

Variational Inference

Foundations and recent advances

(Part 1)

Kamélia Daudel



University of Bristol – 09/03/2022

Outline

- 1 Introduction
- 2 Mean-field Variational Inference
- 3 Black-box Variational Inference
- 4 Alpha-divergence Variational Inference
- 5 Conclusion of Part 1

Outline

- 1 Introduction
- 2 Mean-field Variational Inference
- 3 Black-box Variational Inference
- 4 Alpha-divergence Variational Inference
- 5 Conclusion of Part 1

Bayesian inference

- Goal : model a phenomenon given some **observed data** while taking into account **prior knowledge** on the model parameters.
- Core quantity in Bayesian Inference : **posterior density** of the latent variables y given the data \mathcal{D}

$$p(y|\mathcal{D}) = \frac{p(\mathcal{D}, y)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|y)p_0(y)}{p(\mathcal{D})}$$

$p(\mathcal{D})$: normalisation constant 'marginal likelihood'

- What we would like : **compute / sample** from the posterior density (posterior mean, posterior predictive distribution...)
- Problem : for many important models, we can only evaluate $p(y|\mathcal{D})$ **up to the constant** $p(\mathcal{D})$.

Bayesian inference

- Goal : model a phenomenon given some **observed data** while taking into account **prior knowledge** on the model parameters.
- Core quantity in Bayesian Inference : **posterior density** of the latent variables y given the data \mathcal{D}

$$p(y|\mathcal{D}) = \frac{p(\mathcal{D}, y)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|y)p_0(y)}{p(\mathcal{D})}$$

$p(\mathcal{D})$: normalisation constant 'marginal likelihood'

- What we would like : **compute / sample** from the posterior density (posterior mean, posterior predictive distribution...)
- Problem : for many important models, we can only evaluate $p(y|\mathcal{D})$ **up to the constant** $p(\mathcal{D})$.

Bayesian inference

- Goal : model a phenomenon given some **observed data** while taking into account **prior knowledge** on the model parameters.
- Core quantity in Bayesian Inference : **posterior density** of the latent variables y given the data \mathcal{D}

$$p(y|\mathcal{D}) = \frac{p(\mathcal{D}, y)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|y)p_0(y)}{p(\mathcal{D})}$$

$p(\mathcal{D})$: normalisation constant ‘marginal likelihood’

- What we would like : **compute / sample** from the posterior density (posterior mean, posterior predictive distribution...)
- Problem : for many important models, we can only evaluate $p(y|\mathcal{D})$ **up to the constant** $p(\mathcal{D})$.

Bayesian inference

- Goal : model a phenomenon given some **observed data** while taking into account **prior knowledge** on the model parameters.
- Core quantity in Bayesian Inference : **posterior density** of the latent variables y given the data \mathcal{D}

$$p(y|\mathcal{D}) = \frac{p(\mathcal{D}, y)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|y)p_0(y)}{p(\mathcal{D})}$$

$p(\mathcal{D})$: normalisation constant ‘marginal likelihood’

- What we would like : **compute / sample** from the posterior density (posterior mean, posterior predictive distribution...)
- Problem : for many important models, we can only evaluate $p(y|\mathcal{D})$ **up to the constant** $p(\mathcal{D})$.

Bayesian inference

- Goal : model a phenomenon given some **observed data** while taking into account **prior knowledge** on the model parameters.
- Core quantity in Bayesian Inference : **posterior density** of the latent variables y given the data \mathcal{D}

$$p(y|\mathcal{D}) = \frac{p(\mathcal{D}, y)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|y)p_0(y)}{p(\mathcal{D})}$$

$p(\mathcal{D})$: normalisation constant ‘marginal likelihood’

- What we would like : **compute / sample** from the posterior density (posterior mean, posterior predictive distribution...)
- Problem : for many important models, we can only evaluate $p(y|\mathcal{D})$ **up to the constant** $p(\mathcal{D})$.

Approximate Bayesian Inference

Two broad categories of methods :

① Monte Carlo methods → sampling methods

- Importance Sampling (IS)
- Markov Chain Monte Carlo (MCMC)
- Sequential Monte Carlo (SMC) ...

② Variational Inference methods → optimisation-based methods

- Mean-field Variational Inference (MFVI)
- Black-Box Variational Inference (BBVI)
- Variational Auto-Encoder (VAE) ...

Approximate Bayesian Inference

Two broad categories of methods :

① Monte Carlo methods → sampling methods

- Importance Sampling (IS)
- Markov Chain Monte Carlo (MCMC)
- Sequential Monte Carlo (SMC) ...

