



OPEN

Relationship between gene regulation network structure and prediction accuracy in high dimensional regression

Yuichi Okinaga¹, Daisuke Kyogoku², Satoshi Kondo³, Atsushi J. Nagano^{4,5}✉ & Kei Hirose^{6,7}✉

The least absolute shrinkage and selection operator (lasso) and principal component regression (PCR) are popular methods of estimating traits from high-dimensional omics data, such as transcriptomes. The prediction accuracy of these estimation methods is highly dependent on the covariance structure, which is characterized by gene regulation networks. However, the manner in which the structure of a gene regulation network together with the sample size affects prediction accuracy has not yet been sufficiently investigated. In this study, Monte Carlo simulations are conducted to investigate the prediction accuracy for several network structures under various sample sizes. When the gene regulation network is a random graph, a sufficiently large number of observations are required to ensure good prediction accuracy with the lasso. The PCR provided poor prediction accuracy regardless of the sample size. However, a real gene regulation network is likely to exhibit a scale-free structure. In such cases, the simulation indicates that a relatively small number of observations, such as $N = 300$, is sufficient to allow the accurate prediction of traits from a transcriptome with the lasso.

Technological advancements have enabled the collection of highly multidimensional data from biological systems^{1–4}. For example, RNA sequencing quantifies expression levels of thousands of genes. Such omics data is useful in predicting organismal traits, with empirical applications including diagnosis and classification of diseases and prediction of patient survival^{5–8} and possible future applications in predicting crop yields⁹, insecticide resistance¹⁰, and environmental adaptation¹¹.

A common challenge in predicting traits from omics data is the dimension of the data far exceeding that of the sample size (known as high-dimensional regression). For example, if one is to apply least-squares estimation in multiple regression (e.g. $\text{trait} \approx \beta_0 + \beta_1 \text{gene}_1 + \beta_2 \text{gene}_2 + \dots$) to predict a trait value from a transcriptome, the sample size needs to be (at least) larger than the number of model parameters. However, because transcriptome studies typically observe thousands of genes, a sample size exceeding the number of genes is not realistic at present. In this case, high-dimensional regression modeling must be considered.

The least absolute shrinkage and selection operator (lasso¹²) is one of the most frequently used methods for high-dimensional regression. It simultaneously achieves variable selection and parameter estimation. Theoretically, the prediction accuracy of the lasso is highly dependent on the correlation structure among exploratory variables; it is high under certain strong conditions, such as the compatibility condition¹³. However, in practice, it is not easy to check whether the compatibility condition holds. Another popular estimation method for high-dimensional regression is principal component regression (PCR¹⁴). PCR is a two-stage procedure: first, principal component analysis is conducted for predictors, following which the regression model on which the principal components are used as predictors is fitted. This method may perform well when the exploratory variables are highly correlated.

It is reasonable to assume that gene regulation networks will result in conditional independence among the levels of gene expression^{15–17}. Here, two variables are conditionally independent when they are independent given

¹Graduate School of Mathematics, Kyushu University, 744 Motoooka, Fukuoka 819-0395, Japan. ²The Museum of Nature and Human Activities, 6 Yayoigaoka, Sanda, Hyogo 669-1546, Japan. ³Agriculture and Biotechnology Business Division, Toyota Motor Corporation, Miyoshi, Aichi 470-0201, Japan. ⁴Faculty of Agriculture, Ryukoku University, Otsu, Shiga 520-2194, Japan. ⁵Institute for Advanced Biosciences, Keio University, Tsuruoka, Yamagata 997-0017, Japan. ⁶Institute of Mathematics for Industry, Kyushu University, 744 Motoooka, Fukuoka 819-0395, Japan. ⁷RIKEN Center for Advanced Intelligence Project, 1-4-1 Nihonbashi, Chuo-ku, Tokyo 103-0027, Japan. ✉email: anagano1234@gmail.com; hirose@imi.kyushu-u.ac.jp

other variables (e.g. two focal variables are independently influenced by a third variable¹⁸). When a random vector of exploratory variables follows a multivariate normal distribution, two variables are conditionally independent if and only if the corresponding element of the inverse covariance matrix is zero. Essentially, the networks are characterized by the nonzero pattern of the inverse covariance matrix.

One of the most notable characteristics of biological networks is their scale-free nature, that is, the degree distribution of the network follows a power-law expressed as $p(x) \propto x^{-\gamma}$ ($\gamma > 1$)^{19,20}. Empirical studies suggest that biological networks are often scale-free^{21–23}, although exceptions have also been found²⁴. Therefore, it is reasonable to consider the problem of high-dimensional regression when the networks of exploratory variables are scale-free. Here, it should be noted that the relative performance of different high-dimensional regression techniques may depend on sample sizes. However, to the best of our knowledge, the effect of the gene regulation network structure together with sample size on prediction accuracy has not yet been sufficiently investigated.

This paper provides a general simulation framework to study the effects of correlation structure in explanatory variables. As an example, the prediction of ambient temperature from the transcriptome, for which good empirical data is available^{11,25}, is considered. It should be noted that the implementation of the proposed procedure is independent of the empirical data^{11,25}; the proposed framework may be applied to predict any consequence of gene expression differences. The proposed framework is based on the Monte Carlo simulations. Three datasets of transcriptome and their traits are generated. The datasets are characterized by the covariance structure of exploratory variables; one of the covariance structures corresponds to the scale-free gene regulation network. Both lasso and PCR are applied to these simulated datasets to investigate the prediction accuracy with different types of gene regulation networks. The sample size is also varied to examine its effect on the prediction accuracy.

The remainder of this paper is organized as follows. Section “[Prediction methods for high-dimensional data](#)” describes prediction methods for high-dimensional regression in the given simulation. Section “[Simulation framework](#)” discusses the proposed simulation framework. Finally, Section “[Concluding remarks](#)” presents the concluding remarks.

Prediction methods for high-dimensional data

Suppose that we have n observations $\{(x_i, y_i) \mid i = 1, \dots, n\}$, where x_i are p -dimensional vector of explanatory variables and y_i are responses ($i = 1, \dots, n$). Let $X = (x_1, \dots, x_n)^T$ and $Y = (y_1, \dots, y_n)^T$. Consider the linear regression model:

$$Y = X\beta + \epsilon,$$

where $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T$ is a vector of error variables with $E(\epsilon) = \mathbf{0}$ and $V(\epsilon) = \sigma^2 I_n$.

