# Univariate–Guided Sparse Regression

Sourav Chatterjee[1,3]     Trevor Hastie[2,3]     Robert Tibshirani[2,3]

[1]Department of Mathematics
[2]Department of Biomedical Data Science
[3]Department of Statistics
Stanford University

**Abstract**

In this paper, we introduce "uniLasso"— a novel statistical method for regression. This two-stage approach preserves the signs of the univariate coefficients and leverages their magnitude. Both of these properties are attractive for stability and interpretation of the model. Through comprehensive simulations and applications to real-world datasets, we demonstrate that uniLasso outperforms lasso in various settings, particularly in terms of sparsity and model interpretability. We prove asymptotic support recovery and mean-squared error consistency under a set of conditions different from the well-known irrepresentability conditions for the lasso. Extensions to generalized linear models (GLMs) and Cox regression are also discussed.

## 1  Introduction

High-dimensional regression and classification problems are ubiquitous across fields such as genomics, finance, and social sciences. In these settings, the lasso (Least Absolute Shrinkage and Selection Operator) has emerged as a widely-used methodology due to its ability to simultaneously perform variable selection and regularization, yielding sparse and interpretable models. Despite its widespread use, lasso has certain limitations, including sensitivity to correlated predictors and inclusion of spurious features in the final model. There have been many related proposals including MC+ (Zhang 2010), the elastic net (Zou & Hastie 2005), the adaptive lasso (Zou 2006), and SparseNet (Mazumder et al. 2011).

To address these challenges, we propose uniLasso, a novel regression methodology. In this work, we aim to achieve two primary goals: predictive accuracy and (especially) interpretability. To achieve this, uniLasso integrates marginal (univariate) information into a coherent multivariate framework: our method preserves the signs of the univariate coefficients and leverages their magnitude. This approach not only enhances predictive accuracy but also provides insights into the relative importance of predictors.

Here is an example, taken from unpublished work with collaborators of the third author. (Since the work is unpublished, we do not give background details). The dataset is from cancer proteomics, with 559 biomarkers measured on 81 patients, 20 healthy and 61 with cancer. We divided the data into 70% training and 30% test, and applied both the lasso and uniLasso to the training set. Both methods performed well, giving cross-validation error of about 5%, and 4 and 3 errors respectively out of 30 test samples. Table 1 shows the results of lasso and uniLasso applied to these data, using cross-validation to choose the model. In the top table there are two sign changes from univariate to multivariate: #12 goes from mildly protective to a negative indicator, while #11 goes in the opposite direction.

By design, uniLasso in the bottom table has no sign changes. Further, its univariate coefficients are much larger (in absolute value) than those of the lasso. This makes it a more credible and stable model.

The remainder of this paper is organized as follows. Section 2 provides a detailed description of

the uniLasso methodology. We examine the relationship between uniLasso and the adaptive lasso in Section 3, and in Section 4 we study the underlying reason for the strong sparsity produced by uniLasso. We compare uniLasso to lasso on a car price data set in Section 5: here we show the utility of unregularized uniLasso as an alternative to least squares. In Section 6 we look at the orthonormal design case, where we can derive an explicit expression for the uniLasso solution. Section 7 presents some theoretical results on support recovery and mean square error consistency. Section 8 presents simulation studies comparing uniLasso to lasso and other benchmark methods, while Section 9 examines uniReg, the unregularized ($\lambda = 0$) case. Section 10 applies uniLasso to real-world datasets, and in Section 11 we adapt uniLasso to multiclass classification problems. Section 12 explores the setting where additional data is available: not the complete data, but just the univariate scores. In Section 14 we study whether cross-validation works properly in the uniLasso setting. We discuss uniLasso for GLMs and the Cox Survival Model, as well as computation in Sections 15 and 16. Lastly, Section 17 concludes with a discussion of future directions.

*Model chosen by lasso*

| Biomarker | Univariate LS Coefficient | Lasso Coefficient |
|:---:|:---:|:---:|
| 1 | -49.24 | -18.92 |
| 2 | -37.65 | -5.42 |
| 3 | 25.27 | 37.49 |
| 4 | -24.74 | -11.26 |
| 5 | 22.91 | 14 |
| 6 | -22.38 | -12.75 |
| 7 | -17.97 | -11.91 |
| 8 | 13.02 | 0.79 |
| 9 | -12.71 | -1.85 |
| 10 | -10.14 | -2.32 |
| 11 | 1.81 | -6.27 |
| 12 | -0.34 | 6.25 |

*Model chosen by uniLasso*

| Biomarker | Univariate LS Coefficient | uniLasso Coefficient |
|:---:|:---:|:---:|
| A | 93.13 | 13.82 |
| B | -84.26 | -8.17 |
| C | 77.69 | 15.86 |
| D | 73.84 | 27.66 |
| E | 69.15 | 16.21 |
| F | 58.92 | 1.38 |
| G | -53.46 | -10.5 |
| H | -49.24 | -20.76 |
| I | 25.27 | 10.48 |

Table 1: *Proteomics study. Top table: Univariate least squares coefficients and lasso coefficients, for model chosen by lasso. The two sign changes from univariate to multivariate are shown in blue. Bottom table: same for model chosen by uniLasso. Biomarkers are ordered by decreasing absolute value. Most biomarkers are different in the two tables, with the exception of 1 and 3 in the first table corresponding to H and I in the second table.*

# 2 Univariate-guided lasso

## 2.1 Our proposal

We assume the standard supervised learning setup: we have training features and target $X_{n \times p}$, $y_{n \times 1}$. For now we assume that $y$ is quantitative and we fit a linear model using squared error loss; later we discuss the binomial and other GLM families, as well as the Cox model.

Our procedure has three simple steps, which we motivate here. For interpretability and prediction accuracy, we preprocess the features in Step 1, multiplying them by a robust version of their univariate least-squares coefficient estimates. In Step 2 we fit a non-negative lasso tuned by cross-validation, and Step 3 combines the components of Steps 1 and 2 to produce a final model. Together these ensure that:

(a) the signs of the final coefficients agree with the signs of the univariate coefficients (or they are zero);

(b) features with larger univariate coefficients will tend to have larger coefficients in the final model.

---

### UniLasso algorithm

1. For $j = 1, 2, \ldots, p$ compute the univariate intercepts and slopes $(\hat{\beta}_{0j}, \hat{\beta}_j)$ and their leave-one-out (LOO) counterparts $(\hat{\beta}_{0j}^{-i}, \hat{\beta}_j^{-i})$, $i = 1, \ldots, n$.

2. Fit the lasso — with an intercept, no standardization, and non-negativity constraints — to target $y$ and the univariate LOO fits as features

$$\underset{\theta}{\text{minimize}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \Big( y_i - \theta_0 - \sum_{j=1}^{p} (\hat{\beta}_{0j}^{-i} + \hat{\beta}_j^{-i} x_{ij}) \theta_j \Big)^2 + \lambda \sum_{j=1}^{p} \theta_j \right\} \quad \text{with } \theta_j \geq 0 \; \forall j.$$

Select $\lambda$ by cross-validation.

3. The final model can be written as $\hat{\eta}(x) = \hat{\gamma}_0 + \sum_{j=1}^{p} \hat{\gamma}_j x_j$, with $\hat{\gamma}_j = \hat{\beta}_j \hat{\theta}_j$, and $\hat{\gamma}_0 = \hat{\theta}_0 + \sum_{\ell=1}^{p} \hat{\beta}_{0\ell} \hat{\theta}_\ell$.

---

This procedure is computationally convenient: in Step 1 we can use efficient LOO formulas, and in Step 2 we can apply any efficient $\ell_1$ solver. Here we used the R language package `glmnet`. Specifically, we use the function `cv.glmnet` to estimate the lasso path parameter, and have all of the functionality of `glmnet` at our disposal. We provide a function `cv.UniLasso` in the R package `uniLasso` that implements this approach.

Note that we *do not* standardize the features before applying the non-negative lasso in Step 2; the univariate LOO fits are all on the scale of the response. From our knowledge of multiple linear regression, the first constraint — agreement between univariate and multivariate signs — may seem like an unreasonable restriction. However our belief is that in high-dimensions there are likely to be a multitude of different models that have about the same MSE as the "optimal" one chosen by the lasso. Hence it can make sense to choose one that is interpretable and sparser than that of lasso.

These properties mean, for example, if the feature "age" has a positive univariate coefficient — indicating increasing risk of the outcome variable (such as Alzheimer's disease), it will have a positive (or zero) coefficient in the final lasso model. And if age is strongly significant on its own, it is more likely to be chosen in the multivariate model.

These conjectures are borne out by our numerical studies. In simulations with varying problems sizes and SNR, and a number of real datasets (see Section 8), in almost every case uniLasso did no

worse than the lasso in terms of out-of-sample MSE, and delivered a substantially sparser model. We note that Meinshausen (2012) studies sign-constrained least squares estimation for high-dimensional regression, which relates to our non-negativity constraint in our second step.

**Remark A**. The uniLasso procedure applies seamlessly to other GLMs as well the Cox model. Indeed, all of the models covered by the `glmnet` package are at our disposal. All we need are the separate fitted linear models and their LOO fit vectors in step 1, and then `glmnet` can be directly applied. The only challenge is find an (approximate) LOO formula for the GLM families. We give details of these computations in Section 16.

**Remark B**. UniLasso can be thought of as a version of stacked regression, or *Stacking* (Wolpert 1992, Breiman 1996). Stacking is a two-stage procedure for combining the predictions from a set of models ("learners"), in order to get improved predictions. It works as follows. A set of base models is trained on the training data; these models can be of the same type (e.g. gradient boosting) or different types (e.g., linear regression, decision trees, etc.). Each model generates predictions, capturing specific aspects of the relationship between the predictors and the target variable. A meta-model is then trained to combine the predictions of the base models into a final prediction. This meta-model learns how to weigh and integrate the LOO outputs of the base models to minimize the overall prediction error. UniLasso is a special case of stacking where the individual learners are simple univariate regressions.

**Remark C**. Why do we use the LOO univariate estimates $\hat{\eta}_j^{-i} = \hat{\beta}_{0j}^{-i} + \hat{\beta}_j^{-i} x_{ij}$ as features in step 2, instead of the univariate estimates $\hat{\eta}_j^i = \hat{\beta}_{0j} + \hat{\beta}_j x_{ij}$?[1] In traditional stacking this is essential because the individual learners can be of very different complexity. In uniLasso, one would think that the learners (univariate regression) all have the same complexity so that the LOO versions are not needed. Indeed, the theory in Section 7 is based on the LOO estimates but holds equally well for the usual (non-LOO) estimates. However in practice we have found that the LOO estimates lead to greater sparsity and better performance, and hence we use them. With the use of the usual (non-LOO) univariate estimates in uniLasso, the resulting estimator is closely related to the adaptive lasso and we explore this connection in Section 3.

**Remark D**. The non-negative garotte (Breiman 1995) is another closely related method. It minimizes $\sum_i (y_i - \sum_j c_j \hat{\beta}_j x_{ij})^2$ subject to $c_j \geq 0$ for all $j$, and $\sum c_j \leq s$, where $\hat{\beta}_j$ are the usual (multivariable) LS estimates. This is different from our proposal in an important way: our use of univariate LOO estimates in the first step. The univariate coefficients lead to materially different solutions and also allow application to the $p > n$ scenario.

## 2.2   UniLasso with no-regularization

We get an interesting special case of uniLasso if we set $\lambda = 0$, so that there is no $\ell_1$ regularization in Step 2. That is, we solve the following::

$$\text{minimize}_\theta \ \left\{ \frac{1}{n} \sum_{i=1}^n \left( y_i - \theta_0 - \sum_{j=1}^p \theta_j (\hat{\beta}_{0j}^{-i} + \hat{\beta}_j^{-i} x_{ij}) \right)^2 \right\} \quad \text{with } \theta_j \geq 0 \ \forall j. \tag{1}$$

We call this "uniReg" for univariate-guided regression and it represents an interesting alternative to the usual least squares estimates. For example, the non-negative constraint can still provide sparsity even though there is no $\ell_1$ penalty. If $p > n$ there can be multiple solutions; in this case we look for the sparsest solution by computing the limiting uniLasso solution as $\lambda \downarrow 0$. The standard error and distribution of the estimated coefficients can be estimated via the bootstrap. We study uniReg in detail in Section 9.

## 2.3   Two examples — one good, one bad

**Homecourt.**   We generated data with 100 observations, 30 standard normal features with AR(1)

---

[1]In this case we could simply scale the features by their univariate slope coefficients and ignore the intercepts, which would be handled by the overall intercept in the model.

correlation 0.8, 20% sparsity and non-negative coefficients in two stages, as follows:

$$y' \quad \leftarrow \quad \sum_j x_j \beta_j + \sigma' z' \tag{2}$$

From these data we compute the $p = 30$ univariate least squares coefficients $\hat{\beta}_j^{uni}$ separately for each $j$, and then generate $y$ as

$$y \quad \leftarrow \quad \sum_j x_j \hat{\beta}_j^{uni} \beta_j + \sigma z. \tag{3}$$

The variance terms $\sigma'$ and $\sigma$ were chosen so that at both stages the SNR was 1. The idea here is that (3) roughly mimics the model fit by uniLasso in its second stage.
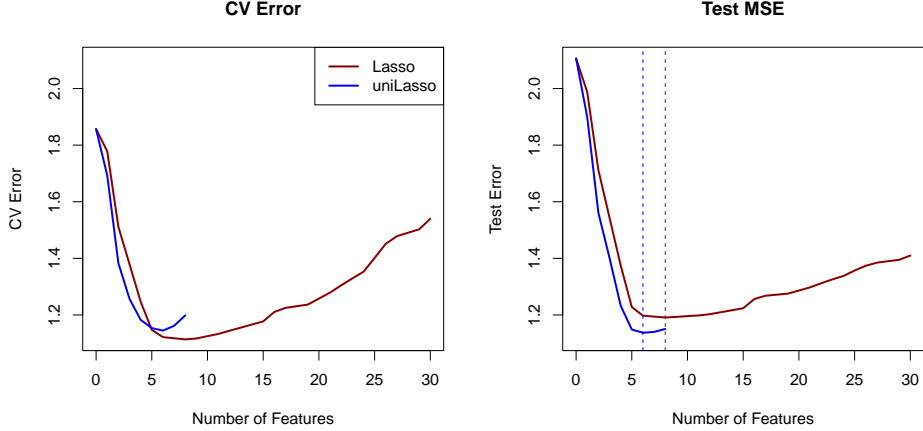


Figure 1: *Results for homecourt example: CV and test set prediction error. The dashed vertical lines for each method correspond to the models chosen by CV.*

|  | MSE-lasso | MSE-uniLasso | Support-lasso | Support-uniLasso |
|---|---|---|---|---|
| Mean | 1.098 | 1.077 | 7.660 | 4.790 |
| se | 0.005 | 0.005 | 0.309 | 0.142 |
|  | TPR-lasso | TPR-uniLasso | FPR-lasso | FPR-uniLasso |
| Mean | 0.737 | 0.700 | 0.135 | 0.025 |
| se | 0.014 | 0.016 | 0.012 | 0.004 |

Table 2: *Test MSE, support, TPR (True Positive Rate) and FPR (False Positive Rate) for 100 simulations from the setting of Figure 1.*

Figure 1 shows the CV and test error curves for the lasso and uniLasso. We see that `uniLasso` has test error a little below that of `lasso`, with a smaller active set. Table 2 shows the result of 100 simulations from this setup. We see that the same pattern emerges, and uniLasso exhibits a slightly lower true positive rate and a much lower false positive rate than the lasso.

**Counter-example.** Here we take $n = 100$, $p = 20$, $x_1 \sim N(0, 1)$, $x_2 = x_1 + N(0, 1)$, $\beta = (1, -.5, 0, 0, \ldots 0)$, and error SD $= 0.5$. The remaining 18 features are standard normal. This is an "Achilles heel" for uniLasso, as the (negative) sign of the population coefficient for $x_2$ differs from its (positive) univariate sign. As a result, the test MSE for uniLasso is about twice that of lasso (detailed results are given in Section 8). Clearly uniLasso fails badly here, but its high CV error will alert the user to this. Motivated by this kind of example, we discuss a post-processing ("polish") for uniLasso in Section 13 that remedies this problem.

# 3   Relationship of uniLasso to the adaptive lasso

The adaptive lasso (Zou 2006) is defined by

$$\hat{\gamma}^* = \text{argmin} = \frac{1}{2} \sum_{i=1}^{n} (y_i - \gamma_0 - \sum_{j=1}^{p} x_{ij}\gamma_j)^2 + \lambda \sum_{j=1}^{p} w_j |\gamma_j| \tag{4}$$

where $w_j = 1/|\hat{\gamma}_j|^\nu$. Here the vector $\hat{\gamma}$ is any root-$n$ consistent estimate of $\gamma$, for example the least squares estimates.

The original paper assumed $n > p$; with $p > n$, one cannot use the least-squares estimates. One possibility is to use an initial estimate (from say ridge or lasso), solve the adaptive lasso, and iterate. Another option is to use the univariate least-squares estimates; (Huang et al. 2008) consider a combination of these to achieve some theoretical guarantees.

Let $\{\hat{\beta}_j\}_1^p$ be the univariate least squares coefficients. Consider the adaptive lasso with weights $w_j = 1/|\hat{\beta}_j|$, $j = 1, \ldots, p$. Then it is easy to show that this procedure is equivalent to uniLasso with $\hat{\eta}_j^i = \hat{\beta}_{0j} + \hat{\beta}_j x_{ij}$ replacing their LOO versions $\hat{\eta}_j^{-i} = \hat{\beta}_{0j}^{-i} + \hat{\beta}_j^{-i} x_{ij}$ in Step 2, and removing the non-negativity constraint.

This procedure does not share some of the properties of uniLasso. In particular, the final coefficients may not have the same signs as the univariate coefficients. Importantly, in our simulations of Section 8, it tends to yield less sparse models and sometimes much higher MSE.

