

The f -Divergence Expectation Iteration Scheme

Kamélia Daudel

Télécom Paris
kamelia.daudel@telecom-paris.fr

November 22, 2019

Joint work with Randal Douc, François Portier and François Roueff



Introduction

- Bayesian statistics : compute / sample from the **posterior density** of the latent variables y given the data \mathcal{D}

$$p(y|\mathcal{D}) = \frac{p(\mathcal{D}, y)}{p(\mathcal{D})} .$$

- Problem : the marginal likelihood $p(\mathcal{D})$ is **untractable**.
→ **Variational Inference** methods :

Goal

Approximate the posterior density $p(\cdot|\mathcal{D})$ by a variational density q_θ , where $\theta \in \mathcal{T}$:

$$\theta^* = \operatorname{arginf}_{\theta \in \mathcal{T}} \mathcal{D}(q_\theta, p(\cdot|\mathcal{D})) ,$$

where \mathcal{D} is a divergence.

↪ Particular case of **density approximation**.

Introduction

- Bayesian statistics : compute / sample from the **posterior density** of the latent variables y given the data \mathcal{D}

$$p(y|\mathcal{D}) = \frac{p(\mathcal{D}, y)}{p(\mathcal{D})} .$$

- Problem : the marginal likelihood $p(\mathcal{D})$ is **untractable**.

→ **Variational Inference** methods :

Goal

Approximate the posterior density $p(\cdot|\mathcal{D})$ by a variational density q_θ , where $\theta \in \mathcal{T}$:

$$\theta^* = \operatorname{arginf}_{\theta \in \mathcal{T}} \mathcal{D}(q_\theta, p(\cdot|\mathcal{D})) ,$$

where \mathcal{D} is a divergence.

↪ Particular case of **density approximation**.

Introduction

- Bayesian statistics : compute / sample from the **posterior density** of the latent variables y given the data \mathcal{D}

$$p(y|\mathcal{D}) = \frac{p(\mathcal{D}, y)}{p(\mathcal{D})} .$$

- Problem : the marginal likelihood $p(\mathcal{D})$ is **untractable**.
→ **Variational Inference** methods :

Goal

Approximate the posterior density $p(\cdot|\mathcal{D})$ by a variational density q_θ , where $\theta \in \mathsf{T}$:

$$\theta^* = \operatorname{arginf}_{\theta \in \mathsf{T}} \mathcal{D}(q_\theta, p(\cdot|\mathcal{D})) ,$$

where \mathcal{D} is a divergence.

↪ Particular case of **density approximation**.

Introduction

- Bayesian statistics : compute / sample from the **posterior density** of the latent variables y given the data \mathcal{D}

$$p(y|\mathcal{D}) = \frac{p(\mathcal{D}, y)}{p(\mathcal{D})} .$$

- Problem : the marginal likelihood $p(\mathcal{D})$ is **untractable**.
→ **Variational Inference** methods :

Goal

Approximate the posterior density $p(\cdot|\mathcal{D})$ by a variational density q_θ , where $\theta \in \mathsf{T}$:

$$\theta^* = \operatorname{arginf}_{\theta \in \mathsf{T}} \mathcal{D}(q_\theta, p(\cdot|\mathcal{D})) ,$$

where \mathcal{D} is a divergence.

↪ Particular case of **density approximation**.

Our approach

Usually in Variational Inference : approximating family

$$\{y \mapsto q_{\theta}(y) : \theta \in \mathcal{T}\} .$$

Let us now consider a broader approximating family

$$\left\{ y \mapsto \int_{\mathcal{T}} \mu(d\theta) q_{\theta}(y) : \mu \in \mathcal{M} \right\} ,$$

\mathcal{M} : subset of $\mathcal{M}_1(\mathcal{T})$, the set of probability measures on $(\mathcal{T}, \mathcal{T})$.

Question : Can we define an **iterative scheme** which diminishes a given objective function at each step ?

→ **Yes** : The f -El(ϕ) algorithm !

Our approach

Usually in Variational Inference : approximating family

$$\{y \mapsto q_{\theta}(y) : \theta \in \mathcal{T}\} .$$

Let us now consider a broader approximating family

$$\left\{ y \mapsto \int_{\mathcal{T}} \mu(d\theta) q_{\theta}(y) : \mu \in \mathcal{M} \right\} ,$$

\mathcal{M} : subset of $\mathcal{M}_1(\mathcal{T})$, the set of probability measures on $(\mathcal{T}, \mathcal{T})$.

Question : Can we define an **iterative scheme** which diminishes a given objective function at each step ?

→ **Yes** : The f -El(ϕ) algorithm !

Our approach

Usually in Variational Inference : approximating family

$$\{y \mapsto q_{\theta}(y) : \theta \in \mathcal{T}\} .$$

Let us now consider a broader approximating family

$$\left\{ y \mapsto \int_{\mathcal{T}} \mu(d\theta) q_{\theta}(y) : \mu \in \mathcal{M} \right\} ,$$

\mathcal{M} : subset of $\mathcal{M}_1(\mathcal{T})$, the set of probability measures on $(\mathcal{T}, \mathcal{T})$.

Question : Can we define an **iterative scheme** which diminishes a given objective function at each step ?

→ **Yes :** The f -El(ϕ) algorithm !

Our approach

Usually in Variational Inference : approximating family

$$\{y \mapsto q_{\theta}(y) : \theta \in \mathcal{T}\} .$$

Let us now consider a broader approximating family

$$\left\{ y \mapsto \int_{\mathcal{T}} \mu(d\theta) q_{\theta}(y) : \mu \in \mathcal{M} \right\} ,$$

\mathcal{M} : subset of $\mathcal{M}_1(\mathcal{T})$, the set of probability measures on $(\mathcal{T}, \mathcal{T})$.

Question : Can we define an **iterative scheme** which diminishes a given objective function at each step ?

→ **Yes :** The f -El(ϕ) algorithm !

