

Using sabinaHSBM for link prediction and network reconstruction using Hierarchical Stochastic Block Models

Example on how to use it to identify missing and spurious links with the
full_reconstruction method (with parallelized computation)

The *sabinaHSBM* package provides tools for link prediction and network reconstruction using hierarchical stochastic block models (HSBM). This document demonstrates a simple use for the `full_reconstruction` method, showcasing how it handles missing and spurious links in an incomplete and error-prone network. To enhance computational efficiency during the prediction links steps we use parallelized computation.

We use the example dataset `dat2` included in the package to show key functionalities, including data preparation, link prediction, and network reconstruction.

Important note:

There are **two ways** to use the package:

- **UNIX users (native installation)** can run *sabinaHSBM* locally **if** their system includes:
 - R (version $\leq 4.0.4$) with all required R packages (listed below)
 - Python $\leq 3.9.2$ with the `graph-tool` library (version ≤ 2.45)
- **All other users**, or those who prefer to avoid manual setup, can use the **ready-to-use Docker image**, which includes everything needed:
 - All R and Python dependencies
 - The *sabinaHSBM* package pre-installed and ready to use

Loading Required Libraries

Note:

The following instructions are **only required for UNIX users** running *sabinaHSBM* *natively* (outside Docker).

These users must ensure their system includes the required dependencies before proceeding.

```
# If the package is not installed, install it from GitHub
if (!requireNamespace("sabinaHSBM", quietly = TRUE)) {
  library(remotes)
  remotes::install_github("h-lima/sabinaHSBM")
}

# Load the sabinaHSBM package
library(sabinaHSBM)
```

Install required R packages if not already available:

```
list.of.packages <- c(
  "dplyr",
  "parallel",
  "reshape2",
  "reticulate",
```

```

"stringr",
"tidyr",
"ROCR",
"data.table"
)

new.packages <- list.of.packages[!(list.of.packages %in% installed.packages()[, "Package"])]
if (length(new.packages) > 0) {
  install.packages(new.packages, dependencies = TRUE)
}

for (package in list.of.packages) {
  library(package, character.only = TRUE)
}

```

If you're using the **Docker image** (see *Supporting Information S3*) **you do not need to install any packages**, simply start the container and load the package in your R session with:

```
library(sabinaHSBM)
```

```

# Record starting time (optional)
start_time <- Sys.time()

```

Load data

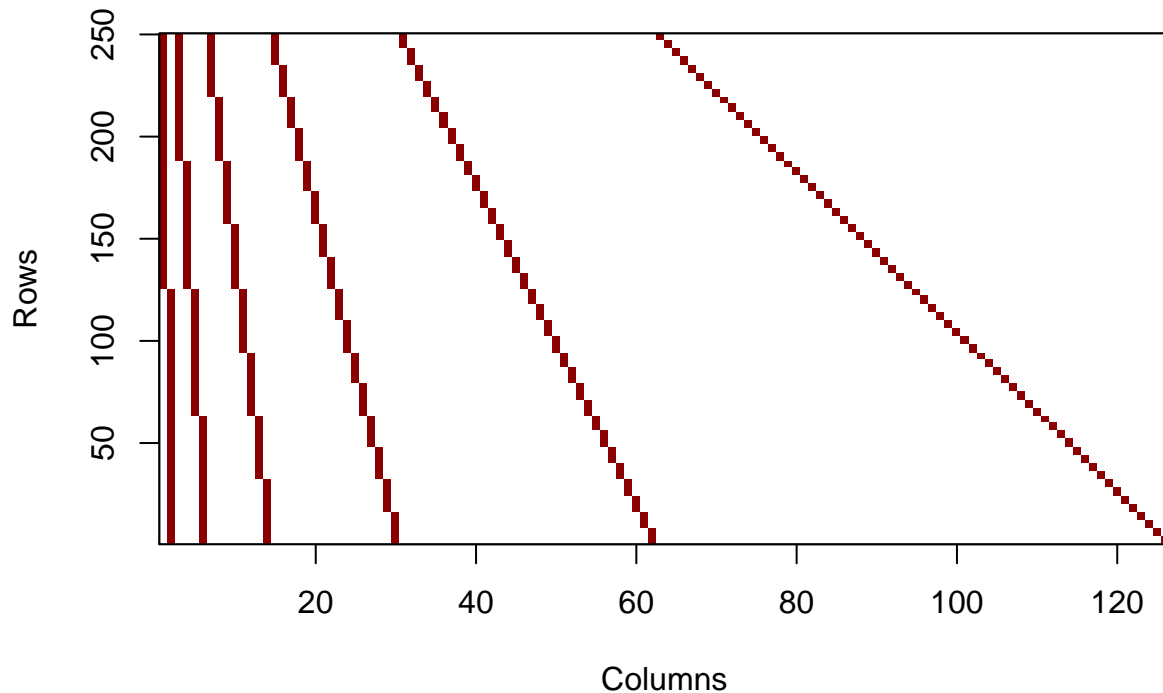
The dataset `dat2` is a binary bipartite matrix representing a hypothetical species interactions network. Columns and rows correspond to two different types of nodes (e.g., hosts and parasites), and links (values of 1, in red) represent interactions between them while 0 (in black) represent lack of an observed interaction. These links are distributed in a structured pattern, with no missing information, to focus the example on identifying spurious links.

```

# Load the dataset
data(dat2, package = "sabinaHSBM")
dat<- dat2

# Plot the original matrix
plot_interaction_matrix(adj_mat = dat, order_mat = FALSE)

