

PREDICTION ERROR AFTER MODEL SEARCH

BY XIAOYING TIAN

Department of Statistics, Stanford University, xtian@alumni.stanford.edu

Estimation of the prediction error of a linear estimation rule is difficult if the data analyst also uses data to select a set of variables and constructs the estimation rule using only the selected variables. In this work, we propose an asymptotically unbiased estimator for the prediction error after model search. Under some additional mild assumptions, we show that our estimator converges to the true prediction error in L^2 at the rate of $O(n^{-1/2})$, with n being the number of data points. Our estimator applies to general selection procedures, not requiring analytical forms for the selection. The number of variables to select from can grow as an exponential factor of n , allowing applications in high-dimensional data. It also allows model misspecifications, not requiring linear underlying models. One application of our method is that it provides an estimator for the degrees of freedom for many discontinuous estimation rules like best subset selection or relaxed Lasso. Connection to Stein's Unbiased Risk Estimator is discussed. We consider in-sample prediction errors in this work, with some extension to out-of-sample errors in low-dimensional, linear models. Examples such as best subset selection and relaxed Lasso are considered in simulations, where our estimator outperforms both C_p and cross validation in various settings.

1. Introduction. In this paper, we consider a homoscedastic model with Gaussian errors. In particular,

$$(1.1) \quad y = \mu(X) + \epsilon, \quad \epsilon \sim N(0, \sigma^2 I),$$

where the feature matrix $X \in \mathbb{R}^{n \times p}$ is considered fixed, $y \in \mathbb{R}^n$ is the response, and the noise level σ^2 is considered known and fixed. Note the mean function $\mu : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^n$ is not necessarily linear in X .

Prediction problems involve finding an estimator $\hat{\mu}$ which fits the data well. We naturally are interested in its performance in predicting a future response vector that is generated from the same mechanism as y . Mallows (1973) provided an unbiased estimator for the prediction error when the estimator is linear

$$\hat{\mu} = Hy,$$

where H is an $n \times n$ matrix independent of the data y . H is often referred to as the hat matrix. But in recent context, it is more and more unrealistic that the data analyst will not use the data to build a linear estimation rule. H , in other words, depends on y . In this case, is there still hope to get an unbiased estimator for the prediction error? In this article, we seek to address this problem.

The following are some examples of applications. In the context of model selection, the data analyst might use some techniques to select a subset of predictors M to build the linear estimation rules. Such techniques can include the more principled methods like LASSO (Tibshirani (1996)), best subset selection, forward stepwise regression and Least Angle Regression (Efron et al. (2004)) or some heuristics or even the combination of both. After the

Received November 2016; revised January 2019.

MSC2010 subject classifications. Primary 62H12, 62F12; secondary 62J07, 62F07.

Key words and phrases. Prediction error, model search, degrees of freedom, SURE.

selection step, we simply project the data onto the column space of X_M , the submatrix of X that consists of M columns, and use that as our estimation rule. Specifically,

$$(1.2) \quad \begin{aligned} \hat{\mu}(y; X) &= H_M \cdot y, & H_M &= P_M = X_M(X_M^T X_M)^{-1} X_M^T, \\ M &= \hat{M}(y), \end{aligned}$$

where \hat{M} can be any selection rule and P_M is the projection matrix onto the column space of X_M . In the case when M is selected by the LASSO, $\hat{\mu} = X_M \hat{\beta}_M(y)$, and $\hat{\beta}_M(y)$ is known as the *relaxed LASSO* solution (Meinshausen (2007)).

We assume the hat matrix $H_{\hat{M}}$ depends on the data y only through \hat{M} . In this sense, \hat{M} is the abstraction of the data-driven part in H . This paper will study the prediction error of

$$\hat{\mu} = H_{\hat{M}} \cdot y.$$

In this paper, we want to estimate the prediction error for $\hat{\mu}$,

$$(1.3) \quad \text{Err} = \mathbb{E}[\|y_{\text{new}} - H_{\hat{M}} \cdot y\|_2^2], \quad y_{\text{new}} \sim N(\mu(X), \sigma^2 I) \perp y.$$

There are several major methods for estimating (1.3) (Efron (2004)).

Penalty methods such as C_p or Akaike's information criterion (AIC) add a penalty term to the loss in training data. The penalty is usually twice the degrees of freedom times σ^2 .

Stein's Unbiased Risk Estimator (Stein (1981)) provides an unbiased estimator for any estimator that is smooth in the data. For nonsmooth estimation rules, Ye (1998) use perturbation techniques to approximate the covariance term for general estimators.

Nonparametric methods like cross validation or related bootstrap techniques provide risk estimators without any model assumption.

Methods like C_p assume a fixed model. Or specifically, the degrees of freedom is defined as $\text{df} = \text{tr}(H)$ for fixed H . Stein's Unbiased Risk Estimator (SURE) only allows risk estimation for almost differentiable estimators. In addition, computing the SURE estimate usually involves calculating the divergence of $\hat{\mu}(y)$. This is difficult when $\hat{\mu}(y)$ does not have an explicit form. Some special cases have been considered. Works by Tibshirani and Taylor (2012), Zou, Hastie and Tibshirani (2007) have computed the "degrees of freedom" for the LASSO estimator, which is Lipschitz. But for general estimators of the form $\hat{\mu} = H_{\hat{M}} y$, where $H_{\hat{M}}$ depends on y , $\hat{\mu}$ might not even be continuous in the data y . Thus, analytical forms of the prediction error are very difficult to derive (Mikkelsen and Hansen (2018), Tibshirani (2015)).

Nonparametric methods like cross validation are probably the most ubiquitous in practice. Cross validation has the advantage of assuming almost no model assumptions. However, Klement, Mamlouk and Martinetz (2008) shows that cross validation is inconsistent for estimating prediction error in high-dimensional scenarios. Moreover, cross validation also includes extra variation from having a different X for the validation set, which is different from the fixed X setup of this work. Efron (2004) also points out that the model-based methods like C_p , AIC, SURE offer substantially better accuracy compared with cross validation, given the model is believable.

In this work, we introduce a method for estimating prediction errors that is applicable to general model selection procedures. Examples include best subset selection for which prediction errors are difficult to estimate beyond X being orthogonal matrices (Tibshirani (2015)). In general, we do not require $H_{\hat{M}}$ to have any analytical forms. The main approach is to apply the selection algorithm \hat{M} to a slightly randomized response vector y^* . This is similar to holding out the validation set in cross validation, with the distinction that we do not have to split the feature matrix X . We can then construct an unbiased estimator for the

prediction error using the holdout information that is analogous to the validation set in cross validation. Note that since y^* would select a different model from y , this estimator will not be unbiased for the prediction error of $\hat{\mu}$. However, If the perturbation in y^* is small and we repeat this process multiple times so the randomization averages out, we will get an asymptotically unbiased and consistent estimator for the prediction error of $\hat{\mu} = H_{\hat{M}(y)}y$, which is the original target of our estimation.

Furthermore, we prove that under mild conditions on the selection procedure, our estimator converges to the true prediction error as in (1.3) in L^2 at the rate of $n^{-\frac{1}{2}}$. This automatically implies consistency of our estimator. Moreover, Li (1989) proves that in general, $n^{-\frac{1}{2}}$ is a lower bound for the estimation of the average squared error. Thus our estimator achieves the lower bound established by Li (1989). The C_p estimator, on the other hand, converges in L^2 at the rate of n^{-1} for fixed hat matrix H . So compared with C_p , our estimator pays a price of $n^{\frac{1}{2}}$ for the protection against any “data-driven” manipulations in choosing the hat matrix $H_{\hat{M}}$ for the linear estimation rules.

1.1. *Organization.* The rest of the paper is organized as follows. In Section 2, we introduce our procedure for unbiased estimation for a slightly different target. This is achieved by first randomizing the data and then constructing an unbiased estimator for the prediction error of this slightly different estimation rule. We then address the question of how randomization affects the accuracy of our estimator for estimating the true prediction error. There is a clear bias-variance tradeoff with respect to the amount of randomization. We derive upper bounds for the bias and the variance in Section 3 and propose an “optimal” scale of randomization that would make our estimator converge to the true prediction error in L^2 . Since the unbiased estimator constructed in Section 2 only uses one instance of randomization. We can further reduce the variance of our estimator by averaging over different randomizations. In Section 4, we propose a simple algorithm to compute the estimator after averaging over different randomizations. We also discuss the condition under which our estimator is equal to the SURE estimator. While SURE is difficult to compute both in terms of analytical formula and simulation, our estimator is easy to compute. Using the relationship between prediction error and degrees of freedom, we also discuss how to compute the “search degrees of freedom,” a term used in Tibshirani (2015) to refer to the degrees of freedom of estimators after model search. Finally, we include some simulation results in Section 5 and conclude with some discussions in Section 6.

2. Method of estimation. First, we assume the homoscedastic Gaussian model in (1.1), $y \sim N(\mu(X), \sigma^2 I)$, and we have a *model selection* algorithm \hat{M} ,

$$\hat{M} : \mathbb{R}^n \times \mathbb{R}^{n \times p} \rightarrow \mathcal{M}, \quad (y, X) \mapsto M.$$

As we assume X is fixed, we often use the shorthand $\hat{M}(y)$, and assume

$$\hat{M} : \mathbb{R}^n \rightarrow \mathcal{M}, \quad y \mapsto M,$$

where \mathcal{M} is a *finite* collection of models we are potentially interested in. The definition of models here is quite general. It can refer to any information we extract from the data. A common model as described in the Introduction can be a subset of predictors of particular interest. In such case, \hat{M} takes a value of the observation y and maps it to a set of selected variables. Note also the inverse image of \hat{M}^{-1} induces a partition on the space of \mathbb{R}^n . We will discuss this partition further in Section 3.

