

EXACT POST-SELECTION INFERENCE, WITH APPLICATION TO THE LASSO

BY JASON D. LEE, DENNIS L. SUN, YUEKAI SUN, AND JONATHAN E.
TAYLOR

Stanford University

We develop a general approach to valid inference after model selection. In a nutshell, our approach produces post-selection inferences with the same frequency guarantees as those given by data splitting but are more powerful. At the core of our framework is a result that characterizes the distribution of a post-selection estimator conditioned on the *selection event*. We specialize the approach to model selection by the lasso to form valid confidence intervals for the selected coefficients and test whether all relevant variables have been included in the model.

1. Introduction. As a statistical technique, linear regression is both simple and powerful. Not only does it provide estimates of the “effect” of each variable, but it also quantifies the uncertainty in those estimates, allowing for inference about the effects. However, in many exploratory investigations, a practitioner starts with a large pool of candidate variables, such as genes or demographic features, and does not know *a priori* which are relevant. This is especially a problem when there are more variables than observations, since then the (full) linear model is unidentified.

We might wish to use the data to select a subset of variables. One approach is to fit a linear model with all variables included (assuming this is possible), observe which ones are significant at level α , and then refit the linear model with only those variables included. The problem with this approach is that the p -values can no longer be trusted, since the variables that are selected will tend to be those that are significant. Intuitively, we are “overfitting” to this realization of the data.

To formalize the problem, consider the standard linear regression setup, where the response $\mathbf{y} \in \mathbb{R}^n$ is generated

$$(1.1) \quad \mathbf{y} \sim N(\boldsymbol{\mu}, \sigma^2 I_n)$$

and $\boldsymbol{\mu}$ is modeled as a linear function of predictors $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^n$. We choose a subset of the predictors $M \subset [p]$ and ask for inferences about the

AMS 2000 subject classifications: Primary 62F03, 62J07; secondary 62E15

Keywords and phrases: lasso, confidence interval, hypothesis test, model selection

submodel coefficients, i.e. the linear combination of predictors in M that minimizes the expected error:

$$(1.2) \quad \beta^M \equiv \arg \min_{\beta} \mathbb{E} \|\mathbf{y} - X_M \beta\|^2 = X_M^+ \boldsymbol{\mu},$$

where $X_M^+ \equiv (X_M^T X_M)^{-1} X_M^T$ is the pseudo-inverse of X_M . Notice that (1.2) implies that the targets β_j^M and $\beta_j^{M'}$ in different models $M \neq M'$ are in general different. This is simply a restatement of the well-known fact that a regression coefficient describes the effect of a predictor, *adjusting for the other predictors in the (sub)model*.

Alternatively, we may ask for inferences about a selected subset of the *full model coefficients*:

$$\beta^{[p]} \equiv \arg \min_{\beta} \mathbb{E} \|\mathbf{y} - X \beta\|^2 = X^+ \boldsymbol{\mu};$$

i.e. we ask for valid selective inferences about the $\beta^{[p]}$. For example, we may ask for inferences about the k most “significant” coefficients, e.g. the k coefficients with the largest absolute Z-score. In this setting, the targets $\beta_j^{[p]}$ remain consistent in different models.

Thus, “inference after selection” is ambiguous in linear regression because the target of inference may change with the selected model (Berk et al., 2013). In the next section, we discuss several ways to resolve this ambiguity.

2. Post-Selection Inference in Linear Regression. At first blush, the fact that the target changes with the model is deeply disturbing, since it seems to imply that the parameters are random. However, it is not so much the target that is random as our selection of targets. This is evidently the case when asking for selective inference about the full model coefficients. When asking for inference about the submodel coefficients, note that there are $p2^{p-1}$ possible submodel coefficient, one for each coefficient in all 2^p possible models:

$$\{\beta_j^M : M \subset [p], j \in M\}.$$

However, we only form inferences for a subset of these parameters—specifically, the parameters $\beta_j^{\hat{M}}$ in the model \hat{M} we select. To avoid repetition, we continue the discussion focusing on inference for the submodel targets.

To be concrete, suppose we want a confidence interval $C_j^{\hat{M}}$ for a submodel coefficient $\beta_j^{\hat{M}}$. What frequency properties should $C_j^{\hat{M}}$ have? By analogy to the classical setting, we might want

$$\mathbb{P}(\beta_j^{\hat{M}} \in C_j^{\hat{M}}) \geq 1 - \alpha,$$

but this is ill-posed because β_j^M does not exist when $j \notin M$. [Berk et al. \(2013\)](#) suggest the following alternatives:

1. Conditional Coverage: Since we form an interval for β_j^M if and only if we select model M , i.e., $\hat{M} = M$, it makes sense to condition on this event. Hence, we might require that our confidence interval C_j^M satisfy

$$(2.1) \quad \mathbb{P}(\beta_j^M \in C_j^M \mid \hat{M} = M) \geq 1 - \alpha.$$

The benefit of this approach is that we avoid ever having to compare coefficients across two different models $M \neq M'$.

2. Simultaneous Coverage: It also makes sense to talk about events that are defined simultaneously over all $j \in \hat{M}$. [Berk et al. \(2013\)](#) propose controlling the familywise error rate

$$(2.2) \quad FWER \equiv \mathbb{P}(\beta_j^{\hat{M}} \notin C_j^{\hat{M}} \text{ for any } j \in \hat{M}),$$

but it may be too stringent when there are many predictors involved. Instead of controlling the probability of making any error, we can control the expected proportion of errors—although the “proportion of errors” is ambiguous when we select no variables. We can simply declare the error to be zero when $|\hat{M}| = 0$ ([Benjamini and Yekutieli, 2005](#)):

$$(2.3) \quad FCR \equiv \mathbb{E} \left[\frac{|\{j \in \hat{M} : \beta_j^{\hat{M}} \notin C_j^{\hat{M}}\}|}{|\hat{M}|}; |\hat{M}| > 0 \right],$$

or condition on $|\hat{M}| > 0$ ([Storey, 2003](#)):

$$(2.4) \quad pFCR \equiv \mathbb{E} \left[\frac{|\{j \in \hat{M} : \beta_j^{\hat{M}} \notin C_j^{\hat{M}}\}|}{|\hat{M}|} \mid |\hat{M}| > 0 \right].$$

Since $FCR = pFCR \cdot \mathbb{P}(|\hat{M}| > 0)$, $pFCR$ control implies FCR control.

Remarkably, conditional coverage (2.1) implies pFCR (2.4) (and hence, FCR) control.

THEOREM 2.1. *Consider a family of intervals $\{C_j^{\hat{M}}\}_{j \in \hat{M}}$ that each have conditional $(1 - \alpha)$ coverage:*

$$\mathbb{P}(\beta_j^{\hat{M}} \notin C_j^{\hat{M}} \mid \hat{M} = M) \leq \alpha, \text{ for all } M \text{ and } j \in M.$$

Then, $FCR \leq pFCR \leq \alpha$.

PROOF. Condition on \hat{M} and iterate expectations.

$$\begin{aligned}
pFCR &= \mathbb{E} \left[\mathbb{E} \left[\frac{|\{j \in \hat{M} : \beta_j^{\hat{M}} \notin C_j^{\hat{M}}\}|}{|\hat{M}|} \middle| \hat{M} \right] \middle| |\hat{M}| > 0 \right] \\
&= \mathbb{E} \left[\frac{\sum_{j \in \hat{M}} \mathbb{P}(\beta_j^{\hat{M}} \notin C_j^{\hat{M}} | \hat{M})}{|\hat{M}|} \middle| |\hat{M}| > 0 \right] \\
&\leq \mathbb{E} \left[\frac{\alpha |\hat{M}|}{|\hat{M}|} \middle| |\hat{M}| > 0 \right] \\
&= \alpha.
\end{aligned}$$

□

3. Outline of Our Approach. We have argued that post-selection intervals for regression coefficients should have $1 - \alpha$ coverage conditional on the model:

$$\mathbb{P}(\beta_j^M \in C_j^M \mid \hat{M} = M) \geq 1 - \alpha,$$

both because this criterion is interesting in its own right and because it implies FCR control. To obtain an interval with this property, we study the conditional distribution

$$(3.1) \quad \boldsymbol{\eta}_M^T \mathbf{y} \mid \{\hat{M} = M\},$$

which will allow, more generally, conditional inference for parameters of the form $\boldsymbol{\eta}_M^T \boldsymbol{\mu}$. In particular, $\beta_j^M = \mathbf{e}_j^T X_M^+ \boldsymbol{\mu}$ can be written in this form, as can many other linear contrasts.

Our paper focuses on the specific case where the lasso is used to select the model \hat{M} . We begin in Section 4 by characterizing the event $\{\hat{M} = M\}$ for the lasso. As it turns out, this event is a union of polytopes. More precisely, the event $\{\hat{M} = M, \hat{\mathbf{s}}_M = \mathbf{s}_M\}$, that specifies the model *and* the signs of the selected variables, is a polytope of the form

$$\{\mathbf{y} \in \mathbb{R}^n : A(M, \mathbf{s}_M) \mathbf{y} \leq \mathbf{b}(M, \mathbf{s}_M)\}.$$

Therefore, if we condition on both the model and the signs, then we only need to study

$$(3.2) \quad \boldsymbol{\eta}^T \mathbf{y} \mid \{A \mathbf{y} \leq \mathbf{b}\}.$$

We do this in Section 5. It turns out that this conditional distribution is essentially a (univariate) truncated Gaussian. We use this to derive a statistic $F^z(\boldsymbol{\eta}^T \mathbf{y})$ whose distribution given $\{A \mathbf{y} \leq \mathbf{b}\}$ is $\text{Unif}(0, 1)$.

