

Approximation variationnelle pour des HMM multivariés en vue de la détection de CNV chez le maïs

Xiaoqiang WANG, Emilie Lebarbier, Julie Aubert & Stéphane Robin

UMR de Génétique Végétale INRA/Paris Sud/CNRS & UMR Mathématiques et
Informatiques Appliquées INRA/AgroParisTech

MSTGA, May 17 2013, Toulouse



1

Introduction

- Biological context
- Methodology
- Problems

2

Multivariate HMM analysis for dependency structure I

- Model
- Inference
- Simulation
- Application

3

Multivariate HMM analysis for genetic dependency structure II

- Model
- Inference
- Simulation

4

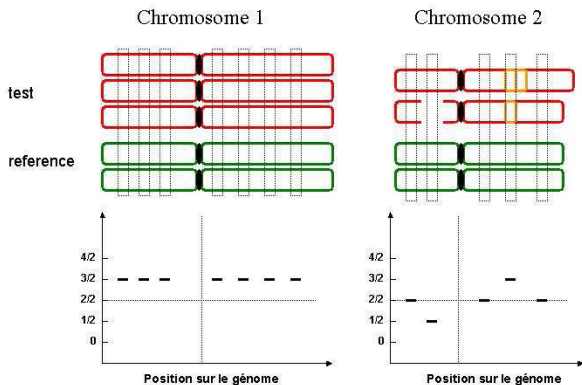
Conclusion

- **Copy Number Variant (CNV)**

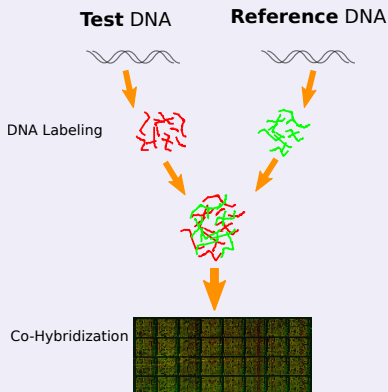
- Alterations of the number of copies of one or more sections of the DNA in genome

- **Technology** : Comparative Genomic Hybridization (**CGH**)

- **CGH in its principle**



CGH: CNV (Test/Reference)



Signal at position t : $Y(t)$



- X: Genomic position
- Y: \log_2 (# Test / # Reference)

- **Objectif:** Classify the genomic regions as, -1 ("deleted"), 0 ("normal"), 1 ("amplified").

- **Notations**

- Hidden status : $\{S_t\} \sim MC(\pi)$, $\pi_{k\ell} = \mathbb{P}(S_t = \ell | S_{t-1} = k)$
- Observed data : $\{Y_t\}$ are independent conditionally to \mathbf{S} .

- **Model**

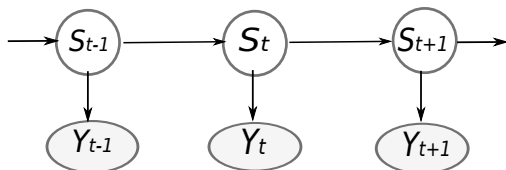
$$Y_t = \sum_k S_t^k \mu_k + \varepsilon_t,$$

where $S_t^k = \mathbb{1}_{\{S_t=k\}}$, $\varepsilon_t \sim \mathcal{N}(0, \sigma^2)$.

- **HMM** denoted by $\theta = (m, \pi, \gamma)$

- m : initial state distribution
- π : transition probability for Markov chain (S_t)
- γ : parameter in emission probability, *i.e.*, $\gamma = (\mu, \sigma^2)$

- **Graphical representation**



- **Complete likelihood**

$$\mathbb{P}(\mathbf{Y}, \mathbf{S}) = \mathbb{P}(\mathbf{S})\mathbb{P}(\mathbf{Y}|\mathbf{S})$$

