

Monotonic Alpha-divergence Minimisation for Variational Inference

Kamélia Daudel

University of Oxford
kamelia.daudel@stats.ox.ac.uk

StatML CDT
09/12/2021

Joint work with Randal Douc and François Roueff

Outline

- 1 Introduction
- 2 Monotonic Alpha-Divergence Minimisation
- 3 Numerical Experiments
- 4 Conclusion

Outline

- 1 Introduction
- 2 Monotonic Alpha-Divergence Minimisation
- 3 Numerical Experiments
- 4 Conclusion

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 : for many important models, we can only evaluate $p(y|\mathcal{D})$ **up to the constant** $p(\mathcal{D})$.

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 : for many important models, we can only evaluate $p(y|\mathcal{D})$ **up to the constant** $p(\mathcal{D})$.

Variational Inference in a nutshell

→ Variational Inference : inference is seen as an **optimisation problem**.

- 1 Posit a *simpler* variational family \mathcal{Q} , where $q \in \mathcal{Q}$.
- 2 Fit q to obtain the best approximation to the posterior density

$$\inf_{q \in \mathcal{Q}} D(\mathbb{Q} \parallel \mathbb{P}_{|\mathcal{D}}) ,$$

where D is a measure of dissimilarity between the variational distribution \mathbb{Q} and the posterior distribution $\mathbb{P}_{|\mathcal{D}}$

Typically, D : exclusive Kullback-Leibler (KL) divergence and \mathcal{Q} : parametric family (e.g. Mean-field)

$$\begin{cases} D_{KL}(\mathbb{Q} \parallel \mathbb{P}) = \int_Y \log \left(\frac{q(y)}{p(y)} \right) q(y) \nu(dy) \\ \mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathbb{T}\} \end{cases}$$

Question : How to choose D and \mathcal{Q} ?

- Can we select **alternative/more general** D ?
- Can we design more **expressive** variational families \mathcal{Q} ?

Variational Inference in a nutshell

→ Variational Inference : inference is seen as an **optimisation problem**.

- 1 Posit a *simpler* variational family \mathcal{Q} , where $q \in \mathcal{Q}$.
- 2 Fit q to obtain the best approximation to the posterior density

$$\inf_{q \in \mathcal{Q}} D(\mathbb{Q} \parallel \mathbb{P}_{|\mathcal{D}}) ,$$

where D is a measure of dissimilarity between the variational distribution \mathbb{Q} and the posterior distribution $\mathbb{P}_{|\mathcal{D}}$

Typically, D : exclusive Kullback-Leibler (KL) divergence and \mathcal{Q} : parametric family (e.g. Mean-field)

$$\begin{cases} D_{KL}(\mathbb{Q} \parallel \mathbb{P}) = \int_Y \log \left(\frac{q(y)}{p(y)} \right) q(y) \nu(dy) \\ \mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathbb{T}\} \end{cases}$$

Question : How to choose D and \mathcal{Q} ?

- Can we select **alternative/more general** D ?
- Can we design more **expressive** variational families \mathcal{Q} ?

Variational Inference in a nutshell

→ Variational Inference : inference is seen as an **optimisation problem**.

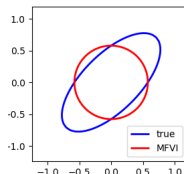
- 1 Posit a *simpler* variational family \mathcal{Q} , where $q \in \mathcal{Q}$.
- 2 Fit q to obtain the best approximation to the posterior density

$$\inf_{q \in \mathcal{Q}} D(\mathbb{Q} || \mathbb{P}_{|\mathcal{D}}) ,$$

where D is a measure of dissimilarity between the variational distribution \mathbb{Q} and the posterior distribution $\mathbb{P}_{|\mathcal{D}}$

Typically, D : exclusive Kullback-Leibler (KL) divergence and \mathcal{Q} : parametric family (e.g. Mean-field)

$$\begin{cases} D_{KL}(\mathbb{Q} || \mathbb{P}) = \int_Y \log \left(\frac{q(y)}{p(y)} \right) q(y) \nu(dy) \\ \mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathbb{T}\} \end{cases}$$



Question : How to choose D and \mathcal{Q} ?

- Can we select **alternative/more general** D ?
- Can we design more **expressive** variational families \mathcal{Q} ?

Variational Inference in a nutshell

→ Variational Inference : inference is seen as an **optimisation problem**.

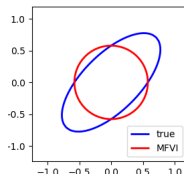
- 1 Posit a *simpler* variational family \mathcal{Q} , where $q \in \mathcal{Q}$.
- 2 Fit q to obtain the best approximation to the posterior density

$$\inf_{q \in \mathcal{Q}} D(\mathbb{Q} || \mathbb{P}_{|\mathcal{D}}) ,$$

where D is a measure of dissimilarity between the variational distribution \mathbb{Q} and the posterior distribution $\mathbb{P}_{|\mathcal{D}}$

Typically, D : exclusive Kullback-Leibler (KL) divergence and \mathcal{Q} : parametric family (e.g. Mean-field)

$$\begin{cases} D_{KL}(\mathbb{Q} || \mathbb{P}) = \int_Y \log \left(\frac{q(y)}{p(y)} \right) q(y) \nu(dy) \\ \mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathbb{T}\} \end{cases}$$



Question : How to choose D and \mathcal{Q} ?

- Can we select **alternative/more general** D ?
- Can we design more **expressive** variational families \mathcal{Q} ?

Variational Inference in a nutshell

→ Variational Inference : inference is seen as an **optimisation problem**.

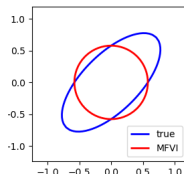
- 1 Posit a *simpler* variational family \mathcal{Q} , where $q \in \mathcal{Q}$.
- 2 Fit q to obtain the best approximation to the posterior density

$$\inf_{q \in \mathcal{Q}} D(\mathbb{Q} || \mathbb{P}_{|\mathcal{D}}) ,$$

where D is a measure of dissimilarity between the variational distribution \mathbb{Q} and the posterior distribution $\mathbb{P}_{|\mathcal{D}}$

Typically, D : exclusive Kullback-Leibler (KL) divergence and \mathcal{Q} : parametric family (e.g. Mean-field)

$$\begin{cases} D_{KL}(\mathbb{Q} || \mathbb{P}) = \int_Y \log \left(\frac{q(y)}{p(y)} \right) q(y) \nu(dy) \\ \mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathbb{T}\} \end{cases}$$



Question : How to choose D and \mathcal{Q} ?

- Can we select **alternative/more general** D ?
- Can we design more **expressive** variational families \mathcal{Q} ?

Variational Inference in a nutshell

→ Variational Inference : inference is seen as an **optimisation problem**.

