

graphsim: An R package for simulating gene expression data from graph structures of biological pathways

S. Thomas Kelly^{1, 2} and Michael A. Black¹

1 Department of Biochemistry, University of Otago, PO Box 56, Dunedin 9054, New Zealand **2**
Present address:RIKEN Center for Integrative Medical Sciences, Suehiro-cho-1-7-22, Tsurumi Ward,
Yokohama, Kanagawa 230-0045, Japan

DOI: [10.21105/joss.0XXXX](https://doi.org/10.21105/joss.0XXXX)

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Statement of Need

Provides a flexible framework to simulate biological pathways from a graph structure based on a statistical model of gene expression.

Editor: [Editor Name](#) ↗

Submitted: 01 January 1900

Published: 01 January 3030

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC-BY](#)).

Summary

Transcriptomic analysis is used to capture the molecular state of a cell or sample in many biological and medical applications. In addition to identifying alterations in activity at the level of individual genes, understanding changes in the gene networks that regulate fundamental biological mechanisms is also an important objective of molecular analysis. As a result, databases that describe biological pathways are increasingly relied on to assist with the interpretation of results from large-scale genomics studies. Incorporating information from biological pathways and gene regulatory networks into a genomic data analysis is a popular strategy, and there are many methods that provide this functionality for gene expression data. When developing or comparing such methods, it is important to gain an accurate assessment of their performance, with simulation-based validation studies a popular choice. This necessitates the use of simulated data that correctly accounts for pathway relationships and correlations. Here we present a versatile statistical framework to simulate correlated gene expression data from biological pathways, by sampling from a multivariate normal distribution derived from a graph structure. This procedure has been released as the `graphsim` R package (<https://github.com/TomKellyGenetics/graphsim>) and is compatible with any graph structure that can be described using the `igraph` package.

Introduction: inference and modelling of biological networks

Network analysis of molecular biological pathways has the potential to lead to new insights into biology and medical genetics (Barabási & Oltvai, 2004; Hu, Thomas, & Brunak, 2016). Since gene expression profiles capture a consistent signature of the regulatory state of a cell (Ozsolak & Milos, 2011; Perou et al., 2000; Svensson, Vento-Tormo, & Teichmann, 2018), they can be used to analyse complex molecular states with genome-scale data. However, biological pathways are often analysed in a reductionist paradigm as amorphous sets of genes involved in particular functions, despite the fact that the relationships defined by pathway structure could further inform gene expression analyses. In many cases, the pathway relationships are well-defined, experimentally-validated, and are available in public databases (Croft et al., 2014). As a result, network analysis techniques could play an important role in furthering our understanding of biological pathways and aiding in the interpretation of genomics studies.

Gene networks provide insights into how cells are regulated, by mapping regulatory interactions between target genes and transcription factors, enhancers, and sites of epigenetic marks or chromatin structures (Barabási & Oltvai, 2004; Yamaguchi, Yoshida, Imoto, Higuchi, & Miyano, 2007). Inference of these regulatory interactions for genomics investigations has the potential to radically expand the range of candidate biological pathways to be further explored, or to improve the accuracy of bioinformatics and functional genomic analysis. A number of methods have already been developed to utilise timecourse gene expression data (Arner et al., 2015; Yamaguchi et al., 2007) using gene regulatory modules in state-space models and recursive vector autoregressive models (Hirose et al., 2008; Shimamura et al., 2009). Various approaches to gene regulation and networks at the genome-wide scale have lead to novel biological insights (Arner et al., 2015; Komatsu et al., 2013). However, inference of regulatory networks has thus far relied on experimental validation or resampling-based approaches to estimate the likelihood of specific network modules being predicted (Hawe, Theis, & Heinig, 2019; Markowitz & Spang, 2007).

There is a need, therefore, for a systematic framework for statistical modelling and simulation of gene expression data derived from hypothetical, inferred or known gene networks. Here we present an package to achieve this, where samples from a multivariate normal distribution are used to generate normally-distributed log-expression data, with correlations between genes derived from the structure of the underlying pathway or gene regulatory network. This methodology enables simulation of expression profiles that approximate the log-transformed and normalised data from microarray and bulk or single-cell RNA-Seq experiments. This procedure has been released as the package to enable the generation of simulated gene expression datasets containing pathway relationships from a known underlying network. These simulated datasets can be used to evaluate various bioinformatics methodologies, including statistical and network inference procedures.