② Variational Inference methods → optimisation-based methods

- Mean-field Variational Inference (MFVI)
- Black-Box Variational Inference (BBVI)
- Variational Auto-Encoder (VAE) ...

Approximate Bayesian Inference

Two broad categories of methods :

① Monte Carlo methods → sampling methods

- Importance Sampling (IS)
- Markov Chain Monte Carlo (MCMC)
- Sequential Monte Carlo (SMC) ...

② Variational Inference methods → optimisation-based methods

- Mean-field Variational Inference (MFVI)
- Black-Box Variational Inference (BBVI)
- Variational Auto-Encoder (VAE) ...

Approximate Bayesian Inference

Two broad categories of methods :

① Monte Carlo methods → sampling methods

- Importance Sampling (IS)
- Markov Chain Monte Carlo (MCMC)
- Sequential Monte Carlo (SMC) ...

② Variational Inference methods → optimisation-based methods

- Mean-field Variational Inference (MFVI)
- Black-Box Variational Inference (BBVI)
- Variational Auto-Encoder (VAE) ...

Approximate Bayesian Inference

Two broad categories of methods :

① Monte Carlo methods → sampling methods

- Importance Sampling (IS)
- Markov Chain Monte Carlo (MCMC)
- Sequential Monte Carlo (SMC) ...

② Variational Inference methods → optimisation-based methods

- Mean-field Variational Inference (MFVI)
- Black-Box Variational Inference (BBVI)
- Variational Auto-Encoder (VAE) ...

Approximate Bayesian Inference

Two broad categories of methods :

① Monte Carlo methods → sampling methods

- Importance Sampling (IS)
- Markov Chain Monte Carlo (MCMC)
- Sequential Monte Carlo (SMC) ...

② Variational Inference methods → optimisation-based methods

- Mean-field Variational Inference (MFVI)
- Black-Box Variational Inference (BBVI)
- Variational Auto-Encoder (VAE) ...

Approximate Bayesian Inference

Two broad categories of methods :

① Monte Carlo methods → sampling methods

- Importance Sampling (IS)
- Markov Chain Monte Carlo (MCMC)
- Sequential Monte Carlo (SMC) ...

② Variational Inference methods → optimisation-based methods

- Mean-field Variational Inference (MFVI)
- Black-Box Variational Inference (BBVI)
- Variational Auto-Encoder (VAE) ...

Approximate Bayesian Inference

Two broad categories of methods :

① Monte Carlo methods → sampling methods

- Importance Sampling (IS)
- Markov Chain Monte Carlo (MCMC)
- Sequential Monte Carlo (SMC) ...

② Variational Inference methods → optimisation-based methods

- Mean-field Variational Inference (MFVI)
- Black-Box Variational Inference (BBVI)
- Variational Auto-Encoder (VAE) ...

Approximate Bayesian Inference

Two broad categories of methods :

① Monte Carlo methods → sampling methods

- Importance Sampling (IS)
- Markov Chain Monte Carlo (MCMC)
- Sequential Monte Carlo (SMC) ...

② Variational Inference methods → optimisation-based methods

- Mean-field Variational Inference (MFVI)
- Black-Box Variational Inference (BBVI)
- Variational Auto-Encoder (VAE) ...

Variational Inference in a nutshell

Variational Inference methodology

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

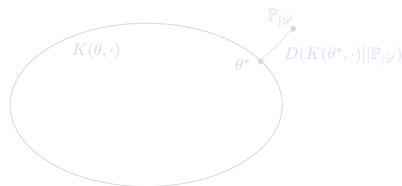
$$\inf_{q \in \mathcal{Q}} D(Q || \mathbb{P}_{|\mathcal{D}}) \quad (1)$$

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

→ D and \mathcal{Q} are key elements in the optimisation problem (1) !

What we want :

- \mathcal{Q} is easy to sample from / optimise over, yet can capture the complexity inside $p(y|\mathcal{D})$ (e.g. well-chosen **parametric** family $\mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathcal{T}\}$)
- $D(Q || \mathbb{P}_{|\mathcal{D}})$ can be optimised efficiently



Variational Inference in a nutshell

Variational Inference methodology

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

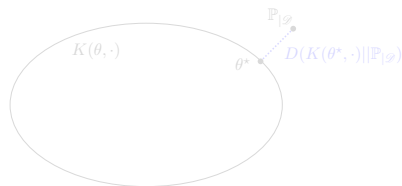
$$\inf_{q \in \mathcal{Q}} D(Q || \mathbb{P}_{|\mathcal{D}}) \quad (1)$$

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

→ D and \mathcal{Q} are key elements in the optimisation problem (1) !