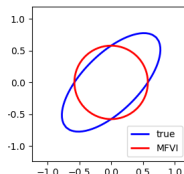
- 1 Posit a *simpler* variational family \mathcal{Q} , where $q \in \mathcal{Q}$.
- 2 Fit q to obtain the best approximation to the posterior density

$$\inf_{q \in \mathcal{Q}} D(\mathbb{Q} || \mathbb{P}_{|\mathcal{D}}) ,$$

where D is a measure of dissimilarity between the variational distribution \mathbb{Q} and the posterior distribution $\mathbb{P}_{|\mathcal{D}}$

Typically, D : exclusive Kullback-Leibler (KL) divergence and \mathcal{Q} : parametric family (e.g. Mean-field)

$$\begin{cases} D_{KL}(\mathbb{Q} || \mathbb{P}) = \int_Y \log \left(\frac{q(y)}{p(y)} \right) q(y) \nu(dy) \\ \mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathbb{T}\} \end{cases}$$



Question : How to choose D and \mathcal{Q} ?

- Can we select **alternative/more general** D ?
- Can we design more **expressive** variational families \mathcal{Q} ?

Variational Inference with the α -divergence family

(Y, \mathcal{Y}, ν) : measured space, ν is a σ -finite measure on (Y, \mathcal{Y}) .

\mathbb{Q} and \mathbb{P} : $\mathbb{Q} \preceq \nu$, $\mathbb{P} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}}{d\nu} = p$.

α -divergence between \mathbb{Q} and \mathbb{P}

$$D_\alpha(\mathbb{Q}||\mathbb{P}) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) ,$$

where

$$f_\alpha = \begin{cases} \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)] , & \text{if } \alpha \in \mathbb{R} \setminus \{0, 1\} , \\ u \log(u) + 1 - u, & \text{if } \alpha = 1 \text{ (Exclusive KL)}, \\ -\log(u) + u - 1, & \text{if } \alpha = 0 \text{ (Inclusive KL)}. \end{cases}$$

Variational Inference with the α -divergence family

(Y, \mathcal{Y}, ν) : measured space, ν is a σ -finite measure on (Y, \mathcal{Y}) .

\mathbb{Q} and \mathbb{P} : $\mathbb{Q} \preceq \nu$, $\mathbb{P} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}}{d\nu} = p$.

α -divergence between \mathbb{Q} and \mathbb{P}

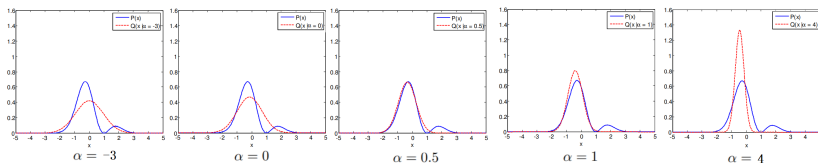
$$D_\alpha(\mathbb{Q}||\mathbb{P}) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) ,$$

where

$$f_\alpha = \begin{cases} \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)] , & \text{if } \alpha \in \mathbb{R} \setminus \{0, 1\} , \\ u \log(u) + 1 - u , & \text{if } \alpha = 1 \text{ (Exclusive KL)} , \\ -\log(u) + u - 1 , & \text{if } \alpha = 0 \text{ (Inclusive KL)} . \end{cases}$$

1 A flexible family of divergences...

Figure: In red, the Gaussian which minimises $D_\alpha(\mathbb{Q}||\mathbb{P})$ for a varying α



Adapted from V. Cevher's lecture notes (2008) <https://www.ece.rice.edu/~vc3/elec633/AlphaDivergence.pdf>

Variational Inference with the α -divergence family

α -divergence between \mathbb{Q} and \mathbb{P}

$$D_{\alpha}(\mathbb{Q}||\mathbb{P}) = \int_{\mathcal{Y}} f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) ,$$

where

$$f_{\alpha} = \begin{cases} \frac{1}{\alpha(\alpha-1)} [u^{\alpha} - 1 - \alpha(u-1)] , & \text{if } \alpha \in \mathbb{R} \setminus \{0, 1\} , \\ u \log(u) + 1 - u, & \text{if } \alpha = 1 \text{ (Forward KL)}, \\ -\log(u) + u - 1, & \text{if } \alpha = 0 \text{ (Reverse KL)}. \end{cases}$$

- ① A **flexible** family of divergences...
- ② ...**suitable** for Variational Inference purposes...

$$\inf_{q \in \mathcal{Q}} D_{\alpha}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \Leftrightarrow \inf_{q \in \mathcal{Q}} \Psi_{\alpha}(q; p)$$

$$\text{with } \Psi_{\alpha}(q; p) = \int_{\mathcal{Y}} f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) \text{ and } p = p(\cdot, \mathcal{D})$$

Black-box alpha divergence minimization. J. Hernandez-Lobato et al. (2016). ICML

Rényi divergence variational inference. Y. Li and R. E Turner (2016). NeurIPS

Variational inference via χ -upper bound minimization A. Dieng et al. (2017). NeurIPS

Variational Inference with the α -divergence family

α -divergence between \mathbb{Q} and \mathbb{P}

$$D_{\alpha}(\mathbb{Q}||\mathbb{P}) = \int_{\mathcal{Y}} f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) ,$$

where

$$f_{\alpha} = \begin{cases} \frac{1}{\alpha(\alpha-1)} [u^{\alpha} - 1 - \alpha(u-1)] , & \text{if } \alpha \in \mathbb{R} \setminus \{0, 1\} , \\ u \log(u) + 1 - u, & \text{if } \alpha = 1 \text{ (Forward KL)}, \\ -\log(u) + u - 1, & \text{if } \alpha = 0 \text{ (Reverse KL)}. \end{cases}$$

- ① A **flexible** family of divergences...
- ② ...**suitable** for Variational Inference purposes...

$$\inf_{q \in \mathcal{Q}} D_{\alpha}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \Leftrightarrow \inf_{q \in \mathcal{Q}} \Psi_{\alpha}(q; p)$$

$$\text{with } \Psi_{\alpha}(q; p) = \int_{\mathcal{Y}} f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) \text{ and } p = p(\cdot, \mathcal{D})$$

Black-box alpha divergence minimization. J. Hernandez-Lobato et al. (2016). ICML

Rényi divergence variational inference. Y. Li and R. E Turner (2016). NeurIPS

Variational inference via χ -upper bound minimization A. Dieng et al. (2017). NeurIPS

Outline

- 1 Introduction
- 2 Monotonic Alpha-Divergence Minimisation**
- 3 Numerical Experiments
- 4 Conclusion

Monotonic Alpha-Divergence Minimisation

Monotonic Alpha-divergence Minimisation.

K. Daudel, R. Douc and F. Roueff (2021). <https://arxiv.org/abs/2103.05684>

Idea : Extend the typical variational parametric family

$$\mathcal{Q} = \{y \mapsto k(\theta, y) : \theta \in \mathbb{T}\}$$

by considering the mixture model variational family

$$\mathcal{Q} = \left\{ q : y \mapsto \mu_{\lambda, \Theta} k(y) := \sum_{j=1}^J \lambda_j k(\theta_j, y) : \lambda \in \mathcal{S}_J, \Theta \in \mathbb{T}^J \right\}$$

and propose an update formula for (λ, Θ) that ensures a systematic decrease in the α -divergence (i.e. Ψ_α) at each step.