Lasso. The lasso minimizes a loss function that consists of quadratic loss with a penalty based on an L_1 norm of a parameter vector:

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1, \quad (1)$$

where $\lambda > 0$ is a regularization parameter. Because of the nature of the L_1 norm in the penalty term, some of the elements of the coefficients are estimated to be exactly zero. Thus, variable selection is conducted, and only variables that correspond to nonzero coefficients affect the responses.

PCR. In some cases, the first few largest eigenvalues of the covariance matrix of predictors (i.e., proportional contributions of principle components) can be considerably large (e.g., spiked covariance model²⁶). In such a case, the lasso may not function effectively in terms of both prediction accuracy and consistency in model selection, because the conditions for its effective performance (e.g., compatibility condition²⁷) may not be satisfied. This issue could be addressed using PCR because it transforms data with a large number of highly correlated variables into a few principal components. In the first stage of PCR, principal component analysis is applied to predictors. The i th observation of predictor, x_i , is linearly mapped onto a d ($< p$)-dimensional vector, $z_i = A^T x_i$, where A is a $p \times d$ matrix. The matrix A is obtained by the following least squares optimization problem²⁸:

$$A = \arg \min_A \sum_{i=1}^n \|(x_i - \bar{x}) - AA^T(x_i - \bar{x})\|_2^2 \quad \text{subject to} \quad A^T A = I_d.$$

here, \bar{x} is the sample mean vector, that is, $\bar{x} = \sum_{i=1}^n x_i/n$. In this work, the number of projected dimensions, d , was chosen such that d principle components collectively explain 90% or more variance (and $d - 1$ principle components do not). Then, in the second stage, regression analysis is conducted, for which the principal components, $\{z_1, \dots, z_n\}$, are used as predictors. Here, the regression coefficients in the second stage are estimated by the lasso.

Simulation framework

An overview of the simulation is presented in Fig. 1. First, the model that defines the relationship between the trait and the levels of gene expression was parameterized. This was done using the empirical data¹¹, which quantified the transcriptome of wild *Arabidopsis halleri* subsp. *gemmifera* weekly for two years in their natural habitat as well as bihourly on the equinoxes and solstices ($p = 17,205$ genes for $n = 835$ observations). Three types of simulated data were generated using different covariance matrices of genes, denoted as Σ_j ($j = 1, 2, 3$). Σ_1 is the sample covariance matrix of genes. Generally, none of the elements of the inverse of sample covariance matrix are exactly zero, implying that each gene interacts with all the other genes. Such a fully connected

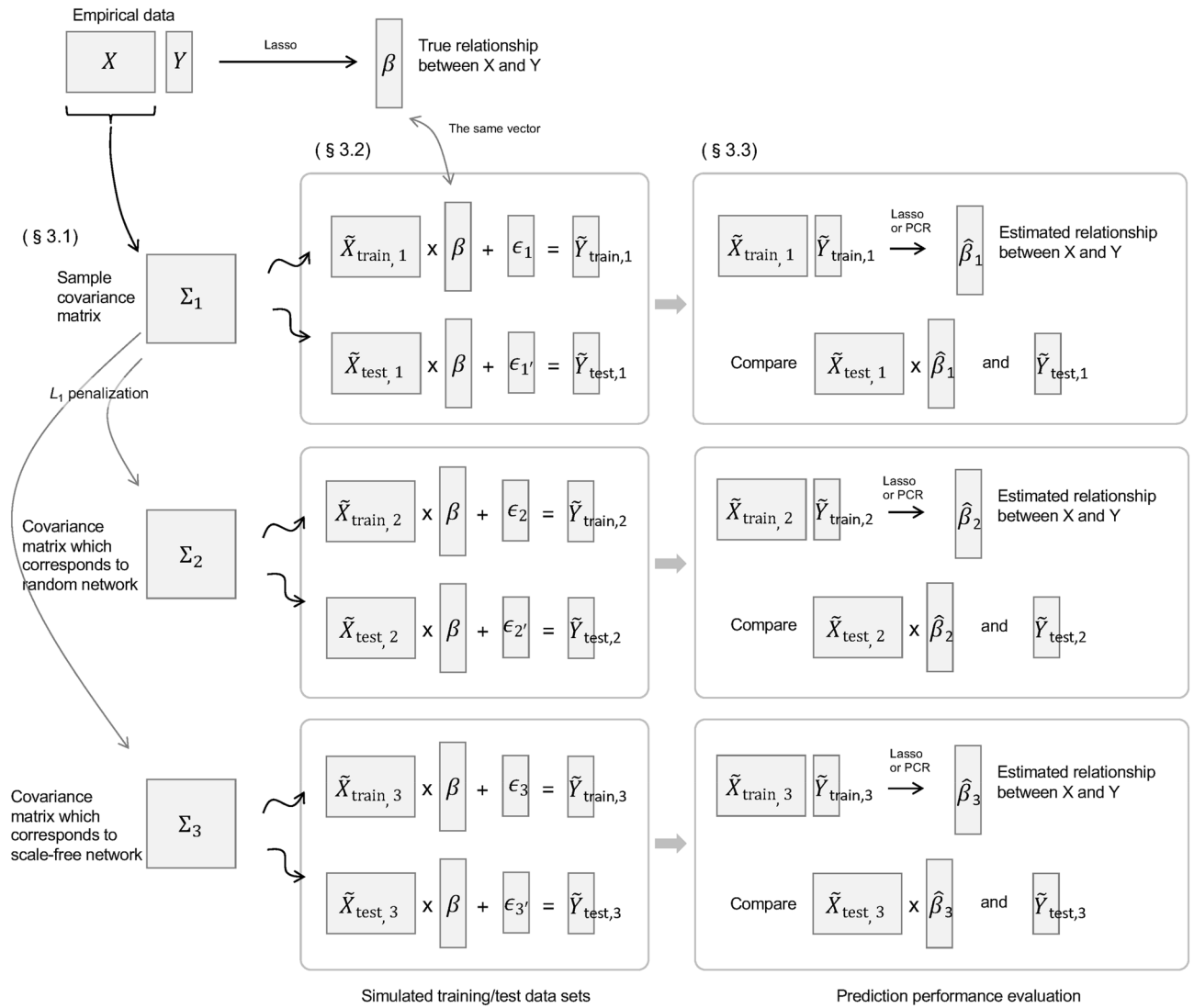


Figure 1. Overview of the simulation.

network is ineffective in terms of interpretation of the mechanism of gene regulation. Thus, two other covariance matrices were produced to simulate sparse networks based on the sample covariance matrix Σ_1 . Σ_2 is generated by the graphical lasso²⁹, which corresponds to the random graph. Although the graphical lasso is widely used because of its computational efficiency, real networks are often scale-free. Therefore, Σ_3 , which corresponds to the scale-free network, was generated here. The estimation of scale-free networks is achieved by the reweighted graphical lasso³⁰. Based on these three covariance matrices Σ_j ($j = 1, 2, 3$), the simulated transcriptome data were generated from the multivariate normal distribution. The simulated trait data were generated from simulated transcriptome data. Finally, lasso and PCR were applied to these simulated data to compare their prediction accuracies. The sample sizes of the simulated data were varied to investigate the relationship between prediction accuracy and sample sizes.

