



Exact Post-Selection Inference for Sequential Regression Procedures

Ryan J. Tibshirani, Jonathan Taylor, Richard Lockhart & Robert Tibshirani

To cite this article: Ryan J. Tibshirani, Jonathan Taylor, Richard Lockhart & Robert Tibshirani (2016) Exact Post-Selection Inference for Sequential Regression Procedures, Journal of the American Statistical Association, 111:514, 600-620, DOI: [10.1080/01621459.2015.1108848](https://doi.org/10.1080/01621459.2015.1108848)

To link to this article: <https://doi.org/10.1080/01621459.2015.1108848>



View supplementary material [↗](#)



Published online: 18 Aug 2016.



Submit your article to this journal [↗](#)



Article views: 1493



View related articles [↗](#)



View Crossmark data [↗](#)



Citing articles: 9 View citing articles [↗](#)

Exact Post-Selection Inference for Sequential Regression Procedures

Ryan J. Tibshirani, Jonathan Taylor, Richard Lockhart, and Robert Tibshirani

ABSTRACT

We propose new inference tools for forward stepwise regression, least angle regression, and the lasso. Assuming a Gaussian model for the observation vector y , we first describe a general scheme to perform valid inference after any selection event that can be characterized as y falling into a polyhedral set. This framework allows us to derive conditional (post-selection) hypothesis tests at any step of forward stepwise or least angle regression, or any step along the lasso regularization path, because, as it turns out, selection events for these procedures can be expressed as polyhedral constraints on y . The p -values associated with these tests are exactly uniform under the null distribution, in finite samples, yielding exact Type I error control. The tests can also be inverted to produce confidence intervals for appropriate underlying regression parameters. The R package `selectiveInference`, freely available on the CRAN repository, implements the new inference tools described in this article. Supplementary materials for this article are available online.

ARTICLE HISTORY

Received October 2014
Revised September 2015

KEYWORDS

Confidence interval; Forward stepwise regression; Inference after selection; Lasso; Least angle regression; p -Value

1. Introduction

Consider observations $y \in \mathbb{R}^n$ drawn from a Gaussian model

$$y = \theta + \epsilon, \quad \epsilon \sim N(0, \sigma^2 I). \quad (1)$$



Given a fixed matrix $X \in \mathbb{R}^{n \times p}$ of predictor variables, our focus is to provide inferential tools for methods that perform variable selection and estimation in an adaptive linear regression of y on X . Unlike much of the related literature on adaptive linear modeling, we *do not assume that the true model is itself linear*, that is, we do not assume that $\theta = X\beta^*$ for a vector of true coefficients $\beta^* \in \mathbb{R}^p$. The particular regression models that we consider in this article are built from sequential procedures that add (or delete) one variable at a time, such as forward stepwise regression (FS), least angle regression (LAR), and the lasso regularization path. However, we stress that the underpinnings of our approach extends well beyond these cases.

To motivate the basic problem and illustrate our proposed solutions, we examine a dataset of 67 observations and 8 variables, where the outcome is the log prostate specific antigen (PSA) level of men who had surgery for prostate cancer. The same dataset was used to motivate the covariance test in Lockhart et al. (2014).¹ The first two numeric columns of Table 1 show the p -values for regression coefficients of variables that enter the model, across steps of FS. The first column shows the results of applying naive, ordinary t -tests to compute the significance of these regression coefficients. We see that the first four variables are apparently significant at the 0.05 level, but this is suspect, as the p -values do not account for the greedy selection of


variables that is inherent to FS. The second column shows our new selection-adjusted p -values for FS, from a *truncated Gaussian (TG) test* developed in Sections 3 and 4. These do properly account for the greediness: they are conditional on the active set at each step, and now just two variables are significant at the 0.05 level.

The last three numeric columns of Table 1 show analogous results for the LAR algorithm applied to the prostate cancer data (the LAR and lasso paths are identical here, as there were no variable deletions). The covariance test (Lockhart et al. 2014), reviewed in the Section 7, measures the improvement in the LAR fit due to adding a predictor at each step, and the third column shows p -values from its $\text{Exp}(1)$ asymptotic null distribution. Our new framework applied to LAR, described in Section 4, produces the results in the rightmost column. We note that this TG test assumes far less than the covariance test. In fact, our TG p -values for both FS and LAR do not require assumptions about the predictors X , or about the true model being linear. They also use a null distribution that is correct in finite samples, rather than asymptotically, under Gaussian errors in (1). The fourth column in the table shows a computationally efficient approximation to the TG test for LAR, that we call the *spacing test*. Later, we establish an asymptotic equivalence between our new spacing test for LAR and the covariance test, and this is supported by the similarity between their p -values in the table.

The R package `selectiveInference` provides an implementation of the TG tests for FS and LAR, and all other inference tools described in this article. This package is available on the CRAN repository, as well as <https://github.com/selective->

CONTACT Ryan J. Tibshirani  ryantibs@cmu.edu  Department of Statistics and Machine Learning Department, 229B Baker Hall, Carnegie Mellon University, Pittsburgh, PA 15213.

Color versions of one or more of the figures in the article can be found online at www.tandfonline.com/r/JASA.

 Supplementary materials for this article are available online. Please go to www.tandfonline.com/r/JASA.

¹The results for the naive FS test and the covariance test differ slightly from those that appear in Lockhart et al. (2014). We use a version of FS that selects variables to maximize the drop in residual sum of squares at each step; Lockhart et al. (2014) used a version based on the maximal absolute correlation of a variable with the residual. Also, our naive FS p -values are one-sided, to match the one-sided nature of the other p -values in the table, whereas Lockhart et al. (2014) used two-sided naive FS p -values. Finally, we use an $\text{Exp}(1)$ limit for the covariance test, and Lockhart et al. (2014) used an F -distribution to account for the unknown variance.

Table 1. Prostate cancer data example: p -values across steps of the forward step-wise (FS) path, computed using naive t -tests that do not account for greedy selection, and our new truncated Gaussian (TG) test for FS; also shown are p -values for the least angle regression (LAR) path, computed using the covariance test of Lockhart et al. (2014), and our new spacing and TG tests for LAR.

| | FS, naive | FS, TG | | LAR, cov | LAR, spacing | LAR, TG |
|---------|-----------|--------|---------|----------|--------------|---------|
| lcavol | 0.000 | 0.000 | lcavol | 0.000 | 0.000 | 0.000 |
| lweight | 0.000 | 0.027 | lweight | 0.047 | 0.052 | 0.052 |
| svi | 0.019 | 0.184 | svi | 0.170 | 0.137 | 0.058 |
| lbph | 0.021 | 0.172 | lbph | 0.930 | 0.918 | 0.918 |
| pgg45 | 0.113 | 0.453 | pgg45 | 0.352 | 0.016 | 0.023 |
| lcp | 0.041 | 0.703 | age | 0.653 | 0.586 | 0.365 |
| age | 0.070 | 0.144 | lcp | 0.046 | 0.060 | 0.800 |
| gleason | 0.442 | 0.800 | gleason | 0.979 | 0.858 | 0.933 |

inference/R-software. A Python implementation is also available at <https://github.com/selective-inference/Python-software>.

A highly nontrivial and important question is to figure out how to combine p -values, such as those in Table 1, to build a rigorous stopping rule, that is, a model selection rule. While we recognize its importance, this topic is *not the focus* of our article. Our focus is to provide a method for computing proper p -values like those in Table 1 in the first place, which we view as a major step in the direction of answering the model selection problem in a practically and theoretically satisfactory manner. Our future work is geared more toward model selection; we also discuss this problem in more detail in Section 2.3.

1.1. Related Work

There is much recent work on inference for high-dimensional regression models. One class of techniques, for example, by Wasserman and Roeder (2009), Meinshausen and Bühlmann (2010), and Minnier, Tian, and Cai (2011) is based on sample-splitting or resampling methods. Another class of approaches, for example, by Zhang and Zhang (2014), Bühlmann (2013), van de Geer et al. (2014), and Javanmard and Montanari (2014a, 2014b) is based on “debiasing” a regularized regression estimator, like the lasso. The inferential targets considered in the aforementioned works are all fixed, and not post-selected, like the targets we study here. As we see it, it is clear (at least conceptually) how to use sample-splitting techniques to accommodate post-selection inferential goals; it is much less clear how to do so with the debiasing tools mentioned above.

Berk et al. (2013) carried out valid post-selection inference (PoSI) by considering all possible model selection procedures that could have produced the given submodel. As the authors state, the inferences are generally conservative for particular selection procedures, but have the advantage that they do not depend on the correctness of the selected submodel. This same advantage is shared by the tests we propose here. Comparisons of our tests, built for specific selection mechanisms, and the PoSI tests, which are much more general, would be interesting to pursue in future work.

Lee et al. (2016), reporting on work concurrent with that of this article, constructed p -values and intervals for lasso coefficients at a fixed value of the regularization parameter λ (instead of a fixed number of steps k along the lasso path, as we consider in Section 4). This article and ours both leverage the same core statistical framework, using truncated Gaussian (TG) distributions, for exact post-selection inference, but differ in

the applications pursued with this framework. After our work was completed, there was further progress on the application and development of exact post-selection inference tools, for example, by Lee and Taylor (2014), Reid, Taylor, and Tibshirani (2014), Loftus and Taylor (2014), Choi, Taylor, and Tibshirani (2014), and Fithian, Sun, and Taylor (2014).

1.2. Notation and Outline

Our notation in the coming sections is as follows. For a matrix $M \in \mathbb{R}^{n \times p}$ and list $S = [s_1, \dots, s_r] \subseteq [1, \dots, p]$, we write $M_S \in \mathbb{R}^{n \times |S|}$ for the submatrix formed by extracting the corresponding columns of M (in the specified order). Similarly for a vector $x \in \mathbb{R}^p$, we write x_S to denote the relevant subvector. We write $(M^T M)^+$ for the (Moore–Penrose) pseudoinverse of the square matrix $M^T M$, and $M^+ = (M^T M)^+ M^T$ for the pseudoinverse of the rectangular matrix M . Finally, we use P_L for the projection operator onto a linear space L .

Here is an outline for the rest of this article. Section 2 gives an overview of our main results. Section 3 describes our general framework for exact conditional inference, with truncated Gaussian (TG) test statistics. Section 4 presents applications of this framework to three sequential regression procedures: FS, LAR, and lasso. Section 5 derives a key approximation to our TG test for LAR, named the *spacing test*, which is considerably simpler (both in terms of form and computational requirements) than its exact counterpart. Section 6 covers empirical examples, and Section 7 draws connections between the spacing and covariance tests. We finish with a discussion in Section 8.

2. Summary of Results

We now summarize our conditional testing framework that yields the p -values demonstrated in the prostate cancer data example, beginning briefly with the general problem setting we consider. Consider testing the hypothesis

$$H_0 : v^T \theta = 0, \quad (2)$$

conditional on having observed $y \in \mathcal{P}$, where \mathcal{P} is a given polyhedral set, and v is a given contrast vector. We derive a test statistic $T(y, \mathcal{P}, v)$ with the property that

$$T(y, \mathcal{P}, v) \stackrel{\mathbb{P}_0}{\sim} \text{Unif}(0, 1), \quad (3)$$

where $\mathbb{P}_0(\cdot) = \mathbb{P}_{v^T \theta = 0}(\cdot | y \in \mathcal{P})$, the probability measure under θ for which $v^T \theta = 0$, conditional on $y \in \mathcal{P}$. The assertion is that $T(y, \mathcal{P}, v)$ is exactly uniform under the null measure, for any finite n and p . This statement assumes nothing about the polyhedron \mathcal{P} , and requires only Gaussian errors in the model (1). As it has a uniform null distribution, the test statistic in (3) serves as its own p -value, and so hereafter we will refer to it in both ways (test statistic and p -value).

Why should we concern ourselves with an event $y \in \mathcal{P}$, for a polyhedron \mathcal{P} ? The short answer: for many regression procedures of interest—in particular, for the sequential algorithms FS, LAR, and lasso—the event that the procedure selects a given model (after a given number of steps) can be represented in this form. For example, consider FS after one step, with $p = 3$ variables total: the FS procedure selects variable 3, and assigns it a

positive coefficient, if and only if

$$\begin{aligned} X_3^T y / \|X_3\|_2 &\geq \pm X_1^T y / \|X_1\|_2, \\ X_3^T y / \|X_3\|_2 &\geq \pm X_2^T y / \|X_2\|_2. \end{aligned}$$

With X considered fixed, these inequalities can be compactly represented as $\Gamma y \geq 0$, where the inequality is meant to be interpreted componentwise, and $\Gamma \in \mathbb{R}^{4 \times n}$ is a matrix with rows $X_3/\|X_3\|_2 \pm X_1/\|X_1\|_2$, $X_3/\|X_3\|_2 \pm X_2/\|X_2\|_2$. Hence if $\hat{j}_1(y)$ and $\hat{s}_1(y)$ denote the variable and sign selected by FS at the first step, then we have shown that

$$\{y : \hat{j}_1(y) = 3, \hat{s}_1(y) = 1\} = \{y : \Gamma y \geq 0\},$$

for a particular matrix Γ . The right-hand side above is clearly a polyhedron (in fact, it is a cone). To test the significance of the third variable, conditional on it being selected at the first step of FS, we consider the null hypothesis H_0 as in (2), with $v = X_3$, and $\mathcal{P} = \{y : \Gamma y \geq 0\}$. The test statistic that we construct in (3) is conditionally uniform under the null. This can be reexpressed as

$$\mathbb{P}_{X_3^T \theta = 0} \left(T_1 \leq \alpha \mid \hat{j}_1(y) = 3, \hat{s}_1(y) = 1 \right) = \alpha, \quad (4)$$

for all $0 \leq \alpha \leq 1$. The conditioning in (4) is important because it properly accounts for the adaptive (i.e., greedy) nature of FS. Loosely speaking, it measures the magnitude of the linear function $X_3^T y$ —not among all y marginally—but among the vectors y that would result in FS selecting variable 3, and assigning it a positive coefficient.

A similar construction holds for a general step k of FS: letting $\hat{A}_k(y) = [\hat{j}_1(y), \dots, \hat{j}_k(y)]$ denote the active list after k steps (so that FS selects these variables in this order) and $\hat{s}_{A_k}(y) = [\hat{s}_1(y), \dots, \hat{s}_k(y)]$ denote the signs of the corresponding coefficients, we have, for any fixed A_k and s_{A_k} ,

$$\{y : \hat{A}_k(y) = A_k, \hat{s}_{A_k}(y) = s_{A_k}\} = \{y : \Gamma y \geq 0\},$$

for another matrix Γ . With $v = (X_{A_k}^+)^T e_k$, where e_k is the k th standard basis vector, the hypothesis in (2) is $e_k^T X_{A_k}^+ \theta = 0$, that is, it specifies that the last partial regression coefficient is not significant, in a projected linear model of θ on X_{A_k} . For $\mathcal{P} = \{y : \Gamma y \geq 0\}$, the test statistic in (3) has the property

$$\mathbb{P}_{e_k^T X_{A_k}^+ \theta = 0} \left(T_k \leq \alpha \mid \hat{A}_k(y) = A_k, \hat{s}_{A_k}(y) = s_{A_k} \right) = \alpha, \quad (5)$$

for all $0 \leq \alpha \leq 1$. We emphasize that the p -value in (5) is exactly (conditionally) uniform under the null, in finite samples. This is true without placing any restrictions on X (besides a general position assumption), and notably, without assuming linearity of the underlying model (i.e., without assuming $\theta = X\beta^*$). Further, though we described the case for FS here, essentially the same story holds for LAR and lasso. The TG p -values for FS and LAR in Table 1 correspond to tests of hypotheses as in (5), that is, tests of $e_k^T X_{A_k}^+ \theta = 0$, over steps of these procedures.

An important point to keep in mind throughout is that our testing framework for the sequential FS, LAR, and lasso procedures is not specific to the choice $v = (X_{A_k}^+)^T e_k$, and allows for the testing of arbitrary linear contrasts $v^T \theta$ (as long as v is fixed by the conditioning event). For concreteness, we will pay close attention to the case $v = (X_{A_k}^+)^T e_k$, since it gives us a test for the

significance of variables as they enter the model, but many other choices of v could be interesting and useful.