What we want :

- \mathcal{Q} is easy to sample from / optimise over, yet can capture the complexity inside $p(y|\mathcal{D})$ (e.g. well-chosen **parametric** family $\mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathcal{T}\}$)
- $D(Q || \mathbb{P}_{|\mathcal{D}})$ can be optimised efficiently



Variational Inference in a nutshell

Variational Inference methodology

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

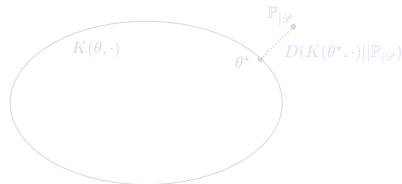
$$\inf_{q \in \mathcal{Q}} D(Q || \mathbb{P}_{|\mathcal{D}}) \quad (1)$$

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

→ D and \mathcal{Q} are key elements in the optimisation problem (1) !

What we want :

- \mathcal{Q} is easy to sample from / optimise over, yet can capture the complexity inside $p(y|\mathcal{D})$ (e.g. well-chosen **parametric** family $\mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathcal{T}\}$)
- $D(Q || \mathbb{P}_{|\mathcal{D}})$ can be optimised efficiently



Variational Inference in a nutshell

Variational Inference methodology

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

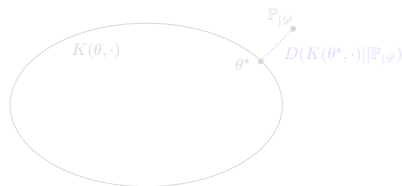
$$\inf_{q \in \mathcal{Q}} D(Q || \mathbb{P}_{|\mathcal{D}}) \quad (1)$$

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

→ D and \mathcal{Q} are key elements in the optimisation problem (1) !

What we want :

- \mathcal{Q} is easy to sample from / optimise over, yet can capture the complexity inside $p(y|\mathcal{D})$ (e.g. well-chosen **parametric** family $\mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathcal{T}\}$)
- $D(Q || \mathbb{P}_{|\mathcal{D}})$ can be optimised efficiently



Variational Inference in a nutshell

Variational Inference methodology

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

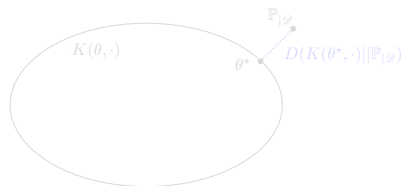
$$\inf_{q \in \mathcal{Q}} D(Q || \mathbb{P}_{|\mathcal{D}}) \quad (1)$$

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

→ D and \mathcal{Q} are key elements in the optimisation problem (1) !

What we want :

- \mathcal{Q} is easy to sample from / optimise over, yet can capture the complexity inside $p(y|\mathcal{D})$ (e.g. well-chosen **parametric** family $\mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathcal{T}\}$)
- $D(Q || \mathbb{P}_{|\mathcal{D}})$ can be optimised efficiently



Variational Inference in a nutshell

Variational Inference methodology

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

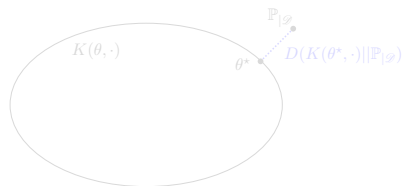
$$\inf_{q \in \mathcal{Q}} D(Q || \mathbb{P}_{|\mathcal{D}}) \quad (1)$$

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

→ D and \mathcal{Q} are key elements in the optimisation problem (1) !

What we want :

- \mathcal{Q} is easy to sample from / optimise over, yet can capture the complexity inside $p(y|\mathcal{D})$ (e.g. well-chosen **parametric** family $\mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathbf{T}\}$)
- $D(Q || \mathbb{P}_{|\mathcal{D}})$ can be optimised efficiently



Variational Inference in a nutshell

Variational Inference methodology

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

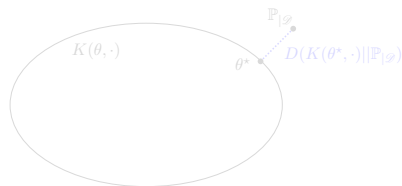
$$\inf_{q \in \mathcal{Q}} D(Q || \mathbb{P}_{|\mathcal{D}}) \quad (1)$$

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

→ D and \mathcal{Q} are key elements in the optimisation problem (1) !

