

Stochastic Block Model for taxonomic identification in (meta)Barcoding

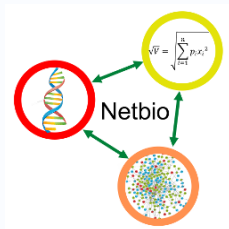
Mohamed Anwar ABOUABDALLAH¹

Directed by : Nathalie Peyrard² Alain Franc¹ Olivier Coulaud³

¹ INRAE , UMR BioGeCo, Pierroton & EPC INRAE / *Inria*-Pleiade, Talence, France

² INRAE , Unité MIAT, Toulouse, France

³ *Inria* , HiePACS, Talence, France



Summary

1 Agreement between botanical and molecular classifications

- Data set
- High taxonomics levels
- General approach
- Results

2 SBM model

- SBM Model parameters
- SBM possible estimation methods

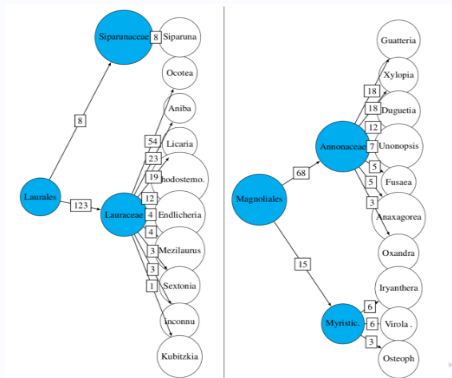
3 Tensor trains approximation for SBM estimation

- The idea
- How it works
- About marginals

Agreement between botanical and molecular classifications

Data set

- 1458 trees from an experimental plot in French Guyana.



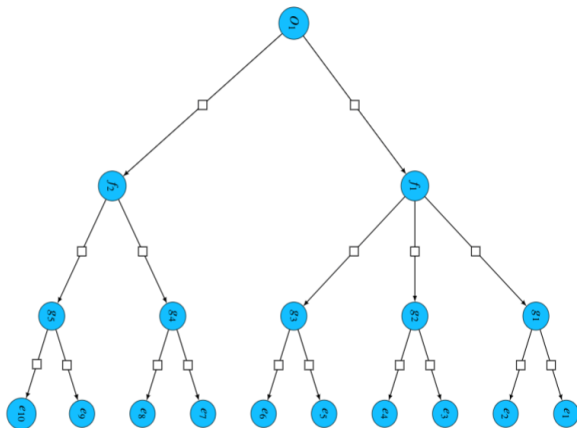
Order, family, genus and species of each individual.

```

>AVChB: 00887: 00504
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAACGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00155: 00049
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00283: 00713
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00255: 01902
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00370: 00050
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00413: 01350
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00412: 01921
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00626: 00400
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAACGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00603: 02589
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00907: 01162
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00975: 02843
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 00990: 02676
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 01122: 02070
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 01153: 00411
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 01157: 02390
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 01209: 00170
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
>AVChB: 01221: 00896
AGGTGAAGTAAAGGTTCTACTTAAACATCACTCGTGTACAATGGGAAGGTTTACACTCGTGCAGAATACGCTAAGGCCCTGGTTCGTAATTGTATGATCGATTTA
  
```

DNA sequence of each individual.