However, instead of using the original response variable y for selection, we use its randomized version y^* ,

$$(2.1) \quad y^* = y + \omega, \quad \omega \sim N(0, \alpha \sigma^2 I) \perp y.$$

For a fixed $\alpha > 0$, after using y^* to select a model M , we can define prediction errors analogous to that defined in (1.3),

$$(2.2) \quad \text{Err}_\alpha = \mathbb{E}[\|y_{\text{new}} - H_{\hat{M}(y+\omega)}y\|_2^2], \quad y_{\text{new}} \sim N(\mu(X), \sigma^2 I) \perp (y, \omega).$$

The subscript α denotes the amount of randomization added to y . Note that although randomization noise ω is added to selection, Err_α integrates over such randomization and thus are *not* random. The prediction error Err as defined in (1.3) corresponds to the case where we set $\alpha = 0$. In this section, we show that can get an unbiased estimator for Err_α for any $\alpha > 0$. Before we introduce the unbiased estimator, we first introduce some background on randomization.

2.1. Randomized selection. It might seem unusual to use y^* for model selection. But actually, using randomization for model selection and fitting is quite common—the common practice of splitting the data into a training set and a test set is a form of randomization. Although not stressed, the split is usually random and thus we are using a random subset of the data instead of the data itself for model selection and training.

The idea of randomization for model selection is not new. The field of differential privacy uses randomized data for database queries to preserve information (Dwork (2008)). This particular additive randomization scheme, $y^* = y + \omega$ is discussed in Tian and Taylor (2018). In this work, we discover that the additive randomization in (2.1) allows us to construct a vector independent of the model selection. This independent vector is analogous to the validation set in data splitting.

To address the question of the effect of randomization, we prove that Err and Err_α are close for small $\alpha > 0$ under mild conditions on the selection procedures. In other words, since we have an unbiased estimator for Err_α for any $\alpha > 0$, when α goes to 0, its bias for Err will diminish as well. For details, see Section 3. In addition, Section 5 also provides some evidence in simulations.

2.2. Unbiased estimation. To construct an unbiased estimator for Err_α , we first construct the following vector that is independent of y^* ,

$$(2.3) \quad y^- = y - \frac{1}{\alpha}\omega.$$

Note this construction is also mentioned in Tian and Taylor (2018). Using the property of Gaussian distributions and calculating the covariance between y^- and $y^* = y + \omega$, it is easy to see y^- is independent of y^* , and thus independent of the selection event $\{\hat{M}(y^*) = M\}$. Now we state our first result that constructs an unbiased estimator for Err_α for any $\alpha > 0$.

THEOREM 2.1 (Unbiased estimator). *Suppose $y \sim N(\mu(X), \sigma^2 I)$ is from the homoscedastic Gaussian model (1.1), then*

$$(2.4) \quad \widehat{\text{Err}}_\alpha = \|y^- - H_{\hat{M}(y^*)}y\|_2^2 + 2 \text{tr}(H_{\hat{M}(y^*)})\sigma^2 - \frac{1}{\alpha}n\sigma^2$$

is unbiased for Err_α for any $\alpha > 0$.

Before we prove the theorem, note that we need knowledge of σ^2 to compute the correction term $n\sigma^2/\alpha$ in $\widehat{\text{Err}}_\alpha$ (see (2.4)). In practice, we often need to estimate σ^2 . If we choose $\alpha = n^{-\frac{1}{4}}$ as suggested in Section 3.3, we need the variance estimate $\hat{\sigma}^2$ to be at least $n^{-\frac{1}{4}}$ consistent to get a consistent estimator for Err_α .

PROOF OF THEOREM 2.1. First notice

$$y = \frac{1}{1+\alpha}y^* + \frac{\alpha}{1+\alpha}y^-, \quad y = \mu(X) + \epsilon,$$

if we let $\epsilon^* = y^* - \mu(X)$ and $\epsilon^- = y^- - \mu(X)$, then

$$(2.5) \quad \epsilon = \frac{1}{1+\alpha}\epsilon^* + \frac{\alpha}{1+\alpha}\epsilon^-.$$

Note $\epsilon^* \perp \epsilon^-$ and $\epsilon^* \sim N(0, (1+\alpha)\sigma^2 I)$, and $\epsilon^- \sim N(0, \frac{1+\alpha}{\alpha}\sigma^2 I)$.

With this, we first define the following estimator for any $\alpha > 0$ and any $M \in \mathcal{M}$:

$$(2.6) \quad \widehat{\text{err}}_\alpha(M) = \|y^- - H_M y\|_2^2 + 2 \text{tr}(H_M)\sigma^2 - \frac{1}{\alpha}n\sigma^2.$$

We claim that $\widehat{\text{err}}_\alpha(M)$ is unbiased for the prediction error conditional on $\{\hat{M}(y^*) = M\}$ for any $M \in \mathcal{M}$ and any $\alpha > 0$. Formally, we prove

$$(2.7) \quad \mathbb{E}[\widehat{\text{err}}_\alpha(M) \mid \hat{M}(y^*) = M] = \mathbb{E}[\|y_{\text{new}} - H_M \cdot y\|^2 \mid \hat{M}(y^*) = M].$$

To see (2.7), we first rewrite

$$(2.8) \quad \begin{aligned} &\mathbb{E}[\|y_{\text{new}} - H_M y\|_2^2 \mid \hat{M}(y^*) = M] \\ &= \mathbb{E}[\|\mu - H_M y\|^2 \mid \hat{M}(y^*) = M] + n\sigma^2. \end{aligned}$$

Now we consider the conditional expectation of $\widehat{\text{err}}_\alpha(M)$. Note

$$\begin{aligned} &\mathbb{E}[\|y^- - H_M y\|_2^2 \mid \hat{M}(y^*) = M] \\ &= \mathbb{E}[\|\mu - H_M y\|^2 \mid \hat{M}(y^*) = M] \\ &\quad + \frac{1+\alpha}{\alpha}n\sigma^2 - 2\mathbb{E}[(\epsilon^-)^T H_M y \mid \hat{M}(y^*) = M] \\ &= \mathbb{E}[\|\mu - H_M y\|^2 \mid \hat{M}(y^*) = M] + \frac{1+\alpha}{\alpha}n\sigma^2 - \frac{2\alpha}{1+\alpha} \text{tr}[H_M \mathbb{E}[y^-(\epsilon^-)^T]] \\ &= \mathbb{E}[\|\mu - H_M y\|^2 \mid \hat{M}(y^*) = M] + \frac{1+\alpha}{\alpha}n\sigma^2 - 2 \text{tr}(H_M)\sigma^2. \end{aligned}$$

The equalities use the decomposition (2.5) as well as the fact that $y^* \perp \epsilon^-$.

Comparing this with (2.8), it is easy to see (2.7). Moreover, marginalizing over $\hat{M}(y^*)$, it is easy to see $\widehat{\text{Err}}_\alpha$ in (2.4) is unbiased for Err_α . \square

In fact, using the proof for Theorem 2.1, we have a even stronger result than the unbiasedness of $\widehat{\text{Err}}_\alpha$.

REMARK 2.2. $\widehat{\text{Err}}_\alpha$ is not only unbiased for the prediction error marginally, but conditional on any selected event $\{\hat{M}(y^*) = M\}$, $\widehat{\text{Err}}_\alpha$ is also unbiased for the prediction error. Formally,

$$\mathbb{E}[\widehat{\text{Err}}_\alpha \mid \hat{M}(y^*) = M] = \mathbb{E}[\|y_{\text{new}} - H_M \cdot y\|^2 \mid \hat{M}(y^*) = M].$$

This is easy to see with (2.7) and

$$\mathbb{E}[\widehat{\text{Err}}_\alpha \mid \hat{M}(y^*) = M] = \mathbb{E}[\widehat{\text{err}}_\alpha(M) \mid \hat{M}(y^*) = M].$$

The simple form of $\widehat{\text{Err}}_\alpha$ in (2.6) has some resemblance to the usual C_p formula for prediction error estimation, with $2 \text{tr}(H_{\hat{M}})\sigma^2$ being the usual correction term for degrees of freedom in the C_p estimator. The additional term $n\sigma^2/\alpha$ helps offset the larger variance in y^- .

3. Randomization and the bias-variance tradeoff. We investigate the effect of randomization in this section. In particular, we are interested in the bias term

$$B_n \stackrel{\text{def}}{=} \frac{\text{Err}_\alpha - \text{Err}}{n}$$

and the variance term

$$B'_n \stackrel{\text{def}}{=} \text{Var} \left[\frac{\widehat{\text{Err}}_\alpha}{n} \right]$$

for small $\alpha > 0$.

There is a simple intuition for the effects of randomization on estimation of prediction error. Since $y^* = y + \omega$, $\omega \sim N(0, \alpha\sigma^2)$, the randomized response vector y^* which we use for model selection will be close to y when the randomization scale α is small. Intuitively, Err_α should be closer to Err when α decreases. On the other hand, the independent vector $y^- = y - \omega/\alpha$ which we use to construct the estimator for the prediction error is more variant when α is small. We seek to find the optimal scale of α for this tradeoff.

Formally, we denote \hat{M} to be a selection procedure that selects an “important” subset of variables. That is

$$\hat{M} : \mathbb{R}^n \rightarrow \mathcal{M} \subseteq 2^{\{1, \dots, p\}}.$$

Without loss of generality, we assume \hat{M} is surjective. Thus, the number of potential models to choose from is $|\mathcal{M}|$ which is finite. Moreover, the map \hat{M} induces a partition of the space of \mathbb{R}^n . In particular, we assume

$$(3.1) \quad U_i = \hat{M}^{-1}(M_i) \subseteq \mathbb{R}^n, \quad i = 1, \dots, |\mathcal{M}|,$$

where $\mathcal{M} = \{M_1, \dots, M_{|\mathcal{M}|}\}$ are different models to choose from. It is easy to see that

$$\bigsqcup_{i=1}^{|\mathcal{M}|} U_i = \mathbb{R}^n,$$

and we further assume $\text{int}(U_i) \neq \emptyset$ and ∂U_i has measure 0 under the Lebesgue measure on \mathbb{R}^n .

Now we assume the hat matrix is a constant matrix in each of the partition U_i . In particular,

$$(3.2) \quad H_{\hat{M}(y)} = \sum_{i=1}^{|\mathcal{M}|} H_{M_i} \mathbb{I}(y \in U_i).$$

The most common matrix is probably the projection matrix onto the column space spanned by a subset of variables. Formally, we assume the following.