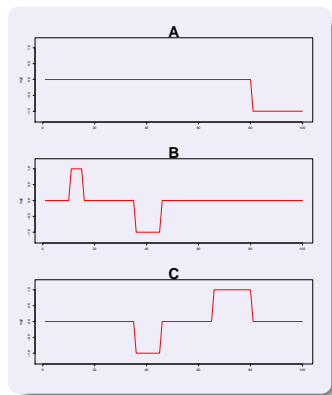
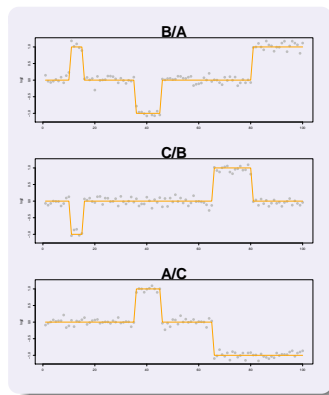
$$= \left\{ \prod_k m_k^{S_1^k} \prod_{t>1} \prod_{k,\ell} \pi_{k\ell}^{S_{t-1}^k S_t^\ell} \right\} \left\{ \prod_t \prod_\ell \phi(Y_t; \gamma_k) S_t^\ell \right\}$$

- **Objectif** : Infer $\mathbf{S}|\mathbf{Y}$

- **Inference** : E-M algorithm

- E-step: calculate $\mathbb{P}(\mathbf{S}|\mathbf{Y})$ by Forward-Backward
- M-step: estimate $\theta = \arg \max_{\theta} \mathbb{E}[\log \mathbb{P}(\mathbf{S}, \mathbf{Y})|\mathbf{Y}]$

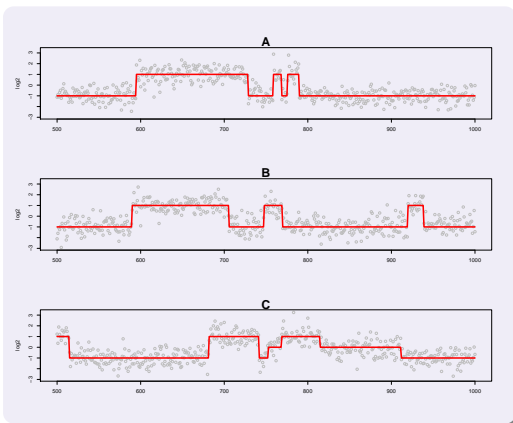
Question: Dependency structure I



Question?

- Perform a joint analysis taking into account the dependency structure due to the experiments
- Detect CNV for lines through their comparisons

Question: Genetic dependency structure II



- The lines are genetically correlated, e.g.,
 $\text{cor}(A, B) = 0.8$
 $\text{cor}(A, C) = 0.2$
 $\text{cor}(B, C) = 0.2$

Question?

- Perform a joint analysis taking into account the genetic dependency structure among lines A, B and C

1

Introduction

- Biological context
- Methodology
- Problems

2

Multivariate HMM analysis for dependency structure I

- Model
- Inference
- Simulation
- Application

3

Multivariate HMM analysis for genetic dependency structure II

- Model
- Inference
- Simulation

4

Conclusion

Multivariate HMM Model

• Notations

- I : number of lines; M : number of comparisons;
- Y_t : observation ;
- $S_{i,t}$: hidden status and $S_{i,t}^q = \mathbb{1}_{\{s_{i,t}=q\}}$

• Model

$$\begin{array}{ccccc} Y_t & = & D & \eta_t & + & \varepsilon_t, \\ M \times 1 & & M \times I & I \times 1 & & \end{array}$$

