ &= \sum_i \int_{U_i^\lambda} n(\psi \operatorname{div}(\hat{\mu}_i))(y) + \langle (\psi D \operatorname{div}(\hat{\mu}_i) + \operatorname{div}(\hat{\mu}_i) D \psi)(y), y \rangle dy \\ &= \sum_i \int_{U_i^\lambda} \psi(y) \operatorname{div}(y \operatorname{div}(\hat{\mu}_i)(y)) + \langle D\psi(y), y \operatorname{div}(\hat{\mu}_i)(y) \rangle dy. \end{aligned}$$

The last line is identified as $\operatorname{df}_S(\tilde{\mu}^\lambda) - \operatorname{df}(\tilde{\mu}^\lambda)$, where

$$\tilde{\mu}^\lambda(y) := \sum_i 1_{U_i^\lambda}(y) y \operatorname{div}(\hat{\mu}_i)(y).$$

Finally (9) follows by applying Theorem 2.4 to $\tilde{\mu}^\lambda$ (which also satisfies Assumption 2.2). \square

Example 3.4. There are naturally occurring examples besides the lasso selection sets that satisfy (8). Consider still a linear regression setup with X an $n \times p$ -matrix. Let ℓ denote the penalized loss function

$$\ell(y, \beta, \lambda) = \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \operatorname{Pen}(\beta),$$

for some penalty function $\operatorname{Pen} : \mathbb{R}^p \rightarrow \mathbb{R}$ and define the sets

$$U_A^\lambda = \operatorname{int} \left\{ y \in \mathbb{R}^n \mid \inf_{\beta: \operatorname{supp}(\beta)=A} \ell(y, \beta, \lambda) = \inf_{\beta} \ell(y, \beta, \lambda) \right\}, \quad (10)$$

for each $A \subseteq \{1, \dots, p\}$. Hence any $y \in U_A^\lambda$ has A as an active set. If Pen is *positive homogeneous* of degree $k \in [0, 2)$ then

$$\ell(\lambda^{\frac{1}{2-k}} y, \lambda^{\frac{1}{2-k}} \beta, \lambda) = \lambda^{\frac{2}{2-k}} \ell(y, \beta, 1).$$

Hence $U_A^\lambda = \lambda^{\frac{1}{2-k}} U_A^1$ holds for all $A \subseteq \{1, \dots, p\}$ and $\lambda > 0$. The (quasi) norms, $\operatorname{Pen}(\beta) = \|\beta\|_k^k$ for $k \in (0, 2)$, and $\operatorname{Pen}(\beta) = \|\beta\|_0 = |\operatorname{supp}(\beta)|$ are examples of positive homogeneous penalties. For these penalties only $k \in [0, 1]$ will result in variable selection. With $\operatorname{Pen}(\cdot) = \|\cdot\|_1$ we see that for lasso the sets U_S^λ in 2.5 satisfy (8) with $q = 1$.

Proof of Theorem 3.2. Let $(U_S^\lambda)_{S \in \mathcal{S}}$ be defined as in Example 2.5, where it was also shown that Assumption 2.2 holds for the lasso-OLS estimator. Moreover, from Example 3.4 we see that $U_S^\lambda = \lambda U_S^1$ for all $\lambda > 0$ and $S \in \mathcal{S}$. By Theorem 2.4 we know that the left hand side of (6) is

$$\begin{aligned} & \operatorname{df}(\hat{\mu}_{1\text{-OLS}}^\lambda) - \operatorname{df}_S(\hat{\mu}_{1\text{-OLS}}^\lambda) \\ &= \frac{1}{2} \sum_{S_1 \neq S_2} \int_{\bar{U}_{S_1}^\lambda \cap \bar{U}_{S_2}^\lambda} \langle (\Pi_{S_2} - \Pi_{S_1})y, \eta_{S_1}(y) \rangle \psi(y) d\mathcal{H}^{n-1}(y). \end{aligned} \quad (11)$$

It will first be established that $\overline{U}_{S_1}^\lambda \cap \overline{U}_{S_2}^\lambda$ for $S_1 \neq S_2$ is a \mathcal{H}^{n-1} null set unless S_1 and S_2 are nested and their dimensions differ by one.

By definition $\hat{\mu}_{\text{lasso}}^\lambda \in S$ on U_S^λ , and by continuity of $\hat{\mu}_{\text{lasso}}^\lambda$ (a consequence of Lemma 3 in [25]) we conclude that the same is true on \overline{U}_S^λ . Hence for $S_1, S_2 \in \mathcal{S}$

$$\hat{\mu}_{\text{lasso}}^\lambda \in S_1 \cap S_2 \quad \text{on } \overline{U}_{S_1}^\lambda \cap \overline{U}_{S_2}^\lambda. \quad (12)$$

For $A \subseteq \{1, \dots, p\}$ and $s \in \{-1, 1\}^{|A|}$ we define the set

$$L_{A,s} := \{u \in \mathbb{R}^n \mid X_A^T u = \lambda s\}.$$

It now follows from the first order subgradient conditions for lasso that

$$y - \hat{\mu}_{\text{lasso}}^\lambda \in \bigcup_{\substack{A \subseteq \{1, \dots, p\}: \\ \text{col}(X_A) = S}} \bigcup_{s \in \{-1, 1\}^{|A|}} L_{A,s} \quad (13)$$

for all $y \in U_S^\lambda$. Note that the dimension of the above set is $n - \dim(S)$. Since the set is closed and $\hat{\mu}_{\text{lasso}}^\lambda$ is continuous, (13) holds for $y \in \overline{U}_S^\lambda$ as well. We therefore conclude that

$$\begin{aligned} y - \hat{\mu}_{\text{lasso}}^\lambda &\in \left(\bigcup_{\substack{A \subseteq \{1, \dots, p\}: \\ \text{col}(X_A) = S_1}} \bigcup_{s \in \{-1, 1\}^{|A|}} L_{A,s} \right) \cap \left(\bigcup_{\substack{A \subseteq \{1, \dots, p\}: \\ \text{col}(X_A) = S_2}} \bigcup_{s \in \{-1, 1\}^{|A|}} L_{A,s} \right) \\ &\subseteq \bigcup_{\substack{A \subseteq \{1, \dots, p\}: \\ \text{col}(X_A) = S_1 + S_2}} \bigcup_{s \in \{-1, 1\}^{|A|}} L_{A,s} \end{aligned} \quad (14)$$

for all $y \in \overline{U}_{S_1}^\lambda \cap \overline{U}_{S_2}^\lambda$ and $S_1, S_2 \in \mathcal{S}$.

From (12) and (14) we deduce that

$$\overline{U}_{S_1}^\lambda \cap \overline{U}_{S_2}^\lambda \subseteq S_1 \cap S_2 + \bigcup_{\substack{A \subseteq \{1, \dots, p\}: \\ \text{col}(X_A) = S_1 + S_2}} \bigcup_{s \in \{-1, 1\}^{|A|}} L_{A,s} \quad (15)$$

for $S_1, S_2 \in \mathcal{S}$. Consequently, if $S_1 \neq S_2$ then $\mathcal{H}^{n-1}(\overline{U}_{S_1}^\lambda \cap \overline{U}_{S_2}^\lambda) = 0$, unless S_1 and S_2 are nested and their dimensions differ by 1.

We can therefore assume $S_1 \subseteq S_2$ and $\dim(S_2) = \dim(S_1) + 1$. Furthermore, $S_2 \ominus S_1 = (S_1 + S_2) \ominus (S_1 \cap S_2)$ is orthogonal to any of the faces $S_1 \cap S_2 + L_{A,s}$ in (15) and thus also orthogonal to $\overline{U}_{S_1}^\lambda \cap \overline{U}_{S_2}^\lambda$. This implies that $\eta_{S_1} = (\Pi_{S_2} - \Pi_{S_1})\eta_{S_1}$ and hence (11) becomes

$$\begin{aligned} &\text{df}(\hat{\mu}_{1\text{-OLS}}^\lambda) - \text{dfs}(\hat{\mu}_{1\text{-OLS}}^\lambda) \\ &= \sum_{\substack{S_1 \subseteq S_2, \\ \dim(S_2) = \dim(S_1) + 1}} \int_{\overline{U}_{S_1}^\lambda \cap \overline{U}_{S_2}^\lambda} \langle y, \eta_{S_1}(y) \rangle \psi(y) d\mathcal{H}^{n-1}(y) \\ &= \sum_{\substack{S_1 \subseteq S_2, \\ \dim(S_2) = \dim(S_1) + 1}} \int_{\overline{U}_{S_1}^\lambda \cap \overline{U}_{S_2}^\lambda} \underbrace{[\text{div}(\Pi_{S_2}) - \text{div}(\Pi_{S_1})]}_{=\dim(S_2) - \dim(S_1) = 1} \langle y, \eta_{S_1}(y) \rangle \psi(y) d\mathcal{H}^{n-1}(y) \\ &= -\lambda \partial_\lambda \text{dfs}(\hat{\mu}_{1\text{-OLS}}^\lambda) \end{aligned}$$

by Proposition 3.3. □

4. Simulation study

We report in this section the results from an extensive simulation study, whose purpose was to quantify how $\widehat{\text{Risk}}_{\text{df}}$ given by (7) performs as an estimate of the risk and in terms of selecting the penalty parameter λ . Its performance was compared to alternatives for risk estimation and tuning, and the resulting lasso-OLS estimator was compared to the lasso estimator. Throughout, the R package *glmnet*, [10], was used to compute the lasso solution path. This section is divided into subsections describing estimators and risk estimates, the design of the simulation study, and the results of the simulation study.