OTU

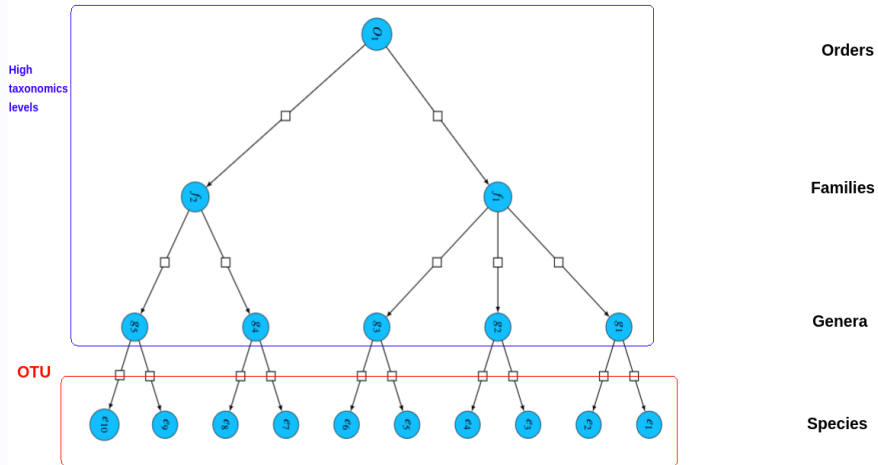


Orders

Families

Genera

Species



General approach

The three steps of the approach

- **Step 1** : Choice of sub-samples to study :
 - Our experiment: Selection

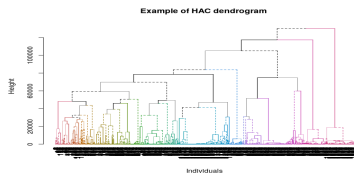
Taxonomic level	Sequences	Number of taxa	Minimal size
Species	313	55	5
Genera	845	36	10
Families	1349	30	10
Orders	1357	11	15

General approach

The three steps of the approach

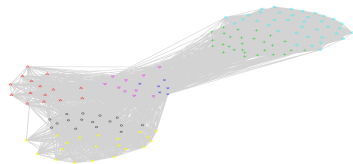
- **Step 2** : For each sub-sample, building partitions with three methods for each sub-sample and with Smith Waterman and kmer dissimilarities :

- **M₁ : Agglomerative Hierarchical Clustering (AHC)**



- **M₂ : Stochastic Block Model (SBM).**

Graph représentant toutes les classes avec layout FL



General approach

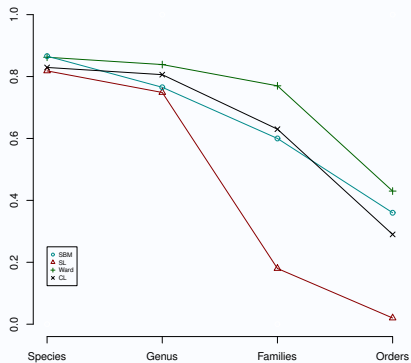
The three steps of the approach

- **Step 3** : Comparing the classifications two by two
 - Using visual tools
 - Using NMI to characterize the adequacy/independence

Results as a function of taxonomic level

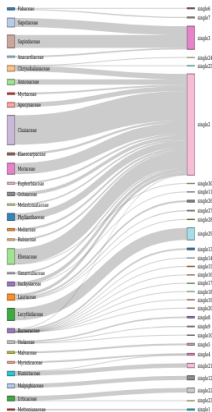
Evolution of NMI index as a function of Taxonomic levels:

kmer based distances

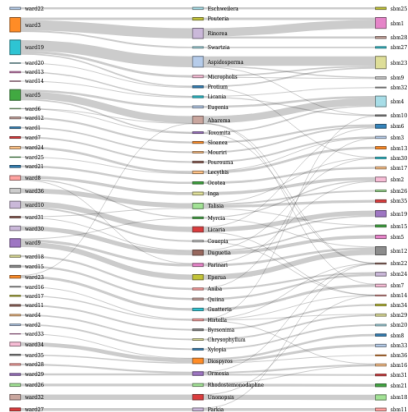


Sankey plots for genera

Botanics WRT SL

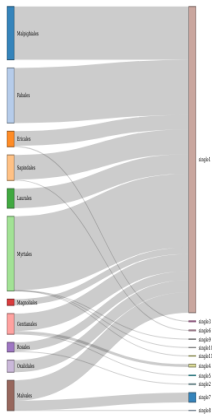


Ward WRT Botanics WRT SBM

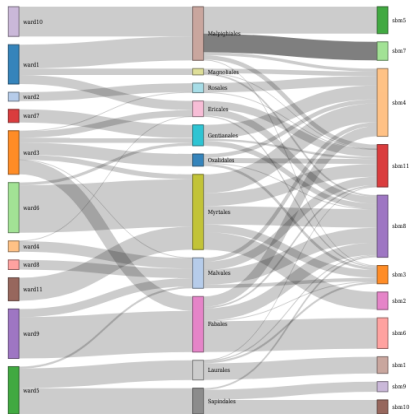


Sankey plots for orders

Botanics WRT SL



Ward WRT Botanics WRT SBM



In conclusion : Interest of SBM models

The main advantage between AHC and SBM :

- AHC produces community w.r.t. SBM produces classes (not necessary communities)
- Outputs of SBM are : Classes and Λ , distance matrix between classes.

Let's talk about Λ :

Case 1 :

$$\Lambda = \begin{pmatrix} 2 & 9 & 11 \\ 6 & 3 & 7 \\ 8 & 5 & 1 \end{pmatrix}$$

There are 3 communities \implies SBM \simeq CAH

In conclusion : Interest of SBM models

Case 2 :

$$\Lambda = \begin{pmatrix} 22 & 9 & 11 \\ 6 & 3 & 7 \\ 8 & 5 & 1 \end{pmatrix} \quad \Lambda = \begin{pmatrix} 2 & 9 & 11 \\ 6 & 23 & 7 \\ 8 & 5 & 1 \end{pmatrix} \quad \Lambda = \begin{pmatrix} 2 & 9 & 11 \\ 6 & 3 & 7 \\ 8 & 5 & 19 \end{pmatrix}$$

There are 2 communities \implies SBM (warning) \neq CAH

Case 3 :

$$\Lambda = \begin{pmatrix} 22 & 9 & 11 \\ 6 & 23 & 7 \\ 8 & 5 & 19 \end{pmatrix}$$

There are no communities \implies SBM (warnings) \neq CAH

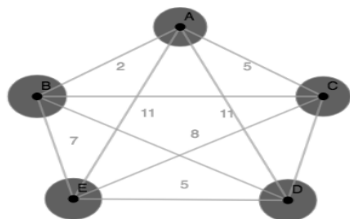
SBM model

Model intuition

Data set

$$\begin{pmatrix} 0 & 2 & 5 & 11 & 11 \\ \vdots & 0 & \ddots & \ddots & 7 \\ \vdots & \vdots & \ddots & \ddots & 8 \\ \vdots & \ddots & \ddots & 0 & 5 \\ \vdots & \dots & \dots & \dots & 0 \end{pmatrix}$$

Representation as graph



Parameters presentation

- $D_{i,j} | z_{i,b} = 1, z_{j,b'} = 1 \sim \text{Pois}(\lambda_{b,b'})$
- $\Lambda \in \mathbb{M}_{q,q}, \lambda_{b,b'}$: The parameter of Poisson probability to have a distance d between a vertex of class b and a vertex of class b' .

$$\forall b, b' = 1, \dots, q, \lambda_{b,b'} = z_b^T \Lambda z_{b'}$$

$$\Lambda = \begin{pmatrix} 2 & 9 & 11 \\ 6 & 3 & 7 \\ 8 & 5 & 1 \end{pmatrix}$$

- $\alpha \in [0, 1]^q, \alpha_i$: The probability to belong to class i .

EM algorithm

$$Z = (z_1, \dots, z_n)^T \in \mathbb{M}_{n,B}([0, 1])$$



- estimating Z needs to obtain $\hat{\theta} = (\hat{\alpha}, \hat{\Lambda})$ we proceed by $\hat{\theta}_{mv} = \operatorname{argmax}(P(D|\theta))$
- The most natural way is the EM algorithm. Each iteration involves two steps :
 - E-step** : Compute : $Q(\theta, \theta^t) = \mathbb{E}_Z[\log P_\theta(D|Z)|\theta^t, D]$
 - M-step** : $\theta^{(t+1)} = \operatorname{argmax}_\theta Q(\theta, \theta^t)$

⚠ The main difficulty of the EM algorithm is to compute the marginals:
It requires q^{n-1} sums.

