

P8160 Project 1 - Designing a simulation study to compare variable Selection methods

Yuqi Miao (ym2771), Jiayi Shen (js5354), Jack Yan (xy2395), Jungang Zou (jz3183)

1 Objectives

In high-dimensional data analysis, variable selection is a common practice to find an optimal model that balances between model fitness and model complexity. Two well-known methods to conduct variable selection are **step-wise forward selection** and **automated LASSO regression**. The aim of this project is to investigate and illustrate how well two variable selection methods in identifying weak and strong predictors in high-dimensional data. To do so, we conducted simulation under different scenarios, and evaluate model performance by measures such as sensitivity, specificity, and F_1 score.

2 Statistical methods to be studied

2.1 Step-wise forward method

Forward selection starts with the empty model, and iteratively adds variable whose inclusion gives the most statistically significant improvement of the fit, and repeats this process until none improves the model to a statistically significant extent. AIC is commonly-used criterion for evaluating the model fit. For linear models,

$$AIC = n \ln \left(\sum_{i=1}^n (y_i - \hat{y}_i)^2 / n \right) + 2p$$

where \hat{y}_i is the fitted values from a model, and p is the dimension of the model (i.e., number of predictors plus 1).

2.2 Automated LASSO regression LASSO

Besides stepwise selection, another popular method for variable selection is called LASSO (least absolute shrinkage and selection operator). It estimates the model parameters by optimizing a penalized loss function:

$$\min_{\beta} \frac{1}{2n} \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \sum_{k=1}^p |\beta_k|$$

where λ is a tuning parameter. By forcing the sum of the absolute value of the regression coefficients to be less than a fixed value, LASSO forces certain coefficients to be set to zero, and thus effectively chooses a simpler model that does not include those coefficients. Therefore, LASSO is able to achieve both prediction accuracy and variable selection at the same time. Cross-validation (CV) is the most common selection criteria for LASSO.

3 Scenarios to be investigated

In the simulation study, we fix the total number of predictors ($p = 200$), among which 100 are null predictors. The correlation coefficient between strong signals and weak but correlated signals is 0.5. Here we change the ratio of signal types r and sample size n .

3.1 Ratio of signal types

The strength of signals are defined by the following criteria:

Definition of strong predictors:

$$S_1 = \{j : |\beta_j| > c\sqrt{\log(p)/n}, \text{some } c > 0, 1 \leq j \leq p\}$$

Definition of weak but correlated predictors:

$$S_2 = \{j : 0 < |\beta_j| \leq c\sqrt{\log(p)/n}, \text{some } c > 0, \text{corr}(X_i, X_{j'}) \neq 0, \text{some } j' \in S_1, 1 \leq j \leq p\}$$

Definition of weak and independent predictors:

$$S_3 = \{j : 0 < |\beta_j| \leq c\sqrt{\log(p)/n}, \text{some } c > 0, \text{corr}(X_i, X_{j'}) = 0, \text{all } j' \in S_1, 1 \leq j \leq p\}$$

Definition of null predictors:

$$S_4 = \{j : |\beta_j| = 0, 1 \leq j \leq p\}$$

Here we consider 2 ratios of predictor types (r: strength: weak_but_correlated: weak_and_independent)
1) r = 5:1:4 2) r = 5:3:2

3.2 Sample size

The following 3 sample sizes are used in the simulation.

- 1) n = 200
- 2) n = 1000
- 3) n = 5000

We use n = 200 to test the robustness of the two methods facing high-dimensional data. We are also interested to see how the performance of the methods would improve with increased sample size.

3.3 Missing weak predictors

For each scenario above, we also investigated the case where those weak but correlated predictors are removed from the original data, thus trying to simulate the situation when weak but correlated predictors are not collected during the data collection process.

4 Methods for generating data

4.1 Data generation

The data matrix \mathbf{X} is generated by R function `MASS::mvrnorm`. A pre-defined covariance matrix is passed to `mvrnorm` function to specify the correlation between predictors. To ensure the positive definite attribute of covariance matrix, we restrict one strong predictor to be correlated with one weak predictor.

4.2 True model Parameter

Linear model parameters are randomly generated from a uniform distribution, subject to the definition of signal strength and ratio of predictor types. All the coefficients are positive values.

4.1 Generating true outcome Y

True distribution of outcome variable is defined as

$$Y \sim N(X^T \beta, \sigma^2)$$

Where X is the data matrix, β is the parameters vector and σ^2 is the constant variance in normal distribution. The variance is fixed at 9.

5 performance measures

5.1 Predictor identification performance

In order to compare the identification performance for these 2 methods, we regard the identification as a classification problem, where the signal predictors are defined as positive and null predictors are defined as negative. 3 indicators have been established:

$$recall = sensitivity = \left(\frac{True\ positive}{True\ positive + False\ negative} \right)$$

$$specificity = \left(\frac{True\ negative}{True\ negative + False\ positive} \right)$$

$$accuracy = \frac{True\ selection}{total\ number\ of\ predictors}$$

$$F_1 = \left(\frac{2}{recall^{-1} + precision^{-1}} \right)$$

Where precision is defined as above, and recall is defined as

$$precision = \left(\frac{True\ positive}{True\ positive + False\ positive} \right)$$

5.2 Parameter Estimation performance

To compare the parameter estimation performance, 3 indicators were calculated to evaluate the estimation: mean bias is calculated as the mean difference between true parameter and estimated parameters for a set of parameters, variance is defined as the variance of the estimated parameter among simulation. mean squared error (MSE) is also used to assess estimation performance.

- bias

$$\frac{1}{p} \sum_{j=1}^p (\hat{\beta}_j - \beta_j)$$

- variance

$$\frac{1}{p} \sum_{j=1}^p (\hat{\beta}_j - \overline{\beta_j})^2$$

- MSE

$$\frac{1}{p} \sum_{j=1}^p (\hat{\beta}_j - \beta_j)^2$$

6 Simulation results

6.1 Identification performance

F1 score (see Fig 1.) is used as an indicator for an overall performance assessment for identification signals. There is a clear trend for both methods that when the ratio of predictors to sample size is decreasing, the overall performance enhances, and there is no significant difference between these two methods when sample size is large enough compared to predictor numbers, while in high dimensional situation (where $n = p$), stepwise method has a poor performance.

In terms of variable selection size (see Fig 2.) , LASSO tends to choose more predictors while stepwise tends to choose less overall, specifically in the settings of this report, where there are 100 true signals and 100 noises, LASSO selects more predictors than true occasion in all scenarios, while stepwise tends to choose less. When comparing in terms of sample size, LASSO performs stably in both high dimensional (where $n = p$) and normal scenario, but stepwise tends to lose its select ability when the ratio between number of parameters and sample size is large.

Sensitivity (see Fig 3.) is an indicator for how many predictors can be captured by the methods. For LASSO, strong predictors can be perfectly captured in all scenario with stable performance, while for weak predictors, difference occurs when correlated-to-independent ratios change. For less correlated scenario, LASSO tends to perform better to capture weak and independent predictors compare to weak but correlated predictors; while in more correlated case, there is no clear difference between these two types of predictors, and performance for capturing weak but correlated predictors is relatively stable. For stepwise methods, except for the same stability for capture weak but correlated predictors in more correlated scenario, there is no clear difference in different settings. When comparing the sensitivity of these 2 methods, LASSO performs better overall but with a comparatively higher variance, and stepwise has a relatively poor performance in aspect of capturing weak predictors in all settings, but the performance stabilises at 0.1 to 0.25.

