# Supplements to: Deterministic Approximate EM algorithm

Application to the Riemann approximation EM and the tempered EM

Thomas Lartigue · Stanley Durrleman · Stéphanie Allassonnière

the date of receipt and acceptance should be inserted later

### 1 Introduction

In these supplementary materials, we give in full details the proofs that were sketched in the main paper detail and a more extensive experimental study of the tempered EM algorithm (tmp-EM). Section 2 is dedicated to the proofs. Section 2.1 details the proof of convergence for the Deterministic Approximate EM algorithm, our central result, Theorem 1 of the paper. Section 2.2 details the proof of convergence for tmp-EM, Theorem 3 of the paper. In Section 3, we perform an in depth experimental study of the behaviour and performances of tmp-EM to demonstrate that it solves the issues raised about the EM. In particular, we illustrate on synthetic data, in the GMM case, that tmp-EM consistently reaches better values of the likelihood than the unmodified EM, in addition to better estimating the GMM parameters. We demonstrate that, as intended, tmp-EM is able to escape bad initialisations, unlike EM, and that more diverse configurations are explored during the procedure before reaching convergence. We confirm these observation on real data from the scikit learn library [6]. Finally, in Section 4, we test the tmp-EM within a more complex pipeline: the Independent Factor Analysis model [3] with a hidden GMM. We illustrate that, with tmp-EM, the identified sources are cleaner, more stable looking, and closer to the real ones when those are known.

S. Durrleman Aramis project-team, Inria E-mail: stanley.durrleman@fr

S. Allassonnière

T. Lartigue

Aramis project-team, Inria and CMAP, CNRS, École polytechnique, I.P. Paris E-mail: thomas.lartigue@inria.fr *Present address:* thomas.lartigue@dzne.de

Centre de Recherche des Cordeliers, Université de Paris, INSERM, Sorbonne Université E-mail: stephanie.allassonniere@parisdescartes.fr

## 2 Proofs of the two main Theorems

2.1 Proof of the general theorem

In this Section, we prove Theorem 1 of the main paper, for the convergence of the Deterministic Approximate EM algorithm.

We use two intermediary results of [5]: their "Proposition 9" and "Proposition 11", which we recall here:

**Proposition 1** ("**Proposition 9**") Let  $\Theta \subseteq \mathbb{R}^l$ , K compact  $\subset \Theta, \mathcal{L} \subseteq \Theta$  such that  $\mathcal{L} \cap K$  compact. Let us assume

- 
$$WC^0$$
 Lyapunov function with regards to  $(T, \mathcal{L})$ .

 $-\exists u_n \in K^{\mathbb{N}} \text{ such that } |W(u_{n+1}) - W \circ T(u_n)| \xrightarrow[n \infty]{} 0$ 

Then

- $\{W(u_n)\}_{n \in \mathbb{N}}$  converges towards a connected component of  $W(\mathcal{L} \cap K)$
- If  $W(\mathcal{L} \cap K)$  has an empty interior, then  $\{W(u_n)\}_n$  converges towards  $w^*$  and  $\{u_n\}_n$  converges towards the set  $\mathcal{L}_{w^*} \cap K$

$$\mathcal{L}_{w^*} = \left\{ \theta \in \mathcal{L} | W(\theta) = w^* \right\}$$

**Proposition 2** ("**Proposition 11**") Let  $\Theta \subseteq \mathbb{R}^l$ , T and  $\{F_n\}_n$  point to point maps on  $\Theta$ . Let  $\{\theta_n\}_n$  be the sequence defined by the stable approximate EM with likelihood f and approximate maps sequence  $\{F_n\}_n$ . Let  $\mathcal{L} \subset \Theta$ . We assume

- the A1 2 conditions of Proposition 10 of [5].
  - (A1) There exists W, a  $C^0$  Lyapunov function with regards to  $(T, \mathcal{L})$  such that  $\forall M > 0, \{\theta \in \Theta, W(\theta) > M\}$  is compact, and:

$$\Theta = \bigcup_{n \in \mathbb{N}} \left\{ \theta \in \Theta | W(\theta) > n^{-1} \right\} \,.$$

- (A2)  $W(\mathcal{L})$  is compact OR (A2')  $W(\mathcal{L} \cap K)$  is finite for all compact  $K \subseteq \Theta$ .  $\begin{array}{l} (\Pi D) & (D) & \text{is compare OR} (\Pi D) & (D + \Pi) & \text{is finite for a} \\ - & \forall u \in K_0, \quad \lim_{n \infty} |W \circ F_n - W \circ T|(u) = 0 \\ - & \forall \text{ compact } K \subseteq \Theta, \quad \lim_{n \infty} |W \circ F_n(u_n) - W \circ T(u_n)| \mathbb{1}_{u_n \in K} = 0 \end{array}$ 

Then

With probability 1,  $\limsup p_n < \infty$  and  $\{u_n\}_n$  compact sequence

Remark 1 In [5], condition (A1) is mistakenly written as:

$$\Theta = \bigcup_{n \in \mathbb{N}} \left\{ \theta \in \Theta | W(\theta) > n \right\} .$$

This is a typo that we have corrected here.

We need to prove that, under the conditions of Theorem 1, we verify the conditions of Proposition Proposition 1 and Proposition 2. Then we will have the results announced in Theorem 1.

#### 2.1.1 Verifying the conditions of 2

f is the likelihood function of a model of the curved exponential family.  $\mathcal{L}$  the set of its critical points:  $\mathcal{L} := \{\theta \in \Theta | \nabla f(\theta) = 0\}$ . Let T be the point to point map describing the transition between  $\theta_n$  and  $\theta_{n+1}$  in the exact EM algorithm. The general properties of the EM tell us that its stationary points are the critical points of  $f: \mathcal{L} = \{\theta \in \Theta | T(\theta) = \theta\}$ . Additionally, f is a  $C^0$  Lyapunov function associated to  $(T, \mathcal{L})$ . Let  $\{\theta_n\}_n$  be the sequence defined by the stable approximate EM with  $\{F_n\}_{n \in \mathbb{N}}$  our sequence of point to point maps.

We verify that under this framework - and with the assumptions of Theorem 1 - we check the conditions of Proposition 2.

### As in [5], M1 - 3 implies A1 - 2.

Let us show that we have the last two conditions for Proposition 2:

$$\forall \theta \in K_0, \quad \lim_{n \to \infty} |f \circ F_n - f \circ T|(\theta) = 0, \qquad (1)$$

and

$$\forall \text{ compact } K \subseteq \Theta, \quad \lim_{n \to \infty} |f \circ F_n(\theta_n) - f \circ T(\theta_n)| \mathbb{1}_{\theta_n \in K} = 0.$$
 (2)

We focus on (2), since (1) is easier to verify and will come from the same reasoning. The first steps are similar to [5]. We underline the most consequent deviations from the proof of [5] when they occur.

Equivalent formulation of the convergence We write Eq. (2) under an equivalent form. First note that  $F_n(\theta_n) = \hat{\theta}(\tilde{S}_n(\theta_n))$  and  $T(\theta_n) = \hat{\theta}(\bar{S}(\theta_n))$ . Hence  $|f \circ F_n(u_n) - f \circ \hat{\theta}(\tilde{S}_n(\theta_n))| = |f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\tilde{S}(\theta_n))|$ . To show Eq. (2):

$$|f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\bar{S}(\theta_n))| \mathbb{1}_{\theta_n \in K} \xrightarrow[n \infty]{} 0,$$

it is sufficient and necessary to have:

$$\forall \epsilon > 0, \, \exists N \in \mathbb{N}, \, \forall n \ge N, \, |f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\bar{S}(\theta_n))| \mathbb{1}_{\theta_n \in K} \le \epsilon \,.$$

An other equivalent formulation is that there are a finite number of integers n such that  $|f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\bar{S}(\theta_n))| \mathbb{1}_{\theta_n \in K} > \epsilon$ , in other words:

$$\forall \epsilon > 0, \sum_{n=1}^{\infty} \mathbb{1}_{|f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\bar{S}(\theta_n))| \mathbb{1}_{\theta_n \in K} > \epsilon} < \infty.$$

Use the uniform continuity We aim to relate the proximity between the images  $f \circ \hat{\theta}$ of to the proximity between the antecedents of  $f \circ \hat{\theta}$ . The function  $f \circ \hat{\theta} : \mathbb{R}^q \to \mathbb{R}$  is continuous, but not necessarily uniformly continuous on  $\mathbb{R}^q$ . As a consequence, we will need to restrict ourselves to a compact to get uniform continuity properties. We already have a given compact K.  $\tilde{S} : \Theta \to \mathbb{R}^l$  is continuous, hence S(K)is a compact as well. Let  $\delta$  be a strictly positive real number. Let  $\bar{S}(K, \delta) :=$  $\left\{s \in \mathbb{R}^q \left| \inf_{t \in K} ||\bar{S}(t) - s|| \le \delta \right\}$ . Where we use any norm ||.|| on  $\mathbb{R}^q$  since they are all equivalent.  $\overline{S}(K, \delta)$  is a compact set as well. As a consequence  $f \circ \theta$  is uniformly continuous on  $\overline{S}(K, \delta)$ , which means that:

$$\forall \epsilon > 0, \ \exists \eta(\epsilon, \delta) > 0, \ \forall x, y \in \bar{S}(K, \delta), \ \|x - y\| \le \eta(\epsilon, \delta) \implies |f \circ \hat{\theta}(x) - f \circ \hat{\theta}(y)| \le \epsilon.$$
(3)

Let us show that, with  $\alpha := \min(\delta, \eta(\epsilon, \delta)), \forall n$ ,

$$|f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\bar{S}(\theta_n))| \mathbb{1}_{\theta_n \in K} > \epsilon \implies \left\| \tilde{S}_n(\theta_n) - \bar{S}(\theta_n) \right\| \mathbb{1}_{\theta_n \in K} > \alpha.$$
(4)