4.1. Estimators and risk estimates

The first alternative risk estimate for lasso-OLS is

$$\widehat{\text{Risk}}_{\text{dfs}} = \|Y - \hat{\mu}_{1\text{-OLS}}^\lambda\|_2^2 - n\sigma^2 + 2\sigma^2 \dim(\hat{S}^\lambda), \quad (16)$$

which does not adjust for the variable selection performed by lasso-OLS. The second alternative is K -fold cross-validation (denoted $\widehat{\text{Risk}}_{\text{CV-K}}$) with $K = 5, 10$. This risk estimate is given by

$$\widehat{\text{Risk}}_{\text{CV-K}} := \sum_{k=1}^K \|Y_k - X_k \hat{\beta}_{1\text{-OLS}}^\lambda(Y_{-k}, X_{-k})\|_2^2 - n\sigma^2, \quad (17)$$

where Y_k and X_k denote the entries of Y and rows of X , respectively, corresponding to the k th fold, and similarly, Y_{-k} and X_{-k} denote the entries and rows not in the k th fold.

The lasso estimator was tuned by minimising the risk estimate

$$\widehat{\text{Risk}}_{\text{lasso}} = \|Y - \hat{\mu}_{\text{lasso}}^\lambda\|_2^2 - n\sigma^2 + 2\sigma^2 \dim(\hat{S}^\lambda). \quad (18)$$

For tuning $\in \{\text{df}, \text{dfs}, \text{CV-5}, \text{CV-10}, \text{lasso}\}$ we let $\hat{\lambda}_{\text{tuning}}$ denote the value of λ that minimises $\widehat{\text{Risk}}_{\text{tuning}}$. The risk of the resulting estimator is denoted

$$\text{Risk}(\text{tuning}) := E \| \mu - \hat{\mu}_{1\text{-OLS}}^{\hat{\lambda}_{\text{tuning}}} \|_2^2$$

for all but the lasso-tuning, whose risk instead is

$$\text{Risk}(\text{lasso}) := E \| \mu - \hat{\mu}_{\text{lasso}}^{\hat{\lambda}_{\text{lasso}}} \|_2^2.$$

When the true mean is $\mu = X\beta$ with $\text{supp}(\beta) = A$ we refer to Π_A as the oracle-OLS estimator. This usage of the oracle terminology is in accordance with e.g. [8]. Its risk is

$$E \| \mu - \Pi_A Y \|_2^2 = \sigma^2 \text{rank}(X_A).$$

The results from the simulation study are reported in terms of $\text{Risk}(\text{tuning})/(\sigma^2 n)$ for each tuning method, which can then be compared to $\text{rank}(X_A)/n$ – the fraction of nonzero parameters.

All simulations were carried out assuming either that σ^2 was known or using the following estimator of σ^2 : first the lasso path $\lambda \mapsto \hat{\mu}_{\text{lasso}}^\lambda$ was calculated, then $\hat{\lambda}$ was selected by minimising the generalized cross-validation criterion

$$\text{gcv}(\lambda) = \frac{\|Y - \hat{\mu}_{\text{lasso}}^\lambda\|_2^2}{(1 - \frac{\dim(\hat{S}^\lambda)}{n})^2},$$

and σ^2 was finally estimated as

$$\hat{\sigma}^2 = \frac{\|Y - \hat{\mu}_{\text{lasso}}^{\hat{\lambda}}\|_2^2}{n - \dim(\hat{S}^{\hat{\lambda}})}.$$

The main reason for choosing this estimator was computational efficiency, as the lasso path must be calculated for lasso-OLS anyway. Thus this variance estimate has virtually no extra computational costs. See also [20] for a comprehensive comparison of variance estimators.

4.2. Simulation study design

In the simulation study the mean was given as $X\beta$ with

$$\beta_i = \begin{cases} \gamma^{i-1} & \text{if } i \leq \lceil n\alpha \rceil, \\ 0 & \text{otherwise,} \end{cases}$$

for different choices of the dimension n , the $n \times p$ design matrix X and the parameters γ and α .

Two simulation designs were implemented with parameters as follows:

Parameter	Values for simulation study I					Values for simulation study II				
σ	0.5					0.1	0.2	0.5	1	2
α	0.1					0	0.05	0.1	0.3	0.5
n	50	100	200	400	800	100	200			
p	200	2000	20,000			n				
γ	1					1	0.9			
X	S					O	S	E		
ρ	0.1					0	0.1	0.4	0.7	

The parameter ρ and the values of the design require some explanation. The three different design types are:

- Orthogonal (O), where $X = I$.
- Simulated (S), where the columns of X are standard normally distributed with one of the following correlation structures:
 - Autoregressive setup: $\text{corr}(X_i, X_j) = \rho^{|i-j|}$ for all $i \neq j$.
 - Constant correlation setup: $\text{corr}(X_i, X_j) = \rho$ for all $i \neq j$.
- Empirical (E), where the rows and columns are randomly selected from the 240×377 matrix of microRNA expression values as used in the earlier study by [26].

The columns of the simulated and empirical designs were standardized to have norm one to obtain a comparable signal-to-noise ratio across the three designs.

The risk estimates were based on 1000 samples for each combination of the parameters, which were generated as follows. For each of the 1000 samples a design matrix X was created/simulated and a single realization of $Y \sim \mathcal{N}(X\beta, \sigma^2 I_n)$ was drawn. For each sample the losses $\|\mu - \hat{\mu}_{\text{lasso}}^{\hat{\lambda}_{\text{lasso}}}\|_2^2$ and $\|\mu - \hat{\mu}_{1\text{-OLS}}^{\hat{\lambda}_{\text{tuning}}}\|_2^2$ for the different tuning methods were computed. The risks were estimated as the average of the losses over the 1000 samples.

In order to assess robustness to deviations from the Gaussian noise assumption, we replicated the second study design with two types of non-Gaussian noise: a t -distribution with 3 degrees of freedom, and a skew normal distribution with shape parameter 3. Location and scale parameters were set so that the noise distribution had mean 0 and variance σ^2 .

4.3. Results from study I

We first report on the accuracy of the risk estimates. Figure 3 shows the risk estimates as a function of λ for 50 samples along with a Monte Carlo estimate of the true risk. Cross-validation appears to give more variable estimates of the risk than $\widehat{\text{Risk}}_{\text{df}}$ across the entire range of λ -values. This is true even when the variance is estimated, though estimation of the variance does appear to degrade the performance of the risk estimates. We note that $\widehat{\text{Risk}}_{\text{df}}$ does not appear to be much more variable than $\widehat{\text{Risk}}_{\text{lasso}}$, though the former relies on the additional smoothed term for the estimation of degrees of freedom.