3.1. *Related Work.* The resulting post-selection test has a similar structure to the pathwise significance tests of [Lockhart et al. \(2014\)](#) and [Taylor et al. \(2014\)](#), which also are conditional tests. However, the intended application of our test is different. Whereas they test the specific hypothesis of whether a newly added coefficient along the LARS path is non-zero, our framework allows more general questions about the model the lasso selects: we can test the model at any value of λ or form confidence intervals for an individual coefficient in the model.

There is also a parallel literature on confidence intervals for coefficients in high-dimensional linear models based on the lasso estimator ([Javanmard and Montanari, 2013](#); [van de Geer et al., 2013](#); [Zhang and Zhang, 2014](#)). The difference between their work and ours is that they do not address post-selection inference; their target is β^0 , the coefficients in the true model, rather than $\beta^{\hat{M}}$, the coefficients in the selected model. The two will not be the same unless \hat{M} happens to contain all non-zero coefficients of β^0 . Although inference for β^0 is appealing, it requires assumptions about correctness of the linear model and sparsity of β^0 . Our approach instead regards the selected model as a linear approximation to the truth, a view shared by [Berk et al. \(2013\)](#).

4. The Lasso and Its Selection Event. In this paper, we apply our post-selection inference procedure to the model selected by the lasso ([Tibshirani, 1996](#)). The lasso estimate is the solution to the usual least squares problem with an additional ℓ_1 penalty on the coefficients:

$$(4.1) \quad \hat{\beta} \in \arg \min_{\beta} \frac{1}{2} \|\mathbf{y} - X\beta\|_2^2 + \lambda \|\beta\|_1.$$

The ℓ_1 penalty shrinks many of the coefficients to exactly zero, and the tradeoff between sparsity and fit to the data is controlled by the penalty parameter $\lambda \geq 0$.

Because the lasso produces sparse solutions, we can define model “selected” by the lasso to be simply the set of predictors with non-zero coefficients:

$$\hat{M} = \{j : \hat{\beta}_j \neq 0\}.$$

Then, post-selection inference seeks to make inferences about β^M , given $\{\hat{M} = M\}$, as defined in [\(1.2\)](#).

The rest of this section focuses on characterizing this event $\{\hat{M} = M\}$. We begin by noting that in order for a vector of coefficients $\hat{\beta}$ and a vector of signs \hat{s} to be solutions to the lasso problem [\(4.1\)](#), it is necessary and

sufficient that they satisfy the Karush-Kuhn-Tucker (KKT) conditions:

$$(4.2) \quad X^T(X\hat{\boldsymbol{\beta}} - \mathbf{y}) + \lambda\hat{\mathbf{s}} = 0,$$

$$\hat{s}_i = \text{sign}(\hat{\beta}_i) \quad \text{if } \hat{\beta}_i \neq 0$$

$$(4.3) \quad \hat{s}_i \in [-1, 1] \quad \text{if } \hat{\beta}_i = 0$$

Following [Tibshirani \(2013\)](#), we consider the *equicorrelation set*

$$(4.4) \quad \hat{M} \equiv \{i \in \{1, \dots, p\} : |\hat{s}_i| = 1\}.$$

Notice that we have implicitly defined the model \hat{M} to be equicorrelation set. Since $|\hat{s}_i| = 1$ for any $\hat{\beta}_i \neq 0$, the equicorrelation set does in fact contain all predictors with non-zero coefficients, although it may also include some predictors with zero coefficients. However, for almost every λ , the equicorrelation set is precisely the set of predictors with non-zero coefficients.

It turns out that it is easier to first characterize $\{(\hat{M}, \hat{\mathbf{s}}) = (M, \mathbf{s})\}$ and obtain $\{\hat{M} = M\}$ as a corollary by taking a union over the possible signs. The next result is an important first step.

LEMMA 4.1. *Assume the columns of X are in general position ([Tibshirani, 2013](#)). Let $M \subset \{1, \dots, p\}$ and $\mathbf{s} \in \{-1, 1\}^{|M|}$ be a candidate set of variables and their signs, respectively. Define the random variables*

$$(4.5) \quad \mathbf{w}(M, \mathbf{s}) := (X_M^T X_M)^{-1} (X_M^T \mathbf{y} - \lambda \mathbf{s})$$

$$(4.6) \quad \mathbf{u}(M, \mathbf{s}) := X_{-M}^T (X_M^T)^+ \mathbf{s} + \frac{1}{\lambda} X_{-M}^T (I - P_M) \mathbf{y}.$$

where $P_M \equiv X_M (X_M^T X_M)^{-1} X_M^T$ is projection onto the column span of X_M . Then the selection procedure can be rewritten in terms of \mathbf{w} and \mathbf{u} as:

$$(4.7) \quad \left\{ (\hat{M}, \hat{\mathbf{s}}) = (M, \mathbf{s}) \right\} = \left\{ \text{sign}(\mathbf{w}(M, \mathbf{s})) = \mathbf{s}, \|\mathbf{u}(M, \mathbf{s})\|_\infty < 1 \right\}$$

PROOF. First, we rewrite the KKT conditions (4.2) by partitioning them according to the equicorrelation set \hat{M} , adopting the convention that $-\hat{M}$ means “variables not in \hat{M} .”

$$\begin{aligned} X_{\hat{M}}^T (X_{\hat{M}} \hat{\boldsymbol{\beta}}_{\hat{M}} - \mathbf{y}) + \lambda \hat{\mathbf{s}}_{\hat{M}} &= 0 \\ X_{-\hat{M}}^T (X_{\hat{M}} \hat{\boldsymbol{\beta}}_{\hat{M}} - \mathbf{y}) + \lambda \hat{\mathbf{s}}_{-\hat{M}} &= 0 \\ \text{sign}(\hat{\boldsymbol{\beta}}_{\hat{M}}) &= \hat{\mathbf{s}} \\ \|\hat{\mathbf{s}}_{-\hat{M}}\|_\infty &< 1. \end{aligned}$$

Since the KKT conditions are necessary and sufficient for a solution, we obtain that $\{(\hat{M}, \hat{\mathbf{s}}) = (M, \mathbf{s})\}$ if and only if there exist \mathbf{w} and \mathbf{u} satisfying:

$$\begin{aligned} X_M^T(X_M\mathbf{w} - \mathbf{y}) + \lambda\mathbf{s} &= 0 \\ X_{-M}^T(X_M\mathbf{w} - \mathbf{y}) + \lambda\mathbf{u} &= 0 \\ \text{sign}(\mathbf{w}) &= \mathbf{s} \\ \|\mathbf{u}\|_\infty &< 1. \end{aligned}$$

We can solve the first two equations for \mathbf{w} and \mathbf{u} to obtain the equivalent set of conditions

$$\begin{aligned} \mathbf{w} &= (X_M^T X_M)^{-1}(X_M^T \mathbf{y} - \lambda \mathbf{s}) \\ \mathbf{u} &= X_{-M}^T (X_M^T)^+ \mathbf{s} + \frac{1}{\lambda} X_{-M}^T (I - P_M) \mathbf{y} \\ \text{sign}(\mathbf{w}) &= \mathbf{s} \\ \|\mathbf{u}\|_\infty &< 1, \end{aligned}$$

where the first two are the definitions of \mathbf{w} and \mathbf{u} given in (4.5) and (4.6), and the last two are the conditions on \mathbf{w} and \mathbf{u} given in (4.7). \square

Lemma 4.1 is remarkable because it says that the event $\{(\hat{M}, \hat{\mathbf{s}}) = (M, \mathbf{s})\}$ can be rewritten as affine constraints on \mathbf{y} . This is because \mathbf{w} and \mathbf{u} are already affine functions of \mathbf{y} , and the constraints $\text{sign}(\cdot) = \mathbf{s}$ and $\|\cdot\|_\infty < 1$ can also be rewritten in terms of affine constraints. The following proposition makes this explicit.

PROPOSITION 4.2. *Let \mathbf{w} and \mathbf{u} be defined as in (4.5) and (4.6). Then:*

$$(4.8) \quad \{\text{sign}(\mathbf{w}) = \mathbf{s}, \|\mathbf{u}\|_\infty < 1\} = \left\{ \begin{pmatrix} A_0(M, \mathbf{s}) \\ A_1(M, \mathbf{s}) \end{pmatrix} \mathbf{y} < \begin{pmatrix} \mathbf{b}_0(M, \mathbf{s}) \\ \mathbf{b}_1(M, \mathbf{s}) \end{pmatrix} \right\}$$

where A_0, \mathbf{b}_0 encode the “inactive” constraints $\{\|\mathbf{u}\|_\infty < 1\}$, and A_1, \mathbf{b}_1 encode the “active” constraints $\{\text{sign}(\mathbf{w}) = \mathbf{s}\}$. These matrices have the explicit forms:

$$\begin{aligned} A_0(M, \mathbf{s}) &= \frac{1}{\lambda} \begin{pmatrix} X_{-M}^T (I - P_M) \\ -X_{-M}^T (I - P_M) \end{pmatrix} & \mathbf{b}_0(M, \mathbf{s}) &= \begin{pmatrix} \mathbf{1} - X_{-M}^T (X_M^T)^+ \mathbf{s} \\ \mathbf{1} + X_{-M}^T (X_M^T)^+ \mathbf{s} \end{pmatrix} \\ A_1(M, \mathbf{s}) &= -\text{diag}(\mathbf{s})(X_M^T X_M)^{-1} X_M^T & \mathbf{b}_1(M, \mathbf{s}) &= -\lambda \text{diag}(\mathbf{s})(X_M^T X_M)^{-1} \mathbf{s} \end{aligned}$$

PROOF. First, substituting expression (4.5) for \mathbf{w} , we rewrite the “active” constraints as

$$\begin{aligned} \{\text{sign}(\mathbf{w}) = \mathbf{s}\} &= \{\text{diag}(\mathbf{s})\mathbf{w} > 0\} \\ &= \{\text{diag}(\mathbf{s})(X_M^T X_M)^{-1}(X_M^T \mathbf{y} - \lambda \mathbf{s}) > 0\} \\ &= \{A_1(M, \mathbf{s})\mathbf{y} < \mathbf{b}_1(M, \mathbf{s})\}. \end{aligned}$$

Next, substituting expression (4.6) for \mathbf{u} , we rewrite the “inactive” constraints as

$$\begin{aligned} \{\|\mathbf{u}\|_\infty < 1\} &= \left\{ -\mathbf{1} < X_{-M}^T (X_M^T)^+ \mathbf{s} + \frac{1}{\lambda} X_{-M}^T (I - P_M) \mathbf{y} < \mathbf{1} \right\} \\ &= \{A_0(M, \mathbf{s})\mathbf{y} < \mathbf{b}_0(M, \mathbf{s})\} \end{aligned}$$

□

Combining Lemma 4.1 with Proposition 4.2, we obtain the following.

