

Graphical Models:

Homework 3

Charles Darmon

Marin Robert De Beauchamp

Master MASH

I HMM - Implementation

(1)

We consider the following HMM model:

The chain q_t which has 4 possible states for $t=1 \dots T$

The observation $u_t = (x_t, y_t) \in \mathbb{R}^2$ for $t=1 \dots T$

An initial probability distribution $\pi \in \mathbb{R}^4$ such that $\sum_{i=1}^4 \pi_i = 1$

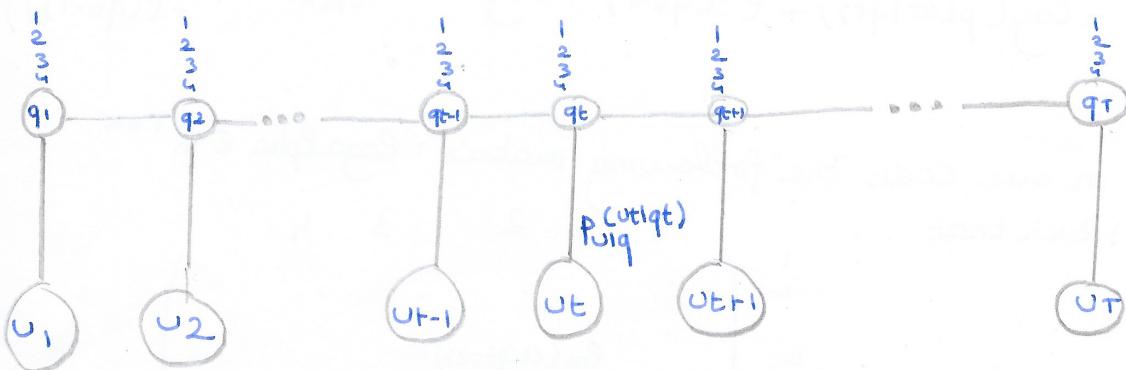
A probability transition matrix $A \in \mathbb{R}^{4 \times 4}$ defined by:

$A_{ij} = P(q_t=i | q_{t-1}=j) \quad \& \quad \sum_i A_{ij} = 1$ such that

$$A = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \left(\begin{array}{cccc} \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ p(q_t=3 | q_{t-1}=2) & \cdot & \cdot & \cdot \end{array} \right) \end{matrix}$$

The observations conditionally on the current states are obtained from Gaussian emission probabilities $u_t | q_t = i \sim N(\mu_i, \Sigma_i)$

The model can be modelled as following:



① We want to implement the α & β -recursions

②

As we have seen on course we have to implement the α & β recursions using logs in order to prevent numerical errors.

The alpha recursion is given by the following formula:

$$\alpha_t(q_t) = p(u_t | q_t) \sum_{q_{t-1}} p(q_t | q_{t-1}) \alpha(q_{t-1}) \text{ where } t=1, \dots, T$$

The first message is initialized as follows:

$$\alpha_1(q_1) = p_{u_1 \rightarrow q_1} c_{q_1} p(q_1) = p(u_1 | q_1) p(q_1) \text{ where:}$$

$$q_1 \sim \pi \text{ & } u_1 | q_1 = i \sim N(C_p, \Sigma_i)$$

- we will denote: $\ell(q_{t-1}) = \log(p(q_t | q_{t-1}) \alpha(q_{t-1}))$

- $c^* = \underset{i \in \{1, \dots, 4\}}{\operatorname{argmax}} \ell(p(q_t | q_{t-1}=i) \alpha(q_{t-1}=i))$

- $e^*(q_{t-1}) = \log(p(q_t | q_{t-1}=c^*) \alpha(q_{t-1}=c^*)) = \ell(q_{t-1}=c^*)$

Then:

$$\log(\alpha_t(q_t)) = \log(p(u_t | q_t)) + \log\left(\sum_{q_{t-1}} p(q_t | q_{t-1}) \alpha(q_{t-1})\right)$$

$$= \log(p(u_t | q_t)) + \log\left(\sum_{i=1}^4 \exp(e^*(q_{t-1}=i))\right)$$

$$= \log(p(u_t | q_t)) + e^*(q_{t-1}) + \log\left(1 + \sum_{i \neq c^*} \exp(e^*(q_{t-1}=i) - e^*(q_{t-1}))\right)$$

We create in our code the following matrix: $\text{logalpha} \in \mathbb{R}^{T \times K}$

(where $K=4$) such that :

$$\begin{matrix} & 1 & 2 & 3 & 4 \\ \vdots & \left(\begin{array}{cccc} \cdot & \cdot & \cdot & \cdot \\ \cdot & \log(\alpha_t(q_t=2)) & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{array} \right) \\ t & & & & \\ \vdots & & & & \\ T & & & & \end{matrix}$$

where $\log(\alpha_t(q_t=2)) = \log(p(u_t | q_t=2)) + e^*(q_{t-1}) + \log\left(1 + \sum_{i \neq c^*} \exp(e^*(q_{t-1}=i) - e^*(q_{t-1}))\right)$

\downarrow

$\sim N(C_2, \Sigma_2)$

The Beta recursion is given by :

③

$$\beta_t(q_t) = \sum_{q_{t+1}} p(u_{t+1}|q_{t+1}) p(q_{t+1}|q_t) \beta_{t+1}(q_{t+1})$$

The last message is given by $\beta_T(q_T) = 1$ and therefore passing on logarithmic scale : $\log(\beta_T(q_T=i)) = 0 \quad \forall i=1,..,4$

We will denote : $e(q_t) = \log(p(u_{t+1}|q_{t+1}) p(q_{t+1}|q_t) \beta_{t+1}(q_{t+1}))$

$$j^* = \underset{j \in \{1,..,4\}}{\operatorname{argmax}} e(q_{t+1}=j)$$

$$e^*(q_{t+1}) = e(q_{t+1}=j^*)$$

Then :

$$\begin{aligned} \log(\beta_t(q_t)) &= \log\left(\sum_{q_{t+1}} p(u_{t+1}|q_{t+1}) p(q_{t+1}|q_t) \beta_{t+1}(q_{t+1})\right) \\ &= \log\left(\sum_{j=1}^4 \exp(e(q_{t+1}=j))\right) \\ &= e^*(q_{t+1}) + \log\left(1 + \sum_{j \neq j^*} \exp(e(q_{t+1}=j) - e^*(q_{t+1}))\right) \end{aligned}$$

We create the following matrix : $\text{logbeta} \in \mathbb{R}^{T \times k}$ such that

$$\begin{matrix} & 1 & 2 & 3 & 4 \\ \vdots & \left(\begin{array}{cccc} \vdots & \vdots & \vdots & \vdots \\ \vdots & \log(\beta_t(q_t=2)) & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{array} \right) \\ t & & & & \end{matrix}$$

$$\begin{aligned} \text{where } \log(\beta_t(q_t=2)) &= e^*(q_{t+1}) + \log\left(1 + \sum_{j \neq j^*} \exp(e(q_{t+1}=j) - e^*(q_{t+1}))\right) \\ \text{and where } e(q_{t+1}=j) &= \log(p(u_{t+1}|q_{t+1}=j)) + \log(p(q_{t+1}=j|q_t=2)) \\ &\quad + \beta_{t+1}(q_{t+1}=j) \end{aligned}$$

$$\downarrow \sim N(\mu_j, \Sigma_j)$$

$$A_{j2}$$

We also need to compute:

(4)

$$\rightarrow p(q_t | u_1, \dots, u_T) = \frac{\alpha_t(q_t) \beta_t(q_t)}{\sum_{q_t} \alpha_t(q_t) \beta_t(q_t)}$$

We compute it with the log and denoting:

$$e(q_t) = \log(\alpha_t(q_t) \beta_t(q_t))$$

$$e^*(q_t) = \log(\alpha_t(q_t=i^*) \beta_t(q_t=i^*)) \text{ where } i^* = \arg\max_{i \in \{1, \dots, 4\}} \log(\alpha_t(q_t=i) \beta_t(q_t=i))$$

We obtain:

$$\log(p(q_t | u_1, \dots, u_T)) = e(q_t) - e^*(q_t) - \log(1 + \sum_{\substack{i=1 \\ i \neq i^*}}^4 \exp(e(q_t=i) - e^*(q_t)))$$

We create a matrix "proba" $\in \mathbb{R}^{T \times K}$ such that:

$$\begin{matrix} & 1 & 2 & 3 & 4 \\ \begin{matrix} 1 \\ \vdots \\ T \end{matrix} & \left(\begin{array}{cccc} \cdot & \cdot & \cdot & \cdot \\ \cdot & \ddots & [p(q_t=2 | u_1, \dots, u_T)] & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{array} \right) \end{matrix}$$

(We pass to the exponential when we add the proba in the matrix)

$$\rightarrow p(q_t, q_{t+1} | u_1, \dots, u_T) = \frac{\alpha_t(q_t) \beta_{t+1}(q_{t+1}) p(q_{t+1} | q_t) p(u_{t+1} | q_{t+1})}{p(u_1, \dots, u_T)}$$

We recall that $p(u_1, \dots, u_T) = \sum_{q \in E} \alpha_t(q_t) \beta_t(q_t)$

We compute the joint probability with the log and we denote as for the computation of $\log(p(q_t | u_1, \dots, u_T))$: $e(q_t) = \log(\alpha_t(q_t) \beta_t(q_t))$

Then we obtain:

$$\begin{aligned} \log(p(q_t, q_{t+1} | u_1, \dots, u_T)) &= \log(\alpha_t(q_t)) + \log(\beta_{t+1}(q_{t+1})) + \log(p(q_{t+1} | q_t)) + \log(p(u_{t+1} | q_{t+1})) \\ &\quad - e^*(q_t) - \log(1 + \sum_{\substack{i=1 \\ i \neq i^*}}^4 \exp(e(q_t=i) - e^*(q_t))) \end{aligned}$$

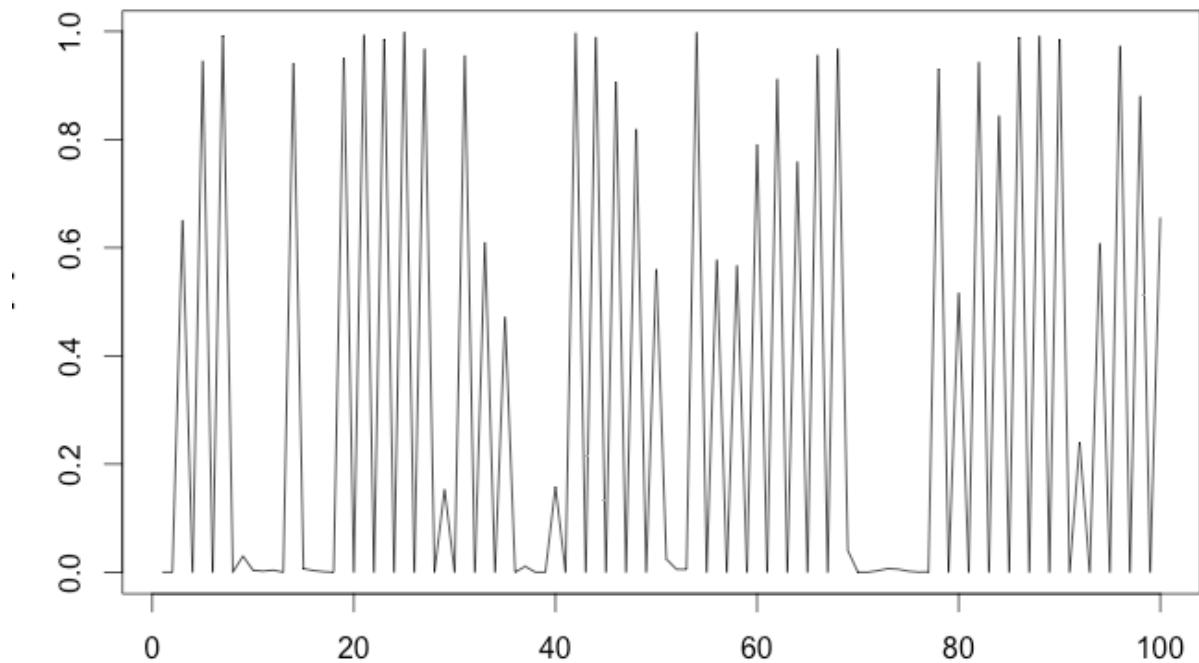
We create a list of matrices where the t th element is $p(q_t, q_{t+1} | u_1, \dots, u_T) \in \mathbb{R}^{4 \times 4} \quad (t \in \{1, \dots, T-1\})$ defined such that:

$$\begin{matrix} & 1 & 2 & 3 & 4 \\ \begin{matrix} 1 \\ 2 \\ 3 \\ \vdots \end{matrix} & \left(\begin{array}{cccc} \cdot & \cdot & \cdot & \cdot \\ \cdot & \ddots & [p(q_2=3, q_1=2 | u_1, \dots, u_T)] & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{array} \right) \end{matrix}$$