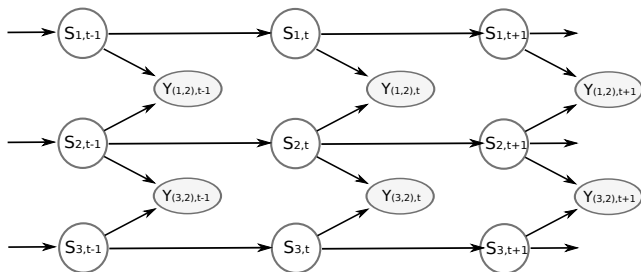
- $\eta_t = (\eta_{1,t}, \dots, \eta_{I,t})^T$ with $\eta_{i,t} = \sum_q \mu_q S_{i,t}^q$,
- D : design matrix, e.g.,

$$\begin{pmatrix} Y_{B/A} \\ Y_{C/B} \\ Y_{A/C} \end{pmatrix} = \begin{pmatrix} -1 & 1 & 0 \\ 0 & -1 & 1 \\ 1 & 0 & -1 \end{pmatrix} \begin{pmatrix} \eta_A \\ \eta_B \\ \eta_C \end{pmatrix}$$

- $\varepsilon_t \sim \mathcal{N}(\mathbf{0}, \sigma^2 I_M)$

Multivariate HMM Model

Graphical representation

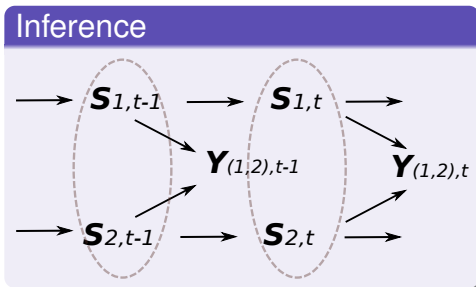


- $S_{1,t}$: hidden status at position t for line 1, e.g. -1, 0, 1.
- $Y_{(1,2),t}$: observation of $\log_2(\# \text{ line 1} / \# \text{ line 2})$ at position t .

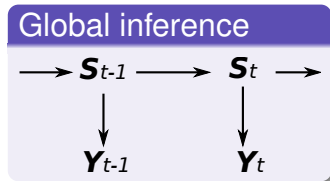
Inference

- Small I : Exact inference by E-M algorithm
- Big I : Approximate inference by variational E-M algorithm

Inference when I is small



\Rightarrow



$$\prod_{i,t>1} \prod_{q,r} \pi_{qr}^{S_{i,t-1}^q S_{i,t}^r} \prod_{(i,j),t} \prod_{q,r} \phi(Y_{(i,j),t}; \mu_{qr}, \sigma^2)^{S_{i,t}^q S_{j,t}^r} \Rightarrow$$

$$\prod_{t>1} \prod_{k,\ell} \pi_{k\ell}^{S_{t-1}^k S_t^\ell} \prod_t \prod_\ell \phi(Y_t; \gamma_k)^{S_t^\ell}$$

- For example

$$(S_{1,t}, S_{2,t}) = (1, 1) \Leftrightarrow S_t = 1$$

$$(S_{1,t}, S_{2,t}) = (1, 2) \Leftrightarrow S_t = 2$$

Inference when I is large

- **Question** : $\mathbb{P}(S_t|Y)$ can not be computed when I is large, e.g, if $I = 10$ and $Q = 3$, S_t has $3^{10} = 59049$ status to be considered.
- **Solution** : $\tilde{\mathbb{P}}(S) \cong \mathbb{P}(S|Y)$
- **Reminder** : Kullback-Leibler divergence

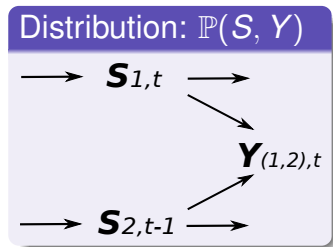
$$KL \left[\tilde{\mathbb{P}}(S) \parallel \mathbb{P}(S|Y) \right] = \int \tilde{\mathbb{P}}(S) \log \frac{\tilde{\mathbb{P}}(S)}{\mathbb{P}(S|Y)} dS$$

- KL is always non-negative
- Null iff $\tilde{\mathbb{P}} = \mathbb{P}$

- $$\tilde{\mathbb{P}}(S) = \arg \min_{\tilde{\mathbb{P}} \in \mathcal{D}} KL \left[\tilde{\mathbb{P}}(S) \parallel \mathbb{P}(S|Y) \right]$$