Alternatively, one could include the non-negativity constraints in this version of the adaptive lasso. The following result shows that we obtain a procedure equivalent to uniLasso, except for the use of LOO estimates in step 1.

**Proposition 1.** Let $(\hat{\beta}_{0j}, \hat{\beta}_j)$ be the univariate least squares coefficients for variable $j$ in a $p$ variable linear model with data $\{(x_i, y_i)\}_1^n$. Let $\eta_j^i = \hat{\beta}_{0j} + \hat{\beta}_j x_{ij}$ be the univariate linear fit for variable $j$ and observation $i$.

Then the following two problems are equivalent:

$$\min_{\theta} \sum_{i=1}^{n} (y_i - \theta_0 - \sum_{j=1}^{p} \eta_j^i \theta_j)^2 + \lambda \sum_{j=1}^{p} |\theta_j| \quad \text{s.t. } \theta_j \geq 0 \ \forall j \tag{5}$$

$$\min_{\gamma} \sum_{i=1}^{n} (y_i - \gamma_0 - \sum_{j=1}^{p} x_{ij}\gamma_j)^2 + \lambda \sum_{j=1}^{p} \frac{|\gamma_j|}{|\hat{\beta}_j|} \quad \text{s.t. } \text{sign}(\gamma_j) = \text{sign}(\hat{\beta}_j) \ \forall j \tag{6}$$

The exact equivalence is obtained with $\hat{\gamma}_j = \hat{\theta}_j \hat{\beta}_j$, and $\hat{\gamma}_0 = \hat{\theta}_0 + \sum_{\ell=1}^{p} \hat{\theta}_\ell \hat{\beta}_{0\ell}$.

The proof is in the Appendix. However our experiments with this version of adaptive lasso produced models that were not nearly as sparse as the ones using LOO. In fact, in Table 3 we see that enforcing sign constraints led to less sparse models than not enforcing sign constraints!

Finally we note that Candes et al. (2008) propose an iterated version of the adaptive lasso, in which the current solutions are used to define weights for the next iteration. They show that the method can enhance sparsity, especially for signal recovery.

Next we examine the relative performance of uniLasso when we modify the three main aspects of its design:

  (a) the use of LOO features $\hat{\eta}_j^{-i}$ versus non-LOO features $\hat{\eta}_j^i$,

  (b) the sign constraints, and

(c) the scaling in Step 2 by the magnitude of the univariate coefficients, rather than just by their signs.

Note that the adaptive lasso is equivalent to uniLasso using non-LOO estimates and no sign constraints.

Table 3 shows the results of 50 simulation runs with $n = 300$, $p = 1000$ and SNR=1, our "medium-SNR" setting detailed in Section 8.

| | | | **LOO** | |
|---|---|---|---|---|
| | lasso | uniLasso | uniLasso-noSign | uniLasso-noMag |
| MSE | 0.55 | 0.59 | 0.61 | 0.60 |
| se | 0.03 | 0.03 | 0.04 | 0.04 |
| Support | 53.78 | 15.32 | 66.38 | 40.12 |
| se | 4.41 | 1.19 | 12.92 | 6.47 |
| | | | **No LOO** | |
| | lasso | uniLasso | uniLasso-noSign | uniLasso-noMag |
| MSE | 0.55 | 0.61 | 0.57 | 0.61 |
| se | 0.03 | 0.04 | 0.03 | 0.05 |
| Support | 53.78 | 36.84 | 24.48 | 48.98 |
| se | 4.41 | 6.34 | 2.71 | 9.66 |

Table 3: *Results for 50 simulations from the medium-SNR setting. The last two columns use variations of uniLasso:* noSign *removes the non-negative constraint, while* noMag *uses just the sign of the univariate estimates in Step 2 of uniLasso (not the magnitude). The top table uses LOO univariate coefficients (as in uniLasso) while the bottom table uses the usual (non-LOO) versions.*

We see that there is not much difference in test error, but uniLasso produces the sparsest models by a large margin. We also see the curious result that with no LOO, the uniLasso *without* sign constraints selects sparser models than *with* sign constraints.

# 4   Why does the uniLasso produce such sparse solutions?

As we saw in the previous section, the combination of LOO univariate estimates in Step (1) and nonnegative constraints in step(2) seem to be the key to the strong sparsity delivered by uniLasso. What is the explanation for this? [2]

The correlation between the responses $y_i$ and the univariate fitted values $\hat{\beta}_{0j} + \hat{\beta}_j x_{ij}$ is always positive, and is the absolute value of the correlation between $y$ and the $j$th feature. When this correlation is large, then the correlation between $y_i$ and the LOO fitted values $\hat{\beta}_{0j}^{-i} + \hat{\beta}_j^{-i} x_{ij}$ will also tend to be positive. But if the first correlation is positive and small then the second correlation is often negative. This could be explained by the fact that leaving an observation out "pushes" the fitted regression line away from that point, and if there wasn't much correlation to start with that will be enough to tip the correlation of the LOO feature to be negative.[3]

Figure 2 shows an example from our "medium-SNR" setting with $n = 300, p = 1000$. The correlations between $y$ and each feature needs to be larger than about 0.08 in order for the correlation between $y$ and each LOO feature to also be positive.

In Appendix D we show that if the absolute correlation between $y$ and a feature is smaller than approximately $\sqrt{2/n}$, then the correlation with the LOO fit will be negative. Since we have $n = 300$ in this example, this evaluates to 0.082. Asymptotically, this means that the two-sided $p$-value for testing $\rho = 0$ has to greater than about 0.16.

In the second step of the uniLasso algorithm in Section 2.1, we fit a multivariate lasso model using all

---

[2]We thank Chris Habron for this analysis.

[3]In the extreme case of zero correlation between $y$ and the $j$th feature, the univariate fit is $\hat{\beta}_{0j} = \bar{y}$, and one can show that the the correlation between $y$ and the LOO version of $\bar{y}$ is -1!
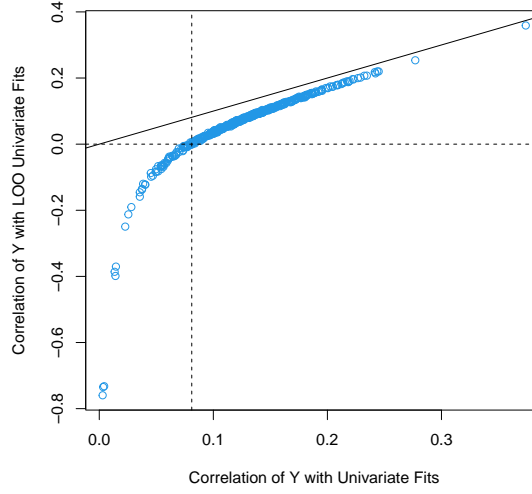
Figure 2: *Medium-SNR setting: correlation between the response and the univariate LOO features, versus the correlation between the response and the usual univariate (non-LOO) features. The solid line represents equality.*

the LOO features, but with a positivity constraint on their coefficients. So even though marginally some LOO features are negatively correlated with the response, this does not guarantee that they will be omitted; however, in practice they do tend to be omitted.

## 5 Application to the car prices data

This dataset consists of 25 features for predicting car prices for 205 cars, taken from `https://www.kaggle.com/datasets/hellbuoy/car-price-prediction` Make of car (22 levels) was first fit to the data and we modeled the residual on the remaining 26 predictors, 11 of which were categorical (and one-hot encoded).

|  | Univariate | Lasso | UniLasso | Polish |
|---|---|---|---|---|
| body.style-3 | -4.768 | -0.224 | 0.000 | 0.000 |
| wheel.base | 0.781 | 0.035 | 0.000 | 0.000 |
| width | 3.082 | 0.475 | 0.378 | 0.000 |
| height | 0.326 | -0.154 | 0.000 | 0.000 |
| curb.weight | 0.014 | 0.004 | 0.002 | -0.001 |
| engine.type-7 | -0.652 | 3.303 | 0.000 | 0.000 |
| num.of.cylinders-2 | 9.867 | -2.546 | 0.000 | 0.000 |
| num.of.cylinders-3 | -13.854 | 0.000 | -0.545 | 0.000 |
| num.of.cylinders-4 | 10.543 | -0.667 | 0.000 | 0.000 |
| num.of.cylinders-7 | -0.652 | 0.009 | 0.000 | 0.000 |
| engine.size | 0.167 | 0.064 | 0.048 | 0.000 |
| fuel.system-2 | -9.143 | 0.208 | 0.000 | 0.000 |
| fuel.system-5 | -0.694 | -1.248 | 0.000 | 0.000 |
| bore | 17.571 | -3.729 | 0.000 | 0.000 |
| horsepower | 0.169 | 0.017 | 0.042 | 0.000 |
| peak.rpm | -0.001 | 0.000 | 0.000 | 0.001 |

Table 4: *Coefficients from various models fit to the car-price data. Coefficients with sign changes relative to the univariate fits are marked in blue.*

Table 4 shows the univariate least squares coefficients on the left, and the lasso and uniLasso coefficients in the middle columns. The right column shows the "uniLasso polish" estimates described

in Section 13. They result from a post-processing of uniLasso in which the lasso is applied to the uniLasso residuals. Only features having a non-zero coefficient in at least one of the three rightmost columns are shown. We see that uniLasso produces a much sparser model than the lasso.

For Figure 3 we took 50 random $2/3 - 1/3$ train-test splits, and computed 4 summaries of model performance. We see that the test errors of all methods are similar, while the support size and number of sign-violations are quite different, as expected. The stability figure is most interesting: for each of the $\binom{50}{2}$ pairs of models, we defined the "stability" as ratio of the number of chosen features in common, divided by the number of unique features in the union of the pair. We see a clear advantage for uniLasso over lasso, where the median proportion of common features is over $80\%$, versus $60\%$ for lasso. The uniLasso "polish", described in Section 13 gives a further improvement in stability, while increasing the support and number of sign violations.
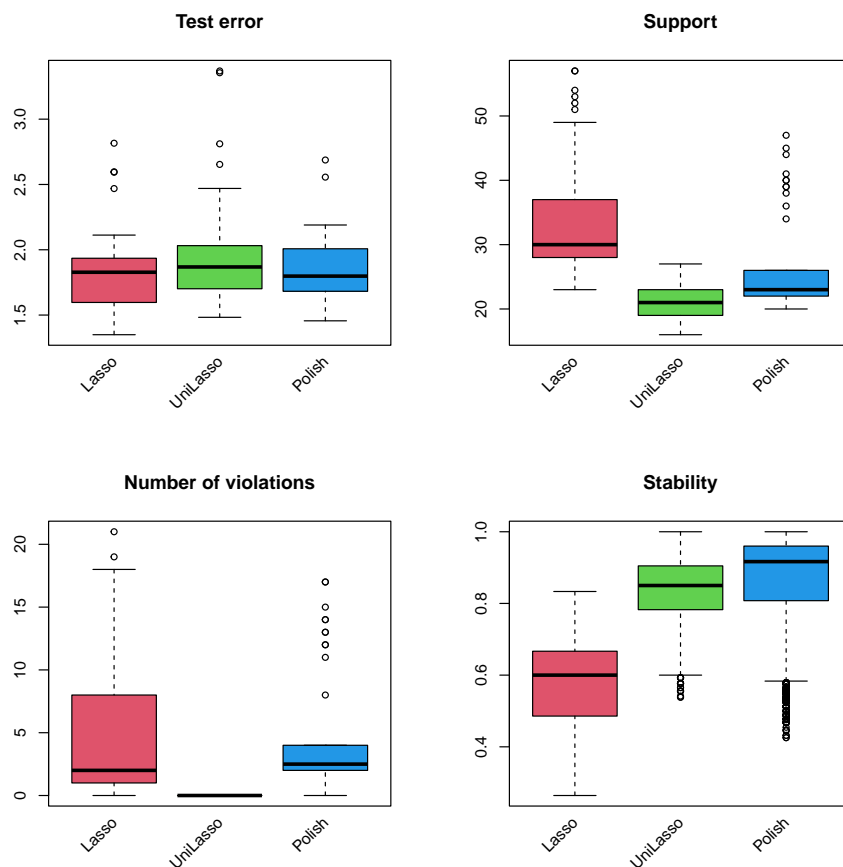


Figure 3: *Results for car prices data over 50 train-test splits. Shown are test set error, size of the chosen model (support), number of sign change violations relative to the univariate signs, and the stability — the average proportion of common features over all model pairs for each method.*

# 6    Analysis of uniLasso with orthogonal features

In this section we derive explicit formula for the uniLasso coefficients in the special case of orthonormal features. Figure 4 shows the lasso and uniLasso paths for a simulated example with an orthonormal feature matrix. They look somewhat different, with the uniLasso path being sparser at any stage along the path (as measured by the $\ell_1$ norm of the coefficients).

It is well-known that in this setting, the lasso coefficients are simply soft-thresholded versions of the univariate least squares estimates. This is (approximately) true of uniLasso, but with a different thresholding function, as we now show.
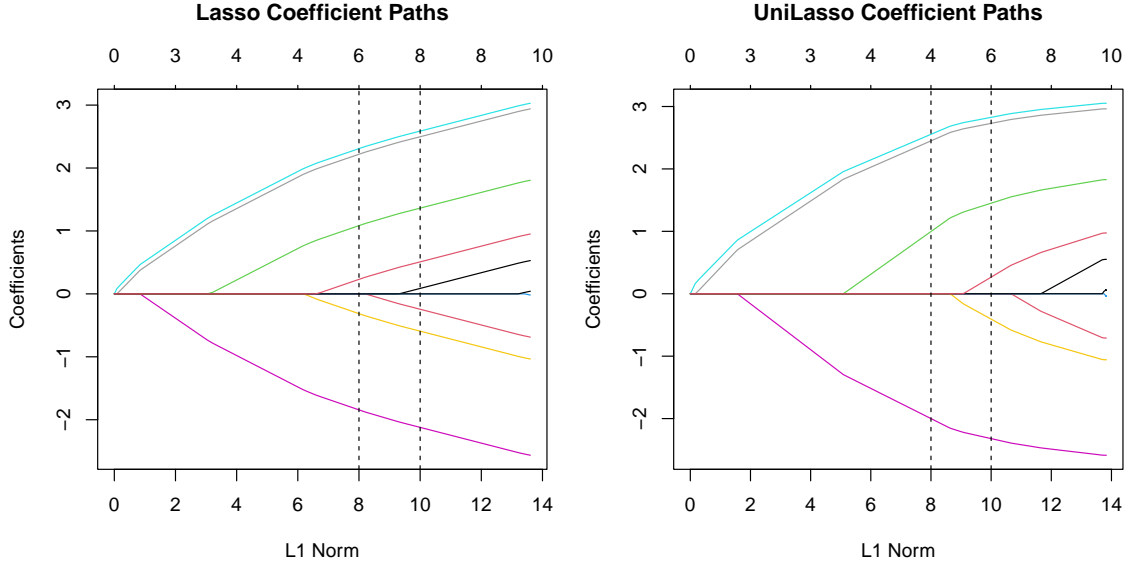
**Figure 4:** *Coefficient paths for lasso and uniLasso, in a simulated example with 10 orthonormal features. The end of the path represents the least squares fit, and in this case also the univariate coefficients. The paths for features with large absolute univariate coefficients look very similar in both plots. Those with small univariate features have a delayed entry in the right plot.*

Suppose $X$ is orthonormal so that each column $x_j$ satisfies $\|x_j\|_2^2 = 1$ and $x_i^T x_j = 0$ for $i \neq j$. Assume also that $\bar{y} = 0$ and $\bar{x}_j = 0$ for each $j$. Note that the least squares coefficients are $\hat{\beta}_{0j} = 0$ and $\hat{\beta}_j = x_j^T y$.

Here we use the actual fits $\hat{\eta}_j^i$ rather than the LOO fits $\hat{\eta}_j^{-i}$ in the second stage — an approximation that simplifies the derivation, and often gives very similar solutions to uniLasso. In this case we have $\hat{\eta}_j^i = \hat{\beta}_j x_{ij}$, since $\hat{\beta}_{0j} = 0$. We can ignore $\theta_0$, which will be zero since $\bar{y} = 0$ and all the $\hat{\eta}_j = \hat{\beta}_j x_j$ have means zero.