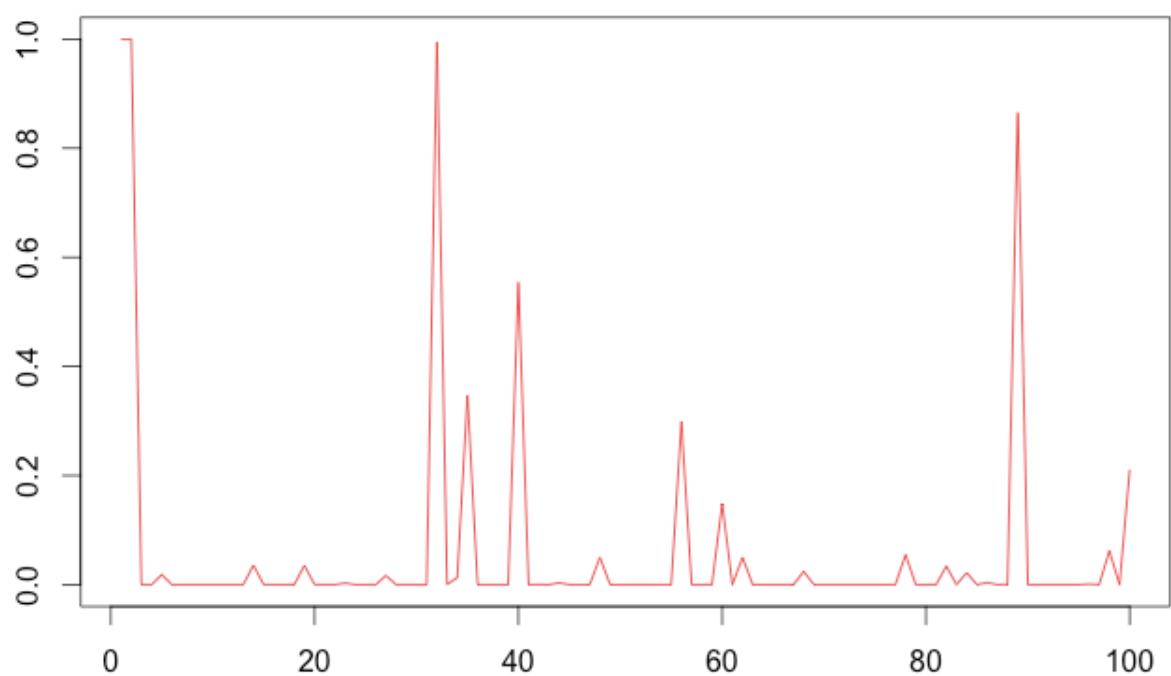
\leftarrow matrix $[E_t]$
of the list

Question 2 :

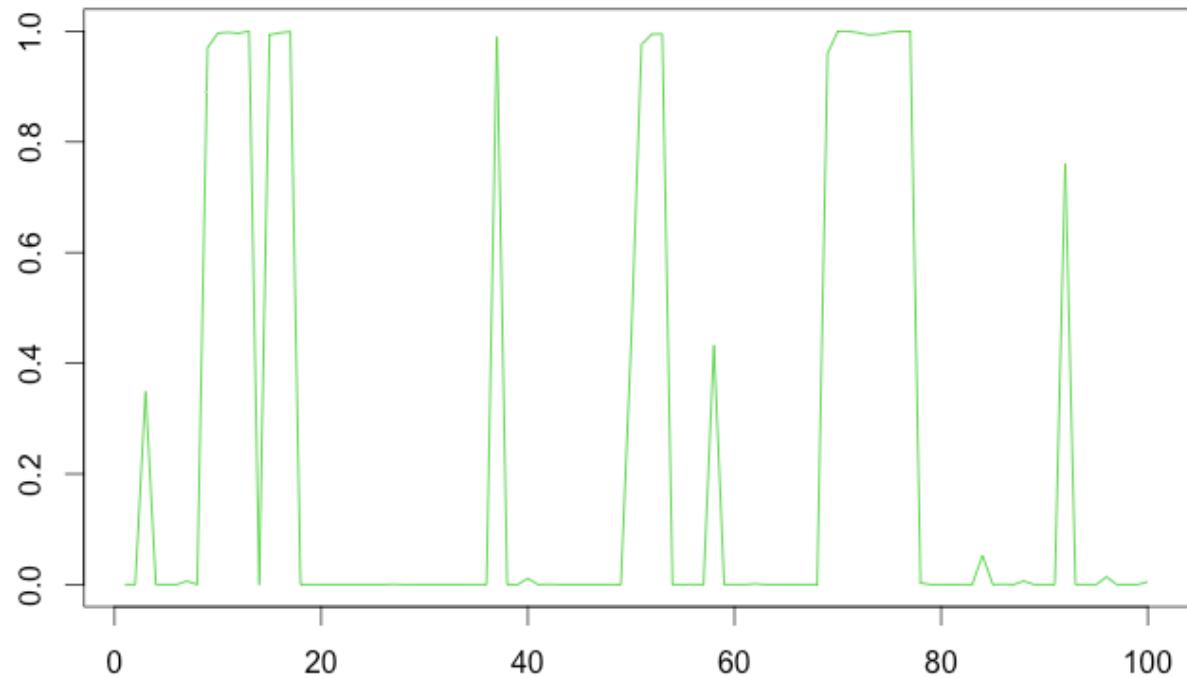
Probability to be in state 1 in function of the time t



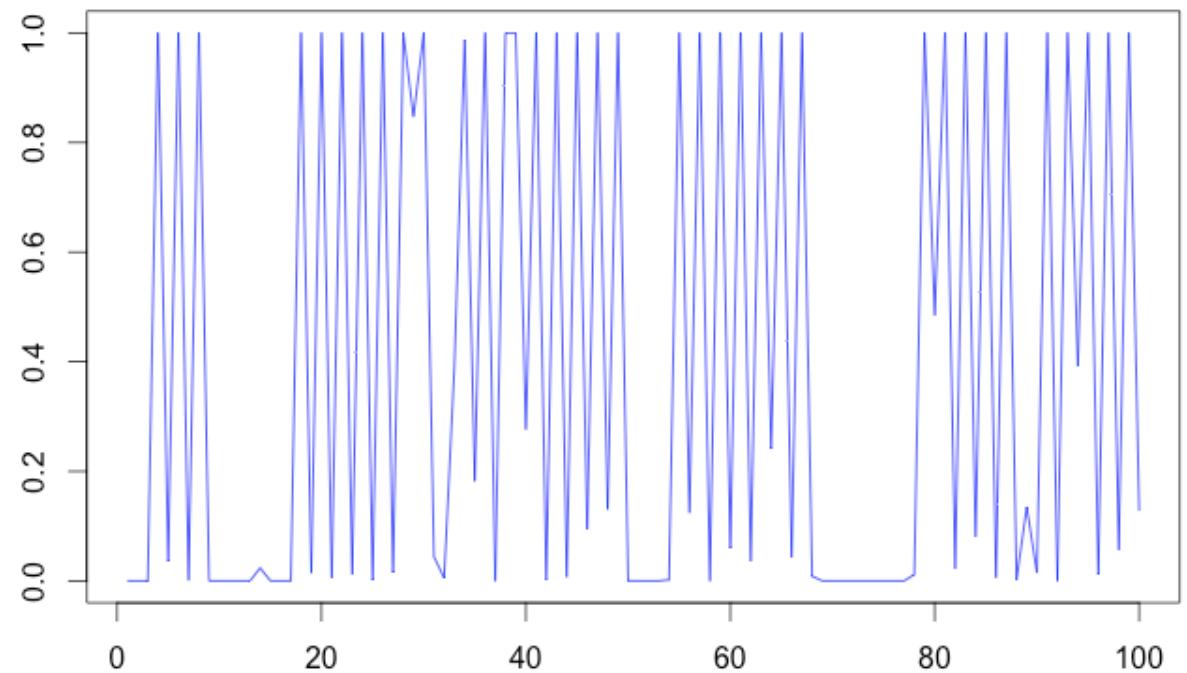
Probability to be in state 2 in function of the time t



Probability to be in state 3 in function of the time t



Probability to be in state 4 in function of the time t



③ We are going to implement the E.M algorithm.