Outline

- 1 Optimisation problem
- 2 The f -Expectation Iteration algorithm $f\text{-EI}(\phi)$
- 3 Application to density approximation
- 4 Conclusion

Outline

- 1 Optimisation problem
- 2 The f -Expectation Iteration algorithm $f\text{-EI}(\phi)$
- 3 Application to density approximation
- 4 Conclusion

Objective function : the f -divergence

- (Y, \mathcal{Y}, ν) : measured space, where ν is a σ -finite measure on (Y, \mathcal{Y})
- f : **convex** function over $(0, \infty)$ that satisfies $f(1) = 0$
- \mathbb{P}_1 and \mathbb{P}_2 : two probability measures on (Y, \mathcal{Y}) such that $\mathbb{P}_1 \preceq \nu$, $\mathbb{P}_2 \preceq \nu$ with $p_1 = \frac{d\mathbb{P}_1}{d\nu}$, $p_2 = \frac{d\mathbb{P}_2}{d\nu}$

Definition 1 : f -divergence between \mathbb{P}_1 and \mathbb{P}_2

$$D_f(\mathbb{P}_1 || \mathbb{P}_2) = \int_Y f\left(\frac{p_1(y)}{p_2(y)}\right) p_2(y) \nu(dy)$$

Objective function : the f -divergence

Definition 1 : f -divergence between \mathbb{P}_1 and \mathbb{P}_2

$$D_f(\mathbb{P}_1 || \mathbb{P}_2) = \int_Y f\left(\frac{p_1(y)}{p_2(y)}\right) p_2(y) \nu(dy)$$

→ a **flexible** family of divergences

$f(u)$	Corresponding divergence
$u \log(u)$	$D_{KL}(\mathbb{P}_1 \mathbb{P}_2) = \int_Y \log\left(\frac{p_1(y)}{p_2(y)}\right) p_1(y) \nu(dy)$
$-\log(u)$	$D_{rKL}(\mathbb{P}_1 \mathbb{P}_2) = \int_Y -\log\left(\frac{p_1(y)}{p_2(y)}\right) p_2(y) \nu(dy)$
$\frac{1}{\alpha(\alpha-1)}[u^\alpha - 1]$	$D_A^{(\alpha)}(\mathbb{P}_1 \mathbb{P}_2) = \frac{1}{\alpha(\alpha-1)} \left[\int_Y \left(\frac{p_1(y)}{p_2(y)}\right)^\alpha p_2(y) \nu(dy) - 1 \right]$

Table 1: Special cases in the f -divergence family

Optimisation problem

- $(\mathsf{T}, \mathcal{T})$: measurable space
- p : measurable positive function on $(\mathsf{Y}, \mathcal{Y})$
- $Q : (\theta, A) \mapsto \int_A q(\theta, y) \nu(dy)$: Markov transition kernel on $\mathsf{T} \times \mathcal{Y}$ with kernel density q

$$\forall \mu \in \mathsf{M}_1(\mathsf{T}), \forall y \in \mathsf{Y}, \mu q(y) = \int_{\mathsf{T}} \mu(d\theta) q(\theta, y)$$

General optimisation problem

$$\operatorname{arginf}_{\mu \in \mathsf{M}} \Psi^{(f)}(\mu)$$

$$\text{where for all } \mu \in \mathsf{M}_1(\mathsf{T}), \Psi^{(f)}(\mu) = \int_{\mathsf{Y}} f\left(\frac{\mu q(y)}{p(y)}\right) p(y) \nu(dy).$$

→ The mapping $\mu \mapsto \Psi^{(f)}(\mu)$ is convex.

Optimisation problem

- $(\mathsf{T}, \mathcal{T})$: measurable space
- p : measurable positive function on $(\mathsf{Y}, \mathcal{Y})$
- $Q : (\theta, A) \mapsto \int_A q(\theta, y) \nu(dy)$: Markov transition kernel on $\mathsf{T} \times \mathcal{Y}$ with kernel density q

$$\forall \mu \in \mathsf{M}_1(\mathsf{T}), \forall y \in \mathsf{Y}, \mu q(y) = \int_{\mathsf{T}} \mu(d\theta) q(\theta, y)$$

General optimisation problem

$$\operatorname{arginf}_{\mu \in \mathsf{M}} \Psi^{(f)}(\mu)$$

$$\text{where for all } \mu \in \mathsf{M}_1(\mathsf{T}), \Psi^{(f)}(\mu) = \int_{\mathsf{Y}} f\left(\frac{\mu q(y)}{p(y)}\right) p(y) \nu(dy).$$

→ The mapping $\mu \mapsto \Psi^{(f)}(\mu)$ is **convex**.

Outline

- 1 Optimisation problem
- 2 The f -Expectation Iteration algorithm $f\text{-EI}(\phi)$
- 3 Application to density approximation
- 4 Conclusion

The f -Expectation Iteration algorithm f -EI(ϕ)

Let $\phi \in \mathbb{R}^*$, $\mu \in \mathcal{M}_1(\mathcal{T})$ such that $\Psi^{(f)}(\mu) < \infty$. We define the sequence of probability measures $(\mu_n)_{n \in \mathbb{N}}$ iteratively by

$$\begin{cases} \mu_0 = \mu, \\ \mu_{n+1} = \mathcal{I}^\phi(\mu_n), \end{cases} \quad n \in \mathbb{N}. \quad (1)$$

Algorithm 1: *Exact f -EI(ϕ) transition*

1. Expectation step : $b_\mu(\theta) = \int_{\mathcal{Y}} q(\theta, y) f' \left(\frac{\mu q(y)}{p(y)} \right) \nu(dy)$
 2. Iteration step : $\mathcal{I}^\phi(\mu)(d\theta) = \frac{\mu(d\theta) \cdot |b_\mu(\theta)|^\phi}{\mu(|b_\mu|^\phi)}$
-

When is the f -EI(ϕ) algorithm well-defined ?

$$b_{\mu}(\theta) = \int_Y q(\theta, y) f' \left(\frac{\mu q(y)}{p(y)} \right) \nu(dy)$$
$$\mathcal{I}^{\phi}(\mu)(d\theta) = \frac{\mu(d\theta) \cdot |b_{\mu}(\theta)|^{\phi}}{\mu(|b_{\mu}|^{\phi})}$$

(A1) For all $(\theta, y) \in \mathcal{T} \times Y$, $q(\theta, y) > 0$, $p(y) > 0$ and $\int_Y p(y) \nu(dy) < \infty$.

(A2) $f : (0, \infty) \rightarrow \mathbb{R}$ is monotonous, strictly convex and continuously differentiable, and $f(1) = 0$.

→ Under (A1) and (A2), b_{μ} is well-defined and $|b_{\mu}| \in (0, \infty]$.

→ The iteration $\mu \mapsto \mathcal{I}^{\phi}(\mu)$ is well-defined if moreover we have

$$0 < \mu(|b_{\mu}|^{\phi}) < \infty . \quad (2)$$

When is the f -EI(ϕ) algorithm well-defined ?

$$b_{\mu}(\theta) = \int_Y q(\theta, y) f' \left(\frac{\mu q(y)}{p(y)} \right) \nu(dy)$$
$$\mathcal{I}^{\phi}(\mu)(d\theta) = \frac{\mu(d\theta) \cdot |b_{\mu}(\theta)|^{\phi}}{\mu(|b_{\mu}|^{\phi})}$$

(A1) For all $(\theta, y) \in \mathsf{T} \times \mathsf{Y}$, $q(\theta, y) > 0$, $p(y) > 0$ and $\int_Y p(y) \nu(dy) < \infty$.

(A2) $f : (0, \infty) \rightarrow \mathbb{R}$ is monotonous, strictly convex and continuously differentiable, and $f(1) = 0$.

→ Under (A1) and (A2), b_{μ} is well-defined and $|b_{\mu}| \in (0, \infty]$.