To that end, we show that:

$$\left\|\tilde{S}_{n}(\theta_{n}) - \bar{S}(\theta_{n})\right\| \mathbb{1}_{\theta_{n} \in K} \leq \alpha \implies |f \circ \hat{\theta}(\tilde{S}_{n}(\theta_{n})) - f \circ \hat{\theta}(\bar{S}(\theta_{n}))| \mathbb{1}_{\theta_{n} \in K} \leq \epsilon.$$

Let us assume that  $\|\tilde{S}_n(\theta_n) - \bar{S}(\theta_n)\| \mathbb{1}_{\theta_n \in K} \leq \alpha$ . If  $\theta_n \notin K$ , then  $|f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\bar{S}(\theta_n))| \mathbb{1}_{\theta_n \in K} = 0 \leq \epsilon$ . If, in contrary,  $\theta_n \in K$ , then  $\bar{S}(\theta_n) \in \bar{S}(K) \subset \bar{S}(K,\delta)$ . Since  $\|\tilde{S}_n(\theta_n) - \bar{S}(\theta_n)\| = \|\tilde{S}_n(\theta_n) - \bar{S}(\theta_n)\| \mathbb{1}_{\theta_n \in K} \leq \alpha \leq \delta$ , then  $\tilde{S}_n(\theta_n) \in \bar{S}(K,\delta)$ . Since  $(\bar{S}(\theta_n), \tilde{S}_n(\theta_n)) \in \bar{S}(K,\delta)^2$  and  $\|\tilde{S}_n(\theta_n) - \bar{S}(\theta_n)\| \leq \alpha \leq \eta(\epsilon, \delta)$ , then we get from Eq. (3)

$$|f \circ \hat{\theta}(\hat{S}_n(\theta_n)) - f \circ \hat{\theta}(\bar{S}(\theta_n))| \mathbb{1}_{\theta_n \in K} \le \epsilon.$$

In both cases, we get that:

$$\left\|\tilde{S}_{n}(\theta_{n})-\bar{S}(\theta_{n})\right\|\mathbb{1}_{\theta_{n}\in K}\leq\alpha\implies|f\circ\hat{\theta}(\tilde{S}_{n}(\theta_{n}))-f\circ\hat{\theta}(\bar{S}(\theta_{n}))|\mathbb{1}_{\theta_{n}\in K}\leq\epsilon,$$

which proves Eq. (4).

Sufficient condition for convergence We use Eq. (4) to find a sufficient condition for (2). This part differs from [5] as our approximation is not defined as a random sum. Eq. (4) is equivalent to

$$\mathbb{1}_{|f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\bar{S}(\theta_n))| \mathbb{1}_{\theta_n \in K} > \epsilon} \leq \mathbb{1}_{\|\tilde{S}_n(\theta_n) - \bar{S}(\theta_n)\| \mathbb{1}_{\theta_n \in K} > \alpha}.$$

From that, we get

$$\forall \epsilon > 0, \, \exists \alpha > 0 \, \sum_{n=1}^{\infty} \mathbbm{1}_{|f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\tilde{S}(\theta_n))| \mathbbm{1}_{\theta_n \in K} > \epsilon} \leq \sum_{n=1}^{\infty} \mathbbm{1}_{\left\|\tilde{S}_n(\theta_n) - \bar{S}(\theta_n)\right\| \mathbbm{1}_{\theta_n \in K} > \alpha}.$$

As a consequence, if

$$\forall \alpha > 0, \ \sum_{n=1}^{\infty} \mathbb{1}_{\left\|\tilde{S}_{n}(\theta_{n}) - \bar{S}(\theta_{n})\right\| \mathbb{1}_{\theta_{n} \in K} > \alpha} < \infty$$

Then

$$\forall \epsilon > 0, \sum_{n=1}^{\infty} \mathbb{1}_{|f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\bar{S}(\theta_n))| \mathbb{1}_{\theta_n \in K} > \epsilon} < \infty$$

In other, equivalent, words:

If 
$$\|\tilde{S}_n(\theta_n) - \bar{S}(\theta_n)\| \mathbb{1}_{\theta_n \in K} \xrightarrow[n\infty]{} 0$$
  
Then  $\|f \circ \hat{\theta}(\tilde{S}_n(\theta_n)) - f \circ \hat{\theta}(\bar{S}(\theta_n))\| \mathbb{1}_{\theta_n \in K} \xrightarrow[n\infty]{} 0.$ 
(5)

Hence, having for all compact sets  $K \subset \Theta$ ,  $\|\tilde{S}_n(\theta_n) - \bar{S}(\theta_n)\| \mathbb{1}_{\theta_n \in K} \xrightarrow[n\infty]{} 0$  is sufficient to have the desired condition (2). Similarly, we find that  $\forall \theta \in K_0$ :

$$\begin{aligned} \left\| \tilde{S}_n(\theta) - \bar{S}(\theta) \right\| &\xrightarrow[n\infty]{} 0 \\ \Longrightarrow \left| f \circ \hat{\theta}(\tilde{S}_n(\theta)) - f \circ \hat{\theta}(\bar{S}(\theta)) \right| &\xrightarrow[n\infty]{} 0 \,, \end{aligned}$$
(6)

which gives us a sufficient condition for (1).

Further simplifications of the desired result with successive sufficient conditions We find another, simpler, sufficient condition for (2) from Eq. (5). This part is unique to our proof and absent from [5]. It is here that we relate the formal conditions of Proposition 2 to the specific assumptions of our Theorem 1.

We first remove the dependency on the terms  $\{\theta_n\}_n$  of the EM sequence:

$$\left\|\tilde{S}_{n}(\theta_{n}) - \bar{S}(\theta_{n})\right\| \mathbb{1}_{\theta_{n} \in K} \leq \sup_{\theta \in K} \left\|\tilde{S}_{n}(\theta) - \bar{S}(\theta)\right\|.$$
(7)

From Eq. (5), (6) and (7) we get that:

Α

$$C \text{ compact } K \subset \Theta, \quad \sup_{\theta \in K} \left\| \tilde{S}_n(\theta) - \bar{S}(\theta) \right\| \underset{n \infty}{\longrightarrow} 0,$$

is a sufficient condition to have both Eq. (1) and (2).

To show that the hypotheses of Theorem 1 imply this sufficient condition, we express it in integral form. Let  $S = \{S_u\}_{u=1,...,q}$ . We recall that  $\tilde{S}_n(\theta) = \{\int_z S_u(z)\tilde{p}_{\theta,n}(z)dz\}_i$  and  $\bar{S}(\theta) = \{\int_z S_u(z)p_{\theta}(z)dz\}_u$ . Hence:

$$\tilde{S}_n(\theta) - \bar{S}(\theta) = \left\{ \int_z S_u(z) \left( \tilde{p}_{\theta,n}(z) - p_{\theta}(z) \right) dz \right\}_u.$$

These q terms can be upper bounded by two different terms depending on the existence of the involved quantities:

$$\int_{z} S_{u}(z) \left( \tilde{p}_{\theta,n}(z) - p_{\theta}(z) \right) dz \leq \left( \int_{z} S_{u}(z)^{2} dz \right)^{\frac{1}{2}} \left( \int_{z} \left( \tilde{p}_{\theta,n}(z) - p_{\theta}(z) \right)^{2} dz \right)^{\frac{1}{2}}$$

and

$$\int_{z} S_{u}(z) \left( \tilde{p}_{\theta,n}(z) - p_{\theta}(z) \right) dz \leq \left( \int_{z} S_{u}(z)^{2} p_{\theta}(z) dz \right)^{\frac{1}{2}} \left( \int_{z} \left( \frac{\tilde{p}_{\theta,n}(z)}{p_{\theta}(z)} - 1 \right)^{2} p_{\theta}(z) dz \right)^{\frac{1}{2}}.$$

As a consequence, if  $\int_z S_u(z)^2 dz$  exists, then it is sufficient to show have:

$$\sup_{\theta \in K} \int_{z} \left( \tilde{p}_{\theta,n}(z) - p_{\theta}(z) \right)^{2} dz \xrightarrow[n \infty]{} 0 \,,$$

and if  $\int_z S_u(z)^2 p_\theta(z) dz$  exists, then it is sufficient to show have:

$$\sup_{\theta \in K} \int_{z} \left( \frac{\tilde{p}_{\theta,n}(z)}{p_{\theta}(z)} - 1 \right)^{2} p_{\theta}(z) dz \xrightarrow[n \infty]{} 0.$$

Among the assumptions of Theorem 1 is one that states that for all compacts  $K \subseteq \Theta$ , one of those scenarios has to be true. Hence our sufficient condition is met.

Conclusion With the hypothesis of Theorem 1, we have

$$\forall \text{ compact } K \subseteq \Theta, \quad \sup_{\theta \in K} \left\| \tilde{S}_n(\theta) - \bar{S}(\theta) \right\| \underset{n \infty}{\longrightarrow} 0,$$

which is a sufficient condition to verify both Eq. (1) and (2). With these two conditions, we can apply Proposition 2.

## 2.1.2 Applying 2

Since we verify all the conditions of Proposition 2, we can apply its conclusions:

With probability 1,  $\limsup_{n \to \infty} p_n < \infty$  and  $\{\theta_n\}_n$  compact sequence,

which is specifically the result (i)(a) of Theorem 1.

### 2.1.3 Verifying the conditions of 1

With Proposition 1, we prove the remaining points of Theorem 1: (i)(b) and (ii).

For the application of Proposition 1:

- $Cl(\{\theta_n\}_n)$  plays the part of the compact K
- $\{\theta \in \Theta | \nabla f(\theta) = 0\} = \{\theta \in \Theta | T(\theta) = \theta\} \text{ plays the part of the set } \mathcal{L}$
- The likelihood f is the  $C^0$  Lyapunov function with regards to  $(T, \mathcal{L})$
- $\{\theta_n\}_n$  is the K valued sequence (since K is  $Cl(\{\theta_n\}_n)$ ).