Evaluation of the estimation procedure. The performance of the estimation procedure is investigated by the following expected prediction error:

$$E \left[\left\| \mathbf{Y}^* - (\mathbf{X}^*)^T \hat{\boldsymbol{\beta}} \right\|_2^2 \right],$$

where \mathbf{X}^* and \mathbf{Y}^* follow $\mathbf{X}^* \sim N(\mathbf{0}, \Sigma_j)$ ($j = 1, 2, \text{ or } 3$) and $\mathbf{Y}^* \sim N((\mathbf{X}^*)^T \boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$, respectively. The estimator $\hat{\boldsymbol{\beta}}$ is obtained using current observations, while \mathbf{X}^* and \mathbf{Y}^* correspond to future observations. The Σ_j ($j = 1, 2, 3$), $\boldsymbol{\beta}$, and σ^2 are true values but unknown. In practice, these parameters are defined by using the actual dataset, (\mathbf{X}, \mathbf{Y}) . Detail of setting of these parameters will be presented in the next subsection.

To estimate the expected prediction error, the Monte Carlo simulation is conducted. We first randomly generate training and test data, $(\tilde{\mathbf{X}}_{train}, \tilde{\mathbf{Y}}_{train})$ and $(\tilde{\mathbf{X}}_{test}, \tilde{\mathbf{Y}}_{test})$, respectively. Here, $\tilde{\mathbf{X}}_{train}$ follows a multivariate normal distribution with mean vector μ_X and variance-covariance matrix Σ_j , where μ_X is the sample mean of X . Then,

\tilde{Y}_{train} is generated by $\tilde{Y}_{train} = \tilde{X}_{train}\beta + \epsilon$, where ϵ is a random sample from $N(0, \sigma^2 I)$ with I being an identity matrix. The test data, $(\tilde{X}_{test}, \tilde{Y}_{test})$, are generated by the same procedure as $(\tilde{X}_{train}, \tilde{Y}_{train})$ but independent of $(\tilde{X}_{train}, \tilde{Y}_{train})$. The number of observations for the training and test data are N ($N = 50, 100, 200, 300, 500, 1000$) and 1000, respectively. The lasso and the PCR are performed with training data $(\tilde{X}_{train}, \tilde{Y}_{train})$, following which RMSE is calculated in (10). The above process, from random generation of data to RMSE calculation, was performed 100 times.

Parameter setting. Covariance structures. Here, the characterization of the network structure of predictors by conditional independence is considered. When the predictors follow a multivariate normal distribution, the network structure based on the conditional independence corresponds to the nonzero pattern of the inverse covariance (precision) matrix. In other words, the network structure is characterized by the inverse covariance matrix of predictors.

Let S be the sample covariance matrix of predictors, that is, $S = \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T / n$. Let $\Omega_j = \Sigma_j^{-1}$ ($j = 1, 2, 3$). Σ_1 is a ridge estimator of the sample variance-covariance matrix, that is, $\Sigma_1 = S + \delta I$. Here δ is a small positive value (in this simulation, $\delta = 10^{-5}$). The term δI allows the existence of Ω_1 . Note that because Ω_1 is not sparse, it leads to the complete graph, which is of no use in interpreting gene regulatory networks. To generate a covariance matrix whose inverse matrix is sparse, L_1 penalization is employed for the estimation of Ω_2 and Ω_3 as follows:

$$\hat{\Omega}_j = \arg \max_{\Omega} \{ \log |\Omega| - \text{tr}(\Omega S) - P_j(\Omega) \} \quad (j = 2, 3), \tag{2}$$

where $P_j(\Omega)$ ($j = 2, 3$) are penalty terms which enhance the sparsity of the inverse covariance matrix. To estimate the sparse inverse covariance matrix, the lasso penalty is typically used as follows:

$$P_2(\Omega) = \rho \sum_{i=1}^p \|\omega_{(-i, \cdot)}\|_1, \tag{3}$$

where $\omega_{(-i, \cdot)} = (\omega_{i1}, \omega_{i2}, \dots, \omega_{i(i-1)}, \omega_{i(i+1)}, \dots, \omega_{ip})^T \in \mathbb{R}^{p-1}$. The problem (3) is referred to as the graphical lasso²⁹, and there exists several efficient algorithms to obtain the solution^{31–33}. The estimator of (2) with (3) corresponds to Ω_2 and $\Sigma_2 = \Omega_2^{-1}$.

The lasso penalty (3) does not enhance scale-free networks. It penalizes all edges equally so that the estimated graph is likely to be a random graph, that is, the degree distribution becomes a binomial distribution. To enhance scale-free networks (i.e., power-law distribution), the log penalty³⁰ is used as follows:

$$P_3(\Omega) = \frac{\rho}{2} \sum_{i=1}^p \{ \log (\|\omega_{(-i, \cdot)}\|_1 + a_i) + \log (\|\omega_{(\cdot, -i)}\|_1 + a_i) \}, \tag{4}$$

where $\omega_{(\cdot, -i)} = (\omega_{1i}, \omega_{2i}, \dots, \omega_{(i-1)i}, \omega_{(i+1)i}, \dots, \omega_{pi})^T$ and $a_i > 0$ are tuning parameters. We note that the penalty (4) is slightly different from original definition³⁰, expressed as

$$P(\Omega) = \rho \sum_{i=1}^p \log (\|\omega_{(-i, \cdot)}\|_1 + a_i). \tag{5}$$

When we do not assume that $\omega_{ij} = \omega_{ji}$, the estimate of the inverse covariance matrix with (5) is not symmetric. Since the original graphical lasso algorithm does not assume that $\omega_{ij} = \omega_{ji}$ ^{31,34}, we slightly modify the penalty as in Eq. (4). Notably, $P_3(\Omega)$ in (4) coincides with (5) when $\omega_{ij} = \omega_{ji}$. From a Bayesian viewpoint, the prior distribution which corresponds to the log penalty becomes the power-law distribution³⁰; thus, the penalty (4) is likely to estimate the scale-free networks. The estimator of (2) with (4) corresponds to Ω_3 .