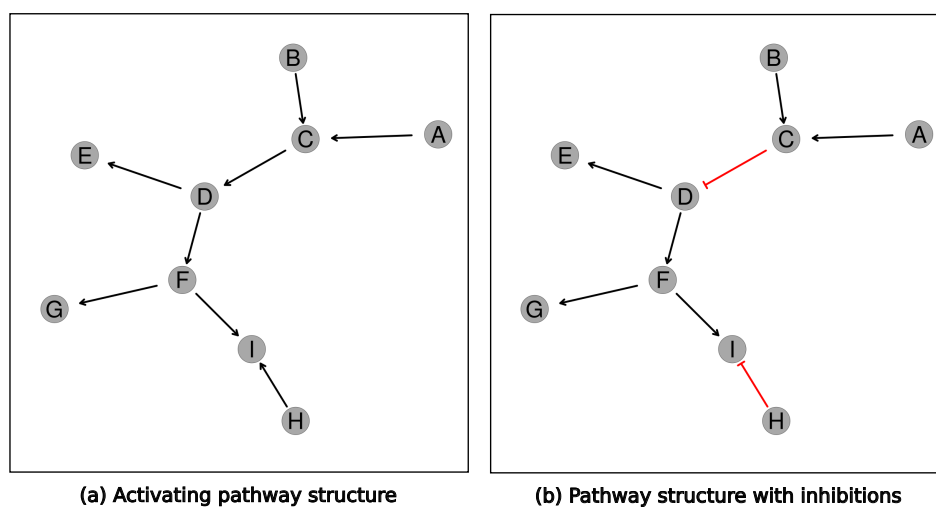


Figure 1: Simulated graph structures. A constructed graph structure used as an example to demonstrate the simulation procedure in Figures 2 and 3. Activating links are denoted by black arrows and inhibiting links by red edges. Inhibiting edges have been highlighted in red.

Methodology and software

Here we present a procedure to simulate gene expression data with correlation structure derived from a known graph structure. This procedure assumes that transcriptomic data have been generated and follow a log-normal distribution (i.e., $\log(X_{ij}) \sim MVN(\mu, \Sigma)$, where μ and Σ are the mean vector and variance-covariance matrix respectively, for gene expression

data derived from a biological pathway) after appropriate normalisation (Law, Chen, Shi, & Smyth, 2014; Li, Piao, Shon, & Ryu, 2015). Log-normality of gene expression matches the assumptions of the popular package, which is often used for the analysis of intensity-based data from gene expression microarray studies and count-based data from RNA-Seq experiments. This approach has also been applied for modelling UMI-based count data from single-cell RNA-Seq experiments in the package (Wang et al., 2018).

In order to simulate transcriptomic data, a pathway is first constructed as a graph structure, using the package (Csardi & Nepusz, 2006), with the status of the edge relationships defined (i.e, whether they activate or inhibit downstream pathway members). This procedure uses a graph structure such as that presented in Figure 1a. The graph can be defined by an adjacency matrix, A (with elements A_{ij}), where

$$A_{ij} = \begin{cases} 1 & \text{if genes } i \text{ and } j \text{ are adjacent} \\ 0 & \text{otherwise} \end{cases}$$

A matrix, R , with elements R_{ij} , is calculated based on distance (i.e., number of edges contained in the shortest path) between nodes, such that closer nodes are given more weight than more distant nodes, to define inter-node relationships. A geometrically-decreasing (relative) distance weighting is used to achieve this:

$$R_{ij} = \begin{cases} 1 & \text{if genes } i \text{ and } j \text{ are adjacent} \\ (\frac{1}{2})^{d_{ij}} & \text{if a path can be found between genes } i \text{ and } j \\ 0 & \text{if no path exists between genes } i \text{ and } j \end{cases}$$

where d_{ij} is the length of the shortest path (i.e., minimum number of edges traversed) between genes (nodes) i and j in graph G . Each more distant node is thus related by $\frac{1}{2}$ compared to the next nearest, as shown in Figure 2b. An arithmetically-decreasing (absolute) distance weighting is also supported in the package which implements this procedure:

$$R_{ij} = \begin{cases} 1 & \text{if genes } i \text{ and } j \text{ are adjacent} \\ 1 - \frac{d_{ij}}{\text{diam}(G)} & \text{if a path can be found between genes } i \text{ and } j \\ 0 & \text{if no path exists between genes } i \text{ and } j \end{cases}$$

Assuming a unit variance for each gene, these values can be used to derive a Σ matrix:

$$\Sigma_{ij} = \begin{cases} 1 & \text{if } i = j \\ \rho R_{ij} & \text{otherwise} \end{cases}$$

where ρ is the correlation between adjacent nodes. Thus covariances between adjacent nodes are assigned by a correlation parameter (ρ) and the remaining off-diagonal values in the matrix are based on scaling these correlations by the geometrically weighted relationship matrix (or the nearest positive definite matrix for Σ with negative correlations).