ASSUMPTION 3.1. For any $M \in \mathcal{M}$, we assume $H_M = X_M(X_M^T X_M)^{-1} X_M^T$, where X_M is the submatrix of X with M as the selected columns. It is easy to see, H_M is symmetric and

$$H_M^2 = H_M \quad \forall M \in \mathcal{M}.$$

Moreover, we also assume that \hat{M} does not select too many variables to include in a model. Specifically, we have the following.

ASSUMPTION 3.2. For any $M \in \mathcal{M}$,

$$\text{rank}(H_M) \leq |M| \leq K.$$

We will see in Section 3.2 that under Assumption 3.2, K and p enter the upper bound for $\text{Var}[\widehat{\text{Err}}_\alpha]$ as $O(\frac{(K \log p)^2}{n^2})$. Thus, so long as $K \log p = o(n)$, this term will vanish as n tends to infinity.

Finally, we also assume the model selection procedure \hat{M} and the resulting solution $H_{\hat{M}}y$ are stable. Specifically, we assume \hat{M} satisfies the following assumption.

ASSUMPTION 3.3.

$$(3.3) \quad \text{Var}[\|H_{\hat{M}(y')}y' - \mu\|_2^2] = O(n\sigma^4), \quad y' \sim N(\mu, \tau),$$

where $\sigma^2 \leq \tau \leq (1 + \delta)\sigma^2$ for some small constant $\delta > 0$.

Under the above assumption, the estimation error of $H_{\hat{M}}y$ cannot be too variable. As an example, we show in the following lemma that the relaxed Lasso estimator satisfies Assumption 3.3 under the setup commonly used in the Lasso literature.

LEMMA 3.4. *We assume a sparse linear model setting where*

$$y \sim N(X\beta^0, \sigma^2 I), \quad \|\beta^0\|_0 = s,$$

and suppose X satisfies the restricted eigenvalue condition with constant ϕ_0 as proposed in Bickel, Ritov and Tsybakov (2009). Then the relaxed Lasso estimator satisfies Assumption 3.3 if we choose the regularization parameter λ as

$$\lambda = \kappa\sigma\mathbb{E}[\|X^T\epsilon'\|_\infty], \quad \epsilon' \sim N(0, I),$$

for $\kappa > 1$ and $s \log p = O(\sqrt{n})$.

With these conditions above, we show in the following that the bias $B_n = O(\alpha)$ (Theorem 3.5) and the variance $B'_n = O(\frac{1}{n\alpha^2} + \frac{(K \log p)^2}{n^2})$ (Theorem 3.7). The proofs of the theorems use some lemmas whose proofs we defer to Section 7.

3.1. *Bias.* The bias B_n is introduced by the fact that selection is performed with y^* , the randomized version of y . However, for small perturbations, the resulting bias will be small as well. Formally, we have the following theorem.

THEOREM 3.5. *Suppose Assumptions 3.1, 3.2 and 3.3 are satisfied, then the bias*

$$B_n \leq C \cdot \alpha\sigma^2 \quad \text{for } \alpha < \delta,$$

where C is a universal constant and $\delta > 0$ is a small constant defined in Assumption 3.3.

Essential to the proof of Theorem 3.5 is that the estimation error for

$$\hat{\mu}(y) = H_{\hat{M}(y)}y$$

is robust to small perturbations on y . This is true under the assumptions introduced at the beginning of Section 3. Formally, we get the following lemma.

LEMMA 3.6. *Under Assumptions 3.1–3.3, we have*

$$\frac{1}{n}|\mathbb{E}[\|\hat{\mu}(y + \omega) - \mu\|_2^2] - \mathbb{E}[\|\hat{\mu}(y) - \mu\|_2^2]| \leq C_1 \cdot \alpha\sigma^2, \quad \omega \sim N(0, \alpha\sigma^2 I),$$

where C_1 is a universal constant. The inequality holds for all $\alpha < \delta$, where δ is the small constant defined in Assumption 3.3. The first expectation is taken over (y, ω) and the second expectation is taken over y .

With Lemma 3.6, it is easy to prove Theorem 3.5.

PROOF. First, notice that

$$H_{\hat{M}(y^*)}y = H_{\hat{M}(y^*)}y^* - H_{\hat{M}(y^*)}\omega = \hat{\mu}(y + \omega) - H_{\hat{M}(y+\omega)}\omega.$$

Thus, we have

$$\begin{aligned} B_n &= \frac{1}{n} |\text{Err}_\alpha - \text{Err}| \\ &= \frac{1}{n} |\mathbb{E}[\|H_{\hat{M}(y+\omega)}y - \mu\|^2] - \mathbb{E}[\|H_{\hat{M}(y)}y - \mu\|^2]| \\ &\leq \frac{1}{n} |\mathbb{E}[\|\hat{\mu}(y + \omega) - \mu\|^2] - \mathbb{E}[\|\hat{\mu}(y) - \mu\|^2]| + \frac{1}{n} \mathbb{E}[\|H_{\hat{M}(y+\omega)}\omega\|^2] + 2\alpha\sigma^2 \\ &\leq \frac{1}{n} |\mathbb{E}[\|\hat{\mu}(y + \omega) - \mu\|^2] - \mathbb{E}[\|\hat{\mu}(y) - \mu\|^2]| + \frac{1}{n} \mathbb{E}[\|\omega\|^2] + 2\alpha\sigma^2 \\ &\leq (C_1 + 3)\alpha\sigma^2. \end{aligned}$$

To see the first inequality, we use the notation as in (2.5) and note that the cross term

$$\begin{aligned} &\mathbb{E}[\omega^T H_{\hat{M}(y^*)}(H_{\hat{M}(y^*)}y^* - \mu)] \\ &= \mathbb{E}[\omega^T H_{\hat{M}(y^*)}\epsilon^*] \\ &= \mathbb{E}\left[\frac{\alpha}{1 + \alpha}(\epsilon^* - \epsilon^-)H_{\hat{M}(y^*)}\epsilon^*\right] \leq \frac{\alpha}{1 + \alpha} \mathbb{E}[\|\epsilon^*\|] = \alpha\sigma^2 \quad \square \end{aligned}$$

3.2. *Variance.* In this section, we discuss the variance of our estimator B'_n . As previously discussed at the beginning of the section, it is intuitive that B'_n will increase as α decreases. Before we establish quantitative results about B'_n with respect to α , recall that the variance of the C_p estimators is of order $O(n)$ when the hat matrix $H_{\hat{M}} = H$ is independent of data. In the following, we show that when we allow model selection, the variance of our estimator $\widehat{\text{Err}}_\alpha$ will increase as a function of α . When the optimal α is chosen, it converges to the average error at a rate of $O(n^{-\frac{1}{2}})$, which is the lower bound established in Li (1989).

In the following, we seek to establish how inflated the variance B'_n is compared to C_p . Theorem 3.7 gives an explicit upper bound for B'_n with respect to α , K and p .

THEOREM 3.7. *Suppose Assumptions 3.1, 3.2, 3.3 are satisfied and $\|\mu\|^2 = O(n)$, then*

$$B'_n = \sigma^4 \cdot O\left(\frac{1}{n\alpha^2} + \frac{(K \log p)^2}{n^2}\right).$$

PROOF. The key to bounding the variance of $\widehat{\text{Err}}_\alpha$ is to bound the variance

$$\text{Var}[\|y^- - H_{\hat{M}(y^*)}y\|^2].$$

First, noting $y^- = y - \frac{1}{\alpha}\omega$, $y^* = y + \omega$, we have

$$\begin{aligned} &\|y^- - H_{\hat{M}(y^*)}y\|^2 \\ &= \|y^- - H_{\hat{M}(y^*)}y^*\|^2 + \|H_{\hat{M}(y^*)}\omega\|^2 - 2\omega^T H_{\hat{M}(y^*)}(y^- - y^*) \\ &= \|y^- - H_{\hat{M}(y^*)}y^*\|^2 + \left(\frac{2}{\alpha} + 3\right)\|H_{\hat{M}(y^*)}\omega\|^2. \end{aligned}$$

The variance of the first term is bounded by

$$\begin{aligned} & \text{Var}[\|y^- - H_{\hat{M}(y^*)}y^*\|^2] \\ &= \mathbb{E}[\text{Var}[\|y^- - H_{\hat{M}(y^*)}y^*\|^2|y^*]] + \text{Var}[\mathbb{E}[\|y^- - H_{\hat{M}(y^*)}y^*\|^2|y^*]] \\ &= 2n\left(1 + \frac{1}{\alpha}\right)^2 \sigma^4 + 4\left(1 + \frac{1}{\alpha}\right)\sigma^2 \mathbb{E}[\|H_{\hat{M}(y^*)}y^* - \mu\|^2] \\ &\quad + \text{Var}[\|H_{\hat{M}(y^*)}y^* - \mu\|^2] \\ &\leq 2n\left(1 + \frac{1}{\alpha}\right)^2 \sigma^4 + 4\left(1 + \frac{1}{\alpha}\right)\sigma^2[(1 + \alpha)n + 2\|\mu\|^2] \\ &\quad + \text{Var}[\|H_{\hat{M}(y^*)}y^* - \mu\|^2]. \end{aligned}$$

The second equality uses the independence of y^* and y^- as well as the variance for a non-central χ^2 distribution.

The variance of the second term is bounded by

$$\begin{aligned} & \text{Var}\left[\left(\frac{2}{\alpha} + 3\right)\|H_{\hat{M}(y^*)}\omega\|^2\right] \\ &\leq \left(\frac{2}{\alpha} + 3\right)^2 \mathbb{E}[\|H_{\hat{M}(y^*)}\omega\|^2]^2 \\ &\leq \left(\frac{2}{\alpha} + 3\right)^2 \cdot \alpha^2 \sigma^4 \cdot O(K^2 + (\log |\mathcal{M}|)^2) = O((K \log p)^2 \sigma^4). \end{aligned}$$

The last inequality uses the tail property of χ^2 random variables, which are summarized in Lemma 7.1 in Section 7.