Conditions for a monotonic decrease

Optimisation problem

$$\inf_{\lambda \in \mathcal{S}_J, \Theta \in \mathcal{T}^J} \Psi_\alpha(\mu_{\lambda, \Theta} k; p) \quad \text{with} \quad \Psi_\alpha(\mu_{\lambda, \Theta} k; p) = \int_Y f_\alpha \left(\frac{\mu_{\lambda, \Theta} k(y)}{p(y)} \right) p(y) \nu(dy)$$

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

Theorem

Assume (A1) and let $\alpha \in [0, 1)$. Then, choosing $(\lambda_n, \Theta_n)_{n \geq 1}$ so that: $\Psi_\alpha(\mu_{\lambda_1, \Theta_1} k; p) < \infty$ and $\forall n \geq 1$,

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{\lambda_{j,n+1}}{\lambda_{j,n}} \right) \nu(dy) \geq 0 \quad (\text{Weights})$$

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

where $\gamma_{j,\alpha}^n(y) = k(\theta_{j,n}, y) \left(\frac{\mu_{\lambda_n, \Theta_n} k(y)}{p(y)} \right)^{\alpha-1}$, yields a systematic decrease in Ψ_α at each step.

Conditions for a monotonic decrease

Optimisation problem

$$\inf_{\lambda \in \mathcal{S}_J, \Theta \in \mathcal{T}^J} \Psi_\alpha(\mu_{\lambda, \Theta} k; p) \quad \text{with} \quad \Psi_\alpha(\mu_{\lambda, \Theta} k; p) = \int_Y f_\alpha \left(\frac{\mu_{\lambda, \Theta} k(y)}{p(y)} \right) p(y) \nu(dy)$$

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

Theorem

Assume (A1) and let $\alpha \in [0, 1)$. Then, choosing $(\lambda_n, \Theta_n)_{n \geq 1}$ so that:
 $\Psi_\alpha(\mu_{\lambda_1, \Theta_1} k; p) < \infty$ and $\forall n \geq 1$,

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{\lambda_{j,n+1}}{\lambda_{j,n}} \right) \nu(dy) \geq 0 \quad (\text{Weights})$$

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

where $\gamma_{j,\alpha}^n(y) = k(\theta_{j,n}, y) \left(\frac{\mu_{\lambda_n, \Theta_n} k(y)}{p(y)} \right)^{\alpha-1}$, yields a systematic decrease in Ψ_α at each step.

Conditions for a monotonic decrease

Optimisation problem

$$\inf_{\lambda \in \mathcal{S}_J, \Theta \in \mathcal{T}^J} \Psi_\alpha(\mu_{\lambda, \Theta} k; p) \quad \text{with} \quad \Psi_\alpha(\mu_{\lambda, \Theta} k; p) = \int_Y f_\alpha \left(\frac{\mu_{\lambda, \Theta} k(y)}{p(y)} \right) p(y) \nu(dy)$$

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

Theorem

Assume (A1) and let $\alpha \in [0, 1)$. Then, choosing $(\lambda_n, \Theta_n)_{n \geq 1}$ so that:

$\Psi_\alpha(\mu_{\lambda_1, \Theta_1} k; p) < \infty$ and $\forall n \geq 1$,

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{\lambda_{j,n+1}}{\lambda_{j,n}} \right) \nu(dy) \geq 0 \quad (\text{Weights})$$

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

where $\gamma_{j,\alpha}^n(y) = k(\theta_{j,n}, y) \left(\frac{\mu_{\lambda_n, \Theta_n} k(y)}{p(y)} \right)^{\alpha-1}$, yields a systematic decrease in Ψ_α at each step.

Conditions for a monotonic decrease (2)

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{\lambda_{j,n+1}}{\lambda_{j,n}} \right) \nu(dy) \geq 0 \quad (\text{Weights})$$

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

- ① (Weights) and (Components) permit **separate/simultaneous** updates
- ② The dependency is **simpler** in (Weights)

→ (Weights) holds for λ_{n+1} such that

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}, \quad j = 1 \dots J$$

where $\eta_n \in (0, 1]$ and κ is such that $(\alpha - 1) \kappa \geq 0$

Conditions for a monotonic decrease (2)

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{\lambda_{j,n+1}}{\lambda_{j,n}} \right) \nu(dy) \geq 0 \quad (\text{Weights})$$

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

- ❶ (Weights) and (Components) permit **separate/simultaneous** updates
- ❷ The dependency is **simpler** in (Weights)

→ (Weights) holds for λ_{n+1} such that

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}, \quad j = 1 \dots J$$

where $\eta_n \in (0, 1]$ and κ is such that $(\alpha - 1) \kappa \geq 0$

Conditions for a monotonic decrease (2)

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{\lambda_{j,n+1}}{\lambda_{j,n}} \right) \nu(dy) \geq 0 \quad (\text{Weights})$$

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

- ❶ (Weights) and (Components) permit **separate/simultaneous** updates
- ❷ The dependency is **simpler** in (Weights)

→ (Weights) holds for λ_{n+1} such that

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}, \quad j = 1 \dots J$$

where $\eta_n \in (0, 1]$ and κ is such that $(\alpha - 1)\kappa \geq 0$

Conditions for a monotonic decrease (2)

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{\lambda_{j,n+1}}{\lambda_{j,n}} \right) \nu(dy) \geq 0 \quad (\text{Weights})$$

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

- ① (Weights) and (Components) permit **separate/simultaneous** updates
- ② The dependency is **simpler** in (Weights)

→ (Weights) holds for λ_{n+1} such that

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}, \quad j = 1 \dots J$$

where $\eta_n \in (0, 1]$ and κ is such that $(\alpha - 1) \kappa \geq 0$

Understanding the mixture weights update

(Weights) and (Components) hold for λ_{n+1} and Θ_{n+1} such that:

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}, \quad j = 1 \dots J$$
$$\Theta_{n+1} = \Theta_n$$

where $\eta_n \in (0, 1]$ and κ is such that $(\alpha - 1)\kappa \geq 0$

→ We recover the **Power Descent** algorithm from

Infinite-dimensional gradient-based descent for alpha-divergence minimisation.

K. Daudel, R. Douc and F. Portier (2021). *Ann. Statist.* 49 (4) 2250 - 2270.

Core insights :

- 1 The mixture weights update is **gradient-based**, η_n plays the role of a **learning rate**
- 2 We can improve on the Power Descent by proposing **simultaneous updates** for Θ with **convergence guarantees!**

Understanding the mixture weights update

(Weights) and (Components) hold for λ_{n+1} and Θ_{n+1} such that:

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}, \quad j = 1 \dots J$$
$$\Theta_{n+1} = \Theta_n$$

where $\eta_n \in (0, 1]$ and κ is such that $(\alpha - 1)\kappa \geq 0$

→ We recover the **Power Descent** algorithm from

Infinite-dimensional gradient-based descent for alpha-divergence minimisation.