Computing the nearest positive definite matrix is necessary to ensure that the variance-covariance matrix could be inverted when used as a parameter in multivariate normal simulations, particularly when negative correlations are included for inhibitions (as shown below). Matrices that could not be inverted occurred rarely with biologically plausible graph structures but this approach allows for the computation of a plausible correlation matrix when the graph structure given is incomplete or contains loops. When required, the nearest positive definite matrix is computed using the `nearPD` function of the package (Bates & Maechler, 2016) to perform Higham's algorithm (Higham, 2002) on variance-covariance matrices. The package gives a warning when this occurs.

Illustrations

Generating a Graph Structure

The graph structure in Figure 1a was used to simulate correlated gene expression data by sampling from a multivariate normal distribution using the package (Genz & Bretz, 2009; Genz et al., 2016). The graph structure visualisation in Figure 1 was specifically developed for (directed) iGraph objects in and is available in the and packages. The `plot_directed` function enables customisation of plot parameters for each node or edge, and mixed (directed) edge types for indicating activation or inhibition. These inhibition links (which occur frequently in biological pathways) are demonstrated in Figure 1b.

A graph structure can be generated and plotted using the following commands in R:

```
#install packages required (once per machine)

install.packages("igraph")
if(! require("devtools") ){
  install.packages("devtools")
  library("devtools")
}
devtools::install_github("TomKellyGenetics/graphsim")

#load required packages (once per R instance)

library("igraph")
library("graphsim")

#generate graph structure

graph_edges <- rbind(c("A", "C"), c("B", "C"), c("C", "D"), c("D", "E"),
                    c("D", "F"), c("F", "G"), c("F", "I"), c("H", "I"))
graph <- graph.edgelist(graph_edges, directed = TRUE)

#plot graph structure (Figure 1)

plot_directed(graph, state = "activating", layout = layout.kamada.kawai,
              cex.node=3, cex.arrow=5, arrow_clip = 0.2)

#generate parameters for inhibitions

state <- c(1, 1, -1, 1, 1, 1, 1, -1, 1)

#plot graph structure with inhibitions (Figure 2)

plot_directed(graph, state=state, layout = layout.kamada.kawai,
              cex.node=3, cex.arrow=5, arrow_clip = 0.2)
```