The last condition that remains to be shown to apply Proposition 1 is that:

$$\lim_{n \to \infty} |f(\theta_{n+1}) - f \circ T(\theta_n)| = 0.$$

We have more or less already proven that, in the previous section of the Proof, with  $F_n(\theta_n)$  in place of  $\theta_{n+1}$ . The only indices where  $F_n(\theta_n) \neq \theta_{n+1}$  are when the value of the sequence  $p_n$  experiences an increment of 1. We have proven with Proposition 2 that there is only a finite number of such increments.

$$|f(\theta_{n+1}) - f \circ T(\theta_n)| = |f(\theta_0) - f \circ T(\theta_n)| \mathbb{1}_{p_{n+1} = p_n + 1} + |f \circ F_n(\theta_n) - f \circ T(\theta_n)| \mathbb{1}_{p_{n+1} = p_n}.$$

Since there is only a finite number of increments of the value of  $p_n$ , then  $\exists N \in \mathbb{N}, \forall n \geq N, \mathbb{1}_{p_{n+1}=p_n+1} = 0$  and  $\mathbb{1}_{p_{n+1}=p_n} = 1$ . In other words:

$$\exists N \in \mathbb{N}, \forall n \ge N, |f(\theta_{n+1}) - f \circ T(\theta_n)| = |f \circ F_n(\theta_n) - f \circ T(\theta_n)| \exists N \in \mathbb{N}, \forall n \ge N, |f(\theta_{n+1}) - f \circ T(\theta_n)| = |f \circ F_n(\theta_n) - f \circ T(\theta_n)| \mathbb{1}_{\theta_n \in Cl(\{\theta_k\}_k)}.$$

Since  $\theta_n$  is always in  $Cl(\{\theta_k\}_k)$  by definition. Additionally Proposition 2 tells us that  $Cl(\{\theta_k\}_k)$  is a compact. Moreover, in order to use Proposition 2 in the first place, we had proven that:

$$\forall \text{ compact } K \subseteq \Theta, \quad \lim_{n \to \infty} |f \circ F_n(\theta_n) - f \circ T(\theta_n)| \mathbb{1}_{\theta_n \in K} = 0$$

We can apply this directly with  $K = Cl(\{\theta_k\}_k)$  to conclude the desired result:

$$\lim_{n \to \infty} |f(\theta_{n+1}) - f \circ T(\theta_n)| = 0$$

Hence we verify all the conditions to apply Proposition 1.

### 2.1.4 Applying 1

Since we verify all we need, we have the conclusions of Proposition 1:

- −  ${f(\theta_n)}_{n \in \mathbb{N}}$  converges towards a connected component of  $f(\mathcal{L} \cap Cl({\theta_n}_n)) \subset f(\mathcal{L})$
- If  $f(\mathcal{L} \cap Cl(\{\theta_n\}_n))$  has an empty interior, then  $\{f(\theta_n)\}_{n \in \mathbb{N}}$  converges towards a  $f^* \in \mathbb{R}$  and  $\{\theta_n\}_n$  converges towards  $\mathcal{L}_{f^*} \cap Cl(\{\theta_n\}_n)$ . Where  $\mathcal{L}_{f^*} := \{\theta \in \mathcal{L} | f(\theta) = f^*\}$

Both points are respectively the statements (i)(b) and (ii) of Theorem 1. Which concludes the proof of the Theorem.

## 2.2 Proof of the tempering theorem

In this Section, we prove Theorem 3 of the main paper, the convergence of the tempered EM algorithm. For that, we need to show that we verify each of the hypothesis of the more general Theorem 1.

We already assumed the conditions M1, M2 and M3 in the hypothesis of Theorem 3. To apply Theorem 1, we need to show that when  $\tilde{p}_{\theta,n}(z) := \frac{p_{\theta}^{\frac{1}{T_n}}(z)}{\int_{z'} p_{\theta}^{\frac{1}{T_n}}(z')dz'}$ , then  $\forall$  compact  $K \subseteq \Theta$ , one of the two following configurations holds:

$$\int_{z} S(z)^{2} dz < \infty \text{ and } \sup_{\theta \in K} \int_{z} \left( \tilde{p}_{\theta,n}(z) - p_{\theta}(z) \right)^{2} dz \underset{n \infty}{\longrightarrow} 0,$$

or

$$\sup_{e \in K} \int_{z} S(z)^{2} p_{\theta}(z) dz < \infty \text{ and } \sup_{\theta \in K} \int_{z} \left( \frac{\tilde{p}_{\theta,n}(z)}{p_{\theta}(z)} - 1 \right)^{2} p_{\theta}(z) dz \underset{n \infty}{\longrightarrow} 0$$

Since we have assumed:

S A

$$\forall \text{ compact } K \in \Theta, \ \forall \alpha \in \overline{\mathcal{B}}(1,\epsilon), \forall u, \quad \sup_{\theta \in K} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz < \infty,$$

then we already verify the first half of the second configuration for all the compacts K. Hence it is sufficient to prove that:

$$\forall \text{ compact } K \in \Theta, \ \sup_{\theta \in K} \int_{z} \left( \frac{\tilde{p}_{\theta,n}(z)}{p_{\theta}(z)} - 1 \right)^{2} p_{\theta}(z) dz \xrightarrow[n \infty]{} 0,$$
(8)

to have the desired result. The rest of the proof is dedicated to this goal.

#### 2.2.1 Taylor development

We use the Taylor's formula of the first order with the mean-value form of the reminder. For a derivable function g:

$$g(x) = g(0) + g'(a)x, \quad a \in [0, x],$$
(9)

where the interval [0, x] has a flexible meaning since x could be negative. We apply it to:

$$g(x) = e^x$$
,  $g'(x) = e^x$ ,  $g(x) = 1 + xe^a$ ,  $a \in [0, x]$ ,

and:

$$g(x) = \frac{1}{1+x}, \quad g'(x) = -\frac{1}{(1+x)^2}, \quad g(x) = 1 - \frac{x}{(1+a)^2}, \quad a \in [0,x].$$

To make the upcoming calculation more readable, we momentarily replace  $p_{\theta}(z)$  by simply p and  $T_n$  by T.

$$p^{\frac{1}{T}} = p\left(p^{\frac{1}{T}-1}\right) = pe^{(\frac{1}{T}-1)ln p} = p + \left(\frac{1}{T}-1\right)p ln p e^{a}, \quad a \in \left[0, \left(\frac{1}{T}-1\right)ln p\right],$$

where  $a = a(z, \theta, T_n)$  since it depends on the value of  $p_{\theta}(z)$  and  $T_n$ . Provided that the following quantities are defined, we have:

$$\int_{z} p^{\frac{1}{T}} = 1 + \left(\frac{1}{T} - 1\right) \int_{z} p \ln p e^{a},$$

Hence:

$$\frac{1}{\int_{z} p^{\frac{1}{T}}} = 1 - \left(\frac{1}{T} - 1\right) \frac{\int_{z} p \ln p e^{a}}{\left(1 + b\right)^{2}}, \quad b \in \left[0, \left(\frac{1}{T} - 1\right) \int_{z} p \ln p e^{a}\right],$$

where  $b = b(\theta, T_n)$  since it depends on the value of  $T_n$  the integral over z of a function of z and  $\theta$ .

In the end, we have:

$$\frac{p^{\frac{1}{T}}}{\int_{z} p^{\frac{1}{T}}} = p + \left(\frac{1}{T} - 1\right) p \ln p e^{a} \left(1 - \left(\frac{1}{T} - 1\right) \frac{\int_{z} p \ln p e^{a}}{\left(1 + b\right)^{2}}\right) - \left(\frac{1}{T} - 1\right) p \frac{\int_{z} p \ln p e^{a}}{\left(1 + b\right)^{2}}.$$
(10)

Since for any real numbers  $(x + y)^2 \le 2(x^2 + y^2)$ , then:

$$\begin{split} \left(\frac{p^{\frac{1}{T}}}{\int_{z} p^{\frac{1}{T}}} - p\right)^{2} &\leq 2\left(\frac{1}{T} - 1\right)^{2} p^{2} \left(\left(\ln p \, e^{a}\right)^{2} \left(1 - \left(\frac{1}{T} - 1\right) \frac{\int_{z} p \ln p \, e^{a}}{\left(1 + b\right)^{2}}\right)^{2} + \left(\frac{\int_{z} p \ln p \, e^{a}}{\left(1 + b\right)^{2}}\right)^{2}\right) \\ &= 2\left(\frac{1}{T} - 1\right)^{2} p^{2} \left(\left(\ln p \, e^{a}\right)^{2} A + B\right) \,. \end{split}$$

where  $A = A(\theta, T_n)$  and  $B = B(\theta, T_n)$ . So far the only condition that has to be verified for all the involved quantities to be defined is that  $\int_z p \ln p e^a$  exists. With this Taylor development on hand, we state, prove and apply two lemmas which allow us to get (8) and conclude the proof of the theorem.

# 2.2.2 Two intermediary lemmas

The two following lemmas provides every result we need to finish the proof.

