

# Principal component-guided sparse regression

J. Kenneth Tay, Jerome Friedman and Robert Tibshirani  
 Department of Statistics, and Department of Biomedical Data Science  
 Stanford University

October 25, 2018

## Abstract

We propose a new method for supervised learning, especially suited to wide data where the number of features is much greater than the number of observations. The method combines the lasso ( $\ell_1$ ) sparsity penalty with a quadratic penalty that shrinks the coefficient vector toward the leading principal components of the feature matrix. We call the proposed method the “*principal components lasso*” (“pcLasso”). The method can be especially powerful if the features are pre-assigned to groups (such as cell-pathways, assays or protein interaction networks). In that case, pcLasso shrinks each group-wise component of the solution toward the leading principal components of that group. In the process, it also carries out selection of the feature groups. We provide some theory for this method and illustrate it on a number of simulated and real data examples.

## 1 Introduction

We consider the usual linear regression model: given  $n$  realizations of  $p$  predictors  $\mathbf{X} = \{x_{ij}\}$  for  $i = 1, 2, \dots, n$  and  $j = 1, 2, \dots, p$ , the response  $\mathbf{y} = (y_1, \dots, y_n)$  is modeled as

$$y_i = \beta_0 \mathbf{1} + \sum_j x_{ij} \beta_j + \epsilon_i, \quad (1)$$

with  $\epsilon \sim (0, \sigma^2)$ . The ordinary least squares (OLS) estimates of  $\beta_j$  are obtained by minimizing the residual sum of squares (RSS). There has been much work on regularized estimators that offer an advantage over the OLS estimates, both in terms of accuracy of prediction on future data and interpretation of the fitted model. One focus has been on the *lasso* (Tibshirani 1996), which minimizes

$$J(\beta_0, \beta) = \frac{1}{2} \|\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1, \quad (2)$$

where  $\beta = (\beta_1, \dots, \beta_p)^T$ , and the tuning parameter  $\lambda \geq 0$  controls the sparsity of the final model. This parameter is often selected by cross-validation (CV). The objective function  $J(\beta_0, \beta)$  is convex, which means that solutions can be found efficiently even for very large  $n$  and  $p$ , in contrast to combinatorial methods like best subset selection.

The lasso is an alternative to ridge regression (Hoerl & Kennard 1970), which uses an objective function with an  $\ell_2$  penalty:

$$\frac{1}{2} \|\mathbf{y} - \beta_0 \mathbf{1} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_2^2. \quad (3)$$

This approach has the disadvantage of producing dense models. However, it also biases the solution vector toward the leading right singular vectors of  $\mathbf{X}$ , which can be effective for improving prediction accuracy.

The elastic net (Zou & Hastie 2005) generalizes the lasso, effectively combining ridge regression and the lasso, by adding an  $\ell_2$  penalty to the lasso’s objective function:

$$\frac{1}{2}\|\mathbf{y} - \beta_0\mathbf{1} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_1 + \frac{\lambda_2}{2}\|\beta\|_2^2. \quad (4)$$

We have found that the bias of ridge regression (and hence the elastic net) toward the leading singular vectors is somewhat mild. Furthermore, if the predictors come in pre-assigned groups— a focus of this paper— these methods do not exploit this information.

In this paper, we propose a new method, the *principal components lasso* (“pcLasso”), that strongly biases the solution vector toward the leading singular vectors and exploits group structure. We give the full definition of the pcLasso for multiple groups in Section 3. For motivation, we describe the simpler single group case in the next section.

## 2 Quadratic penalties based on principal components

Assume that we have centered the columns of  $\mathbf{X}$ , and let  $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$  be the singular value decomposition (SVD) of  $\mathbf{X}$ . Here,  $\mathbf{D}$  is a diagonal matrix with diagonal entries equaling the singular values  $d_1 \geq d_2 \dots \geq d_m > 0$ , with  $m = \text{rank}(\mathbf{X})$ . Assume also that  $\mathbf{y}$  has been centered, so that we may omit the intercept from the model.

In generalizing the  $\ell_2$  penalty, one can imagine a class of penalty functions of the form

$$\beta^T \mathbf{V}\mathbf{Z}\mathbf{V}^T \beta,$$

where  $\mathbf{Z}$  is a diagonal matrix whose diagonal elements are functions of the squared singular values:

$$Z_{11} = f_1(d_1^2, d_2^2, \dots, d_m^2), Z_{22} = f_2(d_1^2, d_2^2, \dots, d_m^2), \dots$$

The ridge penalty chooses  $Z_{jj} = 1$  for all  $j$ , leading to the quadratic penalty  $\beta^T \beta$ . Here we propose a different choice: we minimize

$$J(\beta) = \frac{1}{2}\|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda\|\beta\|_1 + \frac{\theta}{2}\beta^T \mathbf{V}\mathbf{D}_{d_1^2 - d_j^2} \mathbf{V}^T \beta, \quad (5)$$

where  $\mathbf{D}_{d_1^2 - d_j^2}$  is an  $m \times m$  diagonal matrix with diagonal entries equaling  $d_1^2 - d_1^2, d_1^2 - d_2^2, \dots, d_1^2 - d_m^2$ . We call this the “pcLasso penalty” and use it more generally in the next section. The value  $\theta \geq 0$  is a tuning parameter determining the weight given to the second penalty. We see that this penalty term gives no weight (a “free ride”) to the component of  $\beta$  that aligns with the first right singular vector of  $\mathbf{X}$ , i.e, the first principal component (PC). Beyond that, the size of the penalty depends on the gap between the squared singular values  $d_j^2$  and  $d_1^2$ . If we view the quadratic penalty as representing a Bayesian (log-)prior, it puts infinite prior variance on the leading eigenvector.

It is instructive to compare (5) with  $\lambda = 0$  (and with  $\frac{\theta}{2}$  replaced with  $\theta$ ) to ridge regression which uses the penalty  $\theta\beta^T \beta$ . By transforming to principal coordinates, it follows that ridge regression produces the fit

$$\mathbf{X}\hat{\beta}_R = \sum_{j=1}^m \frac{d_j^2}{d_j^2 + \theta} u_j u_j^T \mathbf{y}. \quad (6)$$

In contrast, pcLasso (5) gives

$$\mathbf{X}\hat{\beta}_L = \sum_{j=1}^m \frac{d_j^2}{d_j^2 + \theta(d_1^2 - d_j^2)} u_j u_j^T \mathbf{y}. \quad (7)$$

The latter expression corresponds to a more aggressive form of shrinkage toward the leading singular vectors. To see this, Figure 1 shows the shrinkage factors  $\frac{d_j^2}{d_j^2 + \theta}$  and  $\frac{d_j^2}{d_j^2 + \theta(d_1^2 - d_j^2)}$  for a

100 × 20 design matrix  $\mathbf{X}$  with a strong rank-1 component. In each panel, we have chosen the value  $\theta$  (different for each method) to yield the degrees of freedom indicated at the top of the panel. (For each method, the degrees of freedom is equal to the sum of the shrinkage factors.)

We see that when the degrees of freedom is about “right” (top-left panel), ridge regression shrinks the top component about twice as much as does pcLasso (top-right panel). This contrast is less stark when the degrees of freedom is increased to 5 (bottom-left panel). This aggressiveness can help pcLasso improve prediction accuracy when the signal lines up strongly with the top PCs. More importantly, with pre-specified feature groups, the pcLasso penalty in (5) can be modified in such a way that it shrinks the subvector of  $\beta$  in each group toward the leading PCs of that group. This exploits the group structure to give potentially better predictions and simultaneously selects groups of features. We give details in the next section.

The rest of this paper is organized follows. In Section 3 we give a full definition of our proposal. We illustrate the effectiveness of pcLasso on some real data examples in Section 4, and give details of our computational approach in Section 5. Section 6 derives a formula for the degrees of freedom of the pcLasso fit. In Section 7 we discuss an extensive simulation study, comparing our proposal to the lasso, elastic net and principal components regression. We study the theoretical properties of pcLasso in Section 8, showing that, under certain assumptions, it yields lower  $\ell_2$  and prediction error than the lasso. We end with a discussion and ideas for future work. The Appendix contains further details and proofs, as well as a full description of the set-up and results for our simulation study in Section 7.

### 3 Principal components lasso (“pcLasso”)

#### 3.1 Definition

Suppose that the  $p$  predictors are grouped in  $K$  non-overlapping groups (we discuss the overlapping groups case later in Section 3.2). For example, these groups could be different gene pathways, assays or protein interaction networks. For  $k = 1, \dots, K$ , let  $\mathbf{X}_k$  denote the  $p_k$  columns of  $\mathbf{X}$  corresponding to group  $k$ , let  $m_k = \text{rank}(\mathbf{X}_k)$ , and let  $(\mathbf{V}_k, d_k)$  denote the right singular vectors and singular values of  $\mathbf{X}_k$ . pcLasso minimizes

$$J(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 + \frac{\theta}{2} \sum_k \beta_k^T \left( \mathbf{V}_k \mathbf{D}_{d_{k1}^2 - d_{kj}^2} \mathbf{V}_k^T \right) \beta_k. \quad (8)$$

Here,  $\beta_k$  is the sub-vector of  $\beta$  corresponding to group  $k$ ,  $d_k = (d_{k1}, \dots, d_{km_k})$  are the singular values of  $\mathbf{X}_k$  in decreasing order, and  $\mathbf{D}_{d_{k1}^2 - d_{kj}^2}$  is a diagonal matrix with diagonal entries  $d_{k1}^2 - d_{kj}^2$  for  $j = 1, 2, \dots, m_k$ . Optionally, we might also multiply each penalty term in the sum by a factor  $\sqrt{p_k}$ ,  $p_k$  being the group size, as is standard in related methods such as the group lasso (Ming & Lin 2005).

Some observations on solving (8) are in order:

- The objective function is convex and the non-smooth component is separable. Hence, it can be optimized efficiently by a coordinate descent procedure (see Section 5 for details).
- Our approach is to fix a few values of  $\theta$  (including zero, corresponding to the lasso), and then optimize the objective (8) over a path of  $\lambda$  values. Cross-validation is used to choose  $\hat{\theta}$  and  $\hat{\lambda}$ .
- The parameter  $\theta$  is not unitless, and is difficult to interpret. In our software, instead of asking the user to specify  $\theta$ , we ask the user to specify a function of  $\theta$ , denoted by “rat”. This is the ratio between the shrinkage factors in (7) for  $k = 2$  and  $k = 1$  (the latter being equal to 1). This ratio is between 0 and 1, with 1 corresponding to  $\theta = 0$  (the lasso) and lower values corresponding to more shrinkage. In practice, we find that using the values  $rat = 0.25, 0.5, 0.75, 0.9, 0.95$  and 1 covers a wide range of possible models.

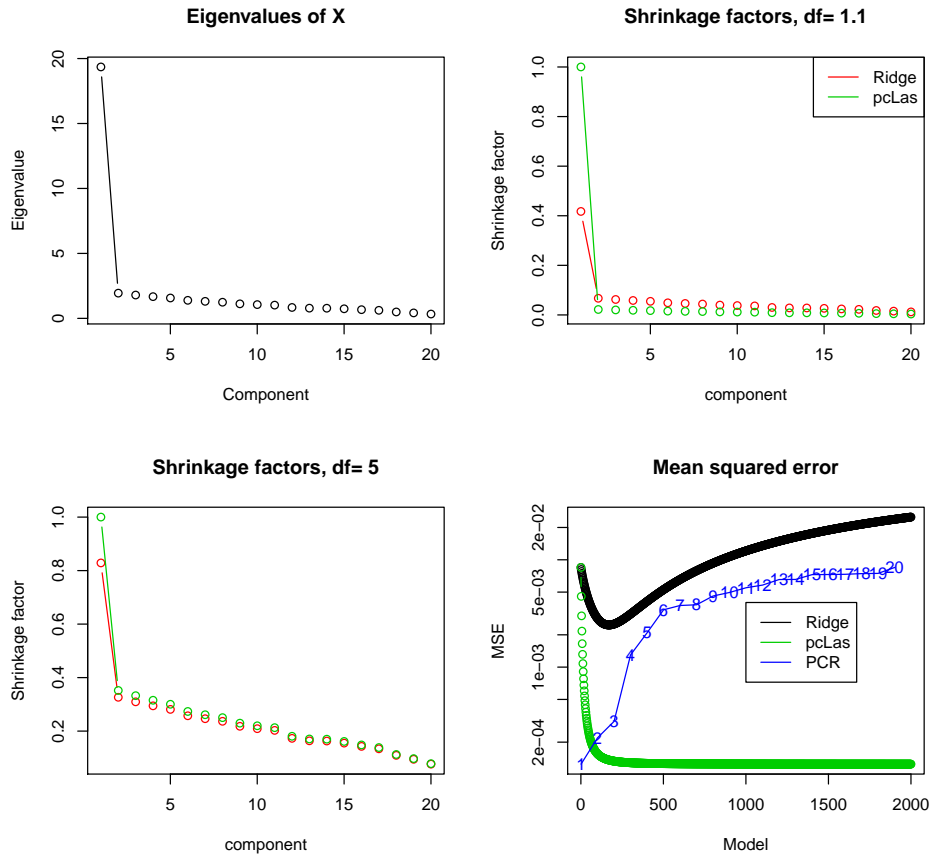


Figure 1: *Simulation example: The design matrix  $\mathbf{X}$  has dimensions  $100 \times 20$ , and has a strong rank-1 component: the eigenvalues of  $\mathbf{X}$  are shown in the top-left panel. Top-right and bottom-left panels show the shrinkage factors from the pcLasso penalty (5) (green) and ridge regression (red) across the PCs, for different degrees of freedom. For small degrees of freedom, pcLasso (5) shrinks more aggressively to the top PCs. For larger degrees of freedom, the two methods have similar shrinkage profiles. The bottom right panel shows the mean-squared error (MSE) of ridge regression and pcLasso as the regularization parameter  $\theta$  varies along the horizontal axis. Also included is the MSE for PC regression for various ranks. We see that pcLasso achieves about the same MSE as PC regression of rank-1, while ridge regression does not.*

- A potentially costly part of the computation is the SVD of each  $\mathbf{X}_k$ . If  $\mathbf{X}_k$  large, for computational ease, we can use an SVD of rank  $< \text{rank}(\mathbf{X}_k)$  as an approximation without much loss of performance.

For some insight into how pcLasso shrinks coefficients, we consider the contours of the penalty in the case of two predictors. Let  $\rho$  denote the correlation between these two predictors. In this case, it is easy to show that the two right singular vectors of  $\mathbf{X}$  are  $\frac{1}{\sqrt{2}}(1, 1)^T$  and  $\frac{1}{\sqrt{2}}(1, -1)^T$ , with the former (latter resp.) being the leading singular vector if  $\rho > 0$  ( $\rho < 0$  resp.). Figure 2 presents the penalty contours for pcLasso along with that for ridge, lasso and elastic net regression for comparison. (See Appendix A for detailed computations on the pcLasso contours.) We can see that pcLasso imposes a smaller penalty in the direction of the leading singular vector, and that this penalty is smaller when the correlation between the two predictors is stronger. We also see that as  $\theta \rightarrow 0$ , the pcLasso contours become like those for the lasso, while as  $\theta \rightarrow \infty$ , they become line segments in the direction of the leading singular vector.

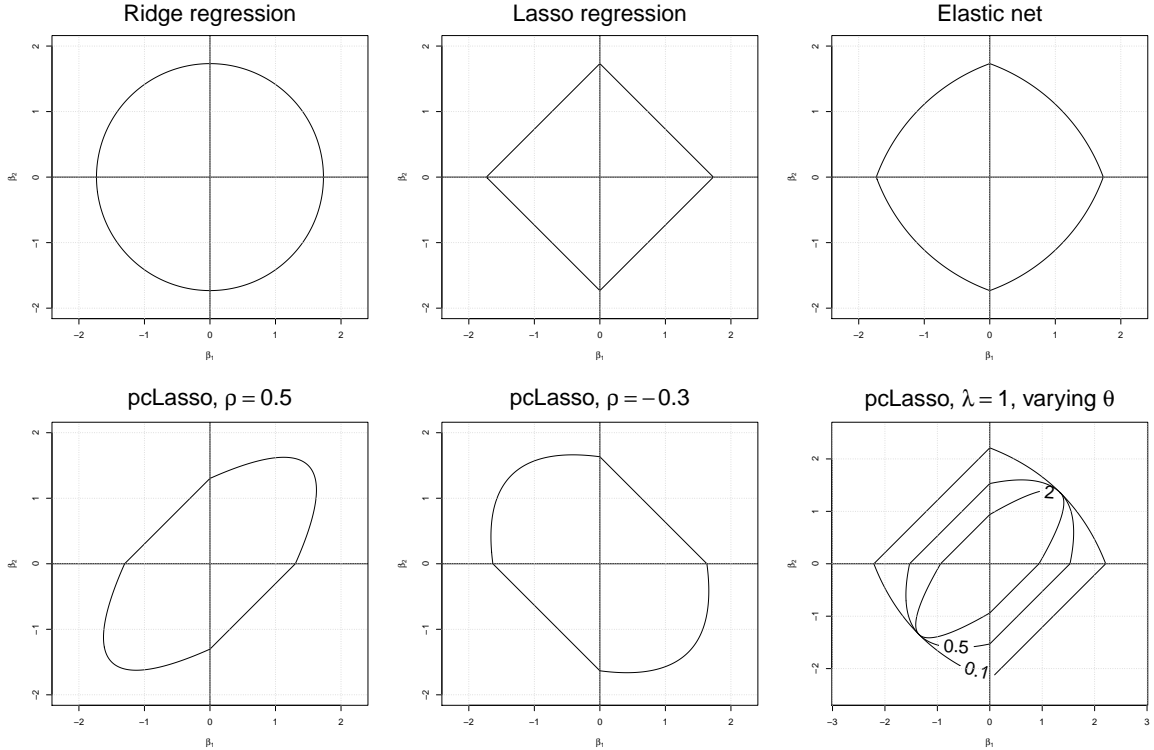


Figure 2: Contours for the pcLasso penalty, compared with those for ridge, lasso and elastic net regression. pcLasso penalizes the coefficient vector less in the direction of the leading singular vector. As  $\theta \rightarrow 0$ , the pcLasso contours become like those for the lasso, while as  $\theta \rightarrow \infty$ , they become line segments in the direction of the leading singular vector.

Figure 3 shows the penalty contours for a particular set of three predictors, with the ratio  $\lambda/\theta$  decreasing as we go from left to right. As this ratio goes from  $\infty$  to 0, the contours move from being the  $\ell_1$  ball (as in the lasso) to becoming more elongated in the direction of the first PC.

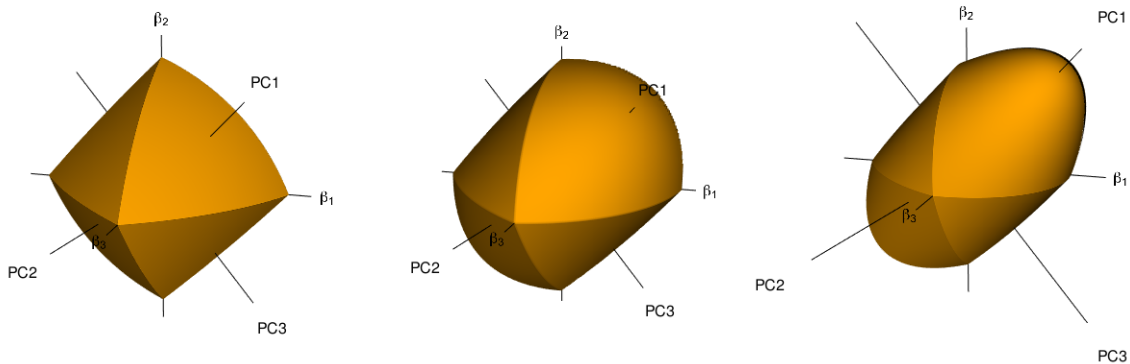


Figure 3: *Contours for the pcLasso penalty for three predictors. For a fixed value of  $\theta$ ,  $\lambda$  decreases as we go from left to right. We see that as we give less relative weight to the  $\ell_1$  penalty, the penalty contours look less and less like that of the lasso and become more elongated in the direction of the first PC.*

## 3.2 The overlapping groups setting

As in the group lasso (Yuan & Lin 2006), there are two approaches one can take to deal with overlapping groups. One approach is to apply the group penalty to the given groups: this however has the undesirable effect of zeroing out a predictor’s coefficient if any group to which it belongs is zeroed out. The overlap also causes additional computational challenges as the penalty is no longer separable. A better approach is the “overlap group lasso” penalty (Jacob et al. 2009), where one replicates columns of  $\mathbf{X}$  to account for multiple memberships, hence creating non-overlapping groups. The model is then fit as if the groups are non-overlapping and the coefficients for each original feature are summed to create the estimated coefficient for that feature. We take this same approach for pcLasso.

# 4 Real data examples

## 4.1 Blog Feedback dataset

This dataset is from the UC Irvine collection. Here is the description of the data and task:

*“This data originates from blog posts. The raw HTML-documents of the blog posts were crawled and processed. The prediction task associated with the data is the prediction of the number of comments in the upcoming 24 hours. In order to simulate this situation, we choose a basetime (in the past) and select the blog posts that were published at most 72 hours before the selected base date/time. Then, we calculate all the features of the selected blog posts from the information that was available at the basetime, therefore each instance various ranks, corresponds to a blog post. The target is the number of comments that the blog post received in the next 24 hours relative to the basetime.*

*In the train data, the basetimes were in the years 2010 and 2011. In the test data the basetimes were in February and March 2012. This simulates the real-world situation in which training data from the past is available to predict events in the future.”*

There are a total of 60,021 observations; we took a random sample of 1000 observations from training set (years 2010 and 2011), and 5 days —2012.03.27 through 2012.03.31— for the test data (a total of 817 observations). There are a total of 281 attributes (features). We applied pcLasso with no pre-specified groups of features. The left panel of Figure 4 shows the eigenvalues of the  $\mathbf{X}^T\mathbf{X}$  matrix, scaled so that the largest eigenvalue is 1. The eigenvalues for a random Gaussian matrix of

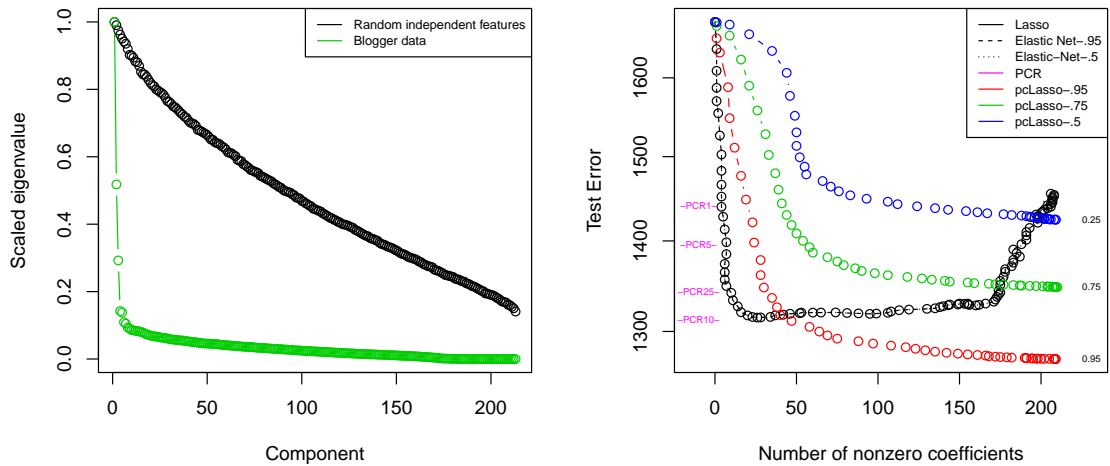


Figure 4: Results for Blogger Feedback dataset: Left panel shows the eigenvalues of the  $\mathbf{X}^T \mathbf{X}$  matrix, scaled so that the largest eigenvalue is 1. The eigenvalues for a random Gaussian matrix of the same size are also shown for reference. Right panel shows the test error for different methods. The lasso and elastic net have give similar results on this dataset. The errors for principal components regression (PCR) do not vary across the horizontal axis and so are marked on the left side of the panel.

the same size are also shown for reference. For the Blog Feedback dataset, there is a sharp drop-off in the eigenvalues after the first 3 or 4 eigenvalues, indicating strong correlation among the features.