Combining the two terms, and using Assumption 3.3 and $\|\mu\|^2 = O(n)$, we have

$$\text{Var}\left[\frac{\widehat{\text{Err}}_\alpha}{n}\right] = \sigma^4 \cdot O\left(\frac{1}{n\alpha^2} + \frac{(K \log p)^2}{n^2}\right). \quad \square$$

Recall the variance of $\frac{C_p}{n}$ is of order $O(n^{-1})$. In comparison, we pay a price of α^{-2} plus $O\left(\frac{(K \log p)^2}{n^2}\right)$, but allows our hat matrix to be dependent on the data y .

3.3. *Bias-variance tradeoff and the choice of α .* Combining Theorem 3.5 and Theorem 3.7, we have the following convergence rate:

$$\mathbb{E}\left[\left(\frac{\widehat{\text{Err}}_\alpha}{n} - \frac{\text{Err}}{n}\right)^2\right] = \sigma^4 \cdot O\left(\alpha^2 + \frac{1}{n\alpha^2} + \frac{(K \log p)^2}{n^2}\right).$$

COROLLARY 3.8. *If we choose $\alpha = n^{-\frac{1}{4}}$ and assume $K \log p = O(n^{\frac{3}{4}})$, then*

$$\mathbb{E}\left[\left(\frac{\widehat{\text{Err}}_\alpha}{n} - \frac{\text{Err}}{n}\right)^2\right] = O(n^{-\frac{1}{2}}).$$

It is easy to see that $\alpha = n^{-\frac{1}{4}}$ achieves the optimal rate for convergence. This should offer some guidance about the choice of α in practice.

Algorithm 1 Algorithm for computing $\widehat{\text{Err}}_\alpha^{(I)}$ for any $\alpha > 0$

```

1: Input:  $X, y$ 
2: Initialize:  $\widehat{\text{Err}}_\alpha^{(I)} \leftarrow 0, N \in \mathbb{Z}_+$ 
3: for  $i$  in  $1:N$  do
4:   Draw  $\omega^{(i)} \sim N(0, \alpha\sigma^2 I)$ 
5:   Compute  $y^* = y + \omega^{(i)}, y^- = y - \frac{1}{\alpha}\omega^{(i)}$ 
6:   Compute  $\hat{M}^* = \hat{M}(y^*)$ 
7:   Use Equation (2.4) to compute  $\widehat{\text{Err}}_\alpha^{(i)}$  from  $y, y^-, \hat{M}^*$ .
8:    $\widehat{\text{Err}}_\alpha^{(I)} = \widehat{\text{Err}}_\alpha^{(I)} + \widehat{\text{Err}}_\alpha^{(i)} / N$ 
return  $\widehat{\text{Err}}_\alpha^{(I)}$ 

```

4. Further properties and applications. In Section 3.3, we show that $\widehat{\text{Err}}_\alpha$ will have diminishing variances if α is chosen properly. However, since $\widehat{\text{Err}}_\alpha$ is computed using only one instance of the randomization ω , its variance can be further reduced if we aggregate over different randomizations ω . Furthermore, in the following section, we will show that after such marginalization over ω , $\widehat{\text{Err}}_\alpha$ is uniform minimum variance unbiased (UMVU) estimators for the prediction error Err_α under some conditions.

4.1. *Variance reduction techniques and UMVU estimators.* We first introduce the following lemma that shows the variance of $\widehat{\text{Err}}_\alpha$ can be reduced at no further assumption.

LEMMA 4.1. *The following estimator is unbiased for Err_α ,*

$$(4.1) \quad \widehat{\text{Err}}_\alpha^{(I)} = \mathbb{E}_\omega[\widehat{\text{Err}}_\alpha \mid y].$$

Furthermore, it has smaller variance,

$$\text{Var}[\widehat{\text{Err}}_\alpha^{(I)}] \leq \text{Var}[\widehat{\text{Err}}_\alpha].$$

The lemma can be easily proved using basic properties of conditional expectation. In practice, we approximate the integration over ω by repeatedly sampling ω and taking the averages. Specifically, Algorithm 1 provides an algorithm for computing $\widehat{\text{Err}}_\alpha^{(I)}$ for any $\alpha > 0$.

Since $\widehat{\text{Err}}_\alpha^{(I)}$ has the same expectation as $\widehat{\text{Err}}_\alpha$ with smaller variances, it is easy to deduce from Corollary 3.8 that $\widehat{\text{Err}}_\alpha^{(I)}$ also converges to Err in L^2 at a rate of at least $O(n^{-\frac{1}{2}})$ (after a proper scaling of n^{-1}). Furthermore, we show that such estimators are UMVU estimators for any $\alpha > 0$ when the parameter space $\mu(X)$ contains a ball in \mathbb{R}^n .

LEMMA 4.2. *If parameter space of $\mu(X)$ contains a ball in \mathbb{R}^n , then $\widehat{\text{Err}}_\alpha^{(I)}$ are UMVU estimators for Err_α for any $\alpha > 0$.*

PROOF. Without loss of generality, assume ω has density g with respect to the Lebesgue measure on \mathbb{R}^n , then the density of (y, ω) with respect to the Lebesgue measure on $\mathbb{R}^n \times \mathbb{R}^n$ is proportional to

$$(4.2) \quad \exp\left[\frac{\mu(X)^T y}{\sigma^2}\right] g(\omega).$$

We note that (4.2) is an exponential family with sufficient statistics y . Moreover, when the parameter space of $\mu(X)$ contains a ball in \mathbb{R}^n , then we have y is sufficient and complete. Thus, taking an unbiased estimator $\widehat{\text{Err}}_\alpha$ and integrating over ω conditional on y , the complete and sufficient statistics, we have the UMVU estimators. \square

4.2. *Relation to the SURE estimator.* In this section, we reveal that our estimator $\widehat{\text{Err}}_\alpha$ is equal to the SURE estimator for the prediction error Err_α if the parameter space of $\mu(X)$ contains a ball in \mathbb{R}^n .

First, we notice that for any $\alpha > 0$, Err_α is the prediction error for

$$\hat{\mu}_\alpha(y) = \mathbb{E}[H_{\hat{M}(y+\omega)}y \mid y], \quad \omega \sim N(0, \alpha\sigma^2I).$$

Although $\hat{\mu}(y)$ might be discontinuous in y , $\hat{\mu}_\alpha(y)$ is actually smooth in the data. To see that, note

$$(4.3) \quad \hat{\mu}_\alpha(y) = \sum_{i=1}^{|\mathcal{M}|} H_i y \int_{U_i} \phi_\alpha(z + y) dz,$$

where ϕ_α is the p.d.f for $N(0, \alpha\sigma^2I)$. Due to the smoothness of ϕ_α and the summation being a finite sum, we have $\hat{\mu}_\alpha(y)$ is smooth in y . Therefore, in theory we can use Stein’s formula to compute an estimate for the prediction error of $\hat{\mu}_\alpha(y)$. Note such estimator would only depend on y , the complete and sufficient statistics for the exponential family in (4.2) when the parameter space of $\mu(X)$ contains a ball in \mathbb{R}^n . Thus it is also the UMVU estimator for Err_α . By Lemma 4.2 and the uniqueness of UMVU estimators, we conclude $\widehat{\text{Err}}_\alpha$ is the same as the SURE estimator.

However, the SURE estimator is quite difficult to compute as the regions U_i ’s may have complex geometry and explicit formulas are hard to derive (Mikkelsen and Hansen (2018)). Moreover, it is difficult to even use Monte Carlo samplers to approximate the integrals in (4.3) since the sets U_i ’s might be hard to describe and there are $|\mathcal{M}|$ integrals to evaluate, making it computationally expensive.

In contrast, $\widehat{\text{Err}}_\alpha$ provides an unbiased estimator for Err_α at a much lower computational cost. That is, we only need to sample ω ’s from $N(0, \alpha\sigma^2I)$ and compute $\widehat{\text{Err}}_\alpha$ at each time and average over them. The major computation involved is reselecting the model with $y^* = y + \omega$. In practice, we choose the number of samples for ω ’s to be less than the number of data points, so the computation involved will be even less than Leave-One-Out cross validation.

4.3. *Prediction error after model selection.* One key message of this work is that we can estimate the prediction error of the estimation rule $\hat{\mu}$ even if we have used some model selection procedure to construct the hat matrix $H_{\hat{M}}$ in $\hat{\mu}$. In practice, however, we need a priori information on σ^2 to compute $\widehat{\text{Err}}_\alpha$. There are several methods for consistent estimation of σ^2 . In the low dimensional setting, we can simply use the residual sum of squares divided by the degrees of freedom to estimate σ^2 . In the high-dimensional setting, the problem is more challenging, but various methods are derived (Reid, Tibshirani and Friedman (2016), Sun and Zhang (2012), Tian, Loftus and Taylor (2018)).

We also want to stress that the prediction error defined in this work is the in-sample prediction error that assumes fixed X . This is the same setup as in C_p (Mallows (1973)), SURE (Stein (1981)) and the prediction errors discussed in Efron (2004). A good estimator for the in-sample prediction error will allow us to evaluate and compare the predictive power of different estimation rules.

However, in other cases, we might be interested in out-of-sample prediction errors. That is, the prediction errors are measured on a new dataset $(X_{\text{new}}, y_{\text{new}})$, $X_{\text{new}} \in \mathbb{R}^{n \times p}$, $y_{\text{new}} \in \mathbb{R}^n$ where $X_{\text{new}} \neq X$. In this case, assuming we observe some new feature matrix X_{new} , and we are interested in the out-of-sample prediction error,

$$(4.4) \quad \text{Err}_{\text{out}} = \mathbb{E}[\| \mu(X_{\text{new}}) - X_{\text{new}}\bar{\beta}(y) \|^2] + n\sigma^2,$$

where

$$\bar{\beta}(y) = (X_M^T X_M)^{-1} X_M^T y, \quad \hat{M}(y) = M,$$

where \hat{M} is the model selection procedure that depends on the data. Analogous to Err_α , we define

$$(4.5) \quad \text{Err}_{\text{out},\alpha} = \mathbb{E}[\|\mu(X_{\text{new}}) - X_{\text{new}}\bar{\beta}^*(y)\|^2] + n\sigma^2,$$

where

$$\bar{\beta}^*(y) = (X_M^T X_M)^{-1} X_M^T y, \quad \hat{M}(y^*) = M.$$

We want to point out that we do not place any assumption on how the feature matrix is sampled. Specifically, we do not need to assume X_{new} is sampled from the same distribution as X . Rather, we condition on the newly observed matrix X_{new} . This is a distinction from cross validation which assumes the rows of the feature matrix X are *i.i.d* samples from some distribution. Such assumption may not be satisfied in practice.