## Lemma 1 With

$$p_{\theta}(z) = exp\left(\psi(\theta) + \langle S(z), \phi(\theta) \rangle\right)$$

then

$$\int_{z} p_{\theta}^{\alpha}(z) \ln^{2} p_{\theta}(z) dz \leq 2\psi(\theta)^{2} \int_{z} p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^$$

and

$$\int_{z} p_{\theta}^{\alpha}(z) \left| \ln p_{\theta}(z) \right| dz \leq \left| \psi(\theta) \right| \int_{z} p_{\theta}^{\alpha}(z) dz + \left\| \phi(\theta) \right\| \cdot \left( \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) \int_{z} p_{\theta}^{\alpha}(z) \right)^{\frac{1}{2}}.$$

*Proof* For the first inequality, using the fact that  $(a + b)^2 \le 2(a^2 + b^2)$ , we have:

$$\int_{z} p_{\theta}^{\alpha}(z) \ln^{2} p_{\theta}(z) dz \leq 2\psi(\theta)^{2} \int_{z} p_{\theta}^{\alpha}(z) dz + 2 \int_{z} p_{\theta}^{\alpha}(z) \left\langle S(z), \phi(\theta) \right\rangle^{2} ,$$

We use Cauchy-Schwartz:

$$\langle S(z), \phi(\theta) \rangle^2 \le \|\phi\|^2 \|S(z)\|^2 = \|\phi\|^2 \sum_u S_u(z)^2,$$

to get the desired result:

$$\int_{z} p_{\theta}^{\alpha}(z) \ln^{2} p_{\theta}(z) dz \leq 2\psi(\theta)^{2} \int_{z} p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) \,.$$

For the second inequality, we start with Cauchy-Schwartz on  $\left\langle \int_{z} S(z) p_{\theta}^{\alpha}(z), \phi(\theta) \right\rangle$ :

$$\int_{z} p_{\theta}^{\alpha}(z) \left| \ln p_{\theta}(z) \right| dz \leq \left| \psi(\theta) \right| \int_{z} p_{\theta}^{\alpha}(z) dz + \left\| \phi(\theta) \right\| \cdot \left\| \int_{z} S(z) p_{\theta}^{\alpha}(z) \right\|$$

Moreover, since:

$$\int_{z} S_{u}(z) p_{\theta}^{\alpha}(z) dz \leq \left( \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz \right)^{\frac{1}{2}} \left( \int_{z} p_{\theta}^{\alpha}(z) dz \right)^{\frac{1}{2}},$$

then

$$\int_{z} p_{\theta}^{\alpha}(z) \left| \ln p_{\theta}(z) \right| dz \leq \left| \psi(\theta) \right| \int_{z} p_{\theta}^{\alpha}(z) dz + \left\| \phi(\theta) \right\| \cdot \left( \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) \int_{z} p_{\theta}^{\alpha}(z) \right)^{\frac{1}{2}}.$$

**Lemma 2** With K compact and  $\epsilon \in \mathbb{R}^*_+$ ,

$$p_{\theta}(z) = exp\left(\psi(\theta) + \langle S(z), \phi(\theta) \rangle\right),$$

and

$$\tilde{p}_{\theta,n}(z):=\frac{p_{\theta}^{\frac{1}{T_n}}(z)}{\int_{z'}p_{\theta}^{\frac{1}{T_n}}(z')dz'},$$

 $i\!f$ 

 $\begin{array}{ll} (i) \ T_n \in \mathbb{R}^*_+ \underset{n \infty}{\longrightarrow} 1 & , \\ (ii) \ \sup_{\theta \in K} \psi(\theta) < \infty & , \\ (iii) \ \sup_{\theta \in K} \|\phi(\theta)\| < \infty & , \\ (iv) \ \forall \alpha \in \overline{\mathcal{B}}(1, \epsilon), & \sup_{\theta \in K} \int_z p_{\theta}^{\alpha}(z) dz < \infty & , \\ (v) \ \forall \alpha \in \overline{\mathcal{B}}(1, \epsilon), \ \forall u, \quad \sup_{\theta \in K} \int_z S_u^2(z) p_{\theta}^{\alpha}(z) dz < \infty & . \end{array}$ 

then

$$\sup_{\theta \in K} \int_{z} \left( \frac{\tilde{p}_{\theta,n}(z)}{p_{\theta}(z)} - 1 \right)^{2} p_{\theta}(z) dz \xrightarrow[n \infty]{} 0.$$

*Proof* Provided that the following integrals exist, we have, thanks to the Taylor development:

$$\int_{z} \frac{1}{p} \left( \frac{p^{\frac{1}{T}}}{\int_{z} p^{\frac{1}{T}}} - p \right)^{2} \leq 2 \int_{z} \left( \frac{1}{T} - 1 \right)^{2} p \left( (\ln p e^{a})^{2} A + B \right)$$

$$= 2 \left( \frac{1}{T} - 1 \right)^{2} A \int_{z} p e^{2a} \ln^{2} p + 2 \left( \frac{1}{T} - 1 \right)^{2} B.$$
(11)

In this proof, we find finite upper bounds independent of  $\theta$  and  $T_n$  for  $A(\theta, T_n)$ ,  $B(\theta, T_n)$  and  $\int_z p e^{2a} ln^2 p$ , then - since  $\left(\frac{1}{T_n} - 1\right) \longrightarrow 0$  - we have the desired result. We start by studying  $A(\theta, T_n) = \left(1 - \left(\frac{1}{T} - 1\right) \frac{\int_z p \ln p e^a}{(1+b)^2}\right)^2$ . The first term of interest here is  $\int_z p \ln p e^a$ . We have:

$$\begin{split} a &\in \left[0, \left(\frac{1}{T} - 1\right) \ln p\right] \\ e^a &\in \left[1, p^{\frac{1}{T} - 1}\right], \\ p \ln p \, e^a &\in \left[p \ln p, p^{\frac{1}{T}} \ln p\right]. \end{split}$$

where we recall that the interval is to be taken in a flexible sense, since we do not now a priory which bound is the largest and which is the smallest. What we have without doubt though is:

$$\left| p \ln p e^{a} \right| \leq max \left( \left| p \ln p \right|, \left| p^{\frac{1}{T}} \ln p \right| \right).$$

We find an upper bound on both those term. Let  $\alpha \in \overline{\mathcal{B}}(1, \epsilon)$ , the second result of Lemma 1 gives us:

$$\int_{z} p_{\theta}^{\alpha}(z) \left| \ln p_{\theta}(z) \right| dz \leq \left| \psi(\theta) \right| \int_{z} p_{\theta}^{\alpha}(z) dz + \left\| \phi(\theta) \right\| \cdot \left( \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) \int_{z} p_{\theta}^{\alpha}(z) \right)^{\frac{1}{2}} \cdot \left( \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) \right)^{\frac{1}{2}} dz$$

10

Thanks to the hypothesises (ii), (iii), (iv) and (v), we have:

$$\begin{split} \int_{z} p_{\theta}^{\alpha}(z) \, \left| \ln p_{\theta}(z) \right| dz &\leq \sup_{\theta \in K} \, \left| \psi(\theta) \right| \cdot \sup_{\theta \in K} \, \int_{z} p_{\theta}^{\alpha}(z) dz \\ &+ \sup_{\theta \in K} \, \left\| \phi(\theta) \right\| \cdot \sum_{u} \left( \sup_{\theta \in K} \, \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) \right)^{\frac{1}{2}} \cdot \left( \sup_{\theta \in K} \, \int_{z} p_{\theta}^{\alpha}(z) \right)^{\frac{1}{2}} \\ &=: C(\alpha) \\ &< \infty \,. \end{split}$$

The upper bound  $C(\alpha)$  in the previous inequality is independent of  $\theta$  and z but still dependent of the exponent  $\alpha$ . However, since  $\overline{\mathcal{B}}(1,\epsilon)$  is closed ball, hypothesises (iv) and (v) can be rephrased as:

$$\begin{aligned} &(iv) \sup_{\alpha \in \overline{\mathcal{B}}(1,\epsilon)} \sup_{\theta \in K} \int_{z} p_{\theta}^{\alpha}(z) dz < \infty \,, \\ &(v) \; \forall u, \; \sup_{\alpha \in \overline{\mathcal{B}}(1,\epsilon)} \sup_{\theta \in K} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz < \infty \end{aligned}$$

Hence we can actually take the supremum in  $\alpha$  as well:

$$\begin{split} \int_{z} p_{\theta}^{\alpha}(z) \, \left| \ln p_{\theta}(z) \right| dz &\leq \sup_{\theta \in K} \, \left| \psi(\theta) \right| \, . \sup_{\alpha \in \overline{\mathcal{B}}(1,\epsilon)} \sup_{\theta \in K} \, \int_{z} p_{\theta}^{\alpha}(z) dz \\ &+ \sup_{\theta \in K} \, \left\| \phi(\theta) \right\| \, . \sum_{u} \left( \sup_{\alpha \in \overline{\mathcal{B}}(1,\epsilon)} \sup_{\theta \in K} \, \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) \right)^{\frac{1}{2}} \, . \left( \sup_{\alpha \in \overline{\mathcal{B}}(1,\epsilon)} \sup_{\theta \in K} \, \int_{z} p_{\theta}^{\alpha}(z) \right)^{\frac{1}{2}} \\ &=: C' \\ &< \infty \, . \end{split}$$

This new upper bound C' is independent of  $\alpha$ .

Since  $T_n \mapsto 1$ , then  $\exists N \in \mathbb{N}, \forall n \geq N, \frac{1}{T_n} \in \overline{\mathcal{B}}(1, \epsilon)$ . Hence for  $n \geq N$ , we can apply the previous inequation to either  $\alpha = 1$  or  $\alpha = \frac{1}{T_n}$ . Which gives us that  $\int_z p_{\theta}(z) |\ln p_{\theta}(z)|, \int_z p_{\theta}^{\frac{1}{T_n}}(z) |\ln p_{\theta}(z)|$  and their supremum in  $\theta$  are all finite, all of them upper bounded by C'.

