Clustering under the slowly mixing Gaussian Hidden Markov Model

Ibrahim Kaddouri
Joint work with Mohamed Ndaoud



Hidden Markov Model

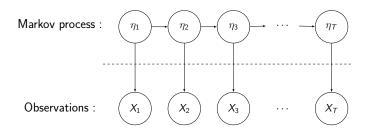


Figure: A Hidden Markov Model.

Latent (unobserved) variables $(\eta_k)_k$ form a Markov chain.

Hidden Markov Model

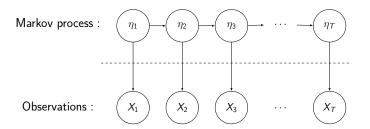


Figure: A Hidden Markov Model.

Latent (unobserved) variables $(\eta_k)_k$ form a Markov chain. Observations $(X_k)_k$ are independent conditionnally to $(\eta_k)_k$.

One observes $X_1,\ldots,X_n\in\mathbb{R}^d$ such that

$$X_i = \theta \eta_i + \xi_i$$

One observes $X_1, \ldots, X_n \in \mathbb{R}^d$ such that

$$X_i = \theta \eta_i + \xi_i$$

where

- ξ_1, \ldots, ξ_n are standard Gaussian random vectors

One observes $X_1, \ldots, X_n \in \mathbb{R}^d$ such that

$$X_i = \theta \eta_i + \xi_i$$

where

- ξ_1, \ldots, ξ_n are standard Gaussian random vectors
- $(\eta_i)_{1 \le i \le n}$ is a Markov chain with two states -1 and 1

One observes $X_1, \ldots, X_n \in \mathbb{R}^d$ such that

$$X_i = \theta \eta_i + \xi_i$$

where

- ξ_1, \ldots, ξ_n are standard Gaussian random vectors
- $(\eta_i)_{1 \le i \le n}$ is a Markov chain with two states -1 and 1
- The transition matrix ensures $Q = egin{pmatrix} 1 \delta & \delta \\ \delta & 1 \delta \end{pmatrix}$

Here, δ is allowed to depend on n and to be very small.

One observes $X_1, \ldots, X_n \in \mathbb{R}^d$ such that

$$X_i = \theta \eta_i + \xi_i$$

where

- ξ_1, \ldots, ξ_n are standard Gaussian random vectors
- $(\eta_i)_{1 \le i \le n}$ is a Markov chain with two states -1 and 1
- The transition matrix ensures $Q = egin{pmatrix} 1 \delta & \delta \\ \delta & 1 \delta \end{pmatrix}$

Here, δ is allowed to depend on n and to be very small.

Some natural questions arise in this setting:

One observes $X_1, \ldots, X_n \in \mathbb{R}^d$ such that

$$X_i = \theta \eta_i + \xi_i$$

where

- ξ_1, \ldots, ξ_n are standard Gaussian random vectors
- $(\eta_i)_{1 \le i \le n}$ is a Markov chain with two states -1 and 1
- The transition matrix ensures $Q = egin{pmatrix} 1 \delta & \delta \\ \delta & 1 \delta \end{pmatrix}$

Here, δ is allowed to depend on n and to be very small.

Some natural questions arise in this setting:

- Estimation of θ

One observes $X_1, \ldots, X_n \in \mathbb{R}^d$ such that

$$X_i = \theta \eta_i + \xi_i$$

where

- ξ_1, \ldots, ξ_n are standard Gaussian random vectors
- $(\eta_i)_{1 \le i \le n}$ is a Markov chain with two states -1 and 1
- The transition matrix ensures $Q = egin{pmatrix} 1 \delta & \delta \\ \delta & 1 \delta \end{pmatrix}$

Here, δ is allowed to depend on n and to be very small.

Some natural questions arise in this setting:

- Estimation of θ
- Estimation of δ

One observes $X_1, \ldots, X_n \in \mathbb{R}^d$ such that

$$X_i = \theta \eta_i + \xi_i$$

where

- ξ_1, \ldots, ξ_n are standard Gaussian random vectors
- $(\eta_i)_{1 \le i \le n}$ is a Markov chain with two states -1 and 1
- The transition matrix ensures $Q = \begin{pmatrix} 1-\delta & \delta \\ \delta & 1-\delta \end{pmatrix}$

Here, δ is allowed to depend on n and to be very small.

