

Infinite-dimensional α -divergence minimisation for Variational Inference

Kamélia Daudel

Télécom Paris, Institut Polytechnique de Paris
kamelia.daudel@telecom-paris.fr

Variational Inference Seminar
24/11/2020

Joint work with Randal Douc and François Portier



Introduction

Goal : build an iterative scheme

$$\mu_{n+1} = \mathcal{I}_\alpha(\mu_n) , \quad n \in \mathbb{N}^\star ,$$

- which **extends** the commonly-used variational approximating family (**Infinite-dimensional Variational Inference**),
- such that one iteration leads to a **systematic decrease** of a certain criterion (**α -divergence**).

Introduction

Goal : build an iterative scheme

$$\mu_{n+1} = \mathcal{I}_\alpha(\mu_n) , \quad n \in \mathbb{N}^\star ,$$

- which extends the commonly-used variational approximating family (Infinite-dimensional Variational Inference),
- such that one iteration leads to a systematic decrease of a certain criterion (α -divergence).

Introduction

Goal : build an iterative scheme

$$\mu_{n+1} = \mathcal{I}_\alpha(\mu_n) , \quad n \in \mathbb{N}^\star ,$$

- which extends the commonly-used variational approximating family (Infinite-dimensional Variational Inference),
- such that one iteration leads to a systematic decrease of a certain criterion (α -divergence).

Outline

- ① Background
- ② The (α, Γ) -descent
- ③ Numerical Experiments
- ④ Take-away message
- ⑤ Proof of the systematic decrease

Outline

- 1 Background
- 2 The (α, Γ) -descent
- 3 Numerical Experiments
- 4 Take-away message
- 5 Proof of the systematic decrease

Variational Inference in a nutshell

- 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 : inference is seen as an **optimisation problem**.

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

$$q^* = \operatorname{arginf}_{q \in \mathcal{Q}} D(Q||\mathbb{P}) ,$$

where D is the a divergence (e.g the Kullback-Leibler).

Variational Inference in a nutshell

- 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 : inference is seen as an **optimisation problem**.

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

$$q^* = \operatorname{arginf}_{q \in \mathcal{Q}} D(Q||\mathbb{P}) ,$$

where D is the a divergence (e.g the Kullback-Leibler).

Variational Inference within the α -divergence family (1)

(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 $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}}{d\nu} = p(\cdot|\mathcal{D})$.

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

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

where

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

Variational Inference within the α -divergence family (1)

(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 $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}}{d\nu} = p(\cdot|\mathcal{D})$.

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

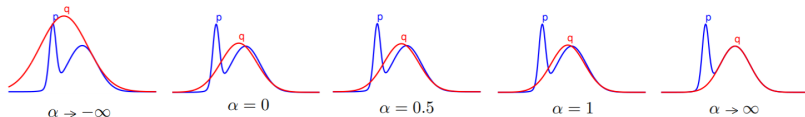
$$D_\alpha(\mathbb{Q}||\mathbb{P}) = \int_Y f_\alpha \left(\frac{q(y)}{p(y|\mathcal{D})} \right) p(y|\mathcal{D}) \nu(dy) ,$$

where

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

❶ A **flexible** family of divergences...

Figure: The Gaussian q which minimizes α -divergence to p (a mixture of two Gaussian), for varying α



[Adapted from T. Minka (2005) Divergence Measures and Message Passing. Technical Report MSR-TR-2005-173]

Variational Inference within the α -divergence family (2)

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

where

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

❶ A **flexible** family of divergences...

❷ ...**suitable** for Variational Inference purposes...

→ We can get rid of $p(\mathcal{D})$ in the optimisation !

$$\begin{aligned} q^* &= \operatorname{arginf}_{q \in \mathcal{Q}} D_{\alpha}(\mathbb{Q}||\mathbb{P}) \\ &= \operatorname{arginf}_{q \in \mathcal{Q}} \int_{\mathcal{Y}} f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) \quad \text{with } p(y) = p(y, \mathcal{D}) . \end{aligned}$$

[J. Hernandez-Lobato, Y. Li, M. Rowland, T. Bui, D. Hernandez-Lobato, and R. E Turner. (2016) Black-box alpha divergence minimization. ICML]

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

[A. Dieng, D. Tran, R. Ranganath, J. Paisley, and D. Blei. (2017) Variational inference via χ -upper bound minimization. NeurIPS]

Variational Inference within the α -divergence family (2)

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

where

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

❶ A **flexible** family of divergences...

❷ ...**suitable** for Variational Inference purposes...

→ We can get rid of $p(\mathcal{D})$ in the optimisation !

$$\begin{aligned} q^{\star} &= \operatorname{arginf}_{q \in \mathcal{Q}} D_{\alpha}(\mathbb{Q}||\mathbb{P}) \\ &= \operatorname{arginf}_{q \in \mathcal{Q}} \int_{\mathcal{Y}} f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) \quad \text{with } p(y) = p(y, \mathcal{D}) . \end{aligned}$$

[J. Hernandez-Lobato, Y. Li, M. Rowland, T. Bui, D. Hernandez-Lobato, and R. E Turner. (2016) Black-box alpha divergence minimization. ICML]

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

[A. Dieng, D. Tran, R. Ranganath, J. Paisley, and D. Blei. (2017) Variational inference via χ -upper bound minimization. NeurIPS]

Variational Inference within the α -divergence family (2)

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

where

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

❶ A **flexible** family of divergences...

❷ ...**suitable** for Variational Inference purposes...

→ We can get rid of $p(\mathcal{D})$ in the optimisation !

$$\begin{aligned} q^{\star} &= \operatorname{arginf}_{q \in \mathcal{Q}} D_{\alpha}(\mathbb{Q}||\mathbb{P}) \\ &= \operatorname{arginf}_{q \in \mathcal{Q}} \int_{\mathcal{Y}} f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) \quad \text{with } p(y) = p(y, \mathcal{D}) . \end{aligned}$$

[J. Hernandez-Lobato, Y. Li, M. Rowland, T. Bui, D. Hernandez-Lobato, and R. E Turner. (2016) Black-box alpha divergence minimization. ICML]

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

[A. Dieng, D. Tran, R. Ranganath, J. Paisley, and D. Blei. (2017) Variational inference via χ -upper bound minimization. NeurIPS]

Variational Inference within the α -divergence family (3)

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

where

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

- ❶ A **flexible** family of divergences...
- ❷ ...**suitable** for Variational Inference purposes...

$$q^\star = \operatorname{arginf}_{q \in \mathcal{Q}} \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) \quad \text{with } p(y) = p(y, \mathcal{D}) .$$

- ❸ ...with good **convexity** properties ! $\rightarrow f_\alpha$ is convex

[D. Wang, H. Liu Q. Liu (2018). Variational Inference with Tail-adaptive f-Divergence. NeurIPS]

Variational Inference within the α -divergence family (3)

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

where

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

- ❶ A **flexible** family of divergences...
- ❷ ...**suitable** for Variational Inference purposes...

$$q^\star = \operatorname{arginf}_{q \in \mathcal{Q}} \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) \quad \text{with } p(y) = p(y, \mathcal{D}) .$$

- ❸ ...with good **convexity** properties ! $\rightarrow f_\alpha$ is convex

[D. Wang, H. Liu Q. Liu (2018). Variational Inference with Tail-adaptive f-Divergence. NeurIPS]

Variational Inference within the α -divergence family (3)

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

where

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

- ❶ A **flexible** family of divergences...
- ❷ ...**suitable** for Variational Inference purposes...

$$q^{\star} = \operatorname{arginf}_{q \in \mathcal{Q}} \int_{\mathcal{Y}} f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) \quad \text{with } p(y) = p(y, \mathcal{D}) .$$

- ❸ ...with good **convexity** properties ! $\rightarrow f_{\alpha}$ is convex

[D. Wang, H. Liu Q. Liu (2018). Variational Inference with Tail-adaptive f-Divergence. NeurIPS]

Approximating family \mathcal{Q}

- Usually in Variational Inference : parametric family

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

- Recently : Hierarchical Variational Inference!