(5)

To do this, we are going to derive the estimation equations of the E.M algorithm.

Our goal is to maximise the incomplete likelihood defined by

$$\log(p(u_1, \dots, u_T)) = \log \left(\sum_{q_t} \lambda_t(q_t) \beta_t(q_t) \right)$$

To do this we are going to compute the complete likelihood

(we denote $\Theta = (\pi, A, \Sigma, p)$).

$$\begin{aligned} \text{Then: } \ell_c(\Theta) &= \log \left(p(q_1) \prod_{t=1}^{T-1} p(q_{t+1}|q_t) \prod_{t=1}^T p(u_t|q_t) \right) \\ &= \sum_{i=1}^K \delta(q_1=i) \log(\pi_i) + \sum_{t=1}^{T-1} \sum_{i,j=1}^K \delta(q_{t+1}=i, q_t=j) \log(A_{ij}) \\ &\quad + \sum_{t=1}^T \sum_{i=1}^K \delta(q_t=i) \log(p(u_t|q_t)) \end{aligned}$$

$$\begin{aligned} \text{E-Step: } \mathbb{E}_{q_1, \dots, q_T} [\ell_c(\Theta)] &= \sum_{i=1}^K p(q_1=i | u_1, \dots, u_T) \log(\pi_i) \\ &\quad + \sum_{t=1}^{T-1} \sum_{i,j=1}^K p(q_{t+1}=i, q_t=j | u_1, \dots, u_T) \log(A_{ij}) \\ &\quad + \sum_{t=1}^T \sum_{i=1}^K p(q_t=i | u_1, \dots, u_T) \log(p(u_t | q_t=i)) \end{aligned}$$

M-Step: We are going to maximise according each parameter one by one since everything decouple.

We first care about π : the associated problem is given by:

$$\max_{\pi_i} \left(\sum_{i=1}^K p(q_1=i | u_1, \dots, u_T) \log(\pi_i) \right) \text{ for } i=1, \dots, K$$

$$\sum_{i=1}^K \pi_i = 1$$

We introduce the lagrangian:

$$\mathcal{L}(\pi, \lambda) = \sum_{i=1}^K p(q_1=i | u_1, \dots, u_T) \log(\pi_i) + \lambda \left(\sum_{i=1}^K \pi_i - 1 \right)$$

$$\text{Then } \frac{\partial \mathcal{L}}{\partial \pi_i} = \frac{p(q_1=i|u_1, \dots, u_T)}{\pi_i} + \lambda \quad (6)$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial \pi_i} = 0 \Rightarrow \pi_i = -\frac{p(q_1=i|u_1, \dots, u_T)}{\lambda}$$

$$\text{Since } \sum_{i=1}^K \pi_i = 1 \Rightarrow \lambda = -\sum_{i=1}^K p(q_1=i|u_1, \dots, u_T)$$

$$\text{Therefore } \pi_i^* = \frac{p(q_1=i|u_1, \dots, u_T)}{\sum_{i=1}^K p(q_1=i|u_1, \dots, u_T)} = \frac{p(q_1=i|u_1, \dots, u_T)}{\sum_{i=1}^K \lambda_i} = \frac{\lambda_i p(q_1=i)}{\sum_{i=1}^K \lambda_i p(q_1=i)}$$

We know care about A : the associated problem is given by

$$\max_{A_{ij}} \sum_{t=1}^{T-1} \sum_{i,j=1}^K p(q_{t+1}=i, q_t=j|u_1, \dots, u_T) \log(A_{ij})$$

The associated Lagrangian is given by :

$$\mathcal{L}(A, \lambda) = \sum_{t=1}^{T-1} \sum_{i,j=1}^K p(q_{t+1}=i, q_t=j|u_1, \dots, u_T) \log(A_{ij}) + \lambda (\sum_{i=1}^K A_{ij} - 1)$$

$$\text{Then } \frac{\partial \mathcal{L}}{\partial A_{ij}} = \sum_{t=1}^{T-1} \frac{p(q_{t+1}=i, q_t=j|u_1, \dots, u_T)}{A_{ij}} + \lambda$$

$$\Rightarrow \frac{\partial \mathcal{L}}{\partial A_{ij}} = 0 \Rightarrow A_{ij}^* = -\frac{\sum_{t=1}^{T-1} p(q_{t+1}=i, q_t=j|u_1, \dots, u_T)}{\lambda}$$

$$\text{Since } \sum_{i=1}^K A_{ij} = 1 \Rightarrow \lambda = -\sum_{i=1}^K \sum_{t=1}^{T-1} p(q_{t+1}=i, q_t=j|u_1, \dots, u_T)$$

$$= -\sum_{t=1}^{T-1} p(q_t=j|u_1, \dots, u_T)$$

$$\Rightarrow A_{ij}^* = \frac{\sum_{t=1}^{T-1} p(q_{t+1}=i, q_t=j|u_1, \dots, u_T)}{\sum_{t=1}^{T-1} p(q_t=j|u_1, \dots, u_T)}$$

We now care about (Σ, μ) .

⑦

We recall that: $u_t | q_{t=i} \sim N(\mu_i, \Sigma_i)$

We have:

$$\begin{aligned} & \sum_{t=1}^T \sum_{i=1}^K p(q_{t=i} | u_1, \dots, u_T) \log(p(u_t | q_{t=i})) \\ = & \sum_{t=1}^T \sum_{i=1}^K p(q_{t=i} | u_1, \dots, u_T) \left\{ \log\left(\frac{1}{(2\pi)^{K/2}}\right) + \log\left(\frac{1}{\sqrt{\Sigma_i}}\right) \right. \\ & \quad \left. - \frac{1}{2} (u_t - \mu_i)^T (\Sigma_i)^{-1} (u_t - \mu_i) \right\} \end{aligned}$$

We have:

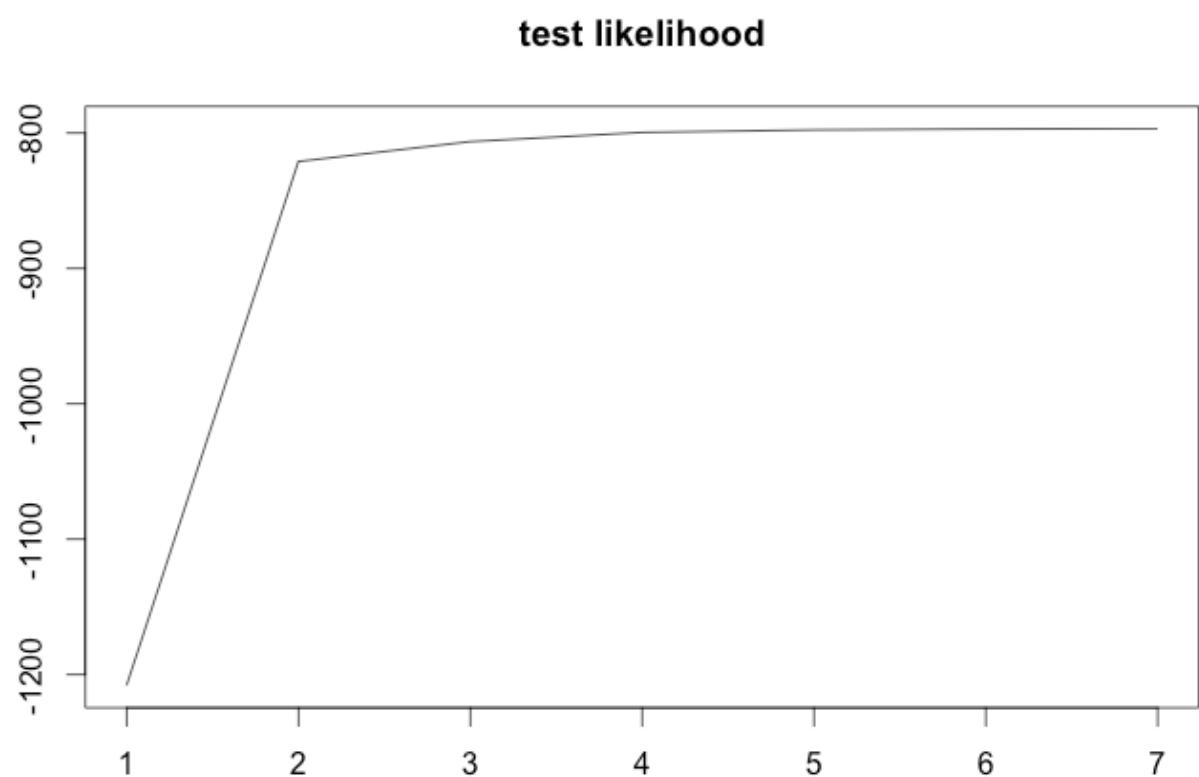
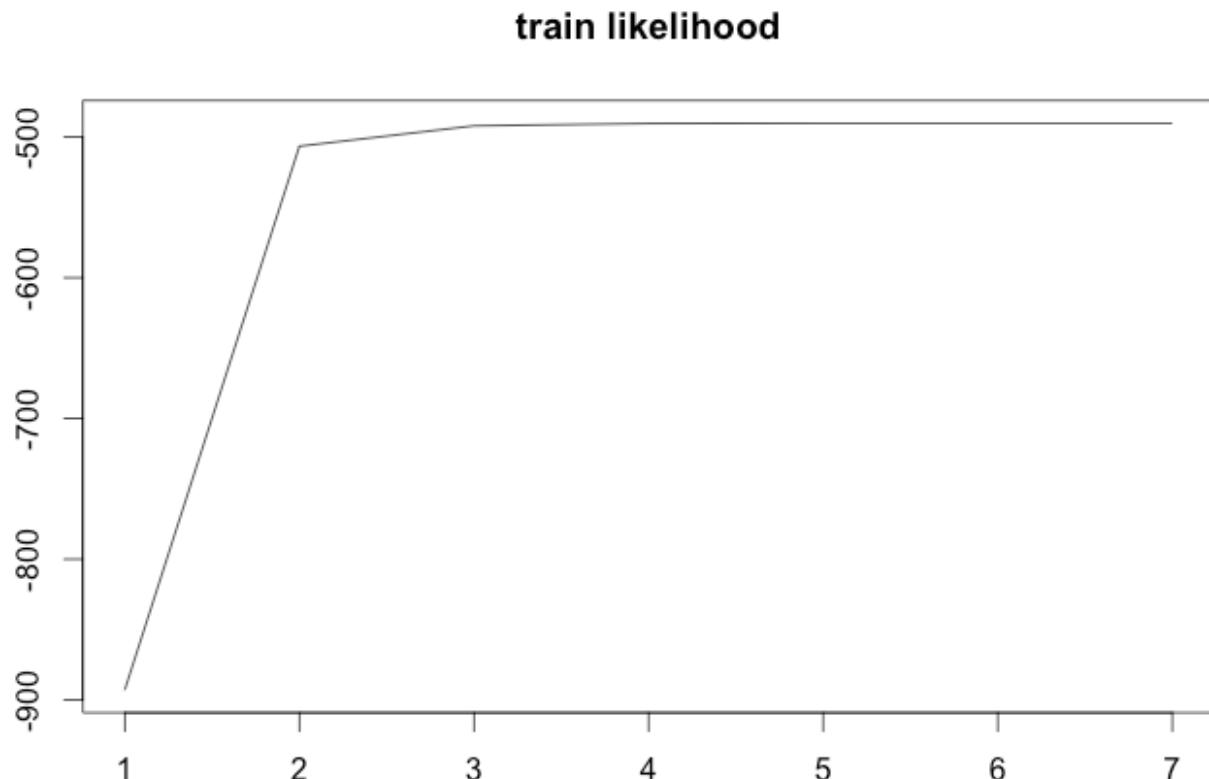
$$\begin{aligned} \frac{\partial}{\partial \mu_i} &= \sum_{t=1}^T p(q_{t=i} | u_1, \dots, u_T) C(\Sigma_i)^{-1} u_t - (\Sigma_i)^{-1} \mu_i \\ \frac{\partial}{\partial \mu_i} &= 0 \Rightarrow \sum_{t=1}^T p(q_{t=i} | u_1, \dots, u_T) u_t = \sum_{t=1}^T p(q_{t=i} | u_1, \dots, u_T) \mu_i \\ \Rightarrow \mu_i^* &= \frac{\sum_{t=1}^T p(q_{t=i} | u_1, \dots, u_T) u_t}{\sum_{t=1}^T p(q_{t=i} | u_1, \dots, u_T)} \end{aligned}$$

For the maximization over Σ , we use the result that we find in
homework 2 adapted here:

$$\Sigma_j^* = \frac{\sum_{t=1}^T p(q_{t=j} | u_1, \dots, u_T) (u_t - \mu_j^*) (u_t - \mu_j^*)^T}{\sum_{t=1}^T p(q_{t=j} | u_1, \dots, u_T)}$$

Question 5:

The likelihood of the train data is higher than the one of the test data because we learned the parameters on the train data.



Question 6:

Comparison of the Log-Likelihood

	Gaussian Mixture	HMM
Train	-921.0602	-490.1315
Test	-1214.876	-796.6534

- We recall that we have scaled the data train & the data test:
- We can replace Train<-scale(Data) by Train<-as.matrix(Data) to modify it.

In the HMM: We consider the the conditional probability of being in state i at position t given the entire observe sequence: $p(q_t = i | u_1, \dots, u_T)$.