Some natural questions arise in this setting:

- Estimation of θ
- Estimation of δ
- Clustering

The problem of estimation

The problem of estimating θ has already been studied recently in Karagulyan and Ndaoud 2024. Let

$$M(n,d,\delta,t) = \inf_{\hat{\theta}(X_{1:n})} \sup_{\|\theta\|=t} \mathbb{E}_{\theta} \left[\min \left\{ \|\hat{\theta}(X_{1:n}) - \theta\|, \|\hat{\theta}(X_{1:n}) + \theta\| \right\} \right]$$

The problem of estimation

The problem of estimating θ has already been studied recently in Karagulyan and Ndaoud 2024. Let

$$M(n,d,\delta,t) = \inf_{\hat{\theta}(X_{1:n})} \sup_{\|\theta\|=t} \mathbb{E}_{\theta} \left[\min \left\{ \|\hat{\theta}(X_{1:n}) - \theta\|, \|\hat{\theta}(X_{1:n}) + \theta\| \right\} \right]$$

It was shown in Karagulyan and Ndaoud 2024 that when $d \leq \delta n$,

$$M(n,d,\delta,t)symp \left\{egin{array}{ll} t & ext{if } t \leq \left(rac{\delta d}{n}
ight)^{1/4} \ rac{1}{t}\sqrt{rac{\delta d}{n}} & ext{if } \left(rac{\delta d}{n}
ight)^{1/4} \leq t \leq \sqrt{\delta} \ \sqrt{rac{d}{n}} & ext{if } \sqrt{\delta} \leq t \end{array}
ight.$$

Let $\eta=(\eta_i)_{1\leq i\leq n}$ and $\ell(\hat{\eta},\eta)=\min_{\nu\in\{-1,1\}}|\hat{\eta}+\nu\eta|$ where |.| is the Hamming distance. We define the risk of a clustering procedure $\hat{\eta}$ as

$$\mathcal{R}(\theta, \delta, \hat{\eta}) = \mathbb{E}\left[\ell\left(\hat{\eta}(X_{1:n}), \eta\right)\right]$$

and the Bayes risk of clustering as

$$\inf_{\hat{\eta}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right)$$

Let $\eta=(\eta_i)_{1\leq i\leq n}$ and $\ell(\hat{\eta},\eta)=\min_{\nu\in\{-1,1\}}|\hat{\eta}+\nu\eta|$ where |.| is the Hamming distance. We define the risk of a clustering procedure $\hat{\eta}$ as

$$\mathcal{R}(\theta, \delta, \hat{\eta}) = \mathbb{E}\left[\ell\left(\hat{\eta}(X_{1:n}), \eta\right)\right]$$

and the Bayes risk of clustering as

$$\inf_{\hat{\eta}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right)$$

This risk was studied in Gassiat, Kaddouri and Naulet 2023 in the strongly mixing setting.

Let $\eta=(\eta_i)_{1\leq i\leq n}$ and $\ell(\hat{\eta},\eta)=\min_{\nu\in\{-1,1\}}|\hat{\eta}+\nu\eta|$ where |.| is the Hamming distance. We define the risk of a clustering procedure $\hat{\eta}$ as

$$\mathcal{R}\left(\theta,\delta,\hat{\eta}\right) = \mathbb{E}\left[\ell\left(\hat{\eta}(X_{1:n}),\eta\right)\right]$$

and the Bayes risk of clustering as

$$\inf_{\hat{\eta}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right)$$

This risk was studied in Gassiat, Kaddouri and Naulet 2023 in the strongly mixing setting. It was shown that for $\tilde{\alpha}_n = \frac{C}{\delta^5 \sqrt{n}}$

$$\frac{\delta^2(1-\tilde{\alpha}_n)}{2(1-\delta)}e^{-2\|\theta\|^2} \leq \inf_{\hat{\eta}} \mathcal{R}\left(\theta,\delta,\hat{\eta}\right) \leq (1-\delta)e^{-\frac{\|\theta\|^2}{2}}$$