```



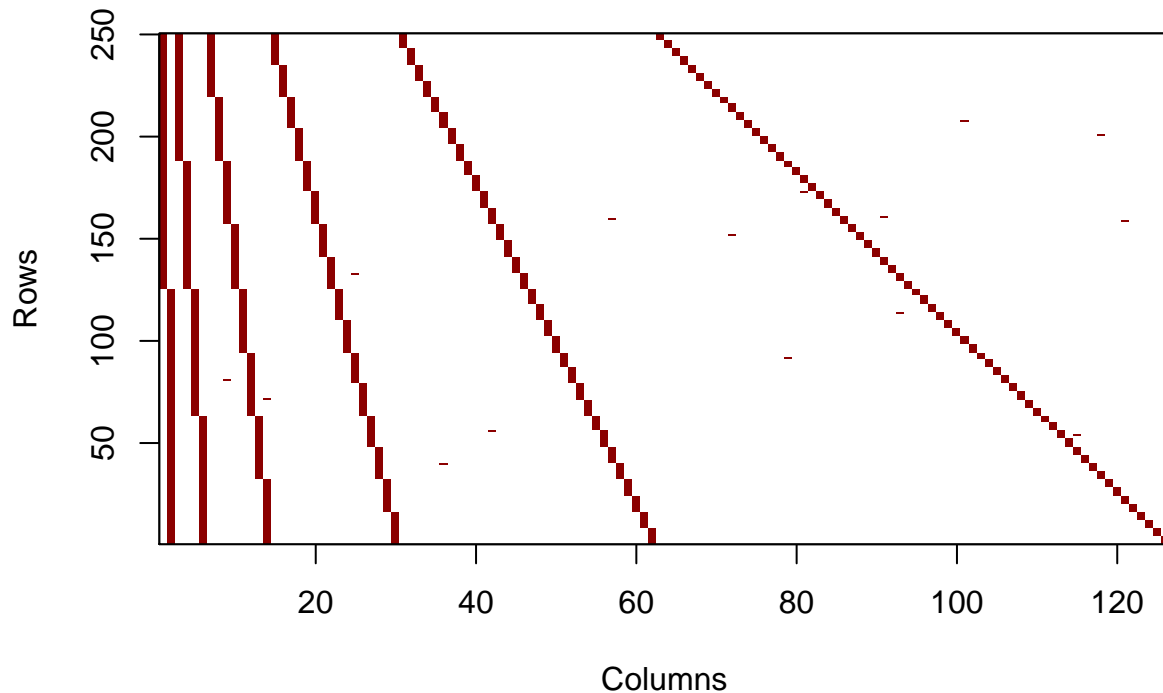
To simulate spurious links and address the `full_reconstruction` method, we randomly add a small number of false positives (1) to the matrix. This controlled modification allows us to demonstrate how the method identifies both spurious links (observed but potentially erroneous interactions, false positives) and missing links (unobserved but likely interactions, false negatives).

```
# Add spurious links artificially
set.seed(123)
num_spurious <- ceiling(sum(dat) * 0.01) # Proportion of spurious links to add
spurious_links <- matrix(nrow = num_spurious, ncol = 2)

for (i in 1:num_spurious) {
  repeat {
    random_row <- sample(1:nrow(dat), 1)
    random_col <- sample(1:ncol(dat), 1)

    if (dat[random_row, random_col] == 0) {
      dat[random_row, random_col] <- 1
      spurious_links[i,] <- c(random_row, random_col)
      break
    }
  }
}

# Plot the matrix with spurious
plot_interaction_matrix(dat, order_mat = FALSE)
```