[R. Ranganath, D. Tran, and D. Blei. (2016) Hierarchical variational models. ICML]

[M. Yin and M. Zhou (2018). Semi-Implicit Variational Inference. ICML]

$$\left\{ y \mapsto \int_{\mathcal{T}} q_{\phi}(\theta) k_{\theta}(y) d\theta : \phi \in \mathcal{A} \right\} .$$

- What if... we consider a **broader** approximating family

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

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

\rightsquigarrow **Mixture models** : $\mu = \sum_{j=1}^J \lambda_j \delta_{\theta_j}$.

Approximating family \mathcal{Q}

- Usually in Variational Inference : parametric family

$$\{y \mapsto k_{\theta}(y) : \theta \in \mathsf{T}\} .$$

- Recently : Hierarchical Variational Inference!

[R. Ranganath, D. Tran, and D. Blei. (2016) Hierarchical variational models. ICML]

[M. Yin and M. Zhou (2018). Semi-Implicit Variational Inference. ICML]

$$\left\{ y \mapsto \int_{\mathsf{T}} q_{\phi}(\theta) k_{\theta}(y) d\theta : \phi \in \mathsf{A} \right\} .$$

- What if... we consider a **broader** approximating family

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

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

\rightsquigarrow **Mixture models** : $\mu = \sum_{j=1}^J \lambda_j \delta_{\theta_j}$.

Approximating family \mathcal{Q}

- Usually in Variational Inference : parametric family

$$\{y \mapsto k_{\theta}(y) : \theta \in \mathsf{T}\} .$$

- Recently : Hierarchical Variational Inference!

[R. Ranganath, D. Tran, and D. Blei. (2016) Hierarchical variational models. ICML]

[M. Yin and M. Zhou (2018). Semi-Implicit Variational Inference. ICML]

$$\left\{ y \mapsto \int_{\mathsf{T}} q_{\phi}(\theta) k_{\theta}(y) d\theta : \phi \in \mathsf{A} \right\} .$$

- What if... we consider a **broader** approximating family

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

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

\rightsquigarrow **Mixture models** : $\mu = \sum_{j=1}^J \lambda_j \delta_{\theta_j}$.

Approximating family \mathcal{Q}

- Usually in Variational Inference : parametric family

$$\{y \mapsto k_{\theta}(y) : \theta \in \mathsf{T}\} .$$

- Recently : Hierarchical Variational Inference!

[R. Ranganath, D. Tran, and D. Blei. (2016) Hierarchical variational models. ICML]

[M. Yin and M. Zhou (2018). Semi-Implicit Variational Inference. ICML]

$$\left\{ y \mapsto \int_{\mathsf{T}} q_{\phi}(\theta) k_{\theta}(y) d\theta : \phi \in \mathsf{A} \right\} .$$

- What if... we consider a **broader** approximating family

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

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

\rightsquigarrow **Mixture models** : $\mu = \sum_{j=1}^J \lambda_j \delta_{\theta_j}$.

Approximating family \mathcal{Q}

- Usually in Variational Inference : parametric family

$$\{y \mapsto k_{\theta}(y) : \theta \in \mathsf{T}\} .$$

- Recently : Hierarchical Variational Inference!

[R. Ranganath, D. Tran, and D. Blei. (2016) Hierarchical variational models. ICML]

[M. Yin and M. Zhou (2018). Semi-Implicit Variational Inference. ICML]

$$\left\{ y \mapsto \int_{\mathsf{T}} q_{\phi}(\theta) k_{\theta}(y) d\theta : \phi \in \mathsf{A} \right\} .$$

- What if... we consider a **broader** approximating family

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

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

\rightsquigarrow **Mixture models** : $\mu = \sum_{j=1}^J \lambda_j \delta_{\theta_j}$.

Our approach

- Let us consider the approximating family...

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

- and minimise the α -divergence w.r.t μ !

Optimisation problem

- $\mu k(y) = \int_{\mathsf{T}} \mu(d\theta) k(\theta, y)$, where $K : (\theta, A) \mapsto \int_A k(\theta, y) \nu(dy)$ is a Markov transition kernel on $\mathsf{T} \times \mathcal{Y}$ with kernel density k
- p : measurable positive function on $(\mathcal{Y}, \mathcal{Y})$

$$\operatorname{arginf}_{\mu \in \mathsf{M}} \underbrace{\int_{\mathcal{Y}} f_{\alpha} \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)}_{:= \Psi_{\alpha}(\mu)}$$

Our approach

- Let us consider the approximating family...

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

- and minimise the α -divergence w.r.t μ !

Optimisation problem

- $\mu k(y) = \int_{\mathsf{T}} \mu(d\theta) k(\theta, y)$, where $K : (\theta, A) \mapsto \int_A k(\theta, y) \nu(dy)$ is a Markov transition kernel on $\mathsf{T} \times \mathcal{Y}$ with kernel density k
- p : measurable positive function on $(\mathcal{Y}, \mathcal{Y})$

$$\operatorname{arginf}_{\mu \in \mathsf{M}} \underbrace{\int_{\mathcal{Y}} f_{\alpha} \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)}_{:= \Psi_{\alpha}(\mu)}$$

Outline

- 1 Background
- 2 The (α, Γ) -descent**
- 3 Numerical Experiments
- 4 Take-away message
- 5 Proof of the systematic decrease

The (α, Γ) -descent

Optimisation problem

$$\operatorname{arginf}_{\mu \in \mathcal{M}} \Psi_{\alpha}(\mu) \quad \text{with} \quad \Psi_{\alpha}(\mu) := \int_{\mathcal{Y}} f_{\alpha} \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)$$

Algorithm

Let $\mu_1 \in \mathcal{M}_1(\mathcal{T})$ such that $\Psi_{\alpha}(\mu_1) < \infty$. We define the sequence of probability measures $(\mu_n)_{n \in \mathbb{N}^*}$ iteratively by

$$\mu_{n+1} = \mathcal{I}_{\alpha}(\mu_n), \quad n \in \mathbb{N}^*. \quad (1)$$

Algorithm 1: *Exact (α, Γ) -descent one-step transition*

- ❶ Expectation step : $b_{\mu, \alpha}(\theta) = \int_{\mathcal{Y}} k(\theta, y) f'_{\alpha} \left(\frac{\mu k(y)}{p(y)} \right) \nu(dy)$
 - ❷ Iteration step : $\mathcal{I}_{\alpha}(\mu)(d\theta) = \frac{\mu(d\theta) \cdot \Gamma(b_{\mu, \alpha}(\theta) + \kappa)}{\mu(\Gamma(b_{\mu, \alpha} + \kappa))}$
-

Monotonicity

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

(A2) The function $\Gamma : \text{Dom}_{\alpha} \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Theorem 1

Assume (A1) and (A2). Let $\mu \in \mathsf{M}_1(\mathsf{T})$ be such that $\Psi_{\alpha}(\mu) < \infty$ and $\mu(\Gamma(b_{\mu, \alpha} + \kappa)) < \infty$. Then, the two following assertions hold.

- 1 We have $\Psi_{\alpha} \circ \mathcal{I}_{\alpha}(\mu) \leq \Psi_{\alpha}(\mu)$.
- 2 We have $\Psi_{\alpha} \circ \mathcal{I}_{\alpha}(\mu) = \Psi_{\alpha}(\mu)$ if and only if $\mu = \mathcal{I}_{\alpha}(\mu)$.