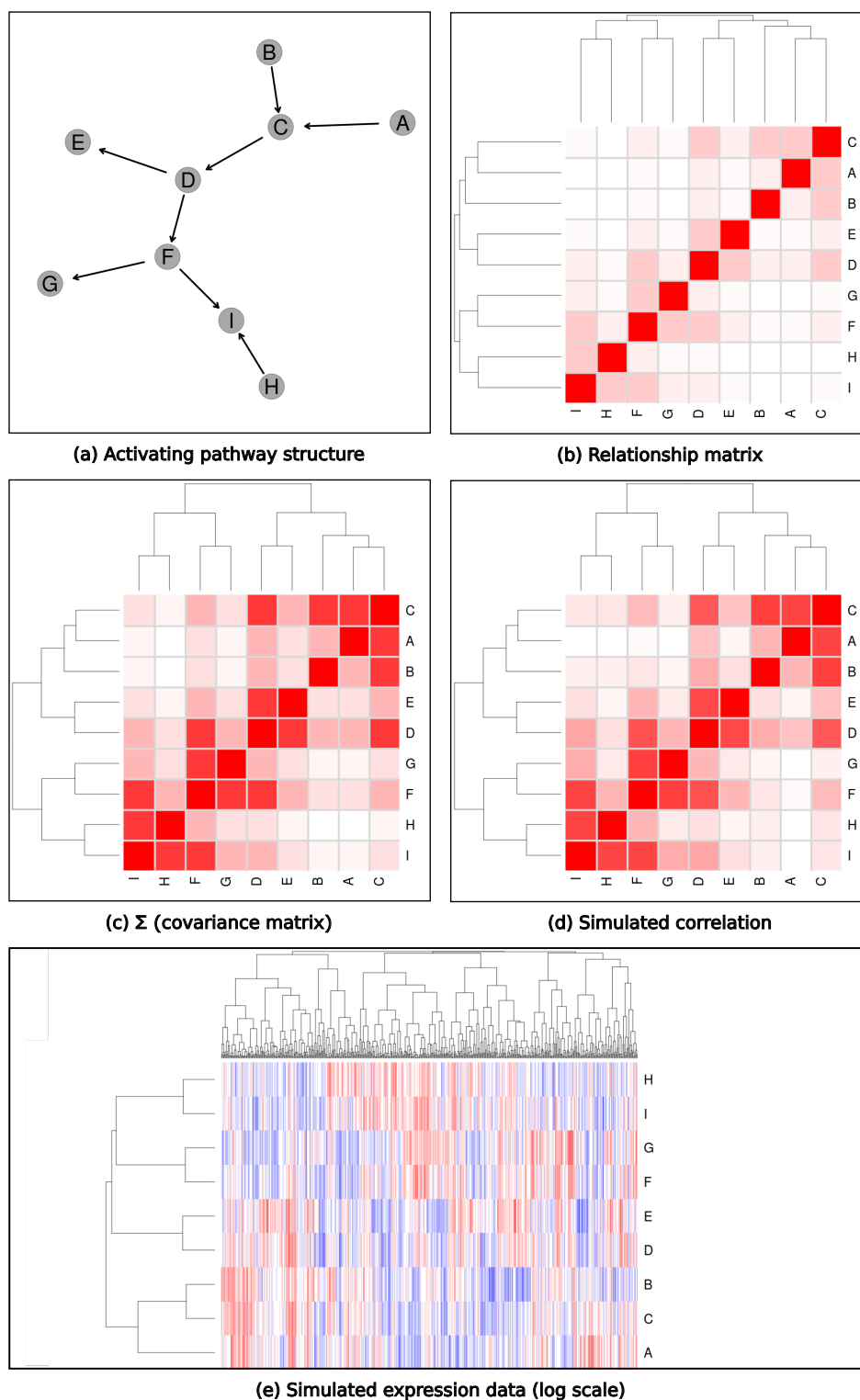


Figure 2: Simulating expression from a graph structure. An example of a graph structure (a) that has been used to derive a relationship matrix (b), Σ matrix (c) and correlation structure (d) from the relative distances between the nodes. Non-negative values are coloured white to red from 0 to 1. This Σ matrix has been used to generate a simulated expression dataset of 100 samples (coloured blue to red from low to high) via sampling from the multivariate normal distribution. Here genes with closer relationships in the pathway structure show higher correlation between simulated values.

Generating a Simulated Expression Dataset

The correlation parameter of $\rho = 0.8$ is used to demonstrate the inter-correlated datasets using a geometrically-generated relationship matrix (as used for the example in Figure 2c). This Σ matrix was then used to sample from a multivariate normal distribution such that each gene had a mean of 0, standard deviation 1, and covariance within the range $[0,1]$ so that the off-diagonal elements of Σ represent correlations. This procedure generated a simulated (continuous normally-distributed) log-expression profile for each node (Figure 2e) with a corresponding correlation structure (Figure 2d). The simulated correlation structure closely resembled the expected correlation structure (Σ in Figure 2c) even for the relatively modest sample size ($N = 100$) illustrated in Figure 2. Once a gene expression dataset comprising multiple pathways has been generated (as in Figure 2e), it can then be used to test procedures designed for analysis of empirical gene expression data (such as those generated by microarrays or RNA-Seq) that have been normalised on a log-scale.

The simulated dataset can be generated using the following code:

```
#adjacency matrix

adj_mat <- make_adjmatrix_graph(graph)

#relationship matrix

dist_mat <- make_distance_graph(graph_test4, absolute = FALSE)

#sigma matrix directly from graph

sigma_mat <- make_sigma_mat_dist_graph(graph, 0.8, absolute = FALSE)

#show shortest paths of graph

shortest_paths <- shortest.paths(graph)

#generate expression data directly from graph

expr <- generate_expression(100, graph, cor = 0.8, mean = 0, comm = F,
                           dist = TRUE, absolute = FALSE, state = state)

#plot adjacency matrix

heatmap.2(make_adjmatrix_graph(graph), scale = "none", trace = "none",
          col = colorpanel(3, "grey75", "white", "blue"),
          colsep = 1:length(V(graph)), rowsep = 1:length(V(graph)))

#plot relationship matrix

heatmap.2(make_distance_graph(graph_test4, absolute = FALSE),
          scale = "none", trace = "none", col = bluered(50),
          colsep = 1:length(V(graph)), rowsep = 1:length(V(graph)))

#plot sigma matrix

heatmap.2(make_sigma_mat_dist_graph(graph, 0.8, absolute = FALSE),
          scale = "none", trace = "none", col = bluered(50),
          colsep = 1:length(V(graph)), rowsep = 1:length(V(graph)))
```