2.1. Conditional Confidence Intervals

A strength of our framework is that our test statistics can be inverted to make coverage statements about arbitrary linear contrasts of θ . In particular, consider the hypothesis test defined by $v = (X_{A_k}^+)^T e_k$, for the k th step of FS (similar results apply to LAR and lasso). By inverting our test statistic in (5), we obtain a conditional confidence interval I_k satisfying

$$\mathbb{P} \left(e_k^T X_{A_k}^+ \theta \in I_k \mid \hat{A}_k(y) = A_k, \hat{s}_{A_k}(y) = s_{A_k} \right) = 1 - \alpha. \quad (6)$$

In words, the random interval I_k traps with probability $1 - \alpha$ the coefficient of the last selected variable, in a regression model that projects θ onto X_{A_k} , conditional on FS having selected variables A_k with signs s_{A_k} , after k steps of the algorithm. As (6) is true conditional on $\Gamma y \geq 0$, we can also marginalize this statement to yield

$$\mathbb{P} \left(e_k^T X_{\hat{A}_k}^+ \theta \in I_k \right) = 1 - \alpha. \quad (7)$$

Note that $\hat{A}_k = \hat{A}_k(y)$ denotes the random active list after k FS steps. Written in the unconditional form (7), we call I_k a *selection interval* for the random quantity $e_k^T X_{\hat{A}_k}^+ \theta$. We use this name to emphasize the difference in interpretation here, versus the conditional case: the selection interval covers a *moving target*, as both the identity of the k th-selected variable, and the identities of all the previously selected variables (which play a role in the k th partial regression coefficient of θ on $X_{\hat{A}_k}$), are random—they depend on y .

We have seen that our intervals can be interpreted conditionally, as in (6), or unconditionally, as in (7). The former is perhaps more aligned with the spirit of post-selection inference, as it guarantees coverage, conditional on the output of our selection procedure. But the latter interpretation is also interesting, and in a way, cleaner. From the unconditional point of view, we can roughly think of the selection interval I_k as covering the project population coefficient of the “ k th most important variable” as deemed by the sequential regression procedure at hand (FS, LAR, or lasso). Figure 1 displays 90% confidence intervals at each step of FS, run on the prostate cancer dataset discussed in the introduction.

2.2. Marginalization

Similar to the formation of selection intervals in the last subsection, we note that any amount of coarsening, that is, marginalization, of the conditioning set in (5) results in a valid interpretation for p -values. For example, by marginalizing over all possible sign lists s_{A_k} associated with A_k , we obtain

$$\mathbb{P}_{e_k^T X_{A_k}^+ \theta = 0} \left(T_k \leq \alpha \mid \hat{A}_k(y) = A_k \right) = \alpha,$$

so that the conditioning event only encodes the observed active list, and not the observed signs. Thus, we have another possible interpretation for the statistic (p -value) T_k : under the null measure, which conditions on FS having selected the variables A_k

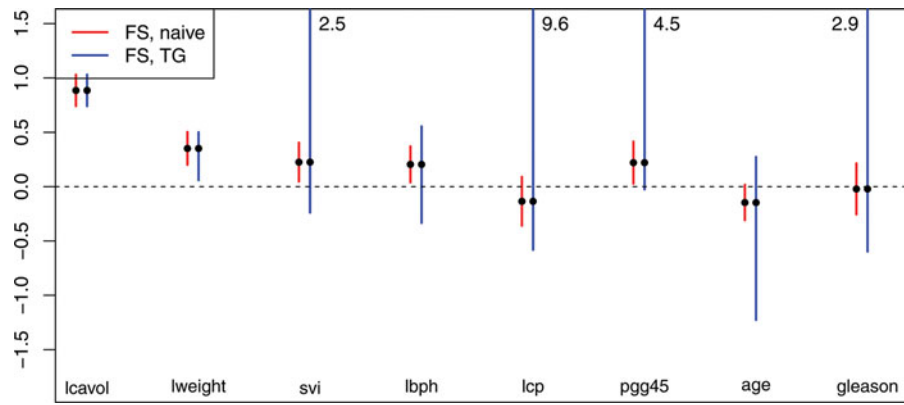


Figure 1. Prostate cancer data example: 90% naive confidence intervals and conditional confidence intervals (or, selection intervals) computed using the TG (truncated Gaussian) statistics, for FS (forward stepwise). Black dots denote the estimated partial regression coefficients for the variable to enter, in a regression on the active submodel. The upper confidence limits for some parameters exceed the range for the y-axis on the plot, and their actual values marked at the appropriate places.

(regardless of their signs), T_k is uniformly distributed. The idea of marginalization will be important when we discuss details of the constructed tests for LAR and lasso.

2.3. Model Selection

How can the inference tools of this article be translated into rigorous rules for model selection? This is of course an important (and difficult) question, and we do not yet possess a complete understanding of the model selection problem, though it is the topic of future work. Below we describe three possible strategies for model selection, using the p -values that come from our inference framework. We do not have extensive theory to explain or evaluate them, but all are implemented in the R package `selectiveInference`.

- *Inference from sequential p -values.* We have advocated the idea of computing p -values across steps of the regression procedure at hand, as exemplified in Table 1. Here, at each step k , the p -value tests $e_k^T X_{A_k}^+ \theta = 0$, that is, tests the significance of the variable to enter the active set A_k , in a projected linear model of the mean θ on the variables in A_k . G'Sell et al. (2016) proposed sequential stopping rules using such p -values, including the “Forward-Stop” rule, which guarantees false discovery rate (FDR) control at a given level. For example, the ForwardStop rule at a nominal 10% FDR level, applied to the TG p -values from the LAR path for the prostate cancer data (the last column of Table 1), yields a model with three predictors. However, it should be noted that the guarantee for FDR control for ForwardStop in G'Sell et al. (2016) assumes that the p -values are independent, and this is not true for the p -values from our inference framework.
- *Inference at a fixed step k .* Instead of looking at p -values across steps, we could instead fix a step k , and inspect the p -values corresponding to the hypotheses $e_j^T X_{A_j}^+ \theta = 0$, for $j = 1, \dots, k$. This tests the significance of every variable, among the rest in the discovered active set A_k , and it still fits within our developed framework: we are just using different linear contrasts $v = (X_{A_j}^+)^T e_j$ of the mean θ , for $j = 1, \dots, k$. The results of these tests are genuinely different, in terms of their statistical meaning,

than the results from testing variables as they enter the model (since the active set changes at each step). Given the p -values corresponding to all active variables at a given step k , we could, for example, perform a Bonferroni correction, and declare significance at the level α/k , to select a model (a subset of A_k) with Type I error controlled at the level α . For example, when we apply this strategy at step $k = 5$ of the LAR path for the prostate cancer data, and examine Bonferroni-corrected p -values at the 0.05 level, only two predictors (lweight and pgg45) end up being significant.

- *Inference at an adaptively selected step k .* Finally, the above scheme for inference could be conducted with a step number k that is adaptively selected, instead of fixed ahead of time, provided the selection event that determines k is a polyhedral set in y . A specific example of this is an Akaike information criterion (AIC)-style rule, which chooses the step k after which the AIC criterion rises, say, twice in a row. We omit the details, but verifying that such a stopping rule defines a polyhedral constraint for y is straightforward (it follows essentially the same logic as the arguments that show the FS selection event is itself polyhedral, which are given in Section 4.1). Hence, by including all the necessary polyhedral constraints—those that determine k , and those that subsequently determine the selected model—we can compute p -values for each of the active variables at an adaptively selected step k , using the inference tools derived in this article. When this method is applied to the prostate cancer dataset, the AIC-style rule (which stops once it sees two consecutive rises in the AIC criterion) chooses $k = 4$. Examining Bonferroni corrected p -values at step $k = 4$, only one predictor (lweight) remains significant at the 0.05 level.

3. Conditional Gaussian Inference After Polyhedral Selection

In this section, we present a few key results on Gaussian contrasts conditional on polyhedral events, which provide a basis for the methods proposed in this article. The same core development appears in Lee et al. (2016); for brevity, we refer the reader to the latter article for formal proofs. We assume $y \sim N(\theta, \Sigma)$,

where $\theta \in \mathbb{R}^n$ is unknown, but $\Sigma \in \mathbb{R}^{n \times n}$ is known. This generalizes our setup in (1) (allowing for a general error covariance matrix). We also consider a generic polyhedron $\mathcal{P} = \{y : \Gamma y \geq u\}$, where $\Gamma \in \mathbb{R}^{m \times n}$ and $u \in \mathbb{R}^m$ are fixed, and the inequality is to be interpreted componentwise. For a fixed $v \in \mathbb{R}^n$, our goal is to make inferences about $v^T \theta$ conditional on $y \in \mathcal{P}$. Next, we provide a helpful alternate representation for \mathcal{P} .

Lemma 1 (Polyhedral selection as truncation). For any Σ, v such that $v^T \Sigma v \neq 0$,

$$\Gamma y \geq u \iff \mathcal{V}^{\text{lo}}(y) \leq v^T y \leq \mathcal{V}^{\text{up}}(y), \mathcal{V}^0(y) \leq 0, \quad (8)$$

where

$$\mathcal{V}^{\text{lo}}(y) = \max_{j: \rho_j > 0} \frac{u_j - (\Gamma y)_j + \rho_j v^T y}{\rho_j}, \quad (9)$$

$$\mathcal{V}^{\text{up}}(y) = \min_{j: \rho_j < 0} \frac{u_j - (\Gamma y)_j + \rho_j v^T y}{\rho_j}, \quad (10)$$

$$\mathcal{V}^0(y) = \max_{j: \rho_j = 0} u_j - (\Gamma y)_j, \quad (11)$$

and $\rho = \Gamma \Sigma v / v^T \Sigma v$. Moreover, the triplet $(\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}, \mathcal{V}^0)(y)$ is independent of $v^T y$.

Remark 1. The result in (8), with $\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}, \mathcal{V}^0$ defined as in (9)–(11), is a deterministic result that holds for all y . Only the last independence result depends on normality of y .

See Figure 2 for a geometric illustration of this lemma. Intuitively, we can explain the result as follows, assuming for simplicity (and without a loss of generality) that $\Sigma = I$. We first decompose $y = P_v y + P_{v^\perp} y$, where $P_v y = v v^T y / \|v\|_2^2$ is the projection of y along v , and $P_{v^\perp} y = y - P_v y$ is the projection onto the orthocomplement of v . Accordingly, we view y as a deviation from $P_{v^\perp} y$, of an amount $v^T y$, along the line determined by v . The quantities \mathcal{V}^{lo} and \mathcal{V}^{up} describe how far we can deviate on either side of $P_{v^\perp} y$, before y leaves the polyhedron. This gives rise to the inequality $\mathcal{V}^{\text{lo}} \leq v^T y \leq \mathcal{V}^{\text{up}}$. Some faces of the polyhedron, however, may be perfectly aligned with v (i.e., their normal vectors may be orthogonal to v), and \mathcal{V}^0 accounts for this by checking that y lies on the correct side of these faces.

From Lemma 1, the distribution of any linear function $v^T y$, conditional on the selection $\Gamma y \geq u$, can be written as the conditional distribution

$$v^T y \mid \mathcal{V}^{\text{lo}}(y) \leq v^T y \leq \mathcal{V}^{\text{up}}(y), \mathcal{V}^0(y) \leq 0. \quad (12)$$

Since $v^T y$ has a Gaussian distribution, the above is a truncated Gaussian distribution (with random truncation limits). A simple transformation leads to a pivotal statistic, which will be critical for inference about $v^T \theta$.

Lemma 2 (Pivotal statistic after polyhedral selection). Let $\Phi(x)$ denote the standard normal cumulative distribution function (CDF), and let $F_{\mu, \sigma^2}^{[a, b]}$ denote the CDF of an $N(\mu, \sigma^2)$ random variable truncated to lie in $[a, b]$, that is,

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

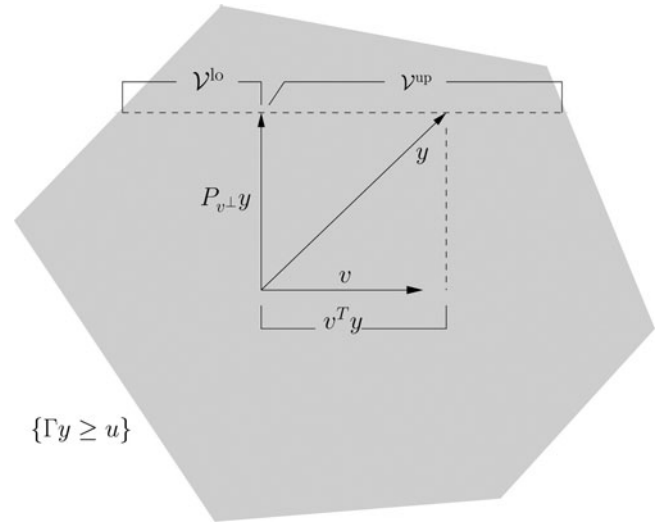


Figure 2. Geometry of polyhedral selection as truncation. For simplicity, we assume that $\Sigma = I$ (otherwise standardize as appropriate). The shaded gray area is the polyhedral set $\{y : \Gamma y \geq u\}$. By breaking up y into its projection onto v and its projection onto the orthogonal complement of v , we see that $\Gamma y \geq u$ holds if and only if $v^T y$ does not deviate too far from $P_{v^\perp} y$, hence trapping it in between bounds $\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}$. Furthermore, these bounds $\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}$ are functions of $P_{v^\perp} y$ alone, so under normality, they are independent of $v^T y$.

For $v^T \Sigma v \neq 0$, the statistic $F_{v^T \theta, v^T \Sigma v}^{[\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}]}(v^T y)$ is a pivotal quantity conditional on $\Gamma y \geq u$:

$$\mathbb{P}\left(F_{v^T \theta, v^T \Sigma v}^{[\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}]}(v^T y) \leq \alpha \mid \Gamma y \geq u\right) = \alpha, \quad (13)$$

for any $0 \leq \alpha \leq 1$, where $\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}$ are as defined in (9), (10).

Remark 2. A referee of this article astutely noted the connection between Lemma 2 and classic results on inference in an exponential family model (e.g., Chapter 4 of Lehmann and Romano 2005), in the presence of nuisance parameters. The analogy is in a rotated coordinate system, the parameter of interest is $v^T \theta$, and the nuisance parameters correspond to $P_{v^\perp} \theta$. This connection is developed in Fithian, Sun, and Taylor (2014).

The pivotal statistic in the lemma leads to valid conditional p -values for testing the null hypothesis $H_0 : v^T \theta = 0$, and correspondingly, conditional confidence intervals for $v^T \theta$. We divide our presentation into two parts, on one-sided and two-sided inference.

3.1. One-Sided Conditional Inference

The result below is a direct consequence of the pivot in Lemma 2.

Lemma 3 (One-sided conditional inference after polyhedral selection). Given $v^T \Sigma v \neq 0$, suppose that we are interested in testing

$$H_0 : v^T \theta = 0 \quad \text{against} \quad H_1 : v^T \theta > 0.$$

Define the test statistic

$$T = 1 - F_{0, v^T \Sigma v}^{[\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}]}(v^T y), \quad (14)$$

where we use the notation of Lemma 2 for the truncated normal CDF. Then T is a valid p -value for H_0 , conditional on $\Gamma y \geq u$:

$$\mathbb{P}_{v^T \theta = 0}(T \leq \alpha \mid \Gamma y \geq u) = \alpha, \quad (15)$$

for any $0 \leq \alpha \leq 1$. Further, define δ_α to satisfy

$$1 - F_{\delta_\alpha, v^T \Sigma v}^{[\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}]}(v^T y) = \alpha. \quad (16)$$

Then $I = [\delta_\alpha, \infty)$ is a valid one-sided confidence interval for $v^T \theta$, conditional on $\Gamma y \geq u$:

$$\mathbb{P}(v^T \theta \geq \delta_\alpha \mid \Gamma y \geq u) = 1 - \alpha. \quad (17)$$

Note that by defining our test statistic in terms of the conditional survival function, as in (14), we are implicitly aligning ourselves to have power against the one-sided alternative $H_1 : v^T \theta > 0$. This is because the truncated normal survival function $1 - F_{\mu, \sigma^2}^{[a, b]}(x)$, evaluated at any fixed point x , is monotone increasing in μ . The same fact (monotonicity of the survival function in μ) validates the coverage of the constructed confidence interval in (16) and (17).