→ The iteration $\mu \mapsto \mathcal{I}^{\phi}(\mu)$ is well-defined if moreover we have

$$0 < \mu(|b_{\mu}|^{\phi}) < \infty . \tag{2}$$

When is the f -EI(ϕ) algorithm well-defined ?

$$b_{\mu}(\theta) = \int_Y q(\theta, y) f' \left(\frac{\mu q(y)}{p(y)} \right) \nu(dy)$$
$$\mathcal{I}^{\phi}(\mu)(d\theta) = \frac{\mu(d\theta) \cdot |b_{\mu}(\theta)|^{\phi}}{\mu(|b_{\mu}|^{\phi})}$$

(A1) For all $(\theta, y) \in \mathsf{T} \times \mathsf{Y}$, $q(\theta, y) > 0$, $p(y) > 0$ and $\int_Y p(y) \nu(dy) < \infty$.

(A2) $f : (0, \infty) \rightarrow \mathbb{R}$ is monotonous, strictly convex and continuously differentiable, and $f(1) = 0$.

→ Under (A1) and (A2), b_{μ} is well-defined and $|b_{\mu}| \in (0, \infty]$.

→ The iteration $\mu \mapsto \mathcal{I}^{\phi}(\mu)$ is well-defined if moreover we have

$$0 < \mu(|b_{\mu}|^{\phi}) < \infty . \tag{2}$$

Monotonicity

Divergence considered		Corresponding range
<i>Reverse KL</i> $f(u) = -\log(u)$		$\phi \in (0, 1]$
α -divergence $f(u) = \frac{1}{\alpha(\alpha-1)}(u^\alpha - 1)$	$\alpha \in (-\infty, -1]$	$\phi \in (0, -1/\alpha]$
	$\alpha \in (-1, 1) \setminus \{0\}$	$\phi \in (0, 1]$
	$\alpha \in (1, \infty)$	$\phi \in (1/(1-\alpha), 0)$

Table 2 : Allowed (f, ϕ) in the f -EI(ϕ) algorithm

Theorem 1

Assume that p and q are as in (A1). Let (f, ϕ) belong to any of the cases in Table 2.

Then (A2) holds. Moreover, let $\mu \in M_1(T)$ be such that $\Psi^{(f)}(\mu) < \infty$. Then the sequence $(\mu_n)_{n \in \mathbb{N}}$ defined by (1) is well-defined and the sequence $(\Psi^{(f)}(\mu_n))_{n \in \mathbb{N}}$ is **non-increasing**.

Limiting behavior

(A3) T is a compact metric space, $\theta \mapsto q(\theta, y)$ is continuous for all $y \in Y$, $\Psi^{(f)}$ and b_μ are uniformly bounded w.r.t μ and θ .

Theorem 2

Assume (A1), (A2) and (A3). Further assume that there exists $\mu, \bar{\mu} \in M_1(T)$ such that the (well-defined) sequence $(\mu_n)_{n \in \mathbb{N}}$ defined by (1) weakly converges to $\bar{\mu}$ as $n \rightarrow \infty$. Then

- ① $\bar{\mu}$ is a fixed point of \mathcal{I}^ϕ ,
- ② $\Psi^{(f)}(\bar{\mu}) = \inf_{\zeta \in M_{1,\mu}(T)} \Psi^{(f)}(\zeta)$,

for f non-increasing and $\phi > 0$ or f non-decreasing and $\phi < 0$.

$M_{1,\mu}(T)$: set of probability measures dominated by μ

→ Compatible with Theorem 1 for (f, ϕ)

Limiting behavior

(A3) T is a compact metric space, $\theta \mapsto q(\theta, y)$ is continuous for all $y \in Y$, $\Psi^{(f)}$ and b_μ are uniformly bounded w.r.t μ and θ .

Theorem 2

Assume (A1), (A2) and (A3). Further assume that there exists $\mu, \bar{\mu} \in M_1(T)$ such that the (well-defined) sequence $(\mu_n)_{n \in \mathbb{N}}$ defined by (1) weakly converges to $\bar{\mu}$ as $n \rightarrow \infty$. Then

- ① $\bar{\mu}$ is a fixed point of \mathcal{I}^ϕ ,
- ② $\Psi^{(f)}(\bar{\mu}) = \inf_{\zeta \in M_{1,\mu}(T)} \Psi^{(f)}(\zeta)$,

for f non-increasing and $\phi > 0$ or f non-decreasing and $\phi < 0$.

$M_{1,\mu}(T)$: set of probability measures dominated by μ

→ Compatible with Theorem 1 for (f, ϕ)

Approximate f -EI(ϕ)

Algorithm 1 typically involves an **intractable** integral in the Expectation step :

$$b_{\mu}(\theta) = \int_Y q(\theta, y) f' \left(\frac{\mu q(y)}{p(y)} \right) \nu(dy) .$$

→ Approximate f -EI(ϕ)

Algorithm 2: *Approximate f -EI(ϕ) transition*

1. Sampling step : Draw independently $Y_1, \dots, Y_K \sim \mu q$
 2. Expectation step : $b_{\mu, K}(\theta) = \frac{1}{K} \sum_{k=1}^K \frac{q(\theta, Y_k)}{\mu q(Y_k)} f' \left(\frac{\mu q(Y_k)}{p(Y_k)} \right)$
 3. Iteration step : $\mathcal{I}_K^{\phi}(\mu)(d\theta) = \frac{\mu(d\theta) \cdot |b_{\mu, K}(\theta)|^{\phi}}{\mu(|b_{\mu, K}|^{\phi})}$
-

Approximate f -EI(ϕ)

Algorithm 1 typically involves an **intractable** integral in the Expectation step :

$$b_{\mu}(\theta) = \int_{\mathcal{Y}} q(\theta, y) f' \left(\frac{\mu q(y)}{p(y)} \right) \nu(\mathrm{d}y) .$$

→ Approximate f -EI(ϕ)

Algorithm 2: *Approximate f -EI(ϕ) transition*

1. Sampling step : Draw independently $Y_1, \dots, Y_K \sim \mu q$
 2. Expectation step : $b_{\mu, K}(\theta) = \frac{1}{K} \sum_{k=1}^K \frac{q(\theta, Y_k)}{\mu q(Y_k)} f' \left(\frac{\mu q(Y_k)}{p(Y_k)} \right)$
 3. Iteration step : $\mathcal{I}_K^{\phi}(\mu)(\mathrm{d}\theta) = \frac{\mu(\mathrm{d}\theta) \cdot |b_{\mu, K}(\theta)|^{\phi}}{\mu(|b_{\mu, K}|^{\phi})}$
-

Total variation convergence

Let Y_1, Y_2, \dots be i.i.d random variables with common density μq w.r.t ν , defined on the same probability space $(\Omega, \mathcal{F}, \mathbb{P})$.