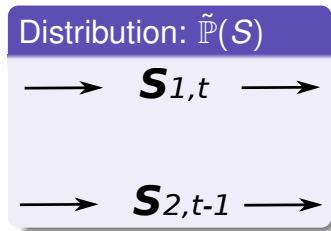
Inference when T is large

- Graphical representation



$KL[\tilde{\mathbb{P}}||\mathbb{P}]$

\iff



$$\prod_t \mathbb{P}(S_t | S_{t-1}) \mathbb{P}(Y_t | S_t)$$

$$\prod_i \prod_t \tilde{\mathbb{P}}(S_{i,t} | S_{i,t-1})$$

- Key** : find $\tilde{\mathbb{P}}(S_{i,t} | S_{i,t-1})$ to break the dependency

Proposition

Let denote $p_{itqr} = \tilde{\mathbb{P}}(S_{i,t} = r | S_{i,t-1} = q)$, then we obtain a set of fixed point equations for p_{itqr} :

$$p_{itqr} \propto \pi_{qr} \prod_{j \in m^-(i), v} \phi(Y_{(i,j),t}; \mu_{rv}, \sigma^2)^{\mathbb{E}_{\tilde{\mathbb{P}}}(S_{j,t}^v)} \\ \times \prod_{j \in m^+(i), u} \phi(Y_{(j,i),t}; \mu_{ur}, \sigma^2)^{\mathbb{E}_{\tilde{\mathbb{P}}}(S_{j,t}^u)}$$

where π is the transition probability and ϕ emission probability.

- Inference by variational E-M algorithm

- **Variational E-step:** find

$$\tilde{\mathbb{P}}^{(h)}(\mathbf{S}) = \arg \min_{\tilde{\mathbb{P}} \in \mathcal{P}} KL \left[\tilde{\mathbb{P}}(\mathbf{S}) \parallel \mathbb{P}(\mathbf{S} | Y, \theta^{(h)}) \right],$$

$$\text{where } \mathcal{P} = \left\{ \tilde{\mathbb{P}}(\mathbf{S}) \mid \tilde{\mathbb{P}}(\mathbf{S}) \propto \prod_i \prod_t \tilde{\mathbb{P}}(\mathbf{S}_{i,t} | \mathbf{S}_{i,t-1}) \right\}$$

- **Variational M-step:** estimate $\theta = (m, \pi, \mu, \sigma^2)$ as

$$\hat{\theta}^{(h+1)} = \arg \max_{\theta} \mathbb{E}_{\tilde{\mathbb{P}}^{(h)}} [\log \mathbb{P}(Y, \mathbf{S} | \theta)].$$

- $\hat{m}_q, \hat{\pi}_{qr}$

$$\hat{m}_q \propto \sum_i \mathbb{E}_{\tilde{\mathbb{P}}}(\mathbf{S}_{i,1}^q) \quad \hat{\pi}_{qr} \propto \sum_{i,t} \mathbb{E}_{\tilde{\mathbb{P}}}(\mathbf{S}_{i,t-1}^q \mathbf{S}_{i,t}^r)$$

- $\hat{\mu}, \hat{\sigma}^2$

$$(\hat{\mu}, \hat{\sigma}^2) = \arg \max_{\mu, \sigma^2} \sum_{t,(i,j)} \sum_{q,r} \mathbb{E}_{\tilde{\mathbb{P}}}(\mathbf{S}_{i,t}^q) \mathbb{E}_{\tilde{\mathbb{P}}}(\mathbf{S}_{j,t}^r) \log \phi(Y_{(i,j),t}; \mu_{qr}, \sigma^2)$$

Simulation and Discussion

1 Input

	Simulation I	Simulation II
I	3	5
M	6	20
μ	(-1,0,1)	(-1,0,1)
σ^2	0.36	0.36
T	5000	5000
Nb. Sim.	100	100