What we want :

- \mathcal{Q} is easy to sample from / optimise over, yet can capture the complexity inside $p(y|\mathcal{D})$ (e.g. well-chosen **parametric** family $Q = \{q : y \mapsto k(\theta, y) : \theta \in \mathbf{T}\}$)
- $D(Q || \mathbb{P}_{|\mathcal{D}})$ can be optimised efficiently



Variational Inference in a nutshell

Variational Inference methodology

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

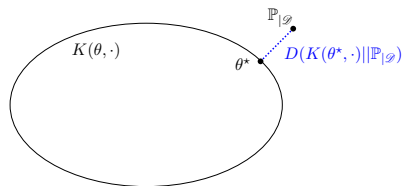
$$\inf_{q \in \mathcal{Q}} D(Q || \mathbb{P}_{|\mathcal{D}}) \quad (1)$$

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

→ D and \mathcal{Q} are key elements in the optimisation problem (1) !

What we want :

- \mathcal{Q} is easy to sample from / optimise over, yet can capture the complexity inside $p(y|\mathcal{D})$ (e.g. well-chosen **parametric** family $\mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathbf{T}\}$)
- $D(Q || \mathbb{P}_{|\mathcal{D}})$ can be optimised efficiently



Variational Inference in a nutshell

Variational Inference methodology

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

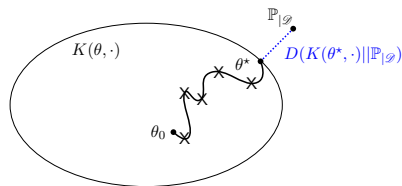
$$\inf_{q \in \mathcal{Q}} D(Q || \mathbb{P}_{|\mathcal{D}}) \quad (1)$$

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

→ D and \mathcal{Q} are key elements in the optimisation problem (1) !

What we want :

- \mathcal{Q} is easy to sample from / optimise over, yet can capture the complexity inside $p(y|\mathcal{D})$ (e.g. well-chosen **parametric** family $\mathcal{Q} = \{q : y \mapsto k(\theta, y) : \theta \in \mathbf{T}\}$)
- $D(Q || \mathbb{P}_{|\mathcal{D}})$ can be optimised efficiently



Outline

- 1 Introduction
- 2 Mean-field Variational Inference**
- 3 Black-box Variational Inference
- 4 Alpha-divergence Variational Inference
- 5 Conclusion of Part 1

Mean-field Variational Inference (MFVI)

→ D : **Kullback-Leibler** (KL) divergence

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

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

$$D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) = \int_Y \log \left(\frac{q(y)}{p(y|\mathcal{D})} \right) q(y) \nu(dy) .$$

$D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \geq 0$ and $D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) = 0$ iff $\mathbb{Q} = \mathbb{P}_{|\mathcal{D}}$

→ \mathcal{Q} : **Mean-field** family

The latent variable y is made of L **independent** latent variables $(y_1, \dots, y_L) \in Y_1 \times \dots \times Y_L$ and

$$\mathcal{Q} = \left\{ q : y \mapsto \prod_{\ell=1}^L q_{\ell}(y_{\ell}) \right\}$$

i.e each latent variable y_{ℓ} is governed by its own variational probability density q_{ℓ} with $\nu(dy) = \bigotimes_{\ell=1}^L \nu_{\ell}(dy_{\ell})$.

Mean-field Variational Inference (MFVI)

→ D : **Kullback-Leibler** (KL) divergence

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

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

$$D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) = \int_Y \log \left(\frac{q(y)}{p(y|\mathcal{D})} \right) q(y) \nu(dy) .$$

$D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \geq 0$ and $D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) = 0$ iff $\mathbb{Q} = \mathbb{P}_{|\mathcal{D}}$

→ \mathcal{Q} : **Mean-field** family

The latent variable y is made of L **independent** latent variables $(y_1, \dots, y_L) \in Y_1 \times \dots \times Y_L$ and

$$\mathcal{Q} = \left\{ q : y \mapsto \prod_{\ell=1}^L q_{\ell}(y_{\ell}) \right\}$$

i.e each latent variable y_{ℓ} is governed by its own variational probability density q_{ℓ} with $\nu(dy) = \bigotimes_{\ell=1}^L \nu_{\ell}(dy_{\ell})$.