Proposition 2

Assume (A1) and (A2). Let $\mu \in M_1(T)$, $\phi \in \mathbb{R}^*$ be such that $\mu(|b_\mu|) \vee \mu(|b_\mu|^\phi) < \infty$ and

$$\int_T \mu(d\theta) \mathbb{E}_{\mu q} \left[\left\{ \frac{q(\theta, Y_1)}{\mu q(Y_1)} \left| f' \left(\frac{\mu q(Y_1)}{p(Y_1)} \right) \right| \right\}^\phi \right] < \infty. \quad (3)$$

Then,

$$\lim_{K \rightarrow \infty} \left\| \mathcal{I}_K^\phi(\mu) - \mathcal{I}^\phi(\mu) \right\|_{TV} = 0, \quad \mathbb{P} - \text{a.s.}$$

Sketch of the proof (1)

- Triangular inequality :

$$\begin{aligned} \left| \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_{\mu,K}|^\phi)} - \frac{|b_\mu(\theta)|^\phi}{\mu(|b_\mu|^\phi)} \right| &= \left| \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_{\mu,K}|^\phi)} - \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_\mu|^\phi)} + \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_\mu|^\phi)} - \frac{|b_\mu(\theta)|^\phi}{\mu(|b_\mu|^\phi)} \right| \\ &\leq \left| \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_{\mu,K}|^\phi)} \right| \left| 1 - \frac{\mu(|b_{\mu,K}|^\phi)}{\mu(|b_\mu|^\phi)} \right| + \frac{||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi|}{\mu(|b_\mu|^\phi)} \end{aligned}$$

which implies :

$$\begin{aligned} \|\mathcal{I}_K^\phi(\mu) - \mathcal{I}^\phi(\mu)\|_{TV} &= \mu \left(\left| \frac{|b_{\mu,K}|^\phi}{\mu(|b_{\mu,K}|^\phi)} - \frac{|b_\mu|^\phi}{\mu(|b_\mu|^\phi)} \right| \right) \\ &\leq \left| 1 - \frac{\mu(|b_{\mu,K}|^\phi)}{\mu(|b_\mu|^\phi)} \right| + \frac{\mu(|b_{\mu,K}|^\phi - |b_\mu|^\phi)}{\mu(|b_\mu|^\phi)} \end{aligned}$$

- First term of the r.h.s : $\left| 1 - \frac{\mu(|b_{\mu,K}|^\phi)}{\mu(|b_\mu|^\phi)} \right|$

Lemma $\lim_{K \rightarrow \infty} \mu(|b_{\mu,K}|^\phi) = \mu(|b_\mu|^\phi), \quad \mathbb{P} - \text{a.s.}$

Sketch of the proof (1)

- Triangular inequality :

$$\begin{aligned} \left| \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_{\mu,K}|^\phi)} - \frac{|b_\mu(\theta)|^\phi}{\mu(|b_\mu|^\phi)} \right| &= \left| \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_{\mu,K}|^\phi)} - \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_\mu|^\phi)} + \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_\mu|^\phi)} - \frac{|b_\mu(\theta)|^\phi}{\mu(|b_\mu|^\phi)} \right| \\ &\leq \left| \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_{\mu,K}|^\phi)} \right| \left| 1 - \frac{\mu(|b_{\mu,K}|^\phi)}{\mu(|b_\mu|^\phi)} \right| + \frac{||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi|}{\mu(|b_\mu|^\phi)} \end{aligned}$$

which implies :

$$\begin{aligned} \|\mathcal{I}_K^\phi(\mu) - \mathcal{I}^\phi(\mu)\|_{TV} &= \mu \left(\left| \frac{|b_{\mu,K}|^\phi}{\mu(|b_{\mu,K}|^\phi)} - \frac{|b_\mu|^\phi}{\mu(|b_\mu|^\phi)} \right| \right) \\ &\leq \left| 1 - \frac{\mu(|b_{\mu,K}|^\phi)}{\mu(|b_\mu|^\phi)} \right| + \frac{\mu(|b_{\mu,K}|^\phi - |b_\mu|^\phi)}{\mu(|b_\mu|^\phi)} \end{aligned}$$

- First term of the r.h.s : $\left| 1 - \frac{\mu(|b_{\mu,K}|^\phi)}{\mu(|b_\mu|^\phi)} \right|$

Lemma $\lim_{K \rightarrow \infty} \mu(|b_{\mu,K}|^\phi) = \mu(|b_\mu|^\phi), \quad \mathbb{P} - \text{a.s.}$

Sketch of the proof (1)

- Triangular inequality :

$$\begin{aligned} \left| \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_{\mu,K}|^\phi)} - \frac{|b_\mu(\theta)|^\phi}{\mu(|b_\mu|^\phi)} \right| &= \left| \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_{\mu,K}|^\phi)} - \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_\mu|^\phi)} + \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_\mu|^\phi)} - \frac{|b_\mu(\theta)|^\phi}{\mu(|b_\mu|^\phi)} \right| \\ &\leq \left| \frac{|b_{\mu,K}(\theta)|^\phi}{\mu(|b_{\mu,K}|^\phi)} \right| \left| 1 - \frac{\mu(|b_{\mu,K}|^\phi)}{\mu(|b_\mu|^\phi)} \right| + \frac{||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi|}{\mu(|b_\mu|^\phi)} \end{aligned}$$

which implies :

$$\begin{aligned} \left\| \mathcal{I}_K^\phi(\mu) - \mathcal{I}^\phi(\mu) \right\|_{TV} &= \mu \left(\left| \frac{|b_{\mu,K}|^\phi}{\mu(|b_{\mu,K}|^\phi)} - \frac{|b_\mu|^\phi}{\mu(|b_\mu|^\phi)} \right| \right) \\ &\leq \left| 1 - \frac{\mu(|b_{\mu,K}|^\phi)}{\mu(|b_\mu|^\phi)} \right| + \frac{\mu(|b_{\mu,K}|^\phi - |b_\mu|^\phi)}{\mu(|b_\mu|^\phi)} \end{aligned}$$

- First term of the r.h.s : $\left| 1 - \frac{\mu(|b_{\mu,K}|^\phi)}{\mu(|b_\mu|^\phi)} \right|$

Lemma $\lim_{K \rightarrow \infty} \mu(|b_{\mu,K}|^\phi) = \mu(|b_\mu|^\phi), \quad \mathbb{P} - \text{a.s.}$

Sketch of the proof (2)

- Second term of the r.h.s : $\frac{\mu(|b_{\mu,K}|^\phi - |b_\mu|^\phi)}{\mu(|b_\mu|^\phi)}$

Generalized Dominated Convergence Theorem :

- ① For all $K \in \mathbb{N}^*$ and for μ -almost all $\theta \in T$,

$$a_K(\theta) \leq b_K(\theta) \leq c_K(\theta) ,$$

and the limits $\lim_{K \rightarrow \infty} a_K(\theta)$, $\lim_{K \rightarrow \infty} b_K(\theta)$, $\lim_{K \rightarrow \infty} c_K(\theta)$ exist.