Correspondingly, the specificity (see Fig 4.) for stepwise method is relatively higher since there is an overall trend for over-screening, but the specificity performance for LASSO is also acceptable especially when the predictor-to-sample-size ratio decreases.

Finally, comparing the accuracy of these two methods (see Fig 5.) , LASSO performs better in every category of predictors in 6 settings. Above all, in terms of variable identification performance, LASSO tends to over-select predictors and performs stably in high dimensional scenarios, while stepwise tends to under-select predictors with a worse overall performance and less accuracy. As compare to the correlation between predictors, there is also a clear trend that a lower correlation ratio of predictors tends to decrease the overall variance of identification performance; this may due to the high correlation between predictors may reduce the overall variance of the sample space and leads a more stable selection results.

6.2 Parameter Estimation

When $n = 200$, the parameter estimate of stepwise selection is highly biased and unstable, while LASSO method has both low bias and variance for the estimation of all types of predictors. When $n = 1000$ and 5000, LASSO underestimates the coefficient of true signals, especially strong signals, while stepwise selection method is approximately unbiased in estimating strong signals. Stepwise selection also overestimates the effect of weak but correlated signals. Both methods overestimates the effect of null predictors.

In terms of RMSE, LASSO has higher RMSE in estimating strong signals, due to its high bias. For predictors other than strong signals, the RMSE of LASSO is much lower than stepwise selection.

6.3. comparison of missing vs no missing

In this part, we will discuss how missing “weak” signals impact the coefficients in linear models. Due to the definition of “missing variables”, we consider 2 situations. The first part is the original “no missing” models, which we have discussed above. The second part is the “missing” models, which were constructed

by deleting the weak_but_correlated data and weak_independent data. After model constructions, we'd like to analyze the performance of "no missing" models and "missing" models.

The procedure to construct models in 2 situations is as follows:

- Use calculated covariance matrix to generate "no missing" data X_{true} , and generate random noise ϵ .
- Use pre-calculated β_{true} , "no missing" data X_{true} and random noise ϵ to calculate response variable $Y = X_{true}\beta_{true} + \epsilon$
- Apply LASSO and stepwise regression on "no missing" data X_{true} , to get the estimated "no missing" model $\hat{Y} = X_{true}\hat{\beta}_{no}$ and estimated coefficient $\hat{\beta}_{no}$.
- Delete the weak_but_correlated and weak_independent variables in "no missing" data X_{true} , to get "missing" data $X_{missing}$.
- Apply LASSO and stepwise regression on "missing" data $X_{missing}$, to get the estimated "missing" model $\hat{Y} = X_{missing}\hat{\beta}_{missing}$ and estimated coefficient $\hat{\beta}_{missing}$.

1) bias

As mentioned above, bias is an important criterion to assess model performance. As in Fig 11, we can see that with the increasing number of samples, the biases for both models are becoming small, because of the average effect on outliers. Another important criterion variance also shows negative relationship when sample size becomes large. So for both "missing" and "no missing" models, the large number of sample size has positive effect to decrease both bias and variance.

For LASSO models and stepwise models, we find different result. For stepwise model, we can find the average of bias has very little difference between "missing" and "no missing" data. However, things become distinct for LASSO models. With the "missing" data, we find LASSO models perform better than "no missing" data, regardless of sample size. This significant difference may result from the different mechanism of model selection. For stepwise model, a variable will be decided to include in the model or not. On the other hand, the LASSO model uses shrinkage method and considers the co-effect among all the variables. If a variable is not important, LASSO model will gradually shrink its coefficient to 0 by considering the relationships for other variables. As a result, in model selection, LASSO model is easy to be influenced by weak signals. On the contrary, stepwise model is robust for weak signals.

2) RMSE

After we draw the conclusion of bias, we need to analyse the rmse for both models. As in Fig 10, we can see that with the increasing number of samples, the rmse for both models are becoming small, due to the decrease for both bias and variance.

For LASSO models and stepwise models, we find the result same as bias. For stepwise model, the average of rmse has very little difference between "missing" and "no missing" data. However, the LASSO model has smaller rmse and large rmse variance with "missing" data. This result shows the disappearance of weak signals will improve the performance of LASSO model, and has little significance for stepwise.

However, as is known to us, LASSO model is a regularized model that can decrease the predictive error on test dataset. This property indicates LASSO models have better generalization ability over all the sample space. Due to the bias nad variance tradeoff, LASSO models increase its bias to decrease its variance. So, on training dataset, the rmse of LASSO will be larger than rmse of stepwise model. However, on the test dataset, the rmse of LASSO will be smaller than rmse of stepwise model. In conclusion, stepwise model is likely to be overfitting on the training dataset.

Finally, we can draw the conclusion, that stepwise model is robust to weak signals, and LASSO model shows a significant improvement with "missing" data.