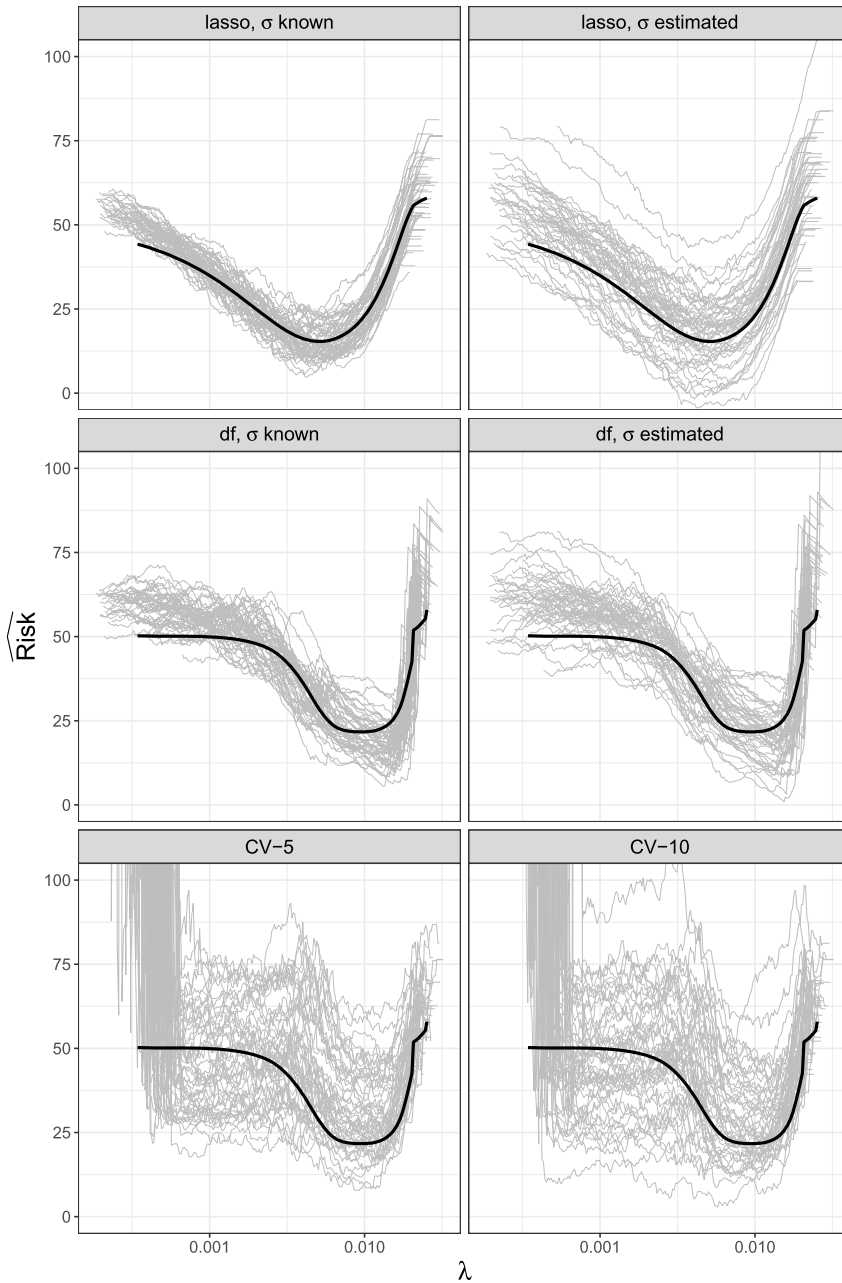


Fig. 3. Risk estimates $\widehat{\text{Risk}}_{\text{df}}$, $\widehat{\text{Risk}}_{\text{CV-5}}$, $\widehat{\text{Risk}}_{\text{CV-10}}$ and $\widehat{\text{Risk}}_{\text{lasso}}$ (gray lines) for 50 samples as a function of λ . The black lines are Monte Carlo estimates of the true risks. The design parameters were: $n = 200$, $p = 2000$, $\sigma = 0.5$, $\gamma = 1$, $\alpha = 0.1$, and the design type was (S) with a constant correlation of $\rho = 0.1$ (see Section 4.2).

Figure 4 shows mean squared errors (MSEs) for the risk estimates. The figure shows the integrated mean squared error as well as the mean squared error in the optimal λ (the λ that minimizes risk as estimated from the Monte Carlo estimate of the risk based on 1000 replications). The cross-validation risk estimates generally have the largest MSEs, while $\widehat{\text{Risk}}_{\text{df}}$ has considerably smaller MSEs. This is true even when the variance is estimated except for $n = 50$ and $p = 2000, 20,000$. From this figure we see that $\widehat{\text{Risk}}_{\text{df}}$ does have a larger MSE than $\widehat{\text{Risk}}_{\text{lasso}}$. Moreover, for n/p large the estimation of σ does not affect the MSE of the risk estimates much.

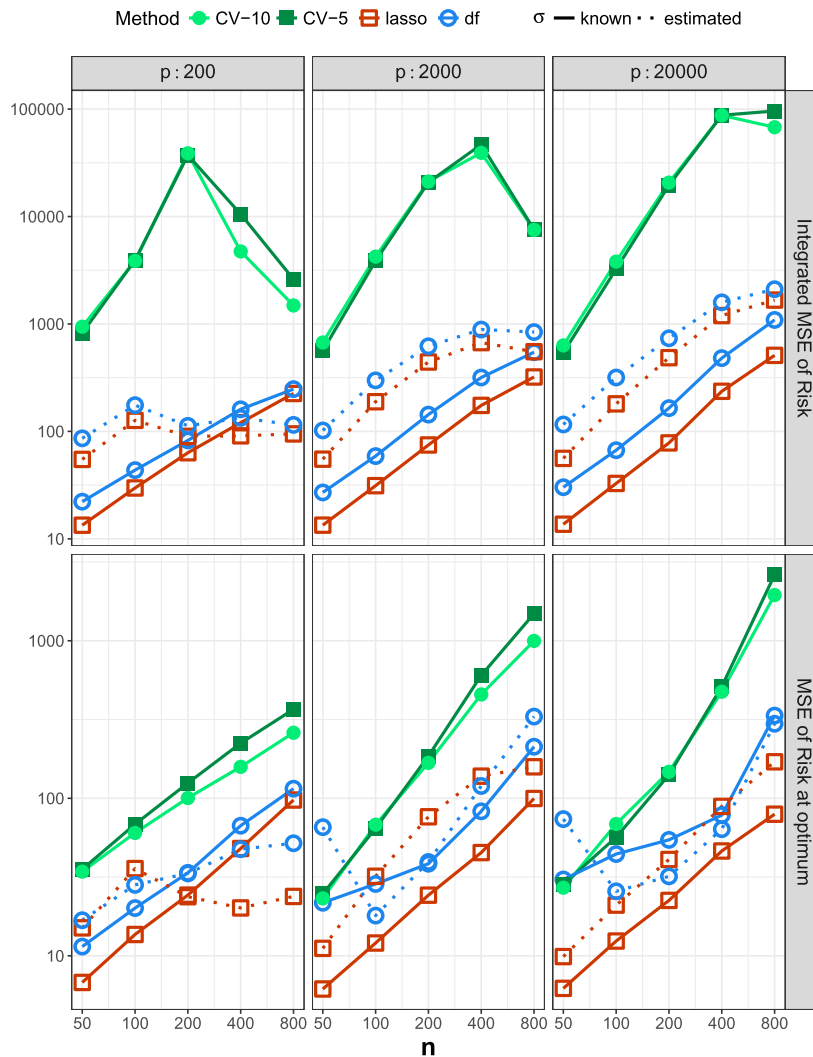


Fig. 4. Integrated mean squared error (top) and mean squared error at the optimal value of λ , $\hat{\lambda}$ (bottom) of the risk estimates $\widehat{\text{Risk}}_{\text{df}}$, $\widehat{\text{Risk}}_{\text{CV-5}}$, $\widehat{\text{Risk}}_{\text{CV-10}}$ and $\widehat{\text{Risk}}_{\text{lasso}}$. The integrated mean squared error was computed over the interval $[\hat{\lambda}/10, 10\hat{\lambda}]$ of $\log(\lambda)$ -values. The design parameters were: $\sigma = 0.5$, $\gamma = 1$, $\alpha = 0.1$, and the design type was (S) with a constant correlation of $\rho = 0.1$ (see Section 4.2).

For this simulation study we also recorded the number of selected predictors as well as the computational time for evaluating and tuning the different estimators. The results can be found as Figure 1 in the supplementary material [19]. The lasso-OLS estimator selects fewer predictors than lasso, but when the variance is estimated, the number of selected predictors is increased – this is particularly so when n/p is small. The lasso estimator using (18) for tuning is fastest, which is unsurprising as the computation of the lasso path is part of all estimators. Moreover, the lasso-OLS estimator using (7) for tuning is about a factor 4 faster than using 5-fold cross-validation for tuning and about a factor 8 faster than 10-fold cross-validation. Thus the added computation of the smoothed term to the estimate of degrees of freedom in (7) has an insignificant effect on the computation time.

