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

## Example on how to use it to identify missing links with the binary_classifier method

The *sabinaHSBM* package provides tools for link prediction and network reconstruction using hierarchical stochastic block models (HSBM). This document demonstrates a simple use case based on the `binary_classifier` method to identify missing links. We will use the example dataset `dat` 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 4.0.4) with all required R packages (listed below)
    - Python 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:
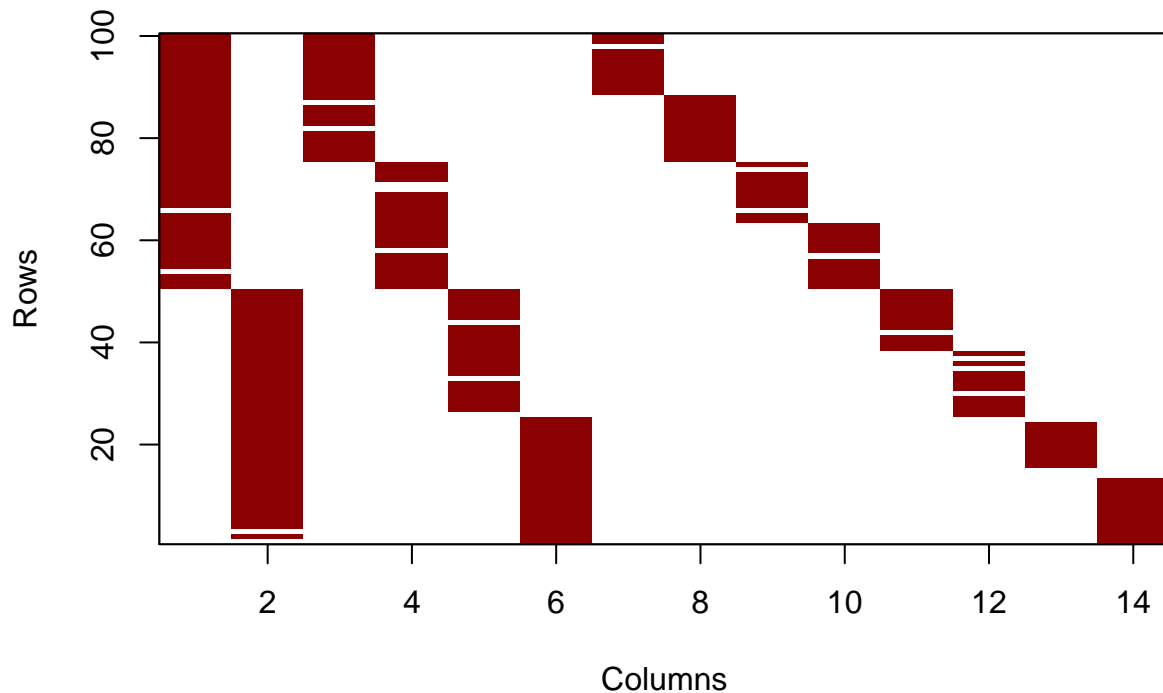
```
library(sabinaHSBM)
```

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

## Load data

The dataset `dat` 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 interaction between them while `0` (in black) represent lack of an observed interaction. The network contains gaps, indicating potential missing links.

```
# Load the dataset
data(dat, package = "sabinaHSBM")
```

```
# Plot the simulated matrix
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.

```r
# Prepare input object
myInput <- hsbm.input(
  dat,            # Binary bipartite matrix of observed links
  n_folds = 5     # Number of folds for cross-validation
)
```

```
## Warning in distribute_zero_sum_folds(folds_lst$com, folds_lst$held_pairs, :
## Nothing to do in distribute_zero_sum_folds. Return original held_pairs.
```

```
## [1] "Actual cross-validation rate is 0.199"
## [2] "Actual cross-validation rate is 0.199"
## [3] "Actual cross-validation rate is 0.199"
## [4] "Actual cross-validation rate is 0.202"
## [5] "Actual cross-validation rate is 0.199"
```

```r
# Summarizes network characteristics
summary(myInput)
```

```
##   n_rows n_cols n_links rows_single_link cols_single_link possible_links
## 1    100     14     277                1                0           1400
```

# Predicting Missing Links

The `hsbm.predict` function applies HSBM to predict link probabilities and group assignments in a network. The function works directly with the processed input created by `hsbm.input`. Here, we use the `binary_classifier` method. This method focuses on predicting probabilities for currently unobserved links (`0s`). Use this method if you want to obtain probabilities for all missing links (unobserved links likely to exist) in partially incomplete networks.

```
# Predict missing links using HSBM
myPred <- hsbm.predict(
  myInput,            # Input data processed by hsbm.input()
  iter = 1000,        # Number of iterations ...
  wait = 1000,        # Number of iterations for MCMC equilibration
  method = "binary_classifier",  # Prediction method
  save_blocks = TRUE,  # Save group assignments
  save_pickle = FALSE, # Save results as pickle files,
  save_plots = FALSE   # Save hierarchical edge bundling plots
)
```

```
## Computing predictions for fold  1
## Computing predictions for fold  2
## Computing predictions for fold  3
## Computing predictions for fold  4
## Computing predictions for fold  5
```

Predicted link probabilities and group assignments are stored for each fold. Below, we extract the probabilities (p) and groups 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  0 101 0.0004558870      sp1      SPB reconstructed
## 2  0 103 0.0325369377      sp1      SPD reconstructed
## 3  0 104 0.0001258867      sp1      SPE reconstructed
## 4  0 105 0.0002862439      sp1      SPF reconstructed
## 5  0 107 0.0123527549      sp1      SPH reconstructed
## 6  0 108 0.0091850525      sp1      SPI reconstructed
```

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 (e.g., 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(groups_fold1)
```

```
##    nodes G1 G2 G3 G4 names
## 1    100 19 65  1  1   SPA
## 2    101 19 65  1  1   SPB
```

```
## 3     102 19 65  1  1    SPC
## 4     103 19 65  1  1    SPD
## 5     104 19 65  1  1    SPE
## 6     105 19 65  1  1    SPF
## 7     106 19 65  1  1    SPG
## 8     107 26 65  1  1    SPH
## 9     108 48 65  1  1    SPI
## 10    109 48 65  1  1    SPJ
## 11    110 48 65  1  1    SPK
## 12    111 78 65  1  1    SPL
## 13    112 78 65  1  1    SPM
## 14    113 78 65  1  1    SPN
```

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 are progressively aggregated, revealing broader patterns and relationships or communities.

## 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.