K. Daudel, R. Douc and F. Portier (2021). *Ann. Statist.* 49 (4) 2250 - 2270.

Core insights :

- 1 The mixture weights update is **gradient-based**, η_n plays the role of a **learning rate**
- 2 We can improve on the Power Descent by proposing **simultaneous updates for Θ with convergence guarantees!**

Understanding the mixture weights update

(Weights) and (Components) hold for λ_{n+1} and Θ_{n+1} such that:

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}, \quad j = 1 \dots J$$
$$\Theta_{n+1} = \Theta_n$$

where $\eta_n \in (0, 1]$ and κ is such that $(\alpha - 1)\kappa \geq 0$

→ We recover the **Power Descent** algorithm from

Infinite-dimensional gradient-based descent for alpha-divergence minimisation.

K. Daudel, R. Douc and F. Portier (2021). *Ann. Statist.* 49 (4) 2250 - 2270.

Core insights :

- 1 The mixture weights update is **gradient-based**, η_n plays the role of a **learning rate**
- 2 We can improve on the Power Descent by proposing **simultaneous updates** for Θ with **convergence guarantees**!

Understanding the mixture weights update

(Weights) and (Components) hold for λ_{n+1} and Θ_{n+1} such that:

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}, \quad j = 1 \dots J$$
$$\Theta_{n+1} = \Theta_n$$

where $\eta_n \in (0, 1]$ and κ is such that $(\alpha - 1)\kappa \geq 0$

→ We recover the **Power Descent** algorithm from

Infinite-dimensional gradient-based descent for alpha-divergence minimisation.

K. Daudel, R. Douc and F. Portier (2021). *Ann. Statist.* 49 (4) 2250 - 2270.

Core insights :

- 1 The mixture weights update is **gradient-based**, η_n plays the role of a **learning rate**
- 2 We can improve on the Power Descent by proposing **simultaneous updates for Θ with convergence guarantees!**

Understanding the mixture weights update

(Weights) and (Components) hold for λ_{n+1} and Θ_{n+1} such that:

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}, \quad j = 1 \dots J$$
$$\Theta_{n+1} = \Theta_n$$

where $\eta_n \in (0, 1]$ and κ is such that $(\alpha - 1)\kappa \geq 0$

→ We recover the **Power Descent** algorithm from

Infinite-dimensional gradient-based descent for alpha-divergence minimisation.

K. Daudel, R. Douc and F. Portier (2021). *Ann. Statist.* 49 (4) 2250 - 2270.

Core insights :

- 1 The mixture weights update is **gradient-based**, η_n plays the role of a **learning rate**
- 2 We can improve on the Power Descent by proposing **simultaneous updates for Θ with convergence guarantees!**

Towards simultaneous updates

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

- Maximisation approach : for all $j = 1 \dots J$,

$$\theta_{j,n+1} = \operatorname{argmax}_{\theta \in \mathcal{T}} \int_Y \gamma_{j,\alpha}^n(y) \log(k(\theta, y)) \nu(dy)$$

- Gradient-based approach : for all $j = 1 \dots J$, $\gamma_{j,n} \in (0, 1]$, $c_{j,n} > 0$

$$\theta_{j,n+1} = \theta_{j,n} - \frac{\gamma_{j,n}}{\beta_{j,n}} \nabla g_{j,n}(\theta)|_{\theta=\theta_{j,n}}$$

where $g_{j,n}$ is assumed to be $\beta_{j,n}$ -smooth on $\mathcal{T} = \mathbb{R}^d$ with

$$g_{j,n}(\theta) = c_{j,n} \int_Y \frac{\gamma_{j,\alpha}^n(y)}{\alpha - 1} \log \left(\frac{k(\theta, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) .$$

Towards simultaneous updates

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

- Maximisation approach : for all $j = 1 \dots J$,

$$\theta_{j,n+1} = \operatorname{argmax}_{\theta \in \mathcal{T}} \int_Y \gamma_{j,\alpha}^n(y) \log(k(\theta, y)) \nu(dy)$$

- Gradient-based approach : for all $j = 1 \dots J$, $\gamma_{j,n} \in (0, 1]$, $c_{j,n} > 0$

$$\theta_{j,n+1} = \theta_{j,n} - \frac{\gamma_{j,n}}{\beta_{j,n}} \nabla g_{j,n}(\theta)|_{\theta=\theta_{j,n}}$$

where $g_{j,n}$ is assumed to be $\beta_{j,n}$ -smooth on $\mathcal{T} = \mathbb{R}^d$ with

$$g_{j,n}(\theta) = c_{j,n} \int_Y \frac{\gamma_{j,\alpha}^n(y)}{\alpha - 1} \log \left(\frac{k(\theta, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) .$$

Towards simultaneous updates

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

- Maximisation approach : for all $j = 1 \dots J$, $a_{j,n} > 0$, $b_{j,n} \geq 0$ and

$$\theta_{j,n+1} = \operatorname{argmax}_{\theta \in T} \int_Y [a_{j,n} \gamma_{j,\alpha}^n(y) + b_{j,n} k(\theta_{j,n}, y)] \log(k(\theta, y)) \nu(dy)$$

- Gradient-based approach : for all $j = 1 \dots J$, $\gamma_{j,n} \in (0, 1]$, $c_{j,n} > 0$

$$\theta_{j,n+1} = \theta_{j,n} - \frac{\gamma_{j,n}}{\beta_{j,n}} \nabla g_{j,n}(\theta)|_{\theta=\theta_{j,n}}$$

where $g_{j,n}$ is assumed to be $\beta_{j,n}$ -smooth on $T = \mathbb{R}^d$ with

$$g_{j,n}(\theta) = c_{j,n} \int_Y \frac{\gamma_{j,\alpha}^n(y)}{\alpha - 1} \log \left(\frac{k(\theta, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) .$$

Towards simultaneous updates

$$\int_Y \sum_{j=1}^J \lambda_{j,n} \gamma_{j,\alpha}^n(y) \log \left(\frac{k(\theta_{j,n+1}, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) \geq 0 \quad (\text{Components})$$

- Maximisation approach : for all $j = 1 \dots J$, $a_{j,n} > 0$, $b_{j,n} \geq 0$ and

$$\theta_{j,n+1} = \operatorname{argmax}_{\theta \in T} \int_Y [a_{j,n} \gamma_{j,\alpha}^n(y) + b_{j,n} k(\theta_{j,n}, y)] \log(k(\theta, y)) \nu(dy)$$

- Gradient-based approach : for all $j = 1 \dots J$, $\gamma_{j,n} \in (0, 1]$, $c_{j,n} > 0$

$$\theta_{j,n+1} = \theta_{j,n} - \frac{\gamma_{j,n}}{\beta_{j,n}} \nabla g_{j,n}(\theta)|_{\theta=\theta_{j,n}}$$

where $g_{j,n}$ is assumed to be $\beta_{j,n}$ -smooth on $T = \mathbb{R}^d$ with

$$g_{j,n}(\theta) = c_{j,n} \int_Y \frac{\gamma_{j,\alpha}^n(y)}{\alpha - 1} \log \left(\frac{k(\theta, y)}{k(\theta_{j,n}, y)} \right) \nu(dy) .$$