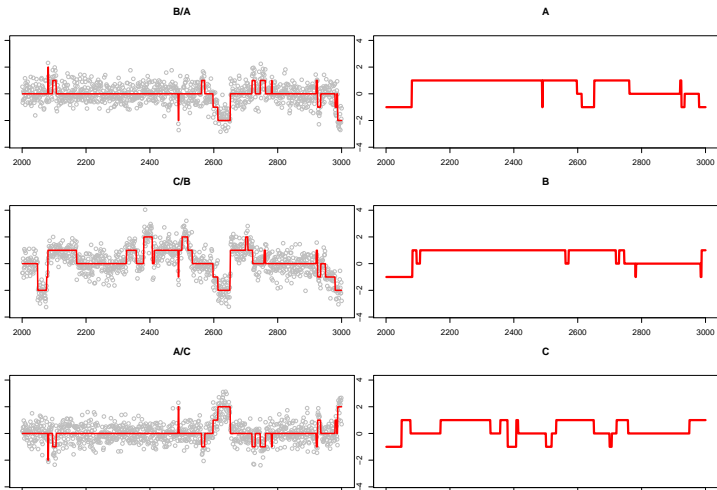
- 2 **Multivariate HMM:** perform a joint analysis with exact inference and variational inference, respectively.

3 Output

$$\hat{m}, \hat{\pi}, \hat{\mu}, \hat{\sigma}, \mathbb{P}(S_{i,t} = k | Y)$$

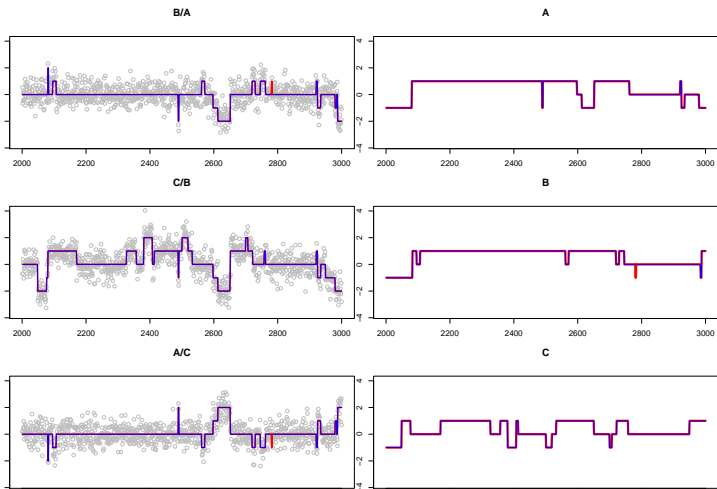
Simulation and Discussion

- One simulation :



Simulation and Discussion

- Result simulation : red (true), blue (estimation)