Mean-field Variational Inference (MFVI)

→ D : **Kullback-Leibler** (KL) divergence

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

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

$$D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) = \int_Y \log \left(\frac{q(y)}{p(y|\mathcal{D})} \right) q(y) \nu(dy) .$$

$D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \geq 0$ and $D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) = 0$ iff $\mathbb{Q} = \mathbb{P}_{|\mathcal{D}}$

→ \mathcal{Q} : **Mean-field** family

The latent variable y is made of L **independent** latent variables

$(y_1, \dots, y_L) \in Y_1 \times \dots \times Y_L$ and

$$\mathcal{Q} = \left\{ q : y \mapsto \prod_{\ell=1}^L q_{\ell}(y_{\ell}) \right\}$$

i.e each latent variable y_{ℓ} is governed by its own variational probability density q_{ℓ} with $\nu(dy) = \bigotimes_{\ell=1}^L \nu_{\ell}(dy_{\ell})$.

Why this choice of D ?

$$\begin{aligned} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) &= \int_Y \log \left(\frac{q(y)}{p(y|\mathcal{D})} \right) q(y) \nu(dy) \\ &= \int_Y q(y) \log \left(\frac{q(y)}{p(y, \mathcal{D})} \right) \nu(dy) + \log p(\mathcal{D}) \\ &:= -\text{ELBO}(q; \mathcal{D}) + \log p(\mathcal{D}) \end{aligned}$$

Evidence Lower Bound (ELBO)

$$\text{ELBO}(q; \mathcal{D}) := \int_Y q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) .$$

→ We deduce :

- ① $\inf_{q \in \mathcal{Q}} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \Leftrightarrow \sup_{q \in \mathcal{Q}} \text{ELBO}(q; \mathcal{D})$
- ② $\text{ELBO}(q; \mathcal{D}) \leq \log p(\mathcal{D})$ with equality iff $\mathbb{Q} = \mathbb{P}_{|\mathcal{D}}$

Why this choice of D ?

$$\begin{aligned} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) &= \int_Y \log \left(\frac{q(y)}{p(y|\mathcal{D})} \right) q(y) \nu(dy) \\ &= \int_Y q(y) \log \left(\frac{q(y)}{p(y, \mathcal{D})} \right) \nu(dy) + \log p(\mathcal{D}) \\ &:= -\text{ELBO}(q; \mathcal{D}) + \log p(\mathcal{D}) \end{aligned}$$

Evidence Lower Bound (ELBO)

$$\text{ELBO}(q; \mathcal{D}) := \int_Y q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) .$$

→ We deduce :

- ① $\inf_{q \in \mathcal{Q}} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \Leftrightarrow \sup_{q \in \mathcal{Q}} \text{ELBO}(q; \mathcal{D})$
- ② $\text{ELBO}(q; \mathcal{D}) \leq \log p(\mathcal{D})$ with equality iff $\mathbb{Q} = \mathbb{P}_{|\mathcal{D}}$

Why this choice of D ?

$$\begin{aligned} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) &= \int_Y \log \left(\frac{q(y)}{p(y|\mathcal{D})} \right) q(y) \nu(dy) \\ &= \int_Y q(y) \log \left(\frac{q(y)}{p(y, \mathcal{D})} \right) \nu(dy) + \log p(\mathcal{D}) \\ &:= -\text{ELBO}(q; \mathcal{D}) + \log p(\mathcal{D}) \end{aligned}$$

Evidence Lower Bound (ELBO)

$$\text{ELBO}(q; \mathcal{D}) := \int_Y q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) .$$

→ We deduce :

- ① $\inf_{q \in \mathcal{Q}} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \Leftrightarrow \sup_{q \in \mathcal{Q}} \text{ELBO}(q; \mathcal{D})$
- ② $\text{ELBO}(q; \mathcal{D}) \leq \log p(\mathcal{D})$ with equality iff $\mathbb{Q} = \mathbb{P}_{|\mathcal{D}}$

Why this choice of D ?

$$\begin{aligned} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) &= \int_Y \log \left(\frac{q(y)}{p(y|\mathcal{D})} \right) q(y) \nu(dy) \\ &= \int_Y q(y) \log \left(\frac{q(y)}{p(y, \mathcal{D})} \right) \nu(dy) + \log p(\mathcal{D}) \\ &:= -\text{ELBO}(q; \mathcal{D}) + \log p(\mathcal{D}) \end{aligned}$$

Evidence Lower Bound (ELBO)

$$\text{ELBO}(q; \mathcal{D}) := \int_Y q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) .$$

→ We deduce :

- ① $\inf_{q \in \mathcal{Q}} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \Leftrightarrow \sup_{q \in \mathcal{Q}} \text{ELBO}(q; \mathcal{D})$
- ② $\text{ELBO}(q; \mathcal{D}) \leq \log p(\mathcal{D})$ with equality iff $\mathbb{Q} = \mathbb{P}_{|\mathcal{D}}$