The coefficients $\hat{\theta}_j$ are determined by solving:

$$\min_{\theta_j} \left\{ \frac{1}{2} \|y - \sum_{\ell=1}^{p} \theta_\ell \hat{\beta}_\ell x_\ell \|_2^2 + \lambda \|\theta\|_1 \right\}$$

We have the derivative w.r.t. $\theta_j$:

$$-\hat{\beta}_j x_j^\top \left( y - \sum_{\ell=1}^{p} \theta_\ell \hat{\beta}_\ell x_\ell \right) + \lambda \cdot \text{sign}(\theta_j) = 0$$

Since the $x_j$ are orthogonal and unit norm, and using $x_j^\top y = \hat{\beta}_j$, we get

$\hat{\beta}_j^2 \theta_j = \hat{\beta}_j^2 - \lambda \cdot \text{sign}(\theta_j)$, and hence

$$\theta_j = \left( 1 - \frac{\lambda}{\hat{\beta}_j^2} \right)_+ .$$

Hence the final coefficients for $x_j$ are

$$\hat{\gamma}_j = \hat{\theta}_j \cdot \hat{\beta}_j = \hat{\beta}_j \left( 1 - \frac{\lambda}{\hat{\beta}_j^2} \right)_+ = \text{sign}(\hat{\beta}_j) \left( |\hat{\beta}_j| - \frac{\lambda}{|\hat{\beta}_j|} \right)_+ \tag{7}$$
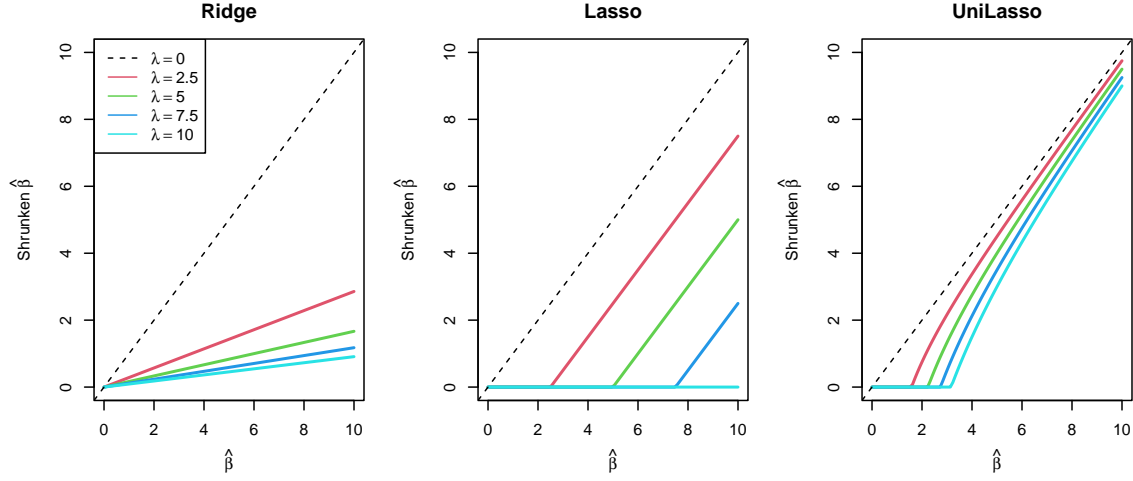
10

Figure 5: *Shrinkage functions for ridge regression, lasso and uniLasso, in a simulated example with 10 orthonormal features.*

The last expression can be compared with the similar expression for the lasso in this situation:

$$\text{sign}(\hat{\beta}_j)\left(|\hat{\beta}_j| - \lambda\right)_+. \tag{8}$$

Figure 5 shows the shrinkage functions for ridge regression, lasso, and uniLasso. Ridge uses proportional shrinkage, while lasso translates all coefficient to zero by the same amount. uniLasso is similar to lasso, except that the larger coefficients are shrunk less than the smaller ones.

The SparseNet procedure (Mazumder et al. 2011) uses a thresholding function that has a roughly similar shape to that of uniLasso, but its objective is not convex. This shrinkage pattern is somewhere between lasso and $\ell_0$ or best subset selection. UniLasso achieves this without losing convexity, a significant computational advantage.

# 7    Theoretical analysis

In this section we study the support recovery and mean-squared error properties of uniLasso. Let $X_1, \ldots, X_p$ be square-integrable random variables with nonzero variance, defined on the same probability space, and let

$$Y = \gamma_0 + \sum_{j \in S} \gamma_j X_j + \epsilon,$$

where $S$ is a subset of $\{1, \ldots, p\}$, $\epsilon$ is a mean zero random variable that is independent of the $X_j$'s, and the $(\gamma_j)_{j \in S}$ are nonzero coefficients. We will refer to $S$ as the *support*. Assume that $Y$ also has nonzero variance. Our data consists of i.i.d. random vectors $(Y_i, X_{i,1}, \ldots, X_{i,p})$, $i = 1, \ldots, n$ (where $n \geq 2$), each having the same distribution as $(Y, X_1, \ldots, X_p)$.

The uniLasso algorithm with penalty parameter $\lambda > 0$ goes as follows. For each $1 \leq i \leq n$ and $1 \leq j \leq p$, let

$$\hat{\beta}_j^{-i} := \frac{\frac{1}{n-1}\sum_{k \neq i} Y_k X_{k,j} - \left(\frac{1}{n-1}\sum_{k \neq i} Y_k\right)\left(\frac{1}{n-1}\sum_{k \neq i} X_{k,j}\right)}{\frac{1}{n-1}\sum_{k \neq i} X_{k,j}^2 - \left(\frac{1}{n-1}\sum_{k \neq i} X_{k,j}\right)^2}$$

be the regression coefficient from the univariate regression of $Y$ on $X_j$ omitting observation $i$, and let

$$\hat{\alpha}_j^{-i} := \frac{1}{n-1}\sum_{k \neq i} Y_k - \frac{\hat{\beta}_j^{-i}}{n-1}\sum_{k \neq i} X_{k,j}$$

be the intercept term. Note that $\hat{\beta}_j^{-i}$ is a consistent estimate of

$$\beta_j := \frac{\text{Cov}(Y, X_j)}{\text{Var}(X_j)},$$

11

and $\hat{\alpha}_j^{-i}$ is a consistent estimate of

$$\alpha_j := \mathrm{E}(Y) - \beta_j \mathrm{E}(X_j).$$

Then, let

$$\hat{Y}_{i,j} := \hat{\alpha}_j^{-i} + \hat{\beta}_j^{-i} X_{i,j}$$

be the predicted value of $Y_i$ from this univariate regression. Next, obtain the estimates $\hat{\theta}_j$, $j = 0, \ldots, p$, by minimizing

$$L(\theta_1, \ldots, \theta_p) = \frac{1}{n} \sum_{i=1}^n (Y_i - \theta_0 - \theta_1 \hat{Y}_{i,1} - \cdots - \theta_p \hat{Y}_{i,p})^2 + \lambda \sum_{j=1}^p \theta_j$$

subject to the constraint that $\theta_j \geq 0$ for each $1 \leq j \leq p$. Finally, define

$$\hat{\gamma}_j := \hat{\theta}_j \hat{\beta}_j,$$

to be the uniLasso estimate of $\gamma_j$ for $1 \leq j \leq p$, where

$$\hat{\beta}_j := \frac{\frac{1}{n} \sum_{k=1}^n Y_k X_{k,j} - (\frac{1}{n} \sum_{k=1}^n Y_k)(\frac{1}{n} \sum_{k=1}^n X_{k,j})}{\frac{1}{n} \sum_{k=1}^n X_{k,j}^2 - (\frac{1}{n} \sum_{k=1}^n X_{k,j})^2}$$

is the regression coefficient from the univariate regression of $Y$ on $X_j$, and let

$$\hat{\gamma}_0 := \hat{\theta}_0 + \sum_{j=1}^p \hat{\theta}_j \hat{\alpha}_j,$$

where

$$\hat{\alpha}_j := \frac{1}{n} \sum_{k=1}^n Y_k - \frac{\hat{\beta}_j}{n} \sum_{k=1}^n X_{k,j}$$

is the intercept term from the univariate regression of $Y$ on $X_j$.

The following theorem shows, roughly speaking, that if (1) $\mathrm{sign}(\gamma_j) = \mathrm{sign}(\beta_j)$ for each $j \in S$, (2) the penalty parameter $\lambda$ is bigger than $|\beta_j|$ for all $j \notin S \cup \{0\}$, and (3) both $\log p$ and $\log n$ are small compared to $n\lambda^2$, then with high probability, $\hat{\gamma}_j = 0$ for all $j \notin S \cup \{0\}$ and $\hat{\gamma}_j = \gamma_j + O(\lambda)$ for all $j \in S \cup \{0\}$.

**Theorem 1.** *Suppose that:*

1. *$\gamma_j \beta_j > 0$ for each $j \in S$.*

2. *The covariance matrix of $(X_j)_{j \in S}$ is nonsingular with minimum eigenvalue $\eta$.*

3. *There is a positive constant $C_0$ such that $\mathrm{Var}(Y) \geq C_0$ and $\mathrm{Var}(X_j) \geq C_0$ for each $j \in S$.*

4. *There are positive constants $C_1$ and $C_2$ such that for each $t \geq 0$ and $1 \leq j \leq p$, $\mathbb{P}(|Y| \geq t)$, $\mathbb{P}(|\epsilon| \geq t)$ and $\mathbb{P}(|X_j| \geq t)$ are bounded above by $C_1 e^{-C_2 t^2}$.*

*Let $M_1 := \max_{j \in S \cup \{0\}} |\gamma_j|$, $M_2 := \min_{j \in S} |\beta_j|$ and $M_3 := \max_{j \in S} |\beta_j|$. Then there are positive constants $K_1, K_2, K_3, K_4, K_5$ depending only on $C_0, C_1, C_2, \eta, M_1, M_2, M_3$ and $|S|$ such that if*

$$K_1 \max_{j \notin S \cup \{0\}} |\beta_j| \leq \lambda \leq K_2,$$

*then*

$$\mathbb{P}(\hat{\gamma}_j = 0 \text{ for all } j \notin S \cup \{0\} \text{ and } |\hat{\gamma}_j - \gamma_j| \leq K_3 \lambda \text{ for all } j \in S \cup \{0\})$$
$$\geq 1 - K_4 p n e^{-K_5 n \lambda^2}.$$

The above theorem is roughly comparable to the available results for the lasso. A close comparison would be, for instance, (Hastie et al. 2015, Theorem 11.3). Like Theorem 1, this theorem also assumes a lower bound on the covariance matrix of the covariates in the support, and the maximum difference between $\hat{\gamma}_j$ and $\gamma_j$ for $j \in S$ is of order $\lambda$. However, there is one key difference. The results about lasso, including the one cited above, require a condition known as *mutual incoherence* or *irrepresentability*. Roughly speaking, it means that if we regress $X_k$ for some $k \notin S$ on $(X_j)_{j \in S}$, the regression coefficients should be small. Notably, our Theorem 1 requires no such relation to hold between the covariates inside and outside the support. All we need is that the univariate regression coefficients of $Y$ on covariates outside the support are small.

Having said that, we make clear that assumption (1) above is a crucial one: namely that $\gamma_j$ and $\beta_j$ have the same sign for each $j \in S$. The next result gives a natural sufficient condition under which this holds.

**Theorem 2.** *Let $\delta_{j,k}$ denote the population value of the univariate regression of $X_j$ on $X_k$. Suppose that for every pair of distinct indices $j, k \in S$, $\delta_{j,k} \geq 0$ if $\gamma_j$ and $\gamma_k$ have the same sign, and $\delta_{j,k} \leq 0$ if $\gamma_j$ and $\gamma_k$ have opposite signs. Then $\beta_j$ is nonzero and has the same sign as $\gamma_j$ for each $j \in S$. Moreover, $|\beta_j| \geq |\gamma_j|$ for each $j \in S$.*

The last result, partly due to Ryan Tibshirani, gives a necessary and sufficient condition for equality of signs.

**Theorem 3.** *Suppose that $Y = \sum_{j \in S} \gamma_j X_j + \epsilon$, where all the $\gamma_j$'s are nonzero, and $\epsilon$ has zero mean, finite variance, and is independent of $(X_j)_{j \in S}$. Let $\beta_j$ be the coefficient of $X_j$ in the univariate (population) regression of $Y$ on $X_j$. Let $\delta_{kj}$ be the coefficient of $X_j$ in the univariate (population) regression of $X_k$ on $X_j$. For each $j$, let $A_j$ be the set of $k \in S$ for which $\mathrm{sign}(\delta_{kj}) = \mathrm{sign}(\gamma_k \gamma_j)$ or $\delta_{kj} = 0$. Then $\beta_j \gamma_j \geq 0$ if and only if*

$$\sum_{k \notin A_j} |\gamma_k \delta_{kj}| \leq \sum_{k \in A_j} |\gamma_k \delta_{kj}|.$$

*In particular, if $A_j = S$, then this holds.*

The proofs of these results are in the Appendix.

# 8  Some simulation results

Figure 6 shows the results of a comparative study of uniLasso with lasso and three other sparse regression methods. Shown are means and standard errors over 100 simulations. The methods are:

- Lasso (using `glmnet`).

- UniLasso

- Polish: a post-processing of the uniLasso solution, described in Section 13.

- Adaptive lasso using penalty weights $1/|\hat{\beta}_j|$ using the univariate coefficients $\hat{\beta}_j$. These first three methods use CV to estimate the value of $\lambda$ that minimizes test error.

- Matching: a variant of the lasso, where we increase the $\lambda$ parameter from the CV-based optimal value, until the support size matches that of uniLasso.

The simulation scenarios are:

1. **Low, medium and high SNR**: Here $n = 300$, $p = 1000$, $N(0,1)$ features with pairwise correlation 0.5. The nonzero regression coefficients are $N(0,1)$; Gaussian errors with SD $\sigma$ chosen to produce low $< 1$, medium $(\approx 1)$, or high $(> 2)$ SNRs.

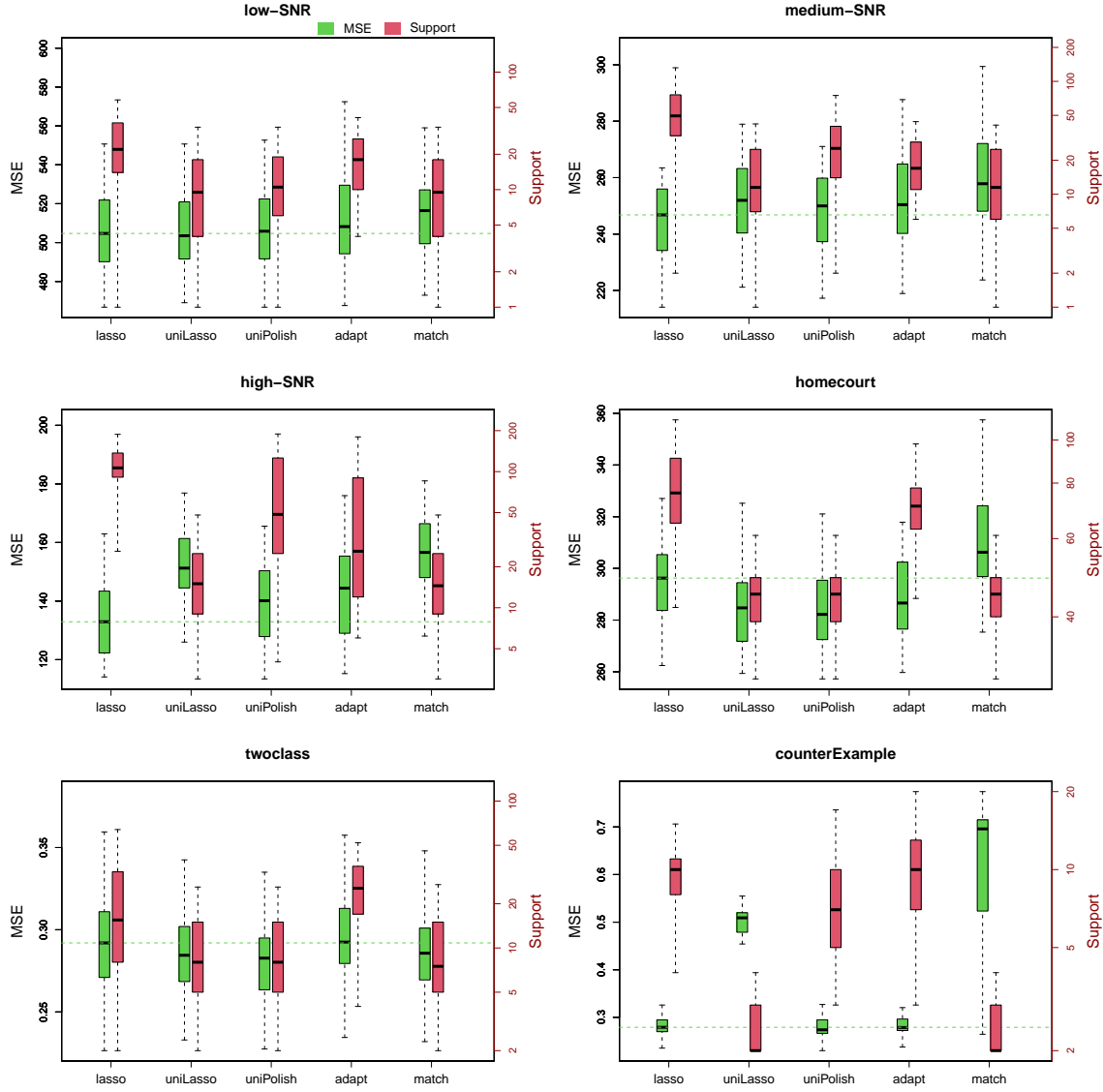2. **Homecourt**: Here $n = 100$, $p = 30$ with details as specified in Section 2.3.

Figure 6: $N = 300$, $p = 1000$: test error (Gaussian) or misclassification rate (Binomial), along with support size.

3. **Two-class**: We set $n = 200$, $p = 500$ with a binary target $y$. The feature covariance within each class is AR(1) with $\rho = 0.8$, and the first 20 features are shifted in the $y = 1$ class by 0.5 units.

4. **Counter-example**: $n = 100$, $p = 20$, $x_1 \sim N(0, 1)$, $x_2 = x_1 + N(0, 1)$, $\beta = (1, -.5, 0, 0, \ldots 0)$, error SD = 0.5.

In the first three scenarios, 10% of the true coefficients are non-zero; for example, in (1), 100 of the 1000 true coefficients are non-zero. Figure 6 shows the test set error (relative to the lasso) and support size (number of non-zero coefficients).

Overall we see that uniLasso shows MSE similar to that of lasso, with often a much smaller support. The same general pattern holds with the real datasets. Exceptions to this are the high-SNR setting where lasso wins handily and the homecourt setting designed to exploit uniLasso's strength (recall Theorem 2). The "counter-example" setting is the classic case where features have positive correlation but the true regression coefficients have opposite signs; not surprisingly uniLasso does poorly. The polish method and adaptive lasso do reasonably well, while the matching method does poorly.

Of course in practice one has cross-validation to help determine which method is preferred in a given example.