4.4. Results from study II

Firstly, we discuss the comparison of the two tuning methods df and df_5 for the lasso-OLS estimator. The purpose of this comparison is to highlight the effect of correctly adjusting for the variable selection in the estimation of degrees of freedom via the term $\hat{\delta}$. Secondly, we discuss the comparison of df to CV-5, CV-10 and lasso. The purpose of this

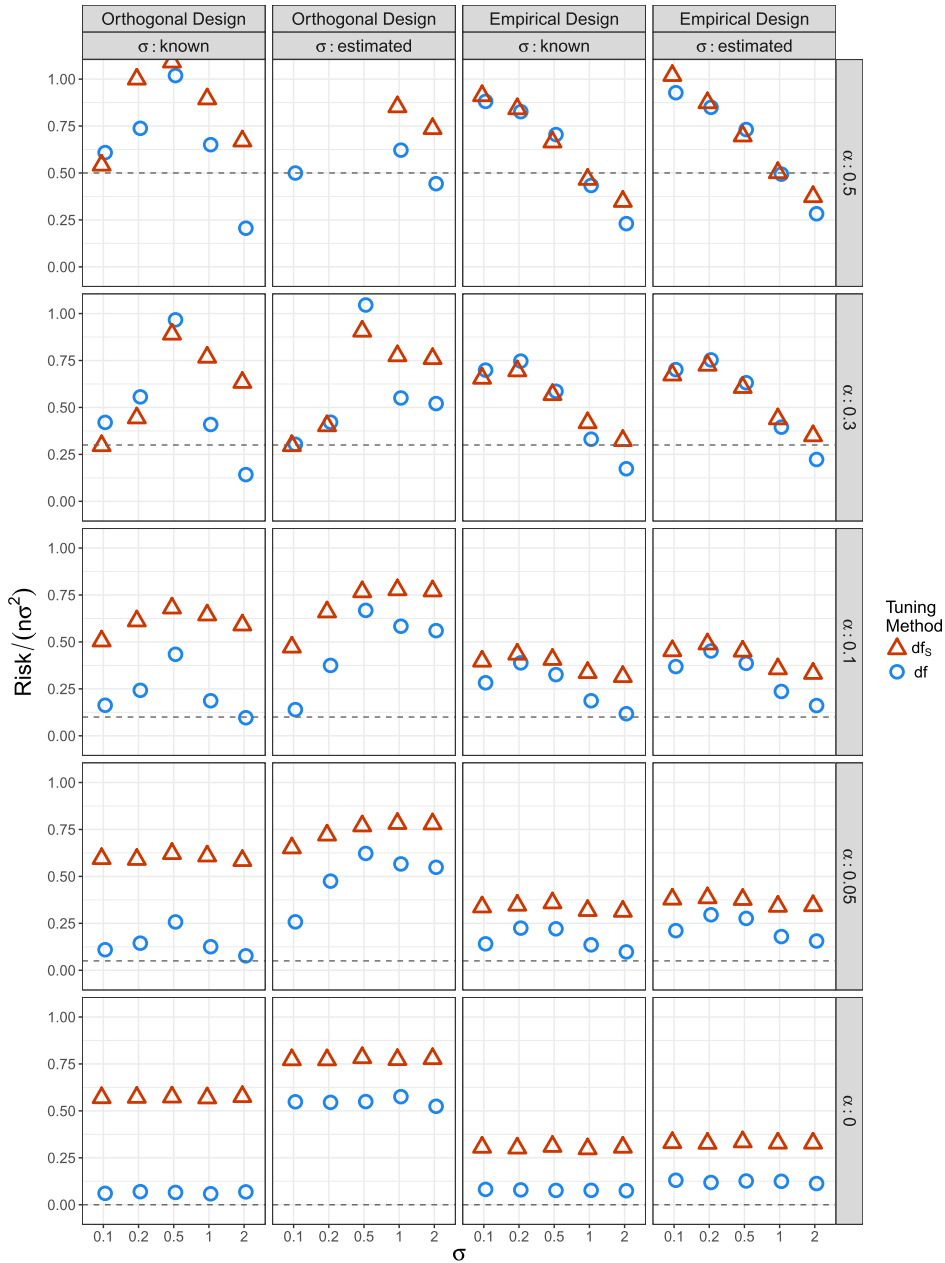


Fig. 5. Risk relative to $\sigma^2 n$ for the estimators $\hat{\mu}_{1-OLS}^{\hat{\lambda}_{df_s}}$ and $\hat{\mu}_{1-OLS}^{\hat{\lambda}_{df}}$ for orthogonal and empirical designs with $n = 100$ and $\gamma = 1$. The dashed line is $\lceil n\alpha \rceil / n \simeq \alpha$, the relative risk for the oracle-OLS estimator.

second comparison is two-fold. It provides a comparison of our proposed tuning method, df , to cross-validation based tuning, and it provides a comparison of lasso-OLS to lasso in terms of predictive performance.

Figure 5 shows the results for the two tuning methods df and df_s in the orthogonal and empirical designs with $\gamma = 1$ and $n = 100$. The results for all the other design parameters can be found in [19]. Tuning λ by using $\dim(\hat{S}^\lambda) + \hat{\varrho}$ as an estimate of degrees of freedom is generally superior to using $\dim(\hat{S}^\lambda)$ and in the worst cases at least comparable. The differences are largest for the lowest signal-to-noise ratios. The benefit of using $\dim(\hat{S}^\lambda) + \hat{\varrho}$ generally increases with the dimension n , and it increases with decreasing signal-to-noise ratio. Furthermore, when the number of non-zero parameters is large and the signal-to-noise ratio is low (specifically, $\gamma = 0.9$, α large and σ large), $\hat{\mu}_{1-OLS}^{\hat{\lambda}_{df}}$ clearly