Why this choice of D ?

$$\begin{aligned} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) &= \int_{\mathbf{Y}} \log \left(\frac{q(y)}{p(y|\mathcal{D})} \right) q(y) \nu(dy) \\ &= \int_{\mathbf{Y}} q(y) \log \left(\frac{q(y)}{p(y, \mathcal{D})} \right) \nu(dy) + \log p(\mathcal{D}) \\ &:= -\text{ELBO}(q; \mathcal{D}) + \log p(\mathcal{D}) \end{aligned}$$

Evidence Lower Bound (ELBO)

$$\text{ELBO}(q; \mathcal{D}) := \int_{\mathbf{Y}} q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) .$$

→ We deduce :

- ① $\inf_{q \in \mathcal{Q}} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \Leftrightarrow \sup_{q \in \mathcal{Q}} \text{ELBO}(q; \mathcal{D})$
- ② $\text{ELBO}(q; \mathcal{D}) \leq \log p(\mathcal{D})$ with equality iff $\mathbb{Q} = \mathbb{P}_{|\mathcal{D}}$

Why this choice of D ?

$$\begin{aligned} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) &= \int_{\mathbf{Y}} \log \left(\frac{q(y)}{p(y|\mathcal{D})} \right) q(y) \nu(dy) \\ &= \int_{\mathbf{Y}} q(y) \log \left(\frac{q(y)}{p(y, \mathcal{D})} \right) \nu(dy) + \log p(\mathcal{D}) \\ &:= -\text{ELBO}(q; \mathcal{D}) + \log p(\mathcal{D}) \end{aligned}$$

Evidence Lower Bound (ELBO)

$$\text{ELBO}(q; \mathcal{D}) := \int_{\mathbf{Y}} q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) .$$

→ We deduce :

- ❶ $\inf_{q \in \mathcal{Q}} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \Leftrightarrow \sup_{q \in \mathcal{Q}} \text{ELBO}(q; \mathcal{D})$
- ❷ $\text{ELBO}(q; \mathcal{D}) \leq \log p(\mathcal{D})$ with equality iff $\mathbb{Q} = \mathbb{P}_{|\mathcal{D}}$

Why this choice of D ?

$$\begin{aligned} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) &= \int_Y \log \left(\frac{q(y)}{p(y|\mathcal{D})} \right) q(y) \nu(dy) \\ &= \int_Y q(y) \log \left(\frac{q(y)}{p(y, \mathcal{D})} \right) \nu(dy) + \log p(\mathcal{D}) \\ &:= -\text{ELBO}(q; \mathcal{D}) + \log p(\mathcal{D}) \end{aligned}$$

Evidence Lower Bound (ELBO)

$$\text{ELBO}(q; \mathcal{D}) := \int_Y q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) .$$

→ We deduce :

- ❶ $\inf_{q \in \mathcal{Q}} D_{KL}(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) \Leftrightarrow \sup_{q \in \mathcal{Q}} \text{ELBO}(q; \mathcal{D})$
- ❷ $\text{ELBO}(q; \mathcal{D}) \leq \log p(\mathcal{D})$ with equality iff $\mathbb{Q} = \mathbb{P}_{|\mathcal{D}}$

Why this choice of \mathcal{Q} ?

Recall that

Mean-field assumption

The latent variable y is made of L **independent** latent variables $(y_1, \dots, y_L) \in Y_1 \times \dots \times Y_L$ and

$$\mathcal{Q} = \left\{ q : y \mapsto \prod_{\ell=1}^L q_{\ell}(y_{\ell}) \right\}$$

i.e each latent variable y_{ℓ} is governed by its own variational probability density q_{ℓ} with $\nu(dy) = \bigotimes_{\ell=1}^L \nu_{\ell}(dy_{\ell})$.

→ Plugging this into the ELBO and keeping all factors but ℓ fixed :

$$q_{\ell}^*(y_{\ell}) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})]) \quad (\text{optimal rule})$$

where $\mathbb{E}_{-\ell}$ is the expectation w.r.t q omitting the factor q_{ℓ}

Why this choice of \mathcal{Q} ?

Recall that

Mean-field assumption

The latent variable y is made of L **independent** latent variables $(y_1, \dots, y_L) \in Y_1 \times \dots \times Y_L$ and

$$\mathcal{Q} = \left\{ q : y \mapsto \prod_{\ell=1}^L q_{\ell}(y_{\ell}) \right\}$$

i.e each latent variable y_{ℓ} is governed by its own variational probability density q_{ℓ} with $\nu(dy) = \bigotimes_{\ell=1}^L \nu_{\ell}(dy_{\ell})$.