The right panel shows the test error for the lasso, elastic net, PC regression of various ranks and pLasso for different values of the “rat” parameter. The results for elastic net are very close to that for the lasso. We see that pLasso achieves a somewhat lower test error by shrinking towards the top principal components.

## 4.2 p53 microarray expression data

Here, we analyze the data from a gene expression study, as described in Zeng & Breheny (2016), taken from Subramanian et al. (2005). It involves the mutational status of the gene p53 in cell lines. The study aims to identify pathways that correlate with the mutational status of p53, which regulates gene expression in response to various signals of cellular stress. The data consists of 50 cell lines, 17 of which are classified as normal and 33 of which carry mutations in the p53 gene. To be consistent with the analysis in Subramanian et al. (2005), 308 gene sets that have size between 15 and 500 are included in our analysis. These gene sets contain a total of 4301 genes and is available in the R package `grpregOverlap`. When the data is expanded to handle the overlapping groups, it contains a total of 13,237 columns.

We divided the data into 10 cross-validation (CV) folds and applied the lasso, pLasso and the group lasso using the R package `grpregOverlap`. (The data was too large for the sparse group lasso package, `SGL`.) The left panel of Figure 5 plots the number of non-zero pathways (i.e. pathways with at least one zero coefficient) against the number of non-zero coefficients for each method, as the complexity parameter of each method is varied. We see that pLasso induces some group-level sparsity, although not as strongly as the group lasso. However, the right panel shows that pLasso

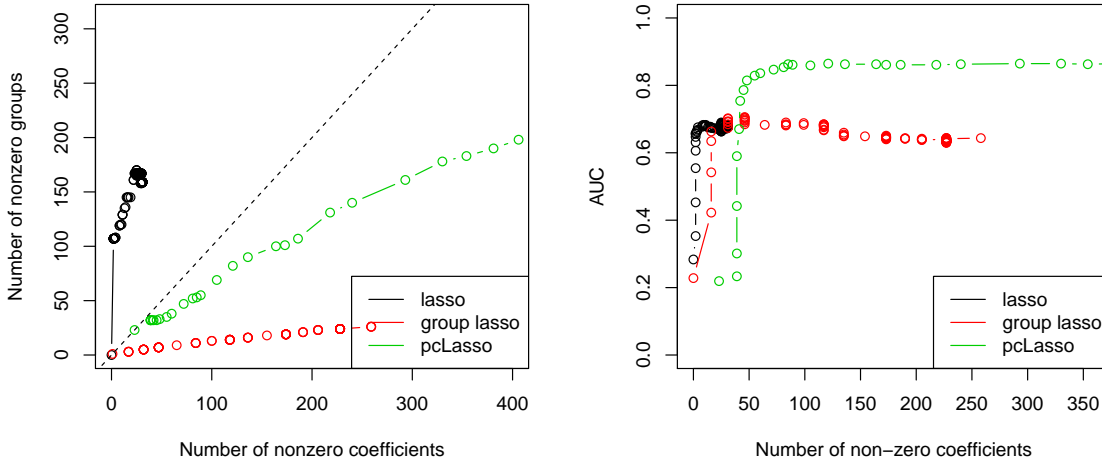


Figure 5: Results for *p53* dataset: Left panel shows the number of non-zero pathways vs. the number of non-zero coefficients for the lasso, group lasso and pcLasso. pcLasso exhibits some group-level sparsity, but not as strongly as the group lasso. The right panel is plot of CV area under the curve (AUC) vs. the number of non-zero coefficients in the selected model. pcLasso achieves the best result here.

achieves higher cross-validation area under the curve (AUC) than the other approaches.

## 5 Computation of the pcLasso solution

Consider first the simpler case of just one group. Let  $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$  be the singular value decomposition of  $\mathbf{X}$  with  $\mathbf{D} = \text{diag}(d_1, \dots, d_m)$ , where  $m = \text{rank}(\mathbf{X})$ , and let  $\mathbf{A} = \mathbf{V}\mathbf{D}_{d_1^2-d_j^2}\mathbf{V}^T$ . The objective function for pcLasso is

$$J(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 + \frac{\theta}{2} \beta^T \mathbf{A} \beta. \quad (9)$$

Since  $J$  is convex in  $\beta$  and the non-smooth part of the penalty is separable, we can minimize  $J$  by applying coordinate descent. The coordinate-wise update has the form

$$\tilde{\beta}_j \leftarrow \frac{\mathcal{S}\left(\sum_i x_{ij} r_i^{(j)} - \theta s_j, \lambda\right)}{\sum_i x_{ij}^2 + \theta A_{jj}},$$

where  $r_i^{(j)} = y_i - \sum_{j' \neq j} x_{ij'} \beta_{j'}$  is the partial residual with predictor  $j$  removed,  $s_j = \sum_\ell A_{j\ell} \beta_\ell - A_{jj} \beta_j$ , and  $\mathcal{S}$  is the soft-thresholding operator.

These updates can be easily generalized to the general case of  $K$  non-overlapping groups (8) and where observations are given different weights (observation  $i$  is given weight  $w_i$ ). Using the notation of Section 3, we apply coordinate descent as shown in Algorithm 1.

We also generalize this procedure to the binomial/logistic model, using the same Newton-Raphson (iteratively reweighted least squares) framework used in `glmnet` (Friedman et al. 2010).



---

**Algorithm 1** *Computation of the pcLasso solution*


---

1. Input: response  $\mathbf{y}_{n \times 1}$ , features  $\mathbf{X}_{n \times p}$ ,  $K$  non-overlapping groups of features each of size  $p_k$ , fixed value of  $\theta \geq 0$ . Observation weights  $w_i \geq 0$ .

Compute the SVD of the columns of  $X$  corresponding to each group:

$$\mathbf{X}_k = \mathbf{U}_k \mathbf{D}_k \mathbf{V}_k^T.$$

Compute  $\mathbf{A}^k = \mathbf{V}_k \mathbf{D}_{d_{k1}^2 - d_{kj}^2} \mathbf{V}_k^T$  for each group  $k$ . Given weights  $w_i$ , let  $v_j = \sum_i \tilde{w}_i x_{ij}^2$ , where  $\tilde{w}_i = n w_i / \sum_i w_i$ . Set  $\hat{\beta} = 0$ .

2. Define a grid of values  $(\lambda_{max}, \dots, \lambda_{min})$ , where  $\lambda_{max}$  is the smallest value of  $\lambda$  yielding  $\hat{\beta} = 0$ . For  $\lambda \in (\lambda_{max}, \dots, \lambda_{min})$ :

For  $j = 1, 2, \dots, p, 1, 2, \dots$  until convergence:

- (a) Compute the partial residual  $r_i^{(j)} = y_i - \sum_{j' \neq j} x_{ij'} \beta_{j'}$ .
- (b) Compute  $s_j = \sum_{\ell} A_{j\ell}^k \beta_{\ell} - A_{jj}^k \beta_j$ , where  $k$  is the group containing predictor  $j$ .
- (c) Update

$$\tilde{\beta}_j \leftarrow \frac{\mathcal{S}\left(\sum_i \tilde{w}_i x_{ij} r_i^{(j)} - \theta s_j, \lambda\right)}{v_j + \theta A_{jj}^k}.$$


---

## 5.1 pcLasso's grouping effect

We saw in Figure 5 that pcLasso can have a sparse grouping effect, even though it does not use a two-norm penalty. We investigate this further here. For Figure 6, we generated  $n = 50$  observations with 750 predictors, arranged in 50 groups of 15 predictors. Within each group, the predictors were either independent (left panel) or had a strong rank-1 component (right panel). The figure plots the number of non-zero groups (i.e. groups with at least one zero coefficient) against the number of non-zero coefficients for the lasso, the sparse group lasso (SGL) (Simon et al. 2013) and pcLasso. The latter two methods were run with different tuning parameters. We see that pcLasso shows no grouping advantage over the lasso in the left panel, but produces moderately strong grouping in the right panel. One might argue that this is reasonable: grouping only occurs when the groups have some strong internal structure. The sparse group lasso, on the other hand, always displays a strong grouping effect.

## 5.2 Connection to the group lasso and a generalization of pcLasso

The group lasso and sparse group lasso objective functions are given by

$$J_{GL}(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda_1 \sum_k \|\beta_k\|_2, \quad (10)$$

$$J_{SGL}(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1 + \lambda_1 \sum_k \|\beta_k\|_2, \quad (11)$$

respectively. In the original group lasso paper of Yuan & Lin (2006) it was assumed that the  $\mathbf{X}_k$ , the block of columns in  $\mathbf{X}$  corresponding to group  $k$ , were orthonormal, i.e.  $\mathbf{X}_k^T \mathbf{X}_k = \mathbf{I}$ . This makes the computation considerable easier and as pointed out by Simon & Tibshirani (2012), is equivalent to a penalty of the form  $\sum_k \|\mathbf{X}_k \beta_k\|_2$ . The R package `grpreg` uses this orthogonalization and states that it is essential to its computational performance (Breheny & Huang 2015). Since

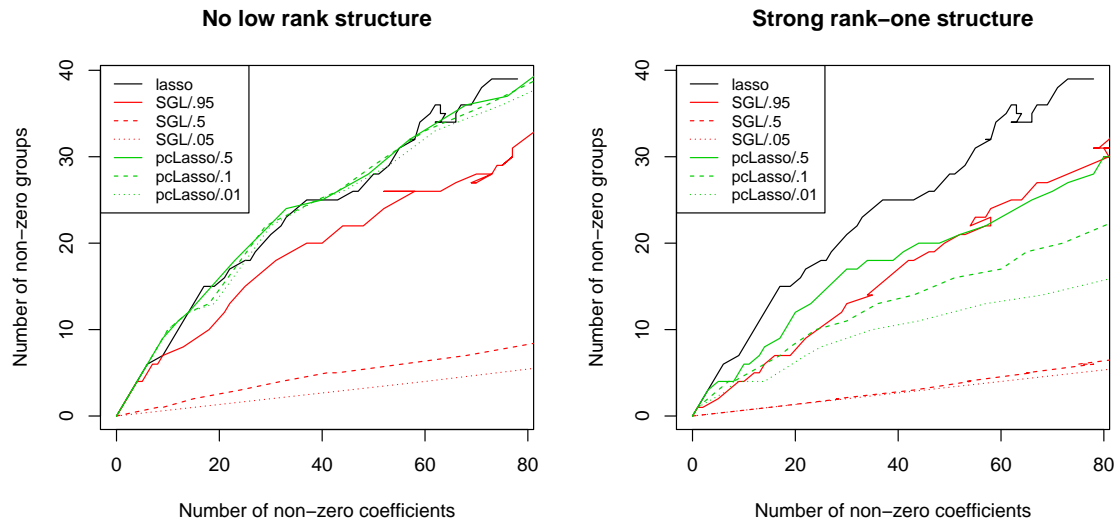


Figure 6: *Grouping results for the lasso, sparse group lasso (SGL) and pcLasso.  $n = 50$  and  $p = 750$ , with the predictors coming in 50 groups of size 15 each. Left panel: Predictors were independent. pcLasso has similar grouping behavior as the lasso. Right panel: Within each group, predictors had a strong rank-1 component. Here, pcLasso shows a stronger grouping effect than the lasso. SGL has a strong grouping effect in both settings.*

$\|\mathbf{X}_k \beta_k\|_2 = \|\mathbf{D}_k \mathbf{V}_k^T \beta_k\|_2$ , this orthogonalization penalizes the top eigenvectors more, in direct contrast to pcLasso which puts *less* penalty on the top eigenvectors. Of course, the group lasso without the  $\ell_1$  term does not deliver sparsity in the features. We also note that this orthogonalization does not work with the sparse group lasso, as the sparsity among the original features would be lost, and does not work whenever  $p_k > n$  for any group  $k$ .

One could envision a variation of pcLasso, namely a version of the sparse group lasso that uses an objective function of the form

$$\frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1 + \lambda_1 \sum_k \|\mathbf{R}_k^{1/2} \beta_k\|_2, \quad (12)$$

where  $\mathbf{R}_k$  is any weight matrix. For example,  $\mathbf{R}_k = \mathbf{V}_k \mathbf{D}_{d_{k1}^2 - d_{kj}^2} \mathbf{V}_k^T$  as in (8) would put less penalty on the higher eigenvectors and induce strong group sparsity. We have not experimented with this yet, as the computation seems challenging due to the presence of  $\ell_2$  norms.

### 5.3 Comparison of timings for model-fitting

Table 1 shows a comparison of timings for Algorithm 1 against that for other methods, all available as R packages. The competitors include the lasso `glmnet`, sparse group lasso SGL and group lasso packages `gglasso` and `grpreg`. As mentioned earlier, the latter uses orthogonalization within groups to speed up the computations, but this changes the problem. Among these competitors, only SGL provides sparsity at the level of individual features. All functions were called with default settings.

We see that SGL and `gglasso` start to slow down considerably as the problem size gets moderately large, while the times for `grpreg` are similar to that for pcLasso. However, when the cost for the initial SVD is separated out, the speed of pcLasso is roughly comparable to that of `glmnet`. In practice one can compute the SVD upfront for the given feature matrix, and then use it for

subsequent applications of pcLasso. For example, in cross-validation, one could compute just one SVD for the entire feature matrix, rather than one SVD in each fold. Table 2 considers larger problem sizes and we see these same general trends.

	n	p	glmnet	SGL	gglasso	grpreg	SVD for pcLasso	Rest	Total pcLasso
1	100	100	0.05	3.59	3.31	0.02	0.00	0.02	0.02
2	100	200	0.02	7.14	0.39	0.03	0.01	0.01	0.01
3	100	500	0.02	29.65	2.40	0.08	0.03	0.02	0.05
4	500	1000	0.19		14.89	1.12	0.25	0.13	0.38
5	1000	2000	0.67			6.38	1.50	0.54	2.04
6	1000	5000	1.10			17.26	11.16	2.35	13.51

Table 1: Comparison of timings for various algorithms. The predictors are pre-assigned to 10 groups. Time in seconds, average over three runs for an entire path of solutions. The last three columns show the pcLasso model-fitting times broken up by the initial SVD(s) and the rest of the computation. SGL refers sparse group lasso. gglasso and grpreg are group lasso packages; the latter does orthogonalization of predictors.

	n	p	glmnet	grpreg	SVD for pcLasso	Rest	Total pcLasso
1	500	1000	0.19	1.92	0.27	0.17	0.43
2	1000	2000	0.76	5.33	1.76	0.54	2.30
3	2000	5000	2.49	36.57	24.15	3.81	27.96
4	2000	10000	4.50	139.63	106.56	12.96	119.52

Table 2: As in the previous table, but for larger problem sizes and focusing on just glmnet, grpreg and pcLasso.

## 6 Degrees of freedom

Given a vector of response values  $\mathbf{y}$  with corresponding fits  $\hat{\mathbf{y}}$ , Efron (1986) defines the degrees of freedom as

$$\text{df}(\hat{\mathbf{y}}) = \frac{\sum_i \text{Cov}(\mathbf{y}_i, \hat{\mathbf{y}}_i)}{\sigma^2}. \quad (13)$$

The degrees of freedom measures the flexibility of the fit: the larger the degrees of freedom, the more closely the fit matches the response values. For fits of the form  $\hat{\mathbf{y}} = \mathbf{M}\mathbf{y}$ , the degrees of freedom is given by  $\text{df}(\hat{\mathbf{y}}) = \text{tr}(\mathbf{M})$ . For the lasso, the number of non-zero elements in the solution is an unbiased estimate of the degrees of freedom (Zou et al. 2007). Zou (2005) derived an unbiased estimate for the degrees of freedom of the elastic net: if the elastic net solves

$$\underset{\beta}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2, \quad (14)$$

then

$$\hat{\text{df}} = \text{tr} [\mathbf{X}_{\mathcal{A}}(\mathbf{X}_{\mathcal{A}}^T \mathbf{X}_{\mathcal{A}} + \lambda_2 \mathbf{I})^{-1} \mathbf{X}_{\mathcal{A}}^T] \quad (15)$$

is an unbiased estimate for the degrees of freedom, where  $\mathcal{A}$  is the active set of the fit.

Using similar techniques as that of Zou (2005), we can derive an unbiased estimate for the degrees of freedom of pcLasso when the number of groups  $K$  is equal to one:

**Theorem 1.** Let pcLasso be the solution to

$$\underset{\beta}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1 + \theta \beta^T \mathbf{V} \mathbf{D}_{d_1^2 - d_2^2} \mathbf{V}^T \beta, \quad (16)$$

where  $\mathbf{D}_{d_1^2-d_j^2}$  is a  $p \times p$  matrix with entries  $d_1^2 - d_j^2$  ( $m = \text{rank}(\mathbf{X}), d_{m+1} = \dots = d_p = 0$ ). Let  $\mathbf{W} = (\mathbf{V}\mathbf{D}_{d_1^2-d_j^2}\mathbf{V}^T)^{1/2}$ . Then an unbiased estimate for the degrees of freedom for pcLasso is

$$\widehat{df} = \text{tr}[\mathbf{X}_{\mathcal{A}}(\mathbf{X}_{\mathcal{A}}^T\mathbf{X}_{\mathcal{A}} + \theta\mathbf{W}_{\mathcal{A}}^T\mathbf{W}_{\mathcal{A}})^{-1}\mathbf{X}_{\mathcal{A}}^T], \quad (17)$$

where  $\mathcal{A}$  is the active set of the fit.

We verify this estimate through a simulation exercise. Consider the model

$$\mathbf{y}^* = \mathbf{X}\beta + \mathcal{N}(0, 1)\sigma. \quad (18)$$

Given this model, for each value of  $\lambda$ , we can evaluate the degrees of freedom of pcLasso via Monte Carlo methods. For  $b = 1, \dots, B$ , we generate a new response vector  $\mathbf{y}^{*b}$  according to (18). We fit pcLasso to this data, generating predictions  $\widehat{\mathbf{y}}^{*b}$ . This gives us the Monte Carlo estimate

$$df(\lambda) \approx \sum_{i=1}^n \widehat{\text{Cov}}(\widehat{\mathbf{y}}_i^*, \mathbf{y}_i^*) / \sigma^2,$$

$$\widehat{\text{Cov}}(\widehat{\mathbf{y}}_i^*, \mathbf{y}_i^*) = \frac{1}{B} \sum_{b=1}^B [\widehat{\mathbf{y}}_i^{*b} - a_i][\mathbf{y}_i^{*b} - \mathbf{y}_i^*],$$

where the  $a_i$ 's can be any fixed known constants (usually taken to be 0). On the other hand, each pcLasso fit gives us an active set  $\mathcal{A}^{*b}$  which allows us to compute an estimate for the degrees of freedom  $\widehat{df}(\lambda)^{*b}$  based on the formula (17). With  $B$  replications, we can estimate  $\mathbb{E}[\widehat{df}(\lambda)] \approx \frac{1}{B} \sum_{b=1}^B \widehat{df}(\lambda)^{*b}$ .

Figure 7 provides evidence for the correctness of (17). The figures on the left are plots of the true value of the degrees of freedom against the mean of the estimate given by (17), with each point corresponding to a value of  $\lambda$ . We can see that for both values of  $\theta$ , there is close agreement between the true value  $df(\lambda)$  and the expectation of the estimate  $\mathbb{E}[\widehat{df}(\lambda)]$ . The figures of the right are plots of the bias  $\mathbb{E}[\widehat{df}(\lambda)] - df(\lambda)$  along with pointwise 95% confidence intervals. The zero horizontal line lies inside these confidence intervals.

## 7 A simulation study

We tested the performance of pcLasso against other methods in an extensive simulation study. The general framework of the simulation is as follows: The training data has a sample size of  $n$  with  $p$  features which fall into  $K$  groups. Denote the design matrix for group  $k$  by  $\mathbf{X}_k$ , and let the SVD of  $\mathbf{X}_k$  be  $\mathbf{X}_k = \mathbf{U}_k\mathbf{D}_k\mathbf{V}_k^T$ . The true signal is a linear combination of the eigenvectors of the  $\mathbf{X}_k$ 's, i.e.  $\text{signal} = \sum_k \mathbf{X}_k\mathbf{V}_k\mathbf{b}_k$ , where the  $\mathbf{b}_k$  are the coefficients of the linear combination. The response  $\mathbf{y}$  is the signal corrupted by additive Gaussian noise. We consider three different scenarios:

- “Home court” for pcLasso: The signal is related to the top eigenvectors of each group, i.e. the non-zero entries of the  $\mathbf{b}_k$  are all at the beginning.
- “Neutral court” for pcLasso: The signal is related to random eigenvectors of each group, i.e. the non-zero entries of the  $\mathbf{b}_k$  are at random positions.
- “Hostile court” for pcLasso: The signal is related to the bottom eigenvectors of each group, i.e. the non-zero entries of the  $\mathbf{b}_k$  are all at the end.

To induce sparsity in the set-up, we also set  $\mathbf{b}_k = \mathbf{0}$  for some  $k$ , corresponding to group  $k$  having no effect on the response. Within each set-up, we looked at a range of signal-to-noise (SNR) ratios in the response, as well as whether there were correlations between predictors in the same group.