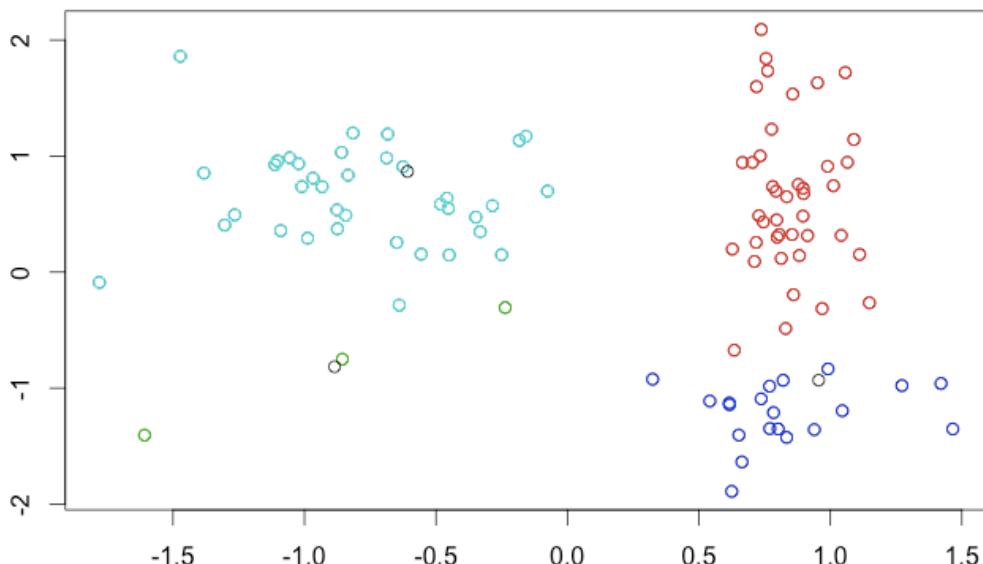
In the Gaussian Mixture Model: We consider the posterior probability of a state given the observation : $p(q_t = i | u_t)$.

The HMM depends on the entire sequence because of the underlying Markov Process whereas the Gaussian Mixture Model can be seen as a special Hidden Markov Model with the underlying state process that are IID.

We will therefore get a higher likelihood in the HMM than in the Gaussian Mixture Model.

Question 8:

The data in 2D with the cluster black centers



7) Viterbi algorithm.

The goal of Viterbi algorithm is to compute

$$\arg \max_{(q_t)_{1 \leq t \leq T}} p((q_t)_{1 \leq t \leq T} \mid (u_t)_{1 \leq t \leq T}) \quad (*),$$

given HMM parameters $p(q_t)$, $p(q_{t+1} \mid q_t)$, $p(u_t \mid q_t)$,
and $(u_t)_{1 \leq t \leq T}$ is the sequence of observation.

$$(*) = \arg \max_{(q_t)_{1 \leq t \leq T}} p\left((q_t)_{1 \leq t \leq T}, (u_t)_{1 \leq t \leq T}\right)$$

$$= \arg \max_{q_T} \max_{(q_t)_{1 \leq t \leq T-1}} p(q_T = h, (q_t)_{1 \leq t \leq T-1}, (u_t)_{1 \leq t \leq T})$$

V_T^h to compute recursively

V_T^h is the probability of the most likely sequence of states ending at state $q_T = h$.

By induction we have:

$$V_T^h = \max_{q_1, \dots, q_{T-1}} p(q_T = h, q_1, \dots, q_{T-1}, u_1, \dots, u_T)$$

$$= p(u_T \mid q_T = h) \max_i p(q_T = h \mid q_{T-1} = i) V_{T-1}^i$$

Hence we can compute the V_T^h & $V_{1 \leq t \leq 4}$,
the probabilities of the most likely sequence of $(t-1)$ states
ending at state $q_{T-t} = h$.

Finally we just have to take the higher probability at each time and the corresponding state:

$$q_T^* = \arg \max_h V_T^h$$

$$q_{t-1}^* = \arg \max_h \{ p(q_T^* | q_{t-1}=i) V_{t-1}^i \}$$

And recursively we find the optimal sequence
 (q_1^*, \dots, q_T^*) .

Algorithm:

- Initialize: $V_1^h = p(u_1 | q_1=h) p(q_1=h) \quad \forall h$

- Iterate : for $t=2, \dots, T$: $\forall h$:

$$V_t^h = p(u_t | q_t=h) \max_i \{ p(q_t=h | q_{t-1}=i) V_{t-1}^i \}$$

- Termination: $\max_{(q_t)_{1 \leq t \leq T}} p((q_t)_{1 \leq t \leq T}, (u_t)_{1 \leq t \leq T}) = \max_h V_T^h$

- Traceback: $q_T^* = \arg \max_h V_T^h$

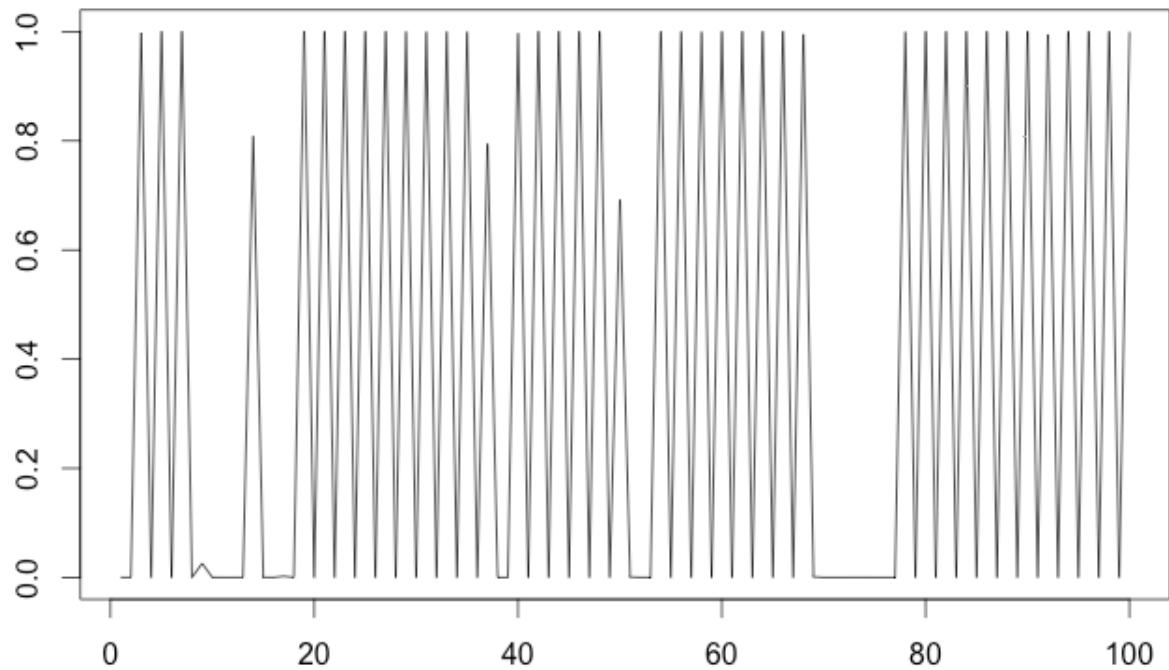
$$q_{t-1}^* = \arg \max_h \{ p(q_T^* | q_{t-1}=h) V_{t-1}^h \}$$

Return (q_1^*, \dots, q_T^*)