14

Figure 7 shows the corresponding results for the simulated examples with $N = 300$ and $p = 100$.
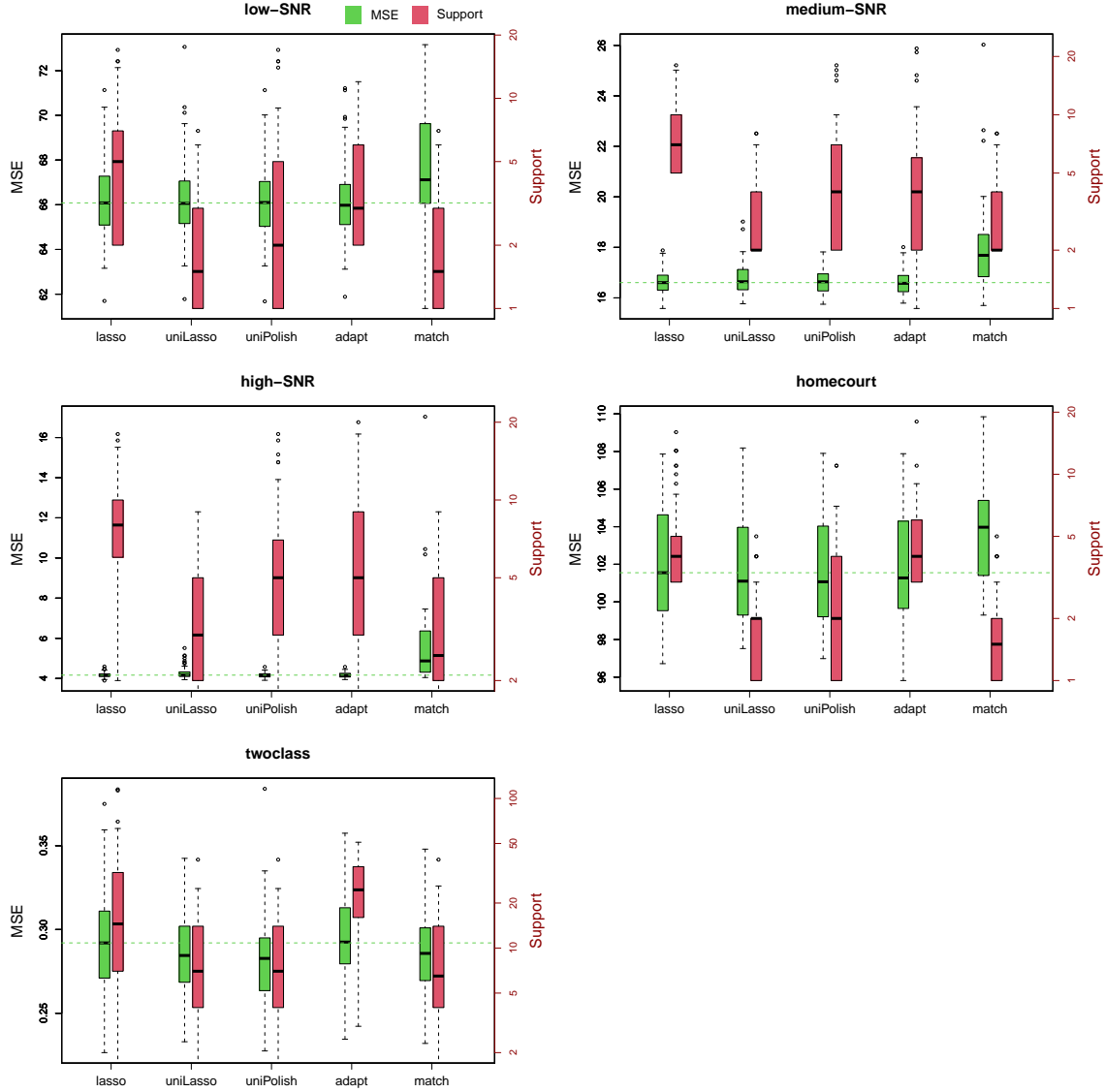


Figure 7: $N = 300 p = 100$: *Test error (Gaussian) or misclassification rate (Binomial), along with support size.*

# 9 The unregularized setting: uniReg

Here we consider the case where $\lambda = 0$, so that there is no $\ell_1$ penalty, as detailed in equation (1). While the main use case has $n > p$, we note that uniReg can be defined for $p > n$ by taking $\lambda \to 0$ in uniLasso. Computationally, this will approximately give the solution with minimum $\ell_1$ norm, in the same way that the limiting lasso fit gives the least squares solution with minimum $\ell_1$ norm. Conveniently even without the $\ell_1$ penalty, the non-negativity constraint still promotes sparsity in the solution.

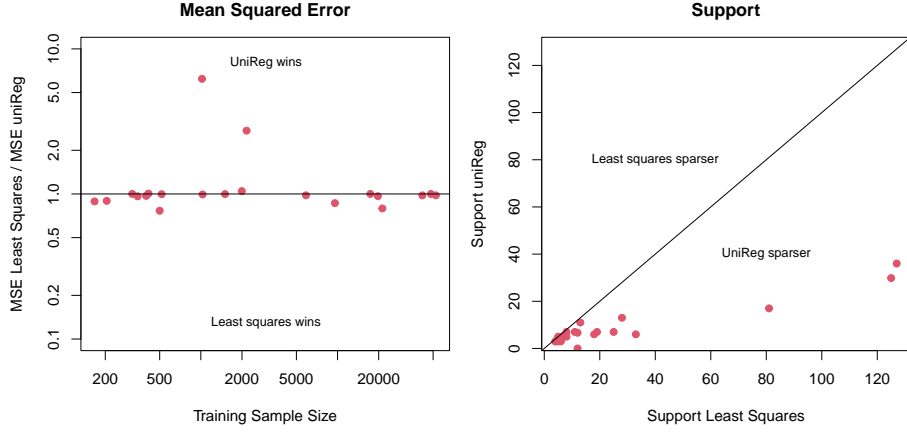In this section we demonstrate the performance of uniReg on simulated and real data.

Figure 8: *LS vs uniReg on regression datasets from the UCI database. All datasets have $n > p$. The left plot shows the MSE ratio on the log scale. The right panel shows the support sizes of the chosen models.*

## 9.1 Simulated data

In Table 5 we compare least squares with uniReg in a subset of the simulated settings defined in Section 8. UniReg performs remarkably well, often winning in MSE and yielding consistently sparser models.

### N=300, p=30

| Setting | p | SNR | MSE-LS | MSE-uniReg | Supp-LS | Supp-uniReg |
| --- | --- | --- | --- | --- | --- | --- |
| low-SNR | 30 | 0.092 | 71.339 | 65.171 | 30 | 4 |
| medium-SNR | 30 | 0.369 | 17.836 | 16.489 | 30 | 4 |
| high-SNR | 30 | 1.475 | 4.459 | 4.178 | 30 | 4 |
| homecourt | 30 | 0.059 | 111.771 | 101.542 | 30 | 4 |

### N=300, p=100

| Setting | p | SNR | MSE-LS | MSE-uniReg | Supp-LS | Supp-uniReg |
| --- | --- | --- | --- | --- | --- | --- |
| low-SNR | 100 | 0.318 | 97.678 | 69.655 | 100 | 8 |
| medium-SNR | 100 | 1.272 | 24.428 | 18.559 | 100 | 9 |
| high-SNR | 100 | 5.087 | 6.104 | 5.82 | 100 | 7 |
| homecourt | 100 | 0.228 | 149.683 | 104.979 | 100 | 13 |

### N=300, p=1000

| Setting | p | SNR | MSE-LS | MSE-uniReg | Supp-LS | Supp-uniReg |
| --- | --- | --- | --- | --- | --- | --- |
| low-SNR | 1000 | 0.578 | 842.191 | 510.837 | 300 | 24.5 |
| medium-SNR | 1000 | 1.48 | 374.915 | 249.076 | 300 | 27 |
| high-SNR | 1000 | 3.613 | 186.171 | 147.276 | 300 | 27 |
| homecourt | 1000 | 1.53 | 577.852 | 420.508 | 300 | 135.5 |

Table 5: *Comparison of least squares with uniReg in a subset of simulated settings defined earlier.*

## 9.2 Real data

We compared least squares and uniReg on the regression datasets from the UCI database[4]. There were 23 datasets in all. Some had multiple targets: we treated each target separately and averaged the results. We randomly sampled each dataset into a 70/30 training/test split, and report the test set MSE and support in Figure 8. One dataset had a small sample ($n = 60$) which produced a very large MSE from least squares; we omit this from the plot for visibility.

The two methods perform similarly in MSE; but uniReg has much smaller support, averaging about

---

[4]https://archive.ics.uci.edu/datasets

1/3 that of the least squares support (which is the number of features $p$).

## 9.3 Theory for uniReg without LOO

Our data consists of $(X_i, Y_i)$, $i = 1, \ldots, n$, where

$$Y_i = \beta_0 + X_i^T \beta + \epsilon_i,$$

where $\beta_0 \in \mathbb{R}$, $\beta = (\beta_1, \ldots, \beta_p) \in \mathbb{R}^p$, $X_i$ are i.i.d. $N_p(0, \Sigma)$ random vectors, where $\Sigma$ is a positive definite matrix with minimum eigenvalue $\lambda_0$ and maximum eigenvalue $\lambda_1$, and $\epsilon_i$ are i.i.d. $N(0, \sigma^2)$ random variables for some $\sigma > 0$. Let $r := p/n$. We assume that $r \in (0, r_0)$ for some $r_0 < 1$.

Let $\hat{\beta}_0^{\text{OLS}}$ and $\hat{\beta}^{\text{OLS}}$ denote the OLS estimates of $\beta_0$ and $\beta$. Let $\hat{\beta}_0^{\text{UR}}$ and $\hat{\beta}^{\text{UR}}$ denote the uniReg estimates of $\beta_0$ and $\beta$ *without leave-one-out*, defined as follows. First, compute the univariate coefficients $\hat{\beta}_j^{\text{uni}}$, $j = 1, \ldots, p$, as the coefficient of $X_j$ when $Y$ is regressed solely on $X_j$ with an intercept term. The estimates $\hat{\beta}_0^{\text{UR}}$ and $\hat{\beta}^{\text{UR}}$ are obtained by minimizing

$$\sum_{i=1}^n (Y_i - a - X_i^T b)^2$$

over $a \in \mathbb{R}$ and $b \in \mathbb{R}^p$ subject to the constraint that $b_j \hat{\beta}_j^{\text{uni}} \geq 0$ for $j = 1, \ldots, p$. Let $\beta_{0,j}^{\text{uni}}$ and $\beta_j^{\text{uni}}$ be the population univariate coefficients; that is,

$$E(Y_i | X_{i,j}) = \beta_{0,j}^{\text{uni}} + \beta_j^{\text{uni}} X_{i,j}.$$

Define three vectors $\mu, \hat{\mu}^{\text{OLS}}, \hat{\mu}^{\text{UR}} \in \mathbb{R}^n$ as

$$\mu_i := E(Y_i | X_i) = \beta_0 + X_i^T \beta,$$
$$\hat{\mu}_i^{\text{OLS}} := \hat{\beta}_0^{\text{OLS}} + X_i^T \hat{\beta}^{\text{OLS}},$$
$$\hat{\mu}_i^{\text{UR}} := \hat{\beta}_0^{\text{UR}} + X_i^T \hat{\beta}^{\text{UR}},$$

for $i = 1, \ldots, n$. Let $q$ be the number of $j$ such that $\beta_j = 0$ and $\beta_j^{\text{uni}} \neq 0$. The following theorem shows that if $\beta_j$ and $\beta_j^{\text{uni}}$ have the same sign for all $j$, then $E\|\mu - \hat{\mu}^{\text{UR}}\|^2$ is less than or equal to $E\|\mu - \hat{\mu}^{\text{OLS}}\|^2$ minus a constant times $q$.

**Theorem 4.** *Suppose that $\beta_j^{\text{uni}} \beta_j \geq 0$ for each $1 \leq j \leq p$. Let $\delta$ be the minimum of $|\beta_j^{\text{uni}}|$ over $1 \leq j \leq p$ such that $\beta_j^{\text{uni}} \neq 0$. Then there are positive constants $C_0, C_1$ depending only on $\lambda_0, \lambda_1, \delta$, $r_0$, and $\sigma$ such that if $n \geq C_0$, then*

$$E\|\mu - \hat{\mu}^{\text{UR}}\|^2 \leq E\|\mu - \hat{\mu}^{\text{OLS}}\|^2 - C_1 q.$$

The theorem is proved in Appendix C. We have not been able to prove the LOO version of uniReg (that is, the actual version that we have proposed). We leave it as an open question.

# 10 Real data examples for uniLasso

Table 6 gives a high level summary of six datasets that we examined in this study. Figure 9 shows the results: these are test error and support size averaged over 100 train/test 50-50 splits.

We see the same general trend as in the simulated examples: uniLasso tends to give test error similar to that of lasso, with often smaller support size.

# 11 Multiclass models

Here we consider a classification problem with more than two classes. The multinomial model is natural for the multiclass problem, but does not generalize easily to uniLasso. The reason is that

| Dataset | n | p | Feature type | Outcome type |
|---|---|---|---|---|
| Spam | 2301 | 57 | mixed | binary |
| NRTI | 1005 | 211 | binary mutations | drug effectiveness (quantitative) |
| Leukemia | 72 | 7129 | gene expression (quantitative) | binary |
| DLBCL | 240 | 1000 | gene expression (quantitative) | survival |
| Breast Cancer | 157 | 1000 | gene expression (quantitative) | survival |
| Ovarian | 190 | 2238 | Mass spec peak intensities | binary |

Table 6: *Summary of datasets used in our study.*



Figure 9: *Real data examples: Blue lines: relative test error (Gaussian) misclassification rate (Binomial) or deviance (Cox), with support size superimposed in gold.*

for each feature there are coefficients for each class, so it is not clear how to replace the feature with a LOO version. Furthermore the coefficients per class are only determined up to a shift, and hence a non-negativity constraint does not make sense. Instead we take a "one-versus rest" (OVR) approach, where the uniLasso binary classifier is used to classify each class from the others. Hence each class gets a predicted probability, and with OVR one classifies to the class with the highest probability as the predicted class.

We applied this to a dataset on childhood cancer, with 63 training and 25 test samples in 4 classes (Khan et al. 2001). With the original train/test split we get zero test errors with `glmnet`. We repeated the train/test split 50 times and obtained these results:

# 12  The use of external univariate scores

Consider the setting where we have our training set `T` and also external data `E` from the same domain (e.g. disease), We assume that `E` does not contain raw data but only summary results. Specifically, `E` contains just univariate coefficients and standard errors for each feature. This setting occurs fairly often in biomedical settings where investigators are not willing to share their raw data, but do publish and share summary results.

We demonstrate here how uniLasso can make productive use of external scores. The idea is to use the univariate coefficients from `E`, rather than computing the LOO estimates from the training set.

| | Ave # Errors | SD | Ave Support | SD |
|---|---|---|---|---|
| Lasso multinomial | 0.96 | 0.03 | 18.68 | 0.08 |
| UniLasso OVR | 0.18 | 0.01 | 27.02 | 0.05 |
| Lasso OVR | 0.22 | 0.01 | 58.08 | 0.08 |

| Number of Errors | 0 | 1 | 2 | 3 | 4 | 7 |
|---|---|---|---|---|---|---|
| Lasso multinomial | 27 | 9 | 8 | 4 | 1 | 1 |
| UniLasso OVR | 41 | 9 | | | | |
| Lasso OVR | 39 | 11 | | | | |

Table 7: *Results for the four-class cancer example. Overall uniLasso OVR makes the fewest misclassification errors averaged over the 50 train-test splits. Of the 50 train(63)/test(25) splits, 41 made zero errors on the 25 test data points, and 9 made 1 error, averaging 0.18. Lasso OVR is sligthly worse, and lasso multinomial makes almost one error on average per train/test split. On the other hand lasso multinomial yields the sparsest model, since the other two select features separately for each class.*

Then in step 2 of uniLasso, we proceed as usual.

To investigate this scheme, we generated data $\mathtt{T}$ with $n = 300$, $p = 1000$, SNR $= 1.5$, with feature covariance AR(1), with $\rho = 0.8$ and the non-zero regression coefficients distributed as $U(0.5, 2)$, and positioned on every other feature $(1, 3, \ldots 99)$. We also generated 600 extra samples from the same distribution to serve as an external data set $\mathtt{E}$, and in addition a large test set. Figure 10 shows the test-set MSE over 200 simulations from the following strategies:

1. lasso and uniLasso, both applied just to the training set;

2. uniLasso+$m$ (green) where the univariate scores are derived from $m$ extra samples from $\mathtt{E}$, and are used in place of the LOO estimates;
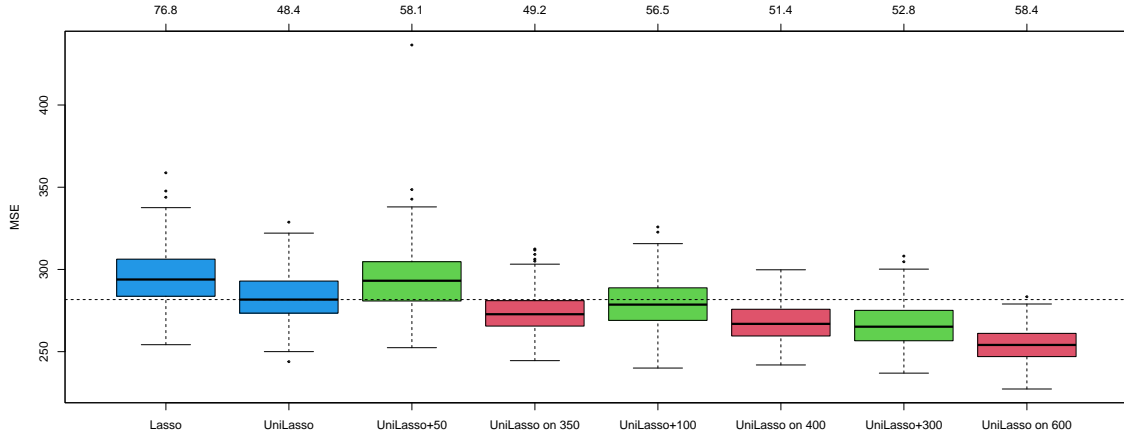
3. unilLasso (red) on all $300 + m$ samples;



Figure 10: *Results from an experiment examining the use of external data with uniLasso. The initial training set has 300 observations and we generated 600 additional samples. Shown in blue are the test set MSE for the lasso, and uniLasso, both applied just to the training set; uniLasso+$m$ (green), where the univariate scores derived from the $m$ extra samples are used in place of the LOO estimates, and unilLasso (red) on $300 + m$ samples for $m = 50, 100 \ldots 300$.*

We see that adding the scores from just 50 external observations does slightly worse than uniLasso on the original data, but builds sparser models than lasso at no cost in test MSE. If instead, we use the extra 50 samples as in strategy 3, the performance improves.