- ②
- $\mu \left| \lim_{K \rightarrow \infty} a_K \right| + \mu \left| \lim_{K \rightarrow \infty} c_K \right| < \infty$
 - $\mu \left(\lim_{K \rightarrow \infty} a_K \right) = \lim_{K \rightarrow \infty} \mu(a_K)$ and $\mu \left(\lim_{K \rightarrow \infty} c_K \right) = \lim_{K \rightarrow \infty} \mu(c_K)$

$$\Rightarrow \mu \left(\lim_{K \rightarrow \infty} b_K \right) = \lim_{K \rightarrow \infty} \mu(b_K)$$

Sketch of the proof (2)

$$b_K(\theta) = ||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi|$$

- Second term of the r.h.s : $\frac{\mu(|b_{\mu,K}|^\phi - |b_\mu|^\phi)}{\mu(|b_\mu|^\phi)}$

Generalized Dominated Convergence Theorem :

- ① For all $K \in \mathbb{N}^*$ and for μ -almost all $\theta \in T$,

$$a_K(\theta) \leq b_K(\theta) \leq c_K(\theta) ,$$

and the limits $\lim_{K \rightarrow \infty} a_K(\theta)$, $\lim_{K \rightarrow \infty} b_K(\theta)$, $\lim_{K \rightarrow \infty} c_K(\theta)$ exist.

- ②
- $\mu|\lim_{K \rightarrow \infty} a_K| + \mu|\lim_{K \rightarrow \infty} c_K| < \infty$
 - $\mu(\lim_{K \rightarrow \infty} a_K) = \lim_{K \rightarrow \infty} \mu(a_K)$ and $\mu(\lim_{K \rightarrow \infty} c_K) = \lim_{K \rightarrow \infty} \mu(c_K)$

$$\Rightarrow \mu(\lim_{K \rightarrow \infty} b_K) = \lim_{K \rightarrow \infty} \mu(b_K)$$

Sketch of the proof (2)

$$b_K(\theta) = ||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi|$$

- Second term of the r.h.s : $\frac{\mu(||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi|)}{\mu(|b_\mu|^\phi)}$

Generalized Dominated Convergence Theorem :

- ① For all $K \in \mathbb{N}^*$, for all $\theta \in \mathbb{T}$,

$$0 \leq ||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi| \leq |b_{\mu,K}(\theta)|^\phi + |b_\mu(\theta)|^\phi$$

and the limits $\lim_{K \rightarrow \infty} a_K(\theta)$, $\lim_{K \rightarrow \infty} b_K(\theta)$, $\lim_{K \rightarrow \infty} c_K(\theta)$ exist.

- ②
 - $\mu|\lim_{K \rightarrow \infty} a_K| + \mu|\lim_{K \rightarrow \infty} c_K| < \infty$
 - $\mu(\lim_{K \rightarrow \infty} a_K) = \lim_{K \rightarrow \infty} \mu(a_K)$ and $\mu(\lim_{K \rightarrow \infty} c_K) = \lim_{K \rightarrow \infty} \mu(c_K)$

$$\Rightarrow \mu(\lim_{K \rightarrow \infty} b_K) = \lim_{K \rightarrow \infty} \mu(b_K)$$

Sketch of the proof (2)

$$b_K(\theta) = ||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi|$$

- Second term of the r.h.s : $\frac{\mu(||b_{\mu,K}|^\phi - |b_\mu|^\phi|)}{\mu(|b_\mu|^\phi)}$

Generalized Dominated Convergence Theorem :

- ① For all $K \in \mathbb{N}^*$, for all $\theta \in \mathbb{T}$,

$$0 \leq ||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi| \leq |b_{\mu,K}(\theta)|^\phi + |b_\mu(\theta)|^\phi$$

$$a_K(\theta)$$

and the limits $\lim_{K \rightarrow \infty} a_K(\theta)$, $\lim_{K \rightarrow \infty} b_K(\theta)$, $\lim_{K \rightarrow \infty} c_K(\theta)$ exist.

$$c_K(\theta)$$

- ②
 - $\mu|\lim_{K \rightarrow \infty} a_K| + \mu|\lim_{K \rightarrow \infty} c_K| < \infty$
 - $\mu(\lim_{K \rightarrow \infty} a_K) = \lim_{K \rightarrow \infty} \mu(a_K)$ and $\mu(\lim_{K \rightarrow \infty} c_K) = \lim_{K \rightarrow \infty} \mu(c_K)$

$$\Rightarrow \mu(\lim_{K \rightarrow \infty} b_K) = \lim_{K \rightarrow \infty} \mu(b_K)$$

Sketch of the proof (2)

$$b_K(\theta) = ||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi|$$

- Second term of the r.h.s : $\frac{\mu(||b_{\mu,K}|^\phi - |b_\mu|^\phi|)}{\mu(|b_\mu|^\phi)}$

Generalized Dominated Convergence Theorem :

- For all $K \in \mathbb{N}^*$, for all $\theta \in \mathbb{T}$,

$$0 \leq ||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi| \leq |b_{\mu,K}(\theta)|^\phi + |b_\mu(\theta)|^\phi$$

$a_K(\theta)$

and the limits $\lim_{K \rightarrow \infty} a_K(\theta)$, $\lim_{K \rightarrow \infty} b_K(\theta)$, $\lim_{K \rightarrow \infty} c_K(\theta)$ exist.

$c_K(\theta)$

- $$\lim_{K \rightarrow \infty} \mu [|b_{\mu,K}|^\phi + |b_\mu|^\phi] = \mu \left[\lim_{K \rightarrow \infty} (|b_{\mu,K}|^\phi + |b_\mu|^\phi) \right] < \infty$$

(since $\mu(|b_\mu|^\phi) < \infty$ and $\lim_{K \rightarrow \infty} \mu(|b_{\mu,K}|^\phi) = \mu(|b_\mu|^\phi)$, \mathbb{P} -a.s.)

$$\Rightarrow \mu(\lim_{K \rightarrow \infty} b_K) = \lim_{K \rightarrow \infty} \mu(b_K)$$

Sketch of the proof (2)

$$b_K(\theta) = ||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi|$$

- Second term of the r.h.s : $\frac{\mu(||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi|)}{\mu(|b_\mu|^\phi)}$

Generalized Dominated Convergence Theorem :

- For all $K \in \mathbb{N}^*$, for all $\theta \in \mathbb{T}$,

$$0 \leq ||b_{\mu,K}(\theta)|^\phi - |b_\mu(\theta)|^\phi| \leq |b_{\mu,K}(\theta)|^\phi + |b_\mu(\theta)|^\phi$$

$$a_K(\theta)$$

and the limits $\lim_{K \rightarrow \infty} a_K(\theta)$, $\lim_{K \rightarrow \infty} b_K(\theta)$, $\lim_{K \rightarrow \infty} c_K(\theta)$ exist.