→ Plugging this into the ELBO and keeping all factors but ℓ fixed :

$$q_{\ell}^*(y_{\ell}) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})]) \quad (\text{optimal rule})$$

where $\mathbb{E}_{-\ell}$ is the expectation w.r.t q omitting the factor q_{ℓ}

CAVI algorithm

Optimal rule **keeping all factors but ℓ fixed** :

$$q_{\ell}^*(y_{\ell}) \propto \exp (\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$$

CAVI algorithm

Optimal rule **keeping all factors but ℓ fixed** :

$$q_{\ell}^*(y_{\ell}) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$$

Algorithm 1: *Coordinate Ascent Variational Inference (CAVI)*

Input: $(q_{\ell})_{1 \leq \ell \leq L}$: initial variational factors.

Output: Return the optimised mean-field variational density q satisfying:

for all $y \in \mathcal{Y}$, $q(y) = \prod_{\ell=1}^L q_{\ell}(y_{\ell})$.

while the ELBO has not converged **do**

for $\ell = 1 \dots L$ **do**

 set $q_{\ell}(y_{\ell}) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$

end

 Compute the ELBO.

end

CAVI algorithm

Optimal rule **keeping all factors but ℓ fixed** :

$$q_{\ell}^*(y_{\ell}) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$$

Algorithm 1: *Coordinate Ascent Variational Inference (CAVI)*

Input: $(q_{\ell})_{1 \leq \ell \leq L}$: initial variational factors.

Output: Return the optimised mean-field variational density q satisfying:

for all $y \in \mathcal{Y}$, $q(y) = \prod_{\ell=1}^L q_{\ell}(y_{\ell})$.

while the ELBO has not converged **do**

for $\ell = 1 \dots L$ **do**

 set $q_{\ell}(y_{\ell}) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$

end

 Compute the ELBO.

end

→ Convergence towards a **local** maximum of the ELBO

CAVI algorithm

Optimal rule **keeping all factors but ℓ fixed** :

$$q_{\ell}^*(y_{\ell}) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$$

Algorithm 1: *Coordinate Ascent Variational Inference (CAVI)*

Input: $(q_{\ell})_{1 \leq \ell \leq L}$: initial variational factors.

Output: Return the optimised mean-field variational density q satisfying:

for all $y \in \mathbf{Y}$, $q(y) = \prod_{\ell=1}^L q_{\ell}(y_{\ell})$.

while the ELBO has not converged **do**

for $\ell = 1 \dots L$ **do**

 set $q_{\ell}(y_{\ell}) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$

end

 Compute the ELBO.

end

→ Convergence towards a **local** maximum of the ELBO

→ Tractable updates for **conditionally conjugate exponential** models
(e.g. Bayesian mixture of Gaussians, Latent Dirichlet Allocation)

CAVI algorithm

Optimal rule **keeping all factors but ℓ fixed** :

$$q_{\ell}^*(y_{\ell}) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$$

Algorithm 1: *Coordinate Ascent Variational Inference (CAVI)*

Input: $(q_{\ell})_{1 \leq \ell \leq L}$: initial variational factors.

Output: Return the optimised mean-field variational density q satisfying:

for all $y \in \mathcal{Y}$, $q(y) = \prod_{\ell=1}^L q_{\ell}(y_{\ell})$.

while the ELBO has not converged **do**

for $\ell = 1 \dots L$ **do**

 set $q_{\ell}(y_{\ell}) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$

end

 Compute the ELBO.

end

→ Convergence towards a **local** maximum of the ELBO

→ Tractable updates for **conditionally conjugate exponential** models
(e.g. Bayesian mixture of Gaussians, Latent Dirichlet Allocation)

Variational Inference: A Review for Statisticians. D. Blei et al. (2017). JASA

The New York Times

music band songs rock album jazz pop song singer night	book life novel story books man stories love children family	art museum show exhibition artist artists paintings painting century works	game knicks nets points team season play games night coach	show film television movie series says life man character know
theater play production show stage street broadway director musical directed	clinton bush campaign gore political republican dole presidential senator house	stock market percent fund investors funds companies stocks investment trading	restaurant sauce menu food dishes street dining dinner chicken served	budget tax governor county mayor billion taxes plan legislature fiscal

Data : 1.8M articles from the New York Times

Model : hierarchical Dirichlet process topic model

Taken from **Stochastic Variational Inference**. M. D. Hoffman et al. (2013). JMRL.