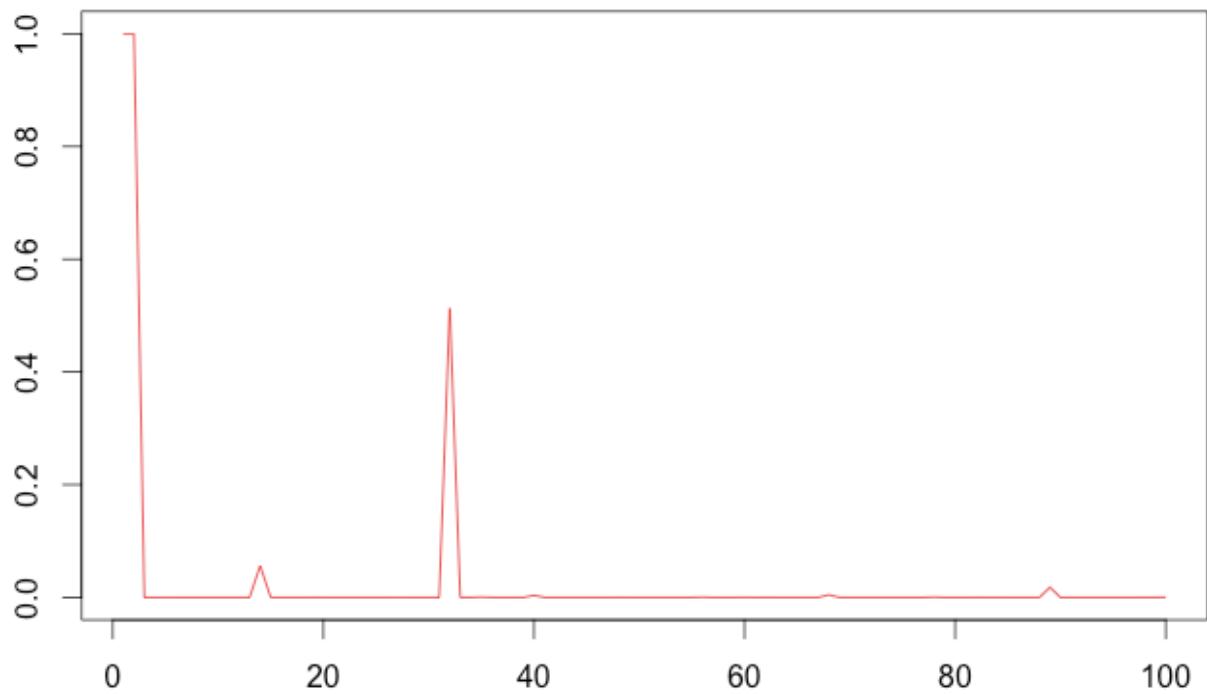
Question 9:

Plot with the parameters learned with the training set.

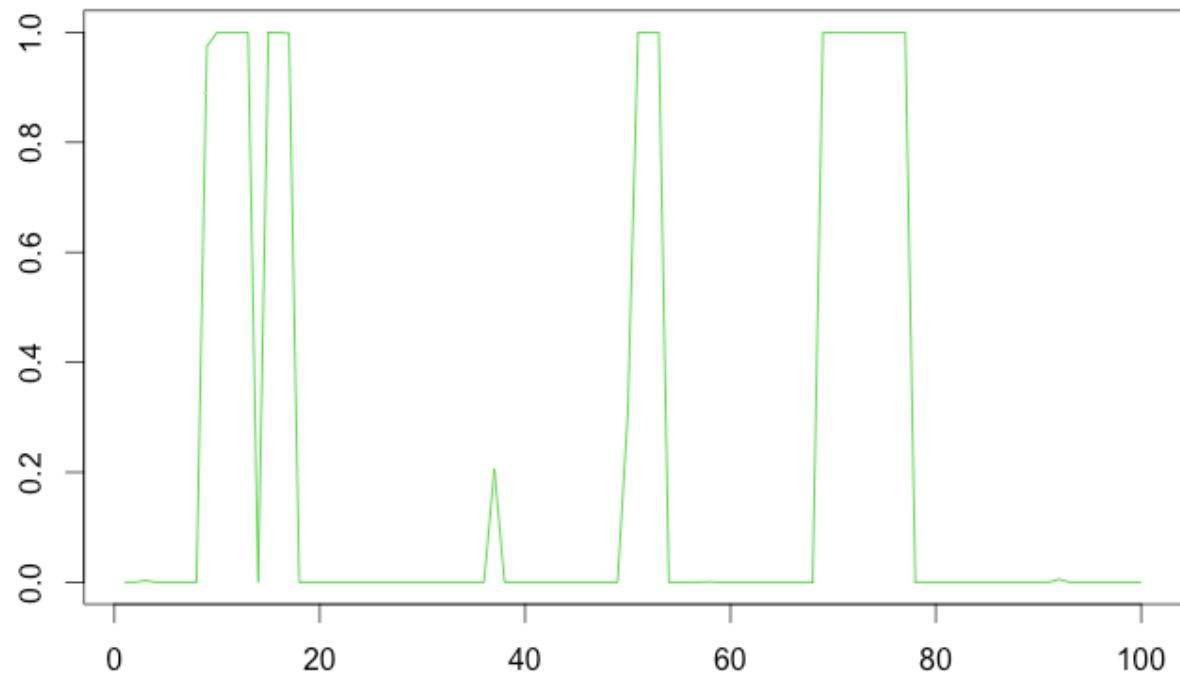
Probability to be in state 1 in function of the time t



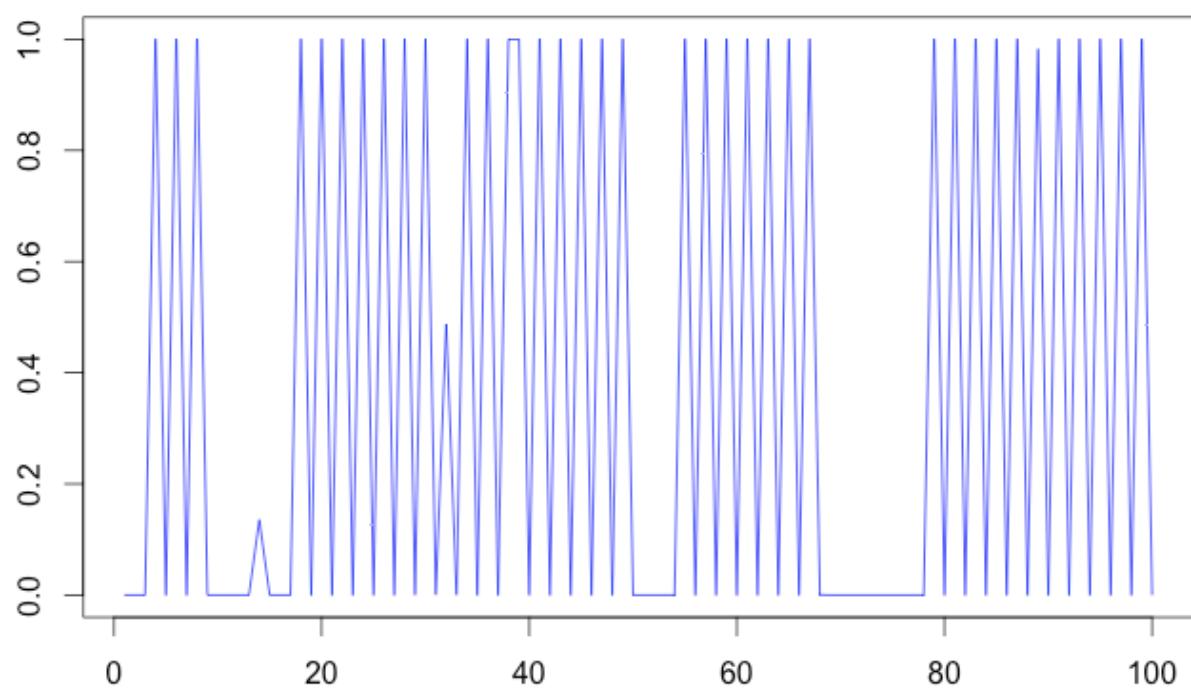
Probability to be in state 2 in function of the time t