Simulation and Discussion

- **Time**

- Mean time for one multivariate analysis

	Exact	Variational
$l = 3$	58 s	2 s
$l = 5$	396 s	13 s

- **Accuracy**

- Error rate for comparison between lines, e.g. $B/A, A/C, \dots$

	Exact	Variational
$l = 3$	0.15%	1.1%

- Error rate for lines, e.g. A, B, \dots

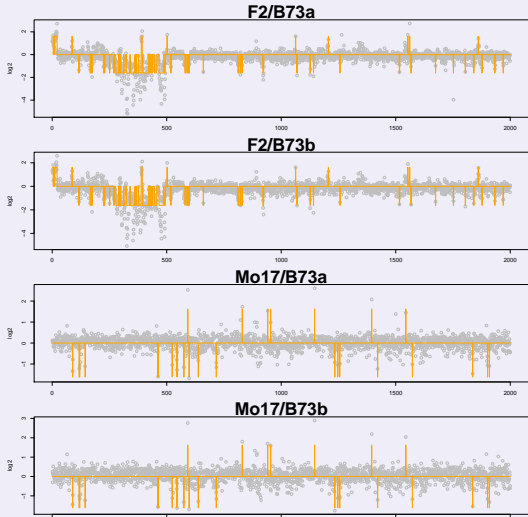
	Exact	Variational
$l = 3$	5.4%	29%

Biogemma data analysis

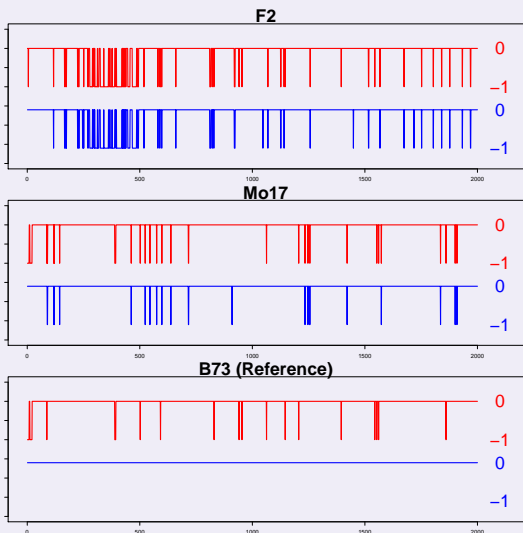
- Experimental design
 - **Lines:** F2, Mo17, B73 (**Reference**)
 - **Comparisons:** F2/B73a, F2/B73b, Mo17/B73a, Mo17/B73b
- Information on Chromosome
 - 2 139 527 probes (data point) on 10 chromosomes

Chromosome	1	2	3	4	5
Probe size	323410	252569	243428	260526	220582
Chromosome	6	7	8	9	10
Probe size	173328	182823	175177	156925	150759

- Joint analyses one by one chromosome



- Chromosome 1
- First 2000 probes
- Orange : mean



- Two joint analyses
 - Red : Multivariate HMM
 - Blue : Multivariate HMM but forcing B73 to be 0

- Different status

	Red \neq Blue
F2	0.9%
Mo17	0.45%
B73	0.65%
Total	0.67%

1

Introduction

- Biological context
- Methodology
- Problems

2

Multivariate HMM analysis for dependency structure I

- Model
- Inference
- Simulation
- Application

3

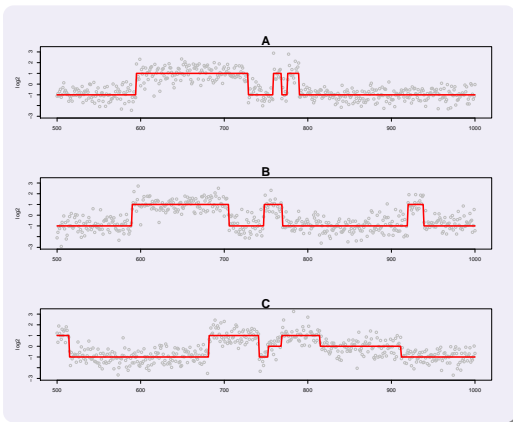
Multivariate HMM analysis for genetic dependency structure II

- Model
- Inference
- Simulation

4

Conclusion

Question: Genetic dependency structure II



- The lines are genetically correlated, e.g.,
 $\text{cor}(A, B) = 0.8$
 $\text{cor}(A, C) = 0.2$
 $\text{cor}(B, C) = 0.2$

Question?

- Perform a joint analysis taking into account the genetic dependency structure among lines A, B and C

Multivariate HMM Model

• Notations

- I : number of lines
- $X_{i,t}$: observation
- $S_{i,t}$: hidden status but $\forall i \neq j, (S_{i,t}, S_{j,t})$ are not independent

• Model

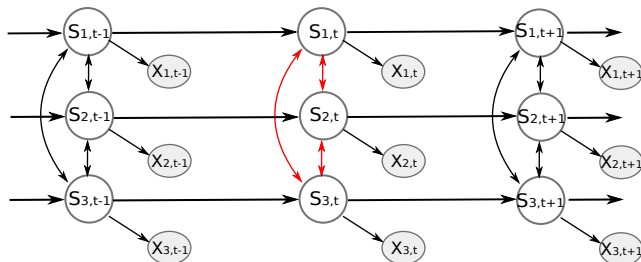
$$X_{i,t} = \sum_q S_{i,t}^q \mu_q + \varepsilon_t$$