THEOREM 4.3. *Let $A(M, \mathbf{s}) = \begin{pmatrix} A_0(M, \mathbf{s}) \\ A_1(M, \mathbf{s}) \end{pmatrix}$ and $\mathbf{b}(M, \mathbf{s}) = \begin{pmatrix} \mathbf{b}_0(M, \mathbf{s}) \\ \mathbf{b}_1(M, \mathbf{s}) \end{pmatrix}$, where A_i and \mathbf{b}_i are defined in Proposition 4.2. Then:*

$$\{\hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}\} = \{A(M, \mathbf{s})\mathbf{y} \leq \mathbf{b}(M, \mathbf{s})\}.$$

As a corollary, $\{\hat{M} = M\}$ is simply the union of the above events over all possible sign patterns.

$$\text{COROLLARY 4.4. } \{\hat{M} = M\} = \bigcup_{\mathbf{s} \in \{-1, 1\}^{|M|}} \{A(M, \mathbf{s})\mathbf{y} \leq \mathbf{b}(M, \mathbf{s})\}.$$

Figure 1 illustrates Theorem 4.3 and Corollary 4.4. The lasso partitions of \mathbb{R}^n into polytopes according to the model it selects and the signs of the coefficients. The shaded area corresponds to the event $\{\hat{M} = \{1, 3\}\}$, which is a union of two polytopes. Notice that the sign patterns $\{+, -\}$ and $\{-, +\}$ are not possible for the model $\{1, 3\}$.

5. Conditioning on Polytopes. In order to obtain inference conditional on the model, we need to understand the distribution of

$$\boldsymbol{\eta}_M^T \mathbf{y} \mid \{\hat{M} = M\}.$$

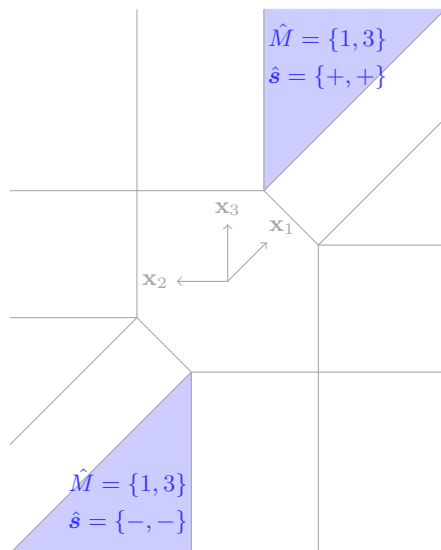


Fig 1: A geometric picture illustrating Theorem 4.3 for $n = 2$ and $p = 3$. The lasso partitions \mathbb{R}^n into polytopes according to the selected model and signs.

However, as we saw in the previous section, $\{\hat{M} = M\}$ is a union of polytopes, so it is easier to condition on both the model *and the signs*,

$$\boldsymbol{\eta}_{\hat{M}}^T \mathbf{y} \mid \{\hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}\},$$

where the conditioning event is a single polytope $\{A(M, \mathbf{s})\mathbf{y} \leq \mathbf{b}(M, \mathbf{s})\}$. Notice that inferences that are valid conditional on this finer event will also be valid conditional on $\{\hat{M} = M\}$. For example, if a confidence interval C_j^M for β_j^M has $(1 - \alpha)$ coverage conditional on the model and signs

$$\mathbb{P}(\beta_j^M \in C_j^M \mid \hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}) \geq 1 - \alpha,$$

then it will also have $(1 - \alpha)$ coverage conditional only on the model:

$$\begin{aligned} \mathbb{P}(\beta_j^M \in C_j^M \mid \hat{M} = M) &= \sum_{\mathbf{s}} \mathbb{P}(\beta_j^M \in C_j^M \mid \hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}) \mathbb{P}(\hat{\mathbf{s}} = \mathbf{s} \mid \hat{M} = M) \\ &\geq \sum_{\mathbf{s}} (1 - \alpha) \mathbb{P}(\hat{\mathbf{s}} = \mathbf{s} \mid \hat{M} = M) \\ &= 1 - \alpha. \end{aligned}$$

This section is divided into two subsections. First, we study how to condition on a single polytope; this will allow us to condition on $\{\hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}\}$. Then, we look at how to extend the framework to condition on a union of polytopes, which will allow us to condition only on the model $\{\hat{M} = M\}$. The inferences obtained by conditioning on the model will in general be more efficient (i.e., narrower intervals, more powerful tests), at the price of more computation.

5.1. *Conditioning on a Single Polytope.* Suppose we observe $\mathbf{y} \sim N(\boldsymbol{\mu}, \Sigma)$, and $\boldsymbol{\eta} \in \mathbb{R}^n$ is some direction of interest. To understand the distribution of

$$(5.1) \quad \boldsymbol{\eta}^T \mathbf{y} \mid \{\mathbf{A}\mathbf{y} \leq \mathbf{b}\},$$

we rewrite $\{\mathbf{A}\mathbf{y} \leq \mathbf{b}\}$ in terms of $\boldsymbol{\eta}^T \mathbf{y}$ and a component \mathbf{z} which is independent of $\boldsymbol{\eta}^T \mathbf{y}$. That component is

$$(5.2) \quad \mathbf{z} \equiv (I_n - \mathbf{c}\boldsymbol{\eta}^T)\mathbf{y},$$

where

$$(5.3) \quad \mathbf{c} \equiv \Sigma\boldsymbol{\eta}(\boldsymbol{\eta}^T\Sigma\boldsymbol{\eta})^{-1}.$$

It is easy to verify that \mathbf{z} is uncorrelated with, and hence independent of, $\boldsymbol{\eta}^T \mathbf{y}$. Although definition (5.2) may seem unmotivated, in the case where $\Sigma = \sigma^2 I_n$, \mathbf{z} is simply the residual $(I_n - P_{\boldsymbol{\eta}})\mathbf{y}$ from projecting \mathbf{y} onto $\boldsymbol{\eta}$.

We can now rewrite $\{\mathbf{A}\mathbf{y} \leq \mathbf{b}\}$ in terms of $\boldsymbol{\eta}^T \mathbf{y}$ and \mathbf{z} .

LEMMA 5.1. *Let \mathbf{z} be defined as in (5.2) and \mathbf{c} as in (5.3). Then, the conditioning set can be rewritten as follows:*

$$\{\mathbf{A}\mathbf{y} \leq \mathbf{b}\} = \{\mathcal{V}^-(\mathbf{z}) \leq \boldsymbol{\eta}^T \mathbf{y} \leq \mathcal{V}^+(\mathbf{z}), \mathcal{V}^0(\mathbf{z}) \geq 0\}$$

where

$$(5.4) \quad \mathcal{V}^-(\mathbf{z}) \equiv \max_{j:(\mathbf{A}\mathbf{c})_j < 0} \frac{b_j - (\mathbf{A}\mathbf{z})_j}{(\mathbf{A}\mathbf{c})_j}$$

$$(5.5) \quad \mathcal{V}^+(\mathbf{z}) \equiv \min_{j:(\mathbf{A}\mathbf{c})_j > 0} \frac{b_j - (\mathbf{A}\mathbf{z})_j}{(\mathbf{A}\mathbf{c})_j}$$

$$(5.6) \quad \mathcal{V}^0(\mathbf{z}) \equiv \min_{j:(\mathbf{A}\mathbf{c})_j = 0} b_j - (\mathbf{A}\mathbf{z})_j.$$

Note that \mathcal{V}^- , \mathcal{V}^+ , and \mathcal{V}^0 refer to functions. Since they are functions of \mathbf{z} only, (5.4)–(5.6) are independent of $\boldsymbol{\eta}^T \mathbf{y}$.

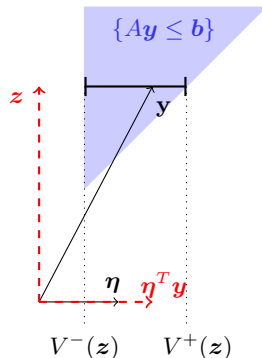


Fig 2: A geometric interpretation of why the event $\{A\mathbf{y} \leq \mathbf{b}\}$ can be characterized as $\{\mathcal{V}^-(\mathbf{z}) \leq \boldsymbol{\eta}^T \mathbf{y} \leq \mathcal{V}^+(\mathbf{z})\}$. Assuming $\Sigma = I$ and $\|\boldsymbol{\eta}\|_2 = 1$, $\mathcal{V}^-(\mathbf{z})$ and $\mathcal{V}^+(\mathbf{z})$ are functions of \mathbf{z} only, which is independent of $\boldsymbol{\eta}^T \mathbf{y}$.