Monotonicity

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

(A2) The function $\Gamma : \text{Dom}_{\alpha} \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Theorem 1

Assume (A1) and (A2). Let $\mu \in \mathsf{M}_1(\mathsf{T})$ be such that $\Psi_{\alpha}(\mu) < \infty$ and $\mu(\Gamma(b_{\mu, \alpha} + \kappa)) < \infty$. Then, the two following assertions hold.

- 1 We have $\Psi_{\alpha} \circ \mathcal{I}_{\alpha}(\mu) \leq \Psi_{\alpha}(\mu)$.
- 2 We have $\Psi_{\alpha} \circ \mathcal{I}_{\alpha}(\mu) = \Psi_{\alpha}(\mu)$ if and only if $\mu = \mathcal{I}_{\alpha}(\mu)$.

Monotonicity

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

(A2) The function $\Gamma : \text{Dom}_{\alpha} \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Theorem 1

Assume (A1) and (A2). Let $\mu \in \mathsf{M}_1(\mathsf{T})$ be such that $\Psi_{\alpha}(\mu) < \infty$ and $\mu(\Gamma(b_{\mu, \alpha} + \kappa)) < \infty$. Then, the two following assertions hold.

- ❶ We have $\Psi_{\alpha} \circ \mathcal{I}_{\alpha}(\mu) \leq \Psi_{\alpha}(\mu)$.
- ❷ We have $\Psi_{\alpha} \circ \mathcal{I}_{\alpha}(\mu) = \Psi_{\alpha}(\mu)$ if and only if $\mu = \mathcal{I}_{\alpha}(\mu)$.

Monotonicity

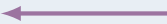
(A1) For all $(\theta, y) \in \mathsf{T} \times \mathsf{Y}$, $k(\theta, y) > 0$, $p(y) > 0$ and $\int_{\mathsf{Y}} p(y) \nu(dy) < \infty$.

(A2) The function $\Gamma : \text{Dom}_{\alpha} \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Theorem 1

Assume (A1) and (A2). Let $\mu \in \mathsf{M}_1(\mathsf{T})$ be such that $\Psi_{\alpha}(\mu) < \infty$ and $\mu(\Gamma(b_{\mu, \alpha} + \kappa)) < \infty$. Then, the two following assertions hold.

- ❶ We have $\Psi_{\alpha} \circ \mathcal{I}_{\alpha}(\mu) \leq \Psi_{\alpha}(\mu)$. 
- ❷ We have $\Psi_{\alpha} \circ \mathcal{I}_{\alpha}(\mu) = \Psi_{\alpha}(\mu)$ if and only if $\mu = \mathcal{I}_{\alpha}(\mu)$.

proof later !

Examples satisfying (A2)

(A2) The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

① Entropic MD : $\eta \in (0, 1]$, $\kappa \in \mathbb{R}$ and $\alpha = 1$

$$\Gamma(v) = e^{-\eta v} .$$

② Power descent : $\eta \in (0, 1]$, $(\alpha - 1)\kappa \geq 0$ and $\alpha \neq 1$

$$\Gamma(v) = [(\alpha - 1)v + 1]^{\frac{\eta}{1-\alpha}} .$$

Examples satisfying (A2)

(A2) The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

❶ Entropic MD : $\eta \in (0, 1]$, $\kappa \in \mathbb{R}$ and $\alpha = 1$

$$\Gamma(v) = e^{-\eta v} .$$

❷ Power descent : $\eta \in (0, 1]$, $(\alpha - 1)\kappa \geq 0$ and $\alpha \neq 1$

$$\Gamma(v) = [(\alpha - 1)v + 1]^{\frac{\eta}{1-\alpha}} .$$

Examples satisfying (A2)

(A2) The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

❶ Entropic MD : $\eta \in (0, 1]$, $\kappa \in \mathbb{R}$ and $\alpha = 1$

$$\Gamma(v) = e^{-\eta v} .$$

❷ Power descent : $\eta \in (0, 1]$, $(\alpha - 1)\kappa \geq 0$ and $\alpha \neq 1$

$$\Gamma(v) = [(\alpha - 1)v + 1]^{\frac{\eta}{1-\alpha}} .$$

Limiting behavior

Table 1: Examples of allowed (Γ, κ) in the (α, Γ) -descent

Divergence considered	Possible choice of (Γ, κ)	
<i>Forward KL</i> ($\alpha = 1$)	$\Gamma(v) = e^{-\eta v}, \eta \in (0, 1)$	any κ
α -divergence with $\alpha \in \mathbb{R} \setminus \{1\}$	$\Gamma(v) = e^{-\eta v}, \eta \in (0, \frac{1}{ \alpha-1 b _{\infty, \alpha+1}})$	any κ
	$\alpha > 1, \Gamma(v) = [(\alpha - 1)v + 1]^{\frac{\eta}{1-\alpha}}, \eta \in (0, 1]$	$\kappa > 0$
	$\alpha < 1, \Gamma(v) = [(\alpha - 1)v + 1]^{\frac{\eta}{1-\alpha}}, \eta \in (0, 1]$	$\kappa \leq 0$

→ Convergence towards the optimum value at a $O(1/N)$ rate

→ Convergence towards the optimum value

Limiting behavior

Table 1: Examples of allowed (Γ, κ) in the (α, Γ) -descent

Divergence considered	Possible choice of (Γ, κ)	
Forward KL ($\alpha = 1$)	$\Gamma(v) = e^{-\eta v}, \eta \in (0, 1)$	any κ
α -divergence with $\alpha \in \mathbb{R} \setminus \{1\}$	$\Gamma(v) = e^{-\eta v}, \eta \in (0, \frac{1}{ \alpha-1 b _{\infty, \alpha+1}})$	any κ
	$\alpha > 1, \Gamma(v) = [(\alpha - 1)v + 1]^{\frac{\eta}{1-\alpha}}, \eta \in (0, 1]$	$\kappa > 0$
	$\alpha < 1, \Gamma(v) = [(\alpha - 1)v + 1]^{\frac{\eta}{1-\alpha}}, \eta \in (0, 1]$	$\kappa \leq 0$

→ Convergence towards the optimum value at a $O(1/N)$ rate

→ Convergence towards the optimum value

Mixture models and (α, Γ) -descent

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

Let $\theta_1, \dots, \theta_J \in \mathcal{T}$ **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, $\mu_n = \underbrace{\mathcal{I}_{\alpha} \circ \dots \circ \mathcal{I}_{\alpha}}_{n \text{ times}}(\mu_{\boldsymbol{\lambda}})$ is of the form $\mu_n = \sum_{j=1}^J \lambda_{j,n} \delta_{\theta_j}$ with

$$\begin{cases} \lambda_1 = \boldsymbol{\lambda} \\ \lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}. \end{cases} \quad (2)$$

- In practice, we will use

$$\hat{b}_{\mu_n, \alpha, M}(\theta_j) = \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_j, Y_{m,n})}{\mu_n k(Y_{m,n})} f'_{\alpha} \left(\frac{\mu_n k(Y_{m,n})}{p(Y_{m,n})} \right),$$

with $Y_{1,n}, \dots, Y_{M,n}$ drawn independently from $\mu_n k$.

- **Exploitation step** which requires no information on the distribution of $\{\theta_1, \dots, \theta_J\}$ (as opposed to Importance Sampling)

Mixture models and (α, Γ) -descent

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

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

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

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

$$\begin{cases} \lambda_1 = \boldsymbol{\lambda} \\ \lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}. \end{cases} \quad (2)$$

- In practice, we will use

$$\hat{b}_{\mu_n, \alpha, M}(\theta_j) = \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_j, Y_{m,n})}{\mu_n k(Y_{m,n})} f'_{\alpha} \left(\frac{\mu_n k(Y_{m,n})}{p(Y_{m,n})} \right),$$

with $Y_{1,n}, \dots, Y_{M,n}$ drawn independently from $\mu_n k$.