Because the log-penalty (4) is nonconvex, it is not easy to directly optimize (2). To implement the maximization problem (2), the minorize-maximization (MM) algorithm³⁵ has been constructed³⁰, in which the weighted lasso penalty $P_M^{(t)}(\Omega)$ with current parameter $\Omega_3^{(t)}$ is used:

$$P_M^{(t)}(\Omega) = \sum_{i=1}^p \sum_{j \neq i} \rho_{ij}^{(t)} |\omega_{ij}|, \tag{6}$$

where $\rho_{ij}^{(t)}$ are the weights

$$\rho_{ij}^{(t)} = \frac{1}{2} \left(\frac{\rho}{\|\omega_{(-i, \cdot)}^{(t)}\|_1 + a_i} + \frac{\rho}{\|\omega_{(\cdot, -j)}^{(t)}\|_1 + a_j} \right). \tag{7}$$

In general, $\hat{\Omega}$ must be symmetric, so that Eq. (7) can be expressed as

$$\rho_{ij}^{(t)} = \frac{1}{2} \left(\frac{\rho}{\|\omega_{(-i, \cdot)}^{(t)}\|_1 + a_i} + \frac{\rho}{\|\omega_{(-j, \cdot)}^{(t)}\|_1 + a_j} \right). \tag{8}$$

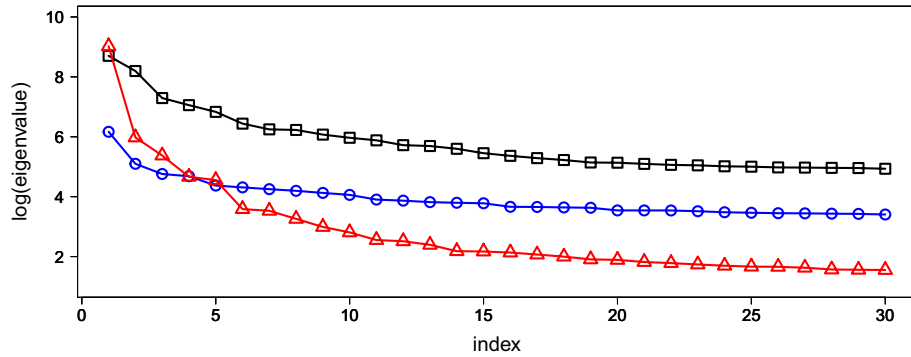


Figure 2. Logarithm graph of the largest 30 eigenvalues of Σ_1 (square), Σ_2 (circle) and Σ_3 (triangle). The horizontal axis expresses the index of eigenvalues arranged in descending order.

Because the weighted graphical lasso can be implemented by a standard graphical lasso algorithm, the estimator is obtained as the following algorithm.

1. Set $t = 0$. Get $\Omega_3^{(0)}$ via ordinary graphical lasso. Repeat 2 to 4 until convergence.
2. Update weights $\rho_{ij}^{(t)}$ using (7).
3. Get $\Omega_3^{(t+1)}$ via the weighted graphical lasso (2) with penalty (6).
4. $t \leftarrow t + 1$.

To obtain Σ_2 and Σ_3 , the tuning parameters a_i ($i = 1 \dots, p$) and ρ must be determined. Following the experiments³⁰, $a_i = 1$ was set for $i = 1 \dots, p$. To select the value of the regularization parameter ρ , several candidates were first prepared. In our simulation, the candidates were $\rho = 0.3, 0.4, 0.5, 0.6, 0.7$. From these, the value of ρ was selected such that the extended Bayesian information criterion (EBIC^{36,37})

$$EBIC = -n\{\log|\Omega_2| - \text{tr}(\Omega_2 S)\} + q \log n + 4q\delta \log p \tag{9}$$

was minimized. Here, q is the number of nonzero parameters of the upper triangular matrix of $\hat{\Omega}$, and $\delta \in [0, 1)$ is a tuning parameter. As the value of δ increases, a sparser graph is generated. Note that $\delta = 0$ corresponds to the ordinary BIC³⁸. We set $\delta = 0.5$ because $\delta = 0.5$ is shown to yield good performance in both simulated and real data analyses³⁷. As a result, the EBIC selected $\rho = 0.5$.

The upper triangular matrix Ω_3 must be estimated with the reweighted graphical lasso problem. A value of $p = 17205$ results in $p(p + 1)/2 \approx 148$ million parameters. As a result, with the machine used in this study (Intel Core Xeon 3 GHz, 128 GB memory), it would take several days to conduct the reweighted graphical lasso approach, even with a small number of iterations such as $T = 5$. For this reason, $T = 5$ iterations were employed to produce Σ_3 here. Finally, Σ_2 and Σ_3 were scaled such that their signal-to-noise ratio became Σ_1 .

Figure 2 depicts the logarithm of the largest 30 eigenvalues of Σ_j ($j = 1, 2, 3$). The first few largest eigenvalues of Σ_3 are significantly larger than those of Σ_2 , implying that the scale-free networks tend to produce predictors with large correlations.

Regression parameters. The values of β and σ^2 are determined as follows. First, 10-fold cross-validation is performed as described below, and the regularization parameter λ in (1) is selected. The data (X, Y) are divided into ten datasets, $(X^{(j)}, Y^{(j)})$ ($j = 1, \dots, 10$), which consist of almost equal sample sizes. Let $X^{(-j)} = (X^{(1)}, \dots, X^{(j-1)}, X^{(j+1)}, \dots, X^{(10)})$, and $Y^{(-j)} = (Y^{(1)}, \dots, Y^{(j-1)}, Y^{(j+1)}, \dots, Y^{(10)})$ ($j = 1, \dots, 10$). For each j ($j = 1, \dots, 10$), the training and test data are defined by $(X^{(-j)}, Y^{(-j)})$ and $(X^{(j)}, Y^{(j)})$, respectively. Then, the parameter $\hat{\beta}^{(j)}$ ($j = 1, \dots, 10$) is found by the lasso:

$$\hat{\beta}^{(j)} = \arg \min_{\beta} \left(\|Y^{(-j)} - X^{(-j)}\beta\|_2^2 + \lambda \|\beta\|_1 \right).$$

For each j ($j = 1, \dots, 10$), the verification error is calculated as follows:

$$CV^{(j)} = \frac{1}{\#Y^{(j)}} \|Y^{(j)} - X^{(j)}\hat{\beta}^{(j)}\|_2^2.$$

Then, λ is adopted such that it minimizes $CV = \frac{1}{10} \sum_{j=1}^{10} CV^{(j)}$, the mean of $CV^{(j)}$. Following this, the dataset (X, Y) is again randomly divided into two datasets: test data (X_{test}, Y_{test}) and training data (X_{train}, Y_{train}) . Lasso estimation (1) is performed using the training data, with λ obtained by the above 10-fold cross-validation. Then, β is defined as the lasso estimator, resulting in the number of nonzero parameters of β being 259. Figure 3 shows the histogram of nonzero parameters of β . It is seen that the majority of the nonzero coefficients were close to zero; only 15 parameters had absolute values larger than 0.1.