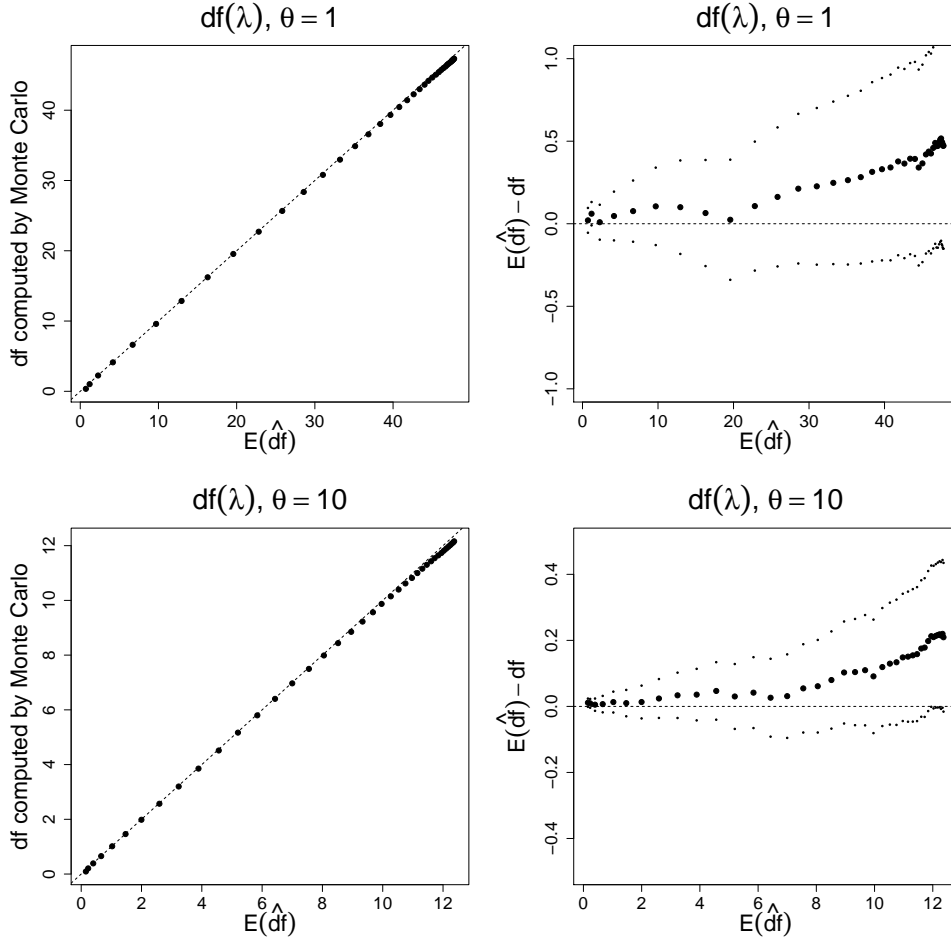


Figure 7: . Degrees of freedom for pcLasso. The synthetic model is  $\mathbf{y}^* = \mathbf{X}\beta + \mathcal{N}(0, 1)\sigma$  with  $\beta = 0$ , i.e. the null model, and  $\sigma = 2$ . We set  $n = 500$ ,  $p = 100$ , and did 500 Monte Carlo replications. The top two plots are for  $\theta = 1$  while the bottom two plots are for  $\theta = 10$ . In the panels on the left, we compare the true  $df(\lambda)$  against  $\mathbb{E}[\hat{df}(\lambda)]$ , with the dotted line being the 45° line (i.e. perfect match). In the panels on the right, the larger dots represent the bias  $\mathbb{E}[\hat{df}(\lambda)] - df(\lambda)$  across the range of  $\mathbb{E}[\hat{df}(\lambda)]$ , and the smaller dots represent the pointwise 95% confidence intervals. Note that the zero horizontal line lies inside these confidence intervals.

For pcLasso, we used the following cross-validation (CV) procedure to select the tuning parameters: For each value of  $rat = 0.25, 0.5, 0.75, 0.9, 0.95$  and 1, we ran pcLasso along a path of  $\lambda$  values with  $rat = rat$ . (We found that these values of  $rat$  covered a good range of models in practice.) For each run, the value of  $\lambda$  which gave the smallest CV error was selected, i.e. the `lambda.min` value as in `glmnet`. We then compared the CV error across the 6 values of  $rat$  and selected the value of  $rat$  with the smallest CV error and its accompanying  $\lambda$  value. A second version of pcLasso was also run using the same procedure, but with the  $\lambda$  values selected by the one standard error rule, i.e. the `lambda.1se` value as in `glmnet`.

To compare the methods, we considered the mean-squared error (MSE) achieved on 5,000 test points, as well as the support size of the fitted model. We benchmarked pcLasso against the following methods:

- The null model, i.e. the mean of the responses in the training data.
- Elastic net with CV across  $\alpha = 0, 0.2, 0.4, 0.6, 0.8$  and 1, with  $\lambda$  values selected both by `lambda.min` and `lambda.1se`.
- Lasso with CV, with  $\lambda$  values selected both by `lambda.min` and `lambda.1se`.
- Principal components (PC) regression with CV across ranks.

Overall, we found that in “home court” settings, pcLasso performs the best in terms of test MSE across the range of SNR and feature correlation settings. The gain in performance appears to be larger in low SNR settings compared to high SNR settings. In “neutral court” and “hostile court” settings, pcLasso performs on par with the elastic net and the lasso. This is because the cross-validation step often picks  $rat = 1$  in these settings, under which pcLasso is equivalent to the lasso. In terms of support size, when there is sparsity pcLasso with  $\lambda$  values selected by the one standard error rule tends to pick models which have support size close to the underlying truth.

Below, we present an illustrative sampling of the results: see Appendix D for more comprehensive results across a wider range of simulation settings.

In the first setting,  $n = 200$ ,  $p = 1,000$  with the features coming in 10 groups of 100 predictors. The response is a linear combination of the top two eigenvectors of just the first group. The performance on test MSE is shown in Figure 8. pcLasso clearly outperforms the other methods when the predictors within each group are uncorrelated. The performance gain is smaller when predictors within each group are correlated with each other. In terms of support size, while pcLasso with  $\lambda$  values selected both by `lambda.min` seems to select models which are too large, pcLasso with  $\lambda$  values selected both by `lambda.1se` has support size closer to the underlying truth.

In the second setting,  $n = 200$ ,  $p = 200$  with the features coming in 10 groups of 20 predictors. The response is a linear combination of two random eigenvectors of just the first group. The performance on test MSE is shown in Figure 9. pcLasso performs comparably to both the elastic net and the lasso, both when the predictors are uncorrelated and when they are correlated. In terms of support size, pcLasso with  $\lambda$  values selected both by `lambda.1se` has support size close to the underlying truth, especially when SNR is high.

In the third setting,  $n = 200$ ,  $p = 50$  with the features coming in 5 groups of 10 predictors. The response is a linear combination of the bottom eigenvectors of the first two groups. The performance on test MSE is shown in Figure 10. pcLasso performs comparably to both the elastic net and the lasso, both when the predictors are uncorrelated and when they are correlated. In terms of support size, pcLasso with  $\lambda$  values selected both by `lambda.1se` has support size close to the underlying truth when the signal-to-noise ratio is high.

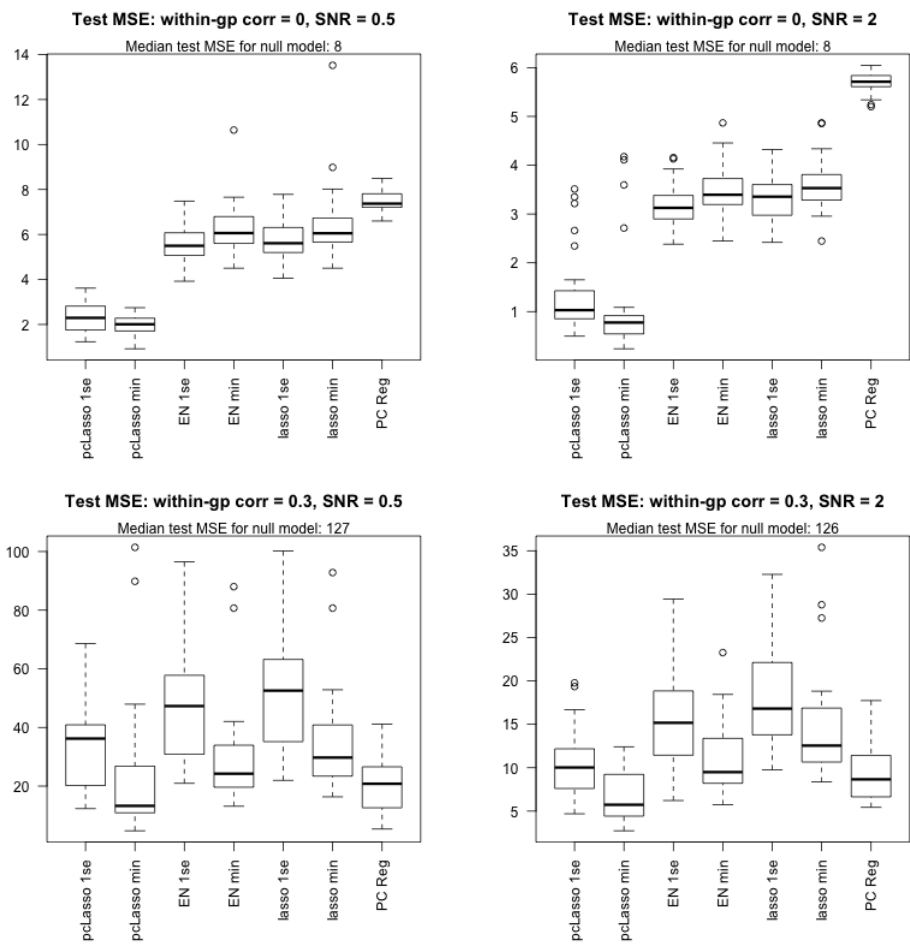


Figure 8: “Home court” simulation:  $n = 200$ ,  $p = 1,000$  with features coming in 10 groups of 100 predictors. The response is a linear combination of the top two eigenvectors of just the first group. For the top two figures, the predictors are uncorrelated while for the bottom two figures, each pair of predictors within each group has correlation 0.3. pcLasso outperforms the other methods across a range of signal-to-noise (SNR) ratios, with the performance gain being larger in the case of uncorrelated predictors.

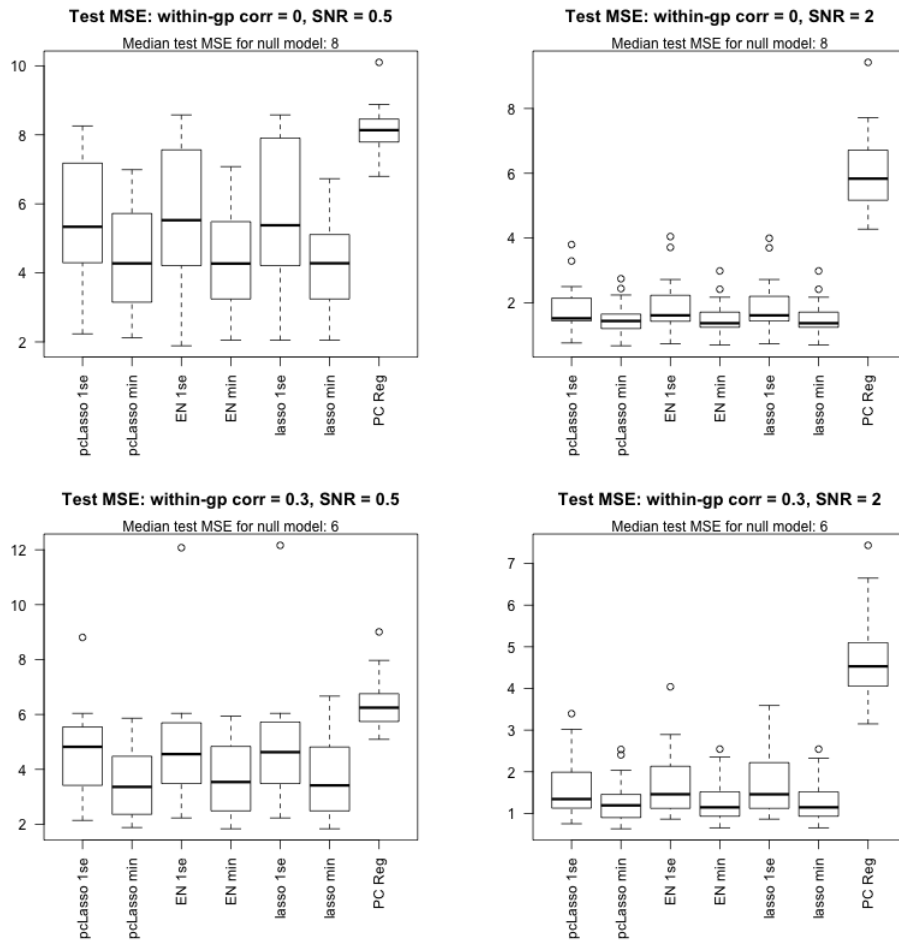


Figure 9: “Neutral court” simulation:  $n = 200$ ,  $p = 200$  with features coming in 10 groups of 20 predictors. The response is a linear combination of two random eigenvectors of just the first group. For the top two figures, the predictors are uncorrelated while for the bottom two figures, each pair of predictors within each group has correlation 0.3. pcLasso performs comparably to the elastic net and the lasso.



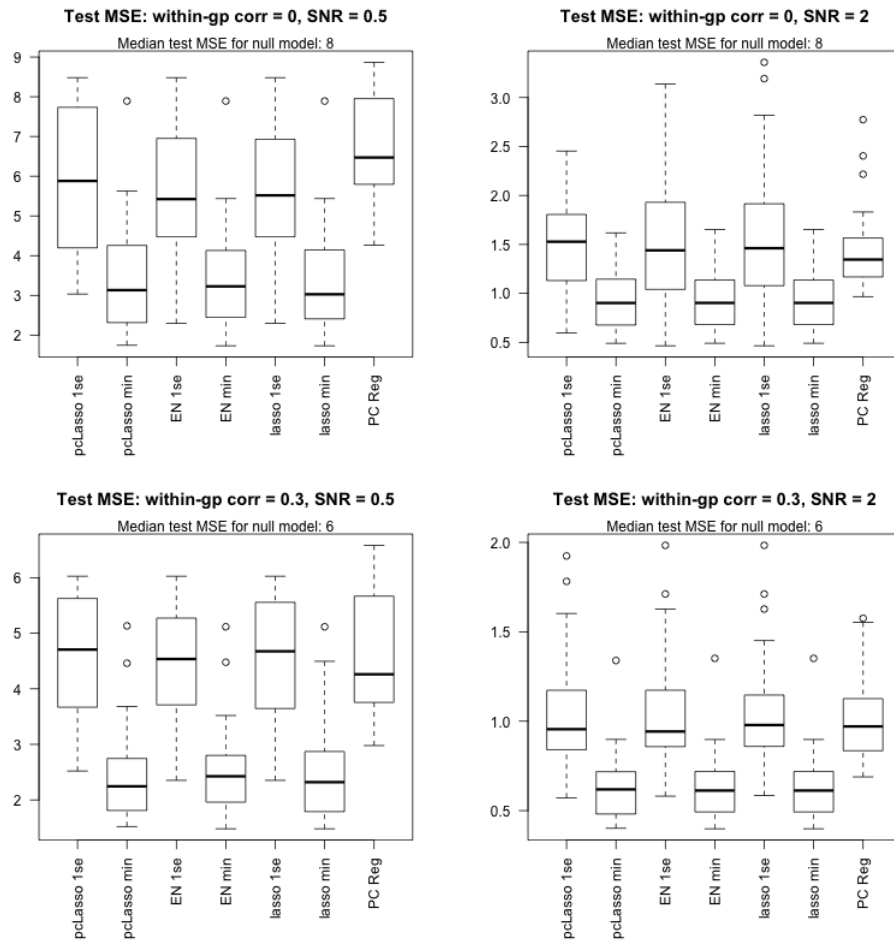


Figure 10: “Hostile court” simulation:  $n = 200$ ,  $p = 50$  with features coming in 5 groups of 10 uncorrelated predictors. The response is a linear combination of the bottom eigenvectors of the first two groups. For the top two figures, the predictors are uncorrelated while for the bottom two figures, each pair of predictors within each group has correlation 0.3. pcLasso performs comparably to the elastic net and the lasso.

## 8 Theoretical properties of pcLasso

In this section, we present some theoretical properties of pcLasso when the number of groups  $K$  is equal to 1, and compare them with that for the lasso as presented in Chapter 11 of Hastie et al. (2015). In particular, we provide non-asymptotic bounds for  $\ell_2$  and prediction error for pcLasso which are better than that for the lasso in certain settings. We note that the results in the section can be extended analogously to pcLasso with  $K$  non-overlapping groups if the  $\mathbf{X}_k$  are orthogonal to each other.

To make the proofs simpler, we consider a slightly different penalty for pcLasso: instead of the second penalty term having  $\mathbf{D}_{d_1^2-d_j^2}$  as an  $m \times m$  matrix, where  $m = \text{rank}(X)$ , we have  $\mathbf{D}_{d_1^2-d_j^2}$  as a  $p \times p$  diagonal matrix with  $d_{m+1} = \dots = d_p = 0$ . Let  $\mathbf{A} = \mathbf{V}\mathbf{D}_{d_1^2-d_j^2}\mathbf{V}^T$ .

As the proofs are rather technical, we state just the results here; proofs can be found in Appendix C.

### 8.1 Set-up

We assume that the true underlying model is

$$\mathbf{y} = \mathbf{X}\beta^* + \mathbf{w}, \quad (19)$$

where  $\mathbf{X} \in \mathbb{R}^{n \times p}$  is the design matrix,  $\mathbf{y}, \mathbf{w} \in \mathbb{R}^n$  are the response and noise vectors respectively, and  $\beta^* \in \mathbb{R}^p$  is the true unknown coefficient vector. We denote the support of  $\beta^*$  by  $S$ . When we solve an optimization problem for  $\beta$ , we denote the estimate by  $\hat{\beta}$  and the error by  $\hat{\nu} = \hat{\beta} - \beta^*$ .

In this section, assume that  $\theta > 0$  is fixed. (If  $\theta = 0$ , pcLasso is equivalent to the lasso.) Define

$$\tilde{\mathbf{y}} = \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix}, \quad \tilde{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \sqrt{n\theta}\sqrt{\mathbf{A}} \end{pmatrix}, \quad \tilde{\mathbf{w}} = \begin{pmatrix} \mathbf{w} \\ -\sqrt{n\theta}\sqrt{\mathbf{A}}\beta^* \end{pmatrix}. \quad (20)$$

Thus,  $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\beta^* + \tilde{\mathbf{w}}$ . The key idea is that with this notation, pcLasso is equivalent to the lasso for the augmented matrices  $\tilde{\mathbf{X}}$  and  $\tilde{\mathbf{y}}$ . Explicitly, the constrained form of pcLasso solves

$$\underset{\beta}{\text{minimize}} \quad \left\| \tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta \right\|_2^2 \quad \text{subject to } \|\beta\|_1 \leq R, \quad (21)$$

while the Lagrangian form (5) solves

$$\underset{\beta}{\text{minimize}} \quad \frac{1}{2n} \left\| \tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta \right\|_2^2 + \lambda \|\beta\|_1. \quad (22)$$

(We have changed the fraction in front of the RSS term to  $\frac{1}{2n}$  so that the results are more directly comparable to those in Hastie et al. (2015).) In the high-dimensional setting, it is standard to impose a *restricted eigenvalue condition* on the design matrix  $\mathbf{X}$ :

$$\frac{1}{n} \nu^T \mathbf{X}^T \mathbf{X} \nu \geq \gamma \|\nu\|_2^2 \quad \text{for all nonzero } \nu \in \mathcal{C}, \quad (23)$$

with  $\gamma > 0$  and  $\mathcal{C}$  an appropriately chosen constraint set. We assume additionally that  $d_1^2 > n\gamma$ . This is not a restrictive assumption: since  $d_1^2$  is the top eigenvalue of  $\mathbf{X}^T \mathbf{X}$ , we automatically have  $d_1^2 \geq n\gamma$ . Equality can only happen if  $\mathcal{C}$  is a subset of the eigenspace associated with  $d_1^2$ .

Since  $\mathbf{A}$  is a positive semidefinite matrix, (23) holds trivially for the augmented matrix  $\tilde{\mathbf{X}}$  as well:

$$\nu^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \nu = \nu^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \nu \geq n\gamma \|\nu\|_2^2.$$

The following key lemma shows that we can improve the constant on the RHS of (23). This will give us better rates of convergence for pcLasso.

**Lemma 2.** *If  $\mathbf{X}$  satisfies the restricted eigenvalue condition (23), then*

$$\nu^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \nu \geq \min [(1 - n\theta)n\gamma + n\theta d_1^2, d_1^2] \|\nu\|_2^2 \quad \text{for all nonzero } \nu \in \mathcal{C}. \quad (24)$$

We note also that the augmented matrix actually satisfies an *unrestricted eigenvalue condition*:

**Lemma 3.** *For any design matrix  $\mathbf{X}$ , the augmented data matrix  $\tilde{\mathbf{X}}$  satisfies*

$$\nu^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \nu \geq \min(n\theta, 1) d_1^2 \|\nu\|_2^2 \quad \text{for all } \nu \in \mathbb{R}^p. \quad (25)$$

This will give us a better “slow rate” for convergence of prediction error, as well as bounds on  $\ell_2$  error without requiring a restricted eigenvalue condition.

## 8.2 Bounds on $\ell_2$ error

The following theorem establishes a bound for  $\ell_2$  error of the constrained form of pcLasso:

**Theorem 4.** *Suppose that  $\mathbf{X}$  satisfies the restricted eigenvalue bound (23) with parameter  $\gamma > 0$  over the set  $\{\nu \in \mathbb{R}^p : \|\nu_{S^c}\|_1 \leq \|\nu_S\|_1\}$ . Then any estimate  $\hat{\beta}$  based on the constrained pcLasso (21) with  $R = \|\beta^*\|_1$  satisfies the bound*

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{4\sqrt{|S|} \|\mathbf{X}^T \mathbf{w} - n\theta \mathbf{A} \beta^*\|_\infty}{\min[(1 - n\theta)n\gamma + n\theta d_1^2, d_1^2]} \leq \frac{4\sqrt{|S|} (n\theta \|\mathbf{A} \beta^*\|_\infty + \|\mathbf{X}^T \mathbf{w}\|_\infty)}{\min[(1 - n\theta)n\gamma + n\theta d_1^2, d_1^2]}. \quad (26)$$

In particular, if  $\hat{\beta}$  is aligned with the first principal component of  $\mathbf{X}$ , then

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{4\sqrt{|S|} \|\mathbf{X}^T \mathbf{w}\|_\infty}{\min[(1 - n\theta)n\gamma + n\theta d_1^2, d_1^2]}, \quad (27)$$

which is a **better** rate of convergence than that for the lasso.

What is interesting is that we can actually obtain a similar (weaker) bound without the restricted eigenvalue condition; the lasso does not have such a bound.