3.2. Two-Sided Conditional Inference

For a two-sided alternative, we use a simple modification of the one-sided test in Lemma 3.

Lemma 4 (Two-sided conditional inference after polyhedral selection). Given $v^T \Sigma v \neq 0$, suppose that we are interested in testing

$$H_0 : v^T \theta = 0 \quad \text{against} \quad H_1 : v^T \theta \neq 0.$$

Define the test statistic

$$T = 2 \cdot \min \left\{ F_{0, v^T \Sigma v}^{[\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}]}(v^T y), 1 - F_{0, v^T \Sigma v}^{[\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}]}(v^T y) \right\}, \quad (18)$$

where we use the notation of Lemma 2 for the truncated normal CDF. Then T is a valid p -value for H_0 , conditional on $\Gamma y \geq u$:

$$\mathbb{P}_{v^T \theta = 0}(T \leq \alpha \mid \Gamma y \geq u) = \alpha, \quad (19)$$

for any $0 \leq \alpha \leq 1$. Further, define $\delta_{\alpha/2}, \delta_{1-\alpha/2}$ to satisfy

$$1 - F_{\delta_{\alpha/2}, v^T \Sigma v}^{[\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}]}(v^T y) = \alpha/2, \quad (20)$$

$$1 - F_{\delta_{1-\alpha/2}, v^T \Sigma v}^{[\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}]}(v^T y) = 1 - \alpha/2. \quad (21)$$

Then

$$\mathbb{P}(\delta_{\alpha/2} \leq v^T \theta \leq \delta_{1-\alpha/2} \mid \Gamma y \geq u) = 1 - \alpha. \quad (22)$$

The test statistic in (18), defined in terms of the minimum of the truncated normal CDF and survival function, has power against the two-sided alternative $H_1 : v^T \theta \neq 0$. The proof of its null distribution in (19) follows from the simple fact that if U is a standard uniform random variable, then so is $2 \cdot \min\{U, 1 - U\}$. The construction of the confidence interval in (20), (21), (22) again uses the monotonicity of the truncated normal survival function in the underlying mean parameter.

4. Exact Selection-Adjusted Tests for FS, LAR, LASSO

Here, we apply the tools of Section 3 to the case of selection in regression using the forward stepwise (FS), least angle regression (LAR), or lasso procedures. We assume that the columns of X are in general position. This means that for any $k < \min\{n, p\}$, any subset of columns X_{j_1}, \dots, X_{j_k} , and any signs $\sigma_1, \dots, \sigma_k \in \{-1, 1\}$, the affine span of $\sigma_1 X_{j_1}, \dots, \sigma_k X_{j_k}$ does not contain any of the remaining columns, up to a sign flip (i.e., does not contain any of $\pm X_j, j \neq j_1, \dots, j_k$). One can check that this implies the sequence of FS estimates is unique. It also implies that the LAR and lasso paths of estimates are uniquely determined (Tibshirani 2013). The general position assumption is not at all stringent, for example, if the columns of X are drawn according to a continuous probability distribution, then they are in general position almost surely.

Next, we show that the model selection events for FS, LAR, and lasso can be characterized as polyhedra (indeed, cones) of the form $\{y : \Gamma y \geq 0\}$. After this, we describe the forms of the exact conditional tests and intervals, as provided by Lemmas 1–4, for these procedures, and discuss some important practical issues.

4.1. Polyhedral Sets for FS Selection Events

Recall that FS repeatedly adds the predictor to the current active model that most improves the fit. After each addition, the active coefficients are recomputed by least-square regression on the active predictors. This process ends when all predictors are in the model, or when the residual error is zero.²

Formally, suppose that $A_k = [j_1, \dots, j_k]$ is the list of active variables selected by FS after k steps, and $s_{A_k} = [s_1, \dots, s_k]$ denotes their signs upon entering. That is, at each step k , the variable j_k and sign s_k satisfy

$$\text{RSS}(y, X_{[j_1, \dots, j_{k-1}, j_k]}) \leq \text{RSS}(y, X_{[j_1, \dots, j_{k-1}, j]}) \quad \text{for all } j \neq j_1, \dots, j_k, \\ \text{and } s_k = \text{sign}(e_k^T (X_{[j_1, \dots, j_k]}^+)^+ y),$$

where $\text{RSS}(y, X_S)$ denotes the residual sum of squares from regressing y onto X_S , for a list of variables S .

The set of all observations vectors y that give active list A_k and sign list s_{A_k} over k steps, denoted

$$\mathcal{P} = \left\{ y : \hat{A}_k(y) = A_k, \hat{s}_{A_k}(y) = s_{A_k} \right\}, \quad (23)$$

is indeed a polyhedron of the form $\mathcal{P} = \{y : \Gamma y \geq 0\}$. The proof of this fact uses induction. The case when $k = 1$ can be seen directly by inspection, as j_1 and s_1 are the variable and sign to be chosen by FS if and only if

$$\left\| \left(I - X_{j_1} X_{j_1}^T / \|X_{j_1}\|_2^2 \right) y \right\|_2^2 \leq \left\| \left(I - X_j X_j^T / \|X_j\|_2^2 \right) y \right\|_2^2 \\ \text{for all } j \neq j_1, \text{ and} \\ s_1 = \text{sign}(X_{j_1}^T y),$$

² Slightly different versions of FS result from different ways of defining the notion of the predictor that “most improves the fit.” Our definition based on RSS drop is equivalent to choosing the variable that, once orthogonalized with respect to the current active model, achieves the largest absolute correlation with the residual. Other common versions of FS—e.g., choosing the variable that achieves the largest absolute correlation with the residual (without orthogonalization)—will clearly also fit into our polyhedral framework.

which is equivalent to

$$s_1 X_{j_1}^T y / \|X_{j_1}\|_2 \geq \pm X_j^T y / \|X_j\|_2 \quad \text{for all } j \neq j_1.$$

Thus, the matrix Γ begins with $2(p-1)$ rows of the form $s_1 X_{j_1} / \|X_{j_1}\|_2 \pm X_j / \|X_j\|_2$, for $j \neq j_1$. Now assume the statement is true for $k-1$ steps. At step k , the optimality conditions for j_k, s_k can be expressed as

$$\begin{aligned} \left\| \left(I - \tilde{X}_{j_k} \tilde{X}_{j_k}^T / \|\tilde{X}_{j_k}\|_2^2 \right) r \right\|_2^2 &\leq \left\| \left(I - \tilde{X}_j \tilde{X}_j^T / \|\tilde{X}_j\|_2^2 \right) r \right\|_2^2 \\ &\quad \text{for all } j \neq j_1, \dots, j_k, \text{ and} \\ s_k &= \text{sign}(\tilde{X}_{j_k}^T r), \end{aligned}$$

where \tilde{X}_j denotes the residual from regressing X_j onto $X_{A_{k-1}}$, and r the residual from regressing y onto $X_{A_{k-1}}$. As in the $k=1$ case, the above is equivalent to

$$s_k \tilde{X}_{j_k}^T r / \|\tilde{X}_{j_k}\|_2 \geq \pm \tilde{X}_j^T r / \|\tilde{X}_j\|_2 \quad \text{for all } j \neq j_1, \dots, j_k,$$

or

$$s_k X_{j_k}^T P_{A_{k-1}}^\perp y / \|P_{A_{k-1}}^\perp X_{j_k}\|_2 \geq \pm X_j^T P_{A_{k-1}}^\perp y / \|P_{A_{k-1}}^\perp X_j\|_2 \quad \text{for all } j \neq j_1, \dots, j_k,$$

where $P_{A_{k-1}}^\perp$ denotes the projection orthogonal to the column space of $X_{A_{k-1}}$. Hence, we append $2(p-k)$ rows to Γ , of the form $P_{A_{k-1}}^\perp (s_k X_{j_k} / \|P_{A_{k-1}}^\perp X_{j_k}\|_2 \pm X_j / \|P_{A_{k-1}}^\perp X_j\|_2)$, for $j \neq j_1, \dots, j_k$. In summary, after k steps, the polyhedral set for the FS selection event (23) corresponds to a matrix Γ with $2pk - k^2 - k$ rows.³

4.2. Polyhedral Sets for LAR Selection Events

The LAR algorithm (Efron et al. 2004) is an iterative method, like FS, that produces a sequence of nested regression models. As before, we keep a list of active variables and signs across steps of the algorithm. Here is a concise description of the LAR steps. At step $k=1$, we initialize the active variable and sign list with $A = [j_1]$ and $s_{A_1} = [s_1]$, where j_1, s_1 satisfy

$$(j_1, s_1) = \underset{j=1, \dots, p, s \in \{-1, 1\}}{\operatorname{argmax}} s X_j^T y. \quad (24)$$

(This is the same selection as made by FS at the first step, provided that X has columns with unit norm.) We also record the first knot

$$\lambda_1 = s_1 X_{j_1}^T y. \quad (25)$$

For a general step $k > 1$, we form the list A_k by appending j_k to A_{k-1} , and form s_{A_k} by appending s_k to $s_{A_{k-1}}$, where j_k, s_k satisfy

$$\begin{aligned} (j_k, s_k) &= \underset{j \notin A_{k-1}, s \in \{-1, 1\}}{\operatorname{argmax}} \frac{X_j^T P_{A_{k-1}}^\perp y}{s - X_j^T (X_{A_{k-1}}^+)^T s_{A_{k-1}}} \\ &\quad \cdot 1 \left\{ \frac{X_j^T P_{A_{k-1}}^\perp y}{s - X_j^T (X_{A_{k-1}}^+)^T s_{A_{k-1}}} \leq \lambda_{k-1} \right\}. \quad (26) \end{aligned}$$

³ We have been implicitly assuming thus far that $k < p$. If $k = p$ (so that necessarily $p \leq n$), then we must add an “extra” row to Γ , this row being $P_{A_{p-1}}^\perp s_p X_{j_p}$, which encodes the sign constraint $s_p X_{j_p}^T P_{A_{p-1}}^\perp y \geq 0$. For $k < p$, this constraint is implicitly encoded due to the constraints of the form $s_k X_{j_k}^T P_{A_{k-1}}^\perp y \geq \pm a$ for some a .

Above, $P_{A_{k-1}}^\perp$ is the projection orthogonal to the column space of $X_{A_{k-1}}$, $1\{\cdot\}$ denotes the indicator function, and λ_{k-1} is the knot value from step $k-1$. We also record the k th knot

$$\lambda_k = \frac{X_{j_k}^T P_{A_{k-1}}^\perp y}{s_k - X_{j_k}^T (X_{A_{k-1}}^+)^T s_{A_{k-1}}}. \quad (27)$$

The algorithm terminates after the k -step model if $k = p$, or if $\lambda_{k+1} < 0$.

LAR is often viewed as “less greedy” than FS. It is also intimately tied to the lasso, as covered in the next subsection. Now, we verify that the LAR selection event

$$\mathcal{P} = \left\{ y : \hat{A}_k(y) = A_k, \hat{s}_{A_k}(y) = s_{A_k}, \hat{S}_\ell(y) = S_\ell, \ell = 1, \dots, k \right\} \quad (28)$$

is a polyhedron of the form $\mathcal{P} = \{y : \Gamma y \geq 0\}$. We can see that the LAR event in (28) contains “extra” conditioning, $\hat{S}_\ell(y) = S_\ell$, $\ell = 1, \dots, k$, when compared to the FS event in (23). Explained in words, $S_\ell \subseteq \{1, \dots, p\} \times \{-1, 1\}$ contains the variable-sign pairs that were “in competition” to become the active variable-sign pair step ℓ . A subtlety of LAR: it is not always the case that $S_\ell = A_{\ell-1}^c \times \{-1, 1\}$, since some variable-sign pairs are automatically excluded from consideration, as they would have produced a knot value that is too large (larger than the previous knot $\lambda_{\ell-1}$). This is reflected by the indicator function in (26). The characterization in (28) is still *perfectly viable for inference*, because any conditional statement over \mathcal{P} in (28) translates into a valid one without conditioning on $\hat{S}_\ell(y)$, $\ell = 1, \dots, k$, by marginalizing over all possible realizations S_ℓ , $\ell = 1, \dots, k$. (Recall the discussion of marginalization in Section 2.2.)

The polyhedral representation for \mathcal{P} in (28) again proceeds by induction. Starting with $k=1$, we can express the optimality of j_1, s_1 in (24) as

$$c(j_1, s_1)^T y \geq c(j, s)^T y, \quad \text{for all } j \neq j_1, s \in \{-1, 1\},$$

where $c(j, s) = s X_j$. Thus Γ has $2(p-1)$ rows, of the form $c(j_1, s_1) - c(j, s)$ for $j \neq j_1, s \in \{-1, 1\}$. (In the first step, $S_1 = \{1, \dots, p\} \times \{-1, 1\}$, and we do not require extra rows of Γ to explicitly represent it.) Further, suppose that the selection set can be represented in the desired manner, after $k-1$ steps. Then the optimality of j_k, s_k in (26) can be expressed as

$$\begin{aligned} c(j_k, s_k, A_{k-1}, s_{A_{k-1}})^T y &\geq c(j, s, A_{k-1}, s_{A_{k-1}})^T y \\ &\quad \text{for all } (j, s) \in S_k \setminus \{(j_k, s_k)\}, \\ c(j_k, s_k, A_{k-1}, s_{A_{k-1}})^T y &\geq 0, \end{aligned}$$

where $c(j, s, A_{k-1}, s_{A_{k-1}}) = (P_{A_{k-1}}^\perp X_j) / (s - X_j^T (X_{A_{k-1}}^+)^T s_{A_{k-1}})$. The set S_k is characterized by

$$\begin{aligned} c(j, s, A_{k-1}, s_{A_{k-1}})^T y &\leq \lambda_{k-1} \quad \text{for } (j, s) \in S_k, \\ c(j, s, A_{k-1}, s_{A_{k-1}})^T y &\geq \lambda_{k-1} \quad \text{for } (j, s) \in (A_{k-1}^c \times \{-1, 1\}) \setminus S_k. \end{aligned}$$

Notice that $\lambda_{k-1} = c(j_{k-1}, s_{k-1}, A_{k-2}, s_{A_{k-2}})^T y$ is itself a linear function of y , by the inductive hypothesis. Therefore, the new Γ matrix is created by appending the following $|S_k| + 2(p-k+1)$ rows to the previous matrix: $c(j_k, s_k, A_{k-1}, s_{A_{k-1}}) - c(j, s, A_{k-1}, s_{A_{k-1}})$, for $(j, s) \in S_k \setminus \{(j_k, s_k)\}$; $c(j_k, s_k, A_{k-1}, s_{A_{k-1}})$; $c(j_{k-1}, s_{k-1}, A_{k-2}, s_{A_{k-2}}) - c(j, s, A_{k-1}, s_{A_{k-1}})$, for $(j, s) \in S_k$; $c(j, s, A_{k-1}, s_{A_{k-1}}) - c(j_{k-1}, s_{k-1}, A_{k-2}, s_{A_{k-2}})$ for $(j, s) \in (A_{k-1}^c \times \{-1, 1\}) \setminus S_k$. In total, the number of rows of Γ at

step k of LAR is bounded above by $\sum_{\ell=1}^k (|S_\ell| + 2(p - \ell + 1)) \leq 3pk - 3k^2/2 + 3k/2$.