- **Exploitation step** which requires no information on the distribution of $\{\theta_1, \dots, \theta_J\}$ (as opposed to Importance Sampling)

Mixture models and (α, Γ) -descent

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

Let $\theta_1, \dots, \theta_J \in \mathcal{T}$ **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, $\mu_n = \underbrace{\mathcal{I}_{\alpha} \circ \dots \circ \mathcal{I}_{\alpha}}_{n \text{ times}}(\mu_{\boldsymbol{\lambda}})$ is of the form $\mu_n = \sum_{j=1}^J \lambda_{j,n} \delta_{\theta_j}$ with

$$\begin{cases} \lambda_1 = \boldsymbol{\lambda} \\ \lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}. \end{cases} \quad (2)$$

- In practice, we will use

$$\hat{b}_{\mu_n, \alpha, M}(\theta_j) = \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_j, Y_{m,n})}{\mu_n k(Y_{m,n})} f'_{\alpha} \left(\frac{\mu_n k(Y_{m,n})}{p(Y_{m,n})} \right),$$

with $Y_{1,n}, \dots, Y_{M,n}$ drawn independently from $\mu_n k$.

- **Exploitation step** which requires no information on the distribution of $\{\theta_1, \dots, \theta_J\}$ (as opposed to Importance Sampling)

Outline

- 1 Background
- 2 The (α, Γ) -descent
- 3 Numerical Experiments**
- 4 Take-away message
- 5 Proof of the systematic decrease

Numerical Experiments

- Framework

Kernel: Gaussian transition kernel k_h with bandwidth h .

$$\left\{ y \mapsto \mu_{\lambda} k_h(y) = \sum_{j=1}^J \lambda_j k_h(y - \theta_j) : \lambda \in \mathcal{S}_J, (\theta_j)_{1 \leq j \leq J} \in \mathbb{T}^J \right\}.$$

At time t ,

- ① **Exploitation step** Optimise λ using the (α, Γ) -descent.
 - ② **Exploration step** Sample $(\theta_{j,t+1})_{1 \leq j \leq J_{t+1}}$ according to $\mu_{\lambda} k_{h_t}$, with $h_t \propto J_t^{-1/(4+d)}$, where d is the dimension of the latent space.
- Toy example
 $p(y) = Z \times [0.5\mathcal{N}(y; -s\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(y; s\mathbf{u}_d, \mathbf{I}_d)]$, $Z = 2$, $s = 2$
 - Bayesian Logistic Regression
Covertypes dataset (581,012 data points and 54 features)

Numerical Experiments

- Framework

Kernel: Gaussian transition kernel k_h with bandwidth h .

$$\left\{ y \mapsto \mu_{\lambda} k_h(y) = \sum_{j=1}^J \lambda_j k_h(y - \theta_j) : \lambda \in \mathcal{S}_J, (\theta_j)_{1 \leq j \leq J} \in \mathbb{T}^J \right\}.$$

At time t ,

- ① **Exploitation step** Optimise λ using the (α, Γ) -descent.
 - ② **Exploration step** Sample $(\theta_{j,t+1})_{1 \leq j \leq J_{t+1}}$ according to $\mu_{\lambda} k_{h_t}$, with $h_t \propto J_t^{-1/(4+d)}$, where d is the dimension of the latent space.
- Toy example
 $p(y) = Z \times [0.5\mathcal{N}(y; -s\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(y; s\mathbf{u}_d, \mathbf{I}_d)]$, $Z = 2$, $s = 2$
 - Bayesian Logistic Regression
Covertypes dataset (581,012 data points and 54 features)

Numerical Experiments

- Framework

Kernel: Gaussian transition kernel k_h with bandwidth h .

$$\left\{ y \mapsto \mu_{\lambda} k_h(y) = \sum_{j=1}^J \lambda_j k_h(y - \theta_j) : \lambda \in \mathcal{S}_J, (\theta_j)_{1 \leq j \leq J} \in \mathbb{T}^J \right\}.$$

At time t ,

- ❶ **Exploitation step** Optimise λ using the (α, Γ) -descent.
- ❷ **Exploration step** Sample $(\theta_{j,t+1})_{1 \leq j \leq J_{t+1}}$ according to $\mu_{\lambda} k_{h_t}$, with $h_t \propto J_t^{-1/(4+d)}$, where d is the dimension of the latent space.

- Toy example

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

- Bayesian Logistic Regression

Covertypes dataset (581,012 data points and 54 features)

Numerical Experiments

- Framework

Kernel: Gaussian transition kernel k_h with bandwidth h .

$$\left\{ y \mapsto \mu_{\lambda} k_h(y) = \sum_{j=1}^J \lambda_j k_h(y - \theta_j) : \lambda \in \mathcal{S}_J, (\theta_j)_{1 \leq j \leq J} \in \mathbb{T}^J \right\}.$$

At time t ,

- ❶ **Exploitation step** Optimise λ using the (α, Γ) -descent.
 - ❷ **Exploration step** Sample $(\theta_{j,t+1})_{1 \leq j \leq J_{t+1}}$ according to $\mu_{\lambda} k_{h_t}$, with $h_t \propto J_t^{-1/(4+d)}$, where d is the dimension of the latent space.
- Toy example
 $p(y) = Z \times [0.5\mathcal{N}(\mathbf{y}; -s\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(\mathbf{y}; s\mathbf{u}_d, \mathbf{I}_d)]$, $Z = 2$, $s = 2$
 - Bayesian Logistic Regression
Covertypes dataset (581,012 data points and 54 features)

Numerical Experiments

- Framework

Kernel: Gaussian transition kernel k_h with bandwidth h .

$$\left\{ y \mapsto \mu_{\lambda} k_h(y) = \sum_{j=1}^J \lambda_j k_h(y - \theta_j) : \lambda \in \mathcal{S}_J, (\theta_j)_{1 \leq j \leq J} \in \mathbb{T}^J \right\}.$$

At time t ,

- ① **Exploitation step** Optimise λ using the (α, Γ) -descent.
 - ② **Exploration step** Sample $(\theta_{j,t+1})_{1 \leq j \leq J_{t+1}}$ according to $\mu_{\lambda} k_{h_t}$, with $h_t \propto J_t^{-1/(4+d)}$, where d is the dimension of the latent space.
- Toy example
 $p(y) = Z \times [0.5\mathcal{N}(\mathbf{y}; -s\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(\mathbf{y}; s\mathbf{u}_d, \mathbf{I}_d)]$, $Z = 2$, $s = 2$
 - Bayesian Logistic Regression
Covertypes dataset (581,012 data points and 54 features)

Toy Example : Mirror Descent vs Power Descent

We compare :

- 0.5-Mirror descent : $\Gamma(v) = e^{-\eta v}$ with $\alpha = 0.5$,
- 0.5-Power descent : $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$ with $\alpha = 0.5$.

$J = M = 100$, initial weights: $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$.

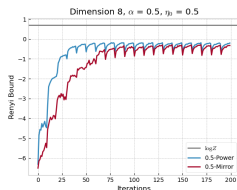
Toy Example : Mirror Descent vs Power Descent

We compare :

- 0.5-Mirror descent : $\Gamma(v) = e^{-\eta v}$ with $\alpha = 0.5$,
- 0.5-Power descent : $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$ with $\alpha = 0.5$.

$J = M = 100$, initial weights: $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$.

Figure: Average Renyi-Bound for the 0.5-Power and 0.5-Mirror descent computed over 100 replicates with $\eta_0 = 0.5$.