Estimation approach

There are two main classes of methods :

Monte-Carlo methods : characterize a distribution by randomly sampling values from the distribution.

- + Precision : Accurate
- Computation time : Slow

Monte-Carlo methods : characterize a distribution by randomly sampling values from the distribution.

Precision : Accurate

Computation time : Slow

Estimation approach

There are two main classes of methods :

Monte-Carlo methods : characterize a distribution by randomly sampling values from of the distribution.

- + Precision : Accurate
- Computation time : Slow

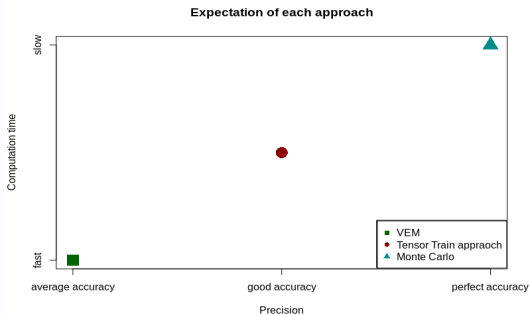
Variational methods: we assume the independence of the nodes knowing the graph to approximate marginal by mean field.

- Precision : Average accuracy
- + Computation time : Fast

Estimation approach

- The ideal approach consists on having good to perfect accuracy and fast computation time

Expectation :



How to get closer of this approach ?

Tensor trains approximation for SBM estimation

The idea

Let's Ψ express a graphical model,

$$\Psi[z_1, \dots, z_n] = \prod_{\ell=1, n} \psi_{A_\ell}(Z)$$

Mean field

$$\Psi : [1, q] \longrightarrow \mathbb{R}$$

$$z_1, \dots, z_n \longrightarrow \Psi[z_1, \dots, z_n]$$

Tensor Train

$$\Psi \in (\mathbb{R}^q)^{\otimes n}$$

$$\Psi[z_1, \dots, z_n] \in \mathbb{R}$$

The idea

- **The mean field** which is approximation by products of scalar (rank =1).

Smartness

$$x_{ij} = a_i b_j$$

$$\sum_{i,j} x_{ij} = \sum_{i,j} a_i b_j = \left(\sum_i a_i \right) \left(\sum_j b_j \right)$$

n^2 products 1 product

- **Tensor Train** which is approximation by matrix products (rank >1).

Smartness

$$\sum_{i,j} \mathbf{a}_i \mathbf{b}_j = \left(\sum_i \mathbf{a}_i \right) \cdot \left(\sum_j \mathbf{b}_j \right)$$

$\mathbf{a}_i | \mathbf{b}_i$ 1 matrix product

The idea

- **The mean field** which is approximation by products of scalar (rank =1).

Smartness

$$x_{ij} = a_i b_j$$

$$\sum_{i,j} x_{ij} = \sum_{i,j} a_i b_j = \left(\sum_i a_i \right) \left(\sum_j b_j \right)$$

n^2 products 1 product

- **Tensor Train** which is approximation by matrix products (rank >1).

Smartness

$$\sum_{i,j} \mathbf{a}_i \mathbf{b}_j = \left(\sum_i \mathbf{a}_i \right) \cdot \left(\sum_j \mathbf{b}_j \right)$$

$\mathbf{a}_i | \mathbf{b}_i$

1 matrix product

The idea

Tensor Train

- The tensor train format of the joint probability distribution :

$$\forall z_1, \dots, z_n, \Psi[z_1, \dots, z_n] = A_1[z_1] \cdot A_2[z_2] \dots A_{n-1}[z_{n-1}] \cdot A_n[z_n]$$

$$\Psi[z_1, \dots, z_n] = \text{---} \cdot \square \dots \square \cdot |$$

$A_i[z_i]$: The cores of Ψ (matrix)

- This format allows variable separation, suitable for marginals

How ?