Then in the low dimensional setting where $p < n$, we are able to construct an unbiased estimator for $\text{Err}_{\text{out},\alpha}$.

LEMMA 4.3. *Suppose $X \in \mathbb{R}^{n \times p}$ and $\text{rank}(X) = p$. Then if we further assume a linear model where*

$$\mu(X) = X\beta^0,$$

where β^0 is the underlying coefficients. Assuming the homoscedastic model in (1.1), we have

$$(4.6) \quad \begin{aligned} \widehat{\text{Err}}_{\text{out},\alpha} &= \|H_0 y^- - H_{\hat{M}(y^*)} y\|^2 + 2 \text{tr}(H_0^T H_{\hat{M}(y^*)})\sigma^2 + n\sigma^2 \\ &\quad - 2 \text{tr}(H_0^T H_0) \left(1 + \frac{1}{\alpha}\right)\sigma^2 \end{aligned}$$

is unbiased for $\text{Err}_{\text{out},\alpha}$, where

$$H_0 = X_{\text{new}}(X^T X)^{-1} X^T, \quad H_{\hat{M}(y^*)} = X_{\text{new},\hat{M}(y^*)}(X_{\hat{M}(y^*)}^T X_{\hat{M}(y^*)})^{-1} X_{\hat{M}(y^*)}^T.$$

The proof of the lemma is analogous to that of Theorem 2.1 noticing that

$$H_0 \mu(X) = X_{\text{new}}(X^T X)^{-1} X^T X \beta^0 = X_{\text{new}} \beta^0 = \mu(X_{\text{new}}).$$

Lemma 4.3 provides an unbiased estimator for $\text{Err}_{\text{out},\alpha}$ for $p < n$ and $\mu(X)$ being a linear function of X . To bound the difference $\text{Err}_{\text{out},\alpha} - \text{Err}_{\text{out}}$, we might need to assume conditions similar to those introduced at the beginning of Section 3. In the case where $p < n$, we might still hope that the matrices $H_0, H_{\hat{M}}$ will be close to projection matrices, and almost satisfy Assumptions 3.1. Thus, intuitively, $\text{Err}_{\text{out},\alpha}$ and Err_{out} will be close and the estimator $\widehat{\text{Err}}_{\text{out},\alpha}$ will be a good estimator of Err_{out} when $p < n$. In simulations, we see that in the low-dimensional setting, the performance of $\text{Err}_{\text{out},\alpha}$ is comparable to that of cross validation. However, in the high-dimensional setting where $n < p$, the estimation of out-of-sample errors remains a very challenging problem that we do not seek to address in the scope of this work.

4.4. *Search degrees of freedom.* There is a close relationship between (in-sample) prediction error and the degrees of freedom of an estimator. In fact, with a consistent estimator for the prediction error Err , we get a consistent estimator for the degrees of freedom.

Under the framework of Stein’s unbiased risk estimator, for any estimation rule $\hat{\mu}$, we have

$$(4.7) \quad \text{Err} = \mathbb{E}[\|y - \hat{\mu}(y)\|^2] + 2 \sum_{i=1}^n \text{Cov}[\hat{\mu}_i(y), y_i],$$

where $\hat{\mu}_i$ is the i th coordinate of $\hat{\mu}$. For almost differentiable $\hat{\mu}$'s, [Stein \(1981\)](#) showed the covariance term is equal to

$$(4.8) \quad \text{Cov}[\hat{\mu}_i(y), y_i] = \sigma^2 \mathbb{E} \left[\frac{\partial \hat{\mu}_i}{\partial y_i} \right].$$

The sum of the covariance terms, properly scaled, is also called the degrees of freedom:

$$(4.9) \quad \text{df} = \sigma^{-2} \sum_{i=1}^n \text{Cov}[\hat{\mu}_i(y), y_i] = \sum_{i=1}^n \mathbb{E} \left[\frac{\partial \hat{\mu}_i}{\partial y_i} \right].$$

However, in many cases, the analytical forms of $\hat{\mu}$ are very hard to compute or there is none. In such cases, the computation of its divergence is only feasible for very special $\hat{\mu}$'s ([Zou, Hastie and Tibshirani \(2007\)](#)). Moreover, for discontinuous $\hat{\mu}$'s which are under consideration in this work, [Mikkelsen and Hansen \(2018\)](#) showed that there are further correction terms for (4.8) to account for the discontinuities. In general, these correction terms do not have analytical forms and are hard to compute. Intuitively, due to the search involved in constructing $\hat{\mu} = H_{\hat{M}}y$, it will have larger degrees of freedom than $\text{tr}(H_{\hat{M}})$ which treats the hat matrix as fixed. We adopt the name used in [Tibshirani \(2015\)](#) to call it ‘‘search degrees of freedom.’’

We circumvent the difficulty in computing $\partial \hat{\mu}_i / \partial y_i$ by providing an asymptotically unbiased estimator for Err . Formally,

$$(4.10) \quad \hat{\text{df}} = \frac{1}{\sigma^2} [\widehat{\text{Err}}_\alpha^{(I)} - \|y - \hat{\mu}\|_2^2],$$

where $\widehat{\text{Err}}_\alpha^{(I)}$ is defined as in (4.1). Using the discussion in Section 3.3, we choose $\alpha = n^{-1/4}$. Notice that such approach as above is not specific to any particular model search procedures involved in constructing $\hat{\mu}$. Thus it offers a unified approach to compute degrees of freedom for any $\hat{\mu} = H_{\hat{M}(y)}y$ satisfying the appropriate assumptions in Section 3. We illustrate this flexibility by computing the search degrees of freedom for the best subset selection where there has been no explicitly computable formula.

Prediction error estimates may also be used for tuning parameters. For example, if the model selection procedure \hat{M} is associated with some regularization parameter λ , we find the optimal λ that minimizes the prediction error of $\hat{\mu}_\lambda$

$$(4.11) \quad \lambda_{\text{optimal}} = \min_{\lambda} \mathbb{E} [\|y_{\text{new}} - H_{\hat{M}_\lambda(y)} \cdot y\|_2^2],$$

where the expectation is taken over both y_{new} and y . [Shen and Ye \(2002\)](#) shows that this model tuning criterion will yield an adaptively optimal model which achieves the optimal prediction error as if the tuning parameter were given in advance.

Using the relationship in (4.7) and (4.9), we easily see that the C_p type criterion (4.11) is equivalent to the AIC criterion using the definition of degrees of freedom (4.9). Analogously, we can also propose the BIC criterion as

$$\text{BIC} = \frac{\|y - \hat{\mu}\|_2^2}{n\sigma^2} + \frac{\log n}{n} \hat{\text{df}}.$$

[Yang \(2005\)](#) points out that compared with the C_p or AIC criterion, BIC tends to recover the true underlying sparse model and recommends it if sparsity is the major concern.

5. Simulations. In this work, we propose a method for risk estimation for a class of ‘‘select and estimate’’ estimators. One remarkable feature of our method is that it provides a consistent estimator of the prediction error for a large class of selection procedures under

general, mild conditions. To demonstrate this strength, we provide simulations for two selection procedure under various setups and datasets. The two estimators are the OLS estimator after best subset selection and relaxed Lasso, which we denote as $\hat{\mu}_{\text{best}}$ and $\hat{\mu}_{\text{RL}}$. In particular,

$$\hat{\mu}(y) = X_{\hat{M}}(X_{\hat{M}}^T X_{\hat{M}})^{-1} X_{\hat{M}}^T y,$$

where \hat{M} is selected by the best subset selection and Lasso at a fixed λ respectively using the original data y . In their Lagrangian forms, best subset selection and Lasso at fixed λ can be written as

$$\min_{\beta} \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_k,$$

where $k = 0$ for best subset selection and $k = 1$ for Lasso. Thus, by showing the good performances (in simulation) of our estimator at both $k = 0$ and $k = 1$, we believe the good performance would persist for all the nonconvex optimization problems with $0 \leq k < 1$. In the simulation, we always marginalize over different randomizations to reduce variance. Specifically, we use Algorithm 1 to compute $\widehat{\text{Err}}_{\alpha}^{(I)}$ which we use in all of the comparisons below.

In the following simulations, we compare both the bias and variances of our estimator $\widehat{\text{Err}}_{\alpha}^{(I)}$ with the C_p estimator, cross validation as well as the parametric bootstrap method proposed in Efron (2004). In particular, to ensure fairness of comparison, we use Leave-One-Out cross validation in all of our simulations. Most of our simulations are for in-sample prediction errors with some exceptions of comparing the out-of-sample estimator $\widehat{\text{Err}}_{\text{out},\alpha}^{(I)}$ in Section 4.3 to cross validation for estimating out-of-sample prediction errors. To establish a “known” truth to compare to, we use mostly synthetic data, with some of the synthetic datasets generated from a diabetes dataset. In the following simulations, we call our estimator “additive” due to the additive randomization used in the estimation. Cross validation is abbreviated as “CV.” The true prediction error is evaluated through Monte Carlo sampling since we have access to the “true” underlying distribution. We assume the variance σ^2 is unknown and estimate it with the OLS residuals when $p < n$. In the high-dimensional setting, we use the methods in Reid, Tibshirani and Friedman (2016) to estimate σ^2 .

5.1. Relaxed Lasso estimator. We perform simulation studies for the prediction error and degrees of freedom estimation for the relaxed Lasso estimator. Unless stated otherwise, the target of prediction error estimation is the in-sample prediction error

$$\text{Err} = \mathbb{E}[\|y_{\text{new}} - \hat{\mu}_{\text{RL}}(y)\|_2^2], \quad y_{\text{new}} \sim N(\mu(X), \sigma^2 I) \perp y.$$

According to the framework of SURE (Stein (1981)), the degrees of freedom of the estimator $\hat{\mu}_{\text{RL}}$ can be defined as

$$\text{df} = \sum_{i=1}^n \frac{\text{Cov}[\hat{\mu}_{\text{RL},i}, y_i]}{\sigma^2},$$

which is the target of our estimation. We first study the performance of the prediction error estimation.