**Theorem 5.** *For any design matrix  $\mathbf{X}$ , any estimate  $\hat{\beta}$  based on the constrained pcLasso (21) with  $R = \|\beta^*\|_1$  satisfies the bound*

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{4\sqrt{|S|} \|\mathbf{X}^T \mathbf{w} - n\theta \mathbf{A} \beta^*\|_\infty}{\min(n\theta, 1) d_1^2}. \quad (28)$$

We can obtain similar results for the Lagrangian form of pcLasso:

**Theorem 6.** *Suppose that  $\mathbf{X}$  satisfies the restricted eigenvalue bound (23) with parameter  $\gamma > 0$  over the set  $\{\nu \in \mathbb{R}^p : \|\nu_{S^c}\|_1 \leq 3\|\nu_S\|_1\}$ . Then any estimate  $\hat{\beta}$  based on the Lagrangian form of pcLasso (22) with  $\lambda \geq \frac{2}{n} \|\mathbf{X}^T \mathbf{w} - n\theta \mathbf{A} \beta^*\|_\infty > 0$  satisfies the bound*

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{3\lambda\sqrt{|S|}}{\min[(1 - n\theta)n\gamma + n\theta d_1^2, d_1^2]/n}. \quad (29)$$

If we remove the restricted eigenvalue condition on  $\mathbf{X}$ , then the bound (29) holds but with the denominator of the RHS being  $\min(n\theta, 1) d_1^2/n$  instead.

### 8.3 Bounds on prediction error

Like the lasso, pcLasso has a “slow rate” convergence for prediction error:

**Theorem 7.** For any design matrix  $\mathbf{X}$ , consider the Lagrangian form of pcLasso (22) with  $\lambda \geq \frac{2}{n} \|\mathbf{X}^T \mathbf{w} - n\theta \mathbf{A} \beta^*\|_\infty$ . We have the following bound on prediction error:

$$\frac{1}{n} \left\| \mathbf{X}(\hat{\beta} - \beta^*) \right\|_2^2 \leq 12\lambda \|\beta^*\|_1. \quad (30)$$

By using Lemma 3, we can get a smaller bound for the RHS, although the expression is much more complicated:

**Theorem 8.** For any design matrix  $\mathbf{X}$ , consider the Lagrangian form of pcLasso (22) with  $\lambda \geq \frac{2}{n} \|\mathbf{X}^T \mathbf{w} - n\theta \mathbf{A} \beta^*\|_\infty$ . Let  $\kappa = \min(n\theta, 1)d_1^2$ . We have the following bound on prediction error:

$$\frac{1}{n} \left\| \mathbf{X}(\hat{\beta} - \beta^*) \right\|_2^2 \leq \frac{3\lambda(-\lambda p + \sqrt{\lambda^2 p^2 + 32\lambda \|\beta^*\|_1 \kappa p})}{4\kappa}. \quad (31)$$

This means that pcLasso has a better “slow rate” convergence than the lasso.

The preceding two theorems hold for all design matrices  $\mathbf{X}$ . If we are willing to assume that  $\mathbf{X}$  satisfies the restricted eigenvalue condition, we recover the same “fast rate” convergence for prediction error as that for the lasso.

**Theorem 9.** Suppose that  $\mathbf{X}$  satisfies the restricted eigenvalue bound (23) with parameter  $\gamma > 0$  over the set  $\{\nu \in \mathbb{R}^p : \|\nu_{S^c}\|_1 \leq 3\|\nu_S\|_1\}$ . For the Lagrangian form of pcLasso with  $\lambda \geq \frac{2}{n} \|\mathbf{X}^T \mathbf{w} - n\theta \mathbf{A} \beta^*\|_\infty$ , we have

$$\frac{1}{n} \left\| \mathbf{X}(\hat{\beta} - \beta^*) \right\|_2^2 \leq \frac{9|S|\lambda^2}{\min[(1 - n\theta)n\gamma + n\theta d_1^2, d_1^2]/n}. \quad (32)$$

## 9 Discussion

We have proposed a new method for supervised learning that exploits the correlation and group structure of predictors to potentially improve prediction performance. It combines the power of PC regression with the sparsity of the lasso. There are several interesting avenues for further research:

- *Efficient screening rules:* To speed up the computation of the lasso solution, `glmnet` uses strong rules (Tibshirani et al. 2012) to discard predictors which are likely to have a coefficient of zero. These strong rules can greatly reduce the number of variables entering the optimization, thus speeding up computation. Using similar techniques, we can derive strong rules for pcLasso to improve computational efficiency. These rules are provided in Appendix B, and merit further investigation.
- *Use of different kinds of correlation or grouping information.* Rather than use the covariance of the observed covariates in a group, one could use other kinds of external information to form this covariance, for example gene or protein pathways. Alternatively, one could use unsupervised clustering to generate the feature groups.
- *Optimality of the quadratic penalty term:* Considering the general class of penalty functions of the form  $\beta^T \mathbf{V} \mathbf{Z} \mathbf{V}^T \beta$ , one could ask if pcLasso’s choice of  $\mathbf{Z} = \mathbf{D}_{d_1^2 - d_2^2}$  is “optimal” in any sense. One could also look for  $\mathbf{Z}$  which satisfy certain notions of optimality.

An R language package `pcLasso` will soon be made available on the CRAN repository.

**Acknowledgements:** We’d like to thank Trevor Hastie, Jacob Bien and Noah Simon for helpful comments. Robert Tibshirani was supported by NIH grant 5R01 EB001988-16 and NSF grant 19 DMS1208164.

## A Details of the pcLasso penalty contours for two predictors

Assume that we only have two predictors which are standardized to have mean 0 and sum of squares

1. Then  $\mathbf{X}^T \mathbf{X} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$ , where  $\rho$  is the sample correlation between the two predictors. Let  $\mathbf{A} = \mathbf{V} \mathbf{D}_{d_1^2 - d_2^2} \mathbf{V}^T$ . If  $\rho > 0$ , the expressions for the SVD of  $\mathbf{X}$  and  $\mathbf{A}$  are

$$\mathbf{X} = \mathbf{U} \cdot \begin{pmatrix} 1+\rho & 0 \\ 0 & 1-\rho \end{pmatrix} \cdot \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix}, \quad \mathbf{A} = 2\rho \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}.$$

If  $\rho < 0$ , the corresponding expressions are

$$\mathbf{X} = \mathbf{U} \cdot \begin{pmatrix} 1-\rho & 0 \\ 0 & 1+\rho \end{pmatrix} \cdot \begin{pmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{pmatrix}, \quad \mathbf{A} = -2\rho \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}.$$

With these expressions, the pcLasso penalty can be written explicitly as

$$\begin{cases} \lambda(|\beta_1| + |\beta_2|) + 2\theta\rho(\beta_1 - \beta_2)^2 & \text{if } \rho > 0, \\ \lambda(|\beta_1| + |\beta_2|) - 2\theta\rho(\beta_1 + \beta_2)^2 & \text{if } \rho < 0. \end{cases}$$

Considering the signs of  $\beta_1$  and  $\beta_2$ , we can get an even more explicit description of the contours. We provide the description for  $\rho > 0$  below:

- $\beta_1 \geq 0, \beta_2 \geq 0$ : The parabola  $\beta_2 = -\frac{2\sqrt{2}\theta\rho}{\lambda}\beta_1^2 + \frac{C}{\sqrt{2}\lambda}$ , rotated 45° clockwise.
- $\beta_1 \geq 0, \beta_2 \leq 0$ : The line  $\beta_2 = \beta_1 + \frac{\lambda - \sqrt{\lambda^2 + 8\theta\rho C}}{4\theta\rho}$ .
- $\beta_1 \leq 0, \beta_2 \leq 0$ : The parabola  $\beta_2 = \frac{2\sqrt{2}\theta\rho}{\lambda}\beta_1^2 - \frac{C}{\sqrt{2}\lambda}$ , rotated 45° clockwise.
- $\beta_1 \leq 0, \beta_2 \geq 0$ : The line  $\beta_2 = \beta_1 - \frac{\lambda - \sqrt{\lambda^2 + 8\theta\rho C}}{4\theta\rho}$ .

## B Strong rules for pcLasso

Recall that in computing the pcLasso solution, we typically hold  $\theta$  fixed and compute the solution for a path of  $\lambda$  values. By casting the pcLasso optimization problem as a lasso problem with a different response and design matrix, we can derive strong rules for pcLasso from that for the lasso.

We first derive strong rules for pcLasso where predictors do not have preassigned groups, i.e. the solution to (5). If we define

$$\tilde{\mathbf{y}} = \begin{pmatrix} \mathbf{y} \\ 0 \end{pmatrix}, \quad \tilde{\mathbf{X}} = \begin{pmatrix} \mathbf{X} \\ \sqrt{\theta} \mathbf{A}^{1/2} \end{pmatrix}, \quad (33)$$

where  $\mathbf{A} = \mathbf{V} \mathbf{D}_{d_1^2 - d_2^2} \mathbf{V}^T$ , then pcLasso is equivalent to the lasso with response vector  $\tilde{\mathbf{y}}$  and design matrix  $\tilde{\mathbf{X}}$ . Fix  $\theta$ , fix a path of  $\lambda$  values  $\lambda_1 \geq \lambda_2 \geq \dots$ , and let  $\hat{\beta}(\lambda_i)$  denote the pcLasso solution at  $\lambda = \lambda_i$ . Applying the lasso strong rules to this set-up, the sequential strong rule for pcLasso for discarding predictor  $j$  at  $\lambda = \lambda_i$  is

$$\begin{aligned} & |\mathbf{X}_j^T (\mathbf{y} - \mathbf{X} \hat{\beta}(\lambda_{i-1})) - \theta (\mathbf{A}^{1/2})_j \mathbf{A}^{1/2} \hat{\beta}(\lambda_{i-1})| < 2\lambda_i - \lambda_{i-1}, \\ \Leftrightarrow & |\mathbf{X}_j^T (\mathbf{y} - \mathbf{X} \hat{\beta}(\lambda_{i-1})) - \theta (\mathbf{A} \hat{\beta}(\lambda_{i-1}))_j| < 2\lambda_i - \lambda_{i-1}. \end{aligned} \quad (34)$$

Next we derive strong rules for pcLasso where predictors come in preassigned non-overlapping groups, i.e. the solution to (8). If we define  $\mathbf{A}_k = \mathbf{V}_k \mathbf{D}_{d_{k1}^2 - d_{kj}^2} \mathbf{V}_k^T$  for  $k = 1, \dots, K$ ,  $\mathbf{A}$  as the block diagonal matrix

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_1 & & 0 \\ & \ddots & \\ 0 & & \mathbf{A}_K \end{pmatrix},$$

and  $\tilde{\mathbf{y}}$  and  $\tilde{\mathbf{X}}$  as in (33), then pcLasso is again equivalent to the lasso with  $\tilde{\mathbf{y}}$  and  $\tilde{\mathbf{X}}$ . This results in the same sequential strong rule as in the single group case (34), although the matrix  $\mathbf{A}$  is defined differently.

## C Proofs for Section 8

Before presenting proofs for Section 8, we first prove two technical lemmas which show that the pcLasso solution lies in the constraint set which we will use for the restricted eigenvalue condition (23).

**Lemma 10.** *If we solve the constrained form of pcLasso (21) with  $R = \|\beta^*\|_1$ , then  $\|\hat{\nu}_{S^c}\|_1 \leq \|\hat{\nu}_S\|_1$ .*

*Proof.* Since  $\hat{\beta}$  is feasible for (21), by the triangle inequality,

$$\|\beta_S^*\|_1 = R \geq \|\beta^* + \hat{\nu}\|_1 = \|\beta_S^* + \hat{\nu}_S\|_1 + \|\hat{\nu}_{S^c}\|_1 \geq \|\beta_S^*\|_1 - \|\hat{\nu}_S\|_1 + \|\hat{\nu}_{S^c}\|_1.$$

□

**Lemma 11.** *If we solve the Lagrangian form of pcLasso (22) with  $\lambda \geq \frac{2}{n} \|\mathbf{X}^T \mathbf{w} - n\theta \mathbf{A} \beta^*\|_\infty$ , then  $\|\hat{\nu}_{S^c}\|_1 \leq 3 \|\hat{\nu}_S\|_1$ .*

*Proof.* Define  $G(\nu) = \frac{1}{2n} \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}(\beta^* + \nu)\|_2^2 + \lambda \|\beta^* + \nu\|_1$ . By definition,  $G(\hat{\nu}) \leq G(0)$ , so

$$\begin{aligned} \frac{1}{2n} \|\tilde{\mathbf{w}} - \tilde{\mathbf{X}}\hat{\nu}\|_2^2 + \lambda \|\beta^* + \hat{\nu}\|_1 &\leq \frac{1}{2n} \|\tilde{\mathbf{w}}\|_2^2 + \lambda \|\beta^*\|_1, \\ \frac{1}{2n} \hat{\nu}^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \hat{\nu} &\leq \frac{\tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \hat{\nu}}{n} + \lambda (\|\beta^*\|_1 - \|\beta^* + \hat{\nu}\|_1). \end{aligned} \quad (35)$$

Note that  $\|\beta^* + \hat{\nu}\|_1 = \|\beta_S^* + \hat{\nu}_S\|_1 + \|\hat{\nu}_{S^c}\|_1 \geq \|\beta_S^*\|_1 - \|\hat{\nu}_S\|_1 + \|\hat{\nu}_{S^c}\|_1$ ; substituting this in the above gives

$$\frac{1}{2n} \hat{\nu}^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \hat{\nu} \leq \frac{\tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \hat{\nu}}{n} + \lambda (\|\hat{\nu}_S\|_1 - \|\hat{\nu}_{S^c}\|_1). \quad (36)$$

By Hölder's inequality and our choice of  $\lambda$ , we have

$$\frac{\tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \hat{\nu}}{n} \leq \frac{1}{n} \|\tilde{\mathbf{X}}^T \tilde{\mathbf{w}}\|_\infty \|\hat{\nu}\|_1 = \frac{1}{n} \|\mathbf{X}^T \mathbf{w} - n\theta \mathbf{A} \beta^*\|_\infty \|\hat{\nu}\|_1 \leq \frac{1}{2} \lambda \|\hat{\nu}\|_1. \quad (37)$$

Substituting this inequality into (36) and noting that the LHS of (36) is always non-negative, we obtain the desired conclusion.

□

### C.1 Proof of Lemma 2

Let  $\nu \in \mathcal{C}$  be expressed as  $\nu = \sum_{j=1}^p a_j v_j$ , where the  $v_j$ 's are the columns of  $\mathbf{V}$ . Direct computation gives

$$\|\nu\|_2^2 = \sum_{j=1}^p a_j^2, \quad \nu^T \mathbf{X}^T \mathbf{X} \nu = \sum_{j=1}^p d_j^2 a_j^2, \quad \nu^T (n\theta \mathbf{A}) \nu = \sum_{j=2}^p n\theta (d_1^2 - d_j^2) a_j^2. \quad (38)$$

Thus, for any  $\alpha \in [0, 1]$ ,

$$\begin{aligned} \nu^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \nu &= \nu^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \nu \\ &= \sum_{j=1}^p d_j^2 a_j^2 + \sum_{j=2}^p n\theta (d_1^2 - d_j^2) a_j^2 \\ &= \alpha \sum_{j=1}^p d_j^2 a_j^2 + (1 - \alpha) d_1^2 a_1^2 + \sum_{j=2}^p [n\theta d_1^2 + (1 - \alpha - n\theta) d_j^2] a_j^2 \\ &\geq \alpha n\gamma \sum_{j=1}^p a_j^2 + (1 - \alpha) d_1^2 a_1^2 + \min [n\theta d_1^2, n\theta d_1^2 + (1 - \alpha - n\theta) d_2^2] \sum_{j=2}^p a_j^2, \\ \nu^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \nu &\geq [\alpha n\gamma + \min ((1 - \alpha) d_1^2, n\theta d_1^2, n\theta d_1^2 + (1 - \alpha - n\theta) d_2^2)] \|\nu\|_2^2. \end{aligned} \quad (39)$$

If  $n\theta \leq 1$ , setting  $\alpha = 1 - n\theta$  in (39) gives

$$\begin{aligned} \nu^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \nu &\geq [(1 - n\theta) n\gamma + \min (n\theta d_1^2, n\theta d_1^2, n\theta d_1^2)] \|\nu\|_2^2 \\ &= [(1 - n\theta) n\gamma + n\theta d_1^2] \|\nu\|_2^2. \end{aligned}$$

If  $n\theta \geq 1$ , setting  $\alpha = 0$  in (39) gives

$$\begin{aligned} \nu^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \nu &\geq \min (d_1^2, n\theta (d_1^2 - d_2^2) + d_2^2) \|\nu\|_2^2 \\ &\geq \min (d_1^2, (d_1^2 - d_2^2) + d_2^2) \|\nu\|_2^2 \\ &= d_1^2 \|\nu\|_2^2. \end{aligned}$$

### C.2 Proof of Lemma 3

Let  $\nu \in \mathbb{R}^p$  be expressed as  $\nu = \sum_{j=1}^p a_j v_j$ , where the  $v_j$ 's are the columns of  $\mathbf{V}$ . Using the formulas in (38),

$$\begin{aligned} \nu^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \nu &= \nu^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \nu \\ &= \sum_{j=1}^p d_j^2 a_j^2 + \sum_{j=2}^p n\theta (d_1^2 - d_j^2) a_j^2 \\ &= \sum_{j=1}^p [n\theta d_1^2 + (1 - n\theta) d_j^2] a_j^2 \\ &\geq \min_j [n\theta d_1^2 + (1 - n\theta) d_j^2] \sum_{j=1}^p a_j^2 \\ &\begin{cases} \geq n\theta d_1^2 \|\nu\|_2^2 & \text{if } n\theta \leq 1, \\ \geq d_1^2 \|\nu\|_2^2 & \text{if } n\theta \geq 1 \end{cases} \\ &= \min(n\theta, 1) d_1^2 \|\nu\|_2^2. \end{aligned}$$

### C.3 Proof of Theorem 4

Let  $\hat{\nu} = \hat{\beta} - \beta^*$ . By definition of  $\hat{\beta}$  and using the relation  $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\beta^* + \tilde{\mathbf{w}}$ ,

$$\begin{aligned} \left\| \tilde{\mathbf{y}} - \tilde{\mathbf{X}}\hat{\beta} \right\|_2^2 &\leq \left\| \tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta^* \right\|_2^2, \\ \left\| \tilde{\mathbf{w}} - \tilde{\mathbf{X}}\hat{\nu} \right\|_2^2 &\leq \|\tilde{\mathbf{w}}\|_2^2, \\ \hat{\nu}^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \hat{\nu} &\leq 2(\tilde{\mathbf{X}}^T \tilde{\mathbf{w}})^T \hat{\nu}. \end{aligned} \quad (40)$$

By Lemma 2, the LHS is  $\geq n\eta \|\hat{\nu}\|_2^2$ , where  $n\eta = \min[(1-n\theta)n\gamma + n\theta d_1^2, d_1^2]$ . We can get an upper bound on the RHS:

$$\begin{aligned} (\tilde{\mathbf{X}}^T \tilde{\mathbf{w}})^T \hat{\nu} &\leq \left\| \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} \right\|_\infty \|\hat{\nu}\|_1 && \text{(Hölder's inequality)} \\ &= \left\| \mathbf{X}^T \mathbf{w} - n\theta \mathbf{A}\beta^* \right\|_\infty (\|\hat{\nu}_S\|_1 + \|\hat{\nu}_{S^c}\|_1) \\ &\leq 2 \left\| \mathbf{X}^T \mathbf{w} - n\theta \mathbf{A}\beta^* \right\|_\infty \|\hat{\nu}_S\|_1 && \text{(Lemma 10)} \\ &\leq 2 \left\| \mathbf{X}^T \mathbf{w} - n\theta \mathbf{A}\beta^* \right\|_\infty \sqrt{|S|} \|\hat{\nu}\|_2. && \text{(Cauchy-Schwarz)} \end{aligned}$$

Hence,

$$\begin{aligned} n\eta \|\hat{\nu}\|_2^2 &\leq 4 \left\| \mathbf{X}^T \mathbf{w} - n\theta \mathbf{A}\beta^* \right\|_\infty \sqrt{|S|} \|\hat{\nu}\|_2, \\ \|\hat{\nu}\|_2 &\leq \frac{4\sqrt{|S|} \left\| \mathbf{X}^T \mathbf{w} - n\theta \mathbf{A}\beta^* \right\|_\infty}{n\eta}. \end{aligned}$$

The second inequality in (26) is the result of using the triangle inequality on the  $\|\cdot\|_\infty$  norm.

When  $\beta^*$  is aligned with the first principal component of  $\mathbf{X}$ ,  $\mathbf{A}\beta^* = 0$ . Hence, the first inequality in (26) reduces to (27).

### C.4 Proof of Theorem 5

The proof of Theorem 5 is exactly the same as that for Theorem 4, except that we bound the LHS of (40) from below by  $\min(n\theta, 1) d_1^2 \|\hat{\nu}\|_2^2$ . We can do this due to Lemma 3.

### C.5 Proof of Theorem 6

As in the proof of Lemma 11, we have (36):

$$\frac{1}{2n} \hat{\nu}^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \hat{\nu} \leq \frac{\tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \hat{\nu}}{n} + \lambda(\|\hat{\nu}_S\|_1 - \|\hat{\nu}_{S^c}\|_1).$$

By Lemma 2, the LHS is  $\geq \frac{\eta}{2} \|\hat{\nu}\|_2^2$ , where  $n\eta = \min[(1-n\theta)n\gamma + n\theta d_1^2, d_1^2]$ . By our choice of  $\lambda$ , we may use (37) to obtain

$$\begin{aligned} \frac{\eta}{2} \|\hat{\nu}\|_2^2 &\leq \frac{1}{2} \lambda \|\hat{\nu}\|_1 + \lambda(\|\hat{\nu}_S\|_1 - \|\hat{\nu}_{S^c}\|_1) \\ &\leq \frac{3\lambda}{2} \|\hat{\nu}_S\|_1 \leq \frac{3\lambda}{2} \sqrt{|S|} \|\hat{\nu}\|_2, \\ \|\hat{\nu}\|_2 &\leq \frac{3\lambda \sqrt{|S|}}{\eta}, \end{aligned}$$

as required.