How to compute $A_i[z_i]$ matrix ?

Novikov approach

- compute the TT-approximation of each factor
- then, it uses the Kronecker mixed-product property to compute the TT-approximation of Ψ from the TT approximation of each factor

Improvement of this approach :

- The matrix $A_i[z_i]$ are sizes $O(q^n)$

A way to deal with this is to use the TT Format that allow us a storage

- + who requires much less memory space
- + who can be used for matrices
- + with efficient operations

Computation of marginals

- Partition function :

$$\begin{aligned} W &= \sum_{z_1, \dots, z_n} A_1[z_1] \dots A_n[z_n] \\ &= \left(\sum_{z_1} A_1[z_1] \right), \dots, \left(\sum_{z_n} A_n[z_n] \right) \end{aligned}$$

Let's $B_i = \sum_{z_i} A_i[z_i]$

$$W = B_1 \times \dots \times B_n$$

Computation of marginals

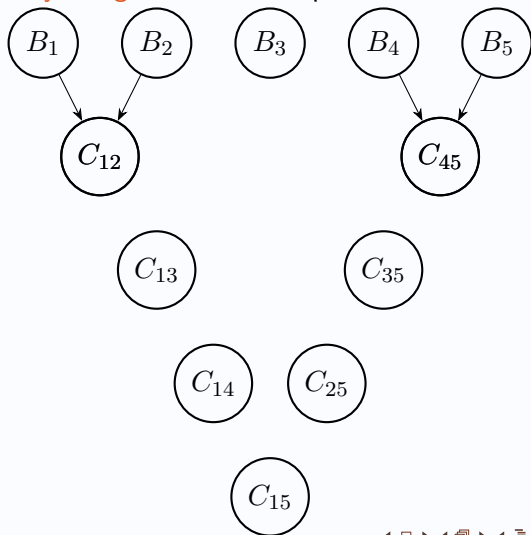
Let's $B_i = \sum_{z_i} A_i[z_i]$

- **Binary marginals** : $\forall i, j \in [1, n]^2, p_{i,j}(z_i, z_j)$

$$\underbrace{B_1 \times \dots \times B_{i-1}}_{\text{left}} A_i[z_i] \underbrace{B_{i+1} \times \dots \times B_{j-1}}_{\text{center}} A_j[z_j] \underbrace{B_{j+1} \times \dots \times B_n}_{\text{right}}$$

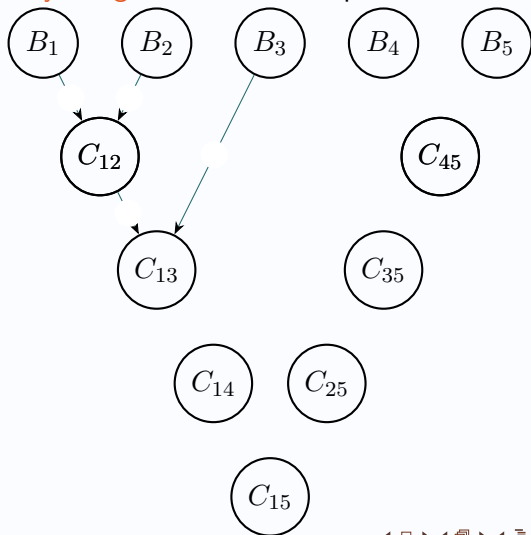
Computation of marginals

- Unary and Binary marginals : First step :