Let $\eta = (\eta_i)_{1 \leq i \leq n}$ and $\ell(\hat{\eta}, \eta) = \min_{\nu \in \{-1,1\}} |\hat{\eta} + \nu \eta|$ where $|\cdot|$ is the Hamming distance. We define the risk of a clustering procedure $\hat{\eta}$ as

$$\mathcal{R}\left(\theta,\delta,\hat{\eta}\right) = \mathbb{E}\left[\ell\left(\hat{\eta}(X_{1:n}),\eta\right)\right]$$

and the Bayes risk of clustering as

$$\inf_{\hat{\eta}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right)$$

This risk was studied in Gassiat, Kaddouri and Naulet 2023 in the strongly mixing setting. It was shown that for $\tilde{\alpha}_n = \frac{C}{\delta^5 \sqrt{n}}$

$$\frac{\delta^2(1-\tilde{\alpha}_n)}{2(1-\delta)}e^{-2\|\theta\|^2} \leq \inf_{\hat{\eta}} \mathcal{R}\left(\theta,\delta,\hat{\eta}\right) \leq (1-\delta)e^{-\frac{\|\theta\|^2}{2}}$$

When δ is allowed to be very small, this fails to provide a precise understanding of the interplay between the strength of the signal and the Markovian dependence measured by δ .

Two frameworks

We study the risk of clustering observations in two frameworks:

• Offline setting: All observations are used in the clustering procedures. Clustering rules are of the form:

$$\hat{\eta}(Y_{1:n}) = (\hat{\eta}_i(Y_{1:n}))_{1 \le i \le n}.$$

The Bayes risk of online clustering will be denoted $\inf_{\hat{\eta}^{O_n}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right)$.

Two frameworks

We study the risk of clustering observations in two frameworks:

• Offline setting: All observations are used in the clustering procedures. Clustering rules are of the form:

$$\hat{\eta}(Y_{1:n}) = (\hat{\eta}_i(Y_{1:n}))_{1 \leq i \leq n}.$$

The Bayes risk of online clustering will be denoted $\inf_{\hat{\eta}^{O_n}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right)$.

• Online setting: Observations are clustered sequentially, with access to the past observations only at each step. Clustering rules are of the form: $\hat{\eta}(Y_{1:n}) = (\hat{\eta}_i(Y_{1:i}))_{1 < i < n}$.

The Bayes risk of offline clustering will be denoted $\inf_{\hat{\eta}} \mathcal{R}(\theta, \delta, \hat{\eta})$.

Theorem

When $\delta \geq \frac{1}{n}$,

$$\inf_{\hat{\eta}^{\mathit{On}}}\mathcal{R}\left(heta,\delta,\hat{\eta}
ight)\leq\left\{
ight.$$

Theorem

When
$$\delta \geq \frac{1}{n}$$
,

$$\inf_{\hat{\eta}^{\mathcal{O}n}}\mathcal{R}\left(heta,\delta,\hat{\eta}
ight)\leq\left\{egin{array}{c} rac{1}{2} \end{array}
ight.$$

if
$$\|\theta\|^2 \le 2\delta$$

Theorem

When $\delta \geq \frac{1}{n}$,

$$\inf_{\hat{\eta}^{On}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right) \leq \begin{cases} \frac{1}{2} & \text{if } \|\theta\|^2 \leq 2\delta \\ \frac{4\delta}{\|\theta\|^2} \left(\log\left(\frac{\|\theta\|^2}{2\delta}\right) + 1\right) & \text{if } 2\delta < \|\theta\|^2 \leq \log\left(\frac{1}{\delta}\right) \end{cases}$$

Theorem

When $\delta \geq \frac{1}{n}$,

$$\inf_{\hat{\eta}^{On}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right) \leq \begin{cases} \frac{1}{2} & \text{if } \|\theta\|^2 \leq 2\delta \\ \frac{4\delta}{\|\theta\|^2} \left(\log\left(\frac{\|\theta\|^2}{2\delta}\right) + 1\right) & \text{if } 2\delta < \|\theta\|^2 \leq \log\left(\frac{1}{\delta}\right) \\ e^{-\frac{\|\theta\|^2}{2}} & \text{if } \|\theta\|^2 > \log\left(\frac{1}{\delta}\right) \end{cases}$$