$$c_K(\theta)$$

- $$\lim_{K \rightarrow \infty} \mu [|b_{\mu,K}|^\phi + |b_\mu|^\phi] = \mu \left[\lim_{K \rightarrow \infty} (|b_{\mu,K}|^\phi + |b_\mu|^\phi) \right] < \infty$$

(since $\mu(|b_\mu|^\phi) < \infty$ and $\lim_{K \rightarrow \infty} \mu(|b_{\mu,K}|^\phi) = \mu(|b_\mu|^\phi)$, \mathbb{P} -a.s.)

$$\Rightarrow \lim_{K \rightarrow \infty} \mu(||b_{\mu,K}|^\phi - |b_\mu|^\phi|) = \mu(\lim_{K \rightarrow \infty} ||b_{\mu,K}|^\phi - |b_\mu|^\phi|) = 0, \quad \mathbb{P} - \text{a.s.}$$

Outline

- 1 Optimisation problem
- 2 The f -Expectation Iteration algorithm $f\text{-EI}(\phi)$
- 3 Application to density approximation**
- 4 Conclusion

f -EI(ϕ) applied to density approximation

Let \tilde{p} be a probability density function on (Y, \mathcal{Y}) and assume that we only have access to an **unnormalized** version p^* of the density \tilde{p} , that is for all $y \in Y$,

$$\tilde{p}(y) = \frac{p^*(y)}{Z}, \quad (4)$$

where $Z := \int_Y p^*(y) \nu(dy)$ is called the *normalizing constant* or *partition function*.

→ Posterior density approximation : $\tilde{p} = p(\cdot | \mathcal{D})$, $p^* = p(\mathcal{D}, \cdot)$ and $Z = p(\mathcal{D})$.

Reformulation of the optimisation problem

- $\tilde{\mathbb{P}}$: probability measure on (Y, \mathcal{Y}) with density \tilde{p} with respect to ν
- for all $\mu \in M_1(T)$, μQ : probability measure on (Y, \mathcal{Y}) with density μq with respect to ν

Lemma 3

Assume (A1). Then, for both the reverse Kullback-Leibler and the α -divergence, optimising the objective

$$D_f(\mu Q || \tilde{\mathbb{P}})$$

(with respect to μ) is equivalent to optimising the objective

$$\Psi^{(f)}(\mu; p) \text{ with } p = p^*.$$

Particular case of the α -divergence

$$\rightarrow \alpha\text{-bound} : \tilde{q} \mapsto \xi^{(\alpha)}(\tilde{q}) := \left[\int_Y \left(\frac{\tilde{q}(y)}{p^*(y)} \right)^\alpha p^*(y) \nu(dy) \right]^{\frac{1}{1-\alpha}}$$

Then,

$$\Psi^{(f)}(\mu; p) = \frac{1}{\alpha(\alpha - 1)} \left(\xi^{(\alpha)}(\mu q)^{1-\alpha} - Z \right) \quad \text{with } p = p^* .$$

Lemma 4

Assume (A1). Let $\mu \in M_1(\mathcal{T})$. Then, for all $\alpha_+ \in (0, 1) \cup (1, +\infty)$ and all $\alpha_- < 0$, we have

$$\xi^{(\alpha_+)}(\mu q) \leq Z \leq \xi^{(\alpha_-)}(\mu q) . \quad (5)$$

\rightarrow We can observe the convergence / monotonicity and obtain a bound on the normalising constant Z .

Particular case of the α -divergence

$$\rightarrow \text{\textcolor{teal}{\alpha-bound}} : \tilde{q} \mapsto \xi^{(\alpha)}(\tilde{q}) := \left[\int_Y \left(\frac{\tilde{q}(y)}{p^*(y)} \right)^\alpha p^*(y) \nu(dy) \right]^{\frac{1}{1-\alpha}}$$

Then,

$$\Psi^{(f)}(\mu; p) = \frac{1}{\alpha(\alpha - 1)} \left(\xi^{(\alpha)}(\mu q)^{1-\alpha} - Z \right) \quad \text{with } p = p^* .$$

Lemma 4

Assume [\(A1\)](#). Let $\mu \in \mathcal{M}_1(\mathcal{T})$. Then, for all $\alpha_+ \in (0, 1) \cup (1, +\infty)$ and all $\alpha_- < 0$, we have

$$\xi^{(\alpha_+)}(\mu q) \leq Z \leq \xi^{(\alpha_-)}(\mu q) . \tag{5}$$

\rightarrow We can observe the [convergence / monotonicity](#) and obtain a [bound](#) on the normalising constant Z .

Particular case of the α -divergence

$$\rightarrow \text{\textcolor{teal}{\(\alpha\)-bound}} : \tilde{q} \mapsto \xi^{(\alpha)}(\tilde{q}) := \left[\int_Y \left(\frac{\tilde{q}(y)}{p^*(y)} \right)^\alpha p^*(y) \nu(dy) \right]^{\frac{1}{1-\alpha}}$$

Then,

$$\Psi^{(f)}(\mu; p) = \frac{1}{\alpha(\alpha - 1)} \left(\xi^{(\alpha)}(\mu q)^{1-\alpha} - Z \right) \quad \text{with } p = p^* .$$

Lemma 4

Assume [\(A1\)](#). Let $\mu \in M_1(\mathcal{T})$. Then, for all $\alpha_+ \in (0, 1) \cup (1, +\infty)$ and all $\alpha_- < 0$, we have

$$\xi^{(\alpha_+)}(\mu q) \leq Z \leq \xi^{(\alpha_-)}(\mu q) . \tag{5}$$

\rightarrow We can observe the [convergence / monotonicity](#) and obtain a [bound](#) on the normalising constant Z .

μ_0 is a weighted sum of Dirac measures

$$\mathcal{S}_J = \left\{ \boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_J) \in \mathbb{R}^J : \forall j \in \{1, \dots, J\}, \lambda_j \geq 0 \text{ and } \sum_{j=1}^J \lambda_j = 1 \right\}.$$

Let $\theta_1, \dots, \theta_J \in \mathbb{T}$ be fixed and denote

$$\mu_{\boldsymbol{\lambda}} = \sum_{j=1}^J \lambda_j \delta_{\theta_j} \quad \text{with } \boldsymbol{\lambda} \in \mathcal{S}_J.$$

Then, $\mu_n = \underbrace{\mathcal{I}^\phi \circ \dots \circ \mathcal{I}^\phi}_{n \text{ times}}(\mu_{\boldsymbol{\lambda}})$ is of the form $\mu_n = \sum_{j=1}^J \lambda_{j,n} \delta_{\theta_j}$

with

$$\begin{cases} \lambda_0 = \boldsymbol{\lambda} \\ \lambda_{j,n+1} = \lambda_{j,n} \frac{|b_{\mu_n, K}(\theta_j)|^\phi}{\mu_n(|b_{\mu_n, K}|^\phi)} , & n \in \mathbb{N}, j \in \{1, \dots, J\} \end{cases}$$