5.1.1. Prediction error estimation. In the following, we describe our data generating distribution as well as the parameters used in the simulation.

- The feature matrix $X \in \mathbb{R}^{n \times p}$ is simulated from an equi-correlated covariance matrix with normal entries. The correlation is $\rho = 0.3$.

- y is generated from a sparse linear model,

$$y = X\beta^0 + \epsilon, \quad \epsilon \sim N(0, \sigma^2 I),$$

where

$$\beta^0 = (\underbrace{\text{snr}, \dots, \text{snr}}_s, 0, \dots, 0)$$

and snr is the signal-to-noise ratio and s is the sparsity of β^0 .

- We fit a Lasso problem with $\lambda = \kappa \lambda_0$, where

$$\lambda_{\min} = \mathbb{E}[\|X^T \epsilon'\|_\infty], \quad \epsilon' \sim N(0, \sigma^2 I),$$

is the level where noise below which noise starts to enter the Lasso path (Negahban et al. (2009)) and we choose $\kappa > 1$.

- The parameter α as defined in (2.1) is taken to be approximately $n^{-\frac{1}{4}}$.

We compare the performances of the estimators for different settings. We take $n = 100$, and $p = 50, 200, 400$ and sparsity to be $s = 10, 20$. Since $n^{-\frac{1}{2}} = 10$, $s = 20$ is the more dense signal situation. We take κ to be 1.1 for the low-dimensional setting and 1.5 for the high-dimensional setting. The randomization parameter $\alpha = 0.25 \approx n^{-1/4}$. We see in Figure 1 that in all settings $\widehat{\text{Err}}_\alpha^{(J)}$ provides an unbiased estimator that has small variance. Remarkably, notice that the variance of our estimator is comparable to the dotted the black lines are the standard error of the true prediction error estimated from Monte Carlo sampling, which is probably the best one can hope for. $\widehat{\text{Err}}_\alpha^{(J)}$ clearly outperforms both C_p and cross validation. Its performance is comparable to the parametric bootstrap estimator in the sparse scenario although parametric bootstrap seems to have more extreme values. Our estimator also performs slightly better in the more dense scenario $s = 20$ in panel 3 of Figure 1. In the dense signal situation, the model selected by Lasso is often misspecified. We suspect that in this situation, that parametric bootstrap overfits the data in this situation, causing a slight bias downwards. The C_p estimator is always biased down because it does not take into account the “degrees of freedom” used for model search. On the other hand, cross validation has an upward bias for in-sample prediction error. However, this bias is twofold. First, the extra randomness in the new feature matrix will cause the out-of-sample prediction error to be higher. However, comparing panels 3 and 4 of Figure 1, we see that when the signal is more dense $s = 20$ in panel 3, cross validation has a much larger bias than when the dimension is higher $p = 400$ in panel 4. This suggests that cross validation might be susceptible to model misspecifications as well. With less sparse signals, the model selected by Lasso is not stable or consistent, causing cross validation to behave wildly even when we only leave out one observation at a time. In contrast, in all of the four settings, our estimator $\widehat{\text{Err}}_\alpha^{(J)}$ provides an unbiased estimator with small variance.

This phenomenon persists when we vary the penalty parameter λ . For a grid of λ 's with varying κ 's from $[0.2, 1.6]$, we see from Figure 2 that cross validation error is always overestimates the in-sample prediction error. Moreover, the amount of over estimation highly depends on the data generating distribution. In both panels of Figure 2, $n = 100$, $p = 200$, $\text{snr} = 7$, and the only difference is the sparsity is $s = 10$ for Figure 2(a) and $s = 20$ for Figure 2(b). Using the same dimensions for X , we seek to control the extra randomness by using a different X for the validation set. However, the change in the sparsity level alone has huge impact for the cross validation estimates of the prediction error. The curve by cross validation is also more kinky due to its bigger variance. However, in both scenarios, $\widehat{\text{Err}}_\alpha^{(J)}$ hugs the true prediction error.

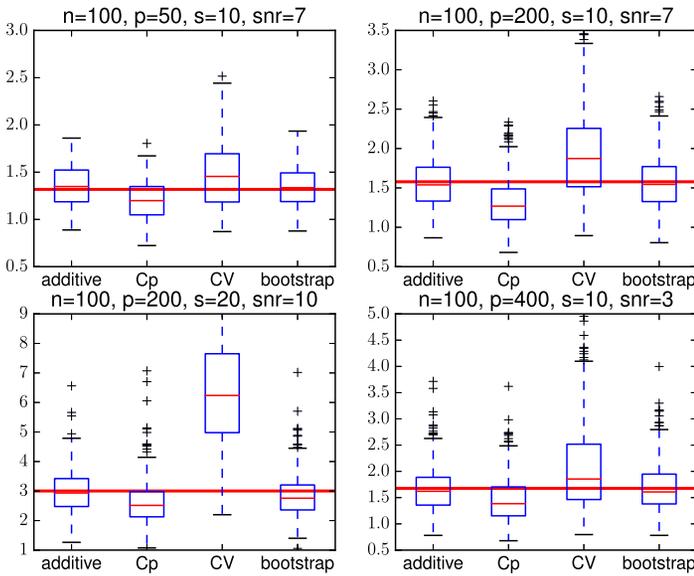


FIG. 1. Comparison of different estimators for different n, p, s, snr . The red horizontal line is the true prediction error estimated by Monte Carlo simulation with the dashed black lines denoting its standard deviation.

5.1.2. *Degrees of freedom.* In this section, we carry out a simulation study for our estimate of the degrees of freedom of the relaxed Lasso estimator $\hat{\mu}_{\text{RL}}$. We take the 64 predictors in the diabetes dataset (Efron et al. (2004)) to be our feature matrix X , which include the interaction terms of the original ten predictors. The positive cone condition is violated on the 64 predictors (Efron et al. (2004), Zou, Hastie and Tibshirani (2007)). We use the response vectors y to compute the OLS estimator $\hat{\beta}_{\text{ols}}$ and $\hat{\sigma}_{\text{ols}}$ and then synthetic data is generated through

$$y = X\hat{\beta}_{\text{ols}} + \epsilon, \quad \epsilon \sim N(0, \sigma_{\text{ols}}^2 I).$$

We choose λ 's to have different ratios $\kappa \in \{0.05, 0.1, 0.15, 0.2, 0.25\}$. Figure 3 shows the estimates of degrees of freedoms by our method as in (4.10) and the naive estimate $\hat{d}_{\text{naive}} = |\hat{M}|$ compared with the truth computed by Monte Carlo sampling. The naive C_p estimator always underestimate the degrees of freedom, not taking into account the inflation

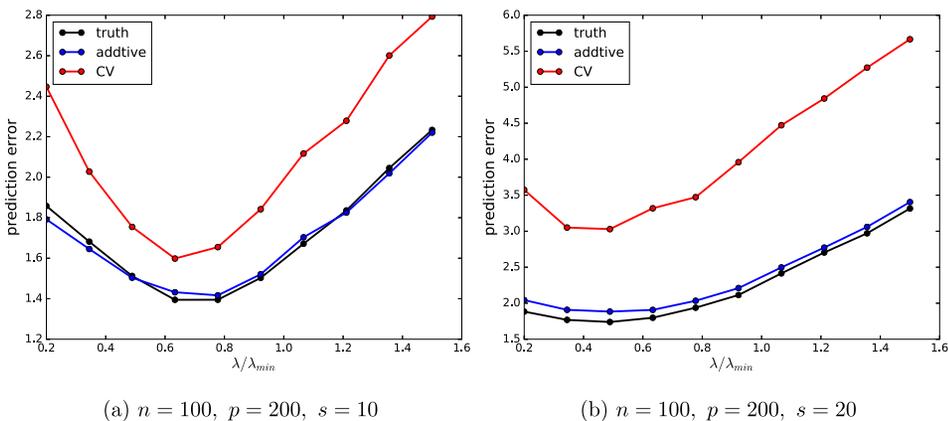


FIG. 2. Estimation of prediction errors for different λ 's. Cross validation is always biased upwards. However, the bias depends on the data generating distribution.

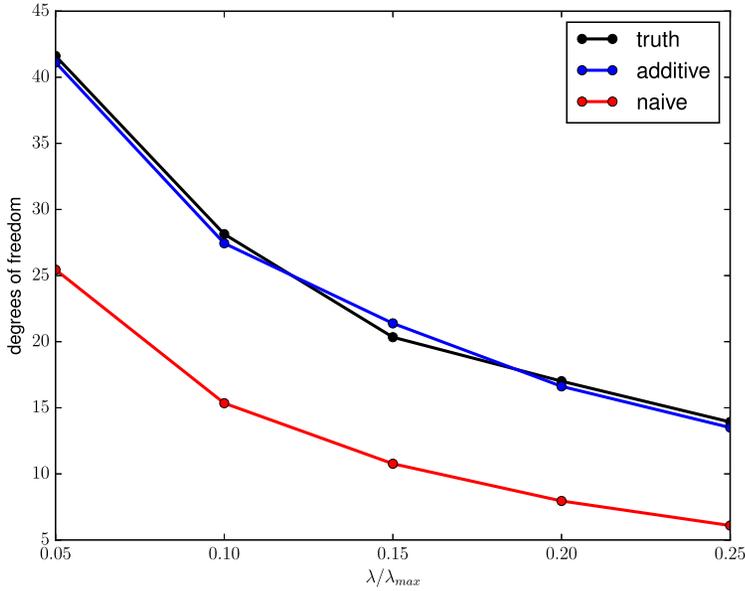


FIG. 3. Comparison of estimates of degrees of freedom by cross validation, \hat{df} in (4.10) and $\hat{df}_{naive} = |\hat{M}|$ at different λ 's. $\alpha = 0.25 \approx n^{-1/4}$.

in degrees of freedom after model search. However, our estimator as defined in (4.9) provides an unbiased estimation for the true degrees of freedom for the relaxed Lasso estimator $\hat{\mu}_{RL}$.