In addition, the root mean squared error (RMSE) is calculated as follows:

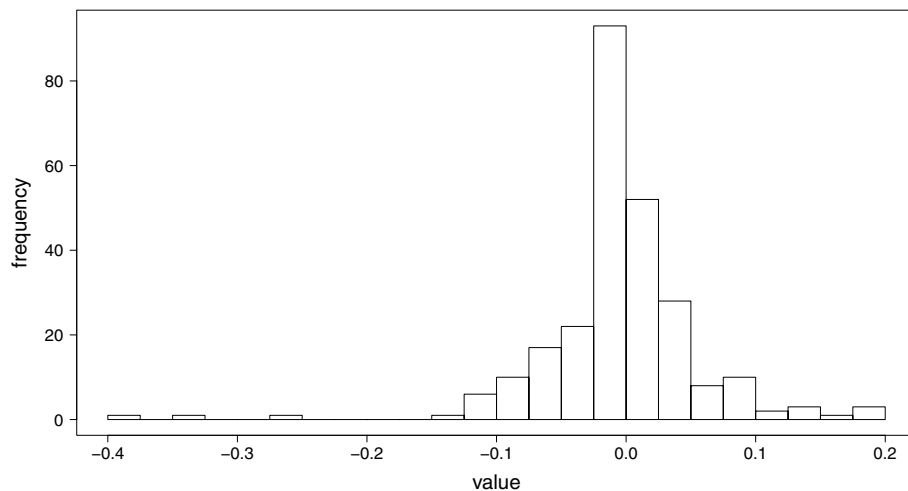


Figure 3. Histogram of 259 nonzero parameters of β .

$$\text{RMSE} = \frac{1}{\sqrt{\#Y_{\text{test}}}} \|Y_{\text{test}} - X_{\text{test}}\hat{\beta}\|_2, \quad (10)$$

and the variance of errors, σ^2 , is defined by $\sigma^2 = (\text{RMSE})^2$.

Results

The box and whisker plot of the RMSE and the coefficient of determination (R^2) are illustrated in Figs. 4 and 5. The horizontal axis is N (the number of observations of training data) and the vertical axis is the RMSE or R^2 based on 1000 observations of test data.

We compared the performance of the lasso with that of the PCR. When Σ_1 and Σ_3 were used, the PCR performed worse than the lasso for small sample sizes. For Σ_2 , the prediction performance with PCR was unsatisfactory even when the sample size N increased. The poor performance of the PCR can be attributed to the predictors associated with small eigenvalues; these predictors affected the prediction performance. Figure 6 depicts a scatter plot of nonzero elements of β and the eigenvector for the maximum eigenvalue of Σ_2 . As can be seen, only a significant amount of correlation existed; in fact, the correlation coefficient was only 0.068.

The prediction accuracy was compared among the three covariance structures. In all the cases except PCR with Σ_2 , the values of RMSE decreased and R^2 increased with the increase in the value of N . Further, R^2 was unstable for small sample sizes for all the cases when the lasso was applied. For large sample sizes, the R^2 of Σ_1 was better than that of Σ_2 and Σ_3 . As described before, Σ_1 was the sample covariance matrix, while Σ_3 (and Σ_2) was estimated using the graphical lasso. As the lasso-type regularization methods shrink parameters toward zero, the correlations among the exploratory variables reduce when the graphical lasso is used. Therefore, Σ_2 and Σ_3 resulted in smaller correlations as compared to Σ_1 . Consequently, the R^2 may increase with stronger correlations. We compared the RMSE results of Σ_2 and Σ_3 . With Σ_2 , we found that a sufficiently large number of observations is required to yield a small RMSE with the lasso. Meanwhile, Σ_3 resulted in a small RMSE with a relatively small number of observations, such as $N = 300$.

Code availability. The proposed simulation is implemented in R package `simrnet`, which is available at <https://github.com/keihirose/simrnet>. Below is a sample code of the `simrnet` in R:

```
library(devtools)
install_github("keihirose/simrnet") #install package
library(simrnet) #load package
data(nagano2019)
attach(nagano2019)
rho <- (1:9) / 10 #tuning parameters for glasso
pars <- genpar(X,Y,rho) #set true parameter
result <- simrnet(pars, times.sim=100) #conduct simulation
plot(result)
```

When $p = 100$, it took less than 12 min to conduct the simulation with 100 replications using the machine employed herein (Intel Core Xeon 3 GHz, 128 GB memory). For high-dimensional data such as $p = 17,205$, which was used in the simulation presented in this paper, several days were required to complete the simulation task.