Related work

Maximisation approach (Gaussian case)

$$\begin{aligned}\lambda_{j,n+1} &= \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}} \\ m_{j,n+1} &= (1 - \gamma_{j,n})m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)} \\ \Sigma_{j,n+1} &= (1 - \gamma_{j,n})\tilde{\Sigma}_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}\end{aligned}$$

NB : $\gamma_{j,n} = (a_{j,n} \int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + b_{j,n})^{-1}$ with $b_{j,n} \geq 0$

Adaptive importance sampling in general mixture classes. O. Cappé, R. Douc, A. Guillin, J-M Marin and C. P Robert (2008). *Statistics and Computing*, 18(4):447–459

→ The M-PMC algorithm a.k.a ‘Integrated EM’ corresponds to

$$\alpha = 0, \eta_n = 1, \kappa = 0 \text{ and } \gamma_{j,n} = 1$$

We have **generalised** an integrated EM algorithm for mixture models optimisation

- 1 We introduce η_n , κ and $\gamma_{j,n}$, where η_n and $\gamma_{j,n}$ act as **learning rates**
- 2 We extend the **systematic** decrease property to $\alpha \in [0, 1)$

Maximisation approach (Gaussian case)

$$\begin{aligned}\lambda_{j,n+1} &= \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}} \\ m_{j,n+1} &= (1 - \gamma_{j,n})m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)} \\ \Sigma_{j,n+1} &= (1 - \gamma_{j,n})\tilde{\Sigma}_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}\end{aligned}$$

NB : $\gamma_{j,n} = (a_{j,n} \int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + b_{j,n})^{-1}$ with $b_{j,n} \geq 0$

Adaptive importance sampling in general mixture classes. O. Cappé, R. Douc, A. Guillin, J-M Marin and C. P Robert (2008). *Statistics and Computing*, 18(4):447–459

→ The M-PMC algorithm a.k.a ‘Integrated EM’ corresponds to

$$\alpha = 0, \eta_n = 1, \kappa = 0 \text{ and } \gamma_{j,n} = 1$$

We have **generalised** an integrated EM algorithm for mixture models optimisation

- 1 We introduce η_n , κ and $\gamma_{j,n}$, where η_n and $\gamma_{j,n}$ act as **learning rates**
- 2 We extend the **systematic** decrease property to $\alpha \in [0, 1)$

Maximisation approach (Gaussian case)

$$\begin{aligned}\lambda_{j,n+1} &= \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}} \\ m_{j,n+1} &= (1 - \gamma_{j,n})m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)} \\ \Sigma_{j,n+1} &= (1 - \gamma_{j,n})\tilde{\Sigma}_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}\end{aligned}$$

NB : $\gamma_{j,n} = (a_{j,n} \int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + b_{j,n})^{-1}$ with $b_{j,n} \geq 0$

Adaptive importance sampling in general mixture classes. O. Cappé, R. Douc, A. Guillin, J-M Marin and C. P Robert (2008). *Statistics and Computing*, 18(4):447–459

→ The M-PMC algorithm a.k.a 'Integrated EM' corresponds to

$$\alpha = 0, \eta_n = 1, \kappa = 0 \text{ and } \gamma_{j,n} = 1$$

We have generalised an integrated EM algorithm for mixture models optimisation

- 1 We introduce η_n , κ and $\gamma_{j,n}$, where η_n and $\gamma_{j,n}$ act as **learning rates**
- 2 We extend the **systematic** decrease property to $\alpha \in [0, 1)$

Maximisation approach (Gaussian case)

$$\begin{aligned}\lambda_{j,n+1} &= \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}} \\ m_{j,n+1} &= (1 - \gamma_{j,n})m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)} \\ \Sigma_{j,n+1} &= (1 - \gamma_{j,n})\tilde{\Sigma}_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}\end{aligned}$$

NB : $\gamma_{j,n} = (a_{j,n} \int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + b_{j,n})^{-1}$ with $b_{j,n} \geq 0$

Adaptive importance sampling in general mixture classes. O. Cappé, R. Douc, A. Guillin, J-M Marin and C. P Robert (2008). *Statistics and Computing*, 18(4):447–459

→ The M-PMC algorithm a.k.a ‘Integrated EM’ corresponds to

$$\alpha = 0, \eta_n = 1, \kappa = 0 \quad \text{and} \quad \gamma_{j,n} = 1$$

We have **generalised** an integrated EM algorithm for mixture models optimisation

- 1 We introduce η_n , κ and $\gamma_{j,n}$, where η_n and $\gamma_{j,n}$ act as **learning rates**
- 2 We extend the **systematic** decrease property to $\alpha \in [0, 1)$

Maximisation approach (Gaussian case)

$$\begin{aligned}\lambda_{j,n+1} &= \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}} \\ m_{j,n+1} &= (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)} \\ \Sigma_{j,n+1} &= (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}\end{aligned}$$

NB : $\gamma_{j,n} = (a_{j,n} \int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + b_{j,n})^{-1}$ with $b_{j,n} \geq 0$

Adaptive importance sampling in general mixture classes. O. Cappé, R. Douc, A. Guillin, J-M Marin and C. P Robert (2008). *Statistics and Computing*, 18(4):447–459

→ The M-PMC algorithm a.k.a ‘Integrated EM’ corresponds to

$$\alpha = 0, \eta_n = 1, \kappa = 0 \quad \text{and} \quad \gamma_{j,n} = 1$$

We have **generalised** an integrated EM algorithm for mixture models optimisation

- 1 We introduce η_n , κ and $\gamma_{j,n}$, where η_n and $\gamma_{j,n}$ act as **learning rates**
- 2 We extend the **systematic** decrease property to $\alpha \in [0, 1)$

Maximisation approach (Gaussian case)

$$\begin{aligned}\lambda_{j,n+1} &= \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}} \\ m_{j,n+1} &= (1 - \gamma_{j,n})m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)} \\ \Sigma_{j,n+1} &= (1 - \gamma_{j,n})\tilde{\Sigma}_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}\end{aligned}$$

NB : $\gamma_{j,n} = (a_{j,n} \int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + b_{j,n})^{-1}$ with $b_{j,n} \geq 0$

Adaptive importance sampling in general mixture classes. O. Cappé, R. Douc, A. Guillin, J-M Marin and C. P Robert (2008). *Statistics and Computing*, 18(4):447–459

→ The M-PMC algorithm a.k.a ‘Integrated EM’ corresponds to

$$\alpha = 0, \eta_n = 1, \kappa = 0 \quad \text{and} \quad \gamma_{j,n} = 1$$

We have **generalised** an integrated EM algorithm for mixture models optimisation

- 1 We introduce η_n , κ and $\gamma_{j,n}$, where η_n and $\gamma_{j,n}$ act as **learning rates**
- 2 We extend the **systematic** decrease property to $\alpha \in [0, 1)$

Maximisation approach (Gaussian case)

$$\begin{aligned}\lambda_{j,n+1} &= \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}} \\ m_{j,n+1} &= (1 - \gamma_{j,n})m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)} \\ \Sigma_{j,n+1} &= (1 - \gamma_{j,n})\tilde{\Sigma}_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}\end{aligned}$$