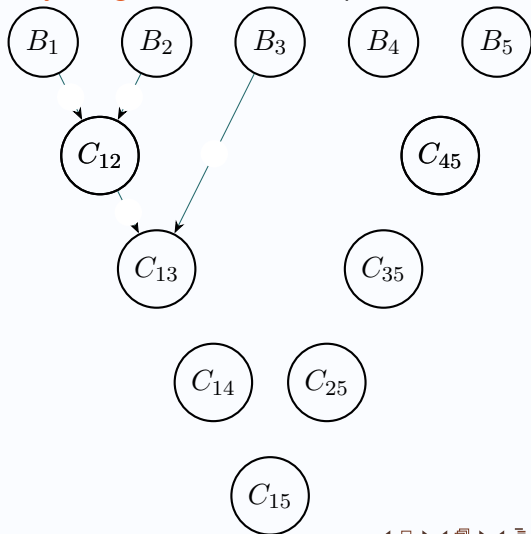
Computation of marginals

- Unary and Binary marginals : Second step : Left side



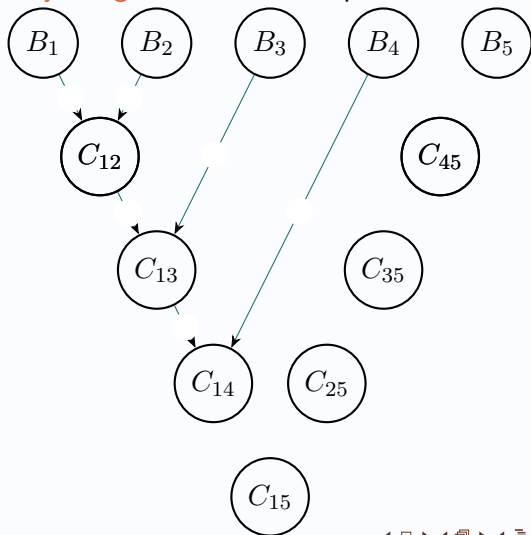
Computation of marginals

- Unary and Binary marginals : Second step : Left side



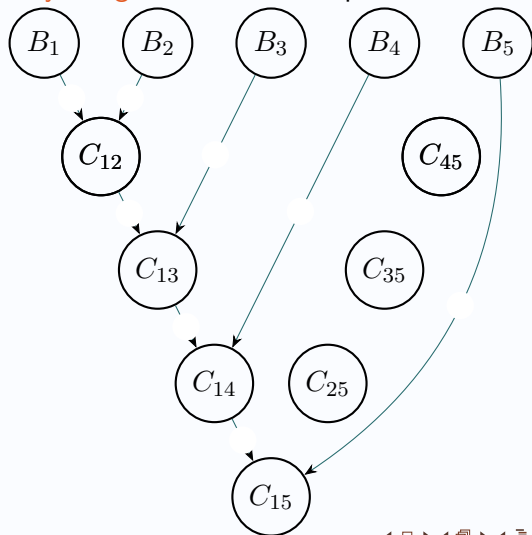
Computation of marginals

- Unary and Binary marginals : Second step : Left side



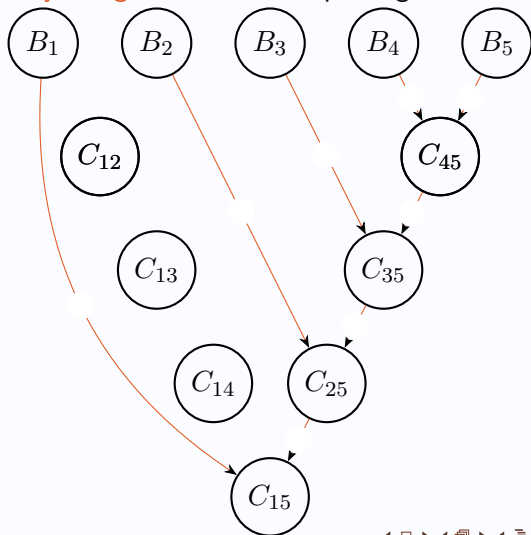
Computation of marginals

- Unary and Binary marginals : Second step : Left side



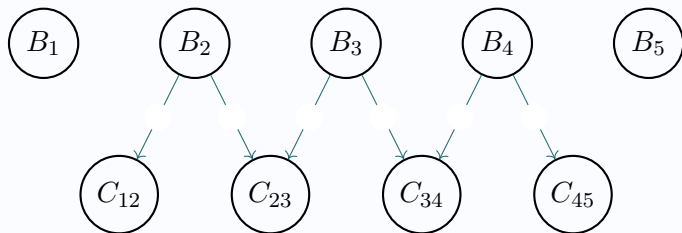
Computation of marginals

- Unary and Binary marginals : Third step : Right side



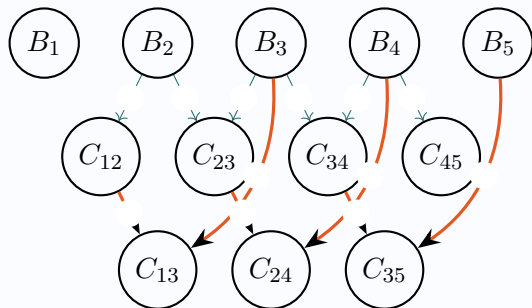
Computation of marginals