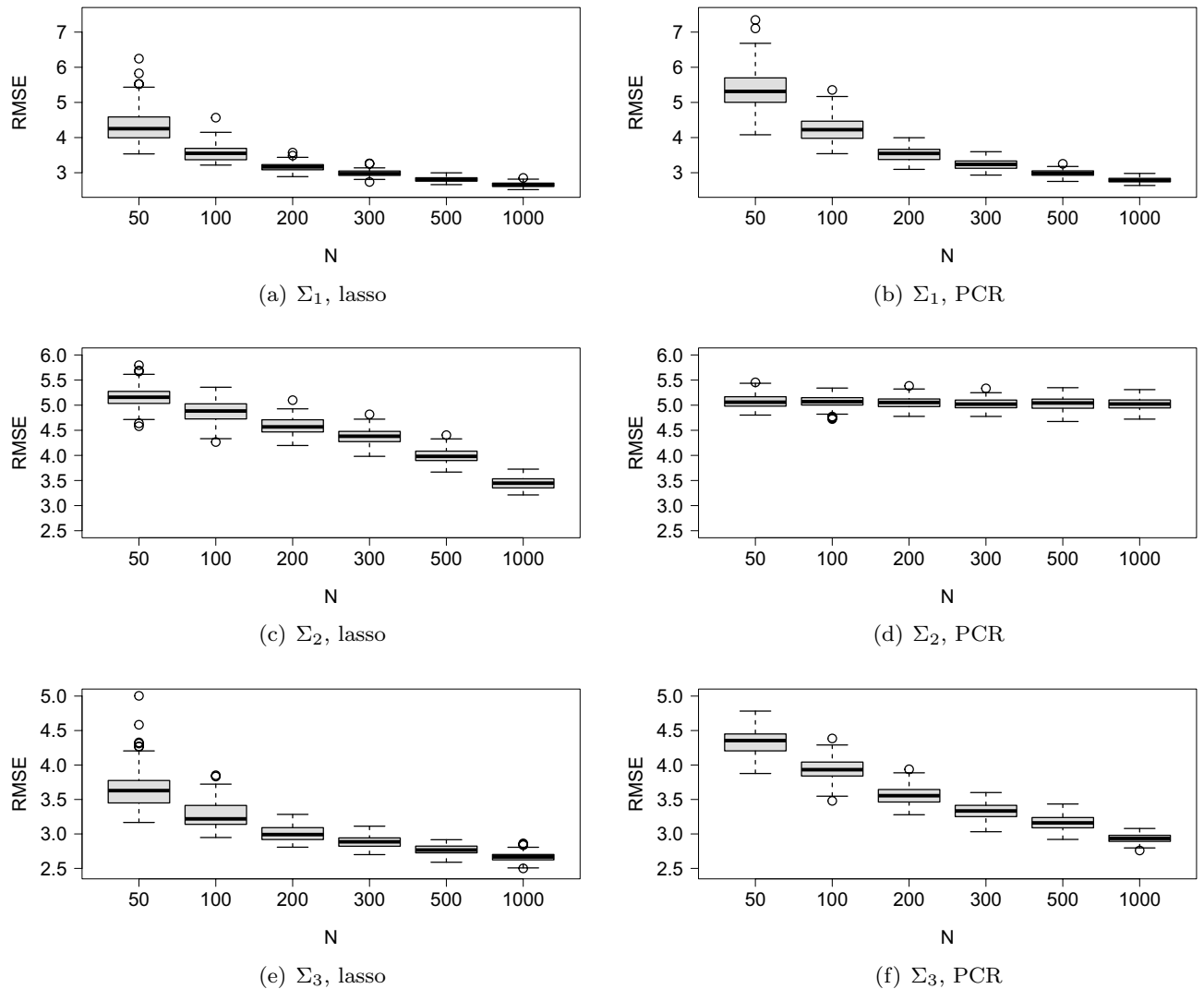


Figure 4. Box and whisker plot of RMSE. The variance-covariance matrix used in the simulations is Σ_1 in (a, b), Σ_2 in (c, d), and Σ_3 in (e, f). The regression model is estimated by the lasso in (a), (c), and (e) and by PCR in (b), (d), and (f).

Concluding remarks

In a gene regulation network, a gene regulates a small portion of a genome, not all the genes in a genome. This indicates that gene regulation network is expected to be a sparse network rather than a complete graph. Therefore, two covariance matrices indicating sparse networks (Σ_2 , Σ_3) were prepared in addition to a covariance matrix derived from empirical data (Σ_1). Generally, although hundreds of genes contribute to defining a trait, their contributions are not equal. It is frequently observed that genes regulating a trait include a few large-effect genes and many small-effect genes. This property was reflected in the distribution of β (Fig. 3). We considered the case where a limited number of regression coefficients significantly contributed to the definition of a trait. The Monte Carlo simulation result indicated that regardless of the network structure, the number of observations should be greater than at least 200 to accurately predict traits from a transcriptome (Σ_1 , Σ_3 , Figs. 4 and 5). We also found that the lasso generally provided better accuracy than the PCR. In particular, when the gene regulation network was random (Σ_2), the prediction accuracy of the PCR was poor even if the sample size increased. In conclusion, it is important to sufficiently secure large sample sizes when performing regression analysis of data that exhibits either the random graph and the scale-free network. Additionally, we concluded that the lasso would be preferable to the PCR to ensure a good prediction accuracy.

Conventional theory on the relationship between RMSE and sample size has been developed under the assumption that the sample size exceeds the number of exploratory variables³⁹. However, omics data, which is rapidly being accumulated, results in high dimensional data with strong correlations. Thus, our simulation study considered more complicated settings than the traditional ones. Our simulation, or its extension, may be used in the future to find clues about theoretical aspects that may ultimately lead to the development of a sample size determination technique for omics data.