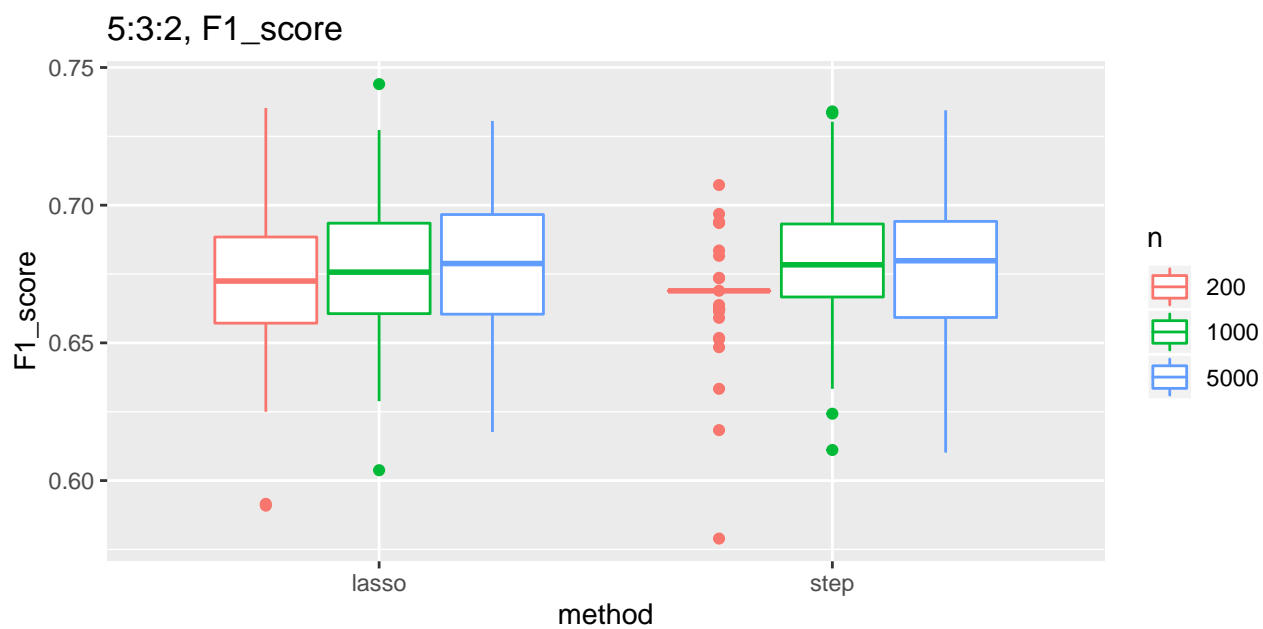
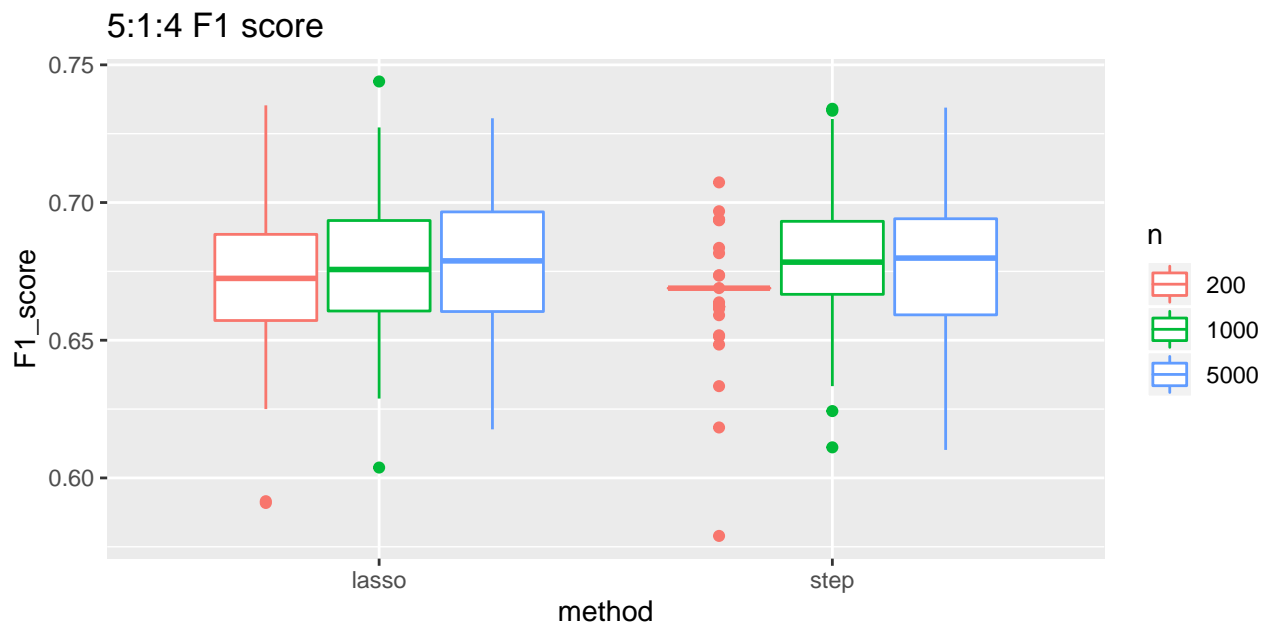