If we increase $m$ to 100, both strategies 2 and 3 improve, but uniLasso using the augmented data performs better each time.

So it appears strategy 2 eventually (as we increase $m$) outperforms lasso and uniLasso both fit on only the original 300 observations. But if we have the raw external data, strategy 3 is the clear winner.

This is one simulation scenario, and it is of course possible that in others the pattern may not be as clear.

# 13  A uniLasso polish

As we have seen earlier, the uniLasso procedure can sometimes perform considerably worse than the lasso, especially in cases where its sign constraint is too restrictive. An example is the "counter-example" problem in Section 8 , where two positively correlated features have different coefficient signs in the population multivariate model. While problematic, the CV estimate of error will alert the user that uniLasso is not working well in a given problem. In this section, we propose a simple "post-processing" of the uniLasso solution that can help to diagnose shortcomings of uniLasso and remedy them.

The idea is very simple: we run the usual lasso (with-cross-validation) to the residual from the uniLasso fit, and then "stitch" together the two solutions. Here is the algorithm in detail, using the `glmnet` R package terminology:

1. Run `cv.uniLasso` on $X, y$ , and extract the fitted linear predictor $\hat{y}$ at $\hat{\lambda}_{min}$.

2. Apply `cv.glmnet` on $X, y$ , with offset $\hat{y}$, giving final predictions $\hat{y}_g$.

3. "Stitch" together the two solutions to produce a single path, starting at the original uniLasso solution.

Note that in the Gaussian model the use of the offset is equivalent to applying the lasso to the residual in step 2; expressed as an offset it allows application of the idea to other GLM families.
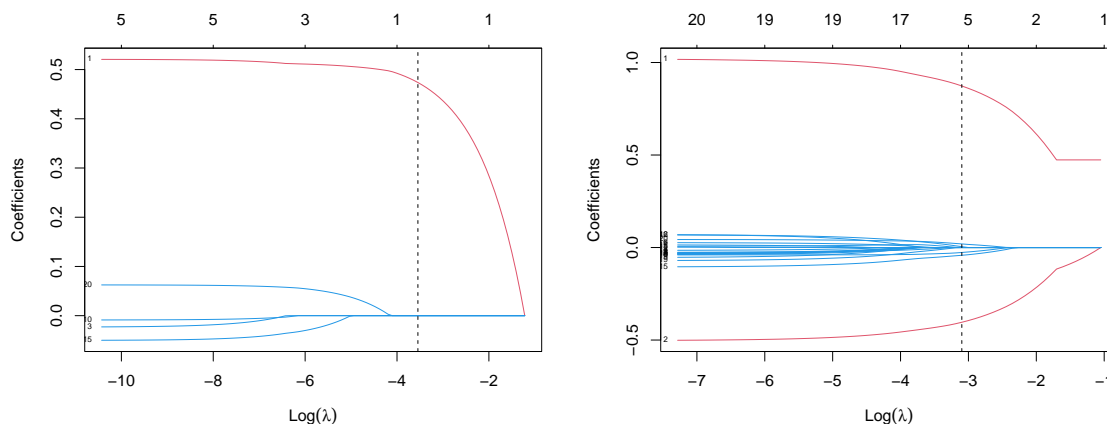


Figure 11: *UniLasso polish applied to the counter-example problem of Section 8. The left panel shows the* `uniLasso` *solution path, with the dashed vertical line indicating the solution chosen by cross-validation. The right panel shows the* `polish.uniLasso` *solution path, which starts at the* `uniLasso`, *and the vertical dashed line indicates the model chosen by cross-validation. The red curves are the support features $x_1$ and $x_2$.*

Figure 11 shows the results of the uniLasso polish applied to the counter-example problem of Section 8. There are two support features ($x_1$ and $x_2$). In the left panel uniLasso enters just feature $x_1$, because $x_2$ has a positive univariate coefficient but a negative multivariate coefficient. The right

panel shows the "polished" path. Via the offset it starts at the uniLasso solution (chosen by CV), and immediately enters $x_2$, and eventually growing the coefficients of both $x_1$ and $x_2$. In the process, the model lets in some of the noise variables.

In Section 8 we included the uniLasso polish in all of the settings, and its performance overall is excellent.

# 14    Does CV work in our two stage uniLasso procedure?

In the uniLasso algorithm we use the target $y$ in both steps, but for (significant) computational convenience, we only apply cross-validation in the 2nd step. Thus one should be concerned that cross-validation may not perform well here. In the simulations of Section 8 we used this form of cross-validation to choose the $\lambda$ tuning parameter, and it seems to have done a reasonably good job.
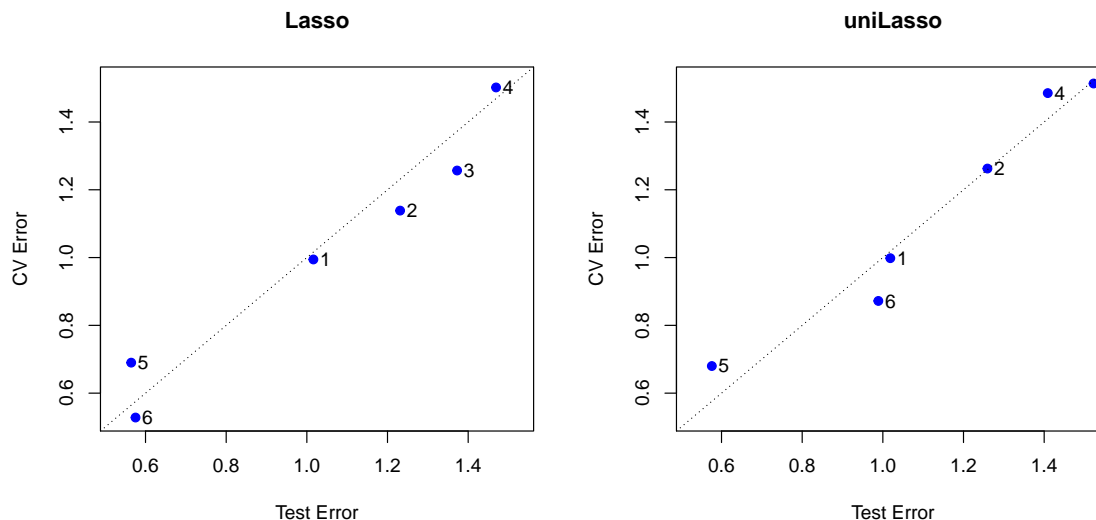


Figure 12: *CV error at the selected value of $\lambda$ versus the test error, for the lasso (left) and uniLasso (right). Each datapoint represents one of our simulated datasets with $p > N$ studied earlier. Both CV estimates track the test error reasonably well.*

A related question is whether the error reported by CV is a good estimate of test error, especially at the selected value of $\lambda$. Figure 12 shows the CV error at the selected value of $\lambda$ versus the test error, for the lasso (left) and uniLasso (right). We see that both CV estimates are quite good. It seems clear that our use of cross-validation in uniLasso is not a major concern.

Finally, a referee asked whether a reason for uniLasso's strong performance was its ability to estimate the the CV tuning parameter $\lambda$. For the car price data, Figure 13 shows the ratio of the achieved test error using CV divided by the minimum test error over the $\lambda$ path, for each of the 50 train/test splits. We see that both lasso and uniLasso do an (equally) good job of estimating the optimal $\lambda$. The average ratio $\lambda_{cv}/\lambda_{opt}$ was 0.87 and 0.38 for the lasso and uniLasso, respectively.

# 15    UniLasso for GLMs and the Cox survival model

Our discussion of uniLasso has focussed on squared-error loss. However, we can use the same idea in any scenario where we are able to fit a lasso model and produce a scalar prediction function $\eta(x)$. This includes in particular all generalized linear models, and the Cox proportional hazards model. These models are included in `glmnet`.

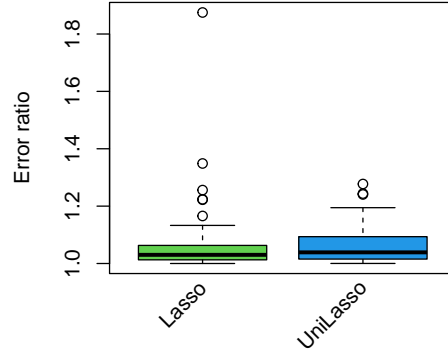The steps are almost the same as before (we will discuss for GLMs):

Figure 13: *Car price data: ratio of the achieved test error using CV divided by the minimum test error over the $\lambda$ path.*

1. Fit the $p$ univariate GLMS, and produce the linear predictor functions $\hat{\eta}_j^i = \hat{\beta}_{0j} + \hat{\beta}_j x_{ij}$. Compute the LOO predictions $\hat{\eta}_j^{-i}$ for the $n$ training observations for each of these.

2. Using these LOO predictions as features, fit a non-negative lasso GLM to the response, yielding coefficients $\hat{\theta} = (\hat{\theta}_0, \hat{\theta}_1, \cdots, \hat{\theta}_p)^\top$.

3. Return the composite estimated linear predictor

$$\hat{\eta}(x) = \hat{\theta}_0 + \sum_{j=1}^p \hat{\theta}_j \hat{\eta}_j(x_j), \tag{9}$$

which collapses as before to a linear model $\hat{\eta}(x) = \hat{\gamma}_0 + \sum_{j=1}^p \hat{\gamma}_j x_j$.

Step 1 can be computationally challenging. Estimating the $p$ univariate functions $\hat{\eta}_j(x_j)$ is fairly straightforward, since each require a few very simple Newton steps. The LOO estimates are more challenging. While for squared-error loss, we have simple formulas for computing these exactly, this is not the case for nonlinear models. However, there are good approximations available that are computationally manageable. We discuss all the computational aspects in the next section.

# 16 Efficient computations for uniLasso

## 16.1 Least squares

We first discuss the computations for squared-error loss. For feature $j$ we have to fit a univariate linear model, and then compute its LOO predictions. We are required to fit the model $\eta_j(x_j) = \beta_{0j} + \beta_j x_j$ using the data $\{x_{ij}, y_i\}_1^n$. This task is simplified if we standardized $x_{ij}$ to have mean zero and unit variance: $z_{ij} = (x_{ij} - \bar{x}_j)/s_j$, where $\bar{x}_j = \frac{1}{n}\sum_{i=1}^n x_{ij}$, and $s_j^2 = \frac{1}{n}\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$. Then the least squares estimates using the $z_{ij}$ are $(\hat{\delta}_{0j}, \hat{\delta}_j) = (\bar{y}, \frac{1}{n}\sum_{i=1}^n z_{ij}y_i)$. These are mapped back to least squares estimates for $x_{ij}$ via $(\hat{\beta}_{0j}, \hat{\beta}_j) = (\hat{\delta}_{oj} - \hat{\delta}_j \bar{x}_j/s_j, \hat{\delta}_j/s_j)$.

This standardized form is particularly useful for computing the LOO fits, since the fits are invariant under affine transformations of the features (so we can use the $z_{ij}$ instead of the $x_{ij}$). We have the classical formula for the LOO residual in linear regression

$$y_i - \hat{\eta}_j^{-i} = \frac{y_i - \hat{\eta}_j^i}{1 - H_{ii}}, \tag{10}$$

where $H_{ii}$ is the $i$th diagonal entry of the *hat* matrix. It is straightforward to show that using the $z_{ij}$ gives $H_{ii} = 1/n + z_{ij}^2/n$. Now simple algebra backs out an expression for $\hat{\eta}_j^{-i}$.

These operations can be performed efficiently for all $i$ and $p$ simultaneously using matrix operations in R, without the need for any loops. In particular we can perform *Hadamard* (elementwise) arithmetic operations (eg multiplication) of like-sized matrices in single operations.

## 16.2   Generalized linear models

The discussion here also applies to the Cox model. Fitting a GLM by maximum likelihood is typically done using a Newton algorithm. This iterative algorithm amounts to making a quadratic approximation to the negative log-likelihood at the current solution, and then solving the quadratic problem to get the updated solution. For GLMs this can be cast as *iteratively reweighted least squares* (IRLS).

Given the fitted linear predictor vector $\eta_j^{(\ell)}$ at iteration $\ell$, one forms a *working* response vector $z_j^{(\ell)}$ that depends on $\eta_j^{(\ell)}$ and $y$ and other properties of the particular GLM family, and an observation weight vector $w_j^{(\ell)}$, and fits an updated model $\eta_j^{(\ell+1)}$ by weighted least squares of $z_j^{(\ell)}$ on $x_j$ with weights $w_j^{(\ell)}$. Usually 4 iterations are sufficient.

Once again we can fit all the univariate GLMs using matrix operations at the same time — i.e. we perform each IRLS step simultaneously for all $p$ univariate GLMS. The expressions are only slightly more complex than in the unweighted case.

What about the LOO fits $\hat\eta_j^{-i}$ for each univariate GLM? Here we make an approximation, as recommended by Rad & Maleki (2020), which amounts to using the final weighted least squares IRLS iteration when fitting the models $\hat\eta_j(x_j)$. A simple formula is available there, similar to what we got before, except a bit more detailed to accommodate observation weights. Again we can use Hadamard matrix operations to do this simultaneously for all $i$ and $j$.

For the Cox model, we make a further approximation, since the implied weight matrix in IRLS is not diagonal. We simply use the diagonal which leads to an approximate Newton algorithm. Note that for the Cox model the intercept is always zero.

## 16.3   Software in R

We provide a R package `uniLasso` for fitting the models described in this paper. It mirrors the behavior of the `glmnet` package for fitting lasso models. The function `cv.uniLasso` has arguments that include all those for `cv.glmnet` plus a few extras. This function does the following, currently for `"binomial"`, `"gaussian"` and `"cox"` families. It does all the work outlined in Sections 16.1–16.2, and performs the second stage `cv.glmnet` using the (approximate) LOO $\hat\eta_j^{-i}$ as features. A `cv.glmnet` object contains information about suggested values of $\lambda$, as well as a fitted `glmnet` model fit to all the data with solutions at all values of $\lambda$. This object is referenced when one predicts from a `cv.glmnet` object. In this case we make sure that the coefficients that are stored on this objects for each $\lambda$ are the collapsed versions as in (9).

Hence `cv.uniLasso` returns a bona-fide `cv.glmnet` object, for which there are a number of methods for plotting, printing and making predictions.

`cv.uniLasso` has three important arguments:

1. `loo = TRUE`: by default it uses the $\hat\eta_j^{-i}$ as features. If this is set to `FALSE` it will use the univariate fitted values $\hat\eta_j^i$ instead.

2. `lower.limits = 0`: this default choice guarantees that the $\hat\theta_j$ are all non-negative. We recommend this in order to get a sparser and more interpretable model.

3. `standardize = FALSE`: this default is the most sensible to use, since part of the point here is to boost strong variables through their scale. Standardizing would undo that.

In addition we have `uniReg` and `cv.uniReg` for fitting the uniReg model, as well as `ci.uniReg` for computing bootstrap intervals for each coefficient. The function `polish.uniLasso` does the polishing as described in Section!13. We also provide simulation functions for generating the data as used in this paper.

# 17 Discussion

In this paper we have introduced a novel method for sparse regression that "stacks" univariate regressions using a non-negative version of the lasso. The procedure has interesting properties, namely that the final features weights have the same signs as the univariate coefficients and tend to be larger when the univariate coefficients are large. The test set MSE of the method is often similar to that of the lasso, with substantially smaller support and lower false positive rate.

There are many possible extensions for this work, for example, to random forests, gradient boosting and neural networks. These are topics for future study.

Current versions of the uniLasso package can be found at

`https://github.com/trevorhastie/uniLasso` and CRAN (R)
`https://github.com/sophial05/uni-lasso` (Python, written by Sophia Lu).

## Disclosure statement

The authors have no conflicts of interest to declare.

## Acknowledgments

# A  Proof of Proposition 1

If you substitute $\eta_j^i = \hat{\beta}_{0j} + \hat{\beta}_j x_{ij}$ in the objective in (5) you get

$$\sum_{i=1}^n (y_i - (\theta_0 + \sum_{j=1}^p \hat{\beta}_{0j}) - \sum_{j=1}^p x_{ij}\hat{\beta}_j\theta_j)^2 + \lambda \sum_{j=1}^p |\theta_j|$$

Now you change variables to $\gamma_0 = \theta_0 + \sum_{j=1}^p \hat{\beta}_{0j}$ and $\gamma_j = \hat{\beta}_j\theta_j$.

The result is immediate, since

$$|\theta_j| = \frac{|\gamma_j|}{|\hat{\beta}_j|}$$

and

$$\{\text{sign}(\gamma_j) = \text{sign}(\hat{\beta}_j) \ \forall j \ \} \equiv \{ \ \theta_j \geq 0 \ \forall j\}$$

$\square$

# B    Proof of theorems in Section 7

We first prove Theorem 1. We will need the following result (see, e.g., (Talagrand 2010, Theorem A.7.1)).

**Lemma 1** (Bernstein's inequality). *Let $Z_1, \ldots, Z_n$ be i.i.d. random variables such that $\mathrm{E}(Z_1) = 0$ and $\mathrm{E}(e^{\beta|Z_1|}) \leq 2$ for some $\beta > 0$. Then for all $t \geq 0$,*

$$\mathbb{P}\left(\left|\frac{1}{n}\sum_{i=1}^{n} Z_i\right| \geq t\right) \leq 2\exp\left(-\min\left\{\frac{n\beta^2 t^2}{4}, \frac{n\beta t}{2}\right\}\right).$$

Define random variables

$$A := \frac{1}{n}\sum_{i=1}^{n} Y_i^2, \quad B_j := \frac{1}{n}\sum_{i=1}^{n} X_{i,j}^2, \quad D_j := \frac{1}{n}\sum_{i=1}^{n} Y_i X_{i,j}.$$

As a corollary of Bernstein's inequality, we obtain the following.