5.1.3. *Out-of-sample prediction errors.* Finally, we test the unbiasedness of the proposed estimator in Section 4.3 for out-of-sample prediction error. We compare with cross validation in the low-dimensional setting where $p = 20$ and $p = 50$, respectively. In this section only, our target is the out-of-sample prediction error

$$Err_{out} = \mathbb{E}[\|X_{new}\beta^0 - X_{new,\hat{M}(y)}\bar{\beta}(y)\|^2] + n\sigma^2,$$

where $\bar{\beta}$ is the relaxed Lasso estimator and $\hat{M}(y)$ is the nonzero set of the Lasso solution at λ . We still abbreviate our estimator as “additive” and compare with the out-of-sample prediction error by cross validation.

We see in Figure 4 that the estimator proposed in Section 4.3 is roughly unbiased for out-of-sample prediction error. Its performance is comparable with cross validation in both settings, with a slightly larger variance. However, as pointed in Section 4.3, our estimator does not assume any assumptions on the underlying distribution of the feature matrix X .

5.2. *Best subset selection.* The C_p estimator was originally proposed for picking the model size in best subset selection. One aspect that often gets neglected is that for any $k < p$, where p is the number of features to choose from, there are more than one models of size k to choose from. And the best subset of size k already includes a selection procedure that needs to be adjusted for. To illustrate this problem, we generate a feature matrix X of dimension 100×6 with i.i.d. standard normal entries, and y is generated from a linear model of X :

$$y = X\beta + N(0, 1), \quad \beta = (1, 2, 3, 4, 5, 6).$$

For each subset of size $k = 1, \dots, 6$, we estimate the prediction error of the best subset of size k using both C_p and $\widehat{Err}_\alpha^{(I)}$. The true prediction error is evaluated using Monte Carlo sampling. From Figure 5, we see that C_p is indeed an under estimate for the prediction error for best subset selection. The bias is bigger when $k = 2, 3, 4$ when there are more potential

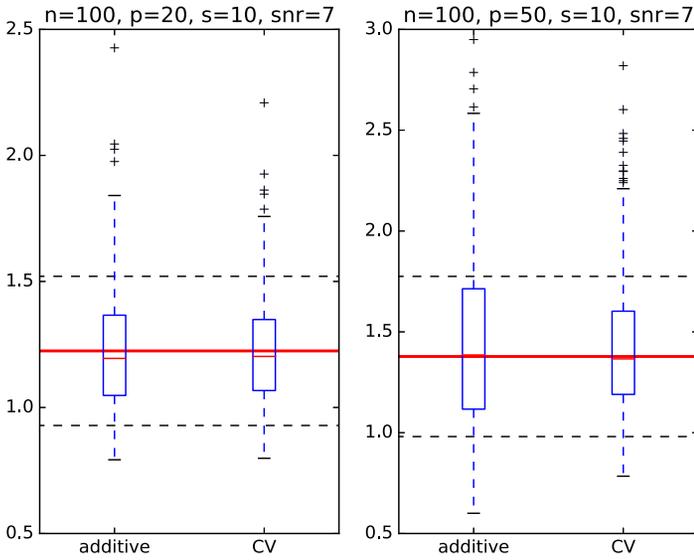


FIG. 4. Out of sample prediction error by $\widehat{\text{Err}}_{\text{out},\alpha}^{(I)}$ and cross validation, respectively.

submodels to select from. In contrast, $\widehat{\text{Err}}_{\alpha}^{(I)}$ hugs the true prediction error at every subset size k .

6. Discussion. In this work, we propose a method for estimating the prediction error after some data snooping in selecting a model. Remarkably, our estimation is not specific to any particular model selection procedures so long as it does not select too many variables to include in the model and it picks up some signals in the data. Different examples are considered.

In the following, we propose two more aspects of the problem that deserve attention but we do not seek to address in this work.

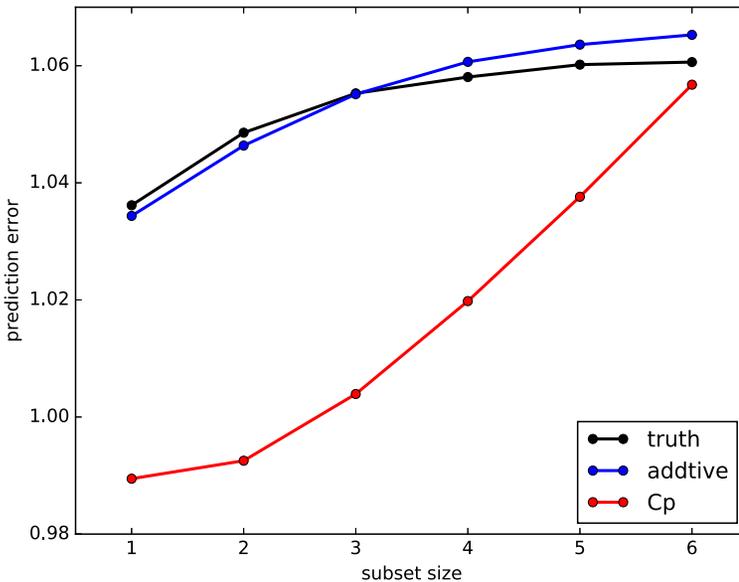


FIG. 5. Comparison of different estimates of prediction errors.

- We mainly focus on “in-sample” prediction errors, with the exception of Section 4.3. But as pointed in Section 4.3, although we can provide a consistent estimator of the (in-sample) prediction error in high dimensions, the same is not true for out-of-sample errors. [Klement, Mamlouk and Martinetz \(2008\)](#) points out that the same difficulty exists for cross validation as well. Under what assumptions can we provide a good estimator for out-of-sample prediction error in high dimensions remains a very interesting question.
- Throughout the work, we assume that the data comes from a homoscedastic normal model (1.1). Some simulations show that the performance of our estimator persists when the noise in the data is sub-Gaussian. The authors of [Tian and Taylor \(2018\)](#) pointed out that it is important that the tail of the randomization noise is heavier than that of the data. Since we add Gaussian noise for randomization, we suspect that the normal assumption on the data can be replaced by a sub-Gaussian assumption. Alternatively, we may investigate what other randomization noise we may add to the data when we have heavier-tailed data.

7. Proof of the lemmas.

7.1. *Proof of Lemma 3.6.* First notice that for hat matrix of the form in (3.2), we have

$$\mathbb{E}[\|\hat{\mu}(y) - \mu\|^2] = \sum_{i=1}^{|\mathcal{M}|} \int_{U_i} \|H_i y - \mu\|^2 \phi(y; \mu, \sigma^2) dy,$$

where H_i is short for H_{M_i} and $\phi(\cdot; \mu, \sigma^2)$ is the density for $N(\mu, \sigma^2 I)$. Let $\sigma^2 \leq \tau \leq (1 + \delta)\sigma^2$, where δ is defined in Assumption 3.3, and we define

$$g(\tau) = \mathbb{E}[\|\hat{\mu}(u) - \mu\|^2], \quad u \sim N(\mu, \tau I).$$

We note that g is differentiable with respect to τ and

$$\begin{aligned} & \frac{1}{n} |\mathbb{E}[\|\hat{\mu}(y + \omega) - \mu\|^2] - \mathbb{E}[\|\hat{\mu}(y) - \mu\|^2]| \\ (7.1) \quad &= \frac{1}{n} |g((1 + \alpha)\sigma^2) - g(\sigma^2)|, \\ & 0 < \alpha \leq \delta. \end{aligned}$$

Moreover, we have

$$\begin{aligned} \partial_\tau g(\tau) &= \sum_{i=1}^{|\mathcal{M}|} \int_{U_i} \|H_i u - \mu\|^2 \partial_\tau \left[\frac{1}{(\sqrt{2\pi\tau})^n} \exp\left(-\frac{\|u - \mu\|_2^2}{2\tau}\right) \right] du \\ (7.2) \quad &= \frac{1}{2\tau} \sum_{i=1}^{|\mathcal{M}|} \int_{U_i} \|H_i u - \mu\|^2 \left[\frac{\|u - \mu\|^2}{\tau} - n \right] \phi(u; \mu, \tau) du \\ &= \frac{1}{2\tau} \mathbb{E} \left[\|H_{\hat{M}(u)} u - \mu\|^2 \left[\frac{\|u - \mu\|^2}{\tau} - n \right] \right], \\ |\partial_\tau g(\tau)| &\leq \frac{1}{2\tau} \text{Var}[\|H_{\hat{M}(u)} u - \mu\|^2]^{\frac{1}{2}} \mathbb{E} \left[\left[\frac{\|u - \mu\|^2}{\tau} - n \right]^2 \right]^{\frac{1}{2}}. \end{aligned}$$

The last inequality holds because $\mathbb{E}[\frac{\|u - \mu\|^2}{\tau} - n] = 0$ and the expectations are taken over $u \sim N(\mu, \tau I)$. Per Assumption 3.3,

$$\text{Var}[\|H_{\hat{M}(u)} u - \mu\|^2]^{\frac{1}{2}} = O(\sqrt{n}\sigma^2).$$

Moreover, note that $\|u - \mu\|^2/\tau$ is a χ_n^2 distribution with mean n , thus

$$\mathbb{E}\left[\left[\frac{\|u - \mu\|^2}{\tau} - n\right]^2\right] = \text{Var}[\chi_n^2] = 2n.$$

Combining the above inequalities with (7.2) and we have

$$|\partial_\tau g(\tau)| = O(n\sigma^2) \quad \forall \tau \in [\sigma^2, (1 + \delta)\sigma^2].$$

Therefore, for any $0 < \alpha \leq \delta$, we have

$$\begin{aligned} & \frac{1}{n} |\mathbb{E}[\|\hat{\mu}(y + \omega) - \mu\|^2] - \mathbb{E}[\|\hat{\mu}(y) - \mu\|^2]| \\ &= \frac{1}{n} |g((1 + \alpha)\sigma^2) - g(\sigma^2)| = O(\alpha\sigma^2). \end{aligned}$$