Figure 1: F1 score for 2 variable selection methods in 6 settings

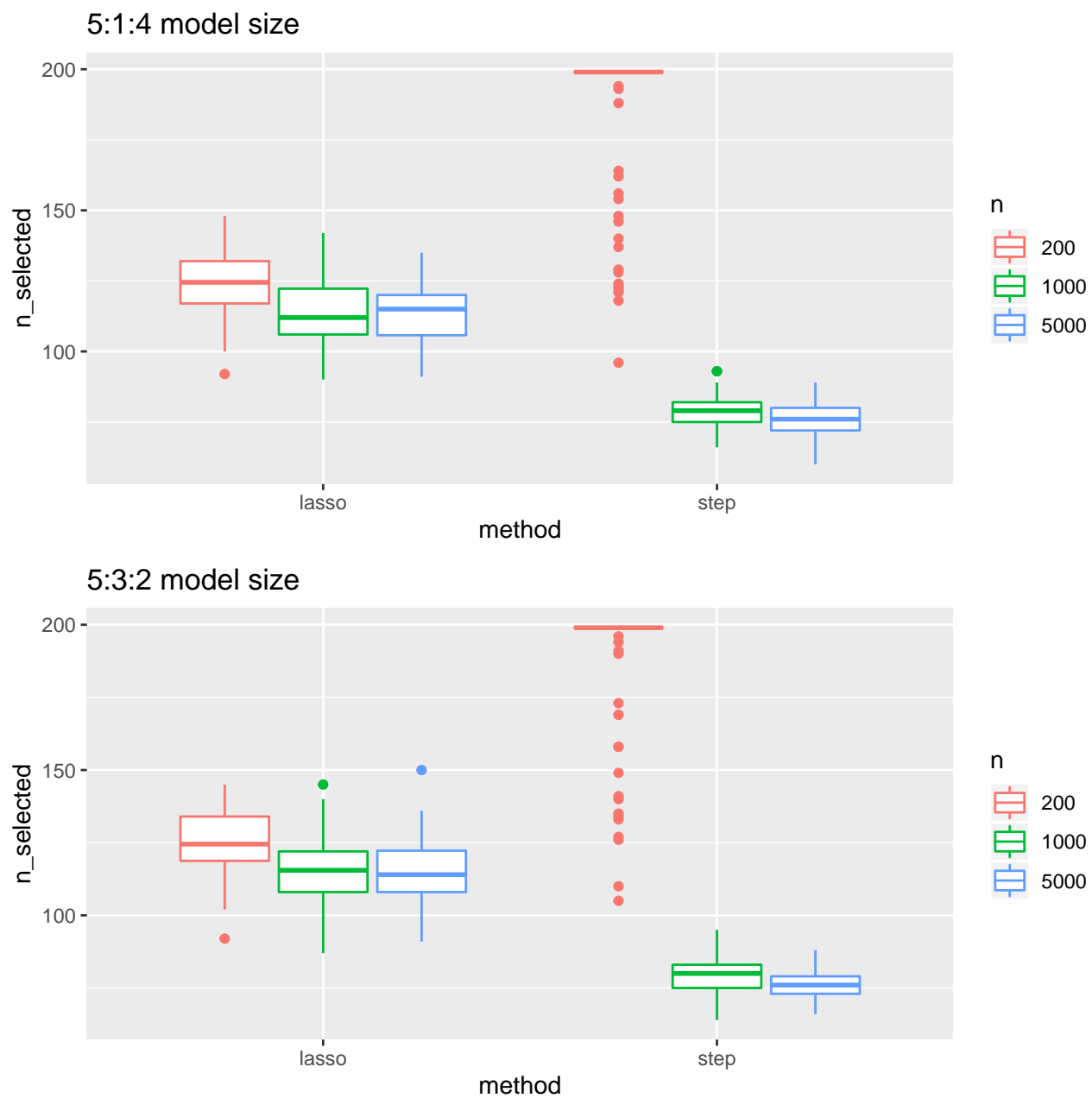


Figure 2: Model size for 2 variable selection methods in 6 settings

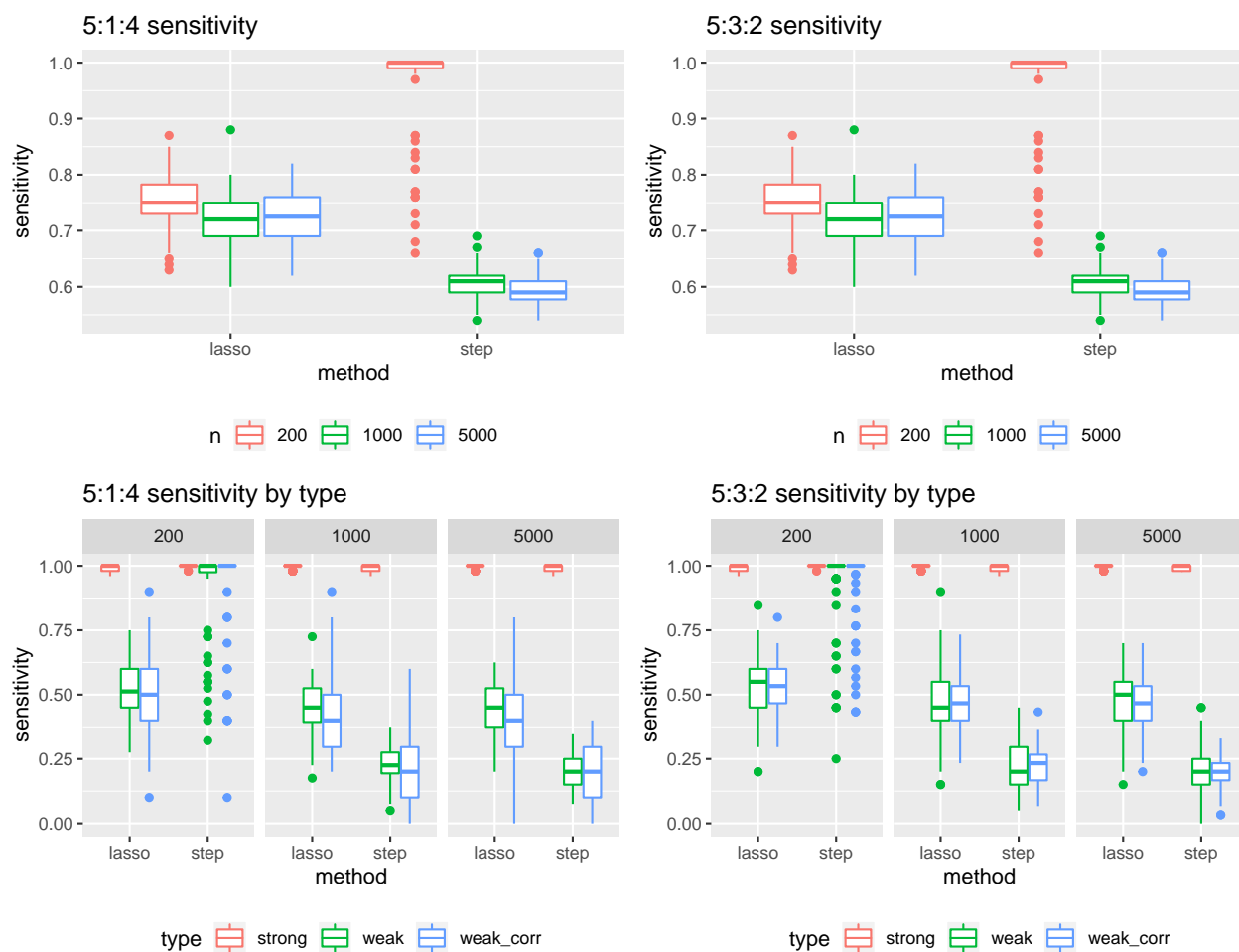


Figure 3: Sensitivity for 2 variable selection methods in 6 settings

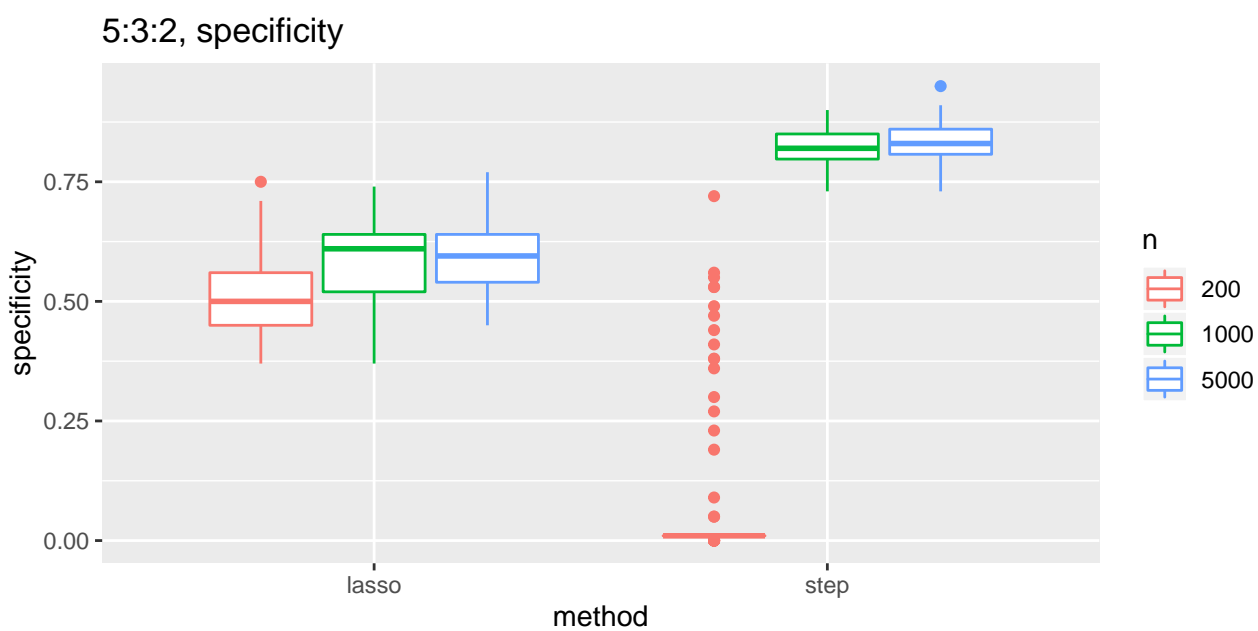
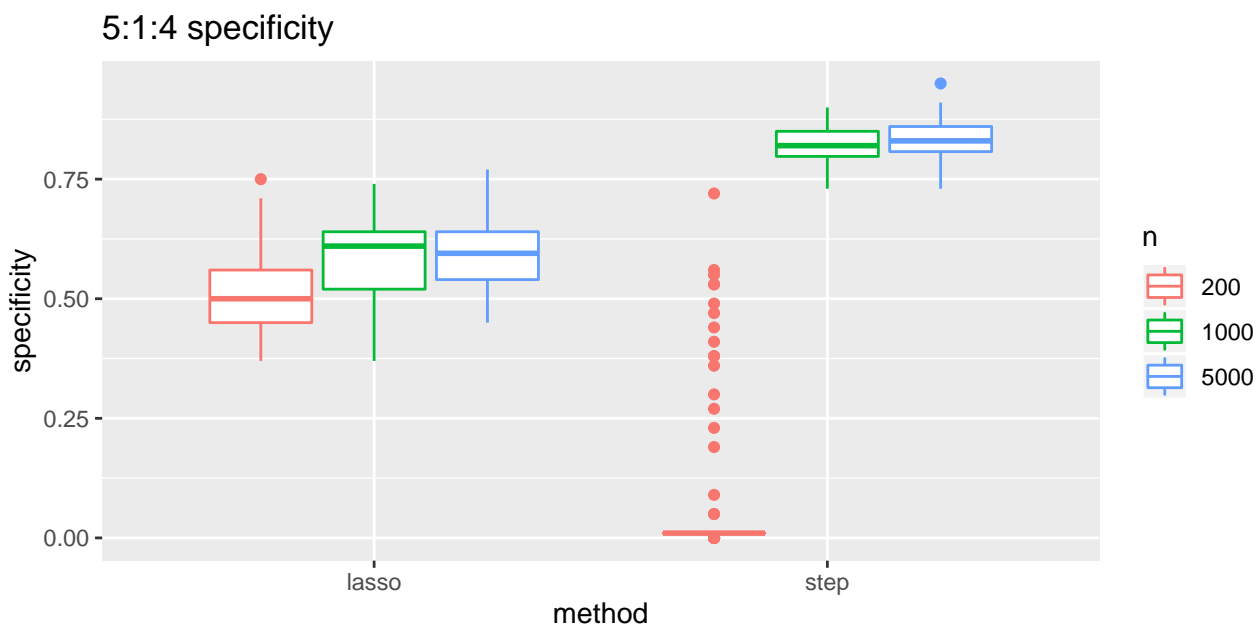