**Corollary 1.** *There is a positive constant $C_3$ depending only on $C_1$ and $C_2$ such that for all $t \geq 0$,*

$$\mathbb{P}(|A - \mathrm{E}(A)| \geq \mathrm{t}) \geq 2\mathrm{e}^{-C_3 n \min\{\mathrm{t}^2, \mathrm{t}\}},$$

*ad the same bound holds for $\mathbb{P}(|B_j - \mathrm{E}(B_j)| \geq \mathrm{t})$ and $\mathbb{P}(|D_j - \mathrm{E}(D_j)| \geq \mathrm{t})$ for all $j$.*

*Proof.* Let $Z_i := Y_i^2 - \mathrm{E}(A)$. Then $\mathrm{E}(Z_i) = 0$, and for any $t \geq 0$,

$$\mathbb{P}(|Z_1| \geq t) \leq \mathbb{P}(Y_1^2 \geq t) \leq C_1 e^{-C_2 t}.$$

Thus, for any $\beta \in (0, C_2)$,

$$\begin{aligned}
\mathrm{E}(e^{\beta|Z_1|}) &= \int_0^\infty \beta e^{\beta t}\mathbb{P}(|Z_1| \geq t)dt \\
&\leq \int_0^\infty \beta e^{\beta t} C_1 e^{-C_2 t} dt \\
&= \frac{C_1 \beta}{C_2 - \beta}.
\end{aligned}$$

Choosing $\beta$ sufficiently small makes the upper bound $\leq 2$, allowing us to apply Bernstein's inequality. A similar argument works for $B_j$. For $D_j$, take $Z_i := Y_i X_{i,j} - \mathrm{E}(D_j)$. Then note that $\mathrm{E}(Z_i) = 0$, and for any $t \geq 2|\mathrm{E}(D_j)|$,

$$\begin{aligned}
\mathbb{P}(|Z_1| \geq t) &\leq \mathbb{P}(|Y_1 X_{1,j}| \geq t/2) \\
&\leq \mathbb{P}(Y_1^2 \geq t/2) + \mathbb{P}(X_{1,j}^2 \geq t/2) \\
&\leq 2C_1 e^{-C_2 t/2}.
\end{aligned}$$

Since $|\mathrm{E}(D_j)|$ can be bounded above by a number that depends only on $C_1$ and $C_2$, this shows that there are constants $C_3$ and $C_4$ depending only on $C_1$ and $C_2$ such that for all $t \geq 0$,

$$\mathbb{P}(|Z_1| \geq t) \leq C_3 e^{-C_4 t}.$$

The proof is now completed by proceeding as before. $\qquad\square$

We also obtain the following second corollary.

**Corollary 2.** *There are positive constants $C_5$, $C_6$ and $C_7$ depending only on $C_0$, $C_1$ and $C_2$ such that for any $t \in [0, C_5]$ and any $i$ and $j$, $\mathbb{P}(|\hat{\beta}_j^{-i} - \beta_j| \geq t)$ and $\mathbb{P}(|\hat{\alpha}_j^{-i} - \alpha_j| \geq t)$ are bounded above by $C_6 e^{-C_7 n t^2}$.*

*Proof.* Fix $i, j$. Let

$$Q_1 := \frac{1}{n-1} \sum_{k \neq i} Y_k X_{k,j}, \quad Q_2 := \frac{1}{n-1} \sum_{k \neq i} Y_k,$$

$$Q_3 := \frac{1}{n-1} \sum_{k \neq i} X_{k,j}, \quad Q_4 := \frac{1}{n-1} \sum_{k \neq i} X_{k,j}^2.$$

Using the same approach as in Corollary 1 and the fact that $n - 1 \geq n/2$ (because $n \geq 2$), we deduce the concentration inequality

$$\mathbb{P}(|Q_l - \mathrm{E}(Q_l)| \geq t) \leq 2e^{-Kn \min\{t^2, t\}} \tag{11}$$

for each $1 \leq l \leq 4$. Now, note that $\mathrm{E}(Q_4) = \mathrm{E}(X_j^2)$ and $\mathrm{E}(Q_3) = \mathrm{E}(X_j)$. Moreover, recall that $\mathrm{Var}(X_j) = \mathrm{E}(X_j^2) - (\mathrm{E}(X_j))^2 \geq C_0 > 0$. Combining these observations with the above inequality, we see that there are positive constants $K_1$ and $K_2$ such that

$$\mathbb{P}(|Q_4 - Q_3^2| < C_0/2) \leq K_1 e^{-K_2 n}. \tag{12}$$

Since

$$\hat{\beta}_j^{-i} = \frac{Q_1 - Q_2 Q_3}{Q_4 - Q_3^2},$$

the inequality (12) gives

$$\mathbb{P}(|\hat{\beta}_j^{-i} - \beta_j| \geq t) = \mathbb{P}\left(\left| \frac{Q_1 - Q_2 Q_3 - (Q_4 - Q_3^2)\beta_j}{Q_4 - Q_3^2} \right| \geq t\right)$$
$$\leq \mathbb{P}(|Q_1 - Q_2 Q_3 - (Q_4 - Q_3^2)\beta_j| \geq C_0 t/2) + K_1 e^{-K_2 n}. \tag{13}$$

Now take any $s > 0$, and suppose that $|Q_l - \mathrm{E}(Q_l)| < s$ for $1 \leq l \leq 4$. Then we also have

$$|Q_3^2 - (\mathrm{E}(Q_3))^2| \leq |Q_3 - \mathrm{E}(Q_3)|^2 + 2|\mathrm{E}(Q_3)||Q_3 - \mathrm{E}(Q_3)| < s^2 + K_3 s,$$

where $K_3$ depends only on $C_1$ and $C_2$. Similarly,

$$|Q_2 Q_3 - \mathrm{E}(Q_2)\mathrm{E}(Q_3)| < K_4 s,$$

where $K_4$ depends only on $C_1$ and $C_2$. Thus, we have

$$|(Q_1 - Q_2 Q_3 - (Q_4 - Q_3^2)\beta_j) - (\mathrm{E}(Q_1) - \mathrm{E}(Q_2)\mathrm{E}(Q_3) - (\mathrm{E}(Q_4) - (\mathrm{E}(Q_3))^2)\beta_j)|$$
$$< K_5 s + K_6 s^2,$$

where $K_5$ and $K_6$ depend only on $C_0, C_1$ and $C_2$. But

$$\mathrm{E}(Q_1) - \mathrm{E}(Q_2)\mathrm{E}(Q_3) - (\mathrm{E}(Q_4) - (\mathrm{E}(Q_3))^2)\beta_j$$
$$= \mathrm{E}(YX_j) - \mathrm{E}(Y)\mathrm{E}(X_j) - (\mathrm{E}(X_j^2) - (\mathrm{E}(X_j))^2)\frac{\mathrm{E}(YX_j) - \mathrm{E}(Y)\mathrm{E}(X_j)}{\mathrm{E}(X_j) - (\mathrm{E}(X_j))^2} = 0.$$

Plugging this into the previous display shows that

$$\mathbb{P}(|Q_1 - Q_2 Q_3 - (Q_4 - Q_3^2)\beta_j| \geq K_5 s + K_6 s^2) \leq \sum_{l=1}^{4} \mathbb{P}(|Q_l - \mathrm{E}(Q_l)| \geq s) \tag{14}$$

Choosing $s$ to solve $K_5 s + K_6 s^2 = C_0 t/2$, and combining equations (11), (13) and (14) yields the proof of the claimed concentration inequality for $\hat{\beta}_j^{-i}$. The proof for $\hat{\alpha}_j^{-i}$ is similar. We omit the details. $\qquad \square$

We are now ready to prove Theorem 1.

*Proof of Theorem 1.* Throughout this proof, $K_1, K_2, \ldots$ will denote positive constants that may depend only on $C_0, C_1, C_2, \theta$ and $|S|$ and nothing else, whose values may change from line to line. Also, we will denote by $o(1)$ any random variable $Z$ such that

$$\mathbb{P}(|Z| \geq t) \leq K_1 n e^{-K_2 n t^2}$$

26

for all $t \in [0, K_3]$ (for some $K_1, K_2, K_3$ according to the above convention), and by $O(1)$ any random variable $Z$ such that $|Z - K_1| = o(1)$ for some $K_1$. Define

$$\tilde{L}(\theta_0, \theta_1, \ldots, \theta_p) := \frac{1}{n} \sum_{i=1}^{n} (Y_i - \theta_0 - \theta_1(\hat{Y}_{i,1} - \hat{\alpha}_1) - \cdots - \theta_p(\hat{Y}_{i,p} - \hat{\alpha}_p))^2 + \lambda \sum_{j=1}^{p} \theta_j.$$

Note that $(\hat{\theta}'_0, \hat{\theta}_1, \ldots, \hat{\theta}_p)$ minimizes $\tilde{L}$ under $\hat{\theta}_j \geq 0$ for all $1 \leq j \leq p$ if and only if $(\hat{\theta}_0, \hat{\theta}_1, \ldots, \hat{\theta}_p)$ minimizes $L$ under the same constraint, where

$$\hat{\theta}_0 = \hat{\theta}'_0 - \sum_{j=1}^{p} \hat{\theta}_j \hat{\alpha}_j.$$

Note that this means $\hat{\gamma}_0 = \hat{\theta}'_0$. We will henceforth consider this modified optimization problem. Note that for each $1 \leq j \leq p$, and for any $\theta_0 \in \mathbb{R}$ and $\theta_1, \ldots, \theta_p \geq 0$,

$$\begin{aligned}
\tilde{L}_j &:= \frac{\partial \tilde{L}}{\partial \theta_j} \\
&= -\frac{2}{n} \sum_{i=1}^{n} (Y_i - \theta_0 - \theta_1(\hat{Y}_{i,1} - \hat{\alpha}_1) - \cdots - \theta_p(\hat{Y}_{i,p} - \hat{\alpha}_p))(\hat{Y}_{i,j} - \hat{\alpha}_j) + \lambda.
\end{aligned} \tag{15}$$

By an application of the Cauchy–Schwarz inequality, we get

$$\begin{aligned}
&\left| \frac{1}{n} \sum_{i=1}^{n} (Y_i - \theta_0 - \theta_1(\hat{Y}_{i,1} - \hat{\alpha}_1) - \cdots - \theta_p(\hat{Y}_{i,p} - \hat{\alpha}_p))(\hat{Y}_{i,j} - \hat{\alpha}_j) \right| \\
&\leq \left( \frac{1}{n} \sum_{i=1}^{n} (Y_i - \theta_0 - \theta_1(\hat{Y}_{i,1} - \hat{\alpha}_1) - \cdots - \theta_p(\hat{Y}_{i,p} - \hat{\alpha}_p))^2 \right)^{1/2} \left( \frac{1}{n} \sum_{i=1}^{n} (\hat{Y}_{i,j} - \hat{\alpha}_j)^2 \right)^{1/2} \\
&\leq \sqrt{\tilde{L}(\theta_0, \ldots, \theta_p)} \left( \frac{1}{n} \sum_{i=1}^{n} (\hat{Y}_{i,j} - \hat{\alpha}_j)^2 \right)^{1/2}.
\end{aligned}$$

Now, note that

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^{n} (\hat{Y}_{i,j} - \hat{\alpha}_j)^2 &\leq \frac{2}{n} \sum_{i=1}^{n} (\hat{Y}_{i,j} - \hat{\alpha}_j^{-i})^2 + \frac{2}{n} \sum_{i=1}^{n} (\hat{\alpha}_j^{-i} - \hat{\alpha}_j)^2 \\
&\leq \frac{2}{n} \sum_{i=1}^{n} \hat{\beta}_{i,j}^2 X_{i,j}^2 + \frac{4}{n} \sum_{i=1}^{n} (\hat{\alpha}_j^{-i} - \alpha_j)^2 + 4(\hat{\alpha}_j - \alpha_j)^2 \\
&\leq \left( \frac{2}{n} \sum_{i=1}^{n} X_{i,j}^2 \right) \max_{1 \leq i \leq n} (\hat{\beta}_j^{-i})^2 + 4 \max\{(\hat{\alpha}_j - \alpha_j)^2, \max_{1 \leq i \leq n} (\hat{\alpha}_j^{-i} - \alpha_j)^2\}.
\end{aligned}$$

Let $M_j^\beta$ and $M_j^\alpha$ denote the two maxima on the right. Combining the previous three displays, we get that for $\theta_0 \in \mathbb{R}$ and $\theta_1, \ldots, \theta_p \geq 0$,

$$\tilde{L}_j(\theta_0, \ldots, \theta_p) \geq \lambda - \sqrt{\tilde{L}(\theta_0, \ldots, \theta_p)(2B_j M_j^\beta + 4M_j^\alpha)}. \tag{16}$$

Take any $j \in \{1, \ldots, p\} \setminus S$. Suppose that $\hat{\gamma}_j \neq 0$. Then $\hat{\theta}_j > 0$. This implies that

$$\tilde{L}_j(\hat{\gamma}_0, \hat{\theta}_1, \ldots, \hat{\theta}_p) = 0,$$

because otherwise, we can slightly perturb $\hat{\theta}_j$ while maintaining the non-negativity constraint and decreasing the value of $L$. By inequality (16), this implies that if $\hat{\theta}_j > 0$, then we must have

$$\tilde{L}(\hat{\gamma}_0, \hat{\theta}_1, \ldots, \hat{\theta}_p)(2B_j M_j^\beta + 4M_j^\alpha) \geq \lambda^2.$$

But note that

$$\tilde{L}(\hat{\gamma}_0, \hat{\theta}_1, \ldots, \hat{\theta}_p) \leq \tilde{L}(0, 0, \ldots, 0) = A.$$

Thus, we conclude that if $\hat{\theta}_j > 0$, then

$$\lambda^2 \le A(2B_j M_j^\beta + 4M_j^\alpha).$$

Thus,

$$\mathbb{P}(\hat{\theta}_j > 0) \le \mathbb{P}(A > 2\mathrm{E}(A)) + \mathbb{P}(\mathrm{B_j} > 2\mathrm{E}(\mathrm{B_j}))$$
$$+ \mathbb{P}\left(M_j^\beta \ge \frac{\lambda^2}{16\mathrm{E}(\mathrm{A})\mathrm{E}(\mathrm{B_j})}\right) + \mathbb{P}\left(M_j^\alpha \ge \frac{\lambda^2}{16\mathrm{E}(\mathrm{A})}\right).$$

By Corollary 1, the first two probabilities on the right are bounded above by $K_1 e^{-K_2 n}$. Next, note that if

$$K_3|\beta_j| \le \lambda \le K_4, \tag{17}$$

for some sufficiently small $K_3$ and $K_4$, then by Corollary 2,

$$\mathbb{P}\left(M_j^\beta \ge \frac{\lambda^2}{16\mathrm{E}(\mathrm{A})\mathrm{E}(\mathrm{B_j})}\right) \le n\mathbb{P}\left(\hat{\beta}_{1,j}^2 \ge \frac{\lambda^2}{16\mathrm{E}(\mathrm{A})\mathrm{E}(\mathrm{B_j})}\right)$$
$$\le n\mathbb{P}(|\hat{\beta}_{1,j} - \beta_j| \ge K_5\lambda) \le K_6 n e^{-K_7 n\lambda^2}.$$

Similarly, the same bound holds for the tail probability of $M_j^\alpha$. Combining, we get that under the condition (17),

$$\mathbb{P}(\hat{\theta}_j > 0) \le K_8 n e^{-K_9 n\lambda^2},$$

and thus,

$$\mathbb{P}(\hat{\gamma}_j \ne 0 \text{ for some } j \notin S \cup \{0\}) \le \mathbb{P}(E^c) \le K_8 p n e^{-K_9 n\lambda^2}, \tag{18}$$

where $E$ denotes the event that $\hat{\theta}_j = 0$ for all $j \notin S \cup \{0\}$. Suppose that $E$ happens. Take any $k \in S$. If $\hat{\theta}_k \ne 0$, then $\tilde{L}_k(\hat{\gamma}_0, \hat{\theta}_1, \ldots, \hat{\theta}_p) = 0$, whereas if $\hat{\theta}_k = 0$, then $\tilde{L}_k(\hat{\gamma}_0, \hat{\theta}_1, \ldots, \hat{\theta}_p) \ge 0$ and $(\gamma_k - \hat{\gamma}_k)/\beta_k = \gamma_k/\beta_k > 0$. Thus, in either case, we have

$$\frac{\gamma_k - \hat{\gamma}_k}{\beta_k} \tilde{L}_k(\hat{\gamma}_0, \hat{\theta}_1, \ldots, \hat{\theta}_p) \ge 0.$$

By the formula (15) for $\tilde{L}_j$, this shows that

$$\frac{\lambda(\gamma_k - \hat{\gamma}_k)}{2\beta_k} \ge \left(\frac{\gamma_k - \hat{\gamma}_k}{\beta_k}\right)\frac{1}{n}\sum_{i=1}^n\left(Y_i - \hat{\gamma}_0 - \sum_{j\in S}\hat{\theta}_j(\hat{Y}_{i,j} - \hat{\alpha}_j)\right)(\hat{Y}_{i,k} - \hat{\alpha}_k)$$
$$= \frac{\gamma_k - \hat{\gamma}_k}{n\beta_k}\sum_{i=1}^n\left[\gamma_0 - \hat{\gamma}_0 + \sum_{j\in S}(\gamma_j X_{i,j} - \hat{\theta}_j(\hat{Y}_{i,j} - \hat{\alpha}_j))\right](\hat{Y}_{i,k} - \hat{\alpha}_k)$$
$$+ \frac{\gamma_k - \hat{\gamma}_k}{n\beta_k}\sum_{i=1}^n\epsilon_i(\hat{Y}_{i,k} - \hat{\alpha}_k). \tag{19}$$