NB : $\gamma_{j,n} = (a_{j,n} \int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + b_{j,n})^{-1}$ with $b_{j,n} \geq 0$

Adaptive importance sampling in general mixture classes. O. Cappé, R. Douc, A. Guillin, J-M Marin and C. P Robert (2008). *Statistics and Computing*, 18(4):447–459

→ The M-PMC algorithm a.k.a ‘Integrated EM’ corresponds to

$$\alpha = 0, \eta_n = 1, \kappa = 0 \quad \text{and} \quad \gamma_{j,n} = 1$$

We have **generalised** an integrated EM algorithm for mixture models optimisation

- 1 We introduce η_n , κ and $\gamma_{j,n}$, where η_n and $\gamma_{j,n}$ act as **learning rates**
- 2 We extend the **systematic** decrease property to $\alpha \in [0, 1]$

Gradient-based approach (Gaussian case)

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}$$

$$(RGD) \quad m_{j,n+1} = m_{j,n} + \gamma_{j,n} \frac{\int_Y \lambda_{j,n} \gamma_{j,\alpha}^n(y) (y - m_{j,n}) \nu(dy)}{\int_Y \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(dy)}$$

$$(MG) \quad m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}$$

- **(RGD)**. Set $p = p(\cdot, \mathcal{D})$, $\gamma_{j,n} := \gamma_n \in (0, 1]$. Usual gradient descent steps on Θ for Rényi's α -divergence minimisation

Rényi divergence variational inference. Y. Li and R. E Turner (2016). NeurIPS

NB : We provide simultaneous updates for λ that preserve the convergence!

- **(MG)**. 'EM-like' : coincides with the mixture means update from the maximisation approach.

NB : Our maximisation approach gives updates for the covariance matrices too!

Gradient-based approach (Gaussian case)

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}$$
$$\text{(RGD)} \quad m_{j,n+1} = m_{j,n} + \gamma_{j,n} \frac{\int_Y \lambda_{j,n} \gamma_{j,\alpha}^n(y) (y - m_{j,n}) \nu(dy)}{\int_Y \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(dy)}$$
$$\text{(MG)} \quad m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}$$

- **(RGD)**. Set $p = p(\cdot, \mathcal{D})$, $\gamma_{j,n} := \gamma_n \in (0, 1]$. Usual gradient descent steps on Θ for Rényi's α -divergence minimisation

Rényi divergence variational inference. Y. Li and R. E Turner (2016). NeurIPS

NB : We provide simultaneous updates for λ that preserve the convergence!

- **(MG)**. 'EM-like' : coincides with the mixture means update from the maximisation approach.

NB : Our maximisation approach gives updates for the covariance matrices too!

Gradient-based approach (Gaussian case)

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}$$
$$\text{(RGD)} \quad m_{j,n+1} = m_{j,n} + \gamma_{j,n} \frac{\int_Y \lambda_{j,n} \gamma_{j,\alpha}^n(y) (y - m_{j,n}) \nu(dy)}{\int_Y \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(dy)}$$
$$\text{(MG)} \quad m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}$$

- **(RGD)**. Set $p = p(\cdot, \mathcal{D})$, $\gamma_{j,n} := \gamma_n \in (0, 1]$. Usual gradient descent steps on Θ for Rényi's α -divergence minimisation

Rényi divergence variational inference. Y. Li and R. E Turner (2016). NeurIPS

NB : We provide simultaneous updates for λ that preserve the convergence!

- **(MG)**. 'EM-like' : coincides with the mixture means update from the maximisation approach.

NB : Our maximisation approach gives updates for the covariance matrices too!

Gradient-based approach (Gaussian case)

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1)\kappa \right]^{\eta_n}}$$
$$\text{(RGD)} \quad m_{j,n+1} = m_{j,n} + \gamma_{j,n} \frac{\int_Y \lambda_{j,n} \gamma_{j,\alpha}^n(y) (y - m_{j,n}) \nu(dy)}{\int_Y \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(dy)}$$
$$\text{(MG)} \quad m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}$$

- **(RGD)**. Set $p = p(\cdot, \mathcal{D})$, $\gamma_{j,n} := \gamma_n \in (0, 1]$. Usual gradient descent steps on Θ for Rényi's α -divergence minimisation

Rényi divergence variational inference. Y. Li and R. E Turner (2016). NeurIPS

NB : We provide simultaneous updates for λ that preserve the convergence!

- **(MG)**. 'EM-like' : coincides with the mixture means update from the maximisation approach.

NB : Our maximisation approach gives updates for the covariance matrices too!

Gradient-based approach (Gaussian case)

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}$$
$$\text{(RGD)} \quad m_{j,n+1} = m_{j,n} + \gamma_{j,n} \frac{\int_Y \lambda_{j,n} \gamma_{j,\alpha}^n(y) (y - m_{j,n}) \nu(dy)}{\int_Y \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(dy)}$$
$$\text{(MG)} \quad m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}$$

- **(RGD)**. Set $p = p(\cdot, \mathcal{D})$, $\gamma_{j,n} := \gamma_n \in (0, 1]$. Usual gradient descent steps on Θ for Rényi's α -divergence minimisation

Rényi divergence variational inference. Y. Li and R. E Turner (2016). NeurIPS

NB : We provide simultaneous updates for λ that preserve the convergence!

- **(MG)**. 'EM-like' : coincides with the mixture means update from the maximisation approach.

NB : Our maximisation approach gives updates for the covariance matrices too!

Gradient-based approach (Gaussian case)

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_Y \gamma_{j,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_Y \gamma_{\ell,\alpha}^n(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}$$
$$\text{(RGD)} \quad m_{j,n+1} = m_{j,n} + \gamma_{j,n} \frac{\int_Y \lambda_{j,n} \gamma_{j,\alpha}^n(y) (y - m_{j,n}) \nu(dy)}{\int_Y \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(dy)}$$
$$\text{(MG)} \quad m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \frac{\int_Y \gamma_{j,\alpha}^n(y) y \nu(dy)}{\int_Y \gamma_{j,\alpha}^n(y) \nu(dy)}$$

- **(RGD)**. Set $p = p(\cdot, \mathcal{D})$, $\gamma_{j,n} := \gamma_n \in (0, 1]$. Usual gradient descent steps on Θ for Rényi's α -divergence minimisation

Rényi divergence variational inference. Y. Li and R. E Turner (2016). NeurIPS

NB : We provide simultaneous updates for λ that preserve the convergence!

- **(MG)**. 'EM-like' : coincides with the mixture means update from the maximisation approach.

NB : Our maximisation approach gives updates for the covariance matrices too!