Toy example : Bayesian Linear Regression (BLR) - 1

- $\mathcal{D} = \{\mathbf{c}, \mathbf{x}\} : I$ 1-D class labels $(c_i)_{1 \leq i \leq I}$, I 2-D covariates $(\mathbf{x}_i)_{1 \leq i \leq I}$
- $y = \{y_1, y_2\} \in \mathbb{R}^2$: regression coefficients
- Model :

$$p(c_i | \mathbf{x}_i, y) = \mathcal{N}(c_i; y^T \mathbf{x}_i, \sigma^2), \quad 1 \leq i \leq I$$
$$p_0(y) = \mathcal{N}(y; \mu_0, \Lambda_0^{-1})$$

μ_0, Λ_0, σ : fixed hyperparameters

In that case,

$$p(y | \mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

with $\Lambda = \Lambda_0 + \sigma^{-2} \sum_{i=1}^I \mathbf{x}_i \mathbf{x}_i^T$ and $\Lambda \mu = \Lambda_0 \mu_0 + \sigma^{-2} \sum_{i=1}^I c_i \mathbf{x}_i$.

Toy example : Bayesian Linear Regression (BLR) - 1

- $\mathcal{D} = \{\mathbf{c}, \mathbf{x}\} : I$ 1-D class labels $(c_i)_{1 \leq i \leq I}$, I 2-D covariates $(\mathbf{x}_i)_{1 \leq i \leq I}$
- $y = \{y_1, y_2\} \in \mathbb{R}^2$: regression coefficients
- Model :

$$p(c_i | \mathbf{x}_i, y) = \mathcal{N}(c_i; y^T \mathbf{x}_i, \sigma^2), \quad 1 \leq i \leq I$$
$$p_0(y) = \mathcal{N}(y; \mu_0, \Lambda_0^{-1})$$

μ_0, Λ_0, σ : fixed hyperparameters

In that case,

$$p(y | \mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

with $\Lambda = \Lambda_0 + \sigma^{-2} \sum_{i=1}^I \mathbf{x}_i \mathbf{x}_i^T$ and $\Lambda \mu = \Lambda_0 \mu_0 + \sigma^{-2} \sum_{i=1}^I c_i \mathbf{x}_i$.

Toy example : Bayesian Linear Regression (BLR) - 2

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

- Mean-field assumption : $q(y) = q_1(y_1)q_2(y_2)$
- Optimal rules : for all $\ell = \{1, 2\}$, $q_\ell(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$
that is,

$$q_1(y_1) \propto \exp(\mathbb{E}_{y_2 \sim q_2}[\log p(y|\mathcal{D})])$$

$$q_2(y_2) \propto \exp(\mathbb{E}_{y_1 \sim q_1}[\log p(y|\mathcal{D})])$$

Notation : $\mu = (\mu_1 \ \mu_2)$, $\Lambda = (\Lambda_{\ell,k})_{1 \leq \ell, k \leq 2}$ with $\Lambda_{1,2} = \Lambda_{2,1}$

$$\log p(y|\mathcal{D}) = -\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(y_2 - \mu_2) \Lambda_{1,2} + (y_2 - \mu_2)^2 \Lambda_{2,2} \} + c_{-y}$$

→ Plugging this in the optimal rule,

$$q_1(y_1) \propto \exp \left(-\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2} \} \right)$$

$$\propto \exp \left(-\frac{1}{2} \{ y_1^2 \Lambda_{1,1} - 2y_1 [\mu_1 \Lambda_{1,1} - (\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}] \} \right)$$

so that : $q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}, \Lambda_{1,1}^{-1})$

Toy example : Bayesian Linear Regression (BLR) - 2

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

- Mean-field assumption : $q(y) = q_1(y_1)q_2(y_2)$
- Optimal rules : for all $\ell = \{1, 2\}$, $q_\ell(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$
that is,

$$q_1(y_1) \propto \exp(\mathbb{E}_{y_2 \sim q_2}[\log p(y|\mathcal{D})])$$

$$q_2(y_2) \propto \exp(\mathbb{E}_{y_1 \sim q_1}[\log p(y|\mathcal{D})])$$

Notation : $\mu = (\mu_1 \ \mu_2)$, $\Lambda = (\Lambda_{\ell,k})_{1 \leq \ell, k \leq 2}$ with $\Lambda_{1,2} = \Lambda_{2,1}$

$$\log p(y|\mathcal{D}) = -\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(y_2 - \mu_2) \Lambda_{1,2} + (y_2 - \mu_2)^2 \Lambda_{2,