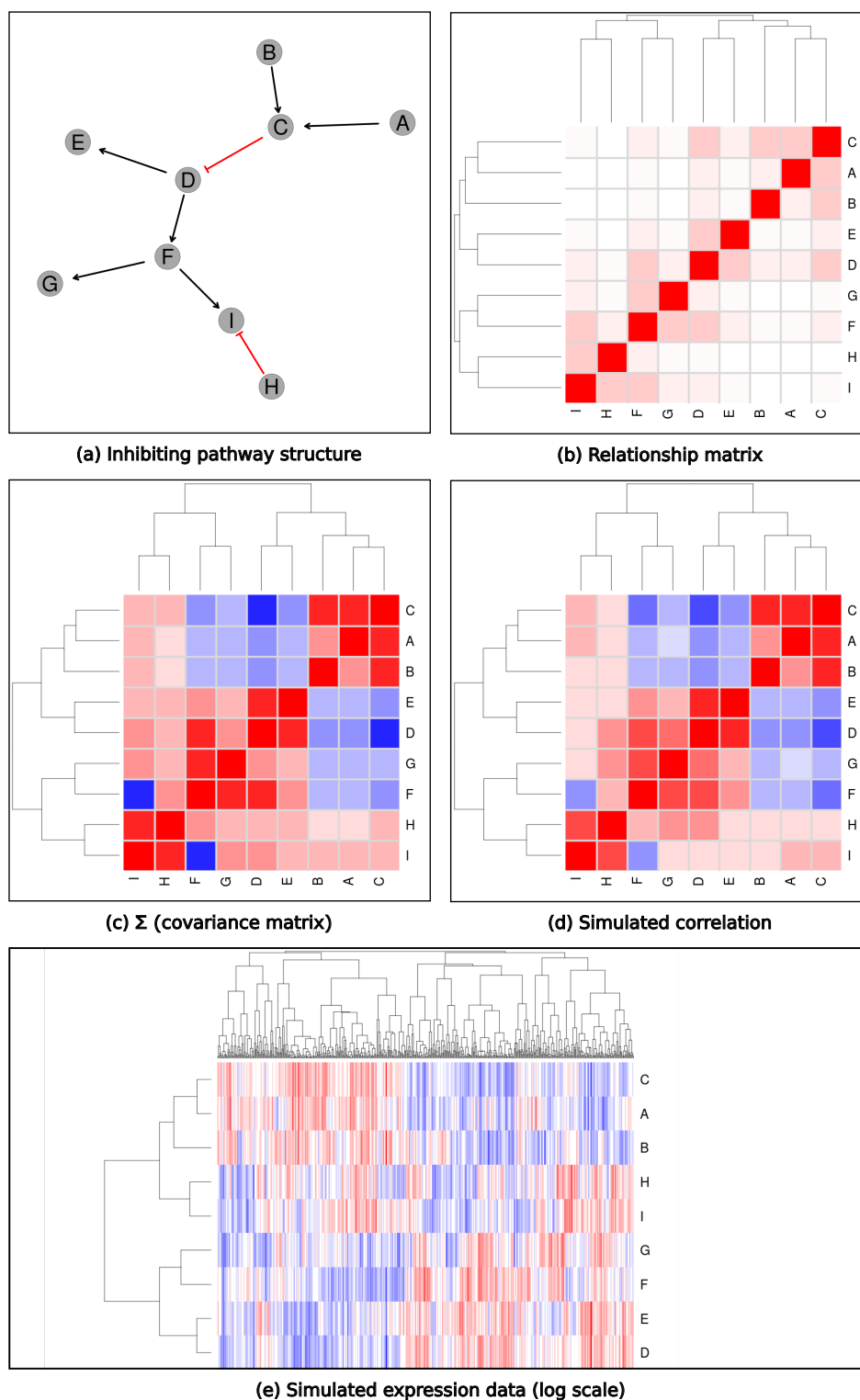


Figure 3: Simulating expression from graph structure with inhibitions. Simulating expression from graph structure with inhibitions.}} An example of a graph structure (a), that has been used to derive a relationship matrix (b), Σ matrix (c), and correlation structure (d), from the relative distances between the nodes. These values are coloured blue to red from -1 to 1 . This has been used to generate a simulated expression dataset of 100 samples (coloured blue to red from low to high) via sampling from the multivariate normal distribution. Here the inhibitory relationships between genes are reflected in negatively correlated simulated values.