## C.6 Proof of Theorem 7

By (35),

$$\begin{aligned}
0 &\leq \frac{1}{2n} \widehat{\nu}^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \widehat{\nu} \leq \frac{\widetilde{\mathbf{w}}^T \widetilde{\mathbf{X}} \widehat{\nu}}{n} + \lambda (\|\beta^*\|_1 - \|\beta^* + \widehat{\nu}\|_1) \\
&\leq \frac{1}{n} \left\| \widetilde{\mathbf{X}}^T \widetilde{\mathbf{w}} \right\|_{\infty} \|\widehat{\nu}\|_1 + \lambda (\|\beta^*\|_1 - \|\beta^* + \widehat{\nu}\|_1) \quad (\text{H\"older's inequality}) \\
&\leq \left( \frac{1}{n} \left\| \mathbf{X}^T \mathbf{w} - n\theta \mathbf{A} \beta^* \right\|_{\infty} - \lambda \right) \|\widehat{\nu}\|_1 + 2\lambda \|\beta^*\|_1 \quad (\text{triangle inequality}) \\
&\leq -\frac{\lambda}{2} \|\widehat{\nu}\|_1 + 2\lambda \|\beta^*\|_1, \quad (\text{by choice of } \lambda),
\end{aligned}$$

so  $\|\widehat{\nu}\|_1 \leq 4 \|\beta^*\|_1$ . Using (35) again,

$$\begin{aligned}
\frac{1}{2n} \left\| \mathbf{X}(\widehat{\beta} - \beta^*) \right\|_2^2 &= \frac{1}{2n} \widehat{\nu}^T \mathbf{X}^T \mathbf{X} \widehat{\nu} \leq \frac{1}{2n} \widehat{\nu}^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \widehat{\nu} \\
&\leq \frac{1}{n} \left\| \widetilde{\mathbf{X}}^T \widetilde{\mathbf{w}} \right\|_{\infty} \|\widehat{\nu}\|_1 + \lambda \|\widehat{\nu}\|_1 \\
&\leq \frac{3\lambda}{2} \|\widehat{\nu}\|_1 \\
&\leq 6\lambda \|\beta^*\|_1.
\end{aligned}$$

Multiplying both sides by 2 establishes the claim.

## C.7 Proof of Theorem 8

Let  $\kappa = \min(n\theta, 1)d_1^2$ . Mimicking the first part of the proof of Theorem 7 and using Lemma 3, we have

$$\begin{aligned}
-\frac{\lambda}{2} \|\widehat{\nu}\|_1 + 2\lambda \|\beta^*\|_1 &\geq \frac{1}{2n} \widehat{\nu}^T (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}}) \widehat{\nu} \geq \kappa \|\widehat{\nu}\|_2^2 \\
&\geq \kappa \frac{\|\widehat{\nu}\|_1^2}{p}, \quad (\text{Cauchy-Schwarz}), \\
\frac{\kappa}{p} \|\widehat{\nu}\|_1^2 + \frac{\lambda}{2} \|\widehat{\nu}\|_1 - 2\lambda \|\beta^*\|_1 &\leq 0, \\
\|\widehat{\nu}\|_1 &\leq \frac{-\lambda/2 + \sqrt{\lambda^2/4 + 8\lambda \|\beta^*\|_1 \kappa/p}}{2\kappa/p} \\
&= \frac{-\lambda p + \sqrt{\lambda^2 p^2 + 32\lambda \|\beta^*\|_1 \kappa p}}{4\kappa}.
\end{aligned}$$

Thus, mimicking the second part of the proof of Theorem 7,

$$\begin{aligned}
\frac{1}{2n} \left\| \mathbf{X}(\widehat{\beta} - \beta^*) \right\|_2^2 &\leq \frac{3\lambda}{2} \|\widehat{\nu}\|_1 \\
&\leq \frac{3\lambda - \lambda p + \sqrt{\lambda^2 p^2 + 32\lambda \|\beta^*\|_1 \kappa p}}{4\kappa}, \\
\frac{1}{n} \left\| \mathbf{X}(\widehat{\beta} - \beta^*) \right\|_2^2 &\leq \frac{3\lambda(-\lambda p + \sqrt{\lambda^2 p^2 + 32\lambda \|\beta^*\|_1 \kappa p})}{4\kappa},
\end{aligned}$$

as required. To show that the RHS above is indeed a better bound than that in Theorem 7:

$$\begin{aligned}
12\lambda \|\beta^*\|_1 &\geq \frac{3\lambda(-\lambda p + \sqrt{\lambda^2 p^2 + 32\lambda \|\beta^*\|_1 \kappa p})}{4\kappa} \\
\Leftrightarrow 16\kappa \|\beta^*\|_1 &\geq -\lambda p + \sqrt{\lambda^2 p^2 + 32\lambda \|\beta^*\|_1 \kappa p} \\
\Leftrightarrow (\lambda p + 16\kappa \|\beta^*\|_1)^2 &\geq \lambda^2 p^2 + 32\lambda \|\beta^*\|_1 \kappa p \\
\Leftrightarrow 256\kappa^2 \|\beta^*\|_1^2 &\geq 0,
\end{aligned}$$

which is obviously true.

## C.8 Proof of Theorem 9

We have

$$\begin{aligned}
\frac{1}{2n} \hat{\nu}^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \hat{\nu} &\leq \frac{\tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \hat{\nu}}{n} + \lambda (\|\hat{\nu}_S\|_1 - \|\hat{\nu}_{S^c}\|_1) && \text{((36) in Lemma 11)} \\
&\leq \frac{1}{n} \left\| \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} \right\|_{\infty} \|\hat{\nu}\|_1 + \lambda \|\hat{\nu}\|_1 && \text{(Hölder's inequality)} \\
&\leq \frac{3\lambda}{2} \|\hat{\nu}\|_1 && \text{(choice of } \lambda) \\
&\leq \frac{3\lambda}{2} \sqrt{|S|} \|\hat{\nu}\|_2. && \text{(Cauchy-Schwarz inequality)}
\end{aligned}$$

By Lemma 11,  $\hat{\nu}$  lies in the set  $\{\nu \in \mathbb{R}^p : \|\nu_{S^c}\|_1 \leq 3\|\nu_S\|_1\}$ , and so the restricted eigenvalue condition and Lemma 2 apply, i.e.  $\|\hat{\nu}\|_2^2 \leq \frac{1}{n\eta} \left\| \tilde{\mathbf{X}} \hat{\nu} \right\|_2^2$ , where  $n\eta = \min [(1 - n\theta)n\gamma + n\theta d_1^2, d_1^2]$ . Hence,

$$\begin{aligned}
\frac{1}{2n} \hat{\nu}^T (\mathbf{X}^T \mathbf{X} + n\theta \mathbf{A}) \hat{\nu} &\leq \frac{3\lambda}{2} \sqrt{|S|} \|\hat{\nu}\|_2 \leq \frac{3\lambda}{2} \sqrt{|S|} \sqrt{\frac{1}{n\eta}} \left\| \tilde{\mathbf{X}} \hat{\nu} \right\|_2, \\
\|\mathbf{X} \hat{\nu}\|_2 &\leq \left\| \tilde{\mathbf{X}} \hat{\nu} \right\|_2 \leq 3n\lambda \sqrt{|S|} \sqrt{\frac{1}{n\eta}}, \\
\frac{\|\mathbf{X} \hat{\nu}\|_2^2}{n} &\leq \frac{9|S|\lambda^2}{\eta},
\end{aligned}$$

as required.

## D Full details of simulation study in Section 7

In this appendix, we present details on how data was generated for the simulation study, as well as results that we obtained for a variety of settings.

### D.1 Data generation details

Let  $n$  be the size of the training data,  $p$  be the number of features for each observation, and  $K$  be the number of groups of features. Let  $size$  be a vector of length  $K$  such that its  $k$ th element denotes the number of features in group  $k$ .

### D.1.1 Generating $\mathbf{X}$

Let  $\mathbf{X}_k$  denote the design matrix for group  $k$ , and let  $\mathbf{X}$  denote the full design matrix.

1. Let  $\rho$  denote the correlation between features in each group. For each group  $k$ , build the covariance matrix  $\Sigma_k = \rho \mathbf{1} + (1 - \rho) \mathbf{I} \in \mathbb{R}^{size_k \times size_k}$ .
2. Generate the rows of  $\mathbf{X}_k$  independently from  $\mathcal{N}(0, \Sigma_k)$ .

### D.1.2 Generating $\mathbf{y}$

1. Perform singular value decomposition (SVD) on each  $\mathbf{X}_k$  to get  $\mathbf{U}_k$ ,  $\mathbf{D}_k$  and  $\mathbf{V}_k$ .
2. Let  $n_{ev}$  be the number of principal components (PCs) to take from each group, and let  $SNR$  be the desired signal-to-noise ratio (SNR).
3. We consider three settings for determining which PCs will form the signal:
  - “*Home court*”: Set  $\mathbf{W}_k$  to be the first  $n_{ev}$  columns of  $\mathbf{V}_k$ .
  - “*Neutral court*”: Set  $\mathbf{W}_k$  to be  $n_{ev}$  random columns of  $\mathbf{V}_k$ .
  - “*Hostile court*”: Set  $\mathbf{W}_k$  to be the bottom  $n_{ev}$  columns of  $\mathbf{V}_k$ .
4. For each  $k$ , specify coefficients  $\mathbf{b}_k \in \mathbb{R}^{n_{ev}}$ . For simplicity, we set  $\mathbf{b}_k = \mathbf{2}$  if the group is related to the response,  $\mathbf{b}_k = \mathbf{0}$  otherwise.
5. Set  $signal = \sum_k \mathbf{X}_k \mathbf{W}_k \mathbf{b}_k$ .
6. Set response  $\mathbf{y} = signal + \epsilon$ , where the entries of  $\epsilon$  are iid  $\mathcal{N}(0, V)$ , where

$$V = \frac{\text{Var}(signal)}{SNR} = \frac{\sum_k \mathbf{b}_k^T \mathbf{W}_k^T \Sigma_k \mathbf{W}_k \mathbf{b}_k}{SNR}.$$

### D.1.3 Generating test data

1. For each  $k$ , generate  $\mathbf{X}_{test,k}$  such that each row is i.i.d.  $\mathcal{N}(\mathbf{0}, \Sigma_k)$ . Generate the  $\mathbf{X}_{test,k}$  independently of each other.
2. Compute  $signal_{test} = \sum_k \mathbf{X}_{test,k} \mathbf{W}_k \mathbf{b}_k$ .

## D.2 Simulation details

We ran simulations for the following settings:

- $n = 200$ ,  $p = 50$ , no groups,  $n_{ev} = 5$ : home court, neutral court and hostile court.
- $n = 200$ ,  $p = 50$ , 5 groups of 10 predictors each,  $n_{ev} = 1$ , signal only depends on the first two groups: home court, neutral court and hostile court.
- $n = 200$ ,  $p = 200$ , 10 groups of 20 predictors each,  $n_{ev} = 2$ , signal only depends on the first group: home court and neutral court.
- $n = 200$ ,  $p = 1,000$ , 10 groups of 100 predictors each,  $n_{ev} = 2$ , signal only depends on the first group: home court and neutral court.

For each setting above, we ran 30 simulations for each combination of within-group correlation  $\rho \in \{0, 0.3\}$  and SNR  $SNR \in \{0.5, 1, 2\}$ .

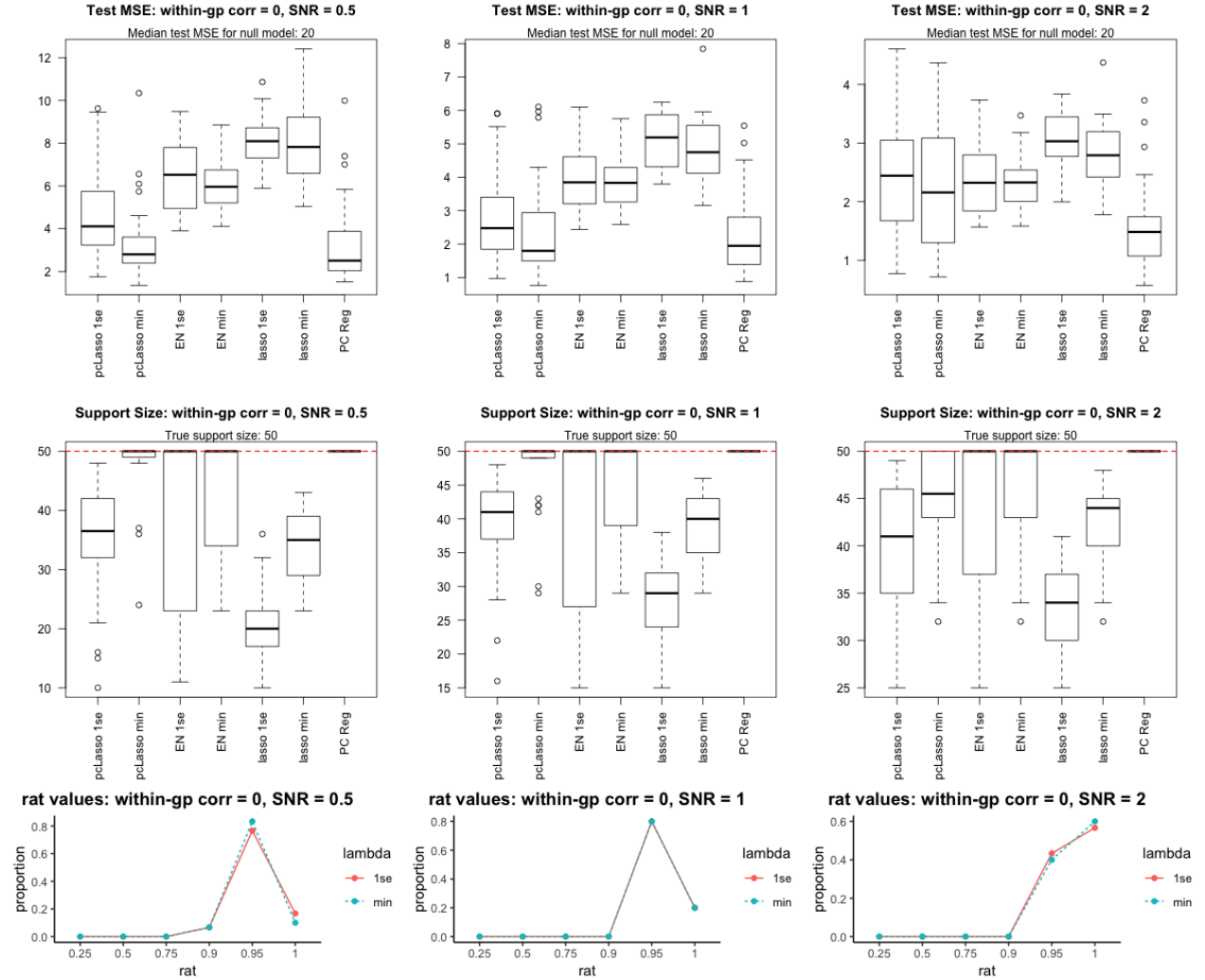
We compared the following models:

- The null model, i.e. the mean of the responses in the training data.
- pcLasso with CV across  $rat = 0.25, 0.5, 0.75, 0.9, 0.95, 1$ , with  $\lambda$  values selected both by `lambda.min` and `lambda.1se`.
- Elastic net with CV across  $\alpha = 0, 0.2, 0.4, 0.6, 0.8$  and  $1$ , with  $\lambda$  values selected both by `lambda.min` and `lambda.1se`.
- Lasso with CV, with  $\lambda$  values selected both by `lambda.min` and `lambda.1se`.
- Principal components (PC) regression with CV across ranks.

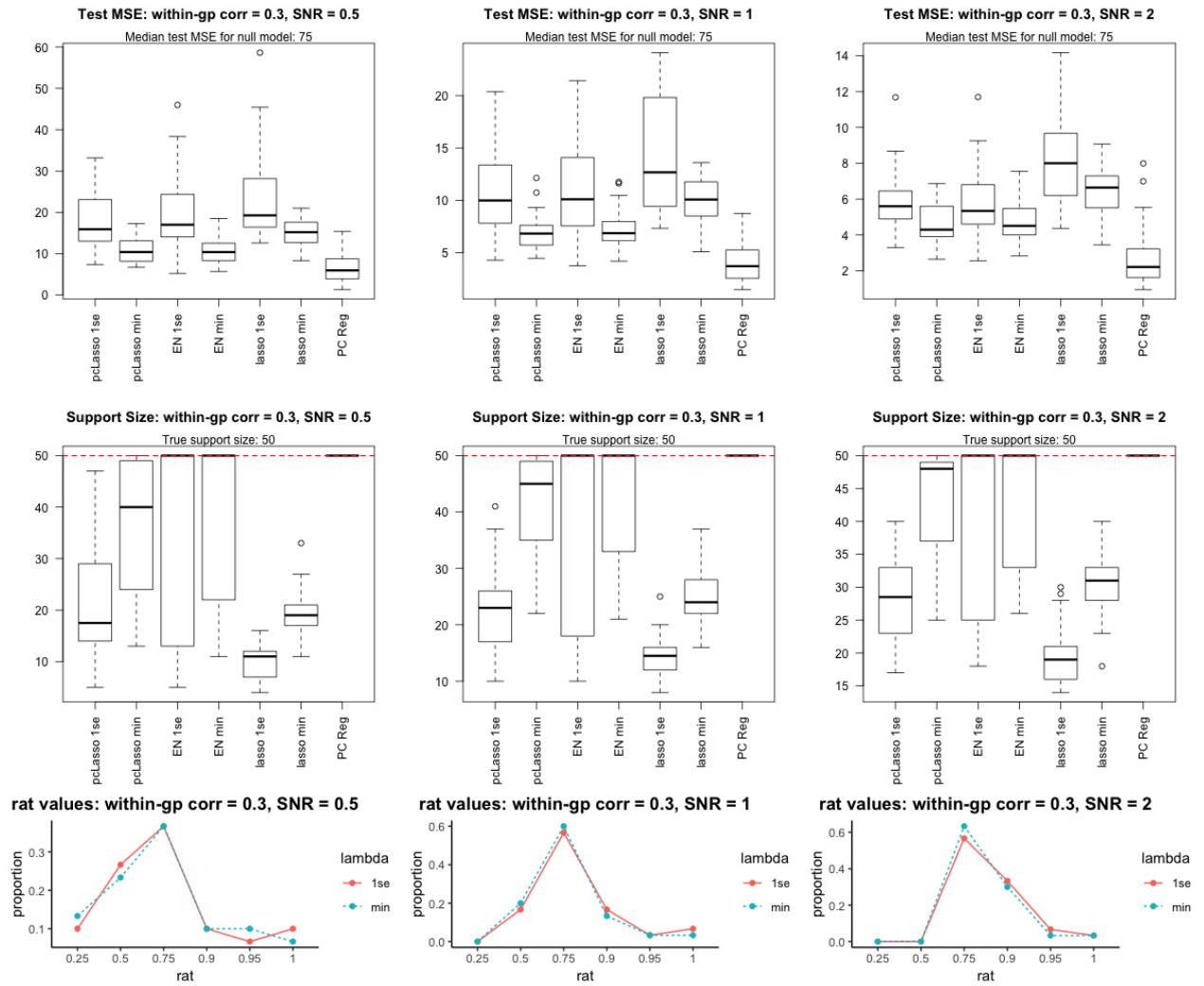
For each method, we present a boxplot of the test mean-squared error (MSE) on 5,000 test points, i.e.  $MSE = \mathbb{E}[(\hat{y}_{test} - signal_{test})^2]$ . We also present boxplots of the support size of the model which the method selects. Finally, we present a line histogram of the values of  $rat$  which the two versions of pcLasso select during model fitting.

### D.2.1 Small $p$ no groups, “home court” for pcLasso

$n = 200, p = 50$ , 1 group of uncorrelated predictors, response related to top 5 eigenvectors.

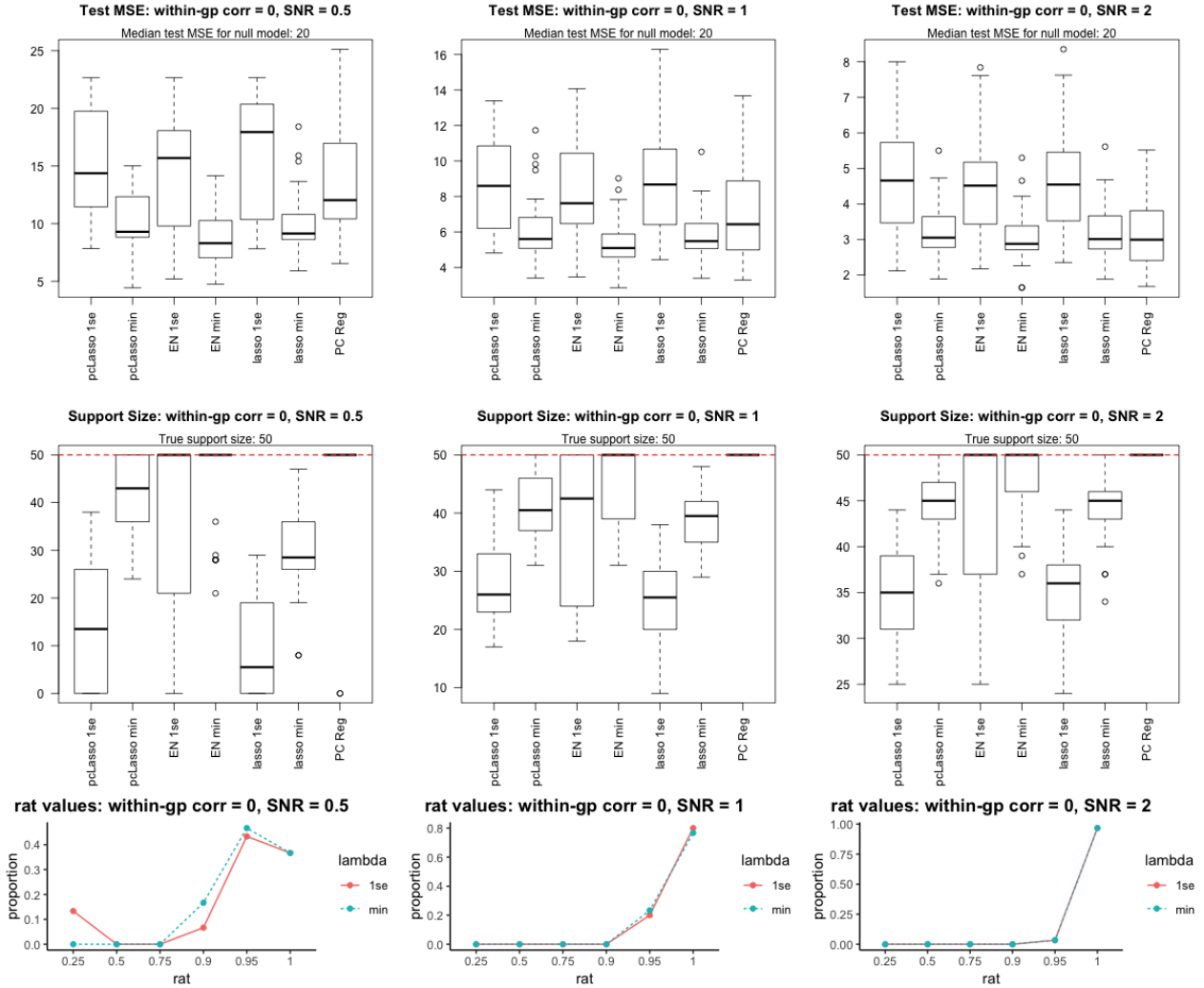


As above, but with predictors having pairwise correlation of 0.3.