In the end, when  $n \ge N$ , we have the control  $\sup_{\theta \in K} \left| \int_z p \ln p \, e^a \right| < C'.$ 

The next term to control is  $\frac{1}{(1+b)^2}$ . Since  $b \in [0, (\frac{1}{T}-1)\int_z p\ln p e^a]$ , then  $|b| \leq (\frac{1}{T}-1)\sup_{\substack{\theta \in K}}\int_z p\ln p e^a$ . We already established that for all  $n \geq N$ ,  $\sup_{\substack{\theta \in K}}|\int_z p\ln p e^a| \leq C' < \infty$ , hence  $\sup_{\substack{\theta \in K}}|b(\theta, T_n)| \underset{T_n \longrightarrow 1}{\longrightarrow} 0$ . In particular,  $\exists N' \in \mathbb{N}, \forall n \geq N', \forall \theta \in K$  we have  $|b(\theta, T_n)| \leq \frac{1}{2}$ . In that case:

$$(1+b)^2 > (1-|b|)^2 \ge \frac{1}{4}$$
  
 $\frac{1}{(1+b)^2} < \frac{1}{(1-|b|)^2} \le 4.$ 

In the end, when  $n \ge max(N, N')$ , for any  $\theta \in K$ :

$$\begin{aligned} A(\theta, T_n) &\leq 2 + 2\left(\frac{1}{T_n} - 1\right)^2 \left(\frac{\int_z p \ln p \, e^a}{\left(1 + b\right)^2}\right)^2 \\ &\leq 2 + 32\left(\frac{1}{T_n} - 1\right)^2 \left(\sup_{\theta \in K} \int_z p \ln p \, e^a\right)^2 \\ &\leq 2 + 32\left(\frac{1}{T_n} - 1\right)^2 C'^2 \\ &\leq 2 + 32\epsilon^2 C'^2 \\ &=: C_1 \,. \end{aligned}$$

This upper bound does not depend en  $\theta$  anymore and the part in  $T_n$  simply converges towards 0 when  $T_n \longrightarrow 1$ .

Treating the term  $B(\theta, T_n) = \left(\frac{\int_z p \ln p e^a}{(1+b)^2}\right)^2 \le 16 \left(\sup_{\theta \in K} \int_z p \ln p e^a\right)^2 \le 16C'^2 =: C_2$ is immediate after having dealt with  $A(\theta, T_n)$ .

We now treat the term  $\int_z p e^{2a} ln^2 p$  in the exact same fashion as we did  $A(\theta, T_n)$ :

$$p \ln p e^{a} \in \left[ p \ln p, p^{\frac{1}{T}} \ln p \right]$$
$$\implies p \left( \ln p e^{a} \right)^{2} \in \left[ p \ln^{2} p, p^{\frac{2}{T} - 1} \ln^{2} p \right]$$
$$\implies p \left( \ln p e^{a} \right)^{2} \le \max(p \ln^{2} p, p^{\frac{2}{T} - 1} \ln^{2} p).$$

We control those two terms as previously. First we apply Lemma 1 (its first result this time) with  $\alpha \in \overline{\mathcal{B}}(1, \epsilon)$ .

$$\int_{z} p_{\theta}^{\alpha}(z) \ln^{2} p_{\theta}(z) dz \leq 2\psi(\theta)^{2} \int_{z} p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z) dz + 2 \left\|\phi(\theta)\right\|^{2} \cdot \sum_{u} \int_{z} S_{u}^{2}(z$$

Thanks to the hypothesis (ii), (iii), (iv) and (v), we can once again take the supremum of the bound over  $\theta \in K$ , then over  $\alpha \in \overline{\mathcal{B}}(1, \epsilon)$  and conserve finite quantities:

.

$$\begin{split} \int_{z} p_{\theta}^{\alpha}(z) \ln^{2} p_{\theta}(z) dz &\leq 2 \sup_{\theta \in K} \psi(\theta)^{2} . \sup_{\alpha \in \overline{\mathcal{B}}(1,\epsilon)} \sup_{\theta \in K} \int_{z} p_{\theta}^{\alpha}(z) dz \\ &+ 2 \sup_{\theta \in K} \|\phi(\theta)\|^{2} . \sum_{u} \sup_{\alpha \in \overline{\mathcal{B}}(1,\epsilon)} \sup_{\theta \in K} \int_{z} S_{u}^{2}(z) p_{\theta}^{\alpha}(z) \\ &=: C_{3} \\ &\leq \infty . \end{split}$$

The previous result is true for  $\alpha = 1$ , and since once again  $\exists N'', \forall n \geq N'', \frac{2}{T_n} - 1 \in \overline{\mathcal{B}}(1,\epsilon) \cap \mathbb{R}^*_+$ , it is also true for  $\alpha = \frac{2}{T_n} - 1$  when *n* is large enough.  $C_3$  is independent of *z*,  $\theta$  and  $T_n$ . In the end  $\forall n \geq N'', \int_z p e^{2a} ln^2 p \leq C_3 < \infty$ .

.

We replace the three terms  $A(\theta, T_n)$ ,  $B(\theta, T_n)$  and  $\int_z p e^{2a} ln^2 p$  by their upper bounds in the inequality (11). When  $n \ge max(N, N', N'')$ :

$$\int_{z} \frac{1}{p} \left( \frac{p^{\frac{1}{T}}}{\int_{z} p^{\frac{1}{T}}} - p \right)^{2} \le 2 \left( \frac{1}{T_{n}} - 1 \right)^{2} C_{1} C_{3} + 2 \left( \frac{1}{T_{n}} - 1 \right)^{2} C_{2}$$

Which converges towards 0 when  $T_n \longrightarrow 1$ , i.e. when  $n \longrightarrow \infty$ . This concludes the proof of the lemma.

#### 2.2.3 Verifying the conditions of Lemma 2

Now that the lemmas are proven, all that remains is to apply Lemma 2.

(i) We have  $T_n \in \mathbb{R}^*_+ \xrightarrow[n\infty]{} 1$  by hypothesis.

(ii) and (iii)  $\sup_{\theta \in K} \psi(\theta) < \infty$  and  $\sup_{\theta \in K} \|\phi(\theta)\| < \infty$  are implied by the fact that  $\psi(\theta) = \psi'(\theta) - \log g(\theta)$  and  $\phi(\theta)$  are continuous

(*iv*) and (*v*) Are also hypothesis of the theorem. Hence we can apply Lemma 2. This means that:

$$\sup_{\theta \in K} \int_{z} \left( \frac{\tilde{p}_{\theta,n}(z)}{p_{\theta}(z)} - 1 \right)^{2} p_{\theta}(z) dz \xrightarrow[n \infty]{} 0$$

With this last condition verified, we can apply Theorem 1. Which concludes the proof.

### 3 Experiments on tmp-EM with Mixtures of Gaussian

In this section, we present more detailed experiments analysing the tempered EM and comparing it to the regular EM. As in the main paper, we focus on likelihood maximisation within the Gaussian Mixture Model. From the optimisation point of view, we demonstrate that tmp-EM does not fall in the first local maximum like EM does but instead consistently finds better one. From the machine learning point of view, we illustrate how tmp-EM is able to better identify the real GMM parameters even when they are ambiguous and when the initialisation is voluntarily tricky.

The only constraints on the temperature profile is that  $T_n \rightarrow 1$  and  $T_n > 0$ . We use two different temperature profiles. First, a decreasing exponential:  $T_n = 1 + (T_0 - 1) \exp(-r.n)$ . We call it the "simple" profile, it works most of the time. Second, we examine the capabilities of a profile with oscillations in addition to the main decreasing trend. These oscillations are meant to momentarily increase the convergence speed to "lock-in" some of the most obviously good decisions of the algorithm, before re-increasing the temperature and continuing the exploration on the other, more ambiguous parameters. Those two regimes are alternated in succession with gradually smaller oscillations, resulting in a multi-scale procedure that "locks-in" gradually harder decisions. The formula is taken from [1]:  $T_n = th(\frac{n}{2r}) + (T_0 - b\frac{2\sqrt{2}}{3\pi}) a^{n/r} + b sinc(\frac{3\pi}{4} + \frac{n}{r})$ . The profile used, as well as the values of the hyper-parameters are specified for each experiment. The hyper parameters are chosen by grid-search.

For the sake of comparison, the following Experiment 1 and 2 are similar to the experiments of [1] on the tmp-SAEM.

### 3.1 Experiment 1: 6 clusters

We start by demonstrating the superior performance of the tempered EM algorithm on an example mixture of K = 6 gaussians in dimension p = 2. The real parameters can be visualised on Figure 1, where the real centroids are represented by black crosses and confidence ellipses help visualise the real covariance matrices. In addition, 500 points were simulated in order to illustrate, among other things, the weights of each class. To quantify the ability of each EM method to increase



Fig. 1 500 sample points from a Mixture of Gaussians with 6 classes. The true centroid of each Gaussian are depicted by black crosses, and their true covariance matrices are represented by the confidence ellipses of level 0.8, 0.99 and 0.999 around the centre.

the likelihood and recover the true parameters, we generate from this model 20 different datasets with n = 500 observations. For each of these datasets, we make 200 EM runs, all of them starting from a different random initialisation. To initialise the mixture parameters, we select uniformly 6 data points to act as centroids. In each run, EM and tmp-EM start with the same initialisation. The number K of clusters is known by the algorithms. For this experiment, the simple tempering profile is used with parameters  $T_0 = 50$  and r = 2.

#### 3.1.1 Illustrative

First, we observe on the left of Figure 2, one example of the final states of the EM algorithm. The observations can be seen in green, the initial centroids are

represented by blue crosses, and the parameters  $\{\hat{\mu}_k\}_{k=1}^K$  and  $\{\hat{\Sigma}_k\}_{k=1}^K$  estimated by the EM are represented in orange. In this EM run, one of the estimated clusters became degenerated and, as counterpart, two different real clusters were fused as one by the method. On the right of Figure 2, we observe the final state of the tmp-EM on the same dataset, from the same initialisation. This time all the clusters were properly identified.



Fig. 2 EM and tmp-EM final states on the same simulation with the same initialisation. tmp-EM positioned correctly the estimated centroides, whereas the regular EM made no distinction between the two bottom classes and ended up with a degenerate class instead.