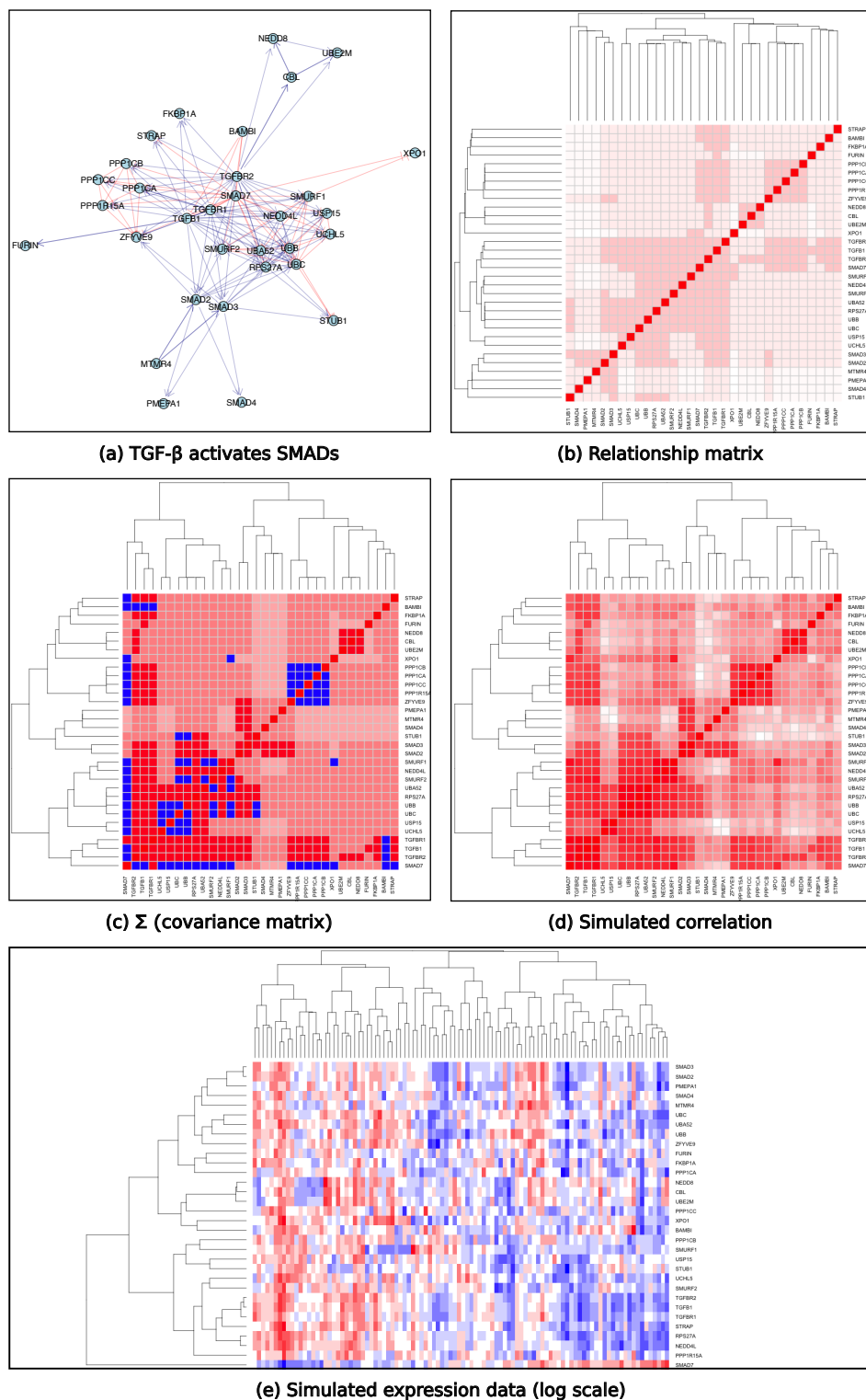


Figure 4: Simulating expression from a biological pathway graph structure. Simulating expression from graph structure with inhibitions.}} The graph structure (a) of a known biological pathway, the TGF- β receptor signaling activates SMADs (R-HSA-2173789), was used to derive a relationship matrix (b), Σ matrix (c) and correlation structure (d) from the relative distances between the nodes. These values are coloured blue to red from -1 to 1 . This has been used to generate a simulated expression dataset of 100 samples (coloured blue to red from low to high) via sampling from the multivariate normal distribution. Here modules of genes with correlated expression can be clearly discerned.


```
expr <- generate_expression(100, graph, cor = 0.8, mean = 0,
  comm = FALSE, dist = TRUE, absolute = FALSE, state = state)

#plot simulated expression data

heatmap.2(expr, scale = "none", trace = "none", col = bluered(50),
  colsep = 1:length(V(graph)), rowsep = 1:length(V(graph)))

#plot simulated correlations

heatmap.2(cor(t(expr)), scale = "none", trace = "none", col = bluered(50),
  colsep = 1:length(V(graph)), rowsep = 1:length(V(graph)))
```