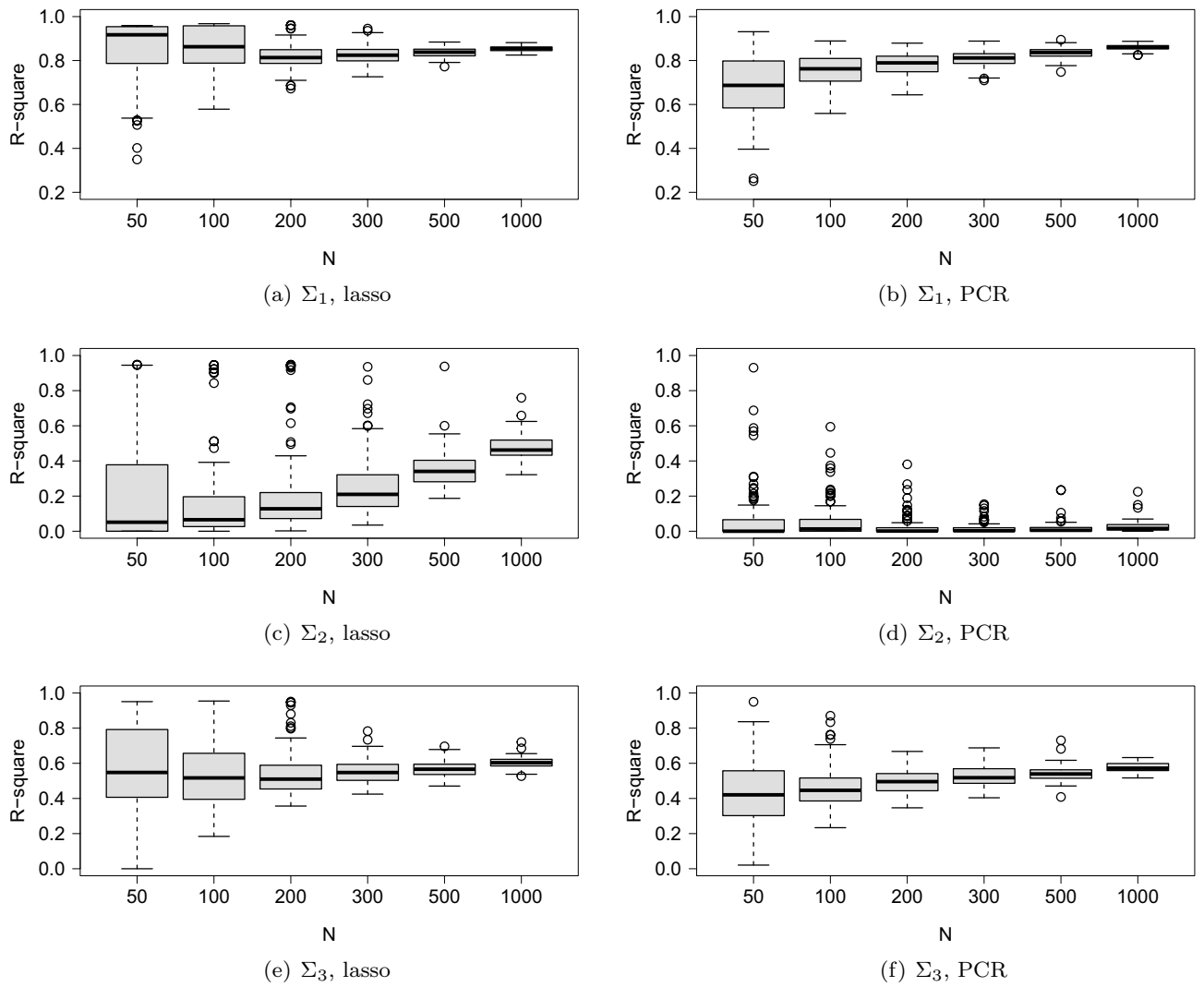


Figure 5. Box and whisker plot of R^2 . The variance-covariance matrix used in the simulations is Σ_1 in (a, b), Σ_2 in (c, d), and Σ_3 in (e, f). The regression model is estimated by the lasso in (a), (c), and (e) and by PCR in (b), (d), and (f).

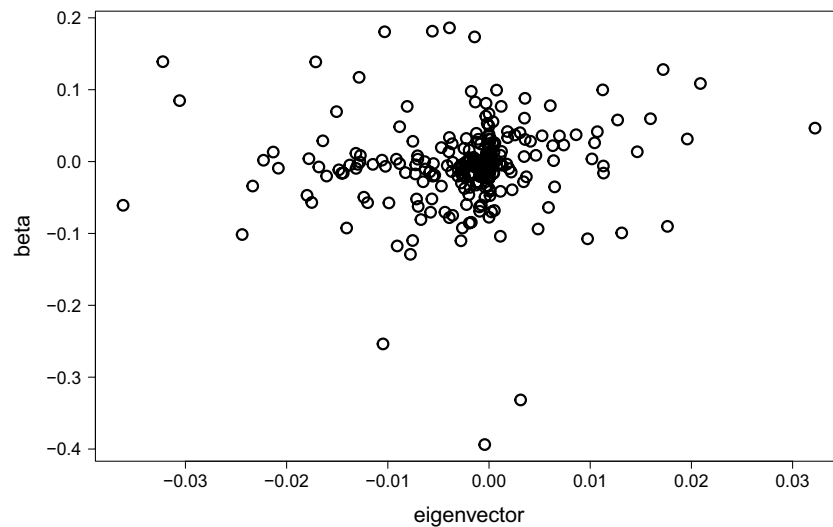


Figure 6. Scatter plot of β and the eigenvector corresponding to the maximum eigenvalue of Σ_2 . The nonzero elements of β are not drawn.

Other than the scale-free network, the small-world network is another notable property in the networks literature⁴⁰. The definition of the small-world networks is that the shortest path length between two randomly chosen variables is proportional to $\log p$; that is, it is considerably small compared with the network size. The small-world networks have been investigated in various fields of research, including the biology^{41–43}. Some statistical properties of the small-world networks have also been studied^{44–46}. The investigation of the prediction accuracy in the small-world networks would be interesting but beyond the scope of this research. We would like to take this as a future research topic. The development of methods that provides better prediction accuracy than the lasso in various network structures with small sample sizes would also be an important future research topic.