Figure 4: Specificity for 2 variable selection methods in 6 settings

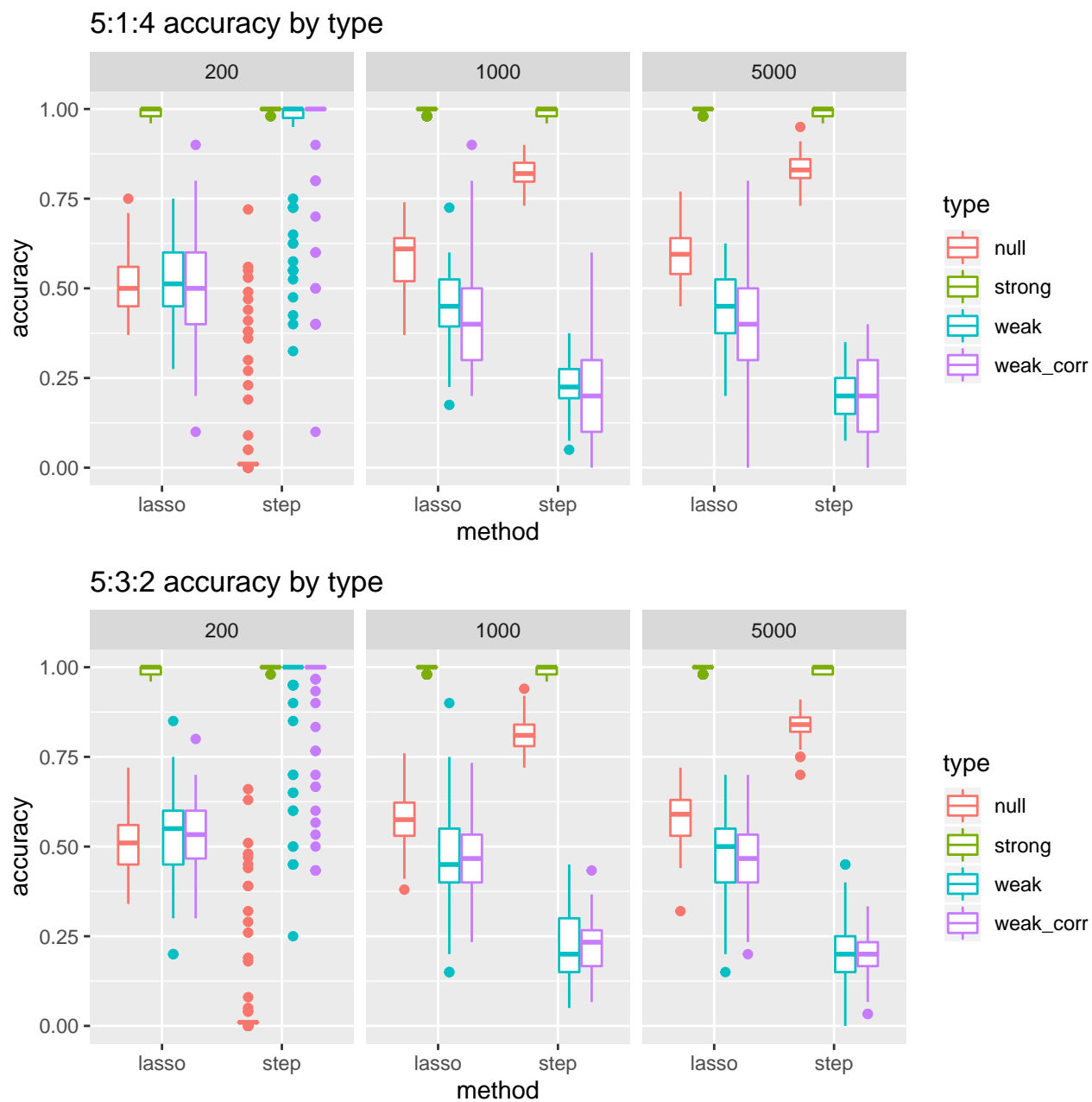


Figure 5: Specificity for 2 variable selection methods in 6 settings