Theorem

When $\delta \geq \frac{1}{n}$,

$$\inf_{\hat{\eta}^{On}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right) \leq \begin{cases} \frac{1}{2} & \text{if } \|\theta\|^2 \leq 2\delta \\ \frac{4\delta}{\|\theta\|^2} \left(\log\left(\frac{\|\theta\|^2}{2\delta}\right) + 1\right) & \text{if } 2\delta < \|\theta\|^2 \leq \log\left(\frac{1}{\delta}\right) \\ e^{-\frac{\|\theta\|^2}{2}} & \text{if } \|\theta\|^2 > \log\left(\frac{1}{\delta}\right) \end{cases}$$

This upper-bound is reached by this online clustering procedure:

- For $i \in [1, k]$:

$$\hat{\eta}_i(X_{1:n}) = \operatorname{sign}(\langle X_i, \theta \rangle)$$

- For i > k:

$$\hat{\eta}_i(X_{1:n}) = \operatorname{sign}\left(\left\langle \frac{1}{k} \sum_{j=i-k+1}^i X_j, \theta \right
angle \right)$$

with
$$k = \left\lceil \frac{2}{\|\theta\|^2} \log \left(\frac{\|\theta\|^2}{2\delta} \right) \right\rceil$$
.

Proposition

$$\inf_{\hat{\eta}^{\mathit{On}}}\mathcal{R}\left(heta,\delta,\hat{\eta}
ight)\gtrsim\left\{
ight.$$

When
$$\delta \gtrsim \frac{1}{n}$$
,

$$\inf_{\hat{\eta}^{On}}\mathcal{R}\left(heta,\delta,\hat{\eta}
ight)\gtrsim\left\{ 1
ight.$$

if
$$\|\theta\|^2 \le 2\delta$$

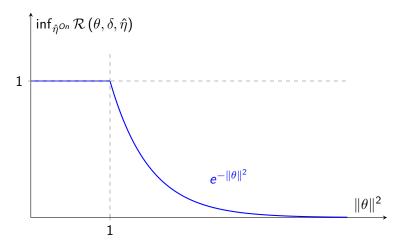
Proposition

$$\inf_{\hat{\eta}^{On}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right) \gtrsim \begin{cases} 1 & \text{if } \|\theta\|^2 \leq 2\delta \\ \frac{\delta}{\|\theta\|^2} \left(\log\left(\frac{\|\theta\|^2}{2\delta}\right) + 1\right) & \text{if } 2\delta < \|\theta\|^2 \leq \log\left(\frac{1}{\delta}\right) \end{cases}$$

When
$$\delta \gtrsim \frac{1}{n}$$
,

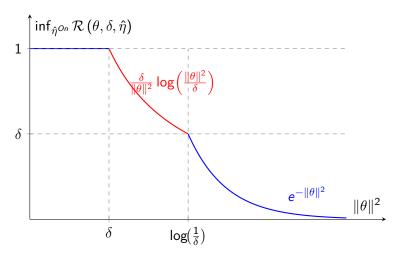
$$\inf_{\hat{\eta}^{On}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right) \gtrsim \begin{cases} 1 & \text{if } \|\theta\|^2 \leq 2\delta \\ \frac{\delta}{\|\theta\|^2} \left(\log\left(\frac{\|\theta\|^2}{2\delta}\right) + 1\right) & \text{if } 2\delta < \|\theta\|^2 \leq \log\left(\frac{1}{\delta}\right) \\ \delta e^{-2\|\theta\|^2} & \text{if } \|\theta\|^2 > \log\left(\frac{1}{\delta}\right) \end{cases}$$

Behavior of the Bayes risk of online clustering



Behavior of the Bayes risk of online clustering in the **strongly** mixing regime.

Behavior of the Bayes risk of online clustering



Behavior of the Bayes risk of online clustering in the **slowly** mixing regime.