Outline

- 1 Introduction
- 2 Monotonic Alpha-Divergence Minimisation
- 3 Numerical Experiments**
- 4 Conclusion

Monte Carlo approximations

Algorithm 1: Gaussian Mixture Models optimisation

At iteration n ,

- ❶ Draw independently M samples $(Y_{m,n})_{1 \leq m \leq M}$ from the proposal q_n . Define for all $j = 1 \dots J$, all $y \in \mathcal{Y}$ and all $n \geq 1$, $\hat{\gamma}_{j,\alpha}^n(y) = k(\theta_{j,n}, y)/q_n(y) (\mu_n k(y)/p(y))^{\alpha-1}$.
- ❷ For all $j = 1 \dots J$, set:

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\sum_{m=1}^M \hat{\gamma}_{j,\alpha}^n(Y_{m,n}) + (\alpha - 1)\kappa_n \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\sum_{m=1}^M \hat{\gamma}_{\ell,\alpha}^n(Y_{m,n}) + (\alpha - 1)\kappa_n \right]^{\eta_n}}$$

$$(RGD) \quad m_{j,n+1} = m_{j,n} + \gamma_n \frac{\lambda_{j,n} \sum_{m=1}^M \hat{\gamma}_{j,\alpha}^n(Y_{m,n}) \cdot (Y_{m,n} - \theta_{j,n})}{\sum_{j=1}^J \sum_{m=1}^M \lambda_{j,n} \hat{\gamma}_{j,\alpha}^n(Y_{m,n})}$$

$$(MG) \quad m_{j,n+1} = (1 - \gamma_n) m_{j,n} + \gamma_n \frac{\sum_{m=1}^M \hat{\gamma}_{j,\alpha}^n(Y_{m,n}) \cdot Y_{m,n}}{\sum_{m=1}^M \hat{\gamma}_{j,\alpha}^n(Y_{m,n})}$$

→ Here $\hat{\gamma}_{j,\alpha}^n(y) = \frac{\gamma_{j,\alpha}^n(y)}{q_n(y)}$.

→ We consider two samplers : $q_n = \mu_{\lambda_n, \Theta_n}$ (IS-n) and $q_n = J^{-1} \sum_{j=1}^J k(\theta_{j,n}, \cdot)$ (IS-unif).

Monte Carlo approximations

Algorithm 1: Gaussian Mixture Models optimisation

At iteration n ,

- ❶ Draw independently M samples $(Y_{m,n})_{1 \leq m \leq M}$ from the proposal q_n . Define for all $j = 1 \dots J$, all $y \in \mathcal{Y}$ and all $n \geq 1$, $\hat{\gamma}_{j,\alpha}^n(y) = k(\theta_{j,n}, y)/q_n(y) (\mu_n k(y)/p(y))^{\alpha-1}$.
- ❷ For all $j = 1 \dots J$, set:

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\sum_{m=1}^M \hat{\gamma}_{j,\alpha}^n(Y_{m,n}) + (\alpha - 1)\kappa_n \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\sum_{m=1}^M \hat{\gamma}_{\ell,\alpha}^n(Y_{m,n}) + (\alpha - 1)\kappa_n \right]^{\eta_n}}$$

$$(RGD) \quad m_{j,n+1} = m_{j,n} + \gamma_n \frac{\lambda_{j,n} \sum_{m=1}^M \hat{\gamma}_{j,\alpha}^n(Y_{m,n}) \cdot (Y_{m,n} - \theta_{j,n})}{\sum_{j=1}^J \sum_{m=1}^M \lambda_{j,n} \hat{\gamma}_{j,\alpha}^n(Y_{m,n})}$$

$$(MG) \quad m_{j,n+1} = (1 - \gamma_n) m_{j,n} + \gamma_n \frac{\sum_{m=1}^M \hat{\gamma}_{j,\alpha}^n(Y_{m,n}) \cdot Y_{m,n}}{\sum_{m=1}^M \hat{\gamma}_{j,\alpha}^n(Y_{m,n})}$$

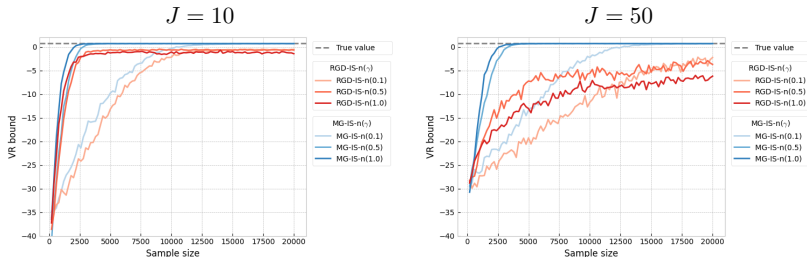
→ Here $\hat{\gamma}_{j,\alpha}^n(y) = \frac{\gamma_{j,\alpha}^n(y)}{q_n(y)}$.

→ We consider two samplers : $q_n = \mu_{\lambda_n, \Theta_n}$ (IS-n) and $q_n = J^{-1} \sum_{j=1}^J k(\theta_{j,n}, \cdot)$ (IS-unif).

Comparing RGD to MG (fixed λ)

Target : $p(y) = 2 \times [0.5\mathcal{N}(y; -2\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(y; 2\mathbf{u}_d, \mathbf{I}_d)]$

- MC estimate of the VR Bound averaged over 30 trials for RGD and MG.
[Here, $\alpha = 0.2$, $d = 16$, $M = 200$, $\kappa_n = 0$, $\eta_n = 0$. and $q_n = \mu_n k$.]



- LogMSE averaged over 30 trials for RGD and MG.

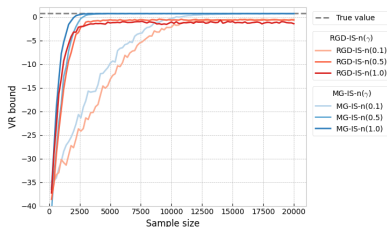
	$J = 10$			$J = 50$		
	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$
RGD-IS-n(γ)	-0.081	-0.076	-0.218	-1.640	-1.673	-1.560
MG-IS-n(γ)	-3.702	-1.875	-2.711	-2.760	-2.771	-2.788

Comparing RGD to MG (fixed λ)

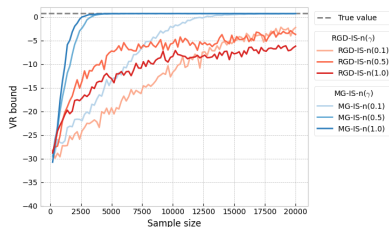
Target : $p(y) = 2 \times [0.5\mathcal{N}(y; -2\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(y; 2\mathbf{u}_d, \mathbf{I}_d)]$

- MC estimate of the VR Bound averaged over 30 trials for RGD and MG.
[Here, $\alpha = 0.2$, $d = 16$, $M = 200$, $\kappa_n = 0$, $\eta_n = 0$. and $q_n = \mu_n k$.]

$J = 10$



$J = 50$



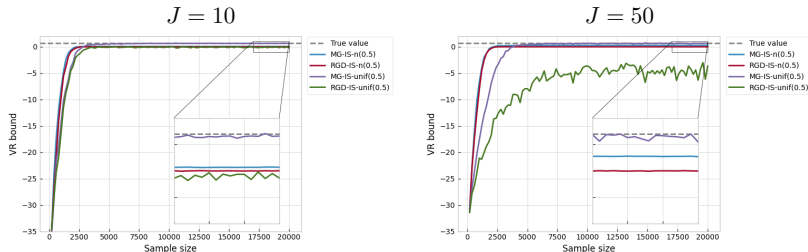
- LogMSE averaged over 30 trials for RGD and MG.