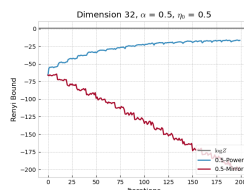
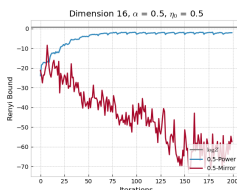
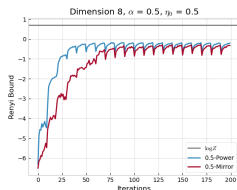
Toy Example : Mirror Descent vs Power Descent

We compare :

- 0.5-Mirror descent : $\Gamma(v) = e^{-\eta v}$ with $\alpha = 0.5$,
- 0.5-Power descent : $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$ with $\alpha = 0.5$.

$J = M = 100$, initial weights: $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$.

Figure: Average Renyi-Bound for the 0.5-Power and 0.5-Mirror descent computed over 100 replicates with $\eta_0 = 0.5$.



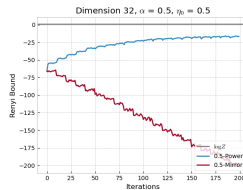
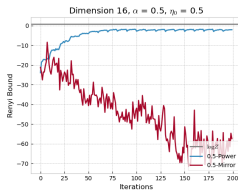
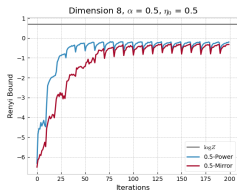
Toy Example : Mirror Descent vs Power Descent

We compare :

- 0.5-Mirror descent : $\Gamma(v) = e^{-\eta v}$ with $\alpha = 0.5$,
- 0.5-Power descent : $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$ with $\alpha = 0.5$.

$J = M = 100$, initial weights: $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$.

Figure: Average Renyi-Bound for the 0.5-Power and 0.5-Mirror descent computed over 100 replicates with $\eta_0 = 0.5$.



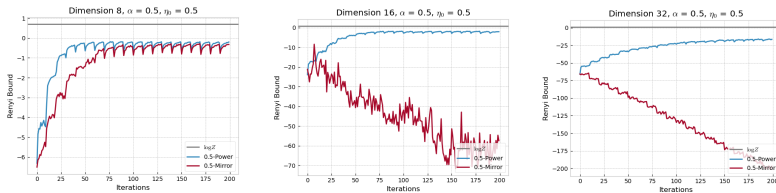
Toy Example : Mirror Descent vs Power Descent

We compare :

- 0.5-Mirror descent : $\Gamma(v) = e^{-\eta v}$ with $\alpha = 0.5$,
- 0.5-Power descent : $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$ with $\alpha = 0.5$.

$J = M = 100$, initial weights: $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$.

Figure: Average Renyi-Bound for the 0.5-Power and 0.5-Mirror descent computed over 100 replicates with $\eta_0 = 0.5$.



$$\begin{aligned} \text{Mirror} \quad \lambda_{j,n} &\propto \exp \left(\frac{\eta}{1-\alpha} ((\alpha-1)b_{\mu_{\lambda_n}, \alpha}(\theta_j) + (\alpha-1)\kappa) \right) \\ \text{Power} \quad \lambda_{j,n} &\propto \exp \left(\frac{\eta}{1-\alpha} \log ((\alpha-1)b_{\mu_{\lambda_n}, \alpha}(\theta_j) + (\alpha-1)\kappa) \right). \end{aligned}$$

Toy Example : $\alpha = 1$

We compare:

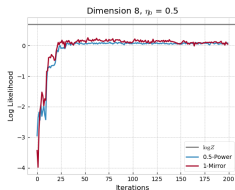
- 1-Mirror descent : $\Gamma(v) = e^{-\eta v}$ with $\alpha = 1$,
- 0.5-Power descent : $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$ with $\alpha = 0.5$.

Toy Example : $\alpha = 1$

We compare:

- 1-Mirror descent : $\Gamma(v) = e^{-\eta v}$ with $\alpha = 1$,
- 0.5-Power descent : $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$ with $\alpha = 0.5$.

Figure: Plotted is the average Log-likelihood for 0.5-Power and 1-Mirror descent in dimension $d = \{8, 16, 32\}$ computed over 100 replicates with $\eta_0 = 0.5$.

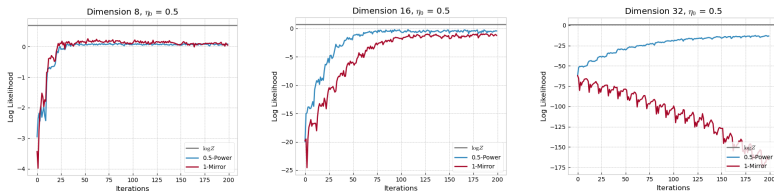


Toy Example : $\alpha = 1$

We compare:

- 1-Mirror descent : $\Gamma(v) = e^{-\eta v}$ with $\alpha = 1$,
- 0.5-Power descent : $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$ with $\alpha = 0.5$.

Figure: Plotted is the average Log-likelihood for 0.5-Power and 1-Mirror descent in dimension $d = \{8, 16, 32\}$ computed over 100 replicates with $\eta_0 = 0.5$.



Outline

- 1 Background
- 2 The (α, Γ) -descent
- 3 Numerical Experiments
- 4 Take-away message**
- 5 Proof of the systematic decrease

Take-away message

The (α, Γ) -descent

- performs an update of probability measures
 - sufficient conditions on (α, Γ) leading to a systematic decrease
 - includes Entropic Mirror Descent
 - convergence to an optimum and $O(1/N)$ convergence rates,
- can be applied to density approximation
 - handles the case of Mixture Models for any kernel K
 - requires no information on the distribution of $\{\theta_1, \dots, \theta_J\}$
 - empirical benefit of using the Power descent.

[Kamélia Daudel, Randal Douc and François Portier (2020). Infinite-dimensional gradient-based descent for alpha-divergence minimisation. To be published in the Annals of Statistics. <https://arxiv.org/abs/2005.10618>]

Take-away message

The (α, Γ) -descent

- performs an update of probability measures
 - sufficient conditions on (α, Γ) leading to a systematic decrease
 - includes Entropic Mirror Descent
 - convergence to an optimum and $O(1/N)$ convergence rates,
- can be applied to density approximation
 - handles the case of Mixture Models for any kernel K
 - requires no information on the distribution of $\{\theta_1, \dots, \theta_J\}$
 - empirical benefit of using the Power descent.

[Kamélia Daudel, Randal Douc and François Portier (2020). Infinite-dimensional gradient-based descent for alpha-divergence minimisation. To be published in the Annals of Statistics. <https://arxiv.org/abs/2005.10618>]

Take-away message

The (α, Γ) -descent

- performs an update of probability measures
 - sufficient conditions on (α, Γ) leading to a systematic decrease
 - includes Entropic Mirror Descent
 - convergence to an optimum and $O(1/N)$ convergence rates,
- can be applied to density approximation
 - handles the case of Mixture Models for any kernel K
 - requires no information on the distribution of $\{\theta_1, \dots, \theta_J\}$
 - empirical benefit of using the Power descent.

[Kamélia Daudel, Randal Douc and François Portier (2020). Infinite-dimensional gradient-based descent for alpha-divergence minimisation. To be published in the Annals of Statistics. <https://arxiv.org/abs/2005.10618>]

Outline

- 1 Background
- 2 The (α, Γ) -descent
- 3 Numerical Experiments
- 4 Take-away message
- 5 Proof of the systematic decrease**

The result we want to prove

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

(A2) The function $\Gamma : \text{Dom}_{\alpha} \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Theorem 1

Assume (A1) and (A2). Let $\mu \in \mathsf{M}_1(\mathsf{T})$ be such that $\Psi_{\alpha}(\mu) < \infty$ and $\mu(\Gamma(b_{\mu, \alpha} + \kappa)) < \infty$. Then, the two following assertions hold.