4.3. Polyhedral Sets for Lasso Selection Events

By introducing a step into the LAR algorithm that deletes variables from the active set if their coefficients pass through zero, the modified LAR algorithm traces out the lasso regularization path (Efron et al. 2004). To concisely describe this modification, at a step $k > 1$, denote by $(j_k^{\text{add}}, s_k^{\text{add}})$ the variable-sign pair to enter the model next, as defined in (26), and denote by λ_k^{add} the value of λ at which they would enter, as defined in (27). Now define

$$j_k^{\text{del}} = \underset{j \in A_{k-1} \setminus \{j_{k-1}\}}{\operatorname{argmax}} \frac{e_j^T X_{A_{k-1}}^+ y}{e_j^T (X_{A_{k-1}}^T X_{A_{k-1}})^{-1} s_{A_{k-1}}} \cdot \mathbf{1} \left\{ \frac{e_j^T X_{A_{k-1}}^+ y}{e_j^T (X_{A_{k-1}}^T X_{A_{k-1}})^{-1} s_{A_{k-1}}} \leq \lambda_{k-1} \right\}, \quad (29)$$

the variable to leave the model next, and

$$\lambda_k^{\text{del}} = \frac{e_{j_k^{\text{del}}}^T X_{A_{k-1}}^+ y}{e_{j_k^{\text{del}}}^T (X_{A_{k-1}}^T X_{A_{k-1}})^{-1} s_{A_{k-1}}}, \quad (30)$$

the value of λ at which it would leave. The lasso regularization path is given by executing whichever action—variable entry, or variable deletion—happens first, when seen from the perspective of decreasing λ . That is, we record the k th knot $\lambda_k = \max\{\lambda_k^{\text{add}}, \lambda_k^{\text{del}}\}$, and we form A_k, s_{A_k} by either adding $j_k^{\text{add}}, s_k^{\text{add}}$ to $A_{k-1}, s_{A_{k-1}}$ if $\lambda_k = \lambda_k^{\text{add}}$, or by deleting j_k^{del} from A_{k-1} and its sign from $s_{A_{k-1}}$ if $\lambda_k = \lambda_k^{\text{del}}$.

We show that the lasso selection event⁴

$$\mathcal{P} = \left\{ y : \hat{A}_\ell(y) = A_\ell, \hat{s}_{A_\ell}(y) = s_{A_\ell}, \hat{S}_\ell^{\text{add}}(y) = S_\ell^{\text{add}}, \hat{S}_\ell^{\text{del}}(y) = S_\ell^{\text{del}}, \ell = 1, \dots, k \right\}, \quad (31)$$

can be expressed in polyhedral form $\{y : \Gamma y \geq 0\}$. A difference between (31) and the LAR event in (28) is that, in addition to keeping track of the set S_ℓ^{add} of variable-sign pairs in consideration to become active (to be added) at step ℓ , we must also keep track of the set S_ℓ^{del} of variables in consideration to become inactive (to be deleted) at step ℓ . As discussed earlier, a valid inferential statement conditional on the lasso event \mathcal{P} in (31) is still valid once we ignore the conditioning on $\hat{S}_\ell^{\text{add}}(y), \hat{S}_\ell^{\text{del}}(y), \ell = 1, \dots, k$, by marginalization.

To build the Γ matrix corresponding to (31), we begin the same construction as we laid out for LAR in the last subsection, and simply add more rows. At a step $k > 1$, the rows we described appending to Γ for LAR now merely characterize the variable-sign pair $(j_k^{\text{add}}, s_k^{\text{add}})$ to enter the model next, as well as the set S_k^{del} . To characterize the variable j_k^{del} to leave the model next, we express its optimality in (29) as

$$d(j_k^{\text{del}}, A_{k-1}, s_{A_{k-1}})^T y \geq d(j, A_{k-1}, s_{A_{k-1}})^T y$$

⁴ The observant reader might notice that the selection event for the lasso in (31), compared to that for FS in (23) and LAR in (28), actually enumerates the assignments of active sets $\hat{A}_\ell(y) = A_\ell, \ell = 1, \dots, k$ across all k steps of the path. This is done because, with variable deletions, it is no longer possible to express an entire history of active sets with a single list. The same is true of the active signs.

for all $j \in S_k^{\text{del}} \setminus \{j_k^{\text{del}}\}$,

$$d(j_k^{\text{del}}, A_{k-1}, s_{A_{k-1}})^T y \geq 0,$$

where $d(j, A_{k-1}, s_{A_{k-1}}) = ((X_{A_{k-1}}^+)^T e_j) / (e_j^T (X_{A_{k-1}}^T X_{A_{k-1}})^{-1} s_{A_{k-1}})$, and S_k^{del} is characterized by

$$\begin{aligned} d(j, s, A_{k-1}, s_{A_{k-1}})^T y &\leq \lambda_{k-1} \quad \text{for } (j, s) \in S_k^{\text{del}}, \\ d(j, s, A_{k-1}, s_{A_{k-1}})^T y &\geq \lambda_{k-1} \quad \text{for } (j, s) \in A_{k-1} \setminus S_k^{\text{del}}. \end{aligned}$$

Recall that $\lambda_{k-1} = b_{k-1}^T y$ is a linear function of y , by the inductive hypothesis. If a variable was added at step $k-1$, then $b_{k-1} = c(j_{k-1}, s_{k-1}, A_{k-2}, s_{A_{k-2}})$; if instead a variable was deleted at step $k-1$, then $b_{k-1} = d(j_{k-1}, A_{k-2}, s_{A_{k-2}})$. Finally, we must characterize step k as either witnessing a variable addition or deletion. The former case is represented by

$$c(j_k^{\text{add}}, s_k^{\text{add}}, A_{k-1}, s_{A_{k-1}})^T y \geq d(j_k^{\text{del}}, A_{k-1}, s_{A_{k-1}})^T y,$$

the latter case reverses the above inequality. Hence, in addition to those described in the previous subsection, we append the following $|S_k^{\text{del}}| + |A_{k-1}| + 1$ rows to Γ : $d(j_k^{\text{del}}, A_{k-1}, s_{A_{k-1}}) - d(j, A_{k-1}, s_{A_{k-1}})$ for $(j, s) \in S_k^{\text{del}} \setminus \{j_k^{\text{del}}\}$; $d(j_k^{\text{del}}, A_{k-1}, s_{A_{k-1}})$; $b_{k-1} - d(j, A_{k-1}, s_{A_{k-1}})$ for $(j, s) \in S_k^{\text{del}}$; $d(j, A_{k-1}, s_{A_{k-1}}) - b_{k-1}$ for $(j, s) \in A_{k-1} \setminus S_k^{\text{del}}$; and either $c(j_k^{\text{add}}, s_k^{\text{add}}, A_{k-1}, s_{A_{k-1}}) - d(j_k^{\text{del}}, A_{k-1}, s_{A_{k-1}})$, or the negative of this quantity, depending on whether a variable was added or deleted at step k . Altogether, the number of rows of Γ at step k is at most $\sum_{\ell=1}^k (|S_\ell^{\text{add}}| + |S_\ell^{\text{del}}| + 2|A_{\ell-1}^c| + |A_{\ell-1}| + 1) \leq 3pk + k$.

4.4. Details of the Exact Tests and Intervals

Given a number of steps k , after we have formed the appropriate Γ matrix for the FS, LAR, or lasso procedures, as derived in the last three subsections, computing conditional p -values and intervals is straightforward. Consider testing a generic null hypothesis $H_0 : v^T \theta = 0$ where v is arbitrary. First, we compute, as prescribed by Lemma 1, the quantities

$$\begin{aligned} \mathcal{V}^{\text{lo}} &= \max_{j: (\Gamma v)_j > 0} -(\Gamma y)_j \cdot \|v\|_2^2 / (\Gamma v)_j + v^T y, \\ \mathcal{V}^{\text{up}} &= \min_{j: (\Gamma v)_j < 0} -(\Gamma y)_j \cdot \|v\|_2^2 / (\Gamma v)_j + v^T y. \end{aligned}$$

Note that the number of operations needed to compute $\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}$ is $O(mn)$, where m is the number of rows of Γ . For testing against a one-sided alternative $H_1 : v^T \theta > 0$, we form the test statistic

$$T_k = 1 - F_{0, \sigma^2 \|v\|_2^2}^{[\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}]}(v^T y) = \frac{\Phi\left(\frac{\mathcal{V}^{\text{up}}}{\sigma \|v\|_2}\right) - \Phi\left(\frac{v^T y}{\sigma \|v\|_2}\right)}{\Phi\left(\frac{\mathcal{V}^{\text{up}}}{\sigma \|v\|_2}\right) - \Phi\left(\frac{\mathcal{V}^{\text{lo}}}{\sigma \|v\|_2}\right)}.$$

By Lemma 3, this serves as valid p -value, conditional on the selection. That is,

$$\mathbb{P}_{v^T \theta = 0} \left(T_k \leq \alpha \mid \hat{A}_k(y) = A_k, \hat{s}_{A_k}(y) = s_{A_k} \right) = \alpha, \quad (32)$$

for any $0 \leq \alpha \leq 1$. Also by Lemma 3, a conditional confidence interval is derived by first computing δ_α that satisfies

$$1 - F_{\delta_\alpha, \sigma^2 \|v\|_2^2}^{[\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}]}(v^T y) = \alpha.$$

Then we let $I_k = [\delta_\alpha, \infty)$, which has the proper conditional coverage, in that

$$\mathbb{P}(v^T \theta \in I_k \mid \hat{A}_k(y) = A_k, \hat{s}_{A_k}(y) = s_{A_k}) = 1 - \alpha. \quad (33)$$

For testing against a two-sided alternative $H_1 : v^T \theta \neq 0$, we instead use the test statistic

$$T'_k = 2 \cdot \min\{T_k, 1 - T_k\},$$

and by Lemma 4, the same results as in (32), (33) follow, but with T'_k in place of T_k , and $I'_k = [\delta_{\alpha/2}, \delta_{1-\alpha/2}]$ in place of I_k .

Recall that the case when $v = (X_{A_k}^+)^T e_k$, and the null hypothesis is $H_0 : e_k^T X_{A_k}^+ \theta = 0$, is of particular interest, as discussed in Section 2. Here, we are testing whether the coefficient of the last selected variable, in the population regression of θ on X_{A_k} , is equal to zero. For this problem, the details of the p -values and intervals follow exactly as above with the appropriate substitution for v . However, as we examine next, the one-sided variant of the test must be handled with care, in order for the alternative to make sense.

4.5. One-Sided or Two-Sided Tests?

Consider testing the partial regression coefficient of the variable to enter, at step k of FS, LAR, or lasso, in a projected linear model of θ on X_{A_k} . With the choice $v = (X_{A_k}^+)^T e_k$, the one-sided setup $H_0 : v^T \theta = 0$ versus $H_1 : v^T \theta > 0$ is not inherently meaningful, since there is no reason to believe ahead of time that the k th population regression coefficient $e_k^T X_{A_k}^+ \theta$ should be positive. By defining $v = s_k (X_{A_k}^+)^T e_k$, where recall s_k is the sign of the k th variable as it enters the (FS, LAR, or lasso) model, the null $H_0 : s_k e_k^T X_{A_k}^+ \theta = 0$ is unchanged, but the one-sided alternative $H_1 : s_k e_k^T X_{A_k}^+ \theta > 0$ now has a concrete interpretation: it says that the population regression coefficient of the last selected variable is nonzero, and *has the same sign* as the coefficient in the fitted (sample) model.

Clearly, the one-sided test here will have stronger power than its two-sided version when the described one-sided alternative is true. It will lack power when the appropriate population regression coefficient is nonzero, and has the opposite sign as the coefficient in the sample model. However, this is not really of concern, because the latter alternative seems unlikely to be encountered in practice, unless the size of the population effect is very small (in which case the two-sided test would not likely reject, as well). For these reasons, we often prefer the one-sided test, with $v = s_k (X_{A_k}^+)^T e_k$, for pure significance testing of the variable to enter at the k th step. The p -values in Table 1, for example, were computed accordingly.

With confidence intervals, the story is different. Informally, we find one-sided (i.e., half-open) intervals, which result from a one-sided significance test, to be less desirable from the perspective of a practitioner. Hence, for coverage statements, we often prefer the two-sided version of our test, which leads to two-sided conditional confidence intervals (selection intervals). The intervals in Figure 1, for example, were computed in this way.

4.6. Models with Intercept

Often, we run FS, LAR, or lasso by first beginning with an intercept term in the model, and then adding predictors. Our selection theory can accommodate this case. It is easiest to simply consider centering y and the columns of X , which is equivalent to including an intercept term in the regression. After centering, the covariance matrix of y is $\Sigma = \sigma^2(I - \mathbb{1}\mathbb{1}^T/n)$, where $\mathbb{1}$ is the vector of all 1s. This is fine, because the polyhedral theory from Section 3 applies to Gaussian random variables with an arbitrary (but known) covariance. With the centered y and X , the construction of the polyhedral set (Γ matrix) carries over just as described in Sections 4.1, 4.2, or 4.3. The conditional tests and intervals also carry over as in Section 4.4, except with the general contrast vector v replaced by its own centered version. Note that when v lies in the column space of X , for example, when $v = (X_{A_k}^+)^T e_k$, no changes at all are needed.

4.7. How Much to Condition on?

In Sections 4.1, 4.2, and 4.3, we saw in the construction of the polyhedral sets in (23), (28), (31) that it was convenient to condition on different quantities to define the FS, LAR, and lasso selection events, respectively. All three of the polyhedra in (23), (28), (31) condition on the active signs s_{A_k} of the selected model, and the latter two condition on more (loosely, the set of variables that were eligible to enter or leave the active model at each step). The decisions here, about what to condition on, were driven entirely by computational convenience. It is important to note that—even though any amount of extra conditioning will still lead to valid inference once we marginalize out part of the conditioning set (recall Section 2.2)—a greater degree of conditioning will generally lead to less powerful tests and wider intervals. This not only refers to the extra conditioning in the LAR and lasso selection events, but also to the specification of active signs s_{A_k} common to all three events. At the price of increased computation, one can eliminate unnecessary conditioning by considering a union of polyhedra (rather than a single one) as determining a selection event. This is done in Lee et al. (2016) and Reid, Taylor, and Tibshirani (2014). In FS regression, one can condition only on the sufficient statistics for the nuisance parameters, and obtain the most powerful selective test. Details are in Fithian, Sun, and Taylor (2014) and Fithian et al. (2016).

5. The Spacing Test for LAR

A computational challenge faced by the FS, LAR, and lasso tests described in the last section is that the matrices Γ computed for the polyhedral representations $\{y : \Gamma y \geq 0\}$ of their selection events can grow very large; in the FS case, the matrix Γ will have $2pk$ after k steps, and for LAR and lasso, it will have roughly $3pk$ rows. This makes it cumbersome to form \mathcal{V}^{lo} , \mathcal{V}^{up} , as the computational cost for these quantities scales linearly with the number of rows of Γ . In this section, we derive a simple approximation to the polyhedral representations for the LAR events, which remedies this computational issue.

5.1. A Refined Characterization of the Polyhedral Set

We begin with an alternative characterization for the LAR selection event, after k steps. The proof draws heavily on results from Lockhart et al. (2014), and is given in Appendix A.1.

Lemma 5. Suppose that the LAR algorithm produces the list of active variables A_k and signs s_{A_k} after k steps. Define $c(j, s, A_{k-1}, s_{A_{k-1}}) = (P_{A_{k-1}}^\perp X_j) / (s - X_j^T (X_{A_{k-1}}^+)^T s_{A_{k-1}})$, with the convention $A_0 = s_{A_0} = \emptyset$, so that $c(j, s, A_0, s_{A_0}) = c(j, s) = sX_j$. Consider the following conditions:

$$c(j_1, s_1, A_0, s_{A_0})^T y \geq c(j_2, s_2, A_1, s_{A_1})^T y \geq \dots \geq c(j_k, s_k, A_{k-1}, s_{A_{k-1}})^T y \geq 0, \quad (34)$$

$$c(j_k, s_k, A_{k-1}, s_{A_{k-1}})^T y \geq M_k^+ \left(j_k, s_k, c(j_{k-1}, s_{k-1}, A_{k-2}, s_{A_{k-2}})^T y \right), \quad (35)$$

$$c(j_\ell, s_\ell, A_{\ell-1}, s_{A_{\ell-1}})^T y \leq M_\ell^- \left(j_\ell, s_\ell, c(j_{\ell-1}, s_{\ell-1}, A_{\ell-2}, s_{A_{\ell-2}})^T y \right), \quad \text{for } \ell = 1, \dots, k, \quad (36)$$

$$0 \geq M_\ell^0 \left(j_\ell, s_\ell, c(j_{\ell-1}, s_{\ell-1}, A_{\ell-2}, s_{A_{\ell-2}})^T y \right), \quad \text{for } \ell = 1, \dots, k, \quad (37)$$

$$0 \leq M_\ell^S y, \quad \text{for } \ell = 1, \dots, k. \quad (38)$$

(Note that for $\ell = 1$ in (36), (37), we are meant to interpret $c(j_0, s_0, A_{-1}, s_{A_{-1}})^T y = \infty$.) The set of all y satisfying the above conditions is the same as the set \mathcal{P} in (28).