PROOF. We can decompose $\mathbf{y} = \mathbf{c}(\boldsymbol{\eta}^T \mathbf{y}) + \mathbf{z}$ and rewrite the polytope as

$$\begin{aligned} \{A\mathbf{y} \leq \mathbf{b}\} &= \{A(\mathbf{c}(\boldsymbol{\eta}^T \mathbf{y}) + \mathbf{z}) \leq \mathbf{b}\} \\ &= \{A\mathbf{c}(\boldsymbol{\eta}^T \mathbf{y}) \leq \mathbf{b} - A\mathbf{z}\} \\ &= \{(A\mathbf{c})_j(\boldsymbol{\eta}^T \mathbf{y}) \leq b_j - (A\mathbf{z})_j \text{ for all } j\} \\ &= \left\{ \begin{array}{l} \boldsymbol{\eta}^T \mathbf{y} \leq \frac{b_j - (A\mathbf{z})_j}{(A\mathbf{c})_j} \text{ for } j : (A\mathbf{c})_j > 0 \\ \boldsymbol{\eta}^T \mathbf{y} \geq \frac{b_j - (A\mathbf{z})_j}{(A\mathbf{c})_j} \text{ for } j : (A\mathbf{c})_j < 0 \\ 0 \leq b_j - (A\mathbf{z})_j \text{ for } j : (A\mathbf{c})_j = 0 \end{array} \right\}, \end{aligned}$$

where in the last step, we have divided the components into three categories depending on whether $(A\mathbf{c})_j \gtrless 0$, since this affects the direction of the inequality (or whether we can divide at all). Since $\boldsymbol{\eta}^T \mathbf{y}$ is the same quantity for all j , it must be at least the maximum of the lower bounds and no more than the minimum of the upper bounds, which is precisely the definition of $\mathcal{V}^-(\mathbf{z})$ and $\mathcal{V}^+(\mathbf{z})$. Finally, $b_j - (A\mathbf{z})_j \geq 0$ for all $j : (A\mathbf{c})_j = 0$ is encoded by $\mathcal{V}^0(\mathbf{z}) \geq 0$. \square

Lemma 5.1 tells us that

$$(5.7) \quad [\boldsymbol{\eta}^T \mathbf{y} \mid \{A\mathbf{y} \leq \mathbf{b}\}] \stackrel{d}{=} [\boldsymbol{\eta}^T \mathbf{y} \mid \{\mathcal{V}^-(\mathbf{z}) \leq \boldsymbol{\eta}^T \mathbf{y} \leq \mathcal{V}^+(\mathbf{z}), \mathcal{V}^0(\mathbf{z}) \geq 0\}]$$

Since $\mathcal{V}^+(\mathbf{z}), \mathcal{V}^-(\mathbf{z}), \mathcal{V}^0(\mathbf{z})$ are independent of $\boldsymbol{\eta}^T \mathbf{y}$, they behave as “fixed” quantities. Thus, $\boldsymbol{\eta}^T \mathbf{y}$ is conditionally like a normal random variable, truncated to be between $\mathcal{V}^-(\mathbf{z})$ and $\mathcal{V}^+(\mathbf{z})$. We would like to be able to say

$$“\boldsymbol{\eta}^T \mathbf{y} \mid \{A\mathbf{y} \leq \mathbf{b}\} \sim TN(\boldsymbol{\eta}^T \boldsymbol{\mu}, \sigma^2 \boldsymbol{\eta}^T \Sigma \boldsymbol{\eta}, \mathcal{V}^-(\mathbf{z}), \mathcal{V}^+(\mathbf{z})),”$$

but this is technically incorrect, since the distribution on the right-hand side changes with \mathbf{z} . By conditioning on the value of \mathbf{z} , $\boldsymbol{\eta}^T \mathbf{y} \mid \{\mathbf{A}\mathbf{y} \leq \mathbf{b}, \mathbf{z} = \mathbf{z}_0\}$ is a truncated normal. We then use the probability integral transform to obtain a statistic $F^{\mathbf{z}}(\boldsymbol{\eta}^T \mathbf{y})$ that has a $\text{Unif}(0, 1)$ distribution for any value of \mathbf{z} . Hence, $F^{\mathbf{z}}(\boldsymbol{\eta}^T \mathbf{y})$ will also have a $\text{Unif}(0, 1)$ distribution marginally over \mathbf{z} . We make this precise in the next theorem.

THEOREM 5.2. *Let $F_{\mu, \sigma^2}^{[a, b]}$ denote the CDF of a $N(\mu, \sigma^2)$ random variable truncated to the interval $[a, b]$, i.e.:*

$$(5.8) \quad F_{\mu, \sigma^2}^{[a, b]}(x) = \frac{\Phi((x - \mu)/\sigma) - \Phi((a - \mu)/\sigma)}{\Phi((b - \mu)/\sigma) - \Phi((a - \mu)/\sigma)}$$

where Φ is the CDF of a $N(0, 1)$ random variable. Then:

$$(5.9) \quad F_{\boldsymbol{\eta}^T \boldsymbol{\mu}, \boldsymbol{\eta}^T \Sigma \boldsymbol{\eta}}^{[\mathcal{V}^-(\mathbf{z}), \mathcal{V}^+(\mathbf{z})]}(\boldsymbol{\eta}^T \mathbf{y}) \mid \{\mathbf{A}\mathbf{y} \leq \mathbf{b}\} \sim \text{Unif}(0, 1)$$

where \mathcal{V}^- and \mathcal{V}^+ are defined in (5.4) and (5.5). Furthermore,

$$[\boldsymbol{\eta}^T \mathbf{y} \mid \mathbf{A}\mathbf{y} \leq \mathbf{b}, \mathbf{z} = \mathbf{z}_0] \sim TN(\boldsymbol{\eta}^T \boldsymbol{\mu}, \sigma^2 \|\boldsymbol{\eta}\|^2, \mathcal{V}^-(\mathbf{z}_0), \mathcal{V}^+(\mathbf{z}_0)).$$

PROOF. First, apply Lemma 5.1:

$$\begin{aligned} [\boldsymbol{\eta}^T \mathbf{y} \mid \mathbf{A}\mathbf{y} \leq \mathbf{b}, \mathbf{z} = \mathbf{z}_0] &\stackrel{d}{=} [\boldsymbol{\eta}^T \mathbf{y} \mid \mathcal{V}^-(\mathbf{z}) \leq \boldsymbol{\eta}^T \mathbf{y} \leq \mathcal{V}^+(\mathbf{z}), \mathcal{V}^0(\mathbf{z}) \geq 0, \mathbf{z} = \mathbf{z}_0] \\ &\stackrel{d}{=} [\boldsymbol{\eta}^T \mathbf{y} \mid \mathcal{V}^-(\mathbf{z}_0) \leq \boldsymbol{\eta}^T \mathbf{y} \leq \mathcal{V}^+(\mathbf{z}_0), \mathcal{V}^0(\mathbf{z}_0) \geq 0, \mathbf{z} = \mathbf{z}_0] \end{aligned}$$

The only random quantities left are $\boldsymbol{\eta}^T \mathbf{y}$ and \mathbf{z} . Now we can eliminate $\mathbf{z} = \mathbf{z}_0$ from the condition using independence:

$$\begin{aligned} [\boldsymbol{\eta}^T \mathbf{y} \mid \mathbf{A}\mathbf{y} \leq \mathbf{b}, \mathbf{z} = \mathbf{z}_0] &\stackrel{d}{=} [\boldsymbol{\eta}^T \mathbf{y} \mid \mathcal{V}^-(\mathbf{z}_0) \leq \boldsymbol{\eta}^T \mathbf{y} \leq \mathcal{V}^+(\mathbf{z}_0)] \\ &\sim TN(\boldsymbol{\eta}^T \boldsymbol{\mu}, \sigma^2 \|\boldsymbol{\eta}\|^2, \mathcal{V}^-(\mathbf{z}_0), \mathcal{V}^+(\mathbf{z}_0)) \end{aligned}$$

Letting $F^{\mathbf{z}}(\boldsymbol{\eta}^T \mathbf{y}) \equiv F_{\boldsymbol{\eta}^T \boldsymbol{\mu}, \boldsymbol{\eta}^T \Sigma \boldsymbol{\eta}}^{[\mathcal{V}^-(\mathbf{z}), \mathcal{V}^+(\mathbf{z})]}(\boldsymbol{\eta}^T \mathbf{y})$, we can apply the probability integral transform to the above result to obtain

$$\begin{aligned} [F^{\mathbf{z}}(\boldsymbol{\eta}^T \mathbf{y}) \mid \mathbf{A}\mathbf{y} \leq \mathbf{b}, \mathbf{z} = \mathbf{z}_0] &\stackrel{d}{=} [F^{\mathbf{z}_0}(\boldsymbol{\eta}^T \mathbf{y}) \mid \mathbf{A}\mathbf{y} \leq \mathbf{b}, \mathbf{z} = \mathbf{z}_0] \\ &\sim \text{Unif}(0, 1) \end{aligned}$$

If we let p_X denote the density of a random variable X given $\{\mathbf{A}\mathbf{y} \leq \mathbf{b}\}$, what we have just shown is that

$$p_{F^{\mathbf{z}}(\boldsymbol{\eta}^T \mathbf{y}) \mid \mathbf{z}}(t \mid \mathbf{z}_0) \equiv \frac{p_{F^{\mathbf{z}}(\boldsymbol{\eta}^T \mathbf{y}), \mathbf{z}}(t, \mathbf{z}_0)}{p_{\mathbf{z}}(\mathbf{z}_0)} = 1_{[0, 1]}(f)$$