#### 3.1.2 Quantitative

To demonstrate the improvements made by tempering, we present aggregated quantitative results over all the simulated datasets and random initialisations.

Likelihood maximisation EM and tmp-EM are optimisation methods whose target function is the likelihood of the estimated mixture parameters. We represent on Figure 3 the empirical distribution of the negative log-likelihoods reached at the end of the two methods, EM in blue, tmp-EM in orange. On those boxplots, the coloured "box" at the centre contains 50% of the distribution, hence it is delimited by the 0.25 and 0.75 quantiles. The median of the distribution is represented by an horizontal black line inside the box. The space between the whiskers on the other end, contain 90% of the distribution, its limits are the 0.05 and 0.95 quantiles. The table provides the numeric values of these statistics.

We note that the negative log-likelihood reached by tmp-EM is lower on average (higher likelihood) than what EM obtains. Moreover, tmp-EM also has a lower variance, its standard deviation being approximately half of the std of EM. More generally, we observe that the distribution of the final loss of tmp-EM is both shifted towards the lower values and less variable. In particular, each of the followed quantiles are lower for tmp-EM, and both the difference Q95-Q5 (space between whiskers) and Q75-Q25 (size of the box) are lower for tmp-EM. This illustrates that it obtains better, more consistent results on our synthetic example.



Fig. 3 Empirical distribution of the negative log-likelihood reached by the EM algorithms. EM is blue and tmp-EM in orange. The boxplot allow us to identify the quantiles 0.05, 0.25, 0.5, 0.75 and 0.95 of each distribution, as well as the outliers. Their numeric values can be found in the table, the better ones being in **bold**. tmp-EM is better overall.

*Parameter recovery* The EM algorithm is an optimisation procedure. Stricto sensu, the optimised metric - the likelihood - should be the only criterion for success. However, in the case of the Mixture of Gaussians, the underlying Machine Learning stakes are always very visible. Hence we dedicate time to assess the relative success parameter recovery of EM and tmp-EM.

The quality of parameter recovery is always dependent on the number of observation. The larger n, the more the likelihood will describe an actual ad-equation with the real parameters behind the simulation. Additionally, as n grows, the situation becomes less and less ambiguous, until all methods yield either the exact same, or at least very similar solutions, with all of them being fairly close to the truth. All of our simulation are done with n = 500 data points. Not a very large number, but since the lowest weight of our K = 6 classes is around 0.09, it is sufficient for all the classes to be guaranteed to contain several points. The three families of parameters in a GMM are the weights  $\{\pi_k\}_{k=1}^K$  of the K classes. We evaluate the error made on  $\mu$  with the relative different in squared norm 2:  $\frac{\|\hat{\mu}_k - \mu_k\|_2^2}{\|\mu_k\|_2^2}$ . For  $\Sigma$ , we compute the KL divergence between the real matrices and the estimates  $KL(\Sigma_k, \hat{\Sigma}_k) = \frac{1}{2} \left( ln \frac{|\Theta_k|}{|\widehat{\Theta_k}|} + tr(\Sigma_k \widehat{\Theta}_k) - p \right)$ , with  $\Theta := \Sigma^{-1}$  for all those matrices. Finally, the analysis on  $\pi$  is harder to interpret and less interesting, but reveals

Finally, the analysis on  $\pi$  is harder to interpret and less interesting, but reveals the same trend, with lower errors for the tempering.

The error on the averages  $\mu_k$  is usually the most informative and easy to interpret metric, quantifying how well each methods position the class centres. Figure 4 and Table 1 represent the distribution of the relative error  $\frac{\|\hat{\mu}_k - \mu_k\|_2^2}{\|\mu_k\|_2^2}$ . The results of tmp-EM are much better with average and median errors often being orders of magnitude below the errors of EM, with similar or lower variance. The other quantiles of the tmp-EM distribution are also either equivalent to or order of magnitudes below the corresponding EM quantiles. The largest errors happen on Class 3 and 6, two of the ambiguous ones, but are always noticeably smaller and less variable with the tempering.

The KL divergences  $KL(\Sigma_k, \widehat{\Sigma}_k)$  assess whether each the covariances  $\Sigma_k$  of each class are properly replicated. Note that since the computation of the KL divergence



**Fig. 4** Empirical distribution of the relative error in squared norm 2  $\frac{\|\hat{\mu}_k - \mu_k\|_2^2}{\|\mu_k\|_2^2}$  between the real centroid positions in  $\mu$  and the estimations by the EM algorithms.

**Table 1** Quantiles and other statistics describing the empirical distribution of the relative error in squared norm  $2 \frac{\|\hat{\mu}_k - \mu_k\|_2^2}{\|\mu_k\|_2^2}$  between the real centroid positions in  $\mu$  and the estimations by the EM algorithms. The error of tmp-EM is always closer to 0 with lower variance (with the exception of class 2 where the variance is similar).

Cl.		mean	$\operatorname{std}$	Q5	Q25	Q50	Q75	Q95
1	EM	0.024	0.119	6.10 <sup>-6</sup>	$6.10^{-5}$	$2.10^{-4}$	0.002	0.065
	tmp-EM	<b>0.002</b>	<b>0.014</b>	6.10 <sup>-6</sup>	4.10 <sup>-5</sup>	1.10 <sup>-4</sup>	<b>4.10<sup>-4</sup></b>	<b>0.005</b>
2	EM	0.038	<b>0.066</b>	5.10 <sup>-5</sup>	$2.10^{-4}$	0.001	0.057	<b>0.169</b>
	tmp-EM	<b>0.032</b>	0.070	5.10 <sup>-5</sup>	2.10 <sup>-4</sup>	<b>5.10<sup>-4</sup></b>	<b>0.013</b>	0.210
3	EM	0.971	1.153	$4.10^{-4}$	0.004	0.297	2.467	2.736
	tmp-EM	<b>0.743</b>	<b>1.072</b>	3.10 <sup>-4</sup>	<b>0.003</b>	<b>0.235</b>	<b>1.500</b>	<b>2.681</b>
4	EM	0.310	0.487	$7.10^{-5}$	$8.10^{-4}$	0.031	0.859	<b>1.158</b>
	tmp-EM	<b>0.287</b>	<b>0.476</b>	3.10 <sup>-5</sup>	5.10 <sup>-4</sup>	<b>0.025</b>	<b>0.076</b>	1.188
5	EM	0.735	1.248	8.10 <sup>-5</sup>	$5.10^{-4}$	0.002	0.814	3.191
	tmp-EM	<b>0.432</b>	<b>1.054</b>	6.10 <sup>-5</sup>	4.10 <sup>-4</sup>	<b>7.10<sup>-4</sup></b>	<b>0.002</b>	<b>3.180</b>
6	$_{\mathrm{tmp-EM}}^{\mathrm{EM}}$	1.940 <b>0.807</b>	2.828 1.735	$7.10^{-4}$ 4.10 <sup>-4</sup>	0.005 <b>0.002</b>	1.158 <b>0.010</b>	2.743 <b>1.066</b>	6.744 <b>3.243</b>

involves the matrix inverse  $\widehat{\Theta}_k = \widehat{\Sigma}_k^{-1}$ , the outliers cases where a class vanishes in an EM have to be removed: they correspond to pathological, non invertible matrices. Figure 5 and Table 2 describe the distribution of the KL divergence. The Figure is cropped and does not show some of the very rare, most upper outliers (less than 1%). Overall, the results are similar to what we get on  $\mu$ : in terms of average KL and median KL, tmp-EM is better than EM, being either similar on some classes and much better on others. Its standard deviation is also lower - sometimes by one order of magnitude - on all classes except Class 4. The other quantiles are also overall better, with one exception on Q95 of class 4.

*Conclusion* We saw that tmp-EM achieved better average and median results with lower variances both on likelihood maximisation and parameter recovery for every

**Table 2** Quantiles and other statistics describing the empirical distribution of the KL divergence  $KL(\Sigma_k, \hat{\Sigma}_k)$  between each covariance matrix estimated by the EMs and the real covariance matrices  $\Sigma$ . On every class but the 4th, the deviation of tmp-EM is closer to 0 with lower or similar variance.

Cl.		mean	std	Q5	Q25	Q50	Q75	Q95
1	$_{\mathrm{tmp-EM}}$	2.741 <b>0.845</b>	39.879 <b>8.683</b>	0.003 <b>0.003</b>	0.009 <b>0.008</b>	0.017 <b>0.013</b>	0.136 <b>0.055</b>	3.222 1.745
2	$_{\mathrm{tmp-EM}}$	0.852 <b>0.412</b>	<b>9.006</b> 9.072	$\begin{array}{c} 0.004 \\ 0.004 \end{array}$	0.015 <b>0.011</b>	0.042 <b>0.027</b>	0.636 <b>0.34</b>	1.015 <b>0.782</b>
3	$_{\mathrm{tmp-EM}}$	1.185 <b>0.648</b>	14.636 <b>4.435</b>	0.015 <b>0.014</b>	0.078 <b>0.066</b>	0.183 <b>0.174</b>	0.414 <b>0.408</b>	1.742 <b>1.331</b>
4	$_{\mathrm{tmp-EM}}$	<b>2.008</b> 2.998	<b>13.156</b> 20.1	0.008 <b>0.006</b>	0.043 <b>0.028</b>	0.386 <b>0.374</b>	1.034 <b>0.637</b>	<b>4.553</b> 5.468
5	$_{\mathrm{tmp-EM}}^{\mathrm{EM}}$	1.772 <b>0.791</b>	12.175 <b>7.088</b>	0.005 <b>0.005</b>	0.015 <b>0.011</b>	0.035 <b>0.026</b>	0.664 <b>0.058</b>	5.813 <b>2.57</b>
6	$_{\mathrm{tmp-EM}}$	2.909 <b>2.072</b>	59.913 <b>25.898</b>	0.012 <b>0.008</b>	0.045 <b>0.023</b>	0.195 <b>0.062</b>	0.676 <b>0.34</b>	4.371 <b>2.883</b>

**Table 3** Synthetic table focusing solely on the average and standard deviation (in parenthesis) of the losses and parameter reconstruction errors made by EM and tmp-EM. We note that the likelihood reached is higher with lower variance, and similarly, the parameter metrics on almost every class are better with lower variance for tmp-EM.