Received: 6 January 2021; Accepted: 17 May 2021

Published online: 01 June 2021

References

- Gehlenborg, N. *et al.* Visualization of omics data for systems biology. *Nat. methods* **7**, S56 (2010).
- Mochida, K. & Shinozaki, K. Advances in omics and bioinformatics tools for systems analyses of plant functions. *Plant Cell Physiol.* **52**, 2017–2038 (2011).
- Li, Z. & Sillanpää, M. J. Overview of lasso-related penalized regression methods for quantitative trait mapping and genomic selection. *Theor. Appl. Genet.* **125**, 419–435 (2012).
- Hasin, Y., Seldin, M. & Lusis, A. Multi-omics approaches to disease. *Genome Biol.* **18**, 83 (2017).
- van't Veer, L. J. *et al.* Gene expression profiling predicts clinical outcome of breast cancer. *Nature* **415**, 530–536 (2002).
- Bøvelstad, H. M. *et al.* Predicting survival from microarray data—A comparative study. *Bioinformatics* **23**, 2080–2087 (2007).
- Chan, A. W. *et al.* 1H-NMR urinary metabolomic profiling for diagnosis of gastric cancer. *Br. J. Cancer* **114**, 59–62 (2016).
- Nandagopal, V., Geeitha, S., Kumar, K. V. & Anbarasi, J. Feasible analysis of gene expression—A computational based classification for breast cancer. *Measurement* **140**, 120–125 (2019).
- Kremling, K. A. *et al.* Dysregulation of expression correlates with rare-allele burden and fitness loss in maize. *Nature* **555**, 520–523 (2018).
- Dermauw, W. *et al.* A link between host plant adaptation and pesticide resistance in the polyphagous spider mite tetranychus urticae. *Proc. Natl. Acad. Sci.* **110**, E113–E122 (2013).
- Nagano, A. J. *et al.* Annual transcriptome dynamics in natural environments reveals plant seasonal adaptation. *Nat. Plants* **5**, 74–83 (2019).
- Tibshirani, R. Regression shrinkage and selection via the Lasso. *J. R. Stat. Soc. Ser. B Methodol.* **58**, 267–288 (1996).
- van de Geer, S. A. & Bühlmann, P. On the conditions used to prove oracle results for the lasso. *Electron. J. Stat.* **3**, 1360–1392. <https://doi.org/10.1214/09-EJS506> (2009).
- Jolliffe, I. T. Principal components in regression analysis. in *Principal Component Analysis*, 129–155 (Springer, 1986).
- Wei, Z. & Li, H. A markov random field model for network-based analysis of genomic data. *Bioinformatics* **23**, 1537–1544 (2007).
- Dobra, A. *et al.* Sparse graphical models for exploring gene expression data. *J. Multivar. Anal.* **90**, 196–212 (2004).
- Yu, D., Kim, M., Xiao, G. & Hwang, T. H. Review of biological network data and its applications. *Genom. Inform.* **11**, 200 (2013).
- Wille, A. & Bühlmann, P. Low-order conditional independence graphs for inferring genetic networks. *Stat. Appl. Genet. Mol. Biol.* **5** (2006).
- Barabási, A.-L. & Albert, R. Emergence of scaling in random networks. *Science* **286**, 509–512 (1999).
- Milo, R. *et al.* Network motifs: simple building blocks of complex networks. *Science* **298**, 824–827 (2002).
- Barabasi, A.-L. & Oltvai, Z. N. Network biology: Understanding the cell's functional organization. *Nat. Rev. Genet.* **5**, 101–113 (2004).
- Albert, R. Scale-free networks in cell biology. *J. Cell Sci.* **118**, 4947–4957 (2005).
- Arita, M. Scale-freeness and biological networks. *J. Biochem.* **138**, 1–4 (2005).
- Broido, A. D. & Clauset, A. Scale-free networks are rare. *Nat. Commun.* **10**, 1017 (2019).
- Nagano, A. *et al.* Deciphering and prediction of transcriptome dynamics under fluctuating field conditions. *Cell* **151**, 1358–1369. <https://doi.org/10.1016/j.cell.2012.10.048> (2012).
- Johnstone, I. M. *et al.* On the distribution of the largest eigenvalue in principal components analysis. *Ann. Stat.* **29**, 295–327 (2001).
- Bühlmann, P. & van de Geer, S. *Statistics for High-Dimensional Data: Methods, Theory and Applications* (Springer Science & Business Media, 2011).
- Hastie, T., Tibshirani, R. & Friedman, J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer Series in Statistics (Springer, 2009).
- Yuan, M. & Lin, Y. Model selection and estimation in the gaussian graphical model. *Biometrika* **94**, 19–35 (2007).
- Liu, Q. & Ihler, A. T. Learning scale free networks by reweighted L1 regularization. *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics in Proceedings of Machine Learning Research*, **15**, 40–48 (2011).
- Friedman, J., Hastie, T. & Tibshirani, R. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* **9**, 432–441 (2008).
- Witten, D. M., Friedman, J. H. & Simon, N. New insights and faster computations for the graphical lasso. *J. Comput. Graph. Stat.* **20**, 892–900 (2011).
- Boyd, S., Parikh, N., Chu, E., Peleato, B. & Eckstein, J. Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers. *Foundations and Trends[®] in Machine Learning*. **3**, 1–122 (2011).
- Rolfs, B. T. & Rajaratnam, B. A note on the lack of symmetry in the graphical lasso. *Comput. Stat. Data Anal.* **57**, 429–434 (2013).
- Hunter, D. R. & Lange, K. A tutorial on mm algorithms. *Am. Stat.* **58**, 30–37 (2004).
- Chen, J. & Chen, Z. Extended bayesian information criteria for model selection with large model spaces. *Biometrika* **95**, 759–771 (2008).
- Foygel, R. & Drton, M. Extended Bayesian information criteria for Gaussian graphical models. *Adv. Neural. Inform. Process. Syst.* **23**, 604–612 (2010).
- Schwarz, G. Estimation of the mean of a multivariate normal distribution. *Ann. Stat.* **9**, 1135–1151 (1978).
- Fahrmeir, L., Kneib, T., Lang, S. & Marx, B. *Regression* (Springer, 2007).
- Watts, D. J. & Strogatz, S. H. Collective dynamics of 'small-world' networks. *Nature* **393**, 440–442. <https://doi.org/10.1038/30918> (1998).
- Girvan, M. & Newman, M. E. J. Community structure in social and biological networks. *Proc. Natl. Acad. Sci.* **99**, 7821–7826. <https://doi.org/10.1073/pnas.122653799> (2002). <https://www.pnas.org/content/99/12/7821.full.pdf>.
- Bassett, D. S. & Bullmore, E. Small-world brain networks. *Neuroscience* **12**, 512–523. <https://doi.org/10.1177/1073858406293182> (2006).
- Bassett, D. S., Meyer-Lindenberg, A., Achard, S., Duke, T. & Bullmore, E. Adaptive reconfiguration of fractal small-world human brain functional networks. *Proc. Natl. Acad. Sci.* **103**, 19518–19523 (2006).

44. Newman, M. & Watts, D. Renormalization group analysis of the small-world network model. *Phys. Lett. A* **263**, 341–346. [https://doi.org/10.1016/S0375-9601\(99\)00757-4](https://doi.org/10.1016/S0375-9601(99)00757-4) (1999).
45. Amara, L., Scala, A., Barthelemy, M. & Stanley, H. *Classes of Small-World Networks* 207–210 (Princeton University Press, 2011).
46. Newman, M. & Watts, D. *Scaling and Percolation in the Small-World Network Model* 310–320 (Princeton University Press, 2011).

Acknowledgements

The authors would like to thank Mr. Kanta Miura for the valuable discussions. We also thank anonymous reviewers for the constructive and helpful comments that improved the quality of the paper.

Author contributions

Y.O. and K.H. created an R package and performed numerical experiments. Y.O. wrote most of this article, and D.K. and A.J.N. significantly revised Introduction and Simulation frameworks. K.H. wrote technical parts. S.K. first proposed to conduct the numerical simulation of high-dimensional regression on plant science.

Funding

This work was partially supported by the Japan Society for the Promotion of Science KAKENHI 19K11862 (KH) and JST CREST Grant number JPMJCR15O2 (AJN).

Competing interests

The authors declare no competing interests.

Additional information

Correspondence and requests for materials should be addressed to A.J.N. or K.H.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021