Preparing Input data for HSBM

The `hsbm.input` function pre-processes the dataset, creating cross-validation folds and edge lists required for modeling. Here, we use 5-fold cross-validation and the `full_reconstruction` method.

```
# Prepare input data
n_folds <- 5 # Number of folds for cross-validation

myInput <- hsbm.input(
  dat, # Binary bipartite matrix of observed links
  n_folds = n_folds,
  add_spurious = FALSE # Simulate spurious links
)
```

```
## [1] "Actual cross-validation rate is 0.200"
## [2] "Actual cross-validation rate is 0.200"
## [3] "Actual cross-validation rate is 0.200"
## [4] "Actual cross-validation rate is 0.200"
## [5] "Actual cross-validation rate is 0.200"
```

For the purposes of this example, artificial spurious links are labeled with the “`spurious_edge`” tag using the code below. This labeling allows for tracking these links throughout the analysis and network reconstruction process, although in typical use, such labeling would not be necessary.

```
myInput$edgelists <- lapply(myInput$edgelists,
  function(x){
    x$x <- 1
  })
```

```

        return(x)
    })

for(i in 1:length(myInput$edgelist)) {
  el <- myInput$edgelist[[i]]
  rows <- as.numeric(el[, 1]) + 1
  cols <- as.numeric(el[, 2]) - nrow(myInput$data) + 1
  rows_cols <- paste0(rows, "_", cols)
  spurious_rows_cols <- paste0(spurious_links[, 1], "_", spurious_links[, 2])
  spurious_loc <- which(rows_cols %in% spurious_rows_cols)
  el[spurious_loc, ]$x <- 1
  el[spurious_loc, ]$edge_type <- "spurious_edge"
  myInput$edgelist[[i]] <- el
}

# Summarizes network characteristics
summary(myInput)

##   n_rows n_cols n_links rows_single_link cols_single_link possible_links
## 1     250   126   1515              0              0          31500

```

Predict link probabilities

The `hsbm.predict` function applies the HSBM to predict probabilities of all links (observed and unobserved links). This step is crucial to identify spurious and missing links within the data, which are often present in incomplete or error-prone networks. The function works directly with the processed input created by `hsbm.input`. This step can be computationally intensive when working with large datasets or numerous folds. To improve performance, we use parallelized computation, distributing tasks across multiple cores.

```

# Define the number of cores to use
nCores <- 2

# Generate HSBM predictions
myPred <- hsbm.predict(
  myInput,          # Input data processed by hsbm.input()
  iter = 10000,     # Number of iterations
  wait = 1000,      # Number of iterations for MCMC equilibration
  method = "full_reconstruction", # Method for link prediction
  save_blocks = TRUE, # Save group assignments
  save_pickle = FALSE, # Save results as pickle files,
  save_plots = FALSE, # Save hierarchical edge bundling plots
  n_cores = nCores
)

```

Predicted link probabilities and group assignments are stored for each fold. Below, we extract the link probabilities (p) and groups assignments of links for fold 1.

```

# View probabilities for fold 1
probabilities_fold1 <- myPred$probs[[1]]
head(probabilities_fold1)

##   v1  v2 p v1_names v2_names edge_type
## 1   0 280 1      sp1      SPeb documented
## 2   0 250 1      sp1      SPa documented
## 3   0 256 1      sp1      SPg documented

```

```
## 4  0 264 1      sp1      SPo documented
## 5  0 312 1      sp1      SPkc documented
## 6  0 252 1      sp1      SPc documented
```