```r
# Network reconstruction
myReconst <- hsbm.reconstructed(
  myPred,                   # Predictions processed by hsbm.predict
  rm_documented = TRUE,   # Use of documented entries during validation
  threshold = "prc_closest_topright", # Binarization threshold
  new_matrix_method = "average_thresholded" # Combine fold predictions
)
```

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

```r
# View the reconstructed network summary and evaluation metrics
summary(myReconst)
```

```
## $`Reconstructed Network Metrics`
##   obs_links unobs_links pred_links kept_links spurious_links missing_links
## 1       277        1123        288        277              0            11
##
## $`Evaluation Metrics`
##   mean_RLRR  mean_auc mean_aucpr  mean_yPRC mean_prec mean_sens mean_spec
## 1 0.4819481 0.9426747  0.5411246 0.04685102 0.6012766 0.5000649 0.9834372
##    mean_ACC   mean_ERR  mean_tss
## 1 0.9607882 0.03921182 0.4835022
```
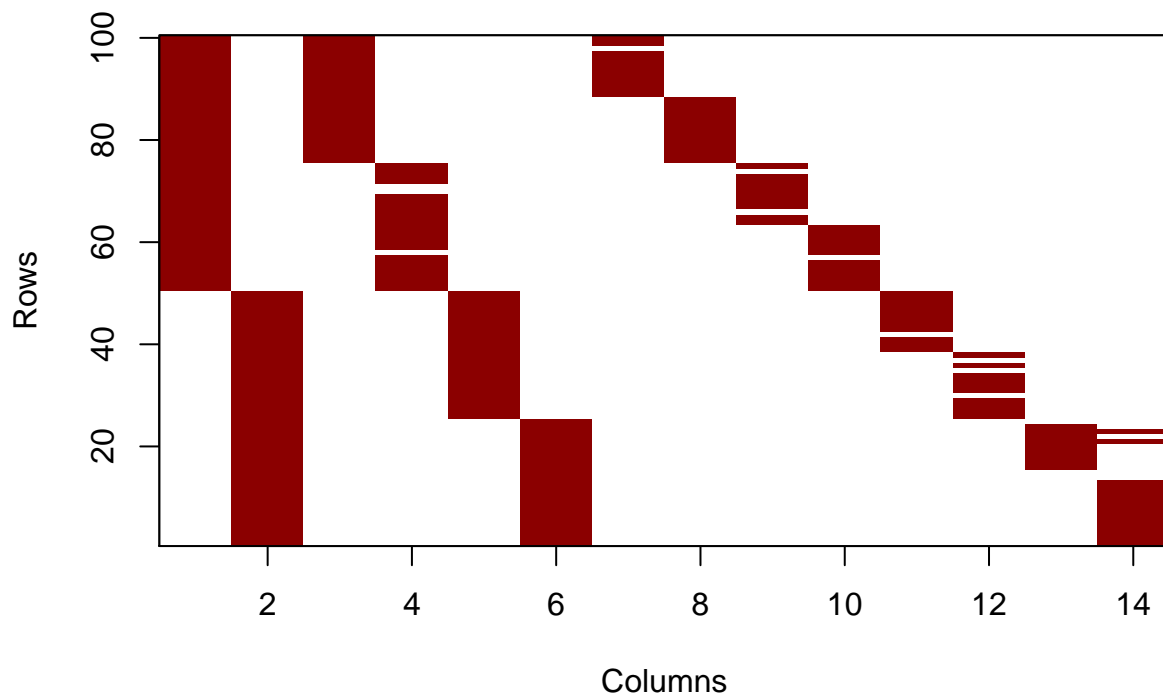
The summary provides the number of missing links, as well as key evaluation metrics, such as the retained link recovery rate (RLRR).

```r
# View the top potential missing links
top_links_df <- top_links(myReconst,
                          n = 10,
                          edge_type = "undocumented") # Type of edge to rank
print(top_links_df)
```

```
##     v1_names v2_names         p        sd
## 1      sp100      SPB 0.2328755 0.2137981
## 2       sp98      SPB 0.2067622 0.2214131
## 3       sp35      SPA 0.1894576 0.1329061
## 4       sp47      SPA 0.1574671 0.1615955
## 5       sp19      SPC 0.1455372 0.1748110
## 6       sp75      SPE 0.1407932 0.1505161
## 7       sp80      SPN 0.1298108 0.1578557
## 8       sp78      SPN 0.1265107 0.1624502
## 9       sp68      SPE 0.1239198 0.1482110
## 10      sp14      SPC 0.1222847 0.1274621
```

The `top-links functions` identifies the undocumented links most likely to be missing links in a network predicted by HSBM.

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



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

# Computing characteristics

```r
# Show processing time and computer characteristics (optional)
end_time <- Sys.time()
```

```r
cat("The processing time of this script took: ", end_time - start_time, "minutes\n")
```

```
## The processing time of this script took:  5.236398 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