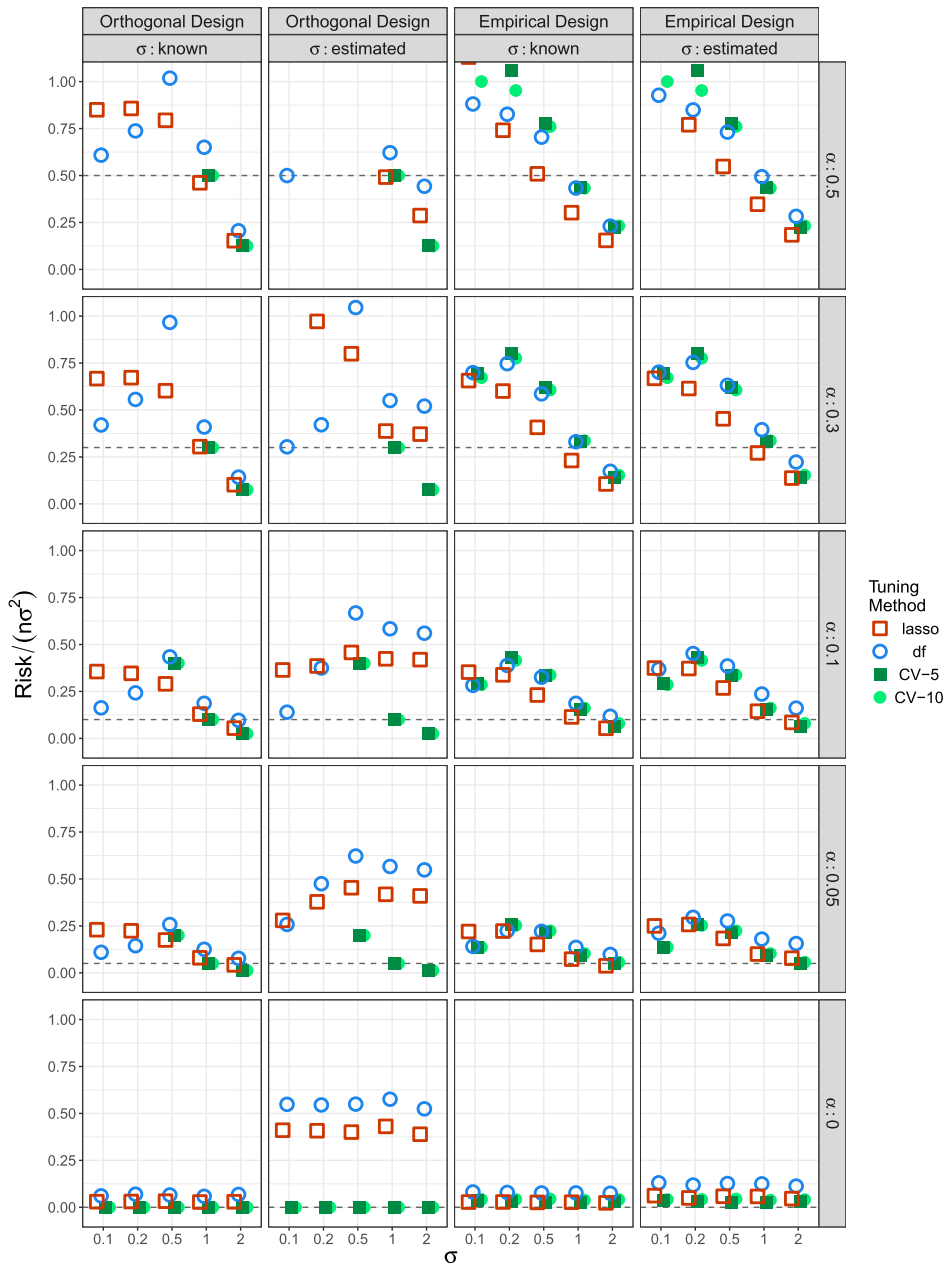


Fig. 6. Risk relative to $\sigma^2 n$ for the estimators $\hat{\mu}_{1-OLS}^{\hat{\lambda}_{df}}$, $\hat{\mu}_{1-OLS}^{\hat{\lambda}_{CV-5}}$, $\hat{\mu}_{1-OLS}^{\hat{\lambda}_{CV-10}}$ and $\hat{\mu}_{1-OLS}^{\hat{\lambda}_{lasso}}$ for orthogonal and empirical designs with $n = 100$ and $\gamma = 1$. The dashed line is $\lceil n\alpha \rceil / n \simeq \alpha$, the relative risk for the oracle-OLS estimator.

outperforms the oracle-OLS estimator, while $\hat{\mu}_{1-OLS}^{\hat{\lambda}_{df}}$ is comparable or worse than the oracle-OLS estimator. Neither of the estimators performs well for small variances and large signal-to-noise ratios. For the orthogonal design the estimation of the variance incurs a clear performance loss, which is not the case for the other designs. We ascribe this to the variance estimator being particularly poor for the orthogonal design.

Figure 6 shows the results for df, CV-5, CV-10 and lasso for the orthogonal and empirical designs with $\gamma = 1$ and $n = 100$. The results for the remaining design parameters are found in [19]. For the orthogonal design cross-validation is not an appropriate tuning method, since $\text{Risk}_{CV-\kappa}$ is constant in λ . This relates to the fact that the folds cannot be considered replications of the same distribution. Consequently, for the orthogonal design, the tuning methods based

on degrees of freedom have clear advantages. On the other hand, the estimation of σ has a quite large negative effect for precisely the orthogonal design.

When restricting attention to the non-orthogonal designs we observe that the tuning methods are quite comparable (see [19]). None of the tuning methods are generally superior or inferior to the others and their performance depends on both design type, signal-to-noise ratio and the signal decay parameter γ . The lasso estimator deviates most from the others, which is mainly due to this being a different estimator. It performs best at low signal-to-noise ratios, while lasso-OLS using either cross-validation or df tuning performs better at high signal-to-noise ratios (α large, σ small and $\gamma = 1$). Cross-validation appears to perform best for highly correlated designs (ρ large).

The results for the non-Gaussian error distributions are included in [19] as well. There are no major differences when compared to the Gaussian error distribution, with the most notable change being that lasso loses some of its performance for the t -distributed noise. The tuning based on df seems to be less affected. Still, all the tuning methods are generally comparable except for orthogonal designs. Since cross-validation does not rely on a Gaussian noise assumption, these results suggest that our proposed tuning method based on df is appropriate even in non-Gaussian settings.

5. Best subset selection

Example 3.4 demonstrates that (8) holds for other estimators than lasso-OLS, and Theorem 3.3 holds, in particular, for best subset selection in the Lagrangian formulation, which corresponds to $\text{Pen}(\cdot) = \|\cdot\|_0$ in Example 3.4. Theorem 3.2 does, however, only partly extend to best subset selection. In this section we demonstrate that this may still provide a practically useful estimate of degrees of freedom.

The best subset selection estimator of μ with tuning parameter $\lambda > 0$, denoted by $\hat{\mu}_{\text{bs}}^\lambda$, is

$$\hat{\mu}_{\text{bs}}^\lambda = X\hat{\beta}^\lambda \quad \text{where } \hat{\beta}^\lambda = \arg \min_{\beta} \frac{1}{2} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_0.$$

It can be written on the form $\hat{\mu}_{\text{bs}}^\lambda = \sum_{A \subseteq \{1, \dots, p\}} 1_{U_A^\lambda} \Pi_A$ (Lebesgue a.e.), where

$$U_A^\lambda := \left\{ y \in \mathbb{R}^n \mid \lambda|A| - \frac{1}{2} \|\Pi_A y\|_2^2 < \min_{B \subseteq \{1, \dots, p\} \setminus A} \lambda|B| - \frac{1}{2} \|\Pi_B y\|_2^2 \right\}, \quad A \subseteq \{1, \dots, p\}. \tag{19}$$

It is straightforward to verify that $\hat{\mu}_{\text{bs}}^\lambda$ fulfils Assumption 2.2 except 2.2(c), which follows by Lemma A.1 in the appendix. Hence Theorem 2.4 applies to $\hat{\mu}_{\text{bs}}^\lambda$.

From (19) we note that the outer unit normal to $\partial U_{A_1}^\lambda$ on $\bar{U}_{A_1}^\lambda \cap \bar{U}_{A_2}^\lambda$ equals $(\Pi_{A_2} - \Pi_{A_1})y$ normalized to have norm 1. Theorem 2.4 yields

$$\begin{aligned} \text{df}(\hat{\mu}_{\text{bs}}^\lambda) - \text{df}_S(\hat{\mu}_{\text{bs}}^\lambda) &= \frac{1}{2} \sum_{A_1 \neq A_2} \int_{\bar{U}_{A_1}^\lambda \cap \bar{U}_{A_2}^\lambda} \frac{\langle (\Pi_{A_2} - \Pi_{A_1})y, (\Pi_{A_2} - \Pi_{A_1})y \rangle}{\|(\Pi_{A_2} - \Pi_{A_1})y\|_2} \psi(y) d\mathcal{H}^{n-1}(y) \\ &= \frac{1}{2} \sum_{A_1 \neq A_2} \int_{\bar{U}_{A_1}^\lambda \cap \bar{U}_{A_2}^\lambda} \|(\Pi_{A_2} - \Pi_{A_1})y\|_2 \psi(y) d\mathcal{H}^{n-1}(y), \end{aligned}$$

which proves that $\text{df} > \text{df}_S$ for best subsection selection. Moreover, Proposition 3.3 and Example 3.4 yields

$$-2\lambda \partial_\lambda \text{df}_S(\hat{\mu}_{\text{bs}}^\lambda) = \frac{1}{2} \sum_{A_1 \neq A_2} \int_{\bar{U}_{A_1}^\lambda \cap \bar{U}_{A_2}^\lambda} \psi(y) \frac{\langle y, (\Pi_{A_2} - \Pi_{A_1})y \rangle}{\|(\Pi_{A_2} - \Pi_{A_1})y\|_2} (|A_2| - |A_1|) d\mathcal{H}^{n-1}(y).$$

For $\text{col}(X_{A_1}) \subseteq \text{col}(X_{A_2})$ and $\text{rank}(X_{A_2}) = \text{rank}(X_{A_1}) + 1$, we see that the integrands in the two identities above coincide. Hence, if we define

$$\begin{aligned} \mathcal{A}_1 &:= \{A_1, A_2 \subseteq \{1, \dots, p\} \mid \text{col}(X_{A_1}) \subseteq \text{col}(X_{A_2}) \text{ and } \text{rank}(X_{A_2}) = \text{rank}(X_{A_1}) + 1\} \quad \text{and} \\ \mathcal{A}_2 &:= \{A_1, A_2 \subseteq \{1, \dots, p\} \mid \text{col}(X_{A_1}) \neq \text{col}(X_{A_2}) \text{ and } (A_1, A_2) \notin \mathcal{A}_1, (A_2, A_1) \notin \mathcal{A}_1\}, \end{aligned}$$