Moreover, the quantity M_k^+ in (35) can be written as a maximum of linear functions of y , each M_ℓ^- in (36) can be written as a minimum of linear functions of y , each M_ℓ^0 in (37) can be written as a maximum of linear functions of y , and each M_ℓ^S in (38) is a matrix. Hence, (34)–(38) can be expressed as $\Gamma y \geq 0$ for a matrix Γ . The number of rows of Γ is bounded above by $4pk - 2k^2 - k$.

At first glance, Lemma 5 seems to have done little for us over the polyhedral characterization in Section 4.2: after k steps, we are now faced with a Γ matrix that has on the order of $4pk$ rows (even more than before!). Meanwhile, at the risk of stating the obvious, the characterization in Lemma 5 is far more succinct (i.e., the Γ matrix is much smaller) without the conditions in (36)–(38). Indeed, in certain special cases (e.g., orthogonal predictors) these conditions are vacuous, and so they do not contribute to the formation of Γ . Even outside of such cases, we have found that dropping the conditions (36)–(38) yields an accurate (and computationally efficient) approximation of the LAR selection set in practice. This is discussed next.

5.2. A Simple Approximation of the Polyhedral Set

It is not hard to see from their definitions in Appendix A.1 that when X is orthogonal (i.e., when $X^T X = I$), we have $M_\ell^- = \infty$ and $M_\ell^0 = -\infty$, and furthermore, the matrix M_ℓ^S has zero rows, for each ℓ . This means that the conditions (36)–(38) are

vacuous. The polyhedral characterization in Lemma 5, therefore, reduces to $\{y : \Gamma y \geq U\}$, where Γ has only $k + 1$ rows, defined by the $k + 1$ constraints (34), (35), and U is a random vector with components $U_1 = \dots = U_k = 0$, and $U_{k+1} = M_k^+(j_k, s_k, c(j_{k-1}, s_{k-1}, A_{k-2}, s_{A_{k-2}})^T y)$.

For a general (nonorthogonal) X , we might still consider ignoring the conditions (36)–(38) and using the compact representation $\{y : \Gamma y \geq U\}$ induced by (34), (35). This is an approximation to the exact polyhedral characterization in Lemma 5, but it is a computationally favorable one, since Γ has only $k + 1$ rows (compared to about $4pk$ rows per the construction of the lemma). Roughly speaking, the constraints in (36)–(38) are often inactive (loose) among the full collection (34)–(38), so dropping them does not change the geometry of the set. Though we do not pursue formal arguments to this end (beyond the orthogonal case), empirical evidence suggests that this approximation is often justified.

Thus, let us suppose for the moment that we are interested in the polyhedron $\{y : \Gamma y \geq U\}$ with Γ, U as defined above, either serving an exact representation, or an approximate one, reducing the full description in Lemma 5. Our focus is the application of our polyhedral inference tools from Section 3 to $\{y : \Gamma y \geq U\}$. Recall that the established polyhedral theory considers sets of the form $\{y : \Gamma y \geq u\}$, where u is fixed. As the equivalence in (8) is a deterministic rather than a distributional result, it holds whether U is random or fixed. But the independence of the constructed $\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}, \mathcal{V}^0$, and $v^T y$ is not as immediate. The quantities $\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}, \mathcal{V}^0$ are now functions of y and U , both of which are random. An important special case occurs when $v^T y$ and the pair $((I - \Sigma v v^T / v^T \Sigma v)y, U)$ are independent. In this case $\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}, \mathcal{V}^0$ —which only depend on the latter pair above—are clearly independent of $v^T y$. To be explicit, we state this result as a corollary.

Corollary 1 (Polyhedral selection as truncation, random U). For any fixed y, Γ, U, v with $v^T \Sigma v \neq 0$,

$$\Gamma y \geq U \iff \mathcal{V}^{\text{lo}}(y, U) \leq v^T y \leq \mathcal{V}^{\text{up}}(y, U), \quad \mathcal{V}^0(y, U) \leq 0,$$

where

$$\begin{aligned} \mathcal{V}^{\text{lo}}(y, U) &= \max_{j: \rho_j > 0} \frac{U_j - (\Gamma y)_j + \rho_j v^T y}{\rho_j}, \\ \mathcal{V}^{\text{up}}(y, U) &= \min_{j: \rho_j < 0} \frac{U_j - (\Gamma y)_j + \rho_j v^T y}{\rho_j}, \\ \mathcal{V}^0(y, U) &= \max_{j: \rho_j = 0} U_j - (\Gamma y)_j, \end{aligned}$$

and $\rho = \Gamma \Sigma v / v^T \Sigma v$. Moreover, assume that y and U are random, and that

$$U \text{ is a function of } (I - \Sigma v v^T / v^T \Sigma v)y, \quad (39)$$

so $v^T y$ and the pair $((I - \Sigma v v^T / v^T \Sigma v)y, U)$ are independent. Then the triplet $(\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}, \mathcal{V}^0)(y, U)$ is independent of $v^T y$.

Under the condition (39) on U , the rest of the inferential treatment proceeds as before, as Corollary 1 ensures that we have the required alternate truncated Gaussian representation of $\Gamma y \geq U$, with the random truncation limits $\mathcal{V}^{\text{lo}}, \mathcal{V}^{\text{up}}$ being independent of the univariate Gaussian $v^T y$. In our LAR problem setup, U is a given random variate (as described in the first

paragraph of this subsection). The relevant question is of course: when does (39) hold? Fortunately, this condition holds with only very minor assumptions on v : this vector must lie in the column space of the LAR active variables at the current step.

Lemma 6. Suppose that we have run k steps of LAR, and represent the conditions (34), (35) in Lemma 5 as $\Gamma y \geq U$. Under our running regression model $y \sim N(\theta, \sigma^2 I)$, if v is in the column space of the active variables A_k , written $v \in \text{col}(X_{A_k})$, then the condition in (39) holds, so inference for $v^T \theta$ can be carried out with the same set of tools as developed in Section 3, conditional on $\Gamma y \geq U$.

The proof is given in Appendix A.2. For example, if we choose the contrast vector to be $v = (X_{A_k}^+)^T e_k$, a case we have revisited throughout the article, then this satisfies the conditions of Lemma 6. Hence, for testing the significance of the projected regression coefficient of the latest selected LAR variable, conditional on $\Gamma y \geq U$, we may use the p -values and intervals derived in Section 3. We walk through this usage in the next subsection.

5.3. The Spacing Test

The (approximate) representation of the form $\{y : \Gamma y \geq U\}$ derived in the last subsection (where Γ is small, having $k+1$ rows), can only be used to conduct inference over $v^T \theta$ for certain vectors v , namely, those lying in the span of current active LAR variables. The particular choice of contrast vector

$$v = c(j_k, s_k, A_{k-1}, s_{A_{k-1}}) = \frac{P_{A_{k-1}}^\perp X_{j_k}}{s_k - X_{j_k}^T (X_{A_{k-1}}^+)^T s_{A_{k-1}}}, \quad (40)$$

paired with the compact representation $\{y : \Gamma y \geq U\}$, leads to a very special test that we name the *spacing test*. From the definition (40), and the well-known formula for partial regression coefficients, we see that the null hypothesis being considered is

$$H_0 : v^T \theta = 0 \iff H_0 : e_k^T X_{A_k}^+ \theta = 0,$$

that is, the spacing test is a test for the k th coefficient in the regression of θ on X_{A_k} , just as we have investigated all along under the equivalent choice of contrast vector $v = (X_{A_k}^+)^T e_k$. The main appeal of the spacing test lies in its simplicity. Letting

$$\omega_k = \|(X_{A_k}^+)^T s_{A_k} - (X_{A_{k-1}}^+)^T s_{A_{k-1}}\|_2, \quad (41)$$

the spacing test statistic is defined by

$$T_k = \frac{\Phi(\lambda_{k-1} \frac{\omega_k}{\sigma}) - \Phi(\lambda_k \frac{\omega_k}{\sigma})}{\Phi(\lambda_{k-1} \frac{\omega_k}{\sigma}) - \Phi(M_k^+ \frac{\omega_k}{\sigma})}. \quad (42)$$

Above, λ_{k-1} and λ_k are the knots at steps $k-1$ and k in the LAR path, and M_k^+ is the random variable from Lemma 5. The statistic in (42) is one-sided, implicitly aligned against the alternative $H_1 : v^T \theta > 0$, where v is as in (40). Since $v^T y = \lambda_k \geq 0$, the denominator in (40) must have the same sign as $X_{j_k}^T P_{A_{k-1}}^\perp y$, that is, the same sign as $e_k^T X_{A_k}^+ y$. Hence,

$$H_1 : v^T \theta > 0 \iff H_1 : \text{sign}(e_k^T X_{A_k}^+ y) \cdot e_k^T X_{A_k}^+ \theta > 0,$$

that is, the alternative hypothesis H_1 is that the population regression coefficient of the last selected variable is nonzero, and shares the sign of the sample regression coefficient of the last

variable. This is a natural setup for a one-sided alternative, as discussed in Section 4.5.

The spacing test statistic falls directly out of our polyhedral testing framework, adapted to the case of a random U (Corollary 1 and Lemma 6). It is a valid p -value for testing $H_0 : v^T \theta = 0$, and has exact conditional size. We emphasize this point by stating it in a theorem.

Theorem 1 (Spacing test). Suppose that we have run k steps of LAR. Represent the conditions (34), (35) in Lemma 5 as $\Gamma y \geq U$. Specifically, we define Γ to have the following $k+1$ rows:

$$\begin{aligned} \Gamma_1 &= c(j_1, s_1, A_0, s_{A_0}) - c(j_2, s_2, A_1, s_{A_1}), \\ \Gamma_2 &= c(j_2, s_2, A_1, s_{A_1}) - c(j_3, s_3, A_2, s_{A_2}), \\ &\dots \\ \Gamma_{k-1} &= c(j_{k-1}, s_{k-1}, A_{k-2}, s_{A_{k-2}}) - c(j_k, s_k, A_{k-1}, s_{A_{k-1}}), \\ \Gamma_k &= \Gamma_{k+1} = c(j_k, s_k, A_{k-1}, s_{A_{k-1}}), \end{aligned}$$

and U to have the following $k+1$ components:

$$\begin{aligned} U_1 &= U_2 = \dots = U_k = 0, \\ U_{k+1} &= M_k^+ \left(j_k, s_k, c(j_{k-1}, s_{k-1}, A_{k-2}, s_{A_{k-2}})^T y \right). \end{aligned}$$

For testing the null hypothesis $H_0 : e_k^T X_{A_k}^+ \theta = 0$, the spacing statistic T_k defined in (41), (42) serves as an exact p -value conditional on $\Gamma y \geq U$:

$$\mathbb{P}_{e_k^T X_{A_k}^+ \theta = 0} (T_k \leq \alpha \mid \Gamma y \geq U) = \alpha,$$

for any $0 \leq \alpha \leq 1$.

Remark 3. The p -values from our polyhedral testing theory depend on the truncation limits \mathcal{V}^{lo} , \mathcal{V}^{up} , and in turn these depend on the polyhedral representation. For the special polyhedron $\{y : \Gamma y \geq U\}$ considered in the theorem, it turns out that $\mathcal{V}^{\text{lo}} = M_k^+$ and $\mathcal{V}^{\text{up}} = \lambda_{k-1}$, which is fortuitous, as it means that no extra computation is needed to form \mathcal{V}^{lo} , \mathcal{V}^{up} (beyond that already needed for the path and M_k^+). Furthermore, for the contrast vector v in (40), it turns out that $\|v\|_2 = 1/\omega_k$. These two facts completely explain the spacing test statistic (42), and the proof of Theorem 1, presented in Appendix A.3, reduces to checking these facts.

Remark 4. The event $\Gamma y \geq U$ is not exactly equivalent to the LAR selection event at the k th step. Recall that, as defined, this only encapsulates the first part (34), (35) of a longer set of conditions (34)–(38) that provides the exact characterization, as explained in Lemma 5. However, in practice, we have found that (34), (35) often provide a very reasonable approximation to the LAR selection event. In most examples, the spacing p -values are either close to those from the exact test for LAR, or exhibit even better power.

Remark 5. A two-sided version of the spacing statistic in (42) is given by $T'_k = 2 \cdot \min\{T_k, 1 - T_k\}$. The result in Theorem 1 holds for this two-sided version, as well.

5.4. Conservative Spacing Test

The spacing statistic in (42) is very simple and concrete, but it still does depend on the random variable M_k^+ . The quantity M_k^+

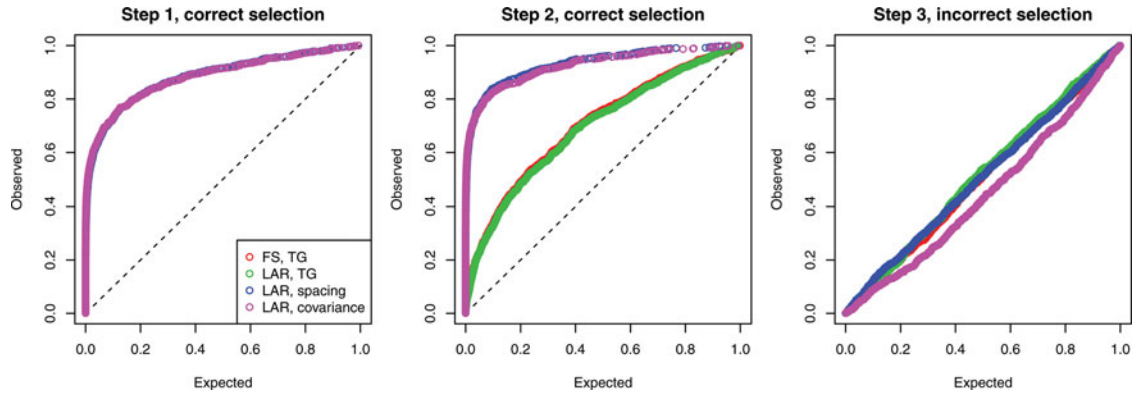


Figure 3. Simulated data with $n = 50$, $p = 100$, and two true active variables. Shown are p -values from the first three steps of FS and LAR, computed using the TG tests of Section 4, the spacing test of Section 5, and the covariance test of Lockhart et al. (2014), across 1000 repetitions (draws of y from the simulation model).

is computable in $O(p)$ operations (see Appendix A.1 for its definition), but it is not an output of standard software for computing the LAR path (e.g., the R package `larss`). To further simplify matters, therefore, we might consider replacing M_k^+ by the next knot in the LAR path, λ_{k+1} . The motivation is that sometimes, but not always, M_k^+ and λ_{k+1} will be equal. In fact, as argued in Appendix A.4, it will always be true that $M_k^+ \leq \lambda_{k+1}$, leading us to a conservative version of the spacing test.

Theorem 2 (Conservative spacing test). After k steps along the LAR path, define the modified spacing test statistic

$$\tilde{T}_k = \frac{\Phi(\lambda_{k-1} \frac{\omega_k}{\sigma}) - \Phi(\lambda_k \frac{\omega_k}{\sigma})}{\Phi(\lambda_{k-1} \frac{\omega_k}{\sigma}) - \Phi(\lambda_{k+1} \frac{\omega_k}{\sigma})}. \quad (43)$$

Here, ω_k is as defined in (41), and λ_{k-1} , λ_k , λ_{k+1} are the LAR knots at steps $k-1$, k , $k+1$ of the path, respectively. Let $\Gamma y \geq U$ denote the compact polyhedral representation of the spacing selection event step k of the LAR path, as described in Theorem 1. Then \tilde{T}_k is conservative, when viewed as a conditional p -value for testing the null hypothesis $H_0 : e_k^T X_{A_k}^+ \theta = 0$:

$$\mathbb{P}_{e_k^T X_{A_k}^+ \theta = 0}(\tilde{T}_k \leq \alpha \mid \Gamma y \geq U) \leq \alpha,$$

for any $0 \leq \alpha \leq 1$.