Metric	class	EM	$\operatorname{tmp-EM}$
$-lnp_{\hat{\theta}}$		$2247.08\;(69.62)$	$2218.80\;(35.21)$
$\frac{\ln p_{\theta_0} - \ln p_{\hat{\theta}}}{\ln p_{\theta_0}}$		$0.12\;(0.04)$	$0.13\;(0.04)$
	$\begin{array}{c} 1 \\ 2 \end{array}$	$egin{array}{c} -0.19 \; (0.36 ) \ 0.11 \; (0.57 ) \end{array}$	$\begin{array}{c} -0.17~(0.29)\\ 0.04~(0.33) \end{array}$
$\frac{\hat{\pi}_k - \pi_k}{\pi_k}$	3	0.56 ( <b>0.81</b> )	0.45(0.83)
	4	<b>0.10</b> (0.57)	0.10 (0.43)
	5	-0.08(0.48)	-0.02(0.31)
	6	-0.20(0.43)	-0.13 (0.40)
	1	0.02(0.12)	$2.10^{-3}\ (0.01)$
	2	0.04 <b>(0.07)</b>	<b>0.03</b> (0.07)
$\frac{\ \hat{\mu}_k - \mu_k\ ^2}{\ \mu_k\ ^2}$	3	$0.97\;(1.15)$	$0.74\;(1.07)$
117.16.11	4	0.31(0.49)	$0.29\;(0.48)$
	5	0.73(1.25)	$0.43\ (1.05\ )$
	6	1.94(2.83)	$0.81\;(1.74)$
	1	2.74 (39.88)	$0.84\ (8.68)$
	2	0.85 <b>(9.01)</b>	0.41 (9.07)
$KL(\Sigma, \widehat{\Sigma})$	3	1.18(14.64)	$0.65\;(4.44)$
	4	$2.01\ (13.16)$	3.00(20.10)
	5	1.77(12.17)	$0.79\ (7.09\ )$
	6	2.91 (59.91)	$2.07\ (25.90)$



**Fig. 5** Empirical distribution of the KL divergence  $KL(\Sigma_k, \widehat{\Sigma}_k)$  between each covariance matrix estimated by the EMs and the real covariance matrices  $\Sigma$ .

Class (with very rare exceptions). A more global look at the overall distributions confirms this trend: the error of tmp-EM are more centred on 0 with less spread than EM. This indicates that the tempering allows the EM algorithm to avoid falling into the first local maximum available and consistently find better ones. From the Machine Learning point of view, we highlighted that with our GMM parameters and n = 500 observations, it was able to better identify the different centroids, despite their ambiguity than the regular EM procedure. Table 3 presents a comparative synthesis of the results of EM and tmp-EM.

#### 3.2 Experiment 2: 3 clusters

In this section, we will assess the capacity of tmp-EM to escape from sub-optimal local maxima near the initialisation. The experimental protocol is the same as in the main paper. Let us recall it here. We confront the algorithm to situations where the true classes have increasingly more ambiguous positions, combined with initialisations designed to be hard to escape from. Even though we still follow the log-likelihood as a critical metric, for illustrative purposes we put more emphasis in this section on visualising whether the clusters were properly identify and following the paths in the 2D space of the estimated centroids towards their final values during the EM procedures.

The setup is the following: we have three clusters of similar shape and same weight. One is isolated and easily identifiable. The other two are next to one another, in a more ambiguous configuration. Figure 6 represents the three, gradually more ambiguous configurations.

We use two different initialisation types to reveal the behaviours of the two EMs. The first - which we call "barycenter" - puts all three initial centroids at the centre of mass of all the observed data points. However, none of the EM procedures would move from this initial state if the three GMM centroids were at the exact same position, hence we actually apply a tiny perturbation to make them all slightly distinct. The blue crosses on Figure 7 represent a typical barycenter initialisation. With this initialisation method, we assess whether the EM procedures are able to correctly estimate the positions of the three clusters, despite the ambiguity, when starting from a fairly neutral position, providing neither direction nor misdirection. On the other hand, the second initialisation type - which we call "2v1" - is voluntarily misguiding the algorithm by positioning two centroids on the isolated right cluster and only one centroid on the side of the two ambiguous left clusters. The blue crosses on Figure 8 represent a typical 2v1 initialisation. This initialisation is intended to assess whether the methods are able to escape the potential well in which they start and make theirs centroids traverse the empty space between the left and right clusters to reach their rightful position. For each of the three parameter families represented on Figure 6, 1000 datasets with 500 observations each are simulated, and the two EMs are ran with both the barycenter and the 2v1 initialisation. In the case of tmp-EM, the oscillating temperature profile is used with parameters  $T_0 = 5$ , r = 2, a = 0.6, b = 20 for the barycenter initialisation, and  $T_0 = 100, r = 1.5, a = 0.02, b = 20$  for the 2v1 initialisation. Although in the case of 2v1, the oscillations are not critical, and the simple temperature profile with  $T_0 = 100$  and r = 1.5 works as well. We have two different sets of tempering hyper-parameters values, one for each of the two very different initialisation types. However, these values then remain the same for the three different parameter families and for every data generation within them. Underlining that the method is not excessively sensitive to the tempering parameters. The experiment with 6 clusters in Section 3.1, already demonstrated that the same hyper parameters could be kept over different initialisation (and different data generations as well) when they were made in a non-adversarial way, by drawing random initial centroids uniformly among the data points.



Fig. 6 500 sample points from a Mixture of Gaussians with 3 classes. The true centroid of each Gaussian are depicted by black crosses, and their true covariance matrices are represented by the confidence ellipses of level 0.8, 0.99 and 0.999 around the centre. There are three different versions of the true parameters. From left to right: the true  $\mu_k$  of the two left clusters ( $mu_1$  and  $mu_2$ ) are getting closer while everything else stays identical.

#### 3.2.1 Illustrative

First we illustrate on unique examples how tmp-EM is able to avoid falling for the tricky initialisations we set up.

As previously stated, the focus will be less on the likelihood optimisation for these illustrative examples. Indeed, they are meant to demonstrate that tmp-EM is able to cross the gaps and put the clusters in the right place even with the disadvantageous initialisation. The more relevant metric to assess success in this task is the error on  $\mu$  (and in a lesser way, the error on  $\Sigma$ ). One reason why the likelihood looses its ability to discriminate between failure and success in escaping the traps set by the initialisations is that there may not be a big likelihood gap between being completely wrong and mostly right. For instance placing two centroids (one of which is linked to an empty class) on the isolated left cluster and putting only one where the two ambiguously close clusters are could have a decent likelihood while being blatantly wrong.

On Figure 7, we represent the results of each EM after convergence for every of the three parameter set, when the start at the barycenter of all data points (blue crosses). The estimated means and covariance matrices of the GMM are represented by orange crosses and confidence ellipses respectively. In those examples, tmp-EM correctly identified the real clusters whereas EM put two centroids on the right, where only the isolated cluster stands, and only one on the left, where the two ambiguous clusters are. Figure 8 shows similar results, with the same conventions in the case of the "2v1" initialisation.

These different outcomes are exactly what one would expect: unlike the classical EM, tmp-EM is by design supposed to avoid the local minima close to the initialisation by taking a more exploratory stance during its first steps. To demonstrate that point, we detail in Figure 9 to 12 the paths taken by the estimated centroids by tmp-EM in those simulations. The paths of the regular EM are straightforward convergences towards their final positions, and are not represented in these supplementary materials. Figure 9 represents the paths of the three cluster centroids during the iterations of tmp-EM. The parameter family is the least ambiguous (the two left cluster are well separated) with the "barycenter" initialisation. On Figure 10, the initialisation is "2v1" instead. The two following Figures, 11 and 12, also features the initialisations "barycenter" and "2v1" respectively, but with the most ambiguous parameter set, where the two left clusters are very close to one another.

These graphs are made of several rows of figures, each row representing a step in the EM procedure. In order to make the Figures informative, the number of steps between each row is not fixed, instead the most interesting steps are represented. Convergence is always achieved within 20 to 50 steps, so there are never big differences between the step gaps anyway. The first row is always the initial stage without any EM step, and the last one is the stage after convergence. Each of the three columns corresponds to one of the three centroids estimated by the EM procedure and represents its evolution in the 2D space, from initialisation to convergence. The corresponding estimated covariance matrix is represented by confidence ellipses. For each of the centroids, the observed data points are coloured accordingly to their (un-tempered) posterior probability of belonging to the associated class at this stage of the the algorithm. Plain blue being a low probability while bright green is a high probability.

We make the following observations on the steps taken by tmp-EM: with a "barycenter" initialisation (Figure 9 and 11), the three centroids gradually converge towards their final position (which correspond to true class centres in these cases) without too much hesitation. We also note that, since the three initial points are slightly distinct, there appears to be preferences at the very beginning, with each class having different high probability points right at the initialisation stage. However those preferences are not respected after a couple EM step, we generally see the centroids directing themselves towards different points than their initial favoured ones. This can be attributed to the tempering reshuffling the positions and preferences at the beginning. The "2v1" initialisation illustrates this phenomenon more clearly and in doing so, showcases the true power of the tempering. The very first steps after this very adversarial initialisation are not very remarkable: the single centroid on the left solidifies its position at the centre of the two ambiguous clusters, while the two centroids on the right try to share the single cluster they started in. However, very quickly this status quo is shattered and every estimated centroid jumps to a completely different position. On both Figure 10 and 12 we see the positions being completely reversed with the lonely centroid moving from the two left clusters to the isolated right one whereas the two close centroids make the inverse trip to reach the two clusters on the left. This jump is an indication that the tempering flattened the likelihood enough to allow each centroid to escape their potential wells. Effectively redoing the initialisation and allowing itself to start from more favourable positions. This behaviour is unattainable with the classical EM.