μ_0 is a weighted sum of Dirac measures

$$S_J = \left\{ \boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_J) \in \mathbb{R}^J : \forall j \in \{1, \dots, J\}, \lambda_j \geq 0 \text{ and } \sum_{j=1}^J \lambda_j = 1 \right\}.$$

Let $\theta_1, \dots, \theta_J \in \mathbb{T}$ be fixed and denote

$$\mu_{\boldsymbol{\lambda}} = \sum_{j=1}^J \lambda_j \delta_{\theta_j} \quad \text{with} \quad \boldsymbol{\lambda} \in S_J.$$

Then, $\mu_n = \underbrace{\mathcal{I}^\phi \circ \dots \circ \mathcal{I}^\phi}_{n \text{ times}}(\mu_{\boldsymbol{\lambda}})$ is of the form $\mu_n = \sum_{j=1}^J \lambda_{j,n} \delta_{\theta_j}$

with

$$\begin{cases} \lambda_0 = \boldsymbol{\lambda} \\ \lambda_{j,n+1} = \lambda_{j,n} \frac{|b_{\mu_n, K}(\theta_j)|^\phi}{\mu_n(|b_{\mu_n, K}|^\phi)} , \quad n \in \mathbb{N}, j \in \{1, \dots, J\} \end{cases}$$

μ_0 is a weighted sum of Dirac measures

$$S_J = \left\{ \boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_J) \in \mathbb{R}^J : \forall j \in \{1, \dots, J\}, \lambda_j \geq 0 \text{ and } \sum_{j=1}^J \lambda_j = 1 \right\}.$$

Let $\theta_1, \dots, \theta_J \in \mathbb{T}$ be fixed and denote

$$\mu_{\boldsymbol{\lambda}} = \sum_{j=1}^J \lambda_j \delta_{\theta_j} \quad \text{with} \quad \boldsymbol{\lambda} \in S_J.$$

Then, $\mu_n = \underbrace{\mathcal{I}^\phi \circ \dots \circ \mathcal{I}^\phi}_{n \text{ times}}(\mu_{\boldsymbol{\lambda}})$ is of the form $\mu_n = \sum_{j=1}^J \lambda_{j,n} \delta_{\theta_j}$ with

$$\begin{cases} \lambda_0 = \boldsymbol{\lambda} \\ \lambda_{j,n+1} = \lambda_{j,n} \frac{|b_{\mu_n, K}(\theta_j)|^\phi}{\mu_n(|b_{\mu_n, K}|^\phi)}, \quad n \in \mathbb{N}, j \in \{1, \dots, J\} \end{cases}$$

Mixing the two

Algorithm 3: Mixture α -Approximate f -EI(ϕ)

Input: p^* : unnormalized version of the density \tilde{p} , Q : Markov transition kernel, K : number of samples, $\Theta_J = \{\theta_1, \dots, \theta_J\} \subset \mathbb{T}$: parameter set.

Output: Optimised weights λ .

Set $\lambda = [\frac{1}{J}, \dots, \frac{1}{J}]$.

while *the α -bound has not converged* **do**

Sampling step : Draw independently K samples Y_1, \dots, Y_K from $\mu_\lambda q$.

Expectation step : Compute $\mathbf{A}_\lambda = (a_j)_{1 \leq j \leq J}$ where

$$a_j = \frac{1}{K} \sum_{k=1}^K q(\theta_j, Y_k) \mu_\lambda q(Y_k)^{\alpha-2} p^*(Y_k)^{1-\alpha}$$

 and deduce $\mathbf{B}_\lambda = (\lambda_j a_j^\phi)_{1 \leq j \leq J}$, $b_\lambda = \sum_{j=1}^J \lambda_j a_j^\phi$ and $c_\lambda = \sum_{j=1}^J \lambda_j a_j$.

Iteration step : Set

$$\xi_K^{(\alpha)}(\mu_\lambda q) \leftarrow c_\lambda^{1/(1-\alpha)}$$

$$\lambda \leftarrow \frac{1}{b_\lambda} \mathbf{B}_\lambda$$

end

Mixing the two

Algorithm 3: Mixture α -Approximate f -EI(ϕ)

Input: p^* : unnormalized version of the density \tilde{p} , Q : Markov transition kernel, K : number of samples, $\Theta_J = \{\theta_1, \dots, \theta_J\} \subset \mathbb{T}$: parameter set.

Output: Optimised weights λ .

Set $\lambda = [\frac{1}{J}, \dots, \frac{1}{J}]$.

while *the α -bound has not converged* **do**

→ Sampling step : Draw independently K samples Y_1, \dots, Y_K from $\mu_\lambda q$.

→ Expectation step : Compute $\mathbf{A}_\lambda = (a_j)_{1 \leq j \leq J}$ where

$$a_j = \frac{1}{K} \sum_{k=1}^K q(\theta_j, Y_k) \mu_\lambda q(Y_k)^{\alpha-2} p^*(Y_k)^{1-\alpha}$$

and deduce $\mathbf{B}_\lambda = (\lambda_j a_j^\phi)_{1 \leq j \leq J}$, $b_\lambda = \sum_{j=1}^J \lambda_j a_j^\phi$ and $c_\lambda = \sum_{j=1}^J \lambda_j a_j$.

→ Iteration step : Set

$$\xi_K^{(\alpha)}(\mu_\lambda q) \leftarrow c_\lambda^{1/(1-\alpha)}$$

$$\lambda \leftarrow \frac{1}{b_\lambda} \mathbf{B}_\lambda$$

end

Mixing the two

Algorithm 3: Mixture α -Approximate f -EI(ϕ)

Input: p^* : unnormalized version of the density \tilde{p} , Q : Markov transition kernel, K : number of samples, $\Theta_J = \{\theta_1, \dots, \theta_J\} \subset \mathbb{T}$: parameter set.

Output: Optimised weights λ .

Set $\lambda = [\frac{1}{J}, \dots, \frac{1}{J}]$.

while *the α -bound has not converged* **do**

→ Sampling step : Draw independently K samples Y_1, \dots, Y_K from $\mu_\lambda q$.

→ Expectation step : Compute $\mathbf{A}_\lambda = (a_j)_{1 \leq j \leq J}$ where

$$a_j = \frac{1}{K} \sum_{k=1}^K q(\theta_j, Y_k) \mu_\lambda q(Y_k)^{\alpha-2} p^*(Y_k)^{1-\alpha}$$

and deduce $\mathbf{B}_\lambda = (\lambda_j a_j^\phi)_{1 \leq j \leq J}$, $b_\lambda = \sum_{j=1}^J \lambda_j a_j^\phi$ and $c_\lambda = \sum_{j=1}^J \lambda_j a_j$.