Define

$$M := \max_{1\le i\le n,\, j\in S}(|\hat{\beta}_j^{-i} - \beta_j| + |\hat{\beta}_j - \beta_j|).$$

Then

$$|(\hat{Y}_{i,k} - \hat{\alpha}_k) - \beta_k X_{i,k}| = |\hat{\beta}_{i,k} - \beta_k||X_{i,k}| \le M|X_{i,k}|, \tag{20}$$

and since $\hat{\gamma}_j = \hat{\theta}_j\hat{\beta}_j$,

$$|\hat{\theta}_j(\hat{Y}_{i,j} - \hat{\alpha}_j) - \hat{\gamma}_j X_{i,j}| = |\hat{\theta}_j\hat{\beta}_j^{-i}X_{i,j} - \hat{\theta}_j\hat{\beta}_j X_{i,j}|$$
$$= |\hat{\beta}_j^{-i} - \hat{\beta}_j||\hat{\theta}_j X_{i,j}| \le M|\hat{\theta}_j||X_{i,j}|. \tag{21}$$

By (20) and (21), we have

$$|\hat{\theta}_j(\hat{Y}_{i,j} - \hat{\alpha}_j)(\hat{Y}_{i,k} - \hat{\alpha}_k) - \hat{\gamma}_j\beta_k X_{i,j}X_{i,k}|$$
$$\le |\hat{\theta}_j(\hat{Y}_{i,j} - \hat{\alpha}_j) - \hat{\gamma}_j X_{i,j}||\hat{Y}_{i,k} - \hat{\alpha}_k| + |\hat{\gamma}_j X_{i,j}||(\hat{Y}_{i,k} - \hat{\alpha}_k) - \beta_k X_{i,k}|$$
$$\le M|\hat{\theta}_j||X_{i,j}|(|\beta_k X_{i,k}| + M|X_{i,k}|) + M|\hat{\gamma}_j X_{i,k}X_{i,j}| \tag{22}$$

28

Let $S' := S \cup \{0\}$. Define $X_0 \equiv 1$ and $X_{i,0} \equiv 1$ for $1 \leq i \leq n$. Combining equations (19), (20), (21) and (22), we get

$$\frac{\lambda(\gamma_k - \hat{\gamma}_k)}{2\beta_k} - \left[(\gamma_k - \hat{\gamma}_k)\sum_{j \in S'}(\gamma_j - \hat{\gamma}_j)\left(\frac{1}{n}\sum_{i=1}^{n}X_{i,j}X_{i,k}\right) + (\gamma_k - \hat{\gamma}_k)\frac{1}{n}\sum_{i=1}^{n}\epsilon_i X_{i,k}\right]$$

$$\geq -K_5(M + M^2)\left|\frac{\gamma_k - \hat{\gamma}_k}{n\beta_k}\right|\sum_{i=1}^{n}\left(|\gamma_0 - \hat{\gamma}_0||X_{i,k}|\right.$$

$$\left. + \sum_{j \in S}(1 + |\hat{\theta}_j| + |\hat{\gamma}_j|)(|X_{i,j}| + |\epsilon_i|)|X_{i,k}|\right). \tag{23}$$

Next note that $\tilde{L}_0(\hat{\gamma}_0, \hat{\theta}_1, \ldots, \hat{\theta}_p) = 0$, where $\tilde{L}_0$ denotes the derivative of $\tilde{L}$ is the zeroth coordinate. This means that

$$0 = \frac{1}{n}\sum_{i=1}^{n}\left(Y_i - \hat{\gamma}_0 - \sum_{j \in S}\hat{\theta}_j(\hat{Y}_{i,j} - \hat{\alpha}_j)\right)$$

$$= \frac{1}{n}\sum_{i=1}^{n}\left[\gamma_0 - \hat{\gamma}_0 + \sum_{j \in S}(\gamma_j X_{i,j} - \hat{\theta}_j(\hat{Y}_{i,j} - \hat{\alpha}_j))\right] + \frac{1}{n}\sum_{i=1}^{n}\epsilon_i.$$

Proceeding as before (and recalling that $X_{i,0} = 1$), this gives

$$(\gamma_0 - \hat{\gamma}_0)\sum_{j \in S'}(\gamma_j - \hat{\gamma}_j)\left(\frac{1}{n}\sum_{i=1}^{n}X_{i,j}X_{i,0}\right) + (\gamma_0 - \hat{\gamma}_0)\frac{1}{n}\sum_{i=1}^{n}\epsilon_i X_{i,0}$$

$$\leq \frac{M|\gamma_0 - \hat{\gamma}_0|}{n}\sum_{j \in S}|\hat{\theta}_j X_{i,j}|. \tag{24}$$

Let $\Delta := \max_{j \in S'}|\hat{\gamma}_j - \gamma_j|$. For each $j \in S$, let $\theta_j := \gamma_j/\beta_j$. Recall the quantity $M_2$ from the statement of the theorem. Note that if

$$M < \frac{1}{2}M_2, \tag{25}$$

we get that for each $j \in S$,

$$|\hat{\theta}_j - \theta_j| = \left|\frac{\hat{\gamma}_j}{\hat{\beta}_j} - \frac{\gamma_j}{\beta_j}\right|$$

$$\leq \frac{|\hat{\gamma}_j - \gamma_j|}{|\hat{\beta}_j|} + \gamma_j\left|\frac{1}{\hat{\beta}_j} - \frac{1}{\beta_j}\right|$$

$$\leq K_6(\Delta + M).$$

Combining this observation with the inequalities (23) and (24), and summing over $k \in S'$, we get that

$$\sum_{j,k \in S'}\hat{\sigma}_{j,k}(\gamma_j - \hat{\gamma}_j)(\gamma_k - \hat{\gamma}_k) \leq K_7\lambda\Delta + K_8\Delta Q_1 + K_9(M + M^2)\Delta(\Delta + M)Q_2, \tag{26}$$

where

$$\hat{\sigma}_{j,k} := \frac{1}{n}\sum_{i=1}^{n}X_{i,j}X_{i,k}, \quad Q_1 := \max_{j \in S'}\left|\frac{1}{n}\sum_{i=1}^{n}\epsilon_i X_{i,j}\right|, \quad Q_2 := \max_{j \in S'}\frac{1}{n}\sum_{i=1}^{n}X_{i,j}^2.$$

Let $\hat{\eta}$ be the smallest eigenvalue of the positive semidefinite matrix $(\hat{\sigma}_{j,k})_{j,k \in S'}$, so that

$$\sum_{j,k \in S'}\hat{\sigma}_{j,k}(\gamma_j - \hat{\gamma}_j)(\gamma_k - \hat{\gamma}_k) \geq \hat{\eta}\sum_{j \in S'}(\hat{\gamma}_j - \gamma_j)^2 \geq \hat{\eta}\Delta^2.$$

Combining this with equation (26) and rearranging, we get

$$(\hat{\eta} - K_9(M + M^2)Q_2)\Delta^2 \leq (K_7\lambda + K_8Q_1 + K_9(M + M^2)MQ_2)\Delta$$

If the coefficient of $\Delta^2$ on the left is positive, this gives

$$\Delta \leq \frac{K_7\lambda + K_8Q_1 + K_9(M + M^2)MQ_2}{\hat{\eta} - K_9(M + M^2)Q_2}.$$

Now, using Bernstein's inequality, it is easy to show that $Q_1 = o(1)$ and $Q_2 = O(1)$ (following the conventions introduced at the beginning of the proof). By Corollary 2, we know that $M = o(1)$. Again by Bernstein's inequality, $\hat{\sigma}_{j,k} = \mathrm{E}(X_jX_k) + o(1)$ for each $j, k \in S$. From this, it is easy to deduce via standard matrix inequalities that $\hat{\eta} = \eta_0 + o(1)$, where $\eta_0$ is the minimum eigenvalue of $(\mathrm{E}(X_jX_k))_{j,k\in S'}$. Recalling that $X_0 = 1$, we see that this is equal to the minimum eigenvalue $\eta$ of the covariance matrix of $(X_j)_{j\in S}$. Combining all of these, it is now straightforward that under the condition (17),

$$\mathbb{P}(E \cap \{\Delta > K_{10}\lambda\}) \leq K_{11}ne^{-K_{12}n\lambda^2}.$$

Together with equation (18), this completes the proof. $\qquad\square$

*Proof of Theorem 2.* Note that since $Y = \gamma_0 + \sum_{j\in S}\gamma_jX_j + \epsilon$, and $\epsilon$ is independent of the $X_j$'s, we have

$$\begin{aligned}
\beta_j &= \frac{\mathrm{Cov}(Y, X_j)}{\mathrm{Var}(X_j)} \\
&= \frac{\sum_{k\in S}\gamma_k\mathrm{Cov}(X_k, X_j)}{\mathrm{Var}(X_j)} \\
&= \gamma_j + \sum_{k\in S\setminus\{j\}}\gamma_k\frac{\mathrm{Cov}(X_k, X_j)}{\mathrm{Var}(X_j)} = \gamma_j + \sum_{k\in S\setminus\{j\}}\gamma_k\delta_{k,j}.
\end{aligned}$$

Since $\gamma_j \neq 0$, we may divide throughout by $\gamma_j$ and get

$$\frac{\beta_j}{\gamma_j} = 1 + \sum_{k\in S\setminus\{j\}}\frac{\gamma_k}{\gamma_j}\delta_{k,j}.$$

Now, by assumption, $\delta_{k,j} \geq 0$ whenever $\gamma_k/\gamma_j > 0$, and $\delta_{k,j} \leq 0$ whenever $\gamma_k/\gamma_j < 0$. Thus,

$$\sum_{k\in S\setminus\{j\}}\frac{\gamma_k}{\gamma_j}\delta_{k,j} \geq 0.$$

This shows that $\beta_j/\gamma_j \geq 1$. In particular, $\beta_j$ is nonzero and has the same sign as $\gamma_j$. $\qquad\square$

*Proof of Theorem 3.* Note that

$$\begin{aligned}
\beta_j &= \frac{\mathrm{Cov}(Y, X_j)}{\mathrm{Var}(X_j)} \\
&= \sum_{k\in S}\gamma_k\frac{\mathrm{Cov}(X_k, X_j)}{\mathrm{Var}(X_j)} \\
&= \sum_{k\in S}\gamma_k\delta_{kj}.
\end{aligned}$$

Thus,

$$\begin{aligned}
\frac{\beta_j}{\gamma_j} &= \sum_{k\in S}\frac{\gamma_k}{\gamma_j}\delta_{kj} \\
&= \sum_{k\in A_j}\left|\frac{\gamma_k}{\gamma_j}\delta_{kj}\right| - \sum_{k\notin A_j}\left|\frac{\gamma_k}{\gamma_j}\delta_{kj}\right| \\
&= \frac{\sum_{k\in A_j}|\gamma_k\delta_{kj}| - \sum_{k\notin A_j}|\gamma_k\delta_{kj}|}{|\gamma_j|}.
\end{aligned}$$

Since $\beta_j\gamma_j \geq 0$ if and only if $\beta_j/\gamma_j \geq 0$, this proves the claim. $\qquad\square$

# C Proof of Theorem 4

Let $\gamma := (\beta_0, \beta) \in \mathbb{R}^{p+1}$ and define $\hat{\gamma}^{\text{OLS}}$ and $\hat{\gamma}^{\text{UR}}$ similarly. Let $X$ denote the $n \times (p+1)$ matrix whose columns are indexed from 0 to $p$; $X_{i,0} = 1$ for each $i$, and for $1 \leq j \leq p$, $X_{i,j}$ is as before. Let $Y \in \mathbb{R}^n$ be the vector whose $i^{\text{th}}$ component is $Y_i$, and define $\epsilon$ similarly. Note that $Y = X\gamma + \epsilon$, and that $\hat{\mu}^{\text{OLS}}$ is the Euclidean projection on $Y$ on the column space $\mathcal{C}$ of $X$. Let

$$\mathcal{D} := \{X\theta : \theta = (\theta_0, \ldots, \theta_p) \in \mathbb{R}^p, \theta_j \hat{\beta}_j^{\text{uni}} \geq 0 \text{ for all } 1 \leq j \leq p\},$$

so that $\mathcal{D}$ is a random closed convex subset of $\mathcal{C}$, and $\hat{\mu}^{\text{UR}}$ is the Euclidean projection of $Y$ on $\mathcal{D}$. Let $E$ be the event that $\hat{\beta}_j^{\text{uni}} \beta_j^{\text{uni}} \geq 0$ for each $1 \leq j \leq p$. Note that if $E$ happens, then $\mu \in \mathcal{D}$.

**Lemma 2.** *There is a positive constant $C_0$ depending only on $\lambda_0$, $\delta$ and $\sigma$ such that*

$$\mathbb{P}(E) \geq 1 - pe^{-C_0 n}.$$

*Proof.* If $\beta_j^{\text{uni}} = 0$, then $\hat{\beta}_j^{\text{uni}} \beta_j^{\text{uni}} \geq 0$ anyway. Thus,

$$\mathbb{P}(E^c) = \mathbb{P}(\hat{\beta}_j^{\text{uni}} \beta_j^{\text{uni}} < 0 \text{ for some } j \text{ such that } \beta_j^{\text{uni}} \neq 0)$$
$$\leq \sum_{j : \beta_j^{\text{uni}} \neq 0} \mathbb{P}(\hat{\beta}_j^{\text{uni}} \beta_j^{\text{uni}} < 0).$$

Take any $j$ such that $\beta_j^{\text{uni}} \neq 0$. Then by standard results, conditional on $X$,

$$\hat{\beta}_j^{\text{uni}} \sim N\left(\beta_j^{\text{uni}}, \frac{\sigma^2}{\sum_{i=1}^n (X_{i,j} - \overline{X}_j)^2}\right),$$

where $\overline{X}_j := \frac{1}{n} \sum_{i=1}^n X_{i,j}$. Recall that $|\beta_j^{\text{uni}}| \geq \delta$ and $\text{Var}(X_{i,j}) \geq \lambda_0$. From this, it is easy to deduce that there is a positive constant $C_0$ depending only on $\lambda_0$, $\delta$ and $\sigma$ such that $\mathbb{P}(E^c) \leq pe^{-C_0 n}$. This completes the proof. $\square$

**Lemma 3.** *If the event $E$ happens, then*

$$\|\mu - \hat{\mu}^{\text{UR}}\|^2 \leq \|\mu - \hat{\mu}^{\text{OLS}}\|^2 - \|\hat{\mu}^{\text{OLS}} - \hat{\mu}^{\text{UR}}\|^2.$$

*Proof.* Since $\hat{\mu}^{\text{OLS}}$ is the projection of $Y$ on $\mathcal{C}$ and $\mathcal{C}$ is a subspace, it follows that $Y - \hat{\mu}^{\text{OLS}}$ is orthogonal to $x - \hat{\mu}^{\text{OLS}}$ for all $x \in \mathcal{C}$. Since $\mathcal{D} \subseteq \mathcal{C}$, this implies that for any $\nu \in \mathcal{D}$,

$$\|Y - \nu\|^2 = \|Y - \hat{\mu}^{\text{OLS}}\|^2 + \|\hat{\mu}^{\text{OLS}} - \nu\|^2.$$

This shows that the projection of $Y$ on $\mathcal{D}$ (that is, $\hat{\mu}^{\text{UR}}$) is the same as the projection of $\hat{\mu}^{\text{OLS}}$ on $\mathcal{D}$. Since $\mathcal{D}$ is convex and $\mu \in \mathcal{D}$ (because $E$ has happened), $t\mu + (1-t)\hat{\mu}^{\text{UR}} \in \mathcal{D}$ for all $t \in [0, 1]$. But then, by the preceding sentence,

$$\|\hat{\mu}^{\text{OLS}} - (t\mu + (1-t)\hat{\mu}^{\text{UR}})\|^2 \geq \|\hat{\mu}^{\text{OLS}} - \hat{\mu}^{\text{UR}}\|^2$$

for all $t \in [0, 1]$. In particular, the derivative of the left side with respect to $t$ should be nonnegative at $t = 0$. This shows that

$$(\hat{\mu}^{\text{OLS}} - \hat{\mu}^{\text{UR}})^T (\hat{\mu}^{\text{UR}} - \mu) \geq 0.$$

From this, we get

$$\|\hat{\mu}^{\text{OLS}} - \mu\|^2 - \|\hat{\mu}^{\text{UR}} - \mu\|^2 = \|\hat{\mu}^{\text{OLS}} - \hat{\mu}^{\text{UR}}\|^2 + 2(\hat{\mu}^{\text{OLS}} - \hat{\mu}^{\text{UR}})^T (\hat{\mu}^{\text{UR}} - \mu)$$
$$\geq \|\hat{\mu}^{\text{OLS}} - \hat{\mu}^{\text{UR}}\|^2.$$

This completes the proof. $\square$

**Lemma 4.** *Let $Z$ be an $n \times p$ matrix of i.i.d. $N(0,1)$ random variables. Let $1 \in \mathbb{R}^n$ be the vector of all 1's, and let $V := Z - \frac{1}{n} 11^T Z$. Let $\mu$ be the minimum eigenvalue of $\frac{1}{n} V^T V$ and $\nu$ be the maximum eigenvalue of $\frac{1}{n} Z^T Z$. There are positive constants $C_1, C_2, \mu_0, \nu_0$ depending only on $r_0$ such that*

$$\mathbb{P}(\mu \geq \mu_0, \nu \leq \nu_0) \geq 1 - C_1 e^{-C_2 n}.$$

*Proof.* By basic facts from statistics, we know that $Z^T Z$ follows the Wishart distribution $W_p(n, I_p)$, and $V^T V \sim W_p(n-1, I_p)$. The claimed bounds now follow from known results in the literature, such as (Rudelson & Vershynin 2010, Proposition 2.4 and Theorem 3.3). $\qquad\square$