### 3.2.2 Quantitative

The quantitative analysis can be found in the main paper.

#### 3.3 Experiment on real data: Wine recognition dataset

To further validate tmp-EM, we compare it once more to the unmodified EM, this time on real observations from the scikit learn [6] classification data base "Wine" [4]. This dataset contains p = 13 chemical measurements of n = 178 wines each belonging to one of K = 3 families. Despite being in high dimension, this dataset is known as not very challenging (the classes are separable) and useful for testing new methods. We expect the unmodified EM to perform quite well already. For tmp-EM, we use the simple decreasing temperature profile, with no oscillations, the tempering parameters are  $T_0 = 100$ , r = 4. Table 4 shows the result of 500 runs of the EMs from different random initial points. We focus on the likelihood and the error on  $\mu_k$ , the other relevant metrics, not presented here, show the same tendencies. We observe, as usual, that tmp-EM reaches in average a lower negative log-likelihood with lower variance. The class centres are also better estimated. As expected, the errors made by the EM are already fairly small, however tmp-EM manages to go further and lower the errors on each class by approximately 17%, 18% and 11% respectively.

The results demonstrate that tmp-EM can improve the EM result on real data. Since this is an easy dataset, the difference is not as drastic as in the hard synthetic cases we ran the EMs by. Still, there was room to improve the EM results, and tmp-EM found those better solutions.

#### 4 Experiments on tmp-EM with Independent Factor Analysis

In this section, we present another application of the tmp-EM with Gaussian Mixture Models, but this time as part of a more complex model. The Indepen-



Fig. 7 Typical final positioning of the centroids by EM (left column) and tmp-EM (right column) when the initialisation is made at the barycenter of all data points (blue crosses). The three rows represent the three gradually more ambiguous parameter sets. Each figure represents the positions of the estimated centroids after convergence of the EM algorithms (orange cross), with their estimated covariance matrices (orange confidence ellipses). In each simulation, 500 sample points were drawn from the real GMM (small green crosses). In those example, tmp-EM managed to correctly identify the position of the three real centroids.

**Table 4** Average and (standard deviation) of the EM and tmp-EM results over 500 random initialisation on the Wine recognition dataset. The classes on this dataset are easily identifiable hence the errors are low. Yet tmp-EM still improved upon the solutions of EM

metric	cl.	EM	tmp-EM
$-lnp_{\hat{ heta}}$		2923 (77)	2905 (71)
$\frac{\ \hat{\mu}_k - \mu_k\ ^2}{\ \mu_k\ ^2}$	$     \begin{array}{c}       1 \\       2 \\       3     \end{array}   $	$\begin{array}{c} 0.017 \ (0.030) \\ 0.026 \ (0.034) \\ 0.089 \ (0.165) \end{array}$	$\begin{array}{c} 0.014 \ (0.028) \\ 0.021 \ (0.033) \\ 0.079 \ (0.156) \end{array}$



Fig. 8 Typical final positioning of the centroids by EM (left column) and tmp-EM (right column) when the initialisation is made by selecting two points in the isolated cluster and one in the lower ambiguous cluster (blue crosses). The three rows represent the three gradually more ambiguous parameter sets. Each figure represents the positions of the estimated centroids after convergence of the EM algorithms (orange cross), with their estimated covariance matrices (orange confidence ellipses). In each simulation, 500 sample points were drawn from the real GMM (small green crosses). In those examples, although EM kept two centroids on the isolated cluster, tmp-EM managed to correctly identify the position of the three real centroids.

dent Factor Analysis (IFA) model was introduced by [3] as an amalgam of Factor Analysis, Principal Component Analysis and Independent Component Analysis to identify and separate independent sources mixed into a single feature vector. From a practical standpoint, the mixing coefficient of each source is assumed to be drawn from a GMM, hence the EM. After estimation of the GMM parameters, the sources are recovered with an optimal non linear estimator. This is a complex model in which the EM plays a key part, works like [2] and [1] use it to assess new variants of the EM on a very practical application. The model is described as follows:

$$\forall i = 1, ..., L', \quad y_i = \sum_{j=1}^{L} H_{ij} x_j + u_i.$$



Fig. 9 Paths of the centroids for tmp-EM with the "barycenter" initialisation. Parameter set 1 (least ambiguous).



Fig. 10 Paths of the centroids for tmp-EM with the "2v1" initialisation. Parameter set 1 (least ambiguous).



Fig. 11 Paths of the centroids for tmp-EM with the "barycenter" initialisation. Parameter set 3 (most ambiguous).



Fig. 12 Paths of the centroids for tmp-EM with the "2v1" initialisation. Parameter set 3 (most ambiguous).

Where  $y \in R^{L'}$  is one vector of observations,  $H \in \mathbb{R}^{L'L}$  is the fixed matrix of the sources,  $u \in R^{L'}$  the vector of noise, and  $x \in R^L$  the random mixing coefficient. Each component  $x_i$  is assumed to be drawn from its own GMM.

An EM that converges too soon towards a local extremum has every chance to yield sub-optimal estimated sources. We demonstrate in this section that an IFA method with tmp-EM can recover sources closer to the original when they are known, and cleaner, more stable looking sources in general.

#### 4.1 Synthetic IFA

We start with a toy example, where the true sources are two easily distinguishable images. As shown on Figure 14, one is a white square on a black background and the other is a white cross on a similar black background but positioned differently. However, once these two sources are mixed and noised, it becomes much harder to identify them with the naked eye - as illustrated by Figure 14 - and a quantitative method is required to properly separate them. To separate the sources,



Fig. 13 The two real sources of a synthetic source mixing model. They are images of size  $20 \times 20$  made of a black background with a white symbol localised either on the bottom left or top right corner.

the identification model assumes that the coefficients used to mix the two sources are drawn from mixtures of gaussian. The outputs were voluntarily generated in a different way to show the generalisation capabilities of the mixture of gaussian assumption. We run an EM and a tmp-EM algorithm to estimate the parameter of those mixtures, recovering in the process an estimation of the mixing matrix H. Figure 15 illustrates the sources typically estimated by each of the two procedure. Although there is noise, tmp-EM essentially identified and corrected the sources correctly. Whereas EM did not manage to completely turn off the square symbol in the estimated sources supposedly dedicated to the cross. Figure 16 displays the quantitative results of several runs over different simulated datasets. It represents the empirical distribution of  $l_2$  errors made on the estimation of the source matrix H by the two EMs. As illustrated by the table in Figure 16, the solutions of tmp-EM have lower mean and median.



Fig. 14 6 typical observation obtained with the source mixing model. With the noise, the sources are harder to identify.



Fig. 15 Estimated sources by EM (up) and tmp-EM (down). The two real sources were correctly identify by tmp-EM, but EM did not fully separate the cross and the square.

## 4.2 ZIP code

We apply this IFA algorithm to the ZIP code dataset from Elements of Statistical learning. This dataset contains handwritten digits between 0 and 9. In this study, we keep only the digits 0,3, 8 (all three being ambiguously similar) and 7 (very different from the three others). We make all classes even by removing half of the 0



Fig. 16 Empirical distribution of the  $l_2$  error on the source matrix H made by EM and tmp-EM. With tmp-EM, we shift the distribution towards the lower errors, with smaller average and median. The numeric values of the quantiles and other statistics can be found in the table, the better ones being in **bold**.

which are originally more numerous. When applying Independent Factor Analysis to such data, one hopes that the distinct digits will be identified as the separable sources making up the signal. We run the IFA model with a Mixture of Gaussians model with a regular and a tempered EM. In the mixing model used, each mixture is composed of two classes. The tempering was made with the oscillating profile, with hyper-parameters:  $T_0 = 50, b = 20, r = 3, a = 0.02$ .

Figure 17 displays the estimated sources by the IFA procedure with either EM or tmp-EM at their core. EM did not really identify an "8" source. Instead, its "3" is a bit ambiguously close to and "8", and the rightmost source in Figure 17 seems like an amalgamation of the four digits. Moreover, the source "7" estimated by EM is actually a mix between a "7" and a "0". On the other hand, the sources estimated by tmp-EM each correspond clearly to a different digit. There is an "8", the "7" is not fused with a "0", the "3" is sharper and more distinct from an "8" then the corresponding EM source, and even the "0" is more symmetrical with tmp-EM than with EM. Tempering the EM within the IFA algorithm allowed for a cleaner separation of the sources. One can infer that tmp-EM was able to identify and reach a better local maximum of the loss function.



Fig. 17 Estimated sources by EM (up) and tmp-EM (down). The "8" and the "7" in particular were much better identified by tmp-EM. Moreover, with tempering the "0" has a more symmetrical shape and the "3" is sharper and less ambiguous.

# References

- 1. Allassonnière, S., Chevallier, J.: A New Class of EM Algorithms. Escaping Local Minima and Handling Intractable Sampling (2019). URL https://hal. archives-ouvertes.fr/hal-02044722. Working paper or preprint
- Allassonniere, S., Younes, L., et al.: A stochastic algorithm for probabilistic independent component analysis. The Annals of Applied Statistics 6(1), 125– 160 (2012)
- 3. Attias, H.: Independent factor analysis. Neural computation **11**(4), 803–851 (1999)
- 4. Dua, D., Graff, C.: UCI machine learning repository (2017). URL http://archive.ics.uci.edu/ml
- Fort, G., Moulines, E., et al.: Convergence of the monte carlo expectation maximization for curved exponential families. The Annals of Statistics 31(4), 1220– 1259 (2003)
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. Journal of Machine Learning Research 12, 2825–2830 (2011)