→ Iteration step : Set

$$\xi_K^{(\alpha)}(\mu_\lambda q) \leftarrow c_\lambda^{1/(1-\alpha)}$$

$$\lambda \leftarrow \frac{1}{b_\lambda} \mathbf{B}_\lambda$$

end

Most of the computing effort

One interesting remark

- Most of the computing effort : compute $(b_{\mu_n, K}(\theta_j))_{1 \leq j \leq J}$ (or equivalently $\mathbf{A}_\lambda = (a_j)_{1 \leq j \leq J}$).
- The **score gradient** of the function

$$\tilde{q} \mapsto \mathcal{L}_A^{(\alpha)}(\tilde{q}) := \int_Y \frac{1}{\alpha(\alpha - 1)} \left(\frac{\tilde{q}(y)}{p^*(y)} \right)^\alpha p^*(y) \nu(dy) ,$$

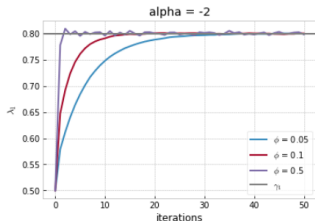
is linked to the quantities approximated in our algorithm

$$\nabla_\lambda \mathcal{L}_A^{(\alpha)}(\mu_\lambda q) = (b_{\mu_\lambda}(\theta_j))_{1 \leq j \leq J} .$$

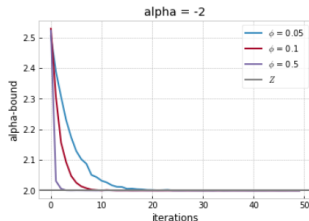
→ similar to computations required in gradient-based methods involving the α -divergence or Renyi's α -divergence.

Numerical experiments : impact of ϕ

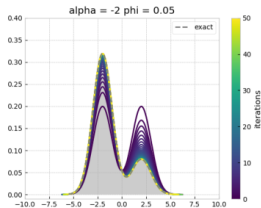
$p^*(y) = Z \times [\gamma_1 \mathcal{N}(y; -s, 1) + \gamma_2 \mathcal{N}(y; s, 1)]$, where $\gamma_1 = 0.8$ $\gamma_2 = 0.2$, $s = 2$ and $Z = 2$



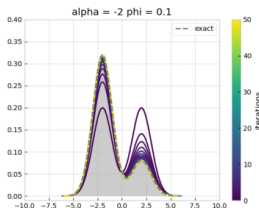
(1)



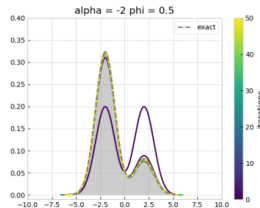
(2)



(3)



(4)



(5)

Towards an adaptive algorithm

- Algorithm 3 leaves $\{\theta_1, \dots, \theta_J\}$ **unchanged** (Exploitation Step)

→ Combine it with an **Exploration step** that modifies the parameter set !

Example : resampling + stochastic perturbation

Towards an adaptive algorithm

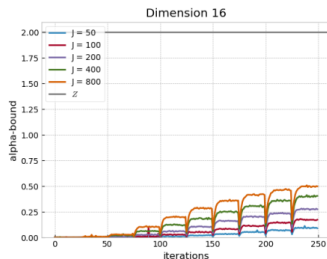
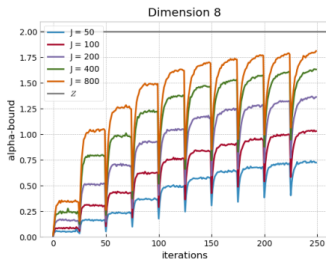
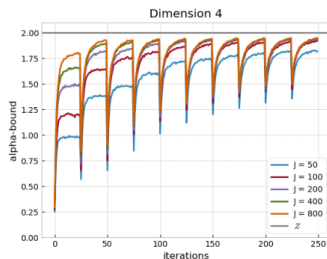
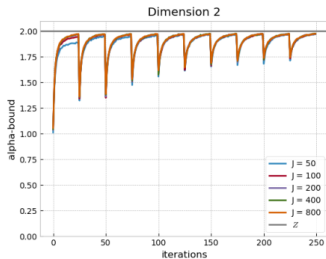
- Algorithm 3 leaves $\{\theta_1, \dots, \theta_J\}$ **unchanged** (Exploitation Step)

→ Combine it with an **Exploration step** that modifies the parameter set !

Example : resampling + stochastic perturbation

Numerical experiments : impact of d and J

$$p^*(y) = Z \times [0.5\mathcal{N}(y; -su_d, I_d) + 0.5\mathcal{N}(y; su_d, I_d)] \text{ with } s = 2 \text{ and } Z = 2$$



Outline

- 1 Optimisation problem
- 2 The f -Expectation Iteration algorithm $f\text{-EI}(\phi)$
- 3 Application to density approximation
- 4 Conclusion**

Conclusion

f -El(ϕ) algorithm : novel iterative scheme that

- performs an update of measures
 - ① Sufficient conditions on (f, ϕ) leading to a systematic decrease
 - ② Convergence to an optimum
 - ③ Approximate version of the algorithm
- can be applied to density approximation
 - ① α -bound: bound on Z , which also measures the convergence
 - ② the computations involved in the Mixture α -Approximate f -El(ϕ) algorithm mostly rely on gradient-based calculations

Conclusion

f -El(ϕ) algorithm : novel iterative scheme that

- performs an update of measures
 - ① Sufficient conditions on (f, ϕ) leading to a systematic decrease
 - ② Convergence to an optimum
 - ③ Approximate version of the algorithm
- can be applied to density approximation
 - ① α -bound: bound on Z , which also measures the convergence
 - ② the computations involved in the Mixture α -Approximate f -El(ϕ) algorithm mostly rely on gradient-based calculations

Perspectives

- ϕ constant in the f -El(ϕ) algorithm \Rightarrow decaying learning rate
- Convergence rate of the f -El(ϕ) algorithm
- Large scale learning
- Try other types of Exploration steps
- Variance reduction schemes in the approximation of b_{μ_n}
- ...

References

- [1] Kamélia Daudel, Randal Douc, François Portier, François Roueff (2019). The f -Divergence Expectation Iteration Scheme. arXiv e-prints, page arXiv:1909.12239.
- [2] Randal Douc, Arnaud Guillin, Jean-Michel Marin, and Christian P. Robert (2007). Convergence of adaptive mixtures of importance sampling schemes. arXiv e-prints, page arXiv:0708.0711.
- [3] José Miguel Hernández-Lobato, Yingzhen Li, Mark Rowland, Daniel Hernández-Lobato, Thang Bui, and Richard E. Turner (2015). Black-box α -divergence Minimization. arXiv e-prints, page arXiv:1511.03243.
- [4] David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational Inference: A Review for Statisticians (2016). arXiv e-prints, page arXiv:1601.00670.
- [5] Yingzhen Li and Richard E. Turner (2016). Rényi Divergence Variational Inference. arXiv e-prints, page arXiv:1602.02311.