**Lemma 5.** *Let $\hat{\lambda}_0$ and $\hat{\lambda}_1$ be the minimum and maximum eigenvalues of $\frac{1}{n} X^T X$. There are positive constants $C_3, C_4, \lambda'_0, \lambda'_1$ depending only on $\lambda_0, \lambda_1,$ and $r_0$ such that*

$$\mathbb{P}(\lambda'_0 \le \hat{\lambda}_0 \le \hat{\lambda}_1 \le \lambda'_1) \ge 1 - C_3 e^{-C_4 n}.$$

*Proof.* Let us write $X = [1 \ \tilde{X}]$, where $\tilde{X}$ consists of columns $1, \ldots, p$ of $X$ and $1$ denotes the first column, which has all 1's. Then for any $y = (y_0, \tilde{y}) \in \mathbb{R}^{p+1}$, where $y_0$ denotes the first coordinate of $y$, we have

$$\|Xy\|^2 = \|y_0 1 + \tilde{X}\tilde{y}\|^2. \tag{27}$$

Let $\Sigma^{1/2}$ denote the positive definite square-root of $\Sigma$. Then we can write $\tilde{X} = Z\Sigma^{1/2}$, where $Z$ is an $n \times p$ matrix with i.i.d. $N(0,1)$ random variables. Let $V := Z - \frac{1}{n} 11^T Z$. Let $\mu$ be the smallest eigenvalue of $\frac{1}{n} V^T V$ and $\nu$ be the largest eigenvalue of $\frac{1}{n} Z^T Z$. By the above identity and the inequality $\|u + v\|^2 \le 2\|u\|^2 + 2\|v\|^2$, we get

$$\begin{aligned}
\|Xy\|^2 &= \|y_0 1 + Z\Sigma^{1/2}\tilde{y}\|^2 \\
&\le 2\|y_0 1\|^2 + 2\|Z\Sigma^{1/2}\tilde{y}\|^2 \\
&\le 2y_0^2 n + 2n\nu\|\Sigma^{1/2}\tilde{y}\|^2 \\
&\le 2y_0^2 n + 2n\nu\lambda_1\|\tilde{y}\|^2 \le 2n(1 + \lambda_1\nu)\|y\|^2.
\end{aligned}$$

This shows that

$$\hat{\lambda}_1 \le 2(1 + \lambda_1\nu). \tag{28}$$

Next, let $z$ be the projection of $Z\Sigma^{1/2}\tilde{y}$ on the span of 1, given by

$$z = n^{-1} 11^T Z\Sigma^{1/2}\tilde{y}.$$

Then by equation (27),

$$\begin{aligned}
\|Xy\|^2 &= \|y_0 1 - z\|^2 + \|\tilde{X}\tilde{y} - z\|^2 \\
&\ge \|\tilde{X}\tilde{y} - z\|^2 = \|Z\Sigma^{1/2}\tilde{y} - n^{-1} 11^T Z\Sigma^{1/2}\tilde{y}\|^2 \\
&= \|V\Sigma^{1/2}\tilde{y}\|^2 \ge n\mu\|\Sigma^{1/2}\tilde{y}\|^2 \ge n\mu\lambda_0\|\tilde{y}\|^2. \tag{29}
\end{aligned}$$

But by the first line of the above display, we also have

$$\|Xy\|^2 \ge \|y_0 1 - z\|^2 \ge (\|y_0 1\| - \|z\|)^2 = (|y_0|\sqrt{n} - \|z\|)^2. \tag{30}$$

Now, note that

$$\|z\| = n^{-1} |1^T Z\Sigma^{1/2}\tilde{y}| \|1\| \le n^{-1/2} \|1\| \|Z\Sigma^{1/2}\tilde{y}\| \le \sqrt{n\nu}\lambda_1\|\tilde{y}\|. \tag{31}$$

So, if $\nu\lambda_1\|\tilde{y}\| \le \frac{1}{2}|y_0|$, then by (30) and (31),

$$\|Xy\|^2 \ge \frac{1}{4} y_0^2 n = \left(\frac{1}{8} y_0^2 + \frac{1}{8} y_0^2\right) n \ge \left(\frac{1}{8} y_0^2 + \frac{4}{8}\nu^2\lambda_1^2\|\tilde{y}\|^2\right) n.$$

Without loss, let us assume that $\nu$ and $\lambda_1$ were chosen so large that $4\nu^2\lambda_1^2 \ge 1$. Then the above inequality gives

$$\|Xy\|^2 \ge \left(\frac{1}{8}(\|y\|^2 - \|\tilde{y}\|^2) + \frac{4}{8}\nu^2\lambda_1^2\|\tilde{y}\|^2\right) n \ge \frac{1}{8}\|y\|^2 n. \tag{32}$$

On the other hand, if $\nu\lambda_1\|\tilde{y}\| > \frac{1}{2}|y_0|$, then by equation (29),

$$\begin{aligned}
\|Xy\|^2 &\ge n\mu\lambda_0\left(\frac{1}{2}\|\tilde{y}\|^2 + \frac{1}{2}\|\tilde{y}\|^2\right) \\
&\ge n\mu\lambda_0\left(\frac{1}{2}\|\tilde{y}\|^2 + \frac{1}{8\nu^2\lambda_1^2} y_0^2\right) \\
&= n\mu\lambda_0\left(\frac{1}{2}\|\tilde{y}\|^2 + \frac{1}{8\nu^2\lambda_1^2}(\|y\|^2 - \|\tilde{y}\|^2)\right) \ge \frac{n\mu\lambda_0}{8\nu^2\lambda_1^2}\|y\|^2, \tag{33}
\end{aligned}$$

where the last inequality follows from the assumption that $4\nu^2\lambda_1^2 \geq 1$. Combining the inequalities (32) and (33), we see that

$$\hat{\lambda}_0 \geq \min\left\{\frac{1}{8}, \frac{\mu\lambda_0}{8\nu^2\lambda_1^2}\right\}. \tag{34}$$

Inequalities (28) and (34) and Lemma 4 complete the proof. □

We are now ready to prove Theorem 4. First, note that

$$\|\hat{\mu}^{\mathrm{OLS}} - \hat{\mu}^{\mathrm{UR}}\|^2 = \|X\hat{\gamma}^{\mathrm{OLS}} - X\hat{\gamma}^{\mathrm{UR}}\|^2 \geq n\hat{\lambda}_0\|\hat{\gamma}^{\mathrm{OLS}} - \hat{\gamma}^{\mathrm{UR}}\|^2$$
$$= n\hat{\lambda}_0(\hat{\beta}_0^{\mathrm{OLS}} - \hat{\beta}_0^{\mathrm{UR}})^2 + n\hat{\lambda}_0\sum_{j=1}^p(\hat{\beta}_j^{\mathrm{OLS}} - \hat{\beta}_j^{\mathrm{UR}})^2.$$

By the definition of $\hat{\beta}^{\mathrm{UR}}$, $\hat{\beta}_j^{\mathrm{UR}}\hat{\beta}_j^{\mathrm{uni}} \geq 0$ for each $j$; that is, $\hat{\beta}_j^{\mathrm{UR}}$ and $\hat{\beta}_j^{\mathrm{uni}}$ are on the same side of zero on the real line (including the possibility that one or both may be equal to zero). Thus, if $E$ happens and $\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}} < 0$, then $\hat{\beta}_j^{\mathrm{OLS}}$ and $\hat{\beta}_j^{\mathrm{UR}}$ are on opposite sides of zero (again, including the possibility that one or both may be equal to zero), and hence,

$$(\hat{\beta}_j^{\mathrm{OLS}} - \hat{\beta}_j^{\mathrm{UR}})^2 \geq (\hat{\beta}_j^{\mathrm{OLS}})^2.$$

(Note that this does not hold if we weaken the condition to $\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}} \leq 0$. This is because this weakening leaves open the possibility that $\beta_j^{\mathrm{uni}} = 0$ and $\hat{\beta}_j^{\mathrm{OLS}}, \hat{\beta}_j^{\mathrm{UR}}$ are on the same side of zero.) Combining, we see that if $E$ happens, then

$$\|\hat{\mu}^{\mathrm{OLS}} - \hat{\mu}^{\mathrm{UR}}\|^2 \geq n\hat{\lambda}_0 \sum_{j\,:\,\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}}<0}(\hat{\beta}_j^{\mathrm{OLS}})^2. \tag{35}$$

Take any $j$ such that $\beta_j = 0$ and $\beta_j^{\mathrm{uni}} \neq 0$. Conditional on $X$,

$$\hat{\beta}_j^{\mathrm{OLS}} \sim N(0, \sigma^2\theta_j), \tag{36}$$

where $\theta_0, \ldots, \theta_p$ are the diagonal elements of $(X^TX)^{-1}$. Since $\theta_j \geq (n\hat{\lambda}_1)^{-1}$, this shows that

$$\mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^2 1_{\{\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}}<0\}}|\mathrm{X}) = \frac{1}{2}\mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^2|\mathrm{X}) \geq \frac{\sigma^2}{2n\hat{\lambda}_1}. \tag{37}$$

Let $F$ be the event that $\lambda_0' \leq \hat{\lambda}_0 \leq \hat{\lambda}_1 \leq \lambda_1'$, where $\lambda_0', \lambda_1'$ are the constants from Lemma 5, and let $H := E \cap F$. Then by the inequality (35),

$$\mathrm{E}\|\hat{\mu}^{\mathrm{OLS}} - \hat{\mu}^{\mathrm{UR}}\|^2 \geq \mathrm{E}(\|\hat{\mu}^{\mathrm{OLS}} - \hat{\mu}^{\mathrm{UR}}\|^2 1_{\mathrm{H}})$$
$$\geq \mathrm{E}\left(n\hat{\lambda}_0 \sum_{j\,:\,\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}}<0}(\hat{\beta}_j^{\mathrm{OLS}})^2 1_{\mathrm{H}}\right)$$
$$\geq n\lambda_0' \sum_{j\,:\,\beta_j=0,\,\beta_j^{\mathrm{uni}}\neq 0} \mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^2 1_{\{\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}}<0\}} 1_{\mathrm{H}}). \tag{38}$$

Take any $j$ such that $\beta_j = 0$ and $\beta_j^{\mathrm{uni}} \neq 0$. Note that

$$\mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^2 1_{\{\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}}<0\}} 1_{\mathrm{H}}) = \mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^2 1_{\{\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}}<0\}} 1_{\mathrm{F}})$$
$$- \mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^2 1_{\{\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}}<0\}} 1_{\mathrm{F}\cap\mathrm{E}^c}). \tag{39}$$

By Lemma 5, the inequality (37), and the fact that the event $F$ depends only on $X$, we get

$$\mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^2 1_{\{\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}}<0\}} 1_{\mathrm{F}}) = \mathrm{E}(\mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^2 1_{\{\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}}<0\}}|\mathrm{X}) 1_{\mathrm{F}})$$
$$\geq \frac{\sigma^2}{2n\lambda_1'}\mathbb{P}(F) \geq \frac{\sigma^2}{2n\lambda_1'}(1 - C_3e^{-C_4n}). \tag{40}$$

33

On the other hand, by the Cauchy–Schwarz inequality, equation (36), and Lemma 2,

$$\mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^2 1_{\{\hat{\beta}_j^{\mathrm{OLS}}\beta_j^{\mathrm{uni}}<0\}} 1_{\mathrm{F}\cap\mathrm{E^c}}) \le [\mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^4 1_{\mathrm{F}})\mathbb{P}(\mathrm{E^c})]^{1/2}$$

$$\le [\mathrm{E}(\mathrm{E}((\hat{\beta}_j^{\mathrm{OLS}})^4|\mathrm{X})1_{\mathrm{F}})]^{1/2}\sqrt{p}e^{-\frac{1}{2}C_0 n}$$

$$\le \left(\frac{3\sigma^4}{n\lambda_0'}\right)^{1/2}\sqrt{p}e^{-\frac{1}{2}C_0 n}$$

$$= \frac{\sqrt{3r}\sigma^2}{\sqrt{\lambda_0'}}e^{-\frac{1}{2}C_0 n}. \tag{41}$$

Combining equations (38), (39), (40) and (41), we get

$$\mathrm{E}\|\hat{\mu}^{\mathrm{OLS}} - \hat{\mu}^{\mathrm{UR}}\|^2 \ge n\lambda_0' q\left(\frac{\sigma^2}{2n\lambda_1'}(1 - C_3 e^{-C_4 n}) - \frac{\sqrt{3r}\sigma^2}{\sqrt{\lambda_0'}}e^{-\frac{1}{2}C_0 n}\right)$$

$$= \frac{\sigma^2 q\lambda_0'}{2\lambda_1'} - C_3\frac{\sigma^2 q\lambda_0'}{2\lambda_1'}e^{-C_4 n} - \sqrt{3r\lambda_0'}q\sigma^2 n e^{-\frac{1}{2}C_0 n}.$$

By Lemma 3, this completes the proof.

# D Derivation supporting Section 4

We wish to establish when the correlations of the univariate LOO fitted values with the response $y$ are likely to be positive. Thanks to Chris Habron for outlining this analysis.

Assume we have standardized both the feature vectors $x_j$ and response $y$ to have zero mean and unit variance. Since the analysis focuses on individual features, we will drop the index $j$.

Formula (10) relates the LOO residuals to the OLS residuals:

$$y_i - \hat{\eta}^{-i} = \frac{y_i - \hat{\eta}^i}{1 - H_{ii}}.$$

Since both $x$ and $y$ are standardized, the vector of OLS fits is given by

$$\hat{\eta} = x\hat{\beta}$$
$$= x(x^\top y)/n. \tag{42}$$

With $D = \mathrm{diag}(\frac{1}{1-H_{11}}, \frac{1}{1-H_{22}}, \ldots, \frac{1}{1-H_{nn}})$ we can write the vector of LOO fits as

$$\hat{\eta}^{loo} = (I - D)y + Dx(x^\top y)/n. \tag{43}$$

We look at the sample covariance between $y$ and $\hat{\eta}^{loo}$

$$\mathrm{Cov}(y, \hat{\eta}^{loo}) = 1 - (y^\top Dy)/n + (y^\top Dx)(x^\top y)/n^2. \tag{44}$$

We would like to know when this covariance is positive. (Note that it is always non-negative if we replace $\hat{\eta}^{loo}$ with $\hat{\eta}$). Since $x$ is standardized,

$$H_{ii} = \frac{1}{n} + \frac{x_i^2}{n}$$

and hence

$$\frac{1}{1 - H_{ii}} = \frac{n}{n - 1 - x_i^2}$$

To identify when this covariance becomes positive, we make the approximation that $H_{ii}$ is constant, which implies $x_i^2 = 1$. This isn't exact, but for small correlations between $x$ and $y$, this introduces minimal error. With this approximation

$$\mathrm{Cov}(y, \hat{\eta}^{loo}) \approx 1 - \frac{n}{n-2} + \mathrm{Cov}(y, x)^2 \cdot \frac{n}{n-2}$$

This becomes positive when

$$\mathrm{Cov}(y, x)^2 > \frac{2}{n}$$

The critical value for a significance test for a correlation using the $t$-distribution is:

$$r = \frac{t}{\sqrt{n - 2 + t^2}}$$

We achieve this when $t = \sqrt{2}$, which corresponds to a two-sided $p$-value of 0.16 for large $n$.

# References

Breiman, L. (1995), 'Better subset selection using the non-negative garotte', *Technometrics* **37**, 738–754.

Breiman, L. (1996), 'Stacked regressions', *Machine Learning* **24**, 51–64.

Candes, E. J., Wakin, M. B. & Boyd, S. P. (2008), 'Enhancing sparsity by reweighted l1 minimization', *Journal of Fourier Analysis and Applications* **14**(5), 877–905.

Hastie, T., Tibshirani, R. & Wainwright, M. (2015), *Statistical Learning with Sparsity: The Lasso and Generalizations*, Chapman & Hall/CRC.

Huang, J., Ma, S. & Zhang, C. (2008), 'Adaptive lasso for sparse high-dimensional regression models', *Statistica Sinica* **18**(4), 1603–1618.

Khan, J., Wei, J., Ringner, M., Saal, L., Ladanyi, M., Westermann, F., Berthold, F., Schwab, M., Antonescu, C., Peterson, C., & Meltzer, P. (2001), 'Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks', *Nature Medicine* **7**, 673–679.

Mazumder, R., Friedman, J. & Hastie, T. (2011), 'Sparsenet: Coordinate descent with non-convex penalties', *J.Amer. Statist. Assoc.* **106**.

Meinshausen, N. (2012), 'Sign-constrained least squares estimation for high-dimensional regression', *Electronic Journal of Statistics* **7**.

Rad, K. R. & Maleki, A. (2020), 'A scalable estimate of the out-of-sample prediction error via approximate leave-one-out cross-validation', *Journal of the Royal Statistical Society Series B: Statistical Methodology* **82**(4), 965–996.
**URL:** *https://doi.org/10.1111/rssb.12374*

Rudelson, M. & Vershynin, R. (2010), Non-asymptotic theory of random matrices: extreme singular values, *in* 'Proceedings of the International Congress of Mathematicians 2010 (ICM 2010) (In 4 Volumes) Vol. I: Plenary Lectures and Ceremonies Vols. II–IV: Invited Lectures', World Scientific, pp. 1576–1602.

Talagrand, M. (2010), *Mean field models for spin glasses: Volume I: Basic examples*, Springer Science & Business Media.

Wolpert, D. (1992), 'Stacked generalization', *Neural Networks* **5**, 241–259.

Zhang, C.-H. (2010), 'Nearly unbiased variable selection under minimax concave penalty', *The Annals of Statistics* **38**(2), 894 – 942.
**URL:** *https://doi.org/10.1214/09-AOS729*

Zou, H. (2006), 'The adaptive lasso and its oracle properties', *Journal of the American Statistical Association* **101**, 1418–1429.

Zou, H. & Hastie, T. (2005), 'Regularization and variable selection via the elastic net', *Journal of the Royal Statistical Society Series B.* **67**(2), 301–320.