- ❶ We have $\Psi_{\alpha} \circ \mathcal{I}_{\alpha}(\mu) \leq \Psi_{\alpha}(\mu)$.
- ❷ We have $\Psi_{\alpha} \circ \mathcal{I}_{\alpha}(\mu) = \Psi_{\alpha}(\mu)$ if and only if $\mu = \mathcal{I}_{\alpha}(\mu)$.

Recall that :

$$\Psi_{\alpha}(\mu) = \int_{\mathsf{Y}} f_{\alpha} \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)$$
$$b_{\mu, \alpha}(\theta) = \int_{\mathsf{Y}} k(\theta, y) f'_{\alpha} \left(\frac{\mu k(y)}{p(y)} \right) \nu(dy)$$
$$\mathcal{I}_{\alpha}(\mu)(d\theta) = \frac{\mu(d\theta) \cdot \Gamma(b_{\mu, \alpha}(\theta) + \kappa)}{\mu(\Gamma(b_{\mu, \alpha} + \kappa))}$$

Step 1 : Proving a general lower bound (1)

Let $\mu, \zeta \in \mathcal{M}_1(\mathcal{T})$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where
$$A_\alpha := \int_{\mathcal{Y}} \nu(dy) \int_{\mathcal{T}} \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$$

Equality holds iif $\zeta = \mu$.

→ By definition $\Psi_\alpha(\mu) = \int_{\mathcal{Y}} f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)$ with f_α convex.

→ By convexity of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{\mu k(y)}{p(y)} [1 - g(\theta)].$$

→ Now integrating first w.r.t to $\frac{\mu(d\theta)k(\theta, y)}{\mu k(y)}$,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq \int_{\mathcal{T}} \frac{\mu(d\theta)k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + \int_{\mathcal{T}} \mu(d\theta)k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{1}{p(y)} [1 - g(\theta)]$$

then w.r.t to $\nu(dy)p(y)$, we deduce

$$\Psi_\alpha(\mu) \geq \int_{\mathcal{Y}} p(y) \nu(dy) \int_{\mathcal{T}} \frac{\mu(d\theta)k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Step 1 : Proving a general lower bound (1)

Let $\mu, \zeta \in M_1(T)$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where
$$A_\alpha := \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$$

Equality holds iif $\zeta = \mu$.

→ By definition $\Psi_\alpha(\mu) = \int_Y f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)$ with f_α **convex**.

→ By convexity of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{\mu k(y)}{p(y)} [1 - g(\theta)] .$$

→ Now integrating first w.r.t to $\frac{\mu(d\theta)k(\theta, y)}{\mu k(y)}$,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq \int_T \frac{\mu(d\theta)k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + \int_T \mu(d\theta)k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{1}{p(y)} [1 - g(\theta)]$$

then w.r.t to $\nu(dy)p(y)$, we deduce

$$\Psi_\alpha(\mu) \geq \int_Y p(y) \nu(dy) \int_T \frac{\mu(d\theta)k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Step 1 : Proving a general lower bound (1)

Let $\mu, \zeta \in \mathcal{M}_1(\mathcal{T})$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where
$$A_\alpha := \int_{\mathcal{Y}} \nu(dy) \int_{\mathcal{T}} \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$$

Equality holds iif $\zeta = \mu$.

→ By definition $\Psi_\alpha(\mu) = \int_{\mathcal{Y}} f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)$ with f_α **convex**.

→ By convexity of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{\mu k(y)}{p(y)} [1 - g(\theta)] .$$

→ Now integrating first w.r.t to $\frac{\mu(d\theta)k(\theta,y)}{\mu k(y)}$,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq \int_{\mathcal{T}} \frac{\mu(d\theta)k(\theta,y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + \int_{\mathcal{T}} \mu(d\theta)k(\theta,y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{1}{p(y)} [1 - g(\theta)]$$

then w.r.t to $\nu(dy)p(y)$, we deduce

$$\Psi_\alpha(\mu) \geq \int_{\mathcal{Y}} p(y) \nu(dy) \int_{\mathcal{T}} \frac{\mu(d\theta)k(\theta,y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Step 1 : Proving a general lower bound (1)

Let $\mu, \zeta \in M_1(T)$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where $A_\alpha := \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$

Equality holds iif $\zeta = \mu$.

→ By definition $\Psi_\alpha(\mu) = \int_Y f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)$ with f_α **convex**.

→ By convexity of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{\mu k(y)}{p(y)} [1 - g(\theta)] .$$

→ Now integrating first w.r.t to $\frac{\mu(d\theta)k(\theta, y)}{\mu k(y)}$,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq \int_T \frac{\mu(d\theta)k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + \int_T \mu(d\theta)k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{1}{p(y)} [1 - g(\theta)]$$

then w.r.t to $\nu(dy)p(y)$, we deduce

$$\Psi_\alpha(\mu) \geq \int_Y p(y) \nu(dy) \int_T \frac{\mu(d\theta)k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Step 1 : Proving a general lower bound (1)

Let $\mu, \zeta \in M_1(T)$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where $A_\alpha := \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$

Equality holds iif $\zeta = \mu$.

→ By definition $\Psi_\alpha(\mu) = \int_Y f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)$ with f_α **convex**.

→ By convexity of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{\mu k(y)}{p(y)} [1 - g(\theta)] .$$

→ Now integrating first w.r.t to $\frac{\mu(d\theta)k(\theta, y)}{\mu k(y)}$,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq \int_T \frac{\mu(d\theta)k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{1}{p(y)} [1 - g(\theta)]$$

then w.r.t to $\nu(dy)p(y)$, we deduce

$$\Psi_\alpha(\mu) \geq \int_Y p(y) \nu(dy) \int_T \frac{\mu(d\theta)k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Step 1 : Proving a general lower bound (1)

Let $\mu, \zeta \in M_1(T)$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where
$$A_\alpha := \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$$

Equality holds iif $\zeta = \mu$.

→ By definition $\Psi_\alpha(\mu) = \int_Y f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)$ with f_α **convex**.

→ By convexity of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{\mu k(y)}{p(y)} [1 - g(\theta)] .$$

→ Now integrating first w.r.t to $\frac{\mu(d\theta)k(\theta,y)}{\mu k(y)}$,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq \int_T \frac{\mu(d\theta)k(\theta,y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{1}{p(y)} [1 - g(\theta)]$$

then w.r.t to $\nu(dy)p(y)$, we deduce

$$\Psi_\alpha(\mu) \geq \int_Y p(y) \nu(dy) \int_T \frac{\mu(d\theta)k(\theta,y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Step 1 : Proving a general lower bound (1)

Let $\mu, \zeta \in M_1(T)$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where $A_\alpha := \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$

Equality holds iif $\zeta = \mu$.

→ By definition $\Psi_\alpha(\mu) = \int_Y f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) p(y) \nu(dy)$ with f_α **convex**.

→ By convexity of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{\mu k(y)}{p(y)} [1 - g(\theta)] .$$

→ Now integrating first w.r.t to $\frac{\mu(d\theta)k(\theta,y)}{\mu k(y)}$,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq \int_T \frac{\mu(d\theta)k(\theta,y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) \frac{1}{p(y)} [1 - g(\theta)]$$

then w.r.t to $\nu(dy)p(y)$, we deduce

$$\Psi_\alpha(\mu) \geq \int_Y p(y) \nu(dy) \int_T \frac{\mu(d\theta)k(\theta,y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Step 1 : Proving a general lower bound (2)

Let $\mu, \zeta \in M_1(\mathcal{T})$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where $A_\alpha := \int_{\mathcal{Y}} \nu(dy) \int_{\mathcal{T}} \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$.

Equality holds iif $\zeta = \mu$.