Proposition $\inf_{\hat{\eta}} \mathcal{R}\left(\theta,\delta,\hat{\eta}\right) \lesssim \left\{$

$$\inf_{\hat{\eta}} \mathcal{R} \left(heta, \delta, \hat{\eta}
ight) \lesssim \left\{ egin{array}{l} rac{1}{2} \end{array}
ight.$$

if
$$\|\theta\|^2 \le 2\delta$$

$$\inf_{\hat{\eta}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right) \lesssim \begin{cases} \frac{1}{2} & \text{if } \|\theta\|^2 \leq 2\delta \\ \frac{2\delta}{\|\theta\|^2} \left(\log\left(\frac{\|\theta\|^2}{2\delta}\right) + 1\right) & \text{if } 2\delta \leq \|\theta\|^2 \leq 1 \end{cases}$$

$$\inf_{\hat{\eta}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right) \lesssim \begin{cases} \frac{1}{2} & \text{if } \|\theta\|^2 \leq 2\delta \\ \frac{2\delta}{\|\theta\|^2} \left(\log\left(\frac{\|\theta\|^2}{2\delta}\right) + 1\right) & \text{if } 2\delta \leq \|\theta\|^2 \leq 1 \\ \delta\log\left(\frac{1}{\delta}\right) e^{-\frac{\|\theta\|^2}{4}} & \text{if } \|\theta\|^2 > 1 \end{cases}$$

Offline clustering procedure

Consider $\tilde{\eta}_{a-1}$ and $\tilde{\eta}_{a+1}$ are the online estimators defined by:

$$\tilde{\eta}_{a-1}(X_{a-k:a-1}) = \langle \frac{1}{k} \sum_{j=a-k}^{a-1} X_j, \theta \rangle, \ \tilde{\eta}_{a+1}(X_{a+1:a+k}) = \langle \frac{1}{k} \sum_{j=a+1}^{a+k} X_j, \theta \rangle$$

where $k = \left\lceil \frac{2}{\|\theta\|^2} \log \left(\frac{\|\theta\|^2}{2\delta} \right) \right\rceil$.

Consider the clustering procedure $\hat{\eta}$ such that for $a \in [1, n]$:

$$\begin{split} \hat{\eta}_{a}(\tilde{\eta}_{a-1}, X_{a-k:a+k}, \tilde{\eta}_{a+1}) &= \\ \left\{ \begin{array}{l} \tilde{\eta}_{a-1} & \text{if } \tilde{\eta}_{a-1} = \tilde{\eta}_{a+1}, \|\theta\|^2 < \log\left(\frac{1}{\delta}\right) \text{ and } a \in [\![k+1, n-k]\!] \\ \text{sign}\left(\langle X_a, \theta \rangle\right) & \text{else.} \end{array} \right. \end{split}$$

Proposition

$$\inf_{\hat{\eta}}\mathcal{R}\left(heta,\delta,\hat{\eta}
ight)\gtrsim\left\{
ight.$$

Proposition

$$\inf_{\hat{\eta}} \mathcal{R} \left(heta, \delta, \hat{\eta}
ight) \gtrsim \left\{ 1
ight.$$

$$\textit{if} \ \|\theta\|^2 \leq 2\delta$$

Proposition

$$\inf_{\hat{\eta}} \mathcal{R} \left(heta, \delta, \hat{\eta}
ight) \gtrsim egin{cases} 1 & & ext{if } \| heta\|^2 \leq 2\delta \ & & ext{if } 2\delta < \| heta\|^2 \leq 1 \end{cases}$$

Proposition

$$\inf_{\hat{\eta}} \mathcal{R}\left(\theta, \delta, \hat{\eta}\right) \gtrsim \begin{cases} 1 & \text{if } \|\theta\|^2 \leq 2\delta \\ \frac{\delta}{\|\theta\|^2} & \text{if } 2\delta < \|\theta\|^2 \leq 1 \\ \delta e^{-2\|\theta\|^2} & \text{if } \|\theta\|^2 > 1 \end{cases}$$

Behavior of the Bayes risk of offline clustering

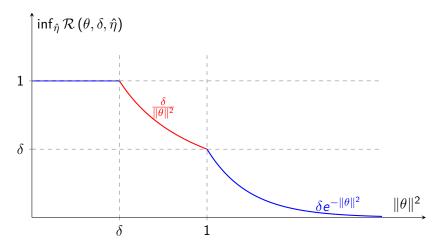


Figure: Behavior of the Bayes risk of offline clustering in the slowly mixing regime

WLOG, assume $n = k\ell$. For each bucket $i \in [1, \ell]$, consider the sample mean of k observations inside the i-th bucket

$$\tilde{X}_i = \frac{1}{k} \sum_{j=k(i-1)+1}^{ki} X_j$$

Note that

$$\tilde{X}_i = \theta \bar{\eta}_i + \frac{\xi_i}{\sqrt{k}}$$

where $\bar{\eta}_i = \frac{1}{k} \sum_{j=k(i-1)+1}^{ki} \eta_j$ and ξ_i is a standard Gaussian vector.