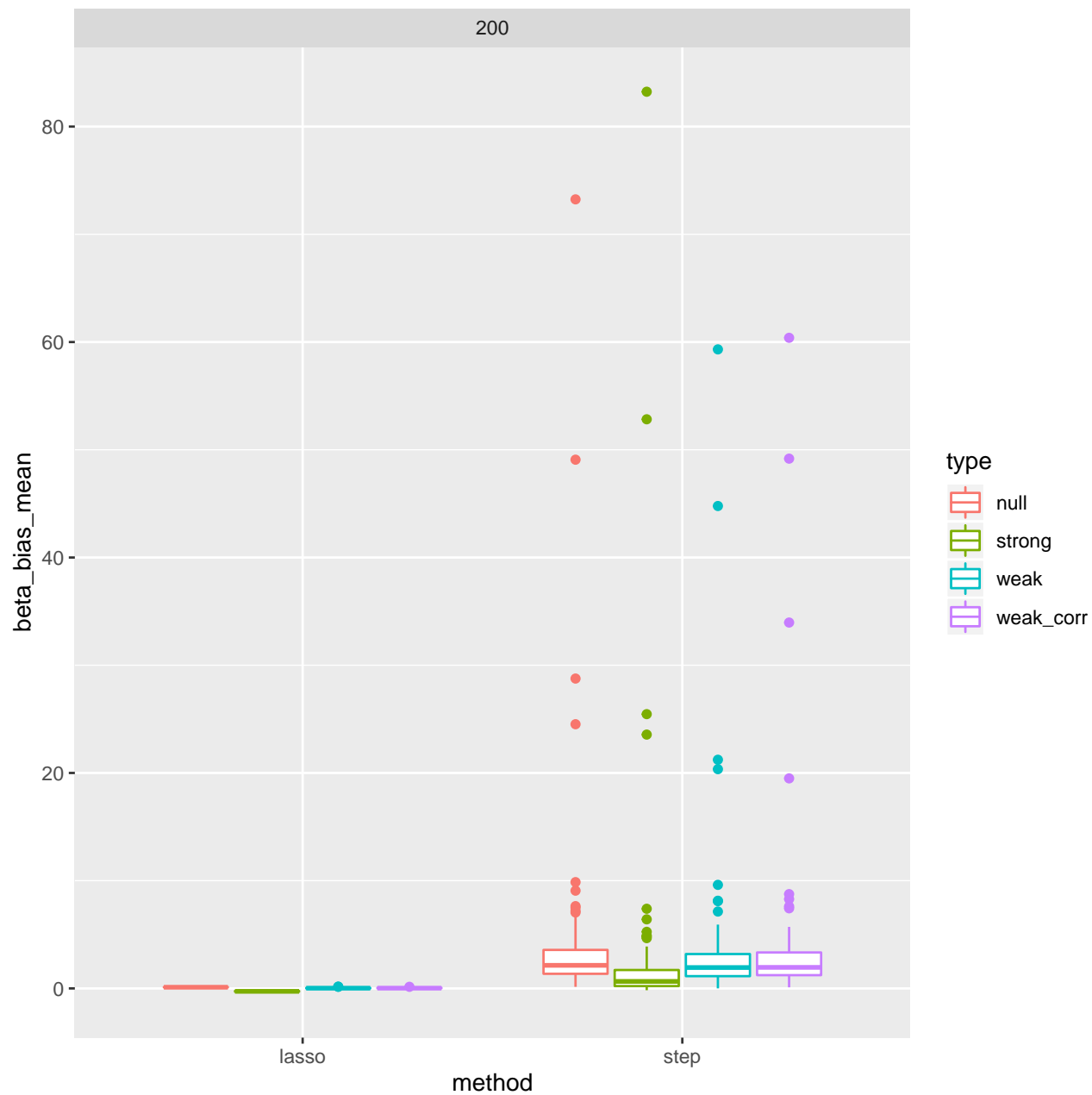


Figure 6: Mean beta bias $n = 200$, 40 weak-but-correlated predictors

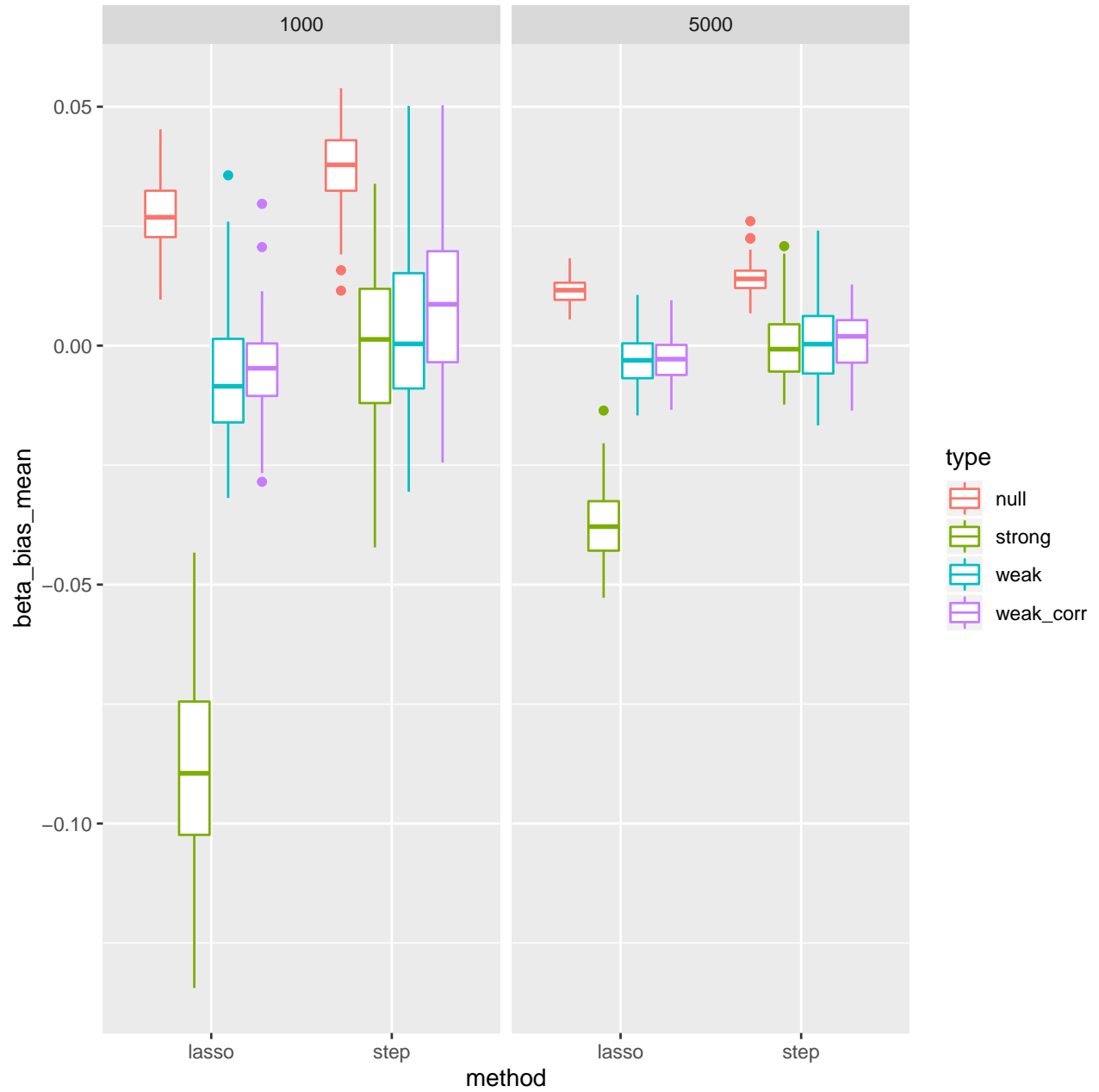


Figure 7: Mean beta bias $n = 200$, 40 weak-but-correlated predictors