- $S_{i,t}^q = \mathbb{1}_{\{S_{i,t}=q\}}$
- $\varepsilon_t \sim \mathcal{N}(0, \sigma^2 I_I)$

• Inference

- Small data: Exact inference by E-M algorithm
- Big data : Approximate inference by variational E-M algorithm

- Graphical representation

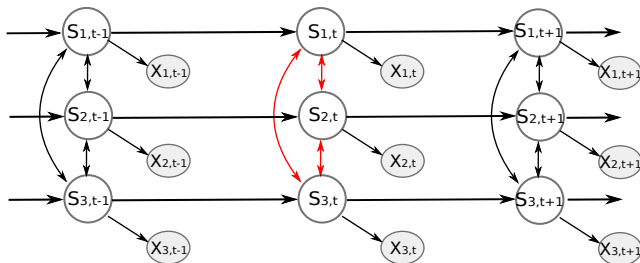


- Complete likelihood

$$\mathbb{P}(\mathbf{X}, \mathbf{S}) \propto \prod_k m_k^{S_1^k} \prod_{t>1} \prod_{k,l} \pi_{kl}^{S_{t-1}^k S_t^l} W_l^{S_t^l} \prod_t \prod_l \phi(X_t; \gamma_l)^{S_t^l}$$

where $W_l = \mathbb{P}(S_t = l)$.

- Graphical representation



- For the part of cyclic

Model

Let $W_\ell = \mathbb{P}(S. = \ell)$,

$$W_\ell \propto \prod_{i,j \neq i} \omega^{s_{ij} \mathbb{1}_{\{q_j^\ell \neq q_i^\ell\}}}$$

with $\omega < 1$, $(s_{ij})_{ij}$ is the similarity matrix.

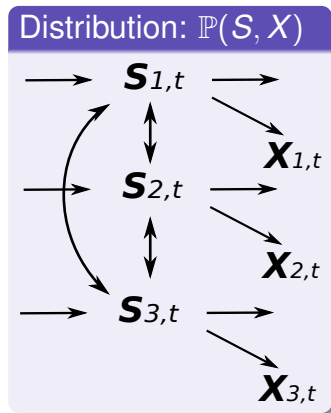
- **Small I** : exact inference with E-M algorithm is similar to that of dependency I.
- **Big I** : $\mathbb{P}(S|X)$ is not computable, find

$$\tilde{\mathbb{P}}(S) \cong \mathbb{P}(S|X)$$

in terms of Kullback-Leibler divergence.

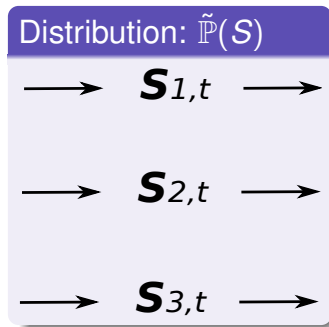
Inference when T is large

- Graphical representation



$KL[\tilde{\mathbb{P}}||\mathbb{P}]$

\iff



$$\prod_i \prod_t \tilde{\mathbb{P}}(S_{i,t} | S_{i,t-1})$$

Proposition

Let denote $p_{itqr} = \tilde{\mathbb{P}}(\mathcal{S}_{i,t} = r | \mathcal{S}_{i,t-1} = q)$, then we obtain a set of fixed point equations for p_{itqr} :

$$p_{itqr} \propto \pi_{qr} \omega^{\sum_{j \neq i} (1 - \mathbb{E}_{\tilde{\mathbb{P}}} S_{j,t}^r)} \phi(X_{i,t}, \mu_r, \sigma^2),$$

where π is the transition probability and ϕ emission probability.