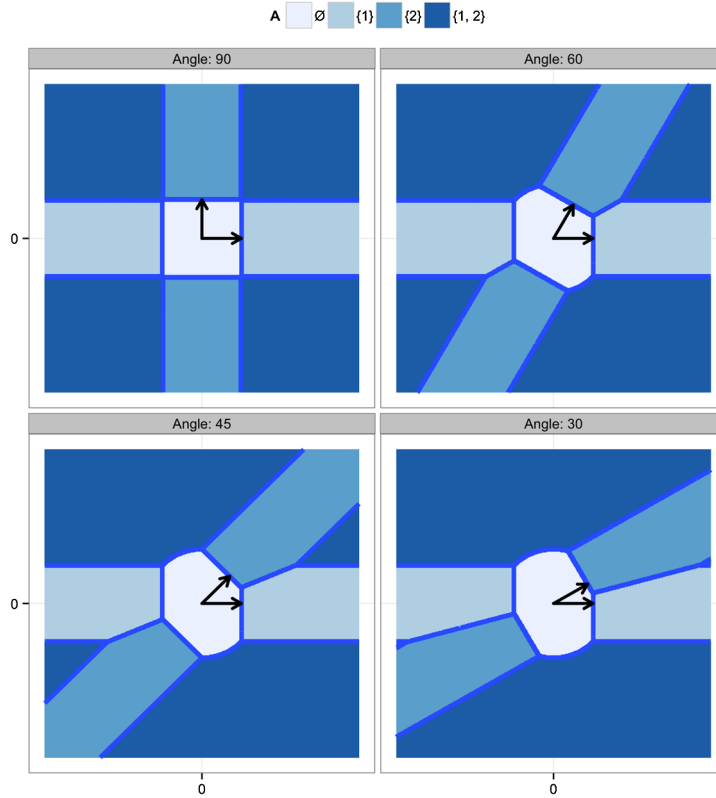


Fig. 7. Illustrations of the decomposition of \mathbb{R}^2 into the four sets U_\emptyset^1 , $U_{\{1\}}^1$, $U_{\{2\}}^1$ and $U_{\{1,2\}}^1$ according to the best subset selection estimator in the Lagrangian formulation with $\lambda = 1$. The set U_\emptyset^1 consists of the points projected onto the 0-dimensional space $\{0\}$, the sets $U_{\{1\}}^1$, $U_{\{2\}}^1$ to the projections onto one of the two 1-dimensional subspaces and $U_{\{1,2\}}^1$ to the identity map. The decomposition depends on the angle between the two columns in X .

then

$$df(\hat{\mu}_{\text{bs}}^\lambda) - df_S(\hat{\mu}_{\text{bs}}^\lambda) = -2\lambda\partial_\lambda df_S(\hat{\mu}_{\text{bs}}^\lambda) + R,$$

where

$$R = \frac{1}{2} \sum_{(A_1, A_2) \in \mathcal{A}_2} \int_{\bar{U}_{A_1}^\lambda \cap \bar{U}_{A_2}^\lambda} \frac{\langle (\Pi_{A_2} - \Pi_{A_1})y, (\Pi_{A_2} - \Pi_{A_1} - (|A_2| - |A_1|)I_n)y \rangle}{\|(\Pi_{A_2} - \Pi_{A_1})y\|_2} \psi(y) d\mathcal{H}^{n-1}(y).$$

The usefulness of this hinges on R being small. For X orthogonal we have already demonstrated that $R = 0$ as $\hat{\mu}_{\text{bs}}^\lambda$ then coincides with lasso-OLS, and in this case $\bar{U}_{A_1}^\lambda \cap \bar{U}_{A_2}^\lambda$ has Hausdorff measure zero for all $(A_1, A_2) \in \mathcal{A}_2$. For non-orthogonal X this is no longer true, see Figure 7. For best subset selection there will generally be boundaries of non-zero Hausdorff measure between many more of the sets \bar{U}_A^λ – boundaries that correspond to including or excluding more than one predictor at the time or replacing predictors. Compare this with lasso-OLS and Figure 1. However, by continuity in X we have $R \rightarrow 0$ for X tending to an orthogonal matrix, and we can expect R to be small for matrices that are not too far from orthogonal matrices. Thus we expect

$$df_S(\hat{\mu}_{\text{bs}}^\lambda) - 2\lambda\partial_\lambda df_S(\hat{\mu}_{\text{bs}}^\lambda) \tag{20}$$

to be a useful approximation for $df(\hat{\mu}_{\text{bs}}^\lambda)$ also for non-orthogonal X .

Using the same procedure for estimating the correction $-2\lambda\partial_\lambda df_S(\hat{\mu}_{\text{bs}}^\lambda)$ as outlined in Section 3 – using $2\hat{\delta}$ instead of $\hat{\delta}$ – we used simulations to investigate if (20) was actually a good approximation of $df(\hat{\mu}_{\text{bs}}^\lambda)$. Figure 8 shows the

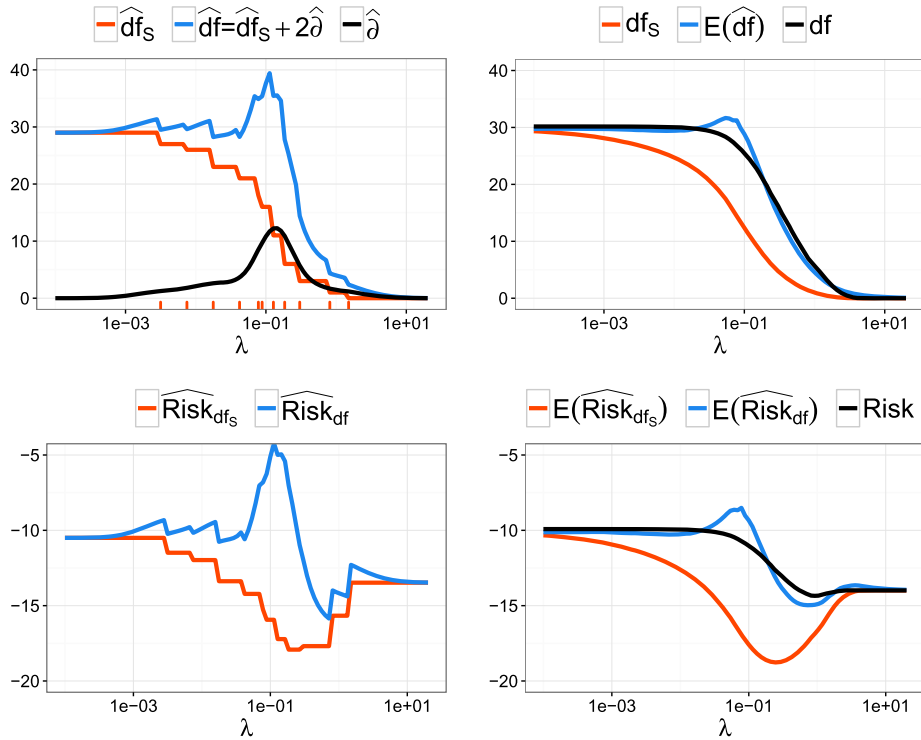


Fig. 8. Left: Realization of the estimates of degrees of freedom $\widehat{df}_S = \dim(\widehat{S}^\lambda)$ and $\widehat{df} = \dim(\widehat{S}^\lambda) + 2\widehat{\partial}$ as well as the correction term $\widehat{\partial}$ as a function of $\log(\lambda)$ for best subset selection (top) and corresponding estimates of the risk (bottom). Right: Similar to the left but mean values of the estimates obtained by averaging over 1000 samples along with the degrees of freedom $df = df(\widehat{\mu}_{bs}^\lambda)$ obtained from the 1000 samples using the covariance definition (1). The design parameters were: $\sigma = 0.5$, $n = p = 30$, $\gamma = 1$, $\alpha = 0.1$ and the design type was (S) with constant correlation of $\rho = 0.1$ (see Section 4).

results using the same configurations as in Figure 2, except that n was lowered to 30 due to computational constraints. The conclusion from this and other similar simulations (not shown) is that even with non-orthogonal designs, (20) is a practically useful approximation. That is, $-2\lambda\partial_\lambda df_S(\widehat{\mu}_{bs}^\lambda)$ accounts for the majority of the increase in the degrees of freedom due to variable selection.