Remark 6. It is not hard to verify that the modified statistic in (43) is a monotone decreasing function of $\lambda_k - \lambda_{k+1}$, the spacing between LAR knots at steps k and $k+1$, hence the name “spacing” test. Similarly, the exact spacing statistic in (42) measures the magnitude of the spacing $\lambda_k - M_k^+$.

6. Empirical Examples

6.1. Conditional Size and Power of FS and LAR Tests

We examine the conditional Type I error and power properties of the truncated Gaussian (TG) tests for FS and LAR, as well as the spacing test for LAR, and the covariance test for LAR. We generated iid standard Gaussian predictors X with $n = 50$, $p = 100$, and then normalized each predictor (column of X) to have unit norm. We fixed true regression coefficients $\beta^* = (5, -5, 0, \dots, 0)$, and we set $\sigma^2 = 1$. For a total of 1000 repetitions, we drew observations according to $y \sim N(X\beta^*, \sigma^2 I)$, ran FS and LAR, and computed p -values across the first three steps.

Figure 3 displays the results in the form of QQ plots. The first plot in the figure shows the p -values at step 1, conditional on the algorithm (FS or LAR) having made a correct selection (i.e., having selected one of the first two variables). The second plot shows the same, but at step 2. The third plot shows p -values at the step 3, conditional on the algorithm having made an incorrect selection.

At step 1, all p -values display very good power, about 73% at a 10% nominal Type I error cutoff. There is an interesting departure between the tests at step 2: we see that the covariance and spacing tests for LAR actually yield much better power than the exact TG tests for FS and LAR: about 82% for the former versus 35% for the latter, again at a nominal 10% Type I error level. At step 3, the TG and spacing tests produce uniform p -values, as desired; the covariance test p -values are super-uniform, showing the conservativeness of this method in the null regime.

Why do the methods display such differences in power at step 2? A rough explanation is as follows. The spacing test, recall, is defined by removing a subset of the polyhedral constraints for the conditioning event for LAR, thus its p -values are based on less conditioning than the exact TG p -values for LAR. Because it conditions on less, that is, it uses a larger portion of the sample space, it can deliver better power; and though it (as well as the covariance test) is not theoretically guaranteed to control Type I error in finite samples, it certainly appears to do so empirically, seen in the third panel of Figure 3. The covariance test is believed to behave more like the spacing test than the exact TG test for LAR; this is based on an asymptotic equivalence between the covariance and spacing tests, given in Section 7, and it explains their similarities in the plots.

6.2. Coverage of LAR Conditional Confidence Intervals

In the same setup as the previous subsection, we computed conditional confidence intervals over the first three steps of LAR, at a 90% coverage level. Figure 4 shows these intervals across the first 100 repetitions. Each interval here is designed to cover the partial regression coefficient of a particular variable, in a population regression of the mean $\theta = X\beta^*$ on the variables in the current active set. These population coefficients are drawn as black dots, and the colors of the intervals reflect the identities of the variables being tested: red for variable 1, green for variable 2, and blue for all other variables. Circles around

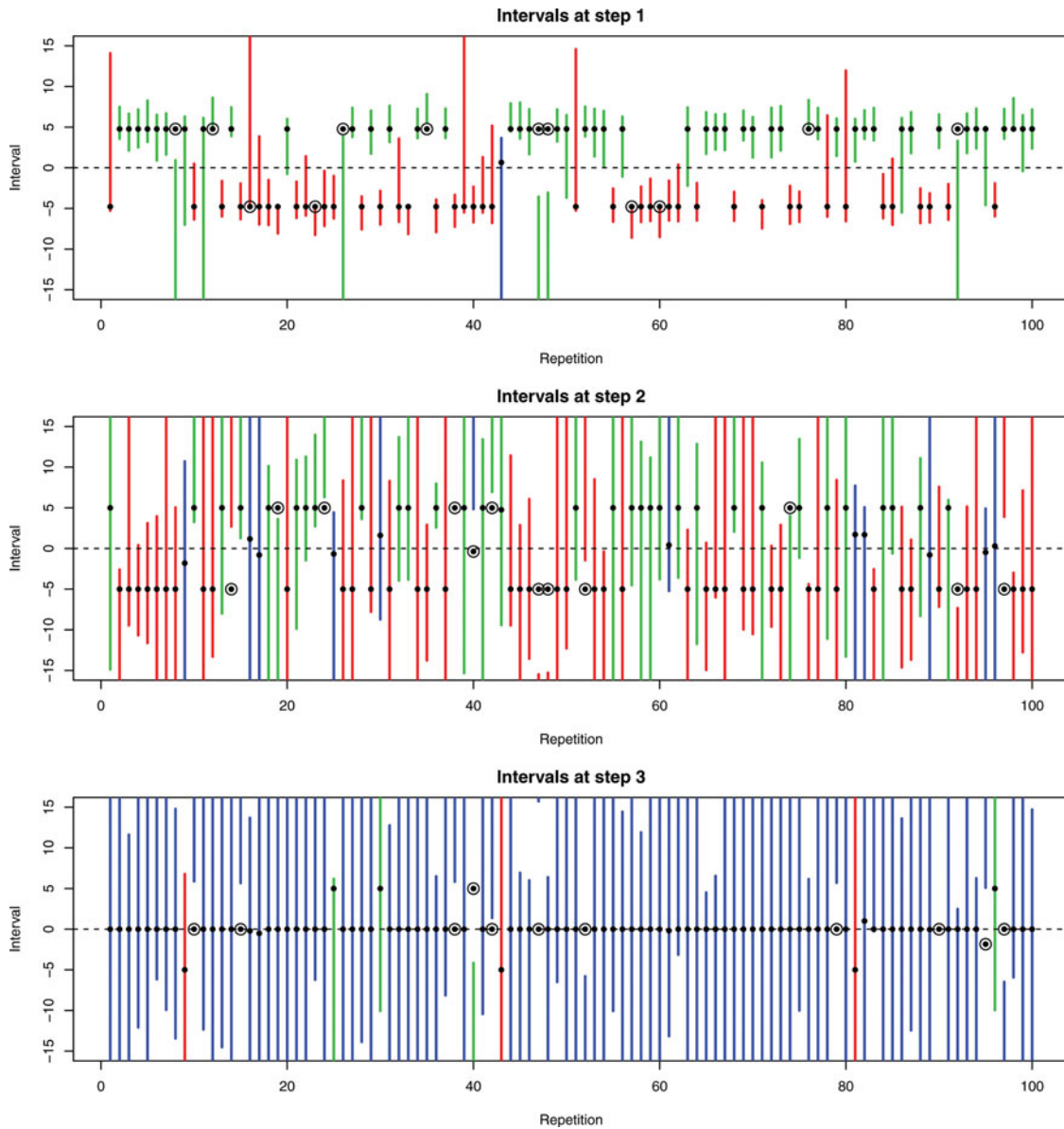


Figure 4. Using the same setup as in Figure 3, shown are 90% confidence intervals for the selections made over the first three steps of LAR, across 100 repetitions (draws of y from the simulation model). The colors assigned to the intervals reflect the identities of the variables whose partial regression coefficients are being tested: red for variable 1, green for variable 2, and blue for all others. The black dots are the true population coefficients, and circles around these dots denote miscoverages. The upper confidence limits for some of the parameters exceed the range for the y -axes on the plots (especially at step 3).

population coefficients indicate that these particular coefficients are not covered by their corresponding intervals. The miscoverage proportion is 12/100 in step 1, 11/100 in step 2, and 11/100 in step 3, all close to the nominal miscoverage level of 10%. An important remark: here we are counting marginal coverage of the intervals. Our theory actually further guarantees *conditional* coverage, for each model selection event, for example, among the red intervals at step 1, the miscoverage proportion is 7/52, and among green intervals, it is 5/47, both close to the nominal 10% level.

6.3. Comparison to the $\max\text{-}|t|$ -Test

The last two subsections demonstrated the unique properties of the exact TG tests for FS and LAR. For testing the significance of variables entered by FS, Buja and Brown (2014) proposed what they call the $\max\text{-}|t|$ -test. Here is a description. At the k th step of

FS, where A_{k-1} is the current active list (with $k-1$ active variables), let

$$t_{\max}(y) = \max_{j \notin A_{k-1}} \frac{|X_j^T P_{A_{k-1}}^\perp y|}{\sigma \|P_{A_{k-1}}^\perp X_j\|_2}.$$

As the distribution of $t_{\max}(y)$ is generally intractable, we simulate $\epsilon \sim N(0, \sigma^2 I)$, and use this to estimate null probability that $t_{\max}(\epsilon) > t_{\max}(y)$, which forms our p -value.

We used the same setup as in the previous two subsections, but with an entirely null signal, that is, we set the mean to be $\theta = X\beta^* = 0$, to demonstrate the following point. As we step farther into the null regime (as we take more and more steps with FS), the $\max\text{-}|t|$ -test becomes increasingly conservative, whereas the exact TG test for FS continues to produce uniform p -values, as expected. The reason is that the TG test for FS at step k properly accounts for all selection events up to and including

step k , but the $\max\text{-}|t|$ -test at step k effectively ignores all selections occurring before this step, creating a conservative bias in the p -value. See Appendix A.5 for the plots.

7. Relationship to the Covariance Test

There is an interesting connection between the LAR spacing test and the covariance test of Lockhart et al. (2014). We first review the covariance test and then discuss this connection.

After k steps of LAR, let A_k denote the list of active variables and s_{A_k} denote the sign list, the same notation as we have been using thus far. The covariance test provides a significance test for the k th step of LAR. More precisely, it assumes an underlying linear model $\theta = X\beta^*$, and tests the null hypothesis

$$H_0 : A_{k-1} \supseteq \text{supp}(\beta^*),$$

where $\text{supp}(\beta^*)$ denotes the support of set of β^* (the true active set). In words, this tests simultaneously the significance of *any variable entered at step k and later*.

Though its original definition is motivated from a difference in the (empirical) covariance between LAR-fitted values, the covariance statistic can be written in an equivalent form that is suggestive of a connection to the spacing test. This form, at step k of the LAR path, is

$$C_k = \omega_k^2 \cdot \lambda_k(\lambda_k - \lambda_{k+1})/\sigma^2, \quad (44)$$

where λ_k, λ_{k+1} are the LAR knots at steps k and $k+1$ of the path, and ω_k is the weight in (41). (All proofs characterizing the null distribution of the covariance statistic in Lockhart et al. (2014) use this equivalent definition.) The main result (Theorem 3) in Lockhart et al. (2014) is that, under correlation restrictions on the predictors X and other conditions, the covariance statistic (44) has a conservative $\text{Exp}(1)$ limiting distribution under the null hypothesis. Roughly, they show that

$$\lim_{n, p \rightarrow \infty} \mathbb{P}_{A_{k-1} \supseteq \text{supp}(\beta^*)} \left(C_k > t \mid \hat{A}_k(y) = A_k, \hat{s}_{A_k}(y) = s_{A_k} \right) \leq e^{-t},$$

for all $t \geq 0$.

A surprising result, perhaps, is that the covariance test in (44) and the spacing test in (43) are asymptotically equivalent. The proof for this equivalence uses relatively straightforward calculations with Mills' inequalities, and is deferred until Appendix A.6.

Theorem 3 (Asymptotic equivalence between spacing and covariance tests). After a fixed number k steps of LAR, the spacing p -value in (43) and the covariance statistic in (44) are asymptotically equivalent, in the following sense. Assume an asymptotic regime in which

$$\omega_k \lambda_{k+1} \xrightarrow{P} \infty, \quad \text{and} \quad \omega_k^2 \cdot \lambda_{k-1}(\lambda_{k-1} - \lambda_k) \xrightarrow{P} \infty,$$

denoting convergence in probability. The spacing statistic, transformed by the inverse $\text{Exp}(1)$ survival function, satisfies

$$-\log \left(\frac{\Phi(\lambda_{k-1} \frac{\omega_k}{\sigma}) - \Phi(\lambda_k \frac{\omega_k}{\sigma})}{\Phi(\lambda_{k-1} \frac{\omega_k}{\sigma}) - \Phi(\lambda_{k+1} \frac{\omega_k}{\sigma})} \right) = \frac{\omega_k^2}{\sigma^2} \lambda_k(\lambda_k - \lambda_{k+1}) + o_P(1).$$

Said differently, the asymptotic p -value of the covariance statistic, under the $\text{Exp}(1)$ limit, satisfies

$$\exp \left(-\frac{\omega_k^2}{\sigma^2} \lambda_k(\lambda_k - \lambda_{k+1}) \right) = \left(\frac{\Phi(\lambda_{k-1} \frac{\omega_k}{\sigma}) - \Phi(\lambda_k \frac{\omega_k}{\sigma})}{\Phi(\lambda_{k-1} \frac{\omega_k}{\sigma}) - \Phi(\lambda_{k+1} \frac{\omega_k}{\sigma})} \right) (1 + o_P(1)).$$

Above, we use $o_P(1)$ to denote terms converging to zero in probability.

Remark 7. The asymptotic equivalence described in this theorem raises an interesting and unforeseen point about the one-sided nature of the covariance test. That is, the covariance statistic is seen to be asymptotically tied to the spacing p -value in (43), which, recall, we can interpret as testing the null $H_0 : e_k^T X_{A_k}^+ \theta = 0$ against the one-sided alternative $H_1 : \text{sign}(e_k^T X_{A_k}^+ y) \cdot e_k^T X_{A_k}^+ \theta > 0$. The covariance test in (44) is hence implicitly aligned to have power when the selected variable at the k th step has a sign that matches that of the projected population effect of this variable.

8. Discussion

In a regression model with Gaussian errors, we have presented a method for exact inference, conditional on a polyhedral constraint on the observations y . Since the FS, LAR, and lasso algorithms admit polyhedral representations for their model selection events, our framework produces exact p -values and confidence intervals post model selection for any of these adaptive regression procedures. One particularly special and simple case arises when we use our framework to test the significance of the projected regression coefficient, in the population, of the latest selected variable at a given step of LAR. This leads to the spacing test, which is asymptotically equivalent to the covariance test of Lockhart et al. (2014). An R language package `selectiveInference`, that implements the proposals in this article, is freely available on the CRAN repository, as well as <https://github.com/selective-inference/R-software>. A Python implementation is also available at <https://github.com/selective-inference/Python-software>.

Supplementary Materials

The supplementary materials contain additional proofs.

Acknowledgment

The authors thank Andreas Buja, Max Grazier G'Sell, Alessandro Rinaldo, and Larry Wasserman for helpful comments and discussion. The authors also thank the editors and referees whose comments led to a complete overhaul of this article.

Funding

Richard Lockhart was supported by the Natural Sciences and Engineering Research Council of Canada; Jonathan Taylor was supported by NSF grant DMS 1208857 and AFOSR grant 113039; Ryan Tibshirani was supported by NSF grant DMS-1309174; and Robert Tibshirani was supported by NSF grant DMS-9971405 and NIH grant N01-HV-28183.