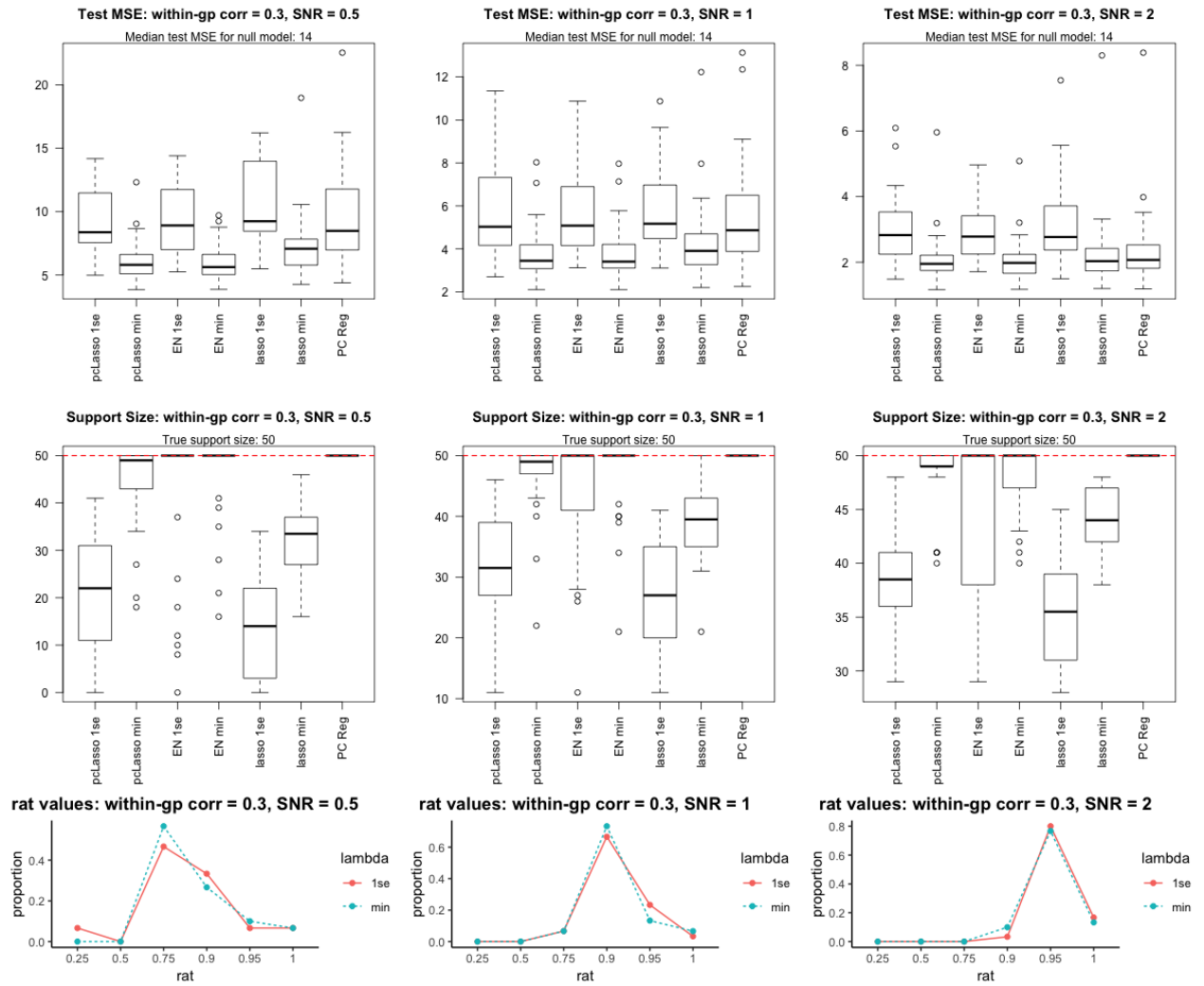


## D.2.2 Small $p$ no groups, “neutral court” for pcLasso

$n = 200$ ,  $p = 50$ , 1 group of uncorrelated predictors, response related to 5 random eigenvectors.



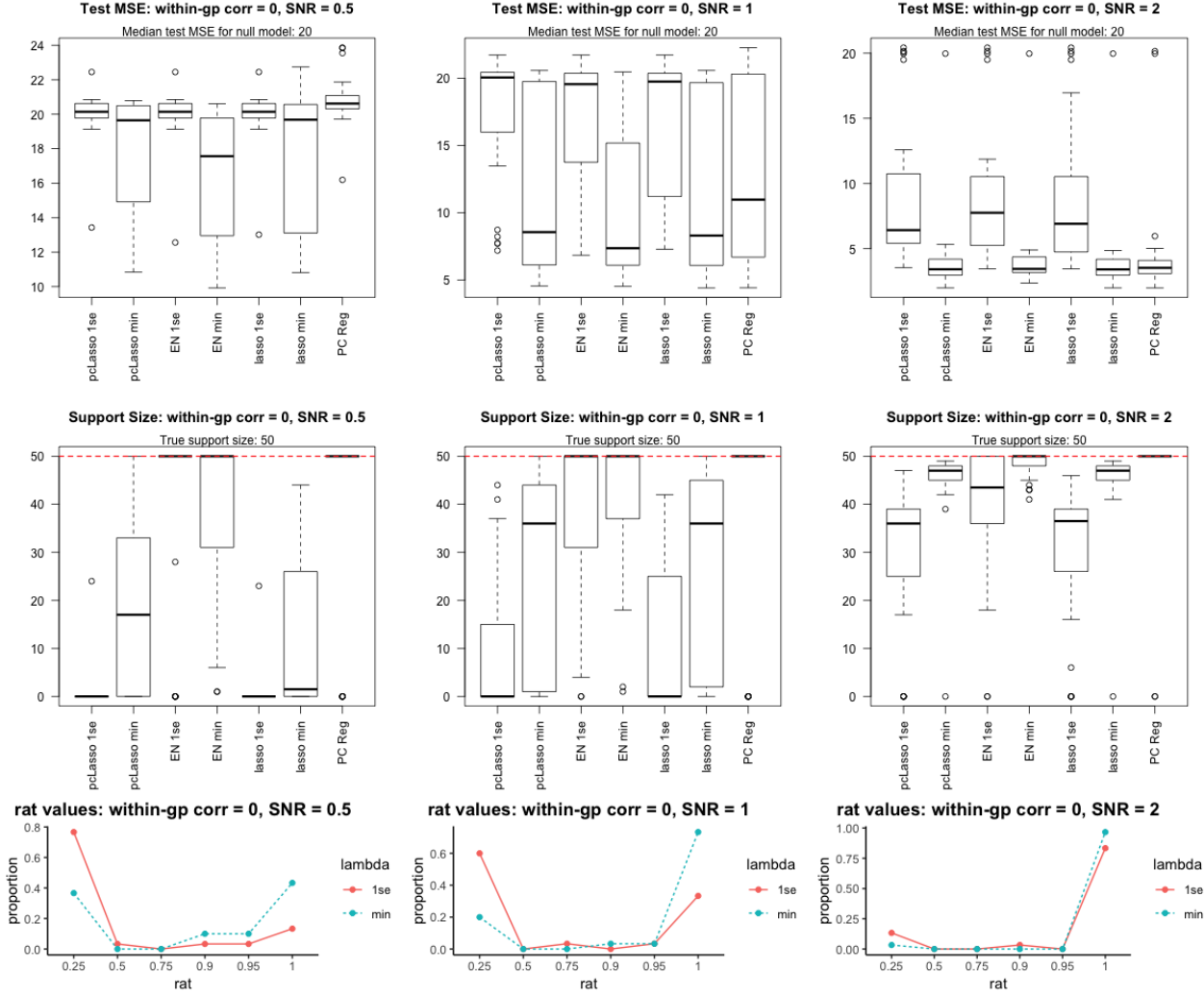
As above, but with predictors having pairwise correlation of 0.3.



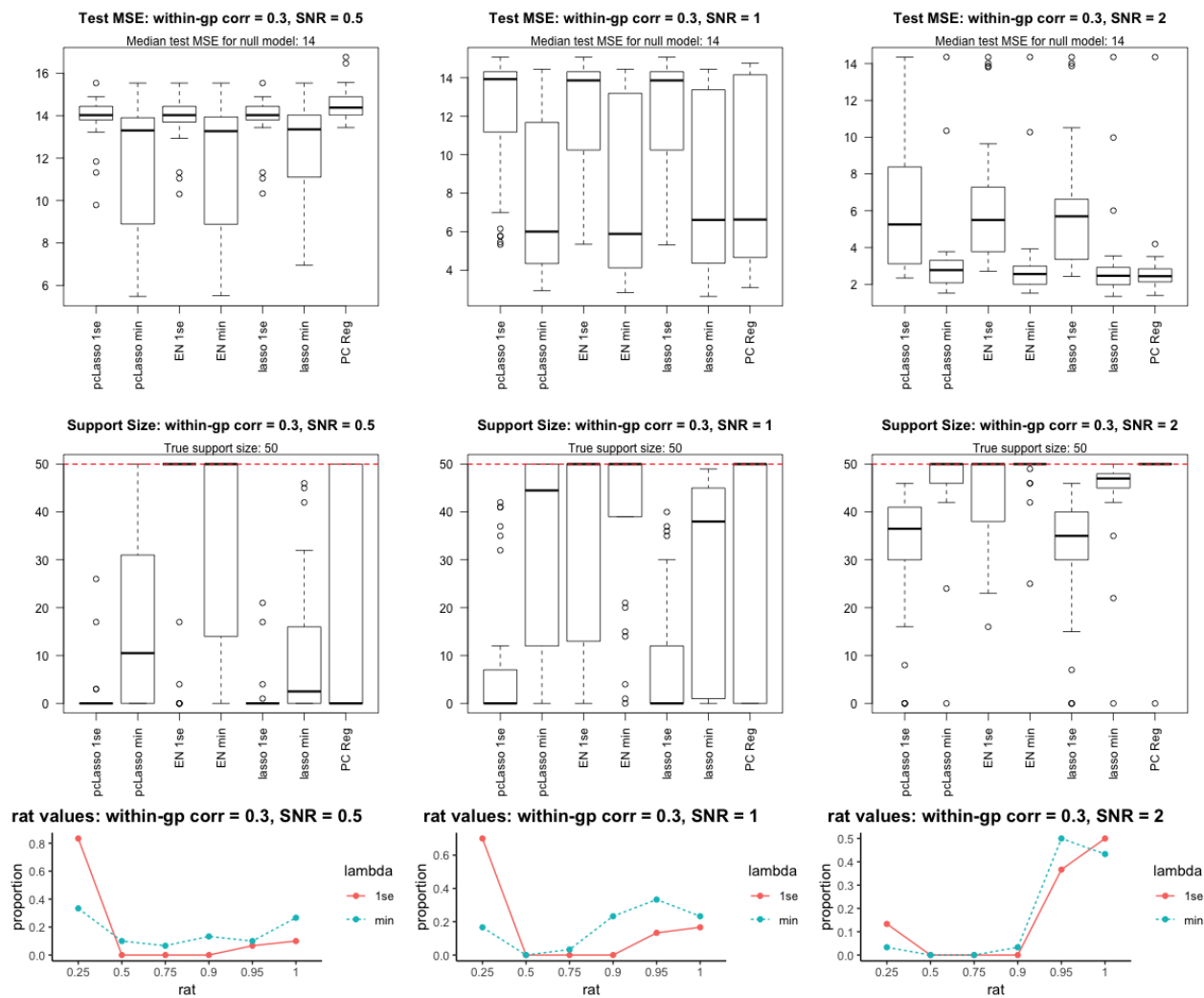


### D.2.3 Small $p$ no groups, “hostile court” for pcLasso

$n = 200, p = 50$ , 1 group of uncorrelated predictors, response related to the bottom 5 eigenvectors.

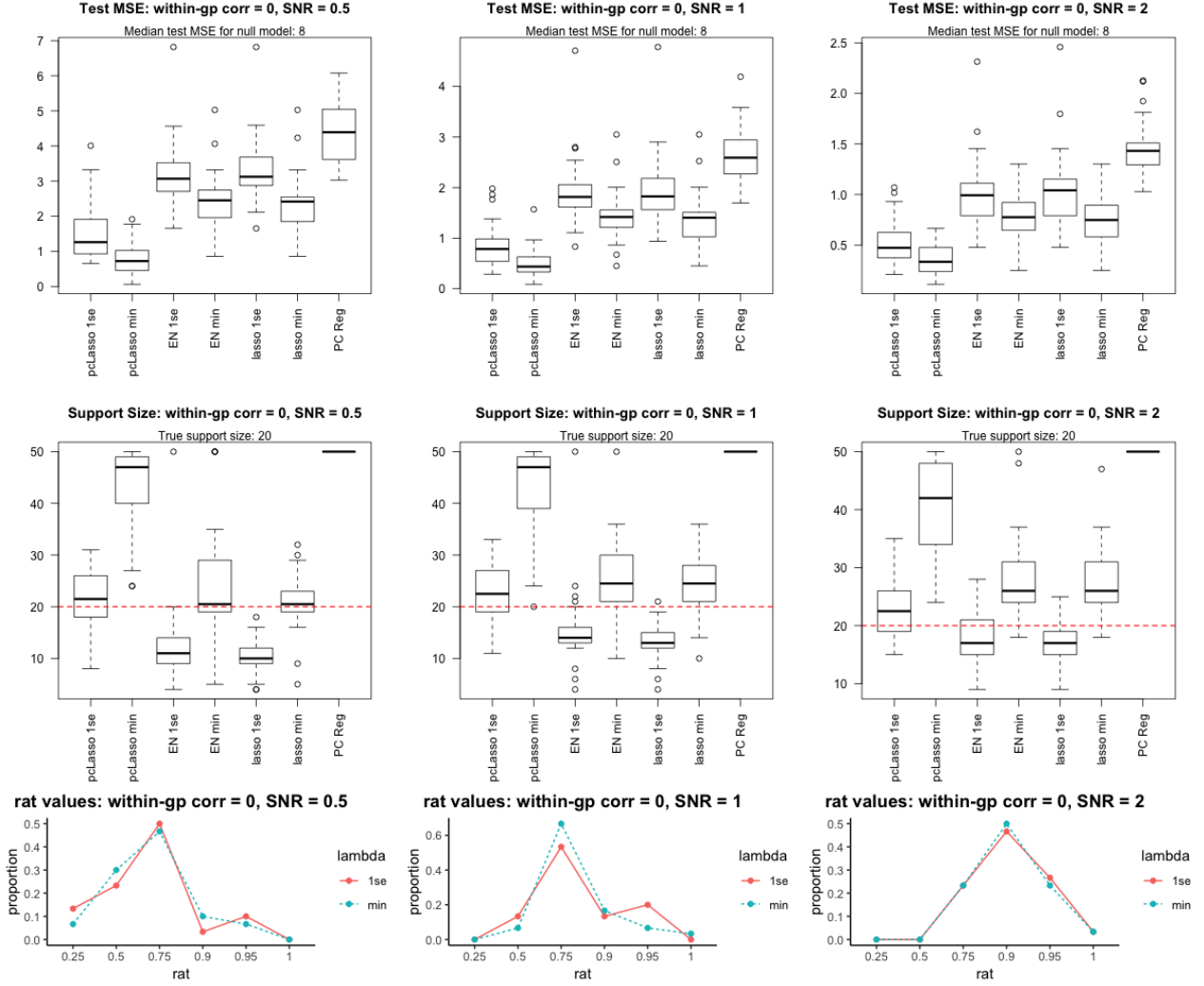


As above, but with predictors having pairwise correlation of 0.3.

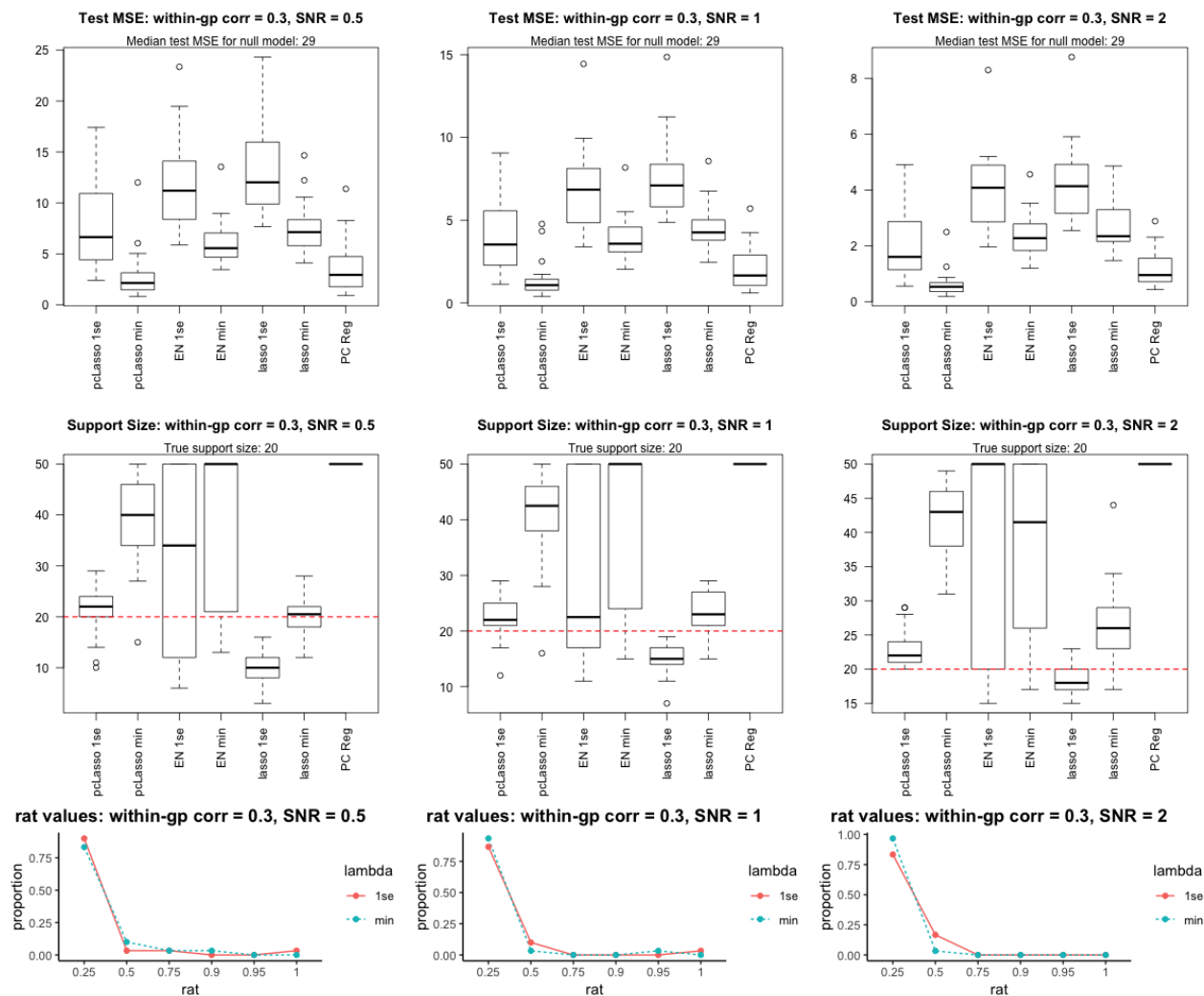


### D.2.4 Small $p$ with groups, “home court” for pcLasso

$n = 200$ ,  $p = 50$ , 5 groups of 10 uncorrelated predictors, response related to the top eigenvector in the first 2 groups.

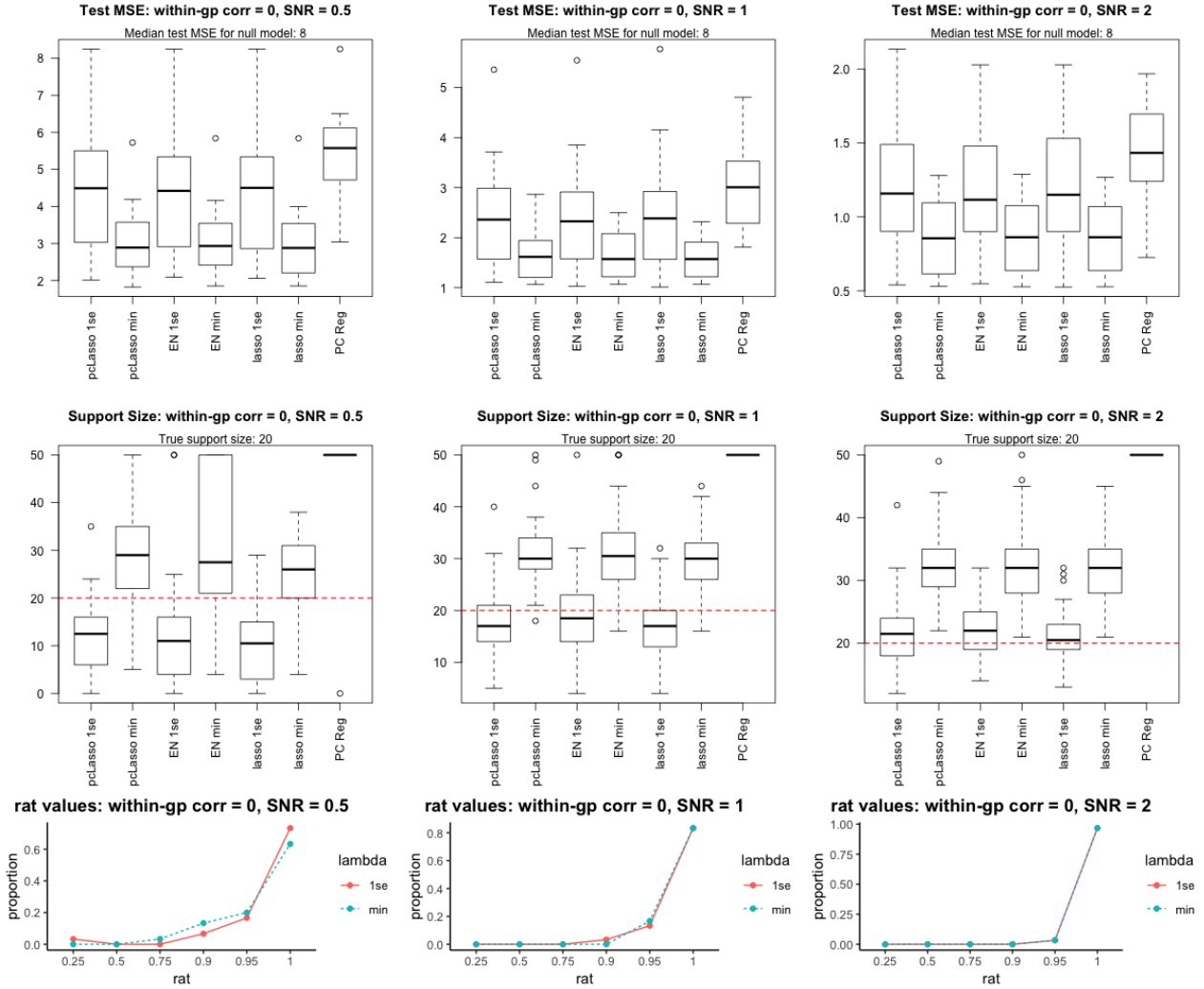


As above, but with predictors in the same group having pairwise correlation of 0.3.