The simulation procedure (Figure 2) can similarly be used for pathways containing inhibitory links (Figure 3) with several refinements. With the inhibitory links (Figure 3a), distances are calculated in the same manner as before (Figure 3b) with inhibitions accounted for by iteratively multiplying downstream nodes by -1 to form modules with negative correlations between them (Figures 3c and 3d). A multivariate normal distribution with these negative correlations can be sampled to generate simulated data (Figure 3e).

The simulation procedure is also demonstrated here (Figure 4) on a pathway structure for a known biological pathway (from reactome R-HSA-2173789) of TGF- β receptor signaling activates SMADs (Figure 4a) derived from the Reactome database version 52 [Reactome]. Distances are calculated in the same manner as before (Figure 4b) producing blocks of correlated genes (Figures 4c and 4d). This shows that multivariate normal distribution can be sampled to generate simulated data to represent expression with the complexity of a biological pathway (Figure 4e). Here *SMAD7* exhibits negative correlations with the other SMADs consistent with its functions as an “inhibitor SMAD” with competitively inhibits *SMAD4*.

These simulated datasets could then be used for simulating synthetic lethal partners of a query gene within a graph network. The query gene was assumed to be separate from the graph network pathway and was added to the dataset using the procedure in Section [methods:simulating_SL]. Thus I can simulate known synthetic lethal partner genes within a synthetic lethal partner pathway structure.

Summary and discussion

Biological pathways are of fundamental importance to understanding molecular biology. In order to translate findings from genomics studies into real-world applications such as improved healthcare, the roles of genes must be studied in the context of molecular pathways. Here we present a statistical framework to simulate gene expression from biological pathways, and provide the package in to generate these simulated datasets. This approach is versatile and can be fine-tuned for modelling existing biological pathways or for testing whether constructed pathways can be detected by other means. In particular, methods to infer biological pathways and gene regulatory networks from gene expression data can be tested on simulated datasets using this framework. The package also enables simulation of complex gene expression datasets to test how these pathways impact on statistical analysis of gene expression data using existing methods or novel statistical methods being developed for gene expression data analysis.

Computational details

The results in this paper were obtained using R 3.6.1 with the packages 1.2.4.1, 1.2-17, 1.0-3, and 1.0-11 packages. itself and all dependent packages used are available from the Comprehensive Archive Network (CRAN) at <https://CRAN.R-project.org>. The packages presented can be installed from <https://github.com/TomKellyGenetics/graphsim> and <https://github.com/TomKellyGenetics/plot.igraph> respectively. These functions can also be installed using the library at <https://github.com/TomKellyGenetics/igraph.extensions> which includes other plotting functions used. This software is cross-platform and compatible with installations on Windows, Mac, and Linux operating systems. The package GitHub repository also contains Vignettes with more information and examples on running functions released in the package. The package (0.1.0) meets CRAN submission criteria and will be released.

Acknowledgements

This package was developed as part of a PhD research project funded by the Postgraduate Tassell Scholarship in Cancer Research Scholarship awarded to STK. We thank members of the Laboratory of Professor Satoru Miyano at the University of Tokyo, Institute for Medical Science, Professor Seiya Imoto, Associate Professor Rui Yamaguchi, and Dr Paul Sheridan (Assistant Professor at Hirosaki University, CSO at Tupac Bio) for helpful discussions in this field. We also thank Professor Parry Guilford at the University of Otago, Professor Cristin Print at the University of Auckland, and Dr Erik Arner at the RIKEN Center for Integrative Medical Sciences for their excellent advice during this project.