References

- Berk, R., Brown, L., Buja, A., Zhang, K., and Zhao, L. (2013), "Valid Post-Selection Inference," *Annals of Statistics*, 41, 802–837. [601]
- Buhlmann, P. (2013), "Statistical Significance in High-Dimensional Linear Models," *Bernoulli*, 19, 1212–1242. [601]
- Buja, A., and Brown, L. (2014), "Discussion: A Significance Test for the Lasso," *Annals of Statistics*, 42, 509–517. [612]
- Choi, Y., Taylor, J., and Tibshirani, R. (2014), "Selecting the Number of Principal Components: Estimation of the True Rank of a Noisy Matrix," arXiv: 1410.8260. [601]
- Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004), "Least Angle Regression," *Annals of Statistics*, 32, 407–499. [606,607]
- Fithian, W., Sun, D., and Taylor, J. (2014), "Optimal Inference After Model Selection," arXiv: 1410.2597. [601,604,608]
- Fithian, W., Taylor, J., Tibshirani, R. J., and Tibshirani, R. (2016), "Selective Sequential Model Selection," arXiv:1512.02565. [608]
- G'Sell, M., Wager, S., Chouldechova, A., and Tibshirani, R. (2016), "Sequential Selection Procedures and False Discovery Rate Control," *Journal of the Royal Statistical Society, Series B*, 78, 423–444. [603]
- Javanmard, A., and Montanari, A. (2014a), "Confidence Intervals and Hypothesis Testing for High-dimensional Regression," *Journal of Machine Learning Research*, 15, 2869–2909. [601]
- (2014b), "Hypothesis Testing in High-dimensional Regression Under the Gaussian Random Design Model: Asymptotic Theory," *IEEE Transactions on Information Theory*, 60, 6522–6554. [601]
- Lee, J., Sun, D., Sun, Y., and Taylor, J. (2016), "Exact Post-Selection Inference With the Lasso," *Annals of Statistics*, 44, 907–927. [601,603,608]
- Lee, J., and Taylor, J. (2014), "Exact Post Model Selection Inference for Marginal Screening," *Advances in Neural Information Processing Systems*, 27, 136–144. [601]
- Lehmann, E., and Romano, J. (2005), *Testing Statistical Hypotheses* (3rd ed.), New York: Springer. [604]
- Lockhart, R., Taylor, J., Tibshirani, R. J., and Tibshirani, R. (2014), "A Significance Test for the Lasso," *Annals of Statistics*, 42, 413–468. [600,609,613]
- Loftus, J., and Taylor, J. (2014), "A Significance Test for Forward Stepwise Model Selection," arXiv: 1405.3920. [601]
- Meinshausen, N., and Buhlmann, P. (2010), "Stability Selection," *Journal of the Royal Statistical Society, Series B*, 72, 417–473. [601]
- Minnier, J., Tian, L., and Cai, T. (2011), "A Perturbation Method for Inference on Regularized Regression Estimates," *Journal of the American Statistical Association*, 106, 1371–1382. [601]
- Reid, S., Taylor, J., and Tibshirani, R. (2014), "Post-selection Point and Interval Estimation of Signal Sizes in Gaussian Samples," arXiv: 1405.3340. [601,608]
- Tibshirani, R. J. (2013), "The Lasso Problem and Uniqueness," *Electronic Journal of Statistics*, 7, 1456–1490. [605]
- van de Geer, S., Buhlmann, P., Ritov, Y., and Dezeure, R. (2014), "On Asymptotically Optimal Confidence Regions and Tests for High-dimensional Models," *Annals of Statistics*, 42, 1166–1201. [601]
- Wasserman, L., and Roeder, K. (2009), "High-Dimensional Variable Selection," *Annals of Statistics*, 37, 2178–2201. [601]
- Zhang, C.-H., and Zhang, S. (2014), "Confidence Intervals for Low Dimensional Parameters in High Dimensional Linear Models," *Journal of the Royal Statistical Society, Series B*, 76, 217–242. [601]

JOURNAL OF THE AMERICAN STATISTICAL ASSOCIATION
2016, VOL. 111, NO. 514, Theory and Methods
<http://dx.doi.org/10.1080/01621459.2016.1182788>

Comment

Lawrence D. Brown and Kory D. Johnson

1. Introduction

The authors provide a novel and exciting framework for analyzing conditional selection. Formalizing the steps of a selection procedure as constraints on the response is applicable beyond the linear models theory discussed here and involves a high degree of technical accomplishment. It also raises interesting questions about different approaches to conditional inference. That being said, the usefulness of these tests appears limited in practice.

For ease of exposition, we focus on the forward stepwise case, though the arguments are also applicable to least angle regressions (LARS). The authors propose an "Exact Forward Stepwise" procedure (FS) that assigns new, "exact" p -values to the

variables in a standard forward selection algorithm based on the usual analysis of variance (ANOVA) forward selection. At each stage, the algorithm includes the variable that creates the largest reduction in error sum of squares. After a variable is added, it is assigned a " p -value" by this "exact" procedure. This is a numerical quantity that has a $U(0,1)$ distribution conditional on the sign of the selected variable and the variables that have been previously chosen.

This extends conditional inference ideas and calculations from other recent articles so as to provide p -values for feature selection algorithms. As these algorithms are commonly used to build multiple regression models, one might think that improved p -values would lead to improved model selection, at least in some circumstances; however, the formulation in

this article involves a serious paradox. One needs to begin with a well-specified model selection algorithm and construct a model independent of the exact p -values described in the article. The exact p -values can be constructed only after the model has been chosen; they cannot validly be used to select the model. If one tries to use them in this way, they become invalid.

While forward stepwise and LARS operate independently of these p -values, one would expect the modeler to want to use the p -values to determine the step at which to “stop” the procedure and provide a final model. Consider the authors’ Table 1 (recreated below), which compares their FS p -values to naive forward stepwise p -values. Identifying a final model using such a table requires considering multiple p -values from separate steps of the procedure. Therein lies the problem: the set of exact p -values cannot be used to make decisions, else they are invalid. Even using these p -values as input into a secondary FDR-controlling procedure as in G’Sell et al. (2016) is inappropriate. Only one exact p -value can validly be used, testing one step of a much larger procedure. Similarly, if a model is selected through other means such as cross-validation, the inferential guarantees of related methods need not hold (Bachoc, Leeb, and Pötscher 2014).

It should be noted that the conventional p -values are single-step values. They do not correct for the multiple testing nature of a stepwise procedure. Later in this commentary we recommend modified versions of the p -value calculations that can be validly and directly used for stepwise selection. See procedures (b) and (c) defined later. The procedure ES, defined later, is built on conditional inference logic and could be used to replace FS. For reasons discussed below, however, we do not favor its use.

The columns in Table 1 labeled “JASA” are taken directly from the authors’ Table 1. The column labeled “Seq. p -value” contains traditional two-sided p -values we calculated from our version of the data. Further information on our computations is in the appendix. Note that our p -values are two-sided, whereas those in the article are one-sided. While providing one-sided p -values may aid numerical comparison to the FS exact p -values, we note that these are one-sided conditional on the sign of the observed effect. Therefore, they are in effect two-sided p -values and should be compared to ordinary two-sided p -values. (We believe our “Seq. p -value” entries should be twice those in the column “FS, naive (JASA).” They are not exactly so but are numerically close to that.)

The paradox in using the FS p -values is rather subtle, and is easiest to explain in the context of an example. Let X_i be independently distributed $N(\theta_i, 1)$, for $i \in \{1, 2\}$. The forward selection problem is equivalent to determining an order for testing $H_{0,i}$: $\theta_i = 0$, while controlling false rejections at level α . Since we are performing model selection, a variable is “included” or “added” to the model when the corresponding null hypothesis is rejected. Allowing correlated variables does not change the basic ideas in our discussion, though it does introduce further complications mentioned in our final paragraph. Without loss of generality, consider the case when $X_1 > X_2 > 0$.

The authors’ FS significance thresholds are given as “FS Step 1” and “FS Step 2” in Figure 1. The conditioning set for both steps of the procedure is the same: $\{X_1 > X_2 > 0\}$. Values to

Table 1. Replicated stepwise table.

| Step | Parameter | Seq. p -value | FS, naive (JASA) | FS, exact (JASA) |
|------|-----------|-----------------|------------------|------------------|
| 1 | lcavol | 0.0000 | 0.000 | 0.000 |
| 2 | lweight | 0.0003 | 0.002 | 0.006 |
| 3 | svi | 0.0424 | 0.024 | 0.425 |
| 4 | lbph | 0.0468 | 0.023 | 0.168 |
| 5 | ppg45 | 0.2304 | 0.116 | 0.423 |
| 6 | lcp | 0.0878 | 0.041 | 0.273 |
| 7 | age | 0.1459 | 0.069 | 0.059 |
| 8 | gleason | 0.8839 | 0.442 | 0.156 |

the right of the dashed curve “FS Step 1” yield p -values below α when testing $H_{0,1}$ while values between “FS Step 1” and the 45° line yield p -values greater than α . Thus, values to the right of FS Step 1 are those for which the statistician using FS p -values would select X_1 with a positive sign at the first step of the selection process. During the second step, values above the dash-dotted curve “FS Step 2” are significant at level α , while values below are not. Note that the calculation at the second step does not change depending on whether or not $H_{0,1}$ was rejected.

In order to use the FS p -values as Table 1 would imply, testing $H_{0,2}$ must account for rejecting $H_{0,1}$. Following the methodology of the authors’ article, this requires updating the conditioning set. We propose the following corrected procedure, “Exact Stepwise” (ES), that terminates on the first step in which a corrected, conditional p -value is above α . If $H_{0,2}$ is only tested when $H_{0,1}$ is rejected by FS, then the conditioning set is the region to the right of FS Step 1. Those points to the right of FS Step 1 and outside the convex, parabolic region whose boundary is the curve “ES Step 2” are those for which the new ES procedure selects X_1 at the first step (with a positive coefficient) and X_2 at the second step (with positive coefficient). It is clear that this correction does not invalidate the authors’ methodology, but it does yield different p -values. Furthermore, the new conditioning sets are not polyhedral and need not be convex. In spite of the authors’ indirect implication, convex polyhedral conditioning regions are not

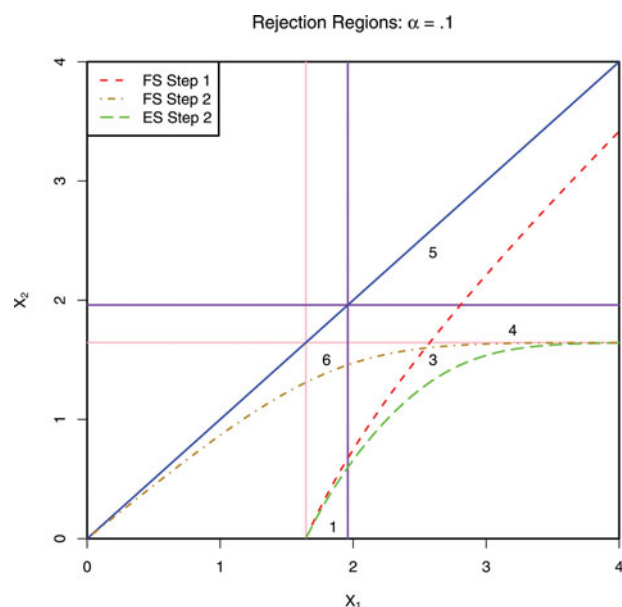


Figure 1. Stepwise rejection regions at $\alpha = 0.1$. The full picture is symmetric around the x - and y -axes. A corresponding image would be drawn if $X_2 > X_1 > 0$, in which case the graph would be rotated 90° and maintain its symmetries.

required for their methodology, although computations are simpler for such regions.

2. Stopping Procedures Using ES p -Values

We are also concerned with the counter-intuitive results given when using conditional p -values, even when they are corrected as already discussed. The problem has obvious symmetries such as relabeling variables 1 and 2 or changing their signs. While our new proposal, ES, preserves those symmetries, it does not preserve the natural monotonicity of the problem. For example, there exist values (x_1, x_2) and (x'_1, x'_2) for which $x_1 \leq x'_1$ and $x_2 \leq x'_2$, but for which ES selects both variables at (x_1, x_2) and no variables at (x'_1, x'_2) . The authors' FS procedure does not produce as extreme an example since $H_{0,2}$ is tested regardless of the result of testing $H_{0,1}$; however, the significance of the test of $H_{0,1}$ depends on the value of X_2 . This is particularly troubling given that X_1 and X_2 are independent.

It is also instructive to compare the rejection regions of the FS and ES procedures to those of more traditional methods (again, see Figure 1). The conventional procedure first adds X_1 if $|X_1| > |X_2|$ and $|X_1| > \Phi^{-1}(1 - \alpha/2)$. It then adds X_2 if, also, $|X_2| > \Phi^{-1}(1 - \alpha/2)$. If $|X_2| > |X_1|$ and $|X_2| > \Phi^{-1}(1 - \alpha/2)$ the first step adds X_2 , etc. Relevant portions of the lines form the boundaries of this region. There are multiple, classically constructed stepwise regions to control for selection and multiple comparisons:

- (a) The thin vertical lines are at $X_1 = 1.645$ and $X_2 = 1.645$. To the right of the 45° line, they show regions where the fully classical stepwise procedure would first choose X_1 (to the right of $X_1 = 1.645$) and then X_2 (above $X_2 = 1.645$). This region is not adjusted for multiple comparisons, and hence has an error probability of choosing nonempty models under the null hypothesis that exceeds the nominal level α .
- (b) Similarly, the thick vertical lines at $X_1 = 1.96$ and $X_2 = 1.96$ bound regions that provide conservative multiple-comparison adjustments based on conventional single-coordinate p -values. This uses the Bonferroni approximation to control for multiple-comparisons. The exact numerical value can be computed from a maximum modulus calculation and is 1.948.
- (c) A better stepwise procedure controlling for multiple comparisons can be constructed as follows: at each step, choose among the remaining k variables using a p -value threshold such that the null probability of choosing any model is less than or equal to α . As in (b), Bonferroni yields the conservative threshold α/k , though an exact calculation is possible when k is small. In the figure, one would include X_1 to the right of $X_1 = 1.96$ (or 1.948 for an exact calculation) and then would include X_2 when X_2 is above $X_2 = 1.645$. At each step of the procedure, the conditional probability under the null hypothesis of continuing with an incorrect rejection is α . This type of procedure was briefly proposed in Buja and Brown (2014). If the goal is to preserve FDR, then one can improve the procedure, and we are currently working on an article to explain how to do so. Johnson, Stine, and Foster (2015a)

provided a related procedure that controls mFDR with much higher power.

Comparison of Procedures

Many interesting comparisons can be made between FS, ES, and the more conventionally motivated, multiplicity-corrected procedures (b) and (c). These regions are labeled (numerically) in Figure 1.

1. Consider the triangular region to the right of FS Step 1 and to the left of $X_1 = 1.96$. This is where the ES procedure selects X_1 and the conventionally motivated procedures choose no variables. Heuristically, this seems to be a success for the ES procedure.
2. Within the region described in 1, there is a sliver between FS Step 1 and ES Step 2. Here, ES selects both X_1 and X_2 , while the conventional procedures select neither. While this maintains a significance guarantee, this may not be an advantage. These points do have conventional (two-sided) p -values for X_1 that are below α , but the conventional p -value for X_2 is quite large. Selecting X_2 appears to be a mistake.
3. There is a more noticeable triangular region bounded by FS Step 1, ES Step 2, and $X_2 = 1.645$. In this region, the ES procedure selects both X_1 and X_2 but (b) and (c) select only X_1 . For reasons similar to those in 2, the second step of ES appears undesirable. The disadvantage is not as clear, however, since X_2 can have a two-sided p -value as small as $\alpha = 0.1$ in this region. The uncorrected FS p -value yields intuitively more satisfactory results in this region.
4. Consider the region between $X_2 = 1.645$ and $X_2 = 1.96$, and to the right of FS Step 1. ES and (c) select both variables, but the simpler, conservative procedure (b) does not include X_2 . The advantage here goes to ES and (c).
5. The area between the 45° line and FS Step 1 and to the right of $X_1 = 1.96$ is where (b) and (c) have a clear advantage in power relative to ES or to a procedure based on FS. In those regions, ES and FS have first step p -values above α and hence do not select any variable, while (b) and (c) always select X_1 and often select X_2 .
6. In the region above FS Step 2 and below $X_2 = 1.645$, X_2 has a significant FS p -value even though its conventional p -value can be close to 1 near the origin.

In summary, (b) or (c) seem preferable to the ES procedure. The latter does better if the data fall in the small, but not negligible region 1; however, (b) and (c) produce much more reasonable models in the more noticeable area 5. Procedure (c) is preferred to (b) because of the difference noted in region 4. The regions 2 and 3 are quite small and nearly negligible in probability. While the ES procedure seems undesirable on these regions, the concern is not important.

As a further comparison, consider the point (4, 3.8). The ES p -values for Steps 1 and 2 are ≈ 0.44 and ≈ 0.0001 , respectively, while the naive p -values are approximately 0.0001 for both variables. The decision of ES to stop at Step 1 and declare an empty model might well be viewed as embarrassing and subjectively undesirable. This is consistent with the claimed ES

p -values though. Similarly, while methods of G'Sell et al. (2016) are not required to stop at Step 1, the penalty for continuing is extremely large given the unconventionally large p -value.