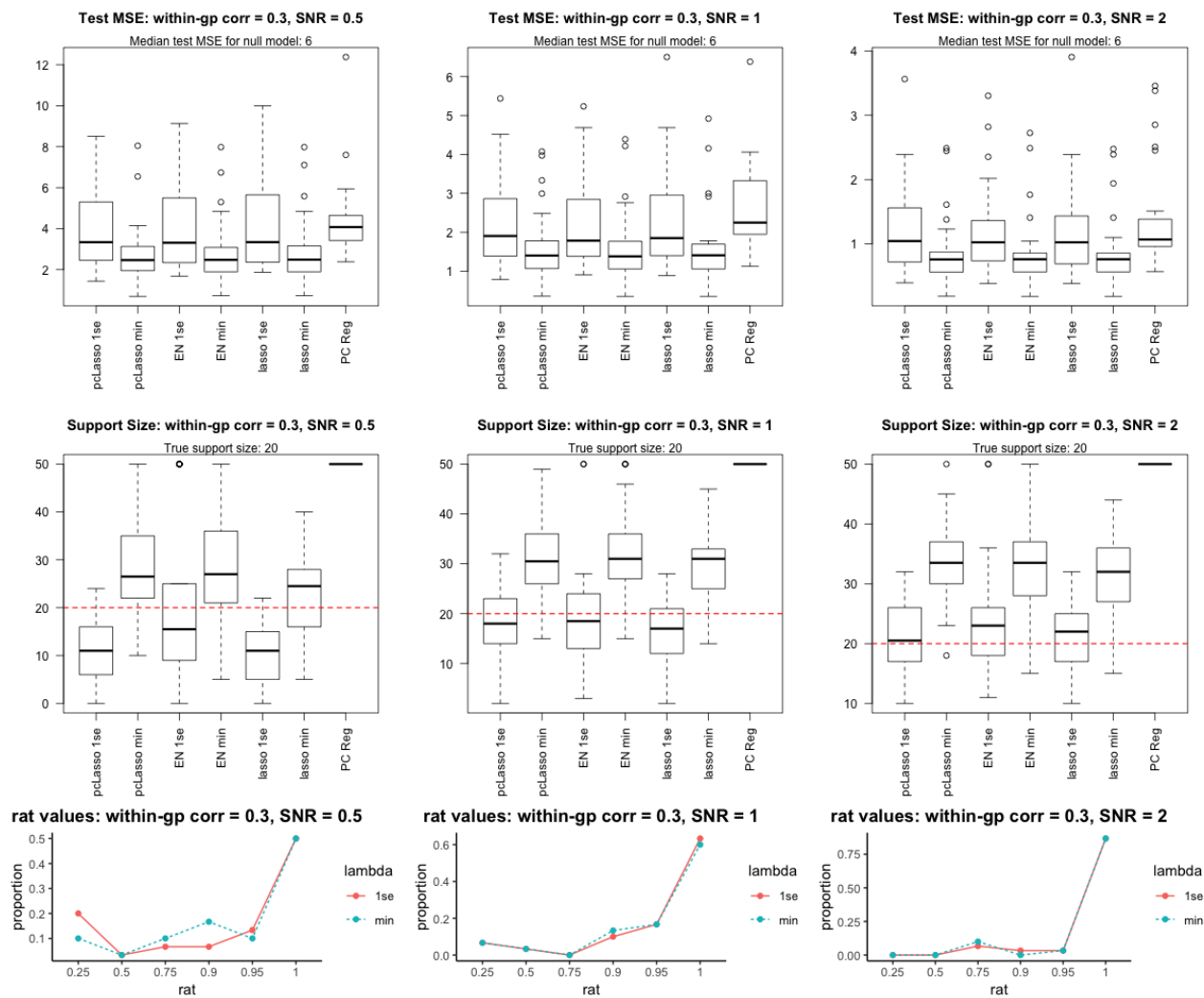


### D.2.5 Small $p$ with groups, “neutral court” for pcLasso

$n = 200$ ,  $p = 50$ , 5 groups of 10 uncorrelated predictors, response related to a random eigenvector in the first 2 groups.

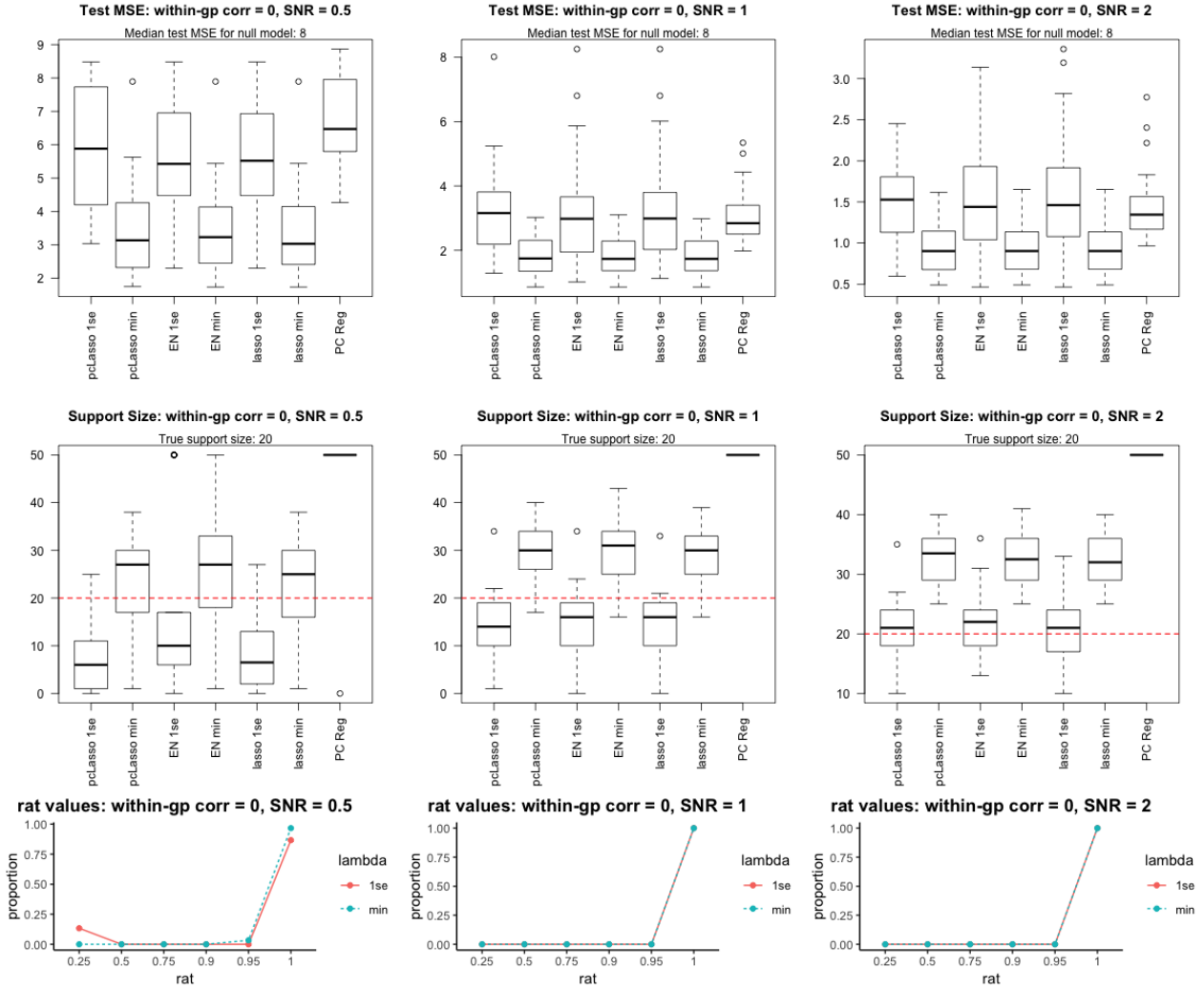


As above, but with predictors in the same group having pairwise correlation of 0.3.

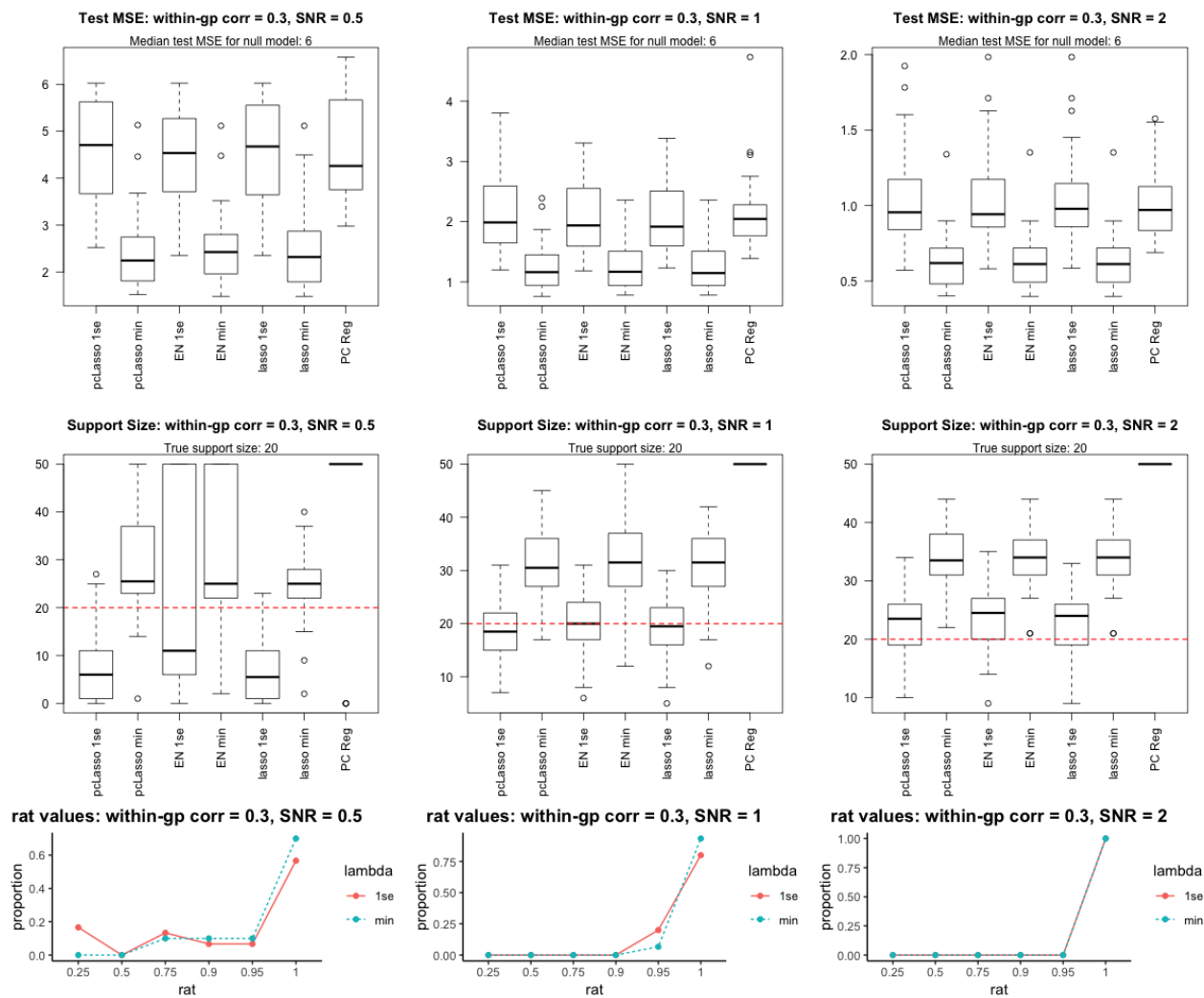


### D.2.6 Small $p$ with groups, “hostile court” for pLasso

$n = 200$ ,  $p = 50$ , 5 groups of 10 uncorrelated predictors, response related to the bottom eigenvector in the first 2 groups.



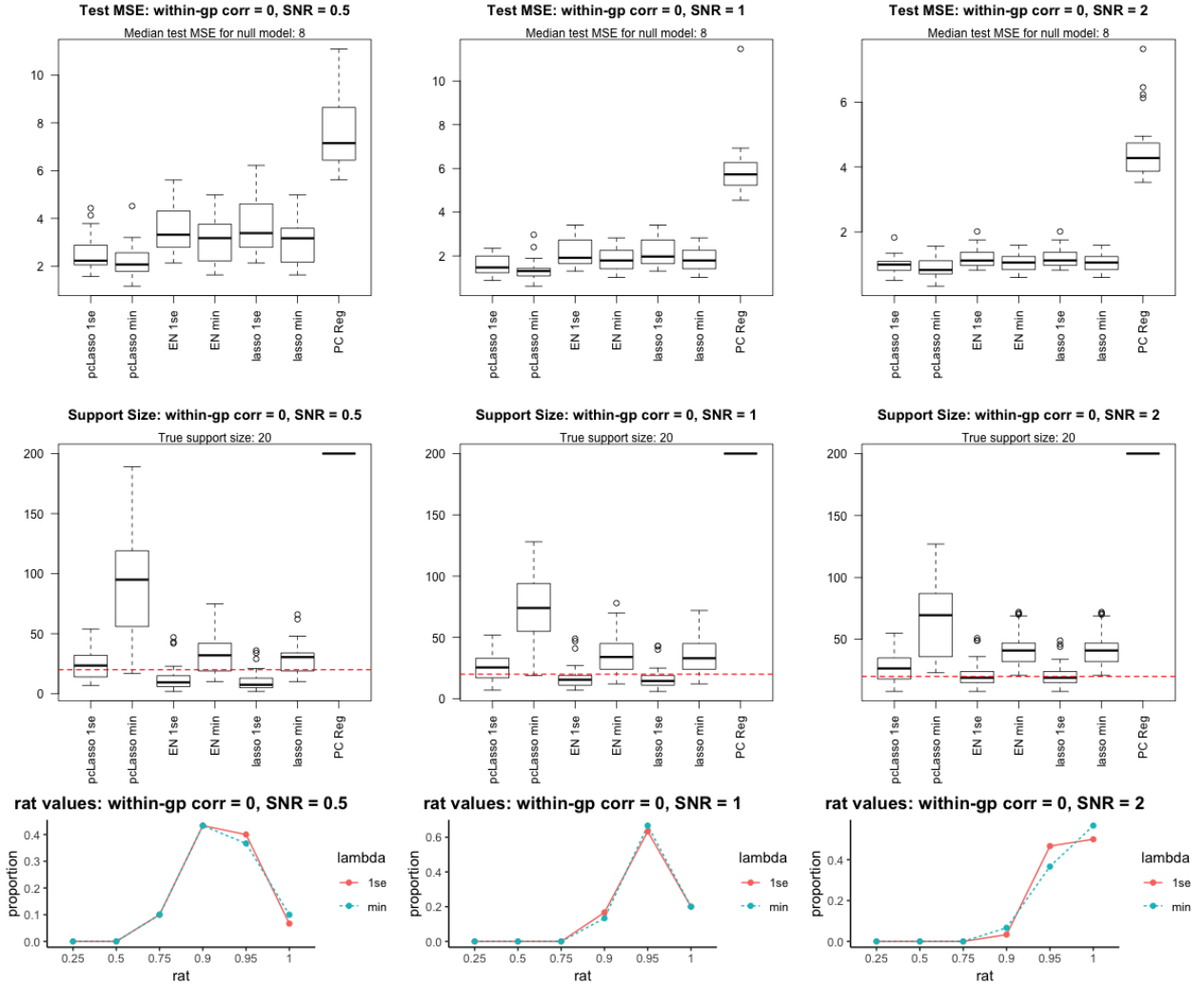
As above, but with predictors in the same group having pairwise correlation of 0.3.



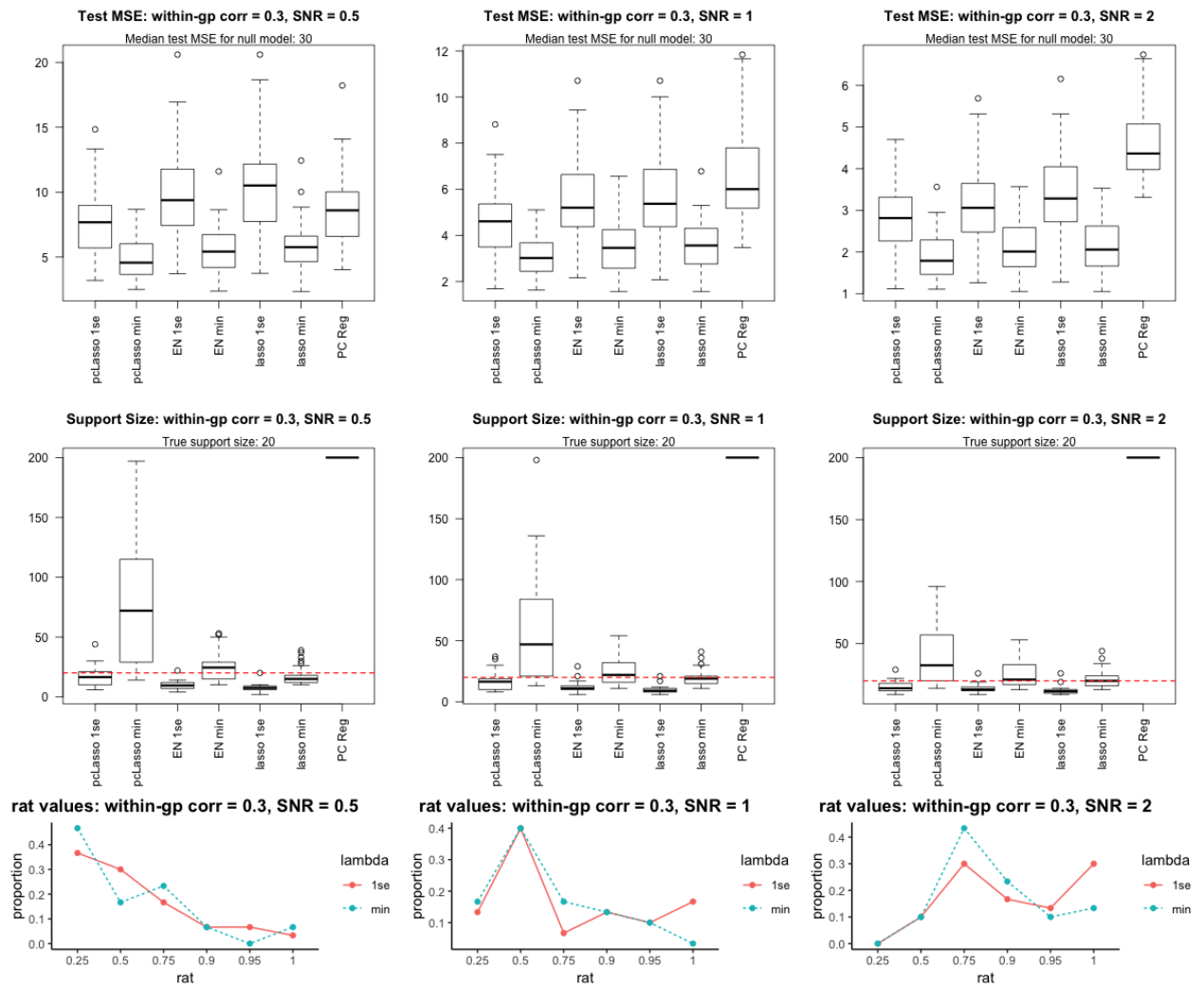


### D.2.7 Medium $p$ with groups, “home court” for pcLasso

$n = 200, p = 200$ , 10 groups of 20 uncorrelated predictors, response related to the top 2 eigenvectors in the first group.

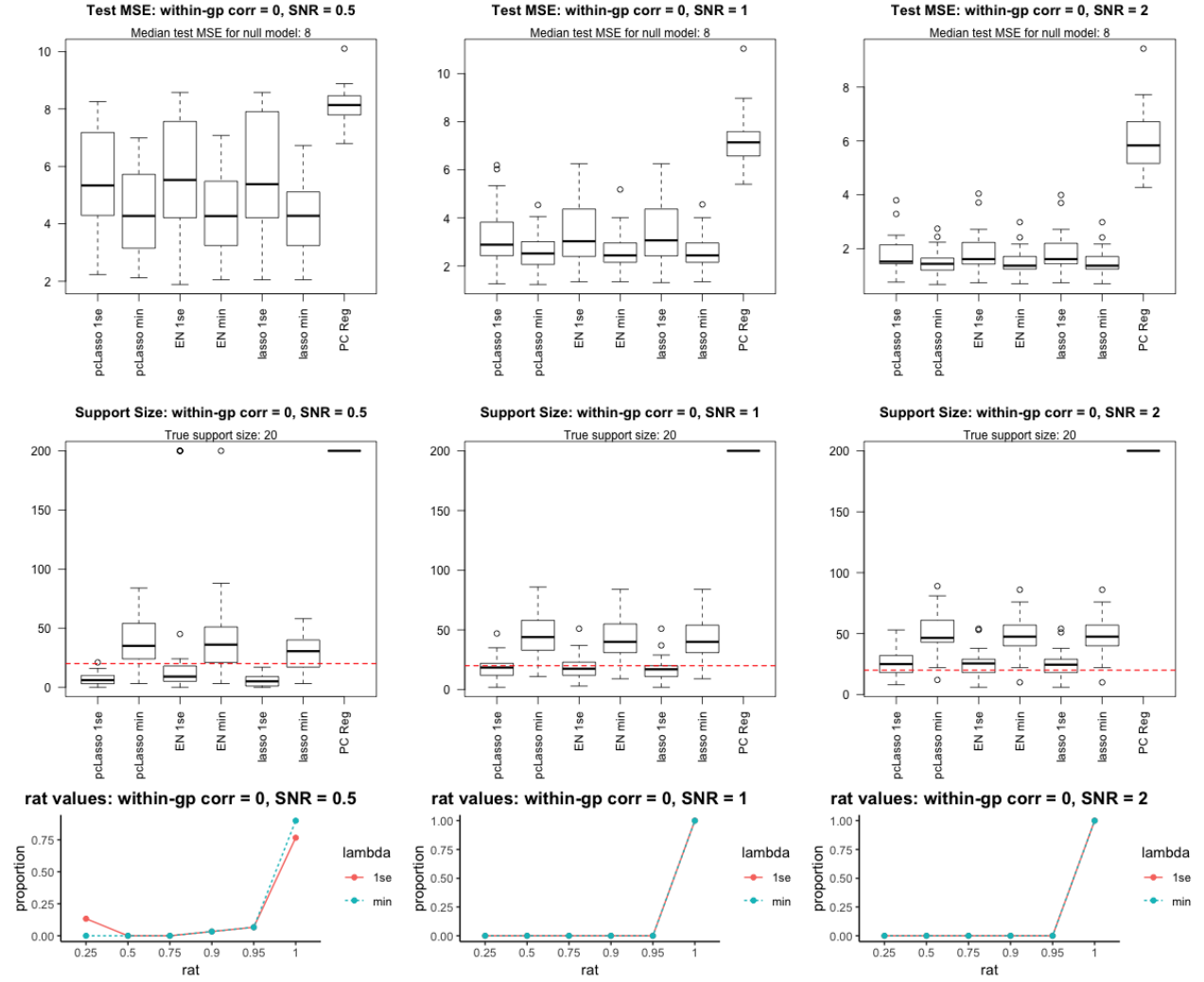


As above, but with predictors in the same group having pairwise correlation of 0.3.