6. Discussion

We have provided a new representation of degrees of freedom for a broad class of discontinuous, piecewise Lipschitz estimators. This representation provides us with a deeper insight into the effect of variable selection, among other things, on the effective dimension of the statistical model and the estimator used. We have demonstrated that for lasso-OLS it was, moreover, possible to derive a practically useful estimator of the degrees of freedom based on the general representation, and we also suggest that a similar estimator can be useful for best subset selection. The estimator was based on relating the derivative of $\lambda \mapsto df_S(\widehat{\mu}^\lambda)$ to the discontinuities of the estimator $\widehat{\mu}^\lambda$ as expressed via the integral representation of $df(\widehat{\mu}^\lambda) - df_S(\widehat{\mu}^\lambda)$. This does, indeed, make some intuitive sense as the first expresses the mean jump of degrees of freedom per unit change of λ and the other (in some sense) the mean discontinuity of degrees of freedom per unit change of y . Changing λ for fixed y or changing y for fixed λ are dual operations, and it is not surprising that we can relate the numbers.

A simulation study demonstrated that the risk of the lasso-OLS estimator can be estimated effectively by using our proposed estimate of degrees of freedom. Our proposal did not incur any substantial computational penalty, nor did it incur a substantial increase in the variance of the risk estimate. The simulation study also showed that lasso-OLS can be effectively tuned by minimising our proposed risk estimate, and that the resulting computations are faster

than using cross-validation. The resulting lasso-OLS estimator selects fewer predictors than lasso with a comparable predictive performance, but it is computationally more expensive.

If we were to generalize our results to other estimators that include a tuning parameter, we expect that it is only the derivative of the part of $\text{df}_S(\hat{\mu}^\lambda)$ that corresponds to jumps that can be related to $\text{df}(\hat{\mu}^\lambda) - \text{df}_S(\hat{\mu}^\lambda)$. That is, in general, $\lambda \mapsto \text{div}(\hat{\mu}^\lambda)$ will have jumps as well as smooth but non-constant pieces, and it is only the expectation of the jump part that we expect can be related to $\text{df}(\hat{\mu}^\lambda) - \text{df}_S(\hat{\mu}^\lambda)$. We believe that our suggested estimator of degrees of freedom may actually be generalizable to a number of discontinuous estimators involving variable selection as well as shrinkage. The requirement will be that the estimator has one or more tuning parameters and that it is computed on a grid or along a path of these. Then we can potentially estimate the derivative of the divergence of the estimator as a function of the tuning parameter(s). It is an ongoing research project to investigate this in detail.

For best subset selection we did not provide any bounds on the residual R in the approximation of $\text{df}(\hat{\mu}^\lambda) - \text{df}_S(\hat{\mu}^\lambda)$. It would, indeed, be very interesting to investigate this approximation in more detail. It would, in particular, be interesting to understand if it in any way can be seen as a “first order approximation” and whether there are higher order terms worth including in some cases.

Finally, we have restricted attention to Gaussian noise in the theoretical derivations. Like Stein’s classical lemma, Theorem 2.4 crucially relies on this assumption. Our simulation study demonstrated some robustness towards deviations from this assumption. However, extensions of Stein’s lemma to non-Gaussian distributions do exist (see, e.g., [3]), but further investigations are required to determine if similar extensions can be made in the more general framework presented in this paper.

7. Supplementary material

The results from the entire simulation study as well as the R-code are available online <http://doi.org/10.5281/zenodo.321847>, [19].

Appendix: Additional results and proofs

A.1. Semialgebraic sets

Observe that for A and B subsets of \mathbb{R}^n it holds that

$$\begin{aligned}\partial A &= \partial(A^c), \\ \partial(A \cup B) &\subseteq \partial A \cup \partial B, \\ \partial(A \cap B) &\subseteq \partial A \cup \partial B.\end{aligned}\tag{21}$$

Especially, the family of sets

$$\{E \in \mathcal{B}(\mathbb{R}^n) \mid r \mapsto \mathcal{H}^{n-1}(\partial E \cap B(0, r)) \text{ is polynomially bounded}\}\tag{22}$$

is stable under complement, finite union and finite intersection. This is a useful observation when we want to verify Assumption 2.2(c).

The following Lemma shows that *semialgebraic sets* belong to the family given by (22). A semialgebraic set is a finite union of finite intersections of sets of the form $(P = 0)$ and $(Q > 0)$, where P and Q are polynomials. A multivariate polynomial is of the form (using multi-index notation)

$$P(x) = \sum_{\alpha \in A} a_\alpha x^\alpha, \quad a_\alpha \in \mathbb{R} \text{ for each } \alpha \in A,$$

with $A \subseteq \mathbb{N}^n$ finite.

Lemma A.1. *If E is semialgebraic then $r \mapsto \mathcal{H}^{n-1}(\partial E \cap B(0, r))$ is polynomially bounded.*

Proof. By the stability under finite set operations of the family given by (22) it suffices to show that $r \mapsto \mathcal{H}^{n-1}((P = 0) \cap B(0, r))$ is polynomially bounded for any nonzero polynomial P . But this follows from Corollary 1 in [16], which implies that

$$\mathcal{H}^{n-1}((P = 0) \cap B(0, r)) \leq \frac{\deg(P)\pi^{\frac{n+1}{2}}}{\Gamma(\frac{n}{2})} r^{n-1}$$

for any nonzero polynomial P with $\deg(P) = \max_{\alpha \neq 0} |\alpha|$ denoting the degree of P . □

A.2. Proof of Theorem 2.4

The following Lemma characterizes the outer unit normal vectors η_i for $i = 1, \dots, N$.

Lemma A.2. *Under Assumption 2.2 the following holds:*

- (a) $\eta_i = 0$ \mathcal{H}^{n-1} a.e. on $\partial U_i \setminus \bigcup_{j \neq i} \overline{U}_j$ for each $i = 1, \dots, N$.
- (b) $\eta_i = -\eta_j$ \mathcal{H}^{n-1} a.e. on $\partial U_i \cap \partial U_j$ with $i \neq j$.
- (c) $\eta_i = 0$ \mathcal{H}^{n-1} a.e. on $\partial U_i \cap \partial U_j \cap \partial U_k$ with i, j, k distinct.

Proof. Firstly, note that the unit outer normal η_i on ∂U_i vanishes outside *the measure theoretic boundary* $\partial_* U_i$, see Definition 5.8 in [7]. Moreover, these two types of boundaries relates to *the reduced boundary* $\partial^* U_i$ (see Definition 5.7 in [7]) by the inclusions:

$$\partial^* U_i \subseteq \partial_* U_i \subseteq \partial U_i.$$

Furthermore, $\mathcal{H}^{n-1}(\partial_* U_i \setminus \partial^* U_i) = 0$ (see Lemma 5.8.1 in [7]). All in all, we see that the lemma holds if we can show the following claims:

$$\begin{aligned} \partial^* U_i &\subseteq \bigcup_{l \neq i} \overline{U}_l, \\ \eta_i &= -\eta_j \quad \text{on } \partial^* U_i \cap \partial^* U_j, \\ \partial^* U_i \cap \partial^* U_j \cap \partial^* U_k &= \emptyset \end{aligned} \tag{23}$$

holds for all i, j, k distinct.

To prove the claims, define for each i and $r > 0$ the sets

$$\begin{aligned} U_i^r(x) &= \{y \mid r(y - x) + x \in U_i\}, \\ H_i(x) &= \{y \mid \langle \eta_i, y - x \rangle \leq 0\}. \end{aligned}$$

Note that $\{U_i^r(x)\}_i$ are still disjoint. By Theorem 5.7.1 in [7]

$$1_{U_i^r(x)} \xrightarrow{r \rightarrow 0} 1_{H_i(x)} \quad \text{in } L^1_{\text{loc}}(\mathbb{R}^n) \text{ for all } x \in \partial^* U_i.$$

Therefore, if there existed $x \in \partial^* U_i \cap \partial^* U_j \cap \partial^* U_k$ for i, j, k distinct, then

$$1_{U_i^r(x) \cup U_j^r(x) \cup U_k^r(x)} \xrightarrow{r \rightarrow 0} 1_{H_i(x)} + 1_{H_j(x)} + 1_{H_k(x)} \quad \text{in } L^1_{\text{loc}}(\mathbb{R}^n), \tag{24}$$

which is impossible as the right hand side is not Lebesgue a.e. an indicator. By the same argument one can deduce that $\eta_i = -\eta_j$ must hold for $x \in \partial^* U_i \cap \partial^* U_j$ and that any $x \in \partial^* U_i$ cannot belong to the open set $(\bigcup_{l \neq i} \overline{U}_l)^c$. □

Proof of Theorem 2.4. For $i = 1, \dots, N$ Gauss–Green’s formula (see Theorem 5.8.1 in [7] and Theorem 4.5.6 in [9]) gives that

$$\int_{U_i} \operatorname{div}(f) \, dm = \int_{\partial U_i} \langle f, \eta_i \rangle \, d\mathcal{H}^{n-1} \quad (25)$$

for all Lipschitz continuous vector fields f with compact support. Here η_i denotes the outer unit normal of ∂U_i , which is well defined and nonzero on a subset of ∂U_i and zero everywhere else by definition.