At this stage,

$$\Psi_\alpha(\mu) \geq \int_{\mathcal{Y}} p(y) \nu(dy) \int_{\mathcal{T}} \frac{\mu(d\theta) k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Finally, applying Jensen's inequality to the convex function f_α

$$\Psi_\alpha(\mu) \geq \int_{\mathcal{Y}} p(y) \nu(dy) f_\alpha \left(\int_{\mathcal{T}} \frac{\mu(d\theta) k(\theta, y)}{\mu k(y)} \frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

that is

$$\Psi_\alpha(\mu) \geq \Psi_\alpha(\zeta) + A_\alpha$$

Step 1 : Proving a general lower bound (2)

Let $\mu, \zeta \in M_1(\mathcal{T})$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where $A_\alpha := \int_{\mathcal{Y}} \nu(dy) \int_{\mathcal{T}} \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$.

Equality holds iff $\zeta = \mu$.

At this stage,

$$\Psi_\alpha(\mu) \geq \int_{\mathcal{Y}} p(y) \nu(dy) \int_{\mathcal{T}} \frac{\mu(d\theta) k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Finally, applying Jensen's inequality to the convex function f_α

$$\Psi_\alpha(\mu) \geq \int_{\mathcal{Y}} p(y) \nu(dy) f_\alpha \left(\int_{\mathcal{T}} \frac{\mu(d\theta) k(\theta, y)}{\mu k(y)} \frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

that is

$$\Psi_\alpha(\mu) \geq \Psi_\alpha(\zeta) + A_\alpha$$

Step 1 : Proving a general lower bound (2)

Let $\mu, \zeta \in M_1(\mathcal{T})$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where $A_\alpha := \int_{\mathcal{Y}} \nu(dy) \int_{\mathcal{T}} \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$.

Equality holds iff $\zeta = \mu$.

At this stage,

$$\Psi_\alpha(\mu) \geq \int_{\mathcal{Y}} p(y) \nu(dy) \int_{\mathcal{T}} \frac{\mu(d\theta) k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Finally, applying Jensen's inequality to the convex function f_α

$$\Psi_\alpha(\mu) \geq \int_{\mathcal{Y}} p(y) \nu(dy) f_\alpha \left(\int_{\mathcal{T}} \frac{\mu(d\theta) k(\theta, y)}{\mu k(y)} \frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

that is

$$\Psi_\alpha(\mu) \geq \Psi_\alpha(\zeta) + A_\alpha$$

Step 1 : Proving a general lower bound (2)

Let $\mu, \zeta \in M_1(\mathsf{T})$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

where $A_\alpha := \int_{\mathsf{Y}} \nu(dy) \int_{\mathsf{T}} \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$.

Equality holds iif $\zeta = \mu$.

At this stage,

$$\Psi_\alpha(\mu) \geq \int_{\mathsf{Y}} p(y) \nu(dy) \int_{\mathsf{T}} \frac{\mu(d\theta) k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Finally, applying Jensen's inequality to the convex function f_α

$$\Psi_\alpha(\mu) \geq \int_{\mathsf{Y}} p(y) \nu(dy) f_\alpha \left(\int_{\mathsf{T}} \frac{\mu(d\theta) k(\theta, y)}{\mu k(y)} \frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

that is

$$\Psi_\alpha(\mu) \geq \Psi_\alpha(\zeta) + A_\alpha$$

Step 1 : Proving a general lower bound (2)

Let $\mu, \zeta \in M_1(\mathsf{T})$ s.t $\zeta \preceq \mu$ and $\Psi_\alpha(\mu) < \infty$. Denote by g the density of ζ w.r.t μ .

We have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta)$$

$$\text{where } A_\alpha := \int_{\mathsf{Y}} \nu(dy) \int_{\mathsf{T}} \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)] .$$

Equality holds iif $\zeta = \mu$.

At this stage,

$$\Psi_\alpha(\mu) \geq \int_{\mathsf{Y}} p(y) \nu(dy) \int_{\mathsf{T}} \frac{\mu(d\theta) k(\theta, y)}{\mu k(y)} f_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

Finally, applying Jensen's inequality to the convex function f_α

$$\Psi_\alpha(\mu) \geq \int_{\mathsf{Y}} p(y) \nu(dy) f_\alpha \left(\int_{\mathsf{T}} \frac{\mu(d\theta) k(\theta, y)}{\mu k(y)} \frac{g(\theta) \mu k(y)}{p(y)} \right) + A_\alpha$$

that is

$$\Psi_\alpha(\mu) \geq \Psi_\alpha(\zeta) + A_\alpha$$

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (1)

Recall that :

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

$$\mathcal{I}_\alpha(\mu)(d\theta) = \frac{\mu(d\theta) \cdot \Gamma(b_{\mu,\alpha}(\theta) + \kappa)}{\mu(\Gamma(b_{\mu,\alpha} + \kappa))}$$

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu)$$

where $A_\alpha := \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$.

The proof is complete if we prove that $A_\alpha \geq 0$.

→ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$. In this case $f'_\alpha(u) = \frac{1}{\alpha-1} [u^{\alpha-1} - 1]$ and

$$b_{\mu,\alpha}(\theta) = \int_Y k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] \nu(dy)$$

$$A_\alpha = \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{g(\theta) \mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] [1 - g(\theta)]$$

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (1)

Recall that :

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

$$\mathcal{I}_\alpha(\mu)(d\theta) = \frac{\mu(d\theta) \cdot \Gamma(b_{\mu,\alpha}(\theta) + \kappa)}{\mu(\Gamma(b_{\mu,\alpha} + \kappa))}$$

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu)$$

where $A_\alpha := \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$.

The proof is complete if we prove that $A_\alpha \geq 0$.

→ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$. In this case $f'_\alpha(u) = \frac{1}{\alpha-1}[u^{\alpha-1} - 1]$ and

$$b_{\mu,\alpha}(\theta) = \int_Y k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] \nu(dy)$$

$$A_\alpha = \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{g(\theta) \mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] [1 - g(\theta)]$$

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (1)

Recall that :

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

$$\mathcal{I}_\alpha(\mu)(d\theta) = \frac{\mu(d\theta) \cdot \Gamma(b_{\mu,\alpha}(\theta) + \kappa)}{\mu(\Gamma(b_{\mu,\alpha} + \kappa))}$$

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu)$$

where $A_\alpha := \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta) \mu k(y)}{p(y)} \right) [1 - g(\theta)]$.

The proof is complete if we prove that $A_\alpha \geq 0$.

→ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$. In this case $f'_\alpha(u) = \frac{1}{\alpha-1} [u^{\alpha-1} - 1]$ and

$$b_{\mu,\alpha}(\theta) = \int_Y k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] \nu(dy)$$

$$A_\alpha = \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{g(\theta) \mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] [1 - g(\theta)]$$

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (2)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

→ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$. In this case $f'_\alpha(u) = \frac{1}{\alpha-1}[u^{\alpha-1} - 1]$ and

$$b_{\mu,\alpha}(\theta) = \int_Y k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] \nu(dy)$$

$$\begin{aligned} A_\alpha &= \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{g(\theta) \mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] [1 - g(\theta)] \\ &= \int_T \mu(d\theta) \left(\int_Y \nu(dy) k(\theta, y) \frac{1}{\alpha-1} \left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} g(\theta)^{\alpha-1} - 1 \right) [1 - g(\theta)] \\ &= \int_T \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha-1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)] \end{aligned}$$

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (2)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

→ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$. In this case $f'_\alpha(u) = \frac{1}{\alpha-1}[u^{\alpha-1} - 1]$ and

$$b_{\mu,\alpha}(\theta) = \int_Y k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] \nu(dy)$$

$$\begin{aligned} A_\alpha &= \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{g(\theta) \mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] [1 - g(\theta)] \\ &= \int_T \mu(d\theta) \left(\int_Y \nu(dy) k(\theta, y) \frac{1}{\alpha-1} \left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} g(\theta)^{\alpha-1} - 1 \right) [1 - g(\theta)] \\ &= \int_T \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha-1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)] \end{aligned}$$

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (2)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

→ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$. In this case $f'_\alpha(u) = \frac{1}{\alpha-1}[u^{\alpha-1} - 1]$ and

$$b_{\mu,\alpha}(\theta) = \int_Y k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] \nu(dy)$$

$$\begin{aligned} A_\alpha &= \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{g(\theta) \mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] [1 - g(\theta)] \\ &= \int_T \mu(d\theta) \left(\int_Y \nu(dy) k(\theta, y) \frac{1}{\alpha-1} \left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} g(\theta)^{\alpha-1} \right) [1 - g(\theta)] \\ &= \int_T \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha-1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)] \end{aligned}$$

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (2)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu).$$

The proof is complete if we prove that $A_\alpha \geq 0$.

→ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$. In this case $f'_\alpha(u) = \frac{1}{\alpha-1}[u^{\alpha-1} - 1]$ and

$$b_{\mu,\alpha}(\theta) = \int_Y k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] \nu(dy)$$

$$\begin{aligned} A_\alpha &= \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{g(\theta) \mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] [1 - g(\theta)] \\ &= \int_T \mu(d\theta) \left(\int_Y \nu(dy) k(\theta, y) \frac{1}{\alpha-1} \left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} g(\theta)^{\alpha-1} - 1 \right) [1 - g(\theta)] \\ &= \int_T \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha-1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)] \end{aligned}$$

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (2)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

→ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$. In this case $f'_\alpha(u) = \frac{1}{\alpha-1}[u^{\alpha-1} - 1]$ and

$$b_{\mu,\alpha}(\theta) = \int_Y k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] \nu(dy)$$

$$\begin{aligned} A_\alpha &= \int_Y \nu(dy) \int_T \mu(d\theta) k(\theta, y) \frac{1}{\alpha-1} \left[\left(\frac{g(\theta) \mu k(y)}{p(y)} \right)^{\alpha-1} - 1 \right] [1 - g(\theta)] \\ &= \int_T \mu(d\theta) \left(\int_Y \nu(dy) k(\theta, y) \frac{1}{\alpha-1} \left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} g(\theta)^{\alpha-1} - 1 \right) [1 - g(\theta)] \\ &= \int_T \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha-1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)] \end{aligned}$$

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (3)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

At this stage,

$$A_\alpha = \int_{\mathcal{T}} \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha - 1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)]$$

On the probability space $(\mathcal{T}, \mathcal{T}, \mu)$, we let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$. Set $\tilde{\Gamma}(v) = \Gamma(v)/\mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$. Then, $\mathbb{E}[1 - \tilde{\Gamma}(V)] = 0$ and

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \end{aligned}$$

Time to recall (A2)! The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Conclusion: $A_\alpha \geq 0$!

□

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (3)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

At this stage,

$$A_\alpha = \int_{\mathcal{T}} \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha - 1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)]$$

On the probability space $(\mathcal{T}, \mathcal{T}, \mu)$, we let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$. Set $\tilde{\Gamma}(v) = \Gamma(v)/\mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$. Then, $\mathbb{E}[1 - \tilde{\Gamma}(V)] = 0$ and

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \end{aligned}$$

Time to recall (A2)! The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Conclusion: $A_\alpha \geq 0$!

□

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (3)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

At this stage,

$$A_\alpha = \int_{\mathbb{T}} \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha - 1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)]$$

On the probability space $(\mathbb{T}, \mathcal{T}, \mu)$, we let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$. Set $\tilde{\Gamma}(v) = \Gamma(v)/\mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$. Then, $\mathbb{E}[1 - \tilde{\Gamma}(V)] = 0$ and

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \end{aligned}$$

Time to recall (A2)! The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Conclusion: $A_\alpha \geq 0$!

□

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (3)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

At this stage,

$$A_\alpha = \int_{\mathcal{T}} \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha - 1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)]$$

On the probability space $(\mathcal{T}, \mathcal{T}, \mu)$, we let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$. Set $\tilde{\Gamma}(v) = \Gamma(v)/\mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$. Then, $\mathbb{E}[1 - \tilde{\Gamma}(V)] = 0$ and

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \end{aligned}$$

Time to recall (A2)! The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Conclusion: $A_\alpha \geq 0$!

□

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (3)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

At this stage,

$$A_\alpha = \int_{\mathcal{T}} \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha - 1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)]$$

On the probability space $(\mathcal{T}, \mathcal{T}, \mu)$, we let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$. Set $\tilde{\Gamma}(v) = \Gamma(v)/\mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$. Then, $\mathbb{E}[1 - \tilde{\Gamma}(V)] = 0$ and

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \end{aligned}$$

Time to recall (A2)! The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Conclusion: $A_\alpha \geq 0$!

□

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (3)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

At this stage,

$$A_\alpha = \int_{\mathcal{T}} \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha - 1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)]$$

On the probability space $(\mathcal{T}, \mathcal{T}, \mu)$, we let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$. Set $\tilde{\Gamma}(v) = \Gamma(v)/\mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$. Then, $\mathbb{E}[1 - \tilde{\Gamma}(V)] = 0$ and

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \end{aligned}$$

Time to recall (A2)! The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is decreasing, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Conclusion: $A_\alpha \geq 0$!

□

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (3)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

At this stage,

$$A_\alpha = \int_{\mathcal{T}} \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha - 1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)]$$

On the probability space $(\mathcal{T}, \mathcal{T}, \mu)$, we let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$. Set $\tilde{\Gamma}(v) = \Gamma(v)/\mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$. Then, $\mathbb{E}[1 - \tilde{\Gamma}(V)] = 0$ and

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \end{aligned}$$

Time to recall (A2)! The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is **decreasing**, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Conclusion: $A_\alpha \geq 0$!

□

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (3)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

At this stage,

$$A_\alpha = \int_{\mathcal{T}} \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha - 1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)]$$

On the probability space $(\mathcal{T}, \mathcal{T}, \mu)$, we let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$. Set $\tilde{\Gamma}(v) = \Gamma(v)/\mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$. Then, $\mathbb{E}[1 - \tilde{\Gamma}(V)] = 0$ and

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \end{aligned}$$

Time to recall (A2)! The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is **decreasing**, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Conclusion: $A_\alpha \geq 0$!

□

Step 2 : take $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$ and show that $A_\alpha \geq 0$ (3)

For $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$, we have that

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha \circ \mathcal{I}_\alpha(\mu) .$$

The proof is complete if we prove that $A_\alpha \geq 0$.

At this stage,

$$A_\alpha = \int_{\mathcal{T}} \mu(d\theta) \left[b_{\mu,\alpha}(\theta) + \frac{1}{\alpha - 1} \right] g(\theta)^{\alpha-1} [1 - g(\theta)]$$

On the probability space $(\mathcal{T}, \mathcal{T}, \mu)$, we let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$. Set $\tilde{\Gamma}(v) = \Gamma(v)/\mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$. Then, $\mathbb{E}[1 - \tilde{\Gamma}(V)] = 0$ and

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \end{aligned}$$

Time to recall (A2)! The function $\Gamma : \text{Dom}_\alpha \rightarrow \mathbb{R}_{>0}$ is **decreasing**, continuously differentiable and satisfies the inequality

$$[(\alpha - 1)(v - \kappa) + 1] (\log \Gamma)'(v) + 1 \geq 0 .$$

Conclusion: $A_\alpha \geq 0$!

□

Thank you !