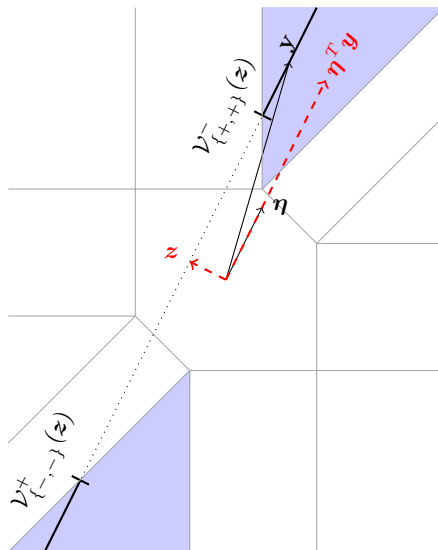


Fig 3: When we take the union over signs, the conditional distribution of $\boldsymbol{\eta}^T \mathbf{y}$ is truncated to a union of disjoint intervals. In this case, the Gaussian is truncated to the set $(-\infty, \mathcal{V}_{\{-,-\}}^+(z)] \cup [\mathcal{V}_{\{+,+\}}^-(z), \infty)$.

for any z_0 . The desired result now follows by integrating over z_0 :

$$\begin{aligned} p_{Fz}(\boldsymbol{\eta}^T \mathbf{y})(t) &= \int p_{Fz}(\boldsymbol{\eta}^T \mathbf{y})|z(t|z_0) p_z(z_0) dz_0 \\ &= \int 1_{[0,1]}(t) p_z(z_0) dz_0 \\ &= 1_{[0,1]}(t). \end{aligned}$$

□

5.2. *Conditioning on a Union of Polytopes.* We have just characterized the distribution of $\boldsymbol{\eta}^T \mathbf{y}$, conditional on \mathbf{y} falling into a single polytope $\{A\mathbf{y} \leq \mathbf{b}\}$. We obtain such a polytope if we condition on both the model and the signs $\{\hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}\}$. If we want to only condition on the model $\{\hat{M} = M\}$, then we will have to understand the distribution of $\boldsymbol{\eta}^T \mathbf{y}$, conditional on \mathbf{y} falling into a union of such polytopes, i.e.,

$$(5.10) \quad \boldsymbol{\eta}^T \mathbf{y} \left| \bigcup_s \{A_s \mathbf{y} \leq \mathbf{b}_s\}.$$

As Figure 3 makes clear, the argument proceeds exactly as before, except that $\boldsymbol{\eta}^T \mathbf{y}$ is now truncated to a union of intervals, instead of a single interval. There is a \mathcal{V}^- and a \mathcal{V}^+ for each possible sign pattern \mathbf{s} , so we index the intervals by the signs. This leads immediately to the next theorem, whose proof is essentially the same as that of Theorem 5.2.

THEOREM 5.3. *Let F_{μ, σ^2}^S denote the CDF of a $N(\mu, \sigma^2)$ random variable truncated to the set S . Then:*

$$(5.11) \quad F_{\boldsymbol{\eta}^T \boldsymbol{\mu}, \boldsymbol{\eta}^T \Sigma \boldsymbol{\eta}}^{\bigcup_{\mathbf{s}} [\mathcal{V}_{\mathbf{s}}^-(\mathbf{z}), \mathcal{V}_{\mathbf{s}}^+(\mathbf{z})]}(\boldsymbol{\eta}^T \mathbf{y}) \left| \bigcup_{\mathbf{s}} \{A_{\mathbf{s}} \mathbf{y} \leq \mathbf{b}_{\mathbf{s}}\} \sim \text{Unif}(0, 1),$$

where $\mathcal{V}_{\mathbf{s}}^-(\mathbf{z})$ and $\mathcal{V}_{\mathbf{s}}^+(\mathbf{z})$ are defined in (5.4) and (5.5) and $A = A_{\mathbf{s}}$ and $b = b_{\mathbf{s}}$.

6. Post-Selection Intervals for Regression Coefficients. In this section, we combine the characterization of the lasso selection event in Section 4 with the results about the distribution of a Gaussian truncated to a (union of) polytope(s) in Section 5 to form post-selection intervals for lasso-selected regression coefficients. The key link is that the lasso selection event can be expressed as a union of polytopes:

$$\begin{aligned} \{\hat{M} = M\} &= \bigcup_{\mathbf{s} \in \{-1, 1\}^{|M|}} \{\hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}\} \\ &= \bigcup_{\mathbf{s} \in \{-1, 1\}^{|M|}} \{A(M, \mathbf{s}) \mathbf{y} \leq \mathbf{b}(M, \mathbf{s})\}, \end{aligned}$$

where $A(M, \mathbf{s})$ and $\mathbf{b}(M, \mathbf{s})$ are defined in Theorem 4.3. Therefore, conditioning on selection is the same as conditioning on a union of polytopes, so the framework of Section 5 applies.

Recall that our goal is to form confidence intervals for $\beta_j^M = \mathbf{e}_j^T X_M^+ \boldsymbol{\mu}$, with $(1 - \alpha)$ -coverage conditional on $\{\hat{M} = M\}$. Taking $\boldsymbol{\eta} = (X_M^+)^T \mathbf{e}_j$, we can use Theorem 5.3 to obtain

$$F_{\beta_j^M, \sigma^2 \|\boldsymbol{\eta}\|^2}^{\bigcup_{\mathbf{s}} [\mathcal{V}_{\mathbf{s}}^-(\mathbf{z}), \mathcal{V}_{\mathbf{s}}^+(\mathbf{z})]}(\boldsymbol{\eta}^T \mathbf{y}) \mid \{\hat{M} = M\} \sim \text{Unif}(0, 1).$$

This gives us a test statistic for testing any hypothesized value of β_j^M . We can invert this test to obtain a confidence set

$$(6.1) \quad C_j^M \equiv \left\{ \beta_j^M : \frac{\alpha}{2} \leq F_{\beta_j^M, \sigma^2 \|\boldsymbol{\eta}\|^2}^{\bigcup_{\mathbf{s}} [\mathcal{V}_{\mathbf{s}}^-(\mathbf{z}), \mathcal{V}_{\mathbf{s}}^+(\mathbf{z})]}(\boldsymbol{\eta}^T \mathbf{y}) \leq 1 - \frac{\alpha}{2} \right\}.$$

In fact, the set C_j^M is an *interval*, as formalized in the next result.

THEOREM 6.1. *Let $\boldsymbol{\eta} = (X_M^+)^T \mathbf{e}_j$. Let L and U be the (unique) values satisfying*

$$F_{L, \sigma^2 \|\boldsymbol{\eta}\|^2}^{\cup_s [\mathcal{V}_s^-(\mathbf{z}), \mathcal{V}_s^+(\mathbf{z})]}(\boldsymbol{\eta}^T \mathbf{y}) = 1 - \frac{\alpha}{2} \quad F_{U, \sigma^2 \|\boldsymbol{\eta}\|^2}^{\cup_s [\mathcal{V}_s^-(\mathbf{z}), \mathcal{V}_s^+(\mathbf{z})]}(\boldsymbol{\eta}^T \mathbf{y}) = \frac{\alpha}{2}$$

Then $[L, U]$ is a $(1 - \alpha)$ confidence interval for β_j^M , conditional on $\{\hat{M} = M\}$, i.e.,

$$(6.2) \quad \mathbb{P}\left(\beta_j^M \in [L, U] \mid \hat{M} = M\right) = 1 - \alpha.$$

PROOF. By construction, $\mathbb{P}_{\beta_j^M}(\beta_j^M \in C_j^M \mid \hat{M} = M) = 1 - \alpha$, where C_j^M is defined in (6.1). The claim is that the set C_j^M is in fact the interval $[L, U]$. To see this, we need to show that the test statistic $F_{L, \sigma^2 \|\boldsymbol{\eta}\|^2}^{\cup_s [\mathcal{V}_s^-(\mathbf{z}), \mathcal{V}_s^+(\mathbf{z})]}(\boldsymbol{\eta}^T \mathbf{y})$ is monotone decreasing in β_j^M so that it crosses $1 - \frac{\alpha}{2}$ and $\frac{\alpha}{2}$ at unique values. This follows from the fact that the truncated Gaussian distribution has monotone likelihood ratio in the mean parameter. See Appendix A for details. \square

Alternatively, we could have conditioned on the signs, in addition to the model, so that we would only have to worry about conditioning on a single polytope. We also showed in Section 5 that

$$F_{\beta_j^M, \sigma^2 \|\boldsymbol{\eta}\|^2}^{\cup_s [\mathcal{V}_s^-(\mathbf{z}), \mathcal{V}_s^+(\mathbf{z})]}(\boldsymbol{\eta}^T \mathbf{y}) \mid \{\hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}\} \sim \text{Unif}(0, 1).$$

Inverting this statistic will produce intervals that have $(1 - \alpha)$ coverage conditional on $\{\hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}\}$, and hence, $(1 - \alpha)$ coverage conditional on $\{\hat{M} = M\}$. However, these intervals will be less efficient; they will in general be wider. However, one may be willing to sacrifice statistical efficiency for computational efficiency. Notice that the main cost in computing intervals according to Theorem 6.1 is determining the intervals $[\mathcal{V}_s^-(\mathbf{z}), \mathcal{V}_s^+(\mathbf{z})]$ for each $\mathbf{s} \in \{-1, 1\}^{|M|}$. The number of such sign patterns is $2^{|M|}$. While this might be feasible when $|M|$ is, say, less than 15, it is not feasible when we select hundreds of variables. Conditioning on the signs means that we only have to compute the interval $[\mathcal{V}_s^-(\mathbf{z}), \mathcal{V}_s^+(\mathbf{z})]$ for the sign pattern \mathbf{s} that was actually observed.