Let $(g_r)_r$ be a sequence of smooth functions with

$$g_r(x) = \begin{cases} 1 & \text{if } x \in B(0, r), \\ 0 & \text{if } x \notin B(0, r+1), \end{cases}$$

and $(g_r)_r$ and $(Dg_r)_r$ uniformly bounded. Since $\hat{\mu}_i$ is Lipschitz continuous on $\overline{U}_i \cap B(0, r+1)$ Kirzbraun’s theorem ensures that $\hat{\mu}_i$ has a Lipschitz extension, $\hat{\mu}_i^r : \mathbb{R}^n \rightarrow \mathbb{R}^n$. Then $f_r = g_r \psi \hat{\mu}_i^r$ is Lipschitz continuous with compact support and $g_r \hat{\mu}_i^r = g_r \hat{\mu}_i$ on U_i . Then (25) applied to f_r yields

$$\int_{\partial U_i} g_r \psi \langle \hat{\mu}_i, \eta_i \rangle \, d\mathcal{H}^{n-1} = \int_{U_i} g_r \psi \operatorname{div}(\hat{\mu}_i) \, dm + \int_{U_i} \langle g_r D\psi + \psi Dg_r, \hat{\mu}_i \rangle \, dm.$$

Due to Assumption 2.2 all integrands above are dominated by integrable functions, and by letting $r \rightarrow \infty$ Lebesgue’s Dominated Convergence Theorem yields

$$\int_{\partial U_i} \psi \langle \hat{\mu}_i, \eta_i \rangle \, d\mathcal{H}^{n-1} = \int_{U_i} \psi \operatorname{div}(\hat{\mu}_i) \, dm + \int_{U_i} \langle D\psi, \hat{\mu}_i \rangle \, dm.$$

By summing over i we get

$$df(\hat{\mu}) = df_S(\hat{\mu}) - \sum_i \int_{\partial U_i} \psi \langle \hat{\mu}_i, \eta_i \rangle \, d\mathcal{H}^{n-1}. \quad (26)$$

By Lemma A.2 we see that

$$\begin{aligned} df(\hat{\mu}) &= df_S(\hat{\mu}) - \sum_{j \neq i} \int_{\partial U_i \cap \partial U_j} \psi \langle \hat{\mu}_i, \eta_i \rangle \, d\mathcal{H}^{n-1} \\ &= df_S(\hat{\mu}) + \frac{1}{2} \sum_{j \neq i} \int_{\partial U_i \cap \partial U_j} \langle \hat{\mu}_j - \hat{\mu}_i, \eta_i \rangle \psi \, d\mathcal{H}^{n-1}. \end{aligned}$$

Since η_i vanishes on $\partial U_i \cap \partial U_j \setminus (\overline{U}_i \cap \overline{U}_j)$ for $i \neq j$ we have proven (4). \square

References

- [1] L. Breiman. The little bootstrap and other methods for dimensionality selection in regression: X-fixed prediction error. *J. Amer. Statist. Assoc.* **87** (419) (1992) 738–754. [MR1185196](#)
- [2] P. Bühlmann and S. van de Geer. *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer Series in Statistics. Springer, Heidelberg, 2011. [MR2807761](#)
- [3] A. Dalalyan and A. Tsybakov. Aggregation by exponential weighting, sharp PAC-Bayesian bounds and sparsity. *Mach. Learn.* **72** (2008) 39–61.
- [4] D. L. Donoho and I. M. Johnstone. Adapting to unknown smoothness via wavelet shrinkage. *J. Amer. Statist. Assoc.* **90** (432) (1995) 1200–1224. [MR1379464](#)
- [5] B. Efron. The estimation of prediction error: Covariance penalties and cross-validation. *J. Amer. Statist. Assoc.* **99** (467) (2004) 619–632. [MR2090899](#)
- [6] B. Efron, T. Hastie, I. Johnstone and R. Tibshirani. Least angle regression. *Ann. Statist.* **32** (2) (2004) 407–499. With discussion, and a rejoinder by the authors. [MR2060166](#)

- [7] L. Evans and R. Gariepy. *Measure Theory and Fine Properties of Functions. Studies in Advanced Mathematics*. Taylor & Francis, London, 1992. [MR1158660](#)
- [8] J. Fan, L. Xue and H. Zou. Strong oracle optimality of folded concave penalized estimation. *Ann. Statist.* **42** (3) (2014) 819–849. [MR3210988](#)
- [9] H. Federer. *Geometric Measure Theory. Grundlehren der mathematischen Wissenschaften*. Springer, Berlin, 1969. [MR0257325](#)
- [10] J. Friedman, T. Hastie and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.* **33** (1) (2010) 1–22.
- [11] G. Givens and J. Hoeting. *Computational Statistics. Wiley Series in Computational Statistics*. John Wiley & Sons, Hoboken, 2012. [MR3236433](#)
- [12] N. R. Hansen and A. Sokol. Degrees of freedom for nonlinear least squares estimation, 2014. Available at <http://arxiv.org/abs/1402.2997>.
- [13] T. J. Hastie and R. J. Tibshirani. *Generalized Additive Models. Monographs on Statistics and Applied Probability*. **43**. Chapman and Hall Ltd., London, 1990. [MR1082147](#)
- [14] K. Kato. On the degrees of freedom in shrinkage estimation. *J. Multivariate Anal.* **100** (7) (2009) 1338–1352. [MR2514133](#)
- [15] J. D. Lee, D. L. Sun, Y. Sun and J. E. Taylor. Exact post-selection inference, with application to the lasso. *Ann. Statist.* **44** (3) (2016) 907–927. [MR3485948](#)
- [16] T. Loi and P. Phien. Bounds of Hausdorff measures of tame sets. *Acta Math. Vietnam.* **39** (4) (2014) 637–647. [MR3292588](#)
- [17] N. Meinshausen. Relaxed lasso. *Comput. Statist. Data Anal.* **52** (1) (2007) 374–393. [MR2409990](#)
- [18] M. Meyer and M. Woodroffe. On the degrees of freedom in shape-restricted regression. *Ann. Statist.* **28** (4) (2000) 1083–1104. [MR1810920](#)
- [19] F. Mikkelsen and N. Hansen. Supplementary material for “Degrees of freedom for piecewise Lipschitz estimators”, 2017. Available at <http://doi.org/10.5281/zenodo.321847>.
- [20] S. Reid, R. Tibshirani and J. Friedman. A study of error variance estimation in lasso regression. *Statist. Sinica* **26** (1) (2016) 35–67. [MR3468344](#)
- [21] C. M. Stein. Estimation of the mean of a multivariate normal distribution. *Ann. Statist.* **9** (6) (1981) 1135–1151. [MR0630098](#)
- [22] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc., B* **58** (1) (1996) 267–288. [MR1379242](#)
- [23] R. J. Tibshirani. The lasso problem and uniqueness. *Electron. J. Stat.* **7** (2013) 1456–1490. [MR3066375](#)
- [24] R. J. Tibshirani. Degrees of freedom and model search. *Statist. Sinica* **25** (3) (2015) 1265–1296. [MR3410308](#)
- [25] R. J. Tibshirani and J. Taylor. Degrees of freedom in lasso problems. *Ann. Statist.* **40** (2) (2012) 1198–1232. [MR2985948](#)
- [26] M. Vincent, K. Perell, F. Nielsen, G. Daugaard and N. Hansen. Modeling tissue contamination to improve molecular identification of the primary tumor site of metastases. *Bioinformatics* **30** (10) (2014) 1417–1423.
- [27] J. Ye. On measuring and correcting the effects of data mining and model selection. *J. Amer. Statist. Assoc.* **93** (441) (1998) 120–131. [MR1614596](#)
- [28] H. Zou, T. Hastie and R. Tibshirani. On the degrees of freedom of the lasso. *Ann. Statist.* **35** (5) (2007) 2173–2192. [MR2363967](#)