The group assignments provide the hierarchical clustering structure of nodes for each fold. Let's extract and examine the group assignments for fold 1:

```
# View the group/block assignments for fold 1
groups_fold1 <- myPred$groups[[1]]

# Filter one type of nodes (nodes in columns)
vnames <- colnames(myPred$data)
groups_fold1 <- groups_fold1 %>% filter(names %in% vnames)
g_cols <- grep("^G", names(groups_fold1))
groups_fold1[g_cols] <- lapply(groups_fold1[g_cols], sort)

print(head(groups_fold1))
```

```
##   nodes G1  G2 G3 G4 G5 names
## 1   250 25 258  0  3  0  SPa
## 2   251 29 258  0  3  0  SPb
## 3   252 29 258  0  3  0  SPc
## 4   253 29 258  0  3  0  SPd
## 5   254 29 258  0  3  0  SPe
## 6   255 29 258  0  3  0  SPf
```

Hierarchical group assignments provide insight into how nodes (e.g., hosts) are organized across multiple levels. At the first level (G1) nodes are divided into specific groups, reflecting fine-scale patterns. Moving to higher levels (G2, G3, G4) these groups (or communities) are progressively aggregated, revealing broader patterns and relationships or communities.

Next, we will explore how the model assigns specific probabilities to different types of links (“documented,” “held out,” and “spurious_edge”) for each fold to evaluate HSBM ability to correctly identify them. Documented links, which are real observations (true positives), should have probabilities close to 1. Held-out links, which are real observations but transformed to 0s during training for validation, should also show high probabilities. In contrast, spurious edges, artificially generated (initially unobserved links transformed to 1s), should have probabilities close to 0 (unless it could be a missing link), as the model should recognize them as false links (false negatives).

```
fold_res <- list()
for(i in 1:5){
  doc <- dplyr::filter(myPred$probs[[i]], edge_type == "documented")
  held <- dplyr::filter(myPred$probs[[i]], edge_type == "held_out")
  spur <- dplyr::filter(myPred$probs[[i]], edge_type == "spurious_edge")
  fold_data <- c(i,
                 nrow(doc), sum(doc$p < 0.95),
                 nrow(held), sum(held$p < 0.95),
                 nrow(spur), sum(spur$p < 0.95))
  fold_res[[i]] <- fold_data
}

res_summary <- do.call(rbind, fold_res)
colnames(res_summary) <- c("Fold", "Nr doc", "Doc p < 0.95",
                          "Nr held", "Held p < 0.95",
                          "Nr spur", "Spur p < 0.95")

kable(res_summary)
```

Fold	Nr doc	Doc p < 0.95	Nr held	Held p < 0.95	Nr spur	Spur p < 0.95
1	1197	0	303	0	15	4
2	1199	0	301	0	15	7
3	1199	0	301	0	15	2
4	1204	0	296	0	15	0
5	1201	0	299	0	15	1

Network Reconstruction

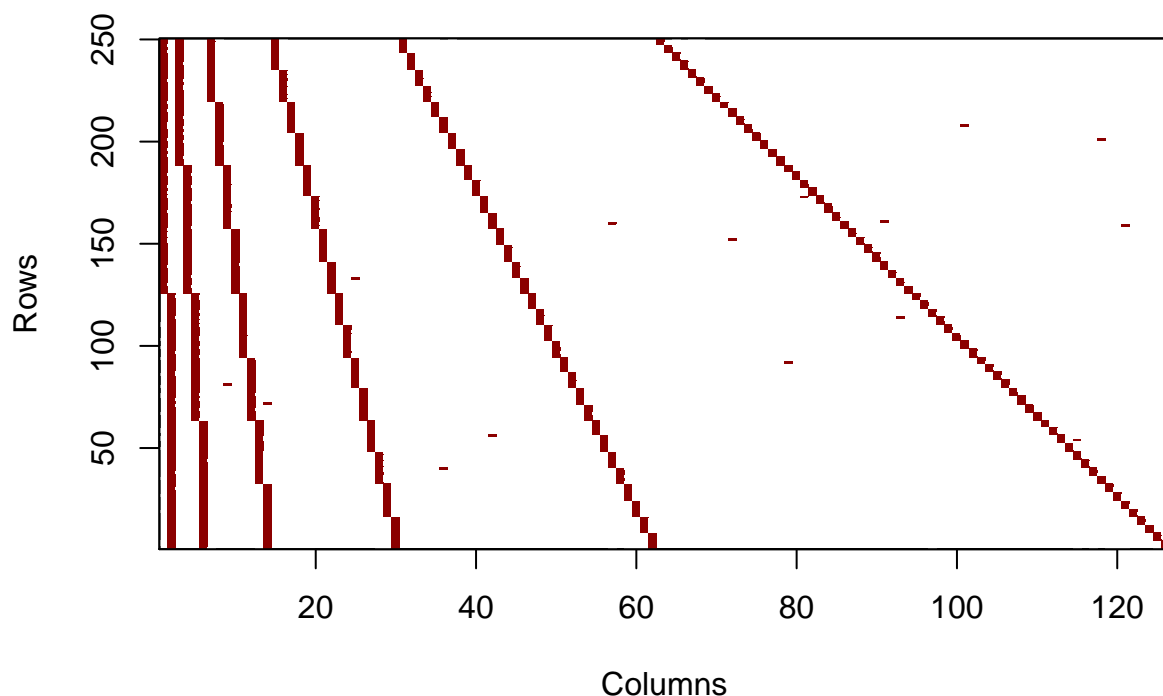
The `hsbm.reconstructed` function generates a reconstructed binary interaction matrix by combining predictions from all folds. The predicted matrix is transformed to binary values using a user-specified threshold.

```
# Network reconstruction
myReconst <- hsbm.reconstructed(
  myPred,                # Predictions processed by hsbm.predict
  rm_documented=FALSE,   # Remove documented links during validation
  na_treatment = "ignore_na", # Handle NA values in predictions
  threshold = 0.5,        # Binarization threshold
  new_matrix_method = "average_thresholded", # Combine fold predictions
  spurious_edges = TRUE,
)
```

```
## Warning in get_hsbm_results(hsbm_out, input_names = TRUE, na_treatment =
## na_treatment): Predictions obtained for 14.96% of the links. Consider
## increasing the number of iterations.
```