7.2. Lemma 7.1 and its proof.

LEMMA 7.1. *If $Z_1, Z_2, \dots, Z_m \sim \chi_k^2$, not necessarily independently distributed, then*

$$\mathbb{E}\left[\max_{i \leq m} Z_i^2\right] = O[k^2 + (\log m)^2].$$

PROOF. Using union bound, we have

$$\mathbb{P}\left(\max_{i \leq m} Z_i > t\right) = m\mathbb{P}(Z_i > t).$$

Using the following bound for χ^2 random variables derived in Laurent and Massart (2000)

$$\mathbb{P}(Z_i \leq k + 2\sqrt{kx} + 2x) \leq \exp(-x) \quad \forall x > 0,$$

we take $x = \log m + t$ for any $t > 0$, and have

$$\begin{aligned} & \mathbb{P}\left(\max_{i \leq m} Z_i > k + 2\sqrt{k \log m + kt} + 2 \log m + 2t\right) \\ & \leq m \times \mathbb{P}(Z_i > k + 2\sqrt{k \log m + kt} + 2 \log m + 2t) \\ & \leq m \exp(-\log m - t) = \exp(-t). \end{aligned}$$

Thus it is easy to see that $\max_{i \leq m} Z_i$ is stochastically dominated by

$$k + 2\sqrt{k \log n} + 2 \log n + \xi,$$

where $\xi \geq 0$ and has exponential tails satisfying

$$\mathbb{P}(\xi > 2t + 2\sqrt{kt}) \leq \exp(-t), \quad \forall t > 0.$$

Since $X_i \geq 0$, we have

$$\mathbb{E}\left[\max_{i \leq m} X_i^2\right] \leq \mathbb{E}[(2 \log m + 2\sqrt{k \log m} + k + \xi)^2] = O[(\log m)^2 + k^2]. \quad \square$$

7.3. *Proof of Lemma 3.4.* We denote $\hat{\mu}_{RL} = H_{\hat{M}}y$ to be the relaxed Lasso solution where \hat{M} is the subset of variables selected by Lasso with regularization parameter λ , and $\hat{\mu}_{lasso}$ to be the Lasso solution. Then the variance can be upper bounded by

$$(7.3) \quad \text{Var}[\|\hat{\mu}_{RL} - \mu\|^2] \leq 2 \text{Var}[\|\hat{\mu}_{RL} - \hat{\mu}_{lasso}\|^2] + 2 \text{Var}[\|\hat{\mu}_{lasso} - \mu\|^2]$$

Using the Karush–Kuhn–Tucker conditions outlined in Tibshirani and Taylor (2012), we have

$$\hat{\mu}_{lasso} = H_{\hat{M}}y - \lambda X_{\hat{M}}(X_{\hat{M}}^T X_{\hat{M}})^{-1} z_{\hat{M}} = \hat{\mu}_{RL} - \lambda X_{\hat{M}}(X_{\hat{M}}^T X_{\hat{M}})^{-1} z_{\hat{M}},$$

where $z_{\hat{M}}$ is the signs of the active variables.

Without the loss of generality, we assume that the columns of X are of length \sqrt{n} . Then the regularization parameter will be chosen as

$$\lambda = \kappa \sigma \sqrt{2n \log p},$$

for some $\kappa > 1$. Therefore the first term in (7.3) can be bounded by

$$\begin{aligned} & \text{Var}[\|\hat{\mu}_{RL} - \hat{\mu}_{lasso}\|^2] \\ & \leq \text{Var}[\|\lambda X_{\hat{M}}(X_{\hat{M}}^T X_{\hat{M}})^{-1} z_{\hat{M}}\|^2] \\ & \leq \lambda^4 \cdot \frac{1}{n^2 \phi_0^4} K^2 = O(K^2 (\log p)^2), \end{aligned}$$

where the last inequality uses the restricted eigenvalue condition with parameter ϕ_0 .

To bound the second term in (7.3), we use the well-known results on the variance of normal distributions derived in Chen (1982), Chernoff (1981) and the fact that the Lasso solution is continuous and almost differentiable,

$$\begin{aligned} & \text{Var}[\|\hat{\mu}_{lasso} - \mu\|^2] \\ & \leq \mathbb{E}[\|H_{\hat{M}}(\hat{\mu}_{lasso} - \mu)\|^2] \leq \mathbb{E}[\|\hat{\mu}_{lasso} - \mu\|^2] \\ & \leq 2\sigma \sqrt{2n \log p} \|\hat{\beta}_{lasso} - \beta^0\|_1 + 2\lambda[\|\beta^0\|_1 - \|\hat{\beta}_{lasso}\|_1] \\ & \leq 2\lambda \|\beta^0\|_1 = O(s \sqrt{n \log p}). \end{aligned}$$

The first inequality uses the fact that $\nabla_y \hat{\mu}_{lasso}$ is $H_{\hat{M}}$ almost everywhere and the last inequality is the standard consistency result for the Lasso solutions.

Finally, Tian and Taylor (2017) showed that under the restrictive eigenvalue conditions, the number of selected variables K is a multiple of s and thus we get the conclusion of the lemma.

Acknowledgments. The author wants to thank Professor Jonathan Taylor, Professor Robert Tibshirani, Frederik Mikkelsen and Professor Ryan Tibshirani for useful discussions during this project.

REFERENCES

BICKEL, P. J., RITOV, Y. and TSYBAKOV, A. B. (2009). Simultaneous analysis of lasso and Dantzig selector. *Ann. Statist.* **37** 1705–1732. MR2533469 <https://doi.org/10.1214/08-AOS620>
 CHEN, L. H. Y. (1982). An inequality for the multivariate normal distribution. *J. Multivariate Anal.* **12** 306–315. MR0661566 [https://doi.org/10.1016/0047-259X\(82\)90022-7](https://doi.org/10.1016/0047-259X(82)90022-7)
 CHERNOFF, H. (1981). A note on an inequality involving the normal distribution. *Ann. Probab.* **9** 533–535. MR0614640

- DWORK, C. (2008). Differential privacy: A survey of results. In *Theory and Applications of Models of Computation. Lecture Notes in Computer Science* **4978** 1–19. Springer, Berlin. MR2472670 https://doi.org/10.1007/978-3-540-79228-4_1
- EFRON, B. (2004). The estimation of prediction error: Covariance penalties and cross-validation. *J. Amer. Statist. Assoc.* **99** 619–642. MR2090899 <https://doi.org/10.1198/016214504000000692>
- EFRON, B., HASTIE, T., JOHNSTONE, I. and TIBSHIRANI, R. (2004). Least angle regression. *Ann. Statist.* **32** 407–499. MR2060166 <https://doi.org/10.1214/009053604000000067>
- KLEMENT, S., MAMLOUK, A. M. and MARTINETZ, T. (2008). Reliability of cross-validation for svms in high-dimensional, low sample size scenarios. In *International Conference on Artificial Neural Networks* 41–50. Springer, Berlin.
- LAURENT, B. and MASSART, P. (2000). Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.* **28** 1302–1338. MR1805785 <https://doi.org/10.1214/aos/1015957395>
- LI, K.-C. (1989). Honest confidence regions for nonparametric regression. *Ann. Statist.* **17** 1001–1008. MR1015135 <https://doi.org/10.1214/aos/1176347253>
- MALLOWS, C. L. (1973). Some comments on c p. *Technometrics* **15** 661–675.
- MEINSHAUSEN, N. (2007). Relaxed Lasso. *Comput. Statist. Data Anal.* **52** 374–393. MR2409990 <https://doi.org/10.1016/j.csda.2006.12.019>
- MIKKELSEN, F. R. and HANSEN, N. R. (2018). Degrees of freedom for piecewise Lipschitz estimators. *Ann. Inst. Henri Poincaré Probab. Stat.* **54** 819–841. MR3795067 <https://doi.org/10.1214/17-AIHP822>
- NEGAHBAN, S. N., RAVIKUMAR, P., WAINWRIGHT, M. J. and YU, B. (2009). A unified framework for high-dimensional analysis of M -estimators with decomposable regularizers. In *Advances in Neural Information Processing Systems* 1348–1356.
- REID, S., TIBSHIRANI, R. and FRIEDMAN, J. (2016). A study of error variance estimation in Lasso regression. *Statist. Sinica* **26** 35–67. MR3468344
- SHEN, X. and YE, J. (2002). Adaptive model selection. *J. Amer. Statist. Assoc.* **97** 210–221. MR1947281 <https://doi.org/10.1198/016214502753479356>
- STEIN, C. M. (1981). Estimation of the mean of a multivariate normal distribution. *Ann. Statist.* **9** 1135–1151. MR0630098
- SUN, T. and ZHANG, C.-H. (2012). Scaled sparse linear regression. *Biometrika* **99** 879–898. MR2999166 <https://doi.org/10.1093/biomet/ass043>
- TIAN, X., LOFTUS, J. R. and TAYLOR, J. E. (2018). Selective inference with unknown variance via the square-root lasso. *Biometrika* **105** 755–768. MR3877864 <https://doi.org/10.1093/biomet/asy045>
- TIAN, X. and TAYLOR, J. (2017). Asymptotics of selective inference. *Scand. J. Stat.* **44** 480–499. MR3658523
- TIAN, X. and TAYLOR, J. (2018). Selective inference with a randomized response. *Ann. Statist.* **46** 679–710. MR3782381 <https://doi.org/10.1214/17-AOS1564>
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B* **58** 267–288. MR1379242
- TIBSHIRANI, R. J. (2015). Degrees of freedom and model search. *Statist. Sinica* **25** 1265–1296. MR3410308
- TIBSHIRANI, R. J. and TAYLOR, J. (2012). Degrees of freedom in lasso problems. *Ann. Statist.* **40** 1198–1232. MR2985948 <https://doi.org/10.1214/12-AOS1003>
- YANG, Y. (2005). Can the strengths of AIC and BIC be shared? A conflict between model identification and regression estimation. *Biometrika* **92** 937–950. MR2234196 <https://doi.org/10.1093/biomet/92.4.937>
- YE, J. (1998). On measuring and correcting the effects of data mining and model selection. *J. Amer. Statist. Assoc.* **93** 120–131. MR1614596 <https://doi.org/10.2307/2669609>
- ZOU, H., HASTIE, T. and TIBSHIRANI, R. (2007). On the “degrees of freedom” of the lasso. *Ann. Statist.* **35** 2173–2192. MR2363967 <https://doi.org/10.1214/009053607000000127>