- Binary marginals : Fourth step Centers



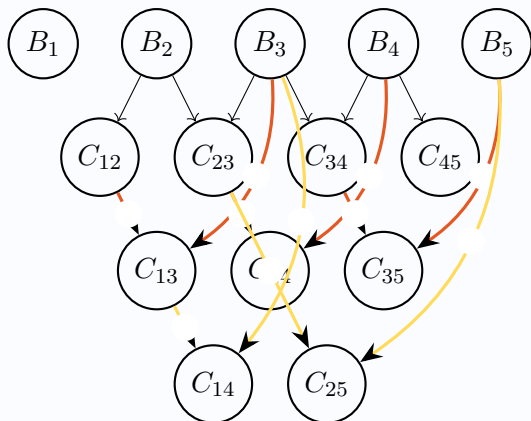
Computation of marginals

- Binary marginals : Fourth step Centers



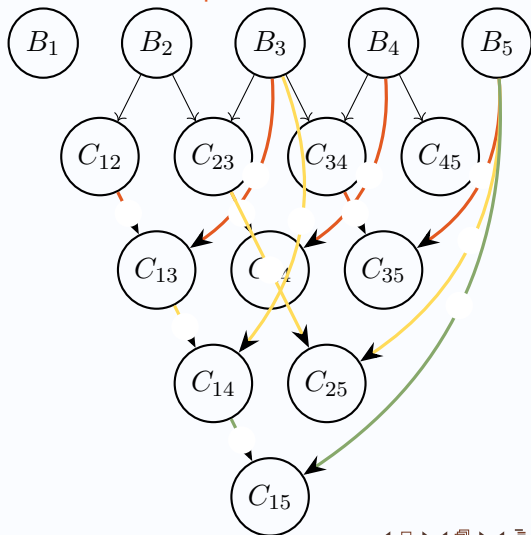
Computation of marginals

- Binary marginals : Fourth step Centers



Computation of marginals

- Binary marginals : Fourth step Centers



Conclusion and prospect

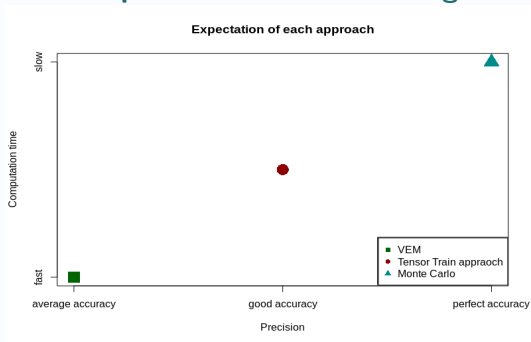
During this presentation, we showed that

- SBM models can be used for metabarcoding
- that we can potentially improve their accuracy thanks to a tensor based approach

Conclusion and prospect

A legitimate question remains unanswered:

How to scale up the SBM model to large datasets ?



E

T

H

N

D

!

A

N

K

Y

U

O