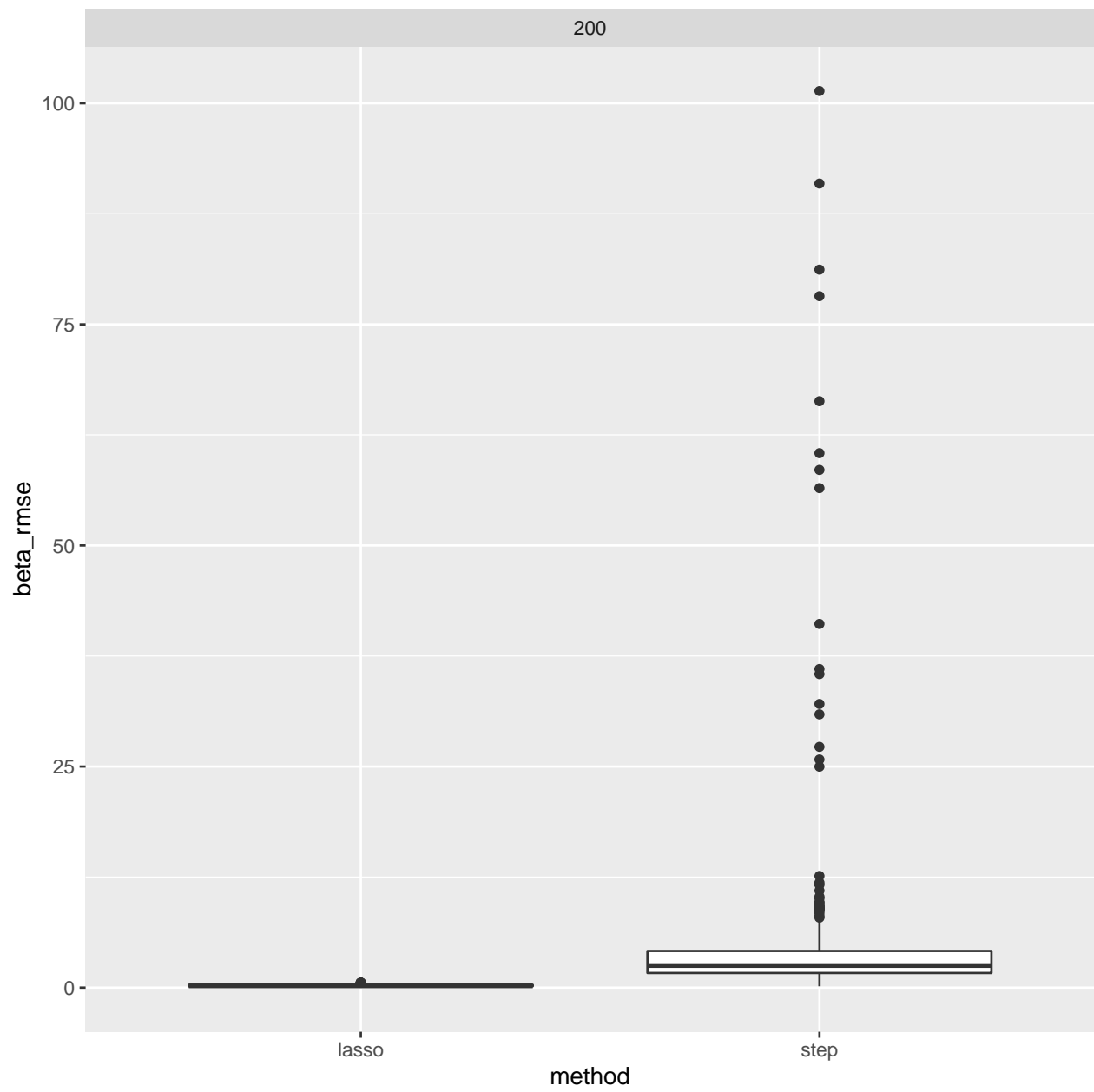


Figure 8: Mean beta bias $n = 1000, 2000$

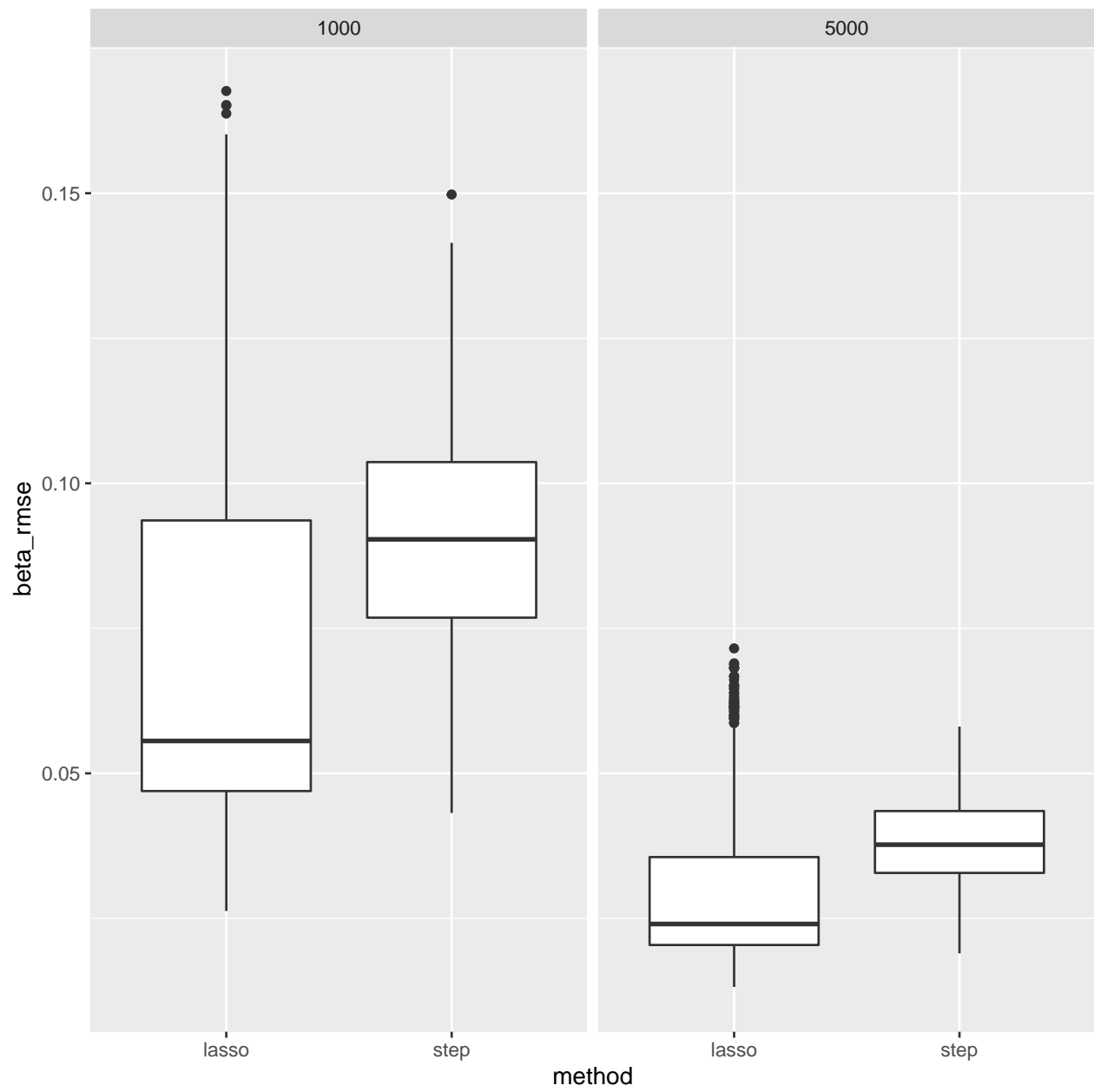


Figure 9: Beta RMSE $n = 1000, 2000$