Probability to be in state 3 in function of the time t

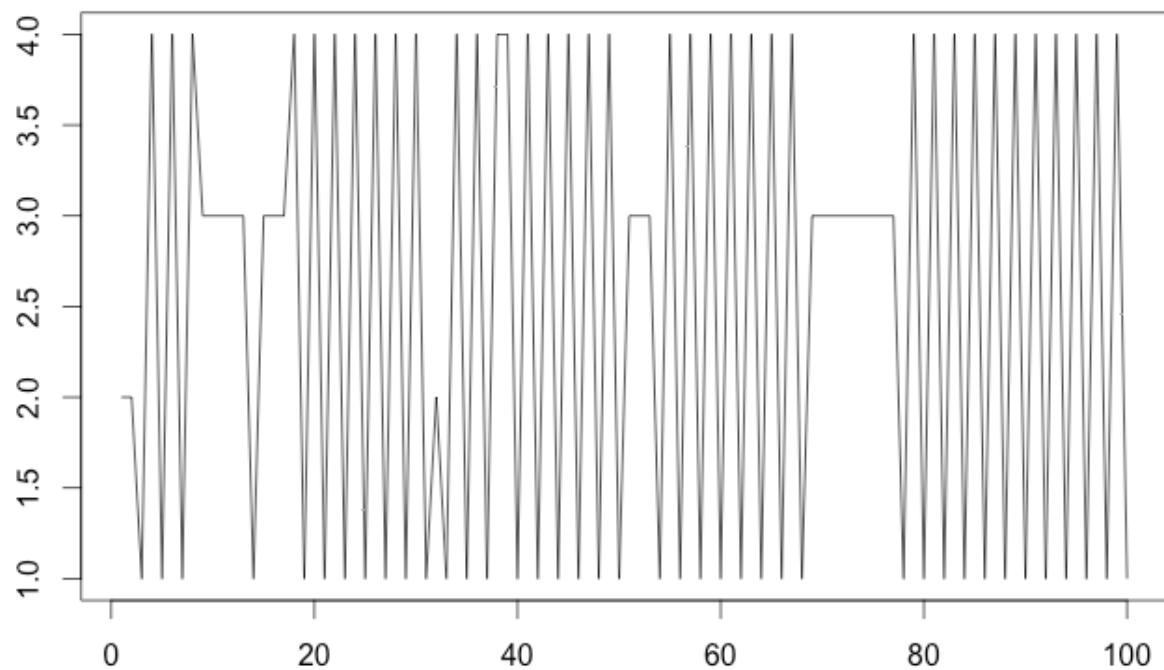


Probability to be in state 4 in function of the time t



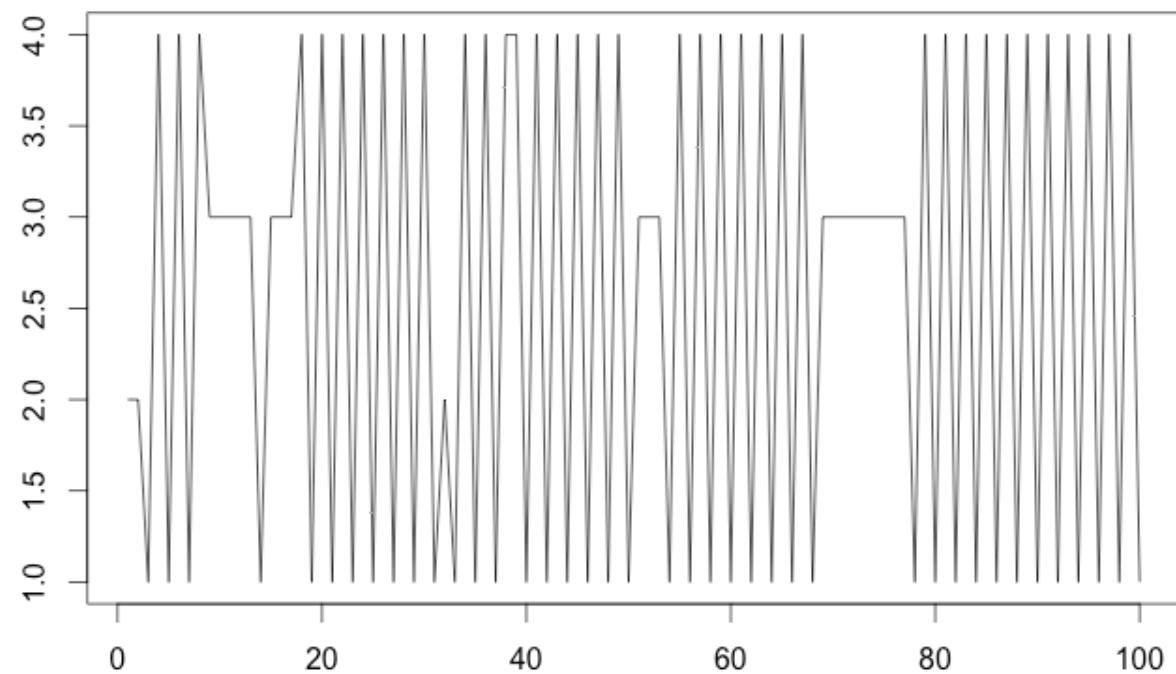
10)

most likely state according to marginal proba



Question 11:

most likely state according to Viterbi



As we can observe, we find the same states according to the Viterbi Algorithm and the marginal probability which seems logical since the most probable states sequence is near to the sequence of more probable states.

Question 12:

If the number of states is not known, we can't not only refer on the likelihood anymore.

If we increase the number of states:

- We will get a higher likelihood.
- We will increase the number of parameters such as transition probabilities and initial state probabilities which lead to slow down our program.
- We even can reach a distortion of 0 when we take the numbers of states equals to the size of the data. In that case, each data point is the center of its own center which doesn't have interest.

We have to make a compromise between the numbers of states and the number of parameters that we will have to compute.

We have to use some criteria as the Bayesian Information Criterion or the Akaike Information Criterion which take into account the likelihood, the number of parameters and also the size of the sample for the (BIC).