This output includes the final averaged probability matrix, the final reconstructed binary matrix, and evaluation metrics. These results highlight the predicted interactions, showcasing the method's ability to detect spurious and missing links effectively. Let's explore some of the key outputs.

```
# Visualize the reconstructed binary matrix
plot_interaction_matrix(myReconst$new_mat, order_mat = FALSE)
```



```
# View the reconstruction summary and evaluation metrics
summary(myReconst)
```

```
## $`Reconstructed Network Metrics`
##  obs_links unobs_links pred_links kept_links spurious_links missing_links
## 1      1515      29985      1515      1515              0              0
##
## $`Evaluation Metrics`
##  mean_RLRR mean_auc mean_aucpr mean_yPRC mean_prec mean_sens mean_spec
## 1          1 0.9999317 0.9999372 0.04809524          1 0.999604          1
##  mean_ACC      mean_ERR mean_tss
## 1 0.999981 1.904762e-05 0.999604
```

The summary provides the number of spurious and missing links, as well as key evaluation metrics, such as the retained link recovery rate (RLRR).

Let's now examine the top 10 predicted links that are most likely to be spurious (false positives) by visualizing the probabilities of "documented" links.

```
# Visualize the top most likely spurious links
top_links_spurious <- top_links(myReconst,
                                n = 10,
                                edge_type = "documented")
print(top_links_spurious)
```

```
##  v1_names v2_names      p      sd
## 1    sp50    SPne 0.7373737 0.21734009
## 2   sp211    SPjb 0.8000000 0.43358967
```



```
## 3      sp43      SPwd 0.8222622 0.38532631
## 4      sp195     SPpb 0.8231623 0.22192967
## 5      sp92      SPqe 0.8691069 0.20669544
## 6      sp91      SPec 0.8812881 0.25736120
## 7      sp137     SPod 0.9022702 0.04370608
## 8      sp118     SPy  0.9359536 0.08191003
## 9      sp99      SPtc 0.9633163 0.07952830
## 10     sp159     SPad 0.9906191 0.02033739
```

This document demonstrates the use of the *sabinaHSBM* package for network reconstruction. By applying the full reconstruction method, we showcased how spurious and missing links can be effectively identified and addressed in complex networks.

Computing characteristics

```
# Show processing time and computer characteristics (optional)
end_time <- Sys.time()
cat("The processing time of this script took: ",
    round(as.numeric(difftime(end_time, start_time, units = "mins", 1))), "minutes\n")
```

```
## The processing time of this script took: 37 minutes
```

This analysis was performed on a Dynabook with the following characteristics:

- Processor (CPU): Intel Core i7-1165G7 @ 2.80GHz (11th Gen)
- Memory (RAM): 32 GB
- Operating System: Windows 11 Pro, Version 24H2
- R Version: 4.3.3