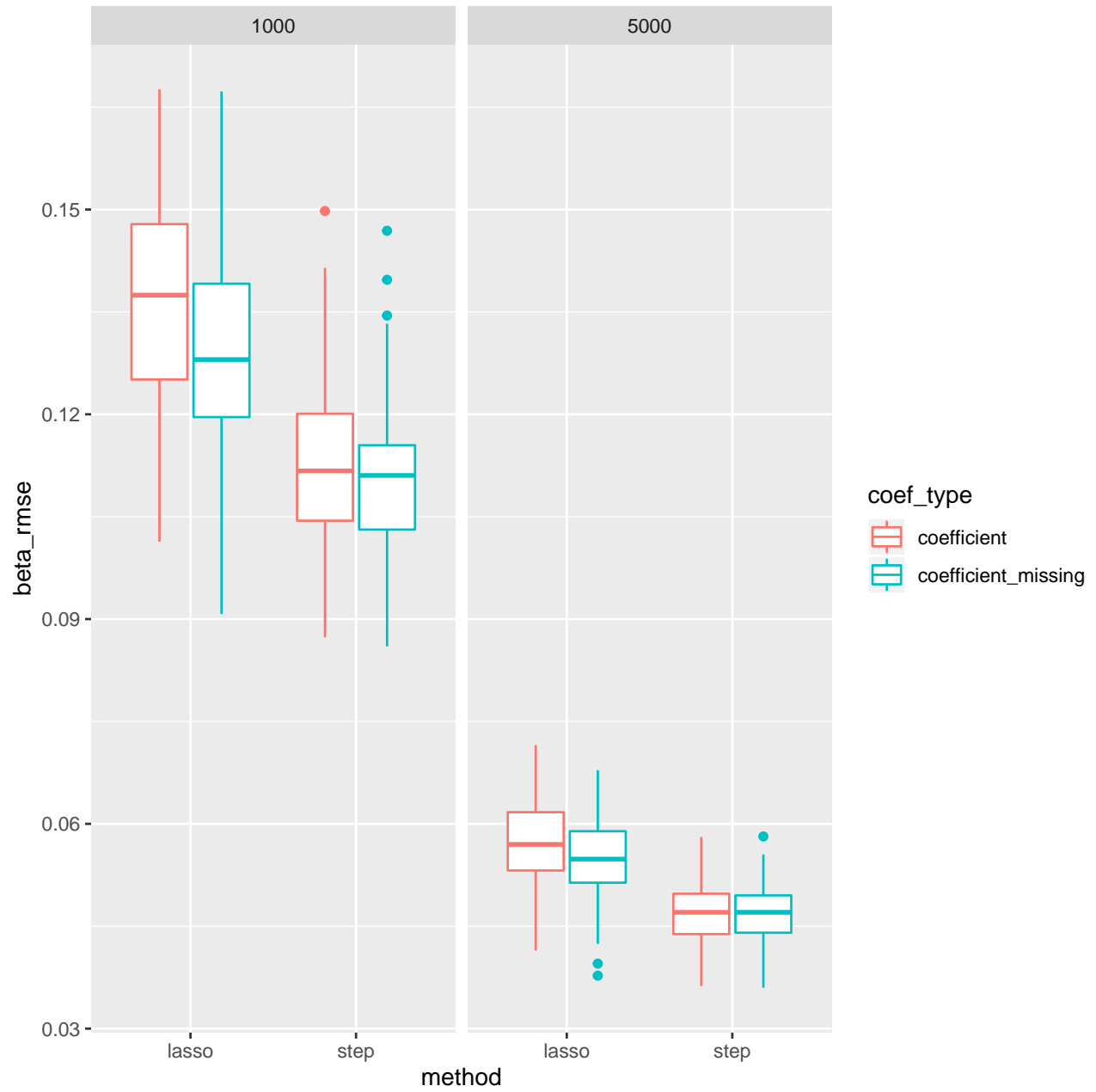


Figure 10: Beta RMSE comparison of with/without weak predictors

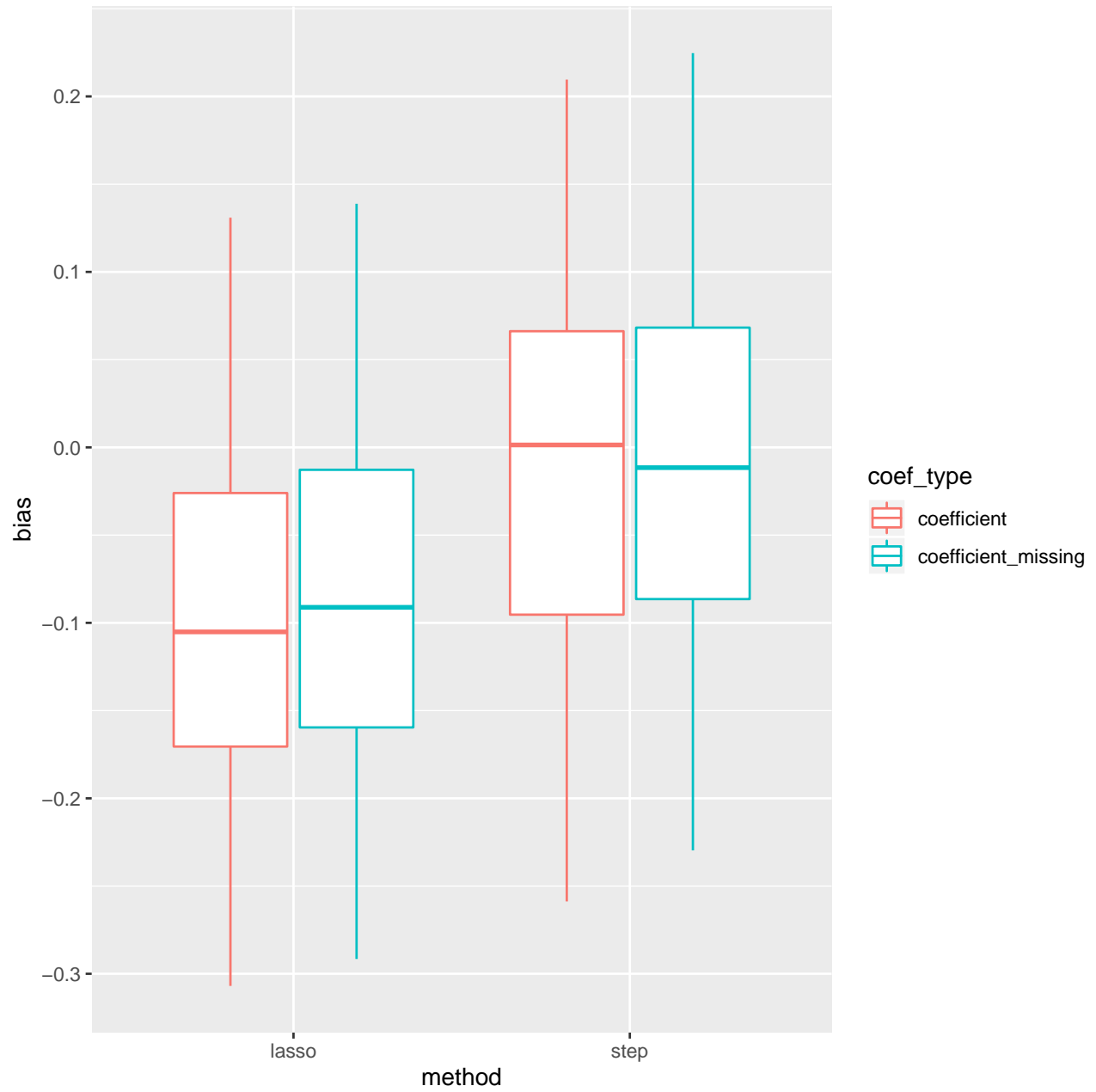


Figure 11: Beta bias comparison of with/without weak predictors