Figure 4 shows the tradeoff in statistical efficiency. When the signal is strong, as in the left-hand plot, there is virtually no difference between the intervals obtained by conditioning on just the model, or the model and signs. On the other hand, in the right-hand plot, we see that we can obtain very

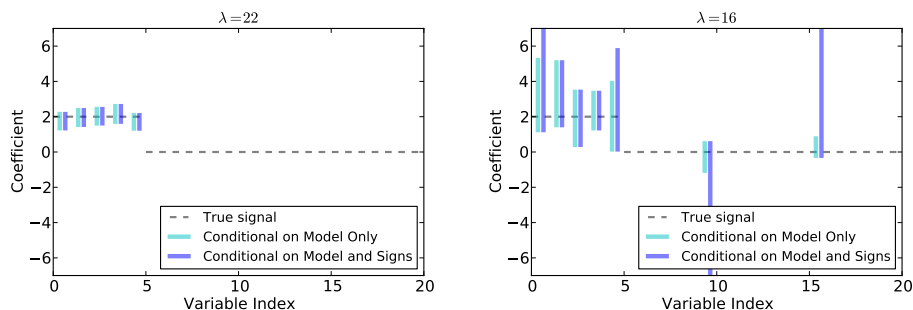


Fig 4: Comparison of the confidence intervals by conditioning on the model only (statistically more efficient, but computationally more expensive) and conditioning on both the model and signs (statistically less efficient, but computationally more feasible). Data were simulated for $n = 25$, $p = 50$, and 5 true non-zero coefficients; only the first 20 coefficients are shown. (Variables with no intervals are included to emphasize that inference is only on the selected variables.) Conditioning on the signs in addition to the model results in no loss of statistical efficiency when the signal is strong (left) but is problematic when the signal is weak (right).

wide intervals when the signal is weak. The widest intervals are for actual noise variables, as expected.

To understand why post-selection intervals are sometimes very wide, notice that when a truncated Gaussian random variable Z is close to the endpoints of the truncation interval $[a, b]$, there are many means μ that would be consistent with that observation—hence, the wide intervals. Figure 5 shows confidence intervals for μ as a function of Z . When Z is far from the endpoints of the truncation interval, we basically recover the nominal OLS intervals (i.e., not adjusted for selection).

The implications are clear. When the signal is strong, $\boldsymbol{\eta}^T \mathbf{y}$ will be far from the endpoints of the truncation region, so we obtain the nominal OLS intervals. On the other hand, when a variable just barely entered the model, then $\boldsymbol{\eta}^T \mathbf{y}$ will be close to the edge of the truncation region, and the interval will be wide.

6.1. *Optimality.* We have derived a confidence interval C_j^M whose conditional coverage, given $\{\hat{M} = M\}$, is $1 - \alpha$. The fact that we have found such an interval is not remarkable, since many such intervals have this property, including the trivial interval $(-\infty, \infty)$. However, given two intervals with the same coverage, we generally prefer the shorter one. We now show that C_j^M

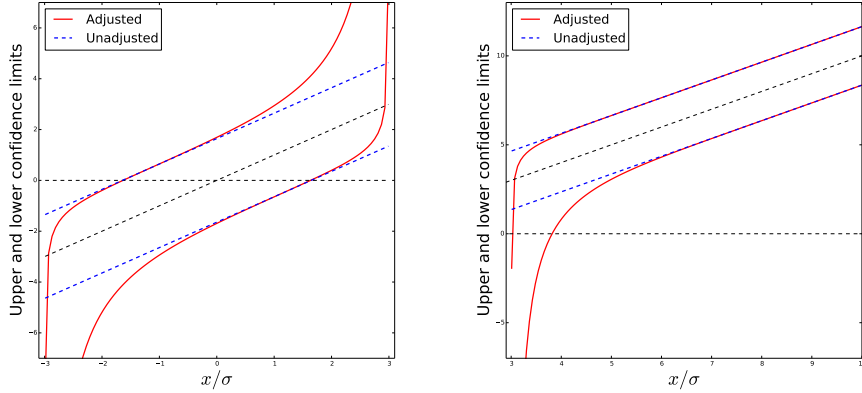


Fig 5: Upper and lower bounds of 90% confidence intervals for μ based on a single observation $x/\sigma \sim TN(0, 1, -3, 3)$. We see that as long as the observation x is roughly 0.5σ away from either boundary, the size of the intervals is comparable to the unadjusted OLS confidence interval.

is, with one small tweak, the interval with expected shortest length among all *unbiased* intervals with $1 - \alpha$ coverage.

An *unbiased* interval C for a parameter θ is one which covers no other parameter θ' with probability more than $1 - \alpha$, i.e.,

$$(6.3) \quad \mathbb{P}_\theta(\theta' \in C) \leq 1 - \alpha, \text{ for all } \theta, \theta' \neq \theta.$$

Unbiasedness is a common restriction to ensure that there is an optimal interval or test at all (Lehmann and Romano, 2005). The interval with expected shortest length for β_j^M , among all intervals with conditional $1 - \alpha$ coverage, is similar to the interval $[L, U]$ in Theorem 6.1. The only difference is that the critical values L and U were chosen symmetrically so that the pivot has $\alpha/2$ area in either tail. However, allocating the α probability equally to either tail may not be optimal in general. The next theorem provides the general recipe for constructing optimal intervals:

THEOREM 6.2. Let $\boldsymbol{\eta} = (X_M^+)^T \mathbf{e}_j$ and define the set

$$(6.4) \quad C_{j,opt}^M \equiv \{\beta_j^M : \boldsymbol{\eta}^T \mathbf{y} \in R(\mathbf{z})\},$$

where the rejection region $R(\mathbf{z})$ is defined in (6.6) below and is a function of \mathbf{z} (5.2). Then $C_{j,opt}^M$ is the interval with expected shortest length for β_j^M

with $1 - \alpha$ coverage, conditional on $\{\hat{M} = M\}$, i.e.,

$$(6.5) \quad \mathbb{P}\left(\beta_j^M \in C_{j,opt}^M \mid \hat{M} = M\right) = 1 - \alpha.$$

PROOF. Without loss of generality, assume $\|\boldsymbol{\eta}\|_2 = 1$ and $\sigma = 1$. Note that the distribution of \mathbf{y} , conditional on $\{\hat{M} = M\}$, is an exponential family. That is, with respect to some base measure ν on $\{\hat{M} = M\}$ and for some log-normalizing function ψ_M , its density is

$$\begin{aligned} \mathbf{y} &\sim \exp\{\boldsymbol{\mu}^T \mathbf{y} - \psi_M(\boldsymbol{\mu})\} d\nu(\mathbf{y}) \\ &= \exp\{\boldsymbol{\mu}^T [P_\boldsymbol{\eta} \mathbf{y} + (I - P_\boldsymbol{\eta}) \mathbf{y}] - \psi_M(\boldsymbol{\mu})\} d\nu(\mathbf{y}) \\ &= \exp\{(\boldsymbol{\eta}^T \boldsymbol{\mu})(\boldsymbol{\eta}^T \mathbf{y}) + ((I - P_\boldsymbol{\eta}) \boldsymbol{\mu})^T \mathbf{z} - \psi_M(\boldsymbol{\mu})\} d\nu(\mathbf{y}). \end{aligned}$$

This is an exponential family with one parameter of interest $(\boldsymbol{\eta}^T \boldsymbol{\mu})$ and $(n-1)$ nuisance parameters, represented by $(I - P_\boldsymbol{\eta}) \boldsymbol{\mu}$.¹ Classical theory (cf. Theorem 4.4.1 in [Lehmann and Romano \(2005\)](#)) says that the uniformly most powerful unbiased (UMPU) test of $H_0 : \boldsymbol{\eta}^T \boldsymbol{\mu} = \beta_j^M$ versus $H_1 : \boldsymbol{\eta}^T \boldsymbol{\mu} \neq \beta_j^M$ is obtained by conditioning on \mathbf{z} and rejecting for values of $\boldsymbol{\eta}^T \mathbf{y}$ that are too large or too small. In other words, we reject for $\boldsymbol{\eta}^T \mathbf{y} \in R(\mathbf{z})$, where

$$(6.6) \quad R(\mathbf{z}) \equiv (-\infty, C_1(\mathbf{z})] \cup [C_2(\mathbf{z}), \infty),$$

where $C_1(\mathbf{z})$ and $C_2(\mathbf{z})$ are chosen to ensure:

1. The test is level α . Letting $\varphi(x) = \frac{e^{-x^2}}{\sqrt{2\pi}}$ and $S(\mathbf{z}) \equiv \bigcup_s [\mathcal{V}_s^-(\mathbf{z}), \mathcal{V}_s^+(\mathbf{z})]$,

$$\frac{\int_{R(\mathbf{z}) \cap S(\mathbf{z})} \varphi(x - \beta_j^M) dx}{\int_{S(\mathbf{z})} \varphi(x - \beta_j^M) dx} = \alpha$$

for almost every \mathbf{z} .

2. Since the test is unbiased, its power function must be minimized under H_0 , so its derivative at β_j^M must be 0, yielding the condition

$$\frac{\int_{R(\mathbf{z}) \cap S(\mathbf{z})} x \varphi(x - \beta_j^M) dx}{\int_{S(\mathbf{z})} \varphi(x - \beta_j^M) dx} = \alpha \frac{\int_{S(\mathbf{z})} x \varphi(x - \beta_j^M) dx}{\int_{S(\mathbf{z})} \varphi(x - \beta_j^M) dx}.$$

This gives us the UMPU test of $H_0 : \boldsymbol{\eta}^T \mathbf{y} = \beta_j^M$ for every β_j^M . Finally, the UMAU interval can be obtained by inverting the UMPU test, i.e.,

$$C_{j,opt}^M = \{\beta_j^M : \boldsymbol{\eta}^T \mathbf{y} \in R(\mathbf{z})\}.$$

The construction above is standard, and the details can be found in Chapter 4 of [Lehmann and Romano \(2005\)](#). By the Ghosh-Pratt theorem, this is the unbiased interval with expected shortest length. \square

¹Although this is technically a vector in \mathbb{R}^n , it is only $(n-1)$ -dimensional.