- **Variational E-step:** find

$$\tilde{\mathbb{P}}(\mathcal{S}) = \arg \min_{\tilde{\mathbb{P}} \in \mathcal{P}} KL \left[\tilde{\mathbb{P}}(\mathcal{S}) \| \mathbb{P}(\mathcal{S} | X) \right],$$

where $\mathcal{P} = \left\{ \tilde{\mathbb{P}}(\mathcal{S}) \mid \tilde{\mathbb{P}}(\mathcal{S}) \propto \prod_i \prod_t \tilde{\mathbb{P}}(\mathcal{S}_{i,t} | \mathcal{S}_{i,t-1}) \right\}$;

Inference when I is large

- **Variational M-step:** estimate $\theta = (m, \pi, \mu, \sigma^2)$ as

$$\hat{\theta}^{(h+1)} = \arg \max_{\theta} \mathbb{E}_{\tilde{\mathbb{P}}^{(h)}} [\log \mathbb{P}(Y, \mathbf{S}|\theta)].$$

- \hat{m}

$$\hat{m}_q \propto \sum_i \mathbb{E}_{\tilde{\mathbb{P}}}(\mathbf{S}_{i,1}^q)$$

- $\hat{\pi}_{qr}$

$$\hat{\pi}_{qr} \propto \sum_{i,t} \mathbb{E}_{\tilde{\mathbb{P}}}(\mathbf{S}_{i,t-1}^q \mathbf{S}_{i,t}^r)$$

- $\hat{\mu}$

$$\hat{\mu}_r = \sum_{i,t} \mathbb{E}_{\tilde{\mathbb{P}}}(\mathbf{S}_{i,t}^r) X_{i,t} / \sum_{i,t} \mathbb{E}_{\tilde{\mathbb{P}}}(\mathbf{S}_{i,t}^r)$$

- $\hat{\sigma}^2$

$$\hat{\sigma}^2 = \sum_{i,t,r} \mathbb{E}_{\tilde{\mathbb{P}}}(\mathbf{S}_{i,t}^r) (X_{i,t} - \hat{\mu}_r)^2 / IT$$

Simulation and Discussion

1 Input

	Simulation I			Simulation II				
I	3			5				
S_{ij}	$\begin{pmatrix} 1 & 0.8 & 0.2 \\ 0.8 & 1 & 0.2 \\ 0.2 & 0.2 & 1 \end{pmatrix}$			$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0.24 & 0.32 & 0.4 \\ 0 & 0.24 & 1 & 0.48 & 0.6 \\ 0 & 0.32 & 0.48 & 1 & 0.8 \\ 0 & 0.4 & 0.6 & 0.8 & 1 \end{pmatrix}$				

- 2 **Multivariate HMM:** perform a joint analysis with exact inference and variational inference, respectively.

3 Output

$$\hat{m}, \quad \hat{\pi}, \quad \hat{\mu}, \quad \hat{\sigma}, \quad \mathbb{P}(S_{i,t} = q|X)$$

- **Time**

- Mean time for one joint analysis

	Exact	Variational
$l = 3$	48 s	0.2 s
$l = 5$	450 s	0.6 s

- **Accuracy**

- Error rate for lines, e.g. A, B, \dots

	Exact	Variational
$l = 3$	0.5%	4.5%

● Multivariate HMM

- Perform the joint analysis taking into account the dependency due to the CGH experiments
- Perform the joint analysis taking into account the dependency among plant lines
- Perform the joint analysis taking into account two above dependencies

● Inference

- When line size is large, the inference becomes impossible.
- Variational technique provides a fast algorithm.

● R package "MHMM" is built, and handles 16 different cases

- Paired / Unpaired for observation
- Dependent / Independent for hidden status
- Exact / Variational for inference
- Yes / No for reference

Thank you !

- UMR AgroParisTech/INRA

Emilie Lebarbier
Julie Aubert
Stéphane Robin

- UMR Moulon

Stéphane Nicolas
.....

- Biogemma

Jean-Philippe Pichon
.....