References

- Arner, E., Daub, C. O., Vitting-Seerup, K., Andersson, R., Lilje, B., Drabløs, F., Lennartsson, A., et al. (2015). Transcribed enhancers lead waves of coordinated transcription in transitioning mammalian cells. *Science*, 347(6225), 1010–1014.
- Barabási, A. L., & Oltvai, Z. N. (2004). Network biology: Understanding the cell's functional organization. *Nat Rev Genet*, 5(2), 101–113.
- Bates, D., & Maechler, M. (2016). *Matrix: Sparse and dense matrix classes and methods*. Retrieved from <https://CRAN.R-project.org/package=Matrix>
- Croft, D., Mundo, A. F., Haw, R., Milacic, M., Weiser, J., Wu, G., Caudy, M., et al. (2014). The Reactome pathway knowledgebase. *Nucleic Acids Res*, 42(database issue), D472–D477. Journal Article. doi:[10.1093/nar/gkt1102](https://doi.org/10.1093/nar/gkt1102)
- Csardi, G., & Nepusz, T. (2006). The igraph software package for complex network research. *InterJournal, Complex Systems*, 1695. Retrieved from <http://igraph.org>
- Genz, A., & Bretz, F. (2009). Computation of multivariate normal and t probabilities. In *Lecture notes in statistics* (Vol. 195). Heidelberg: Springer-Verlag.
- Genz, A., Bretz, F., Miwa, T., Mi, X., Leisch, F., Scheipl, F., & Hothorn, T. (2016). *Mvtnorm: Multivariate normal and t distributions*. Retrieved from <http://CRAN.R-project.org/package=mvtnorm>
- Hawe, J. S., Theis, F. J., & Heinig, M. (2019). Inferring Interaction Networks From Multi-Omics Data. *Front Genet*, 10, 535.
- Higham, N. J. (2002). Computing the nearest correlation matrix—a problem from finance. *IMA Journal of Numerical Analysis*, 22(3), 329–343. doi:[10.1093/imanum/22.3.329](https://doi.org/10.1093/imanum/22.3.329)

- Hirose, O., Yoshida, R., Imoto, S., Yamaguchi, R., Higuchi, T., Charnock-Jones, D. S., Print, C., et al. (2008). Statistical inference of transcriptional module-based gene networks from time course gene expression profiles by using state space models. *Bioinformatics*, 24(7), 932–942. doi:[10.1093/bioinformatics/btm639](https://doi.org/10.1093/bioinformatics/btm639)
- Hu, J. X., Thomas, C. E., & Brunak, S. (2016). Network biology concepts in complex disease comorbidities. *Nat. Rev. Genet.*, 17(10), 615–629.
- Komatsu, M., Yoshimaru, T., Matsuo, T., Kiyotani, K., Miyoshi, Y., Tanahashi, T., Rokutan, K., et al. (2013). Molecular features of triple negative breast cancer cells by genome-wide gene expression profiling analysis. *Int. J. Oncol.*, 42(2), 478–506.
- Law, C. W., Chen, Y., Shi, W., & Smyth, G. K. (2014). voom: Precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol.*, 15(2), R29. doi:[10.1186/gb-2014-15-2-r29](https://doi.org/10.1186/gb-2014-15-2-r29)
- Li, P., Piao, Y., Shon, H. S., & Ryu, K. H. (2015). Comparing the normalization methods for the differential analysis of Illumina high-throughput RNA-Seq data. *BMC Bioinformatics*, 16, 347.
- Markowitz, F., & Spang, R. (2007). Inferring cellular networks—a review. *BMC Bioinformatics*, 8 Suppl 6, S5.
- Ozsolak, F., & Milos, P. M. (2011). RNA sequencing: advances, challenges and opportunities. *Nat. Rev. Genet.*, 12(2), 87–98.
- Perou, C. M., Sørlie, T., Eisen, M. B., Rijn, M. van de, Jeffrey, S. S., Rees, C. A., Pollack, J. R., et al. (2000). Molecular portraits of human breast tumours. *Nature*, 406(6797), 747–752.
- Shimamura, T., Imoto, S., Yamaguchi, R., Fujita, A., Nagasaki, M., & Miyano, S. (2009). Recursive regularization for inferring gene networks from time-course gene expression profiles. *BMC Systems Biology*, 3(1), 41. doi:[10.1186/1752-0509-3-41](https://doi.org/10.1186/1752-0509-3-41)
- Svensson, V., Vento-Tormo, R., & Teichmann, S. A. (2018). Exponential scaling of single-cell RNA-seq in the past decade. *Nat Protoc*, 13(4), 599–604.
- Wang, J., Huang, M., Torre, E., Dueck, H., Shaffer, S., Murray, J., Raj, A., et al. (2018). Gene expression distribution deconvolution in single-cell RNA sequencing. *Proc. Natl. Acad. Sci. U.S.A.*, 115(28), E6437–E6446.
- Yamaguchi, R., Yoshida, R., Imoto, S., Higuchi, T., & Miyano, S. (2007). Finding module-based gene networks with state-space models - Mining high-dimensional and short time-course gene expression data. *IEEE Signal Processing Magazine*, 24(1), 37–46.