7. Data Example. We apply our post-selection intervals to the diabetes data set from Efron et al. (2004). After standardizing all variables, we chose λ according to the strategy in Negahban et al. (2012), $\lambda = 2 \mathbf{E}(\|X^T \epsilon\|_\infty)$, using an estimate of σ from the full model, resulting in $\lambda \approx 190$. The lasso selected four variables: BMI, BP, S3, and S5.

The post-selection intervals are shown in Figure 6, alongside the nominal confidence intervals produced by fitting OLS to the four selected variables, ignoring selection. The nominal intervals do not have $(1 - \alpha)$ coverage conditional on the model and are not valid post-selection intervals. Also depicted are the confidence intervals obtained by *data splitting*; that is, if one splits the n observations into two halves, then uses one half for selection and the other for inference. This is a competitor method that also produces valid confidence intervals conditional on the model. The lasso selected the same four variables on half of the data, and then nominal intervals for these four variables using OLS on the other half of the data.

We can make two observations from Figure 6.

1. The adjusted intervals provided by our method essentially reproduces the OLS intervals for the strong effects, whereas data splitting intervals are wider by a factor of $\sqrt{2}$ (since only $n/2$ observations are used in the inference). For this dataset, the POSI intervals are 1.36 times wider than the OLS intervals. For all the variables, our method produces the shortest intervals among the methods that control selective type 1 error.
2. One variable, S3, which would have been deemed significant using the OLS intervals, is no longer significant after accounting for selection. Data splitting, our selection-adjusted intervals, and POSI intervals conclude that S3 is not significant. This demonstrates that taking model selection into account can have substantive impacts on the conclusions.

8. Testing the Lasso-Selected Model. Having observed that the lasso selected the variables \hat{M} , another relevant question is whether it has captured all of the signal in the model, i.e.,

$$(8.1) \quad H_0 : \beta_{-\hat{M}}^0 = \mathbf{0}.$$

We consider a slightly more general question, which does not assume the correctness of the linear model $\boldsymbol{\mu} = X\boldsymbol{\beta}^0$ and also takes into account whether the non-selected variables can improve the fit:

$$(8.2) \quad H_0 : X_{-\hat{M}}^T (I - P_{-\hat{M}}) \boldsymbol{\mu} = \mathbf{0}.$$

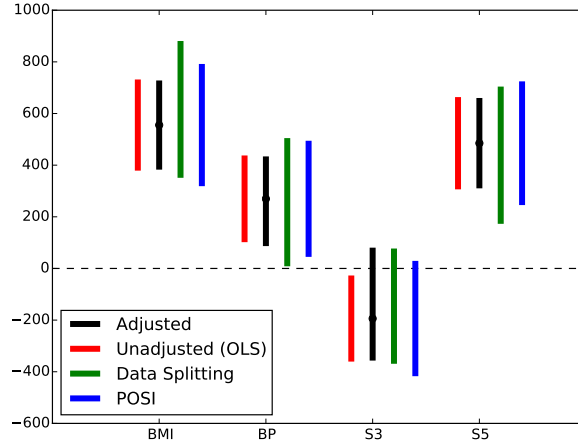


Fig 6: Inference for the four variables selected by the lasso ($\lambda = 190$) on the diabetes data set. The point estimate and adjusted confidence intervals using the approach in Section 6 are shown in black. The OLS intervals, which ignore selection, are shown in red. The green lines show the intervals produced by splitting the data into two halves, forming the interval based on only half of the data. The blue line corresponds to the POSI method of Berk et al. (2013)

This quantity is the partial correlation of the non-selected variables with μ , adjusting for the variables in \hat{M} . This is more general because if we assume $\mu = X\beta^0$ for some β^0 and X is full rank, then rejecting (8.2) implies that there exists $i \in \text{supp}(\beta^0)$ not in \hat{M} , so we would also reject (8.1).

The natural approach is to compare the observed partial correlations $X_{-M}^T(I - P_M)\mathbf{y}$ to $\mathbf{0}$. However, the framework of Section 5 only allows tests of μ in a single direction η . To make use of that framework, we can choose η such that it selects the maximum magnitude of $X_{-M}^T(I - P_M)\mathbf{y}$. In particular, this direction provides the most evidence against the null hypothesis of zero partial correlation, so if the null hypothesis cannot be rejected in this direction, it would not be rejected in any direction.

Letting $\hat{\mathbf{v}} := \text{argmax}_{\mathbf{v} \in B_\infty} \mathbf{v}^T X_{-M}^T(I - P_M)\mathbf{y}$, (B_∞ is the ℓ_∞ ball) we set

$$(8.3) \quad \eta = (I - P_M)X_{-M}\hat{\mathbf{v}},$$

and test $H_0 : \eta^T \mu = 0$. To ensure η is measurable with respect to the

selection event, we condition on not only on $(\hat{M}, \hat{\mathbf{s}})$, but also on $\hat{\mathbf{v}}$:

$$(8.4) \quad \{(\hat{M}, \hat{\mathbf{s}}, \hat{\mathbf{v}}) = (M, \mathbf{s}, \mathbf{v})\}.$$

A test that is level α conditional on (8.4) for all $(\hat{M}, \hat{\mathbf{s}}, \hat{\mathbf{v}})$ is also level α conditional on $(\hat{M}, \hat{\mathbf{s}})$.

In order to use the results of Section 5, we must show that (8.4) can be written in the form $A(M, \mathbf{s}, \mathbf{v})\mathbf{y} \leq b(M, \mathbf{s}, \mathbf{v})$. This is indeed possible, and the following proposition provides an explicit construction.

PROPOSITION 8.1. *Let $A_0, \mathbf{b}_0, A_1, \mathbf{b}_1$ be defined as in Proposition 4.2. Then:*

$$\{(\hat{M}, \hat{\mathbf{s}}, \hat{\mathbf{v}}) = (M, \mathbf{s}, \mathbf{v})\} = \left\{ \begin{pmatrix} A_0(M, \mathbf{s}) \\ A_1(M, \mathbf{s}) \\ A_2(M, \mathbf{v}) \end{pmatrix} \mathbf{y} < \begin{pmatrix} \mathbf{b}_0(M, \mathbf{s}) \\ \mathbf{b}_1(M, \mathbf{s}) \\ \mathbf{0} \end{pmatrix} \right\}$$

where $A_2(M, \mathbf{v})$ is defined as

$$A_2(M, \mathbf{v}) = D(\mathbf{v})X_{-M}^T(I - P_M),$$

where the rows of $D(\mathbf{v}) \in \mathbf{R}^{(2|M|-1) \times |M|}$ are given by $(\mathbf{w} - \mathbf{v})^T$ for all $\mathbf{w} \in \text{ext}(\mathcal{B}_\infty)$.

PROOF. The constraints $\{A_0\mathbf{y} < \mathbf{b}_0\}$ and $\{A_1\mathbf{y} < \mathbf{b}_1\}$ come from Proposition (4.2) and encode the constraints $\{(\hat{M}, \hat{\mathbf{s}}) = (M, \mathbf{s})\}$. We show that the last two sets of constraints encode $\{\hat{\mathbf{v}} = \mathbf{v}\}$.

Let $\mathbf{r} := X_{-M}^T(I - P_M)\mathbf{y}$ denote the vector of partial correlations. Since B_∞ is a polytope, the maximum of $\arg\max_{\mathbf{v} \in B_\infty} \mathbf{v}^T X_{-M}^T(I - P_M)\mathbf{y}$ is attained at an extreme point. Thus $\{\hat{\mathbf{v}} = \mathbf{v}\} = \{D(\mathbf{v})\mathbf{r} < 0\}$. \square

Because of Proposition 8.1, we can now obtain the following result as a simple consequence of Theorem 5.2, which says that $F_{0, \sigma^2 \|\boldsymbol{\eta}\|^2}^{[\mathcal{V}^-, \mathcal{V}^+]}(\boldsymbol{\eta}^T \mathbf{y}) \sim \text{Unif}(0, 1)$, conditional on the set (8.4) and H_0 . We reject when $F_{0, \sigma^2 \|\boldsymbol{\eta}\|^2}^{[\mathcal{V}^-, \mathcal{V}^+]}(\boldsymbol{\eta}^T \mathbf{y})$ is large because $F_{0, \sigma^2 \|\boldsymbol{\eta}\|^2}^{[\mathcal{V}^-, \mathcal{V}^+]}(\cdot)$ is monotone increasing in the argument and $\boldsymbol{\eta}^T \boldsymbol{\mu}$ is likely to be positive under the alternative.