WLOG, assume $n = k\ell$. For each bucket $i \in [1, \ell]$, consider the sample mean of k observations inside the i-th bucket

$$\tilde{X}_i = \frac{1}{k} \sum_{j=k(i-1)+1}^{ki} X_j$$

Note that

$$\tilde{X}_i = \theta \bar{\eta}_i + \frac{\xi_i}{\sqrt{k}}$$

where $\bar{\eta}_i = \frac{1}{k} \sum_{j=k(i-1)+1}^{ki} \eta_j$ and ξ_i is a standard Gaussian vector.

We stack the ℓ terms in an $\mathbb{R}^{d \times \ell}$ matrix form as follows:

$$ilde{X} = heta ar{\eta}^ op + rac{\xi}{\sqrt{k}}$$

We then consider the Gram matrix of observations $\tilde{X}\tilde{X}^{\top}$.



It is easy to see that:

$$\mathbb{E}\left[\frac{1}{\ell}\tilde{X}\tilde{X}^{\top}\right] = \frac{\mathbb{E}\left[\|\bar{\eta}\|^2\right]}{\ell}\theta\theta^{\top} + \frac{1}{k}\mathbf{I}_d$$

Since our goal is adaptation up to multiplicative constants, it is natural to consider the following estimator of $\|\theta\|$:

$$\|\hat{ heta}(\ell)\| = \left\| \frac{1}{\ell} \tilde{X} \tilde{X}^{\top} - \frac{1}{k} \mathbf{I}_d \right\|_{op}^{1/2}$$

Proposition

There exists c>0 and C>0 such that for $\varepsilon>0$ and $\ell\geq n\|\theta\|^2\vee\frac{d}{\varepsilon^2}\vee\frac{2}{3\varepsilon}n\delta$,

$$\mathbb{P}\left(\left|\|\hat{\theta}(\ell)\|^2 - \|\theta\|^2\right| \ge C\varepsilon \frac{\ell}{n}\right) \le e^{-\frac{c\varepsilon^2\ell}{2}}.$$

Define now
$$\mathbf{I}_{\ell} = \left[\|\hat{\theta}(\ell)\|^2 - \frac{5}{6} \frac{\ell}{n}, \ \|\hat{\theta}(\ell)\|^2 + \frac{5}{6} \frac{\ell}{n} \right]$$
 and $\hat{\ell} = \min \left\{ \tilde{\ell} \mid \cap_{\ell \geq \tilde{\ell}} \mathbf{I}_{\ell} \neq \varnothing \right\}$. Let $\tilde{\theta} \in \cap_{\ell \geq \tilde{\ell}} \mathbf{I}_{\ell}$.

Theorem

There exist positive absolute constants c_1, c_2 and c_3 such that in the regime where $\|\theta\|^2 \ge c_1 \left(\frac{d}{n} \lor \delta\right)$,

$$\mathbb{P}\left(\left|\tilde{\theta}^2 - \|\theta\|^2\right| \ge \frac{\|\theta\|^2}{2}\right) \le c_2 e^{-c_3 n \|\theta\|^2}.$$

This shows that $\tilde{\theta}^2$ is an appropriate estimator of $\|\theta\|^2$ up to multiplicative constants.

Conclusion

- Clustering is easier under the slowly mixing regime.
- Unexpected behavior of the Bayes risk in some regimes.
- Many questions are still not answered: full adaptation, high dimension, estimation of δ .

References



Gassiat, Elisabeth, Ibrahim Kaddouri and Zacharie Naulet (2023). "Clustering and Classification Risks in Non-parametric Hidden Markov and I.I.D. Models". inarXiv: 2309.12238 [stat.ST].



Karagulyan, Vahe and Mohamed Ndaoud (2024). "Adaptive Mean Estimation in the Hidden Markov sub-Gaussian Mixture Model". inarXiv: 2406.12446 [stat.ST].