	$J = 10$			$J = 50$		
	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$
RGD-IS-n(γ)	-0.081	-0.076	-0.218	-1.640	-1.673	-1.560
MG-IS-n(γ)	-3.702	-1.875	-2.711	-2.760	-2.771	-2.788

Comparing RGD to MG (varying λ)

Target : $p(y) = 2 \times [0.5\mathcal{N}(y; -2\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(y; 2\mathbf{u}_d, \mathbf{I}_d)]$

- MC estimate of the VR Bound averaged over 30 trials for RGD and MG.
[Here, $\alpha = 0.2$, $d = 16$, $M = 200$, $\eta = 0.1$, $\kappa_n = 0$.]



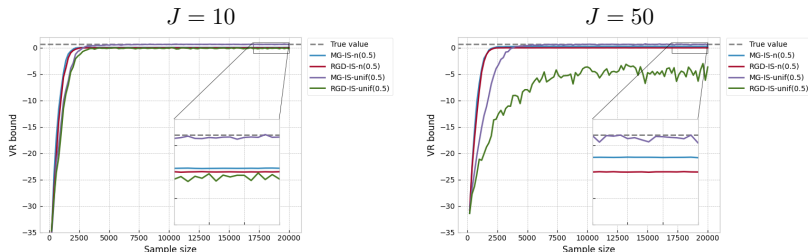
- LogMSE averaged over 30 trials for RGD and MG.

	$J = 10$			$J = 50$		
	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$
RGD-IS-n(γ)	0.372	0.510	0.384	-0.616	-0.713	-0.778
MG-IS-n(γ)	1.104	1.074	0.387	1.135	-0.077	-0.060
RGD-IS-unif(γ)	0.359	0.469	0.458	-0.688	-0.670	-0.583
MG-IS-unif(γ)	-0.200	-0.229	-0.515	-1.500	-1.462	-1.246

Comparing RGD to MG (varying λ)

$$\text{Target : } p(y) = 2 \times [0.5\mathcal{N}(y; -2\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(y; 2\mathbf{u}_d, \mathbf{I}_d)]$$

- MC estimate of the VR Bound averaged over 30 trials for RGD and MG.
[Here, $\alpha = 0.2$, $d = 16$, $M = 200$, $\eta = 0.1$, $\kappa_n = 0$.]



- LogMSE averaged over 30 trials for RGD and MG.

	$J = 10$			$J = 50$		
	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$
RGD-IS-n(γ)	0.372	0.510	0.384	-0.616	-0.713	-0.778
MG-IS-n(γ)	1.104	1.074	0.387	1.135	-0.077	-0.060
RGD-IS-unif(γ)	0.359	0.469	0.458	-0.688	-0.670	-0.583
MG-IS-unif(γ)	-0.200	-0.229	-0.515	-1.500	-1.462	-1.246

Comparing RGD to MG (varying λ , cont'd)

$$\text{Target : } p(y) = 2 \times [0.5\mathcal{N}(\mathbf{y}; -2\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(\mathbf{y}; 2\mathbf{u}_d, \mathbf{I}_d)]$$

- LogMSE averaged over 30 trials for RGD and MG.

[Here, $\alpha = 0.2$, $d = 16$, $M = 200$, $\gamma = 0.5$, $\kappa_n = 0$.]

	$J = 10$			$J = 50$		
	$\eta = 0.05$	$\eta = 0.1$	$\eta = 0.5$	$\eta = 0.05$	$\eta = 0.1$	$\eta = 0.5$
RGD-IS-n(γ)	0.045	0.510	1.299	-1.355	-0.713	0.924
MG-IS-n(γ)	0.087	1.074	1.343	-1.205	-0.077	1.329
RGD-IS-unif(γ)	-0.018	0.469	1.328	-1.385	-0.670	0.928
MG-IS-unif(γ)	-1.244	-0.229	1.100	-2.524	-1.462	0.309

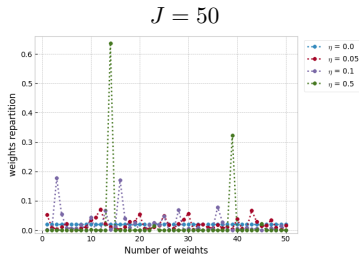
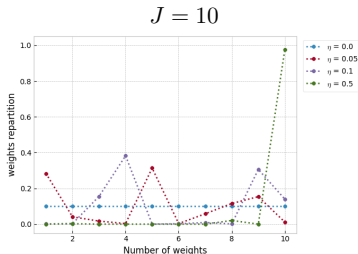
Comparing RGD to MG (varying λ , cont'd)

$$\text{Target : } p(y) = 2 \times [0.5\mathcal{N}(y; -2\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(y; 2\mathbf{u}_d, \mathbf{I}_d)]$$

- LogMSE averaged over 30 trials for RGD and MG.

[Here, $\alpha = 0.2$, $d = 16$, $M = 200$, $\gamma = 0.5$, $\kappa_n = 0$.]

	$J = 10$			$J = 50$		
	$\eta = 0.05$	$\eta = 0.1$	$\eta = 0.5$	$\eta = 0.05$	$\eta = 0.1$	$\eta = 0.5$
RGD-IS-n(γ)	0.045	0.510	1.299	-1.355	-0.713	0.924
MG-IS-n(γ)	0.087	1.074	1.343	-1.205	-0.077	1.329
RGD-IS-unif(γ)	-0.018	0.469	1.328	-1.385	-0.670	0.928
MG-IS-unif(γ)	-1.244	-0.229	1.100	-2.524	-1.462	0.309



Outline

- 1 Introduction
- 2 Monotonic Alpha-Divergence Minimisation
- 3 Numerical Experiments
- 4 Conclusion**

Conclusion

Novel framework for **monotonic α -divergence minimisation**

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **empirical benefits** of our general framework

Perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications

Conclusion

Novel framework for **monotonic α -divergence minimisation**

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **empirical benefits** of our general framework

Perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications

Conclusion

Novel framework for **monotonic α -divergence minimisation**

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **empirical benefits** of our general framework

Perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications

Conclusion

Novel framework for **monotonic α -divergence minimisation**

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **empirical benefits** of our general framework

Perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications

Conclusion

Novel framework for **monotonic α -divergence minimisation**

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **empirical benefits** of our general framework

Perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications

Conclusion

Novel framework for **monotonic α -divergence minimisation**

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **empirical benefits** of our general framework

Perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications

Conclusion

Novel framework for **monotonic α -divergence minimisation**

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **empirical benefits** of our general framework

Perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications

Conclusion

Novel framework for **monotonic α -divergence minimisation**

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **empirical benefits** of our general framework

Perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications

Thank you for your attention!

kamelia.daudel@stats.ox.ac.uk

Monotonic Alpha-divergence Minimisation

K. Daudel, R. Douc and F. Roueff (2021). <https://arxiv.org/abs/2103.05684>

Infinite-dimensional gradient-based descent for alpha-divergence minimisation.

K. Daudel, R. Douc and F. Portier (2020). Ann. Statist. 49 (4) 2250 - 2270.