COROLLARY 8.2. *Let H_0 and $\boldsymbol{\eta}$ be defined as in (8.3). Then, the test which rejects when*

$$\left\{ F_{0, \sigma^2 \|\boldsymbol{\eta}\|^2}^{[\mathcal{V}^-, \mathcal{V}^+]}(\boldsymbol{\eta}^T \mathbf{y}) > 1 - \alpha \right\}$$

is level α , conditional on $\{(\hat{M}, \hat{\mathbf{s}}, \hat{\mathbf{v}}) = (M, \mathbf{s}, \mathbf{v})\}$. That is,

$$\mathbb{P}_0 \left(F_{0, \sigma^2 \|\boldsymbol{\eta}\|^2}^{[\mathcal{V}^-, \mathcal{V}^+]}(\boldsymbol{\eta}^T \mathbf{y}) > 1 - \alpha \mid \{(\hat{M}, \hat{\mathbf{s}}, \hat{\mathbf{v}}) = (M, \mathbf{s}, \mathbf{v})\} \right) = \alpha.$$

In particular, since this holds for every $(M, \mathbf{s}, \mathbf{v})$, this test also controls Type I error conditional only on $(\hat{M}, \hat{\mathbf{s}})$, and unconditionally:

$$\begin{aligned} \mathbb{P}_0 \left(F_{0, \sigma^2 \|\boldsymbol{\eta}\|^2}^{[\mathcal{V}^-, \mathcal{V}^+]}(\boldsymbol{\eta}^T \mathbf{y}) > 1 - \alpha \mid \{(\hat{M}, \hat{\mathbf{s}}) = (M, \mathbf{s})\} \right) &= \alpha \\ \mathbb{P}_0 \left(F_{0, \sigma^2 \|\boldsymbol{\eta}\|^2}^{[\mathcal{V}^-, \mathcal{V}^+]}(\boldsymbol{\eta}^T \mathbf{y}) > 1 - \alpha \right) &= \alpha. \end{aligned}$$

9. Extension to the Elastic Net. One problem with the lasso is that it tends to select one variable out of a set of correlated variables, resulting in estimates that are unstable. One way to stabilize them is to add an ℓ_2 penalty to the lasso objective, resulting in the elastic net (Zou and Hastie, 2005):

$$(9.1) \quad \tilde{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - X\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 + \frac{\gamma}{2} \|\boldsymbol{\beta}\|_2^2.$$

Using a nearly identical argument to Lemma 4.1, we see that $\{\hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}\}$ if and only if there exist $\tilde{\mathbf{w}}$ and $\tilde{\mathbf{u}}$ satisfying

$$\begin{aligned} (X_M^T X_M + \gamma I) \tilde{\mathbf{w}} - X_M^T \mathbf{y} + \lambda \mathbf{s} &= 0 \\ X_{-M}^T (X_M \tilde{\mathbf{w}} - \mathbf{y}) + \lambda \tilde{\mathbf{u}} &= 0 \\ \operatorname{sign}(\tilde{\mathbf{w}}) &= \mathbf{s} \\ \|\tilde{\mathbf{u}}\|_\infty &< 1. \end{aligned}$$

These four conditions differ from those of Lemma 4.1 in only one respect: $X_M^T X_M$ in the first expression is replaced by $X_M^T X_M + \gamma I$. Continuing the argument of Section 4, we see that the selection event can be rewritten

$$(9.2) \quad \{\hat{M} = M, \hat{\mathbf{s}} = \mathbf{s}\} = \left\{ \begin{pmatrix} \tilde{A}_0(M, \mathbf{s}) \\ \tilde{A}_1(M, \mathbf{s}) \end{pmatrix} \mathbf{y} < \begin{pmatrix} \tilde{\mathbf{b}}_0(M, \mathbf{s}) \\ \tilde{\mathbf{b}}_1(M, \mathbf{s}) \end{pmatrix} \right\}$$

where \tilde{A}_k and $\tilde{\mathbf{b}}_k$ are analogous to A_k and \mathbf{b}_k in Proposition 4.2, except with $(X_M^T X_M)^{-1}$ replaced by $(X_M^T X_M + \gamma I)^{-1}$. Notice that this quantity not only appears explicitly in A_1 and \mathbf{b}_1 , but also appears implicitly in A_0 and \mathbf{b}_0 through P_M and $(X_M^T)^+$.

Now that we have rewritten the selection event in the form $\{\mathbf{A}\mathbf{y} \leq \mathbf{b}\}$, we can once again apply the framework in Section 5 to obtain a test for the elastic net conditional on this event.

10. Conclusion. Model selection and inference have long been regarded as conflicting goals in linear regression. Following the lead of Berk et al. (2013), we have proposed a framework for post-selection inference that *conditions on which model was selected*, i.e., the event $\{\hat{M} = M\}$. We characterize this event for the lasso and derive optimal and exact confidence intervals for linear contrasts $\boldsymbol{\eta}^T \boldsymbol{\mu}$, conditional on $\{\hat{M} = M\}$. With this general framework, we can form post-selection intervals for regression coefficients, equipping practitioners with a way to obtain “valid” intervals even after model selection.

Acknowledgements. We thank Will Fithian, Sam Gross, and Josh Loftus for helpful comments and discussions. In particular, Will Fithian provided insights that led to the geometric intuition of our procedure shown in Figure 2. J. Lee was supported by a National Defense Science and Engineering Graduate Fellowship and a Stanford Graduate Fellowship. D. L. Sun was supported by a Ric Weiland Graduate Fellowship and the Stanford Genome Training Program (SGTP; NIH/NHGRI). Y. Sun was partially supported by the NIH, grant U01GM102098. J.E. Taylor was supported by the NSF, grant DMS 1208857, and by the AFOSR, grant 113039.

APPENDIX A: MONOTONICITY OF F

LEMMA A.1. *Let $F_\mu(x) := F_{\mu, \sigma^2}^{[a, b]}(x)$ denote the cumulative distribution function of a truncated Gaussian random variable, as defined as in (5.8). Then $F_\mu(x)$ is monotone decreasing in μ .*

PROOF. First, the truncated Gaussian distribution with CDF $F_\mu := F_{\mu, \sigma^2}^{[a, b]}$ is a natural exponential family in μ , since it is just a Gaussian with a different base measure. Therefore, it has monotone likelihood ratio in μ . That is, for all $\mu_1 > \mu_0$ and $x_1 > x_0$:

$$\frac{f_{\mu_1}(x_1)}{f_{\mu_0}(x_1)} > \frac{f_{\mu_1}(x_0)}{f_{\mu_0}(x_0)}$$

where $f_{\mu_i} := dF_{\mu_i}$ denotes the density. (Instead of appealing to properties of exponential families, this property can also be directly verified.)

This implies

$$f_{\mu_1}(x_1)f_{\mu_0}(x_0) > f_{\mu_1}(x_0)f_{\mu_0}(x_1) \quad x_1 > x_0.$$

Therefore, the inequality is preserved if we integrate both sides with respect

to x_0 on $(-\infty, x)$ for $x < x_1$. This yields:

$$\begin{aligned} \int_{-\infty}^x f_{\mu_1}(x_1)f_{\mu_0}(x_0) dx_0 &> \int_{-\infty}^x f_{\mu_1}(x_0)f_{\mu_0}(x_1) dx_0 && x < x_1 \\ f_{\mu_1}(x_1)F_{\mu_0}(x) &> f_{\mu_0}(x_1)F_{\mu_1}(x) && x < x_1 \end{aligned}$$

Now we integrate both sides with respect to x_1 on (x, ∞) to obtain:

$$(1 - F_{\mu_1}(x))F_{\mu_0}(x) > (1 - F_{\mu_0}(x))F_{\mu_1}(x)$$

which establishes $F_{\mu_0}(x) > F_{\mu_1}(x)$ for all $\mu_1 > \mu_0$. \square

References.

- BENJAMINI, Y. and YEKUTIELI, D. (2005). False discovery rate-adjusted multiple confidence intervals for selected parameters. *Journal of the American Statistical Association*, **100** 71–81.
- BERK, R., BROWN, L., BUJA, A., ZHANG, K. and ZHAO, L. (2013). Valid post-selection inference. *Annals of Statistics*, **41** 802–837.
- EFRON, B., HASTIE, T., JOHNSTONE, I. and TIBSHIRANI, R. (2004). Least angle regression. *Annals of Statistics*, **32** 407–499.
- JAVANMARD, A. and MONTANARI, A. (2013). Confidence intervals and hypothesis testing for high-dimensional regression. *arXiv preprint arXiv:1306.3171*.
- LEHMANN, E. L. and ROMANO, J. P. (2005). *Testing Statistical Hypotheses*. 3rd ed. Springer.
- LOCKHART, R., TAYLOR, J., TIBSHIRANI, R. and TIBSHIRANI, R. (2014). A significance test for the lasso (with discussion). *Annals of Statistics*.
- NEGAHBAN, S. N., RAVIKUMAR, P., WAINWRIGHT, M. J. and YU, B. (2012). A unified framework for high-dimensional analysis of m -estimators with decomposable regularizers. *Statistical Science*, **27** 538–557.
- STOREY, J. D. (2003). The positive false discovery rate: A bayesian interpretation and the q-value. *Annals of statistics* 2013–2035.
- TAYLOR, J., LOCKHART, R., TIBSHIRANI, R. J. and TIBSHIRANI, R. (2014). Post-selection adaptive inference for least angle regression and the lasso. *arXiv preprint arXiv:1401.3889*.
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* 267–288.
- TIBSHIRANI, R. J. (2013). The lasso problem and uniqueness. *Electronic Journal of Statistics*, **7** 1456–1490.
- VAN DE GEER, S., BÜHLMANN, P., RITOV, Y. and DEZEURE, R. (2013). On asymptotically optimal confidence regions and tests for high-dimensional models. *arXiv preprint arXiv:1303.0518*.
- ZHANG, C.-H. and ZHANG, S. (2014). Confidence intervals for low-dimensional parameters with high-dimensional data. *Journal of the Royal Statistical Society: Series B (Methodological)*, **76** 217–242.
- ZOU, H. and HASTIE, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Methodological)*, **67** 301–320.

INSTITUTE FOR COMPUTATIONAL AND MATHE-
MATICAL ENGINEERING
STANFORD UNIVERSITY
STANFORD, CALIFORNIA
E-MAIL: jd117@stanford.edu,
yuekai@stanford.edu

DEPARTMENT OF STATISTICS
STANFORD UNIVERSITY
STANFORD, CALIFORNIA
E-MAIL: dlsun@stanford.edu,
jonathan.taylor@stanford.edu