In the correlated setting, the interesting simulation in Section 6 strongly suggests that the p -values used in procedures (b) and (c) can be extremely conservative. Hence, the naive scheme previously suggested must be modified to achieve desirable performance. To decide whether such modification is possible needs further research, and, if possible, further investigation of the geometric structure and stochastic performance of the resulting tests. One possibility for improved performance is to switch to the sequential selection viewpoint of Johnson, Stine, and Foster (2015a).

For all considered testing methods, when the regressors are correlated, the values of regression coefficients depend on which other coefficients are in the current model. Hence, a coefficient may have a nonzero value within the currently active set of variables; and so be correctly included into that model at that step. Within a later active model it might then have a value of 0. Thus, a correct selection at a given step may become incorrect as the process proceeds, and vice-versa. The related phenomenon of suppression can yield a series of insignificant steps followed by highly significant steps (Johnson, Stine, and Foster 2015b). These issues have important consequences for interpretation of p -values produced in a stepwise routine. Such issues do not occur in the simple model at hand involving independent variables with fixed mean values.

Appendix

The sequential p -values were constructed using data downloaded from Robert Tibshirani's website. The p -values computed in Table 1 are computed from the standard F -test with 1 and $58 = 67 - 9$ degrees of freedom. As some additional numerical details, note that the mean squared error (MSE) from the full model is $\hat{\sigma}^2 = 0.5074$. Thus, for example, the sequential F -value for testing "svi" is $2.1841/0.5074 = 4.305$ with a t -value of $2.075 = \sqrt{4.305}$. This has a p -value with 58 degrees of freedom of 0.0426.

FS Step 1 (dashed curve): If X_1 is chosen before X_2 with a positive sign, the observation lies in the cone $R_1 = \{X_1 > X_2 > 0\}$. To have a level α test of $H_0: \theta_1 = 0$ conditional on $(x_1, x_2) \in R_1$, one must have

$$\theta = \mathbb{P}(X_1 > \tau_1 | (x_1, x_2) \in R_1, X_2 = x_2) \quad \forall x_2.$$

This entails choosing the point via

$$\alpha = \frac{1 - \Phi(x_1)}{1 - \Phi(|x_2|)}. \quad (1)$$

This defines $x_1 = x_1(x_2)$ for the dashed curve.

FS Step 2 (dash-dotted curve): The conditioning region is the same, so the level α test of $H_0: \theta_2 = 0$ conditional on $(x_1, x_2) \in R_1$ requires

$$\theta = \mathbb{P}(X_2 > \tau_2 | (x_1, x_2) \in R_1, X_1 = x_1) \quad \forall x_1.$$

This entails choosing the point via

$$\alpha = \frac{\Phi(x_1) - \Phi(X_2)}{\Phi(x_1) - 1/2}. \quad (2)$$

ES Step 2 (dotted curve): Given $H_0: \theta_1 = 0$ has been rejected, possible values of (X_1, X_2) lie to the right of FS Step 1. Denote this region as R_2 . Now the test $H_0: \theta_2 = 0$ must satisfy

$$\theta = \mathbb{P}(X_2 > \tau_2 | (x_1, x_2) \in R_2, X_1 = x_1)$$

for all x_1 for which the conditioning region is nonempty. The only change from FS Step 2 is that the conditioned region is a function of x_2 . This entails choosing the point $x_2 = x_2(x_1)$ for which

$$\alpha = \frac{\Phi(x_2^*) - \Phi(X_2)}{\Phi(x_2^*) - 1/2}, \quad (3)$$

where x_2^* denotes the value for which $x_1(x_2^*) = x_1$. The computation in Equation (3) is facilitated by noting that Equation (1) implies

$$\Phi(x_2^*(x_1)) = 1 + \frac{\Phi(x_1) - 1}{\theta}. \quad (4)$$

Funding

This research was supported in part by an NSF grant.

References

- Bachoc, F., Leeb, H., and Pötscher, B. M. (2014), "Valid Confidence Intervals for Post-Model-Selection Predictors," *ArXiv preprint*. Available at <https://arxiv.org/abs/1412.4605>. [615]
- BuJa, A., and Brown, L. (2014), Discussion of "A Significance Test for the Lasso," *Annals of Statistics*, 42, 509–517. [616]
- G'Sell, M. G., Wager, S., Chouldechova, A., and Tibshirani, R. (2016), "Sequential Selection Procedures and False Discovery Rate Control," *Journal of the Royal Statistical Society, Series B*, 78, 423–444. [615, 617]
- Johnson, K. D., Stine, R. A., and Foster, D. P. (2015a), "Revisiting Alpha-Investing: Conditionally Valid Stepwise Regression," *ArXiv e-prints*. Available at <http://arxiv.org/abs/1510.06322>. [617]
- (2015b), "Submodularity in Statistics: Comparing the Success of Model Selection Methods," *ArXiv e-prints*. Available at <http://arxiv.org/abs/1510.06301>. [617]

Rejoinder

Ryan J. Tibshirani, Jonathan Taylor, Richard Lockhart, and Robert Tibshirani

We thank Drs. Brown and Johnson for their detailed discussion. We appreciate their comments, and it is helpful to hear honest criticism. But in this case we think that much of the criticism results from misinterpretations. As our thinking has evolved since this article was written, we take the blame for much of their confusion.

1. Usefulness of the Proposed Tools

The discussants say that the “usefulness of [the proposed] tests appears limited in practice.” This seems because they believe we are proposing *purely sequential tests* that should be used to select a model (i.e., choose a stopping point) along the forward stepwise or lasso paths. However, this is *not the focus of our article*. Instead, we propose a much broader set of tools that can be used to conduct inference on any variable, at any step along the forward stepwise, LAR, or lasso paths. Here, we reiterate the discussion from Section 2.3 of our article, which describes different ways of using these tools. We focus on forward stepwise regression (FS), although analogous points apply to the least angle regression (LAR) and lasso paths.

- (a) If the practitioner has a fixed step k in mind, then the framework in our article allows him or her to build (say) confidence intervals for each partial regression coefficient corresponding to an active variable in the k -step model. These can be interpreted simultaneously, but one must apply the usual Bonferroni correction, adjusting here for k tests.

For example, with $p = 200$ variables, a practitioner might decide to inspect the significance of no more than $k = 20$ of these variables, but (of course) would not specify a priori, which ones to inspect. The practitioner would then run FS for 20 steps and examine the resulting confidence intervals from our framework. To adjust for multiplicity, he or she would compute confidence intervals at the level $1 - \alpha/20$, and would declare significance of a variable, in the context of the discovered 20-step model, whenever its corresponding interval excludes 0 (thus controlling the simultaneous Type I error at the level α). In our `selectiveInference` package, one would apply the `fsInf` function with `k=20` and `type="all"`.

- (b) If instead the practitioner does not have a fixed step k in mind, then he or she can use an AIC-like (or BIC-like) stopping rule to choose a step \hat{k} adaptively from the data. Consider the rule that chooses the smallest step for which the Akaike information criterion (AIC) or Bayesian information criterion (BIC) criterion rises (say) twice in a row. We omit the details here, but it is straightforward to show that the event determining \hat{k} here can also be phrased in terms of a polyhedral constraint on y . This means that the practitioner can run this stopping rule, select \hat{k} , and then test the significance of all variables in the \hat{k} -step model. Conditioning both on \hat{k} and on the \hat{k} -step model (which is possible because this is just one large polyhedral constraint on y) yields valid inferences. Again, these can be interpreted simultaneously, but one needs a Bonferroni correction, now for \hat{k} tests.

For example, with $p = 200$ variables, a practitioner might run forward stepwise and use the modified AIC or BIC rule described above to stop after $\hat{k} = 26$ steps. He or she would then compute confidence intervals for each partial regression coefficient in the 26-step model, at the level $1 - \alpha/26$, to properly adjust for multiplicity. In our `selectiveInference` package, this workflow is given by using the `fsInf` function with `k=NULL` and `type="aic"`.

- (c) Finally, using our framework, it is possible to compute a p -value or confidence interval for the variable that enters the FS active model at each step. G'Sell et al. (2015) developed a stopping rule “ForwardStop” that can be applied to such a sequence of p -values, and this rule is guaranteed to control the FDR as long as the p -values are independent under the null. The sequential p -values from our framework are not generally independent, and so these guarantees do not apply, but the ForwardStop rule can still be used as an approximation. In our `selectiveInference` package, this workflow is given by using the `fsInf` function with `k=NULL` and `type="active"`, followed by a call to the function `forwardStop`.

The discussants point out the shortcomings of this approach, namely, that (i) this strategy can have low

power (i.e., produce large p -values), for example, when we have entered only the first of two strong variables; and (ii) it is not possible to use these sequential p -values to select a stopping point \hat{k} and then perform inference (since in this case we must condition on the event determining \hat{k} , which may be very complicated). Both of these are valid criticisms. Indeed we have been aware of them since the near-beginnings of our work. Section 4.3 of Fithian, Sun, and Taylor (2014), for example, discusses precisely the issue (i) above.

In our article, we did not intend to endorse the purely sequential workflow—as described in point (c) above—as the main use case of our tools, but in hindsight, we see that the language may well have been confusing. We would now like to make it clear that, given the framework in the article under discussion, inspecting sequential p -values is *not the ideal usage* for the reasons raised above, and points (a) and (b) describe better strategies.

We would also like to point out that in Fithian et al. (2015), we focus primarily on the sequential aspect of the FS, LAR, and lasso procedures, and develop much improved sequential hypothesis tests for their paths. The sequential p -values constructed in this new article are *independent under the null*, meaning that ForwardStop can be applied to determine a stopping point \hat{k} , with guaranteed FDR control. Further, the sequential p -values from this new article also have *much better power*, and they circumvent the issue (i) above. See Section 5.3 of Fithian et al. (2015) for a discussion.

Figure 1 shows the results of a simulation study with $n = 40$, $p = 8$, and iid $N(0, 1)$ errors. The data were generated from a linear model with the first four coefficients being nonzero, drawn from $N(0, 1)$. We computed the p -values from this article, called “saturated model p -values” in Fithian et al. (2015), and those called “selective model p -values” in the same article. We then applied the ForwardStop rule to estimate the stopping point, that is, the number of predictors chosen by FS. The target FDR was 0.1. The predictors were normally distributed with equal pairwise correlation ρ from -0.1 to 0.8 , indicated along the horizontal axis of each panel. The left panel shows the achieved FDR and the right panel depicts the average number of predictors chosen, both computed over 2000 repetitions

from the described simulation setup. The standard errors are about 0.003 on the left and 0.02 on the right. We see that both procedures seem to control the FDR, while the selective model p -values yield a larger model on average. As we said earlier, FDR control can only be proven for the selective model p -values, since they are independent under the null. On the other hand, the selective model p -values require Markov chain Monte Carlo (MCMC) sampling and for large problems, this could be a limitation.

2. Comments on Their ES Procedure and Analysis

As explained above, the discussants have interpreted our framework to be entirely sequential in nature, and their comments and analysis are directed toward the sequential model selection problem as a result. Though this is not truly the focus of our work, we would still like to make several comments in reply to their ES procedure and analysis.

- Their proposed ES procedure is interesting, but it does not seem generally feasible computationally, as the conditioning regions appear to become quite complicated after several steps.
- We did not intend to “indirectly” imply that our framework is limited to polyhedral conditioning events; in Fithian, Sun, and Taylor (2014), we consider very general selection events, though polyhedral events are of course computationally convenient and make the computation feasible in essentially arbitrary dimensions.
- The conditioning event used by the discussants in Figure 1 is not quite correct, that is, it is not what would be constructed out of our framework. In the “normal means” problem they describe, with $X_i \sim N(\theta_i, 1)$ for $i = 1, 2$, if FS were to select X_1 and X_2 , in that order, both with positive signs, then the conditioning event at Step 1 would be $\{X_1 \geq X_2, X_1 \geq -X_2\}$, and at Step 2 it would be $\{X_1 \geq X_2 \geq 0\}$. The discussants merely consider the latter event for both steps. This means that for step 1, the regions in their Figure 1 should have counterparts that are given by reflections around the x -axis.
- The Bonferroni correction described in their point (b) in Figure 1 would need to be applied, in general, across p

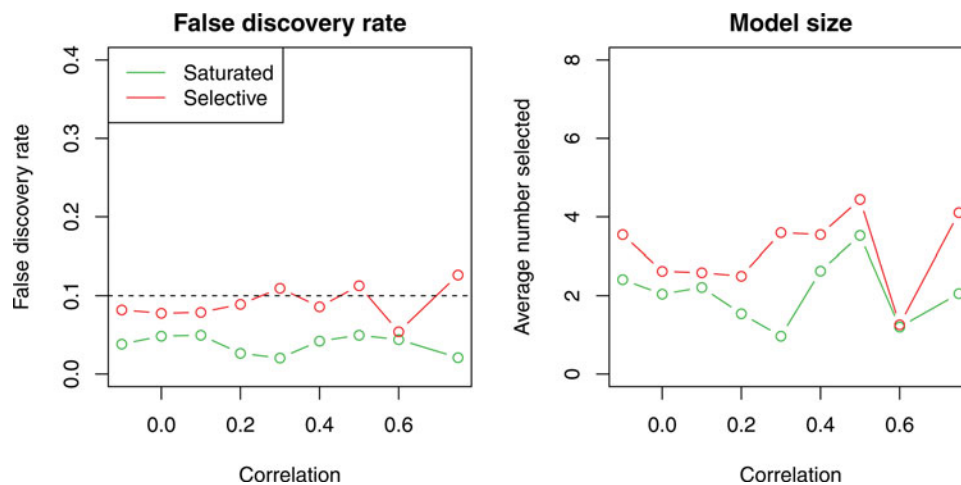


Figure 1. Results of a simulation study to assess p -value stopping rules. Details are in the text.

tests. The discussants consider $p = 2$. In larger problems, of course, such a Bonferroni correction will become very conservative. This should be compared to the strategies we described in the first two bullet points in the previous section, which need only perform a Bonferroni correction at the level k or \hat{k} .

- The point $(X_1, X_2) = (4, 3.8)$ in their Figure 1, which results in a large p -value of about 0.44 at Step 1 of FS and a small p -value of about 0.0001 at Step 2 (that the discussants describe as “embarrassing and subjectively undesirable”) is simply due to using the saturated model p -values in a sequential manner. As we explained at the end of the previous section, the selected model p -values from Fithian et al. (2015) are much better-suited to the sequential problem. In this instance $(X_1, X_2) = (4, 3.8)$, for example, the selected model p -values are each less than 0.0001 for both FS steps.

3. Summary

To summarize, the best use of the tools developed in the discussed article is to compute p -values or confidence intervals for all selected variables in the FS (or LAR or lasso) model after some number of steps. We might write this number as k to indicate that it has been fixed a priori, or as \hat{k} to indicate that it has been chosen by a procedure akin to AIC or BIC. These strategies were described in points (a) and (b) in Section 1. The strategy described in point (c), in which we sequentially compute p -values of the variable to enter the FS (or LAR or lasso) active

set, and use these p -values to determine a final model size \hat{k} of interest, has its flaws, as pointed out by the discussants. First, these p -values can be low-powered, and second, after they have been used to select a stopping point \hat{k} , it is unclear how to perform valid inference (since the conditioning event determining \hat{k} is itself very complicated). We should, however, emphasize the fact that our newer work in Fithian et al. (2015) produces sequential p -values that overcome the first problem: they display much better power along the path. Using these p -values to determine a stopping point, and then performing valid inference (in a computationally tractable manner), remains an open issue.

Finally, we remind the reader of our R language package `selectiveInference`, that implements the proposals in our article, is freely available on CRAN, as well as <https://github.com/selective-inference/R-software>. A Python implementation is also available, at <https://github.com/selective-inference/Python-software>. We welcome feedback from users, to help us improve this package. We plan to actively develop and support it.

References

- Fithian, W., Sun, D., and Taylor, J. (2014), “Optimal Inference After Model Selection,” arXiv: 1410.2597. [xxxx]
- Fithian, W., Taylor, J., Tibshirani, R., and Tibshirani, R. J. (2015), “Selective Sequential Model Selection,” arXiv: 1512.02565. [619]
- G’Sell, M., Wager, S., Chouldechova, A., and Tibshirani, R. (2015), “Sequential Selection Procedures and False Discovery Rate Control,” *Journal of the Royal Statistical Society, Series B*, 78, 426–444. [618]