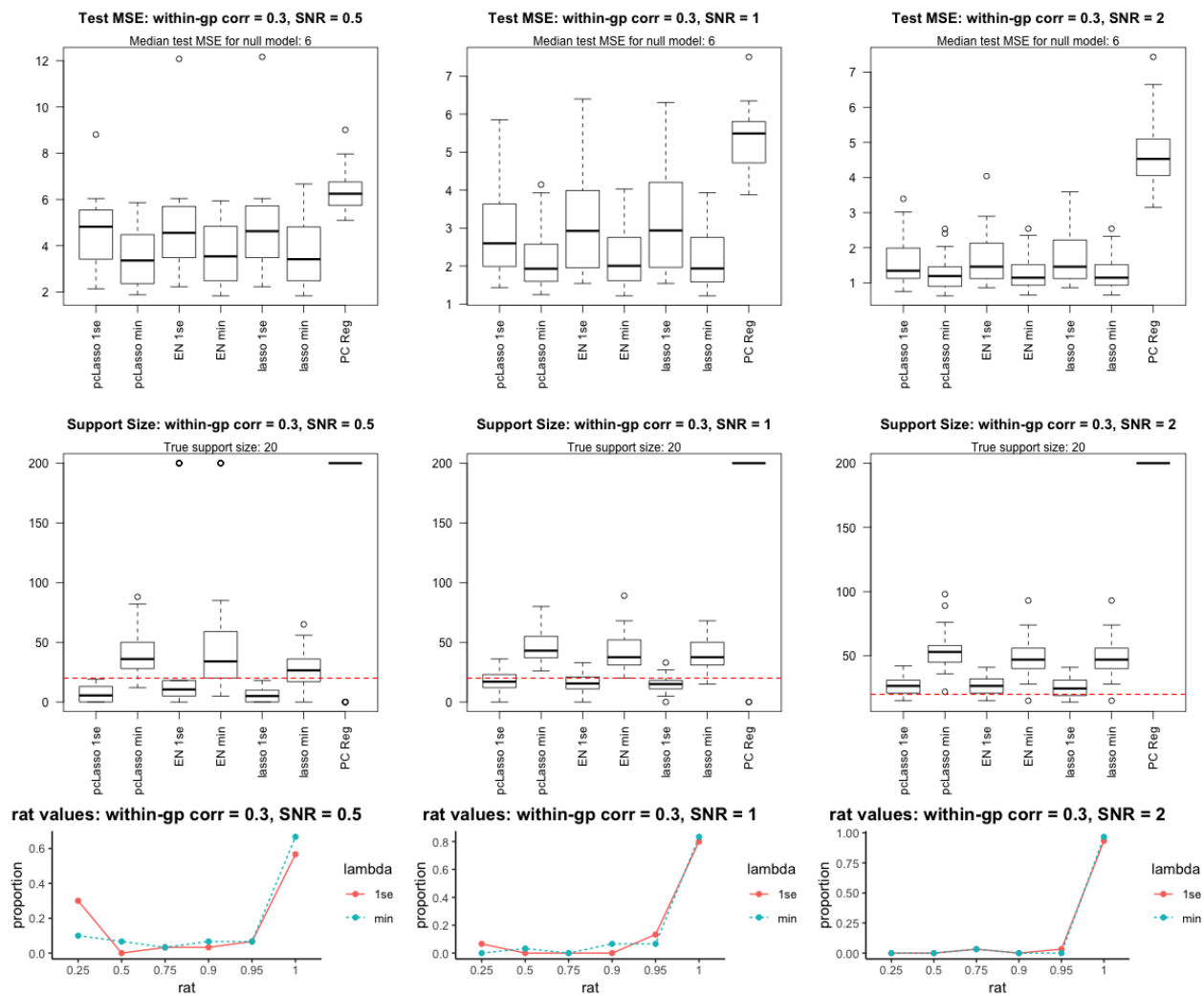


### D.2.8 Medium $p$ with groups, “netural court” for pcLasso

$n = 200, p = 200$ , 10 groups of 20 uncorrelated predictors, response related to 2 random eigenvectors in the first group.

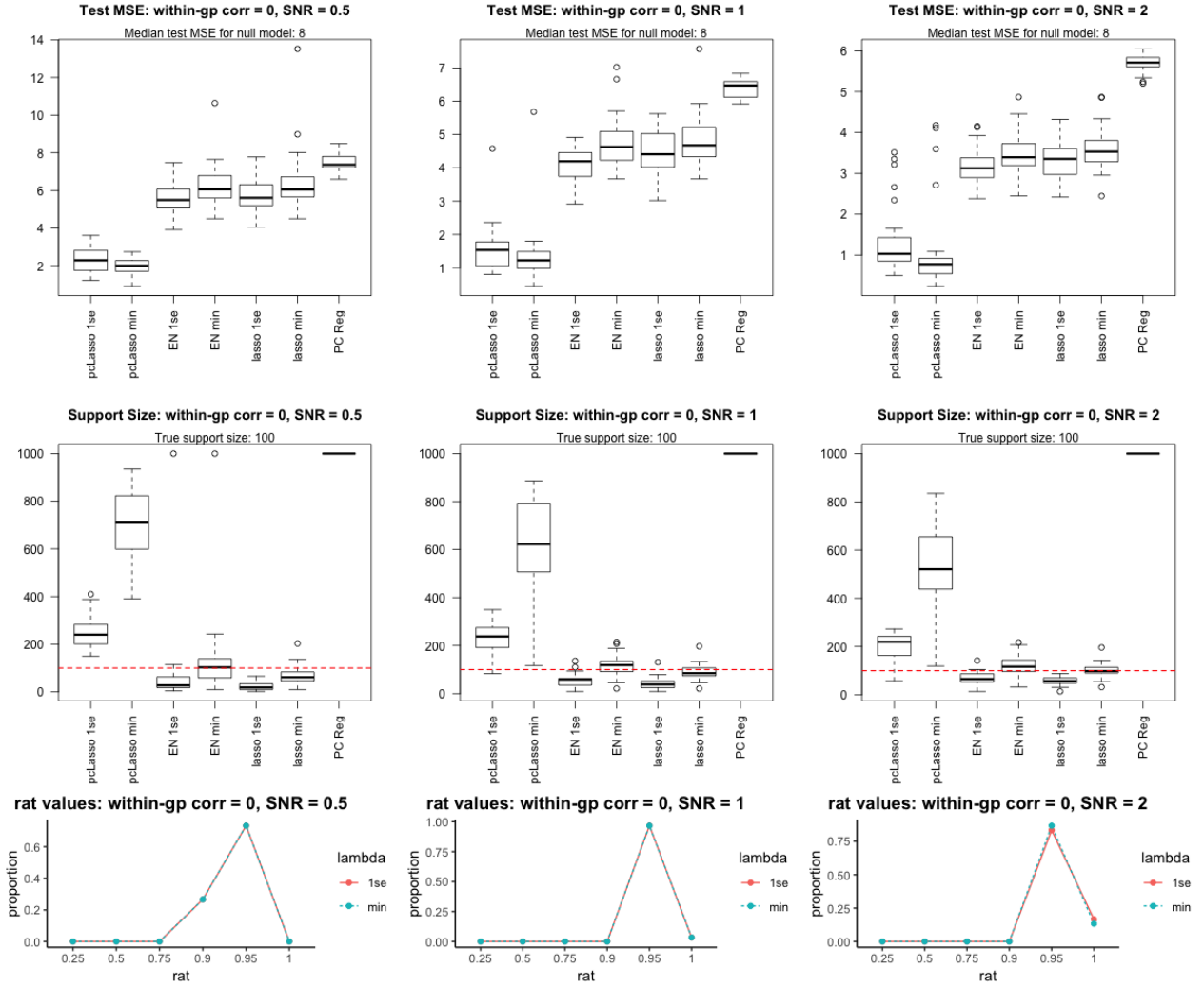


As above, but with predictors in the same group having pairwise correlation of 0.3.

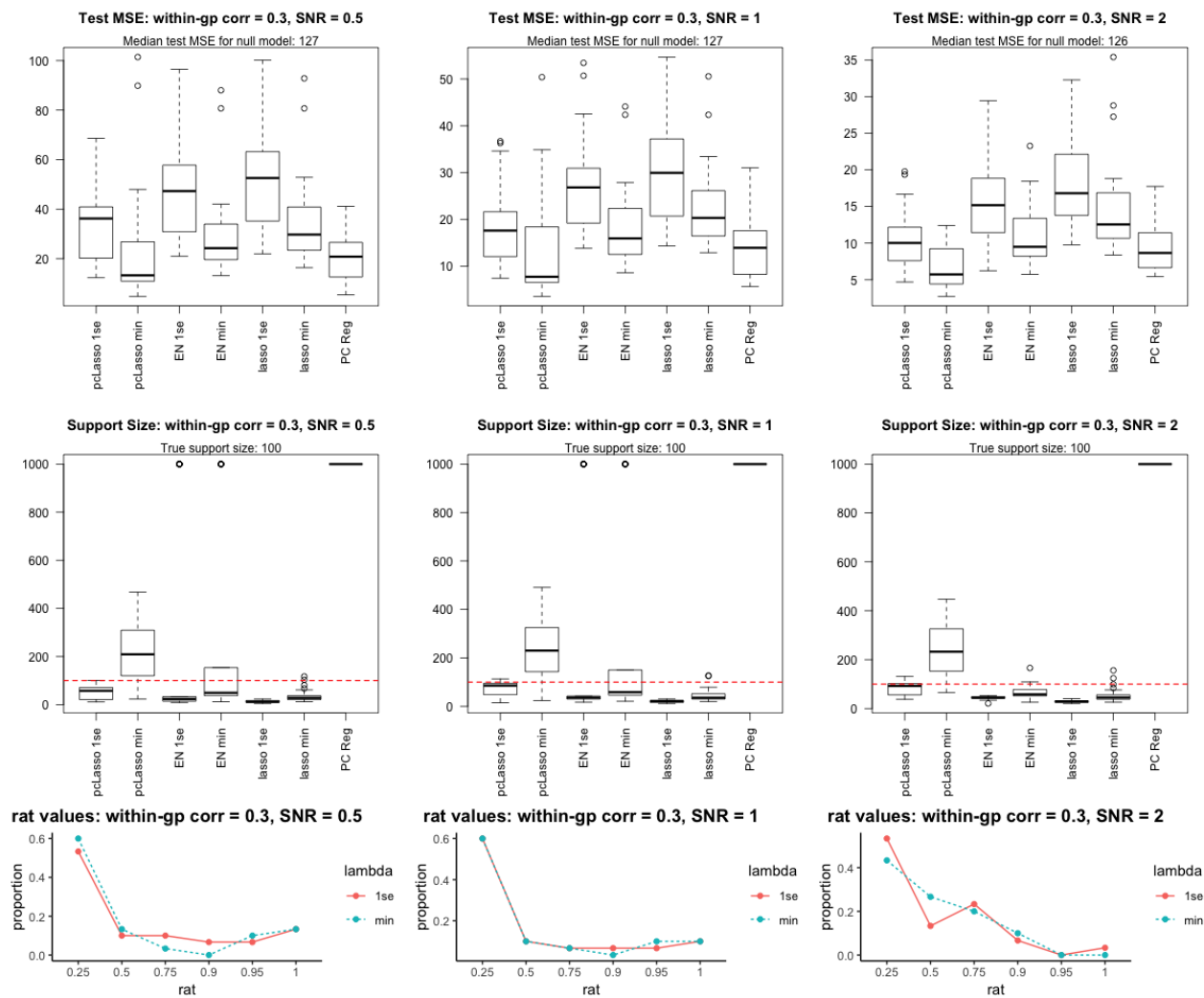


### D.2.9 Large $p$ with groups, “home court” for pcLasso

$n = 200$ ,  $p = 1000$ , 10 groups of 100 uncorrelated predictors, response related to the first 2 eigenvectors in the first group.

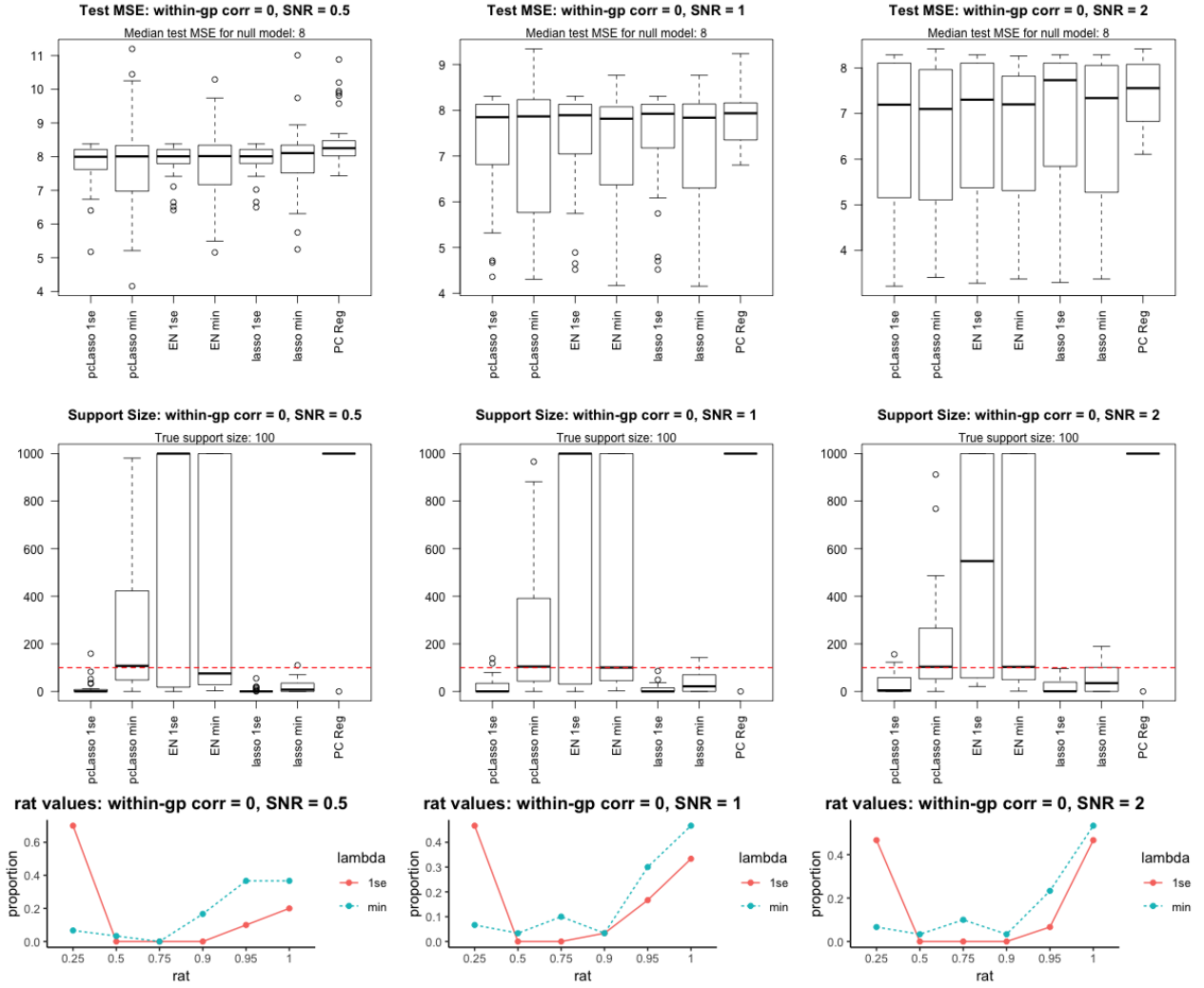


As above, but with predictors in the same group having pairwise correlation of 0.3.

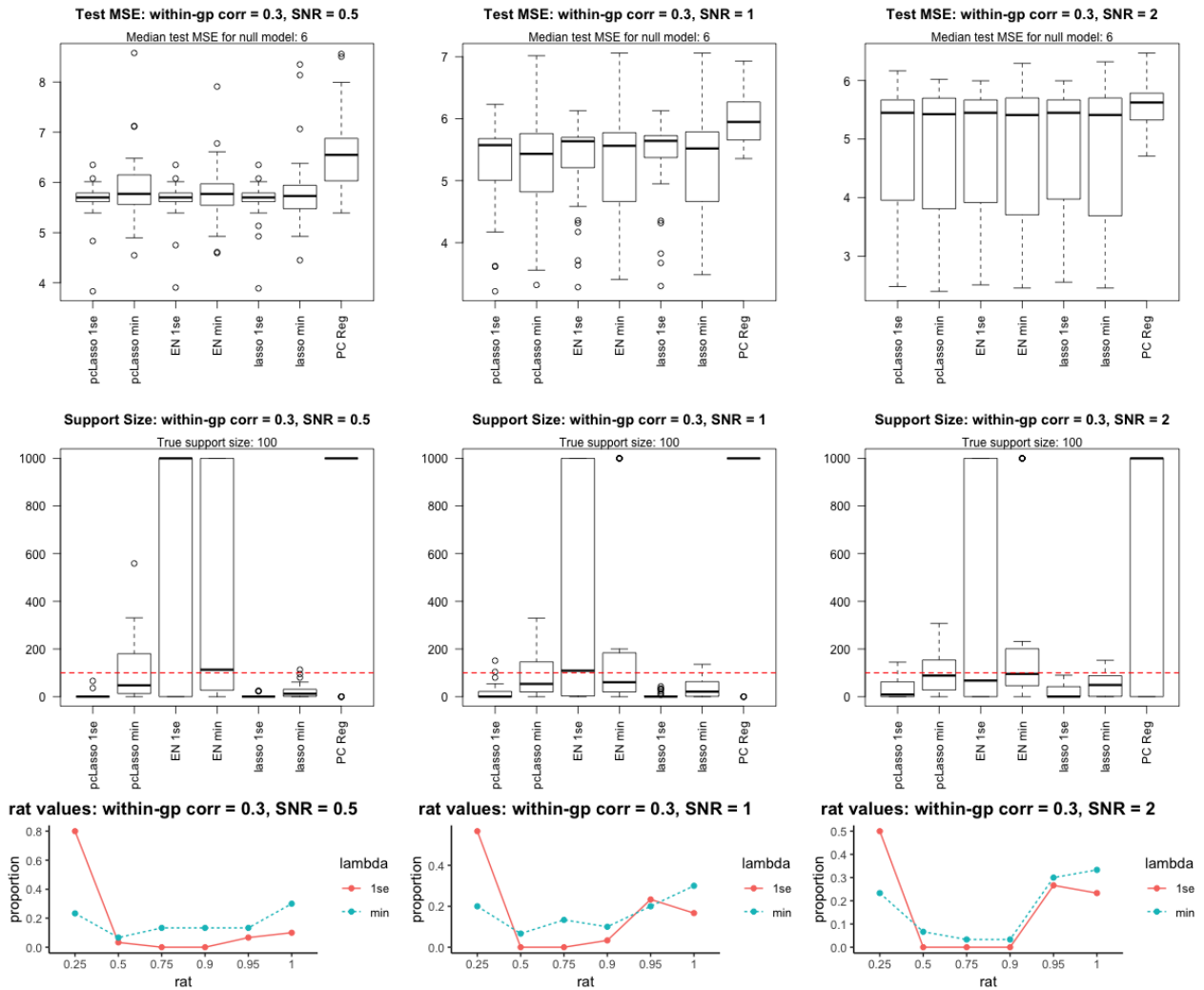


### D.2.10 Large $p$ with groups, “neutral court” for pcLasso

$n = 200$ ,  $p = 1000$ , 10 groups of 100 uncorrelated predictors, response related to 2 random eigenvectors in the first group.



As above, but with predictors in the same group having pairwise correlation of 0.3.



## References

- Brehereny, P. & Huang, J. (2015), ‘Group descent algorithms for nonconvex penalized linear and logistic regression models with grouped predictors’, *Statistics and Computing* **25**(2), 173–187.
- Efron, B. (1986), ‘How biased is the apparent error rate of a prediction rule?’, *Journal of the American Statistical Association* **81**, 461–70.
- Friedman, J., Hastie, T. & Tibshirani, R. (2010), ‘Regularization paths for generalized linear models via coordinate descent’, *Journal of Statistical Software* **33**, 1–22.
- Hastie, T., Tibshirani, R. & Wainwright, M. (2015), *Statistical learning with sparsity: the lasso and generalizations*, CRC press.



- Hoerl, A. E. & Kennard, R. (1970), ‘Ridge regression: biased estimation for nonorthogonal problems’, *Technometrics* **12**, 55–67.
- Jacob, L., Obozinski, G. & Vert, J.-P. (2009), ‘Group lasso with overlap and graph lasso’, *Proceedings of the 26th International Conference on Machine Learning* .
- Ming, Y. & Lin, Y. (2005), ‘Model selection and estimation in regression with grouped variables’, *Journal of the Royal Statistical Society Series B* **68**(1), 49–67.
- Simon, N., Friedman, J., Hastie, T. & Tibshirani, R. (2013), ‘A sparse-group Lasso’, *Journal of Computational and Graphical Statistics* **22**(2), 231–245.
- Simon, N. & Tibshirani, R. (2012), ‘Standardization and the group lasso penalty’, *Statistica Sinica* **22**(3), 983.
- Subramanian, A., Tamayo, P. Mootha, V. K., Mukherjee, S., Ebert, B. L., Gillette, M. A., Paulovich, A., Pomeroy, S. L., Golub, T. R., Lander, E. S. & Mesirov, J. P. (2005), ‘Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles’, *Proc. Natl. Acad. Sci. USA* **102**, 15545–15550.
- Tibshirani, R. (1996), ‘Regression shrinkage and selection via the lasso’, *Journal of the Royal Statistical Society Series B* **58**, 267–288.
- Tibshirani, R., Bien, J., Friedman, J. Hastie, T., Simon, N. Taylor, J. & Tibshirani, R. (2012), ‘Strong rules for discarding predictors in lasso-type problems’, *Journal of the Royal Statistical Society Series B* **74**, 245–266.
- Yuan, M. & Lin, Y. (2006), ‘Model selection and estimation in regression with grouped variables’, *Journal of the Royal Statistical Society, Series B* **68**, 49–67.
- Zeng, Y. & Breheny, P. (2016), ‘Overlapping group logistic regression with applications to genetic pathway selection’, *Cancer informatics* **15**, CIN–S40043.
- Zou, H. (2005), Some Perspectives on Sparse Statistical Modeling, PhD thesis, Department. of Statistics, Stanford.
- 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.
- Zou, H., Hastie, T. & Tibshirani, R. (2007), ‘On the “degrees of freedom” of the lasso’, *Annals of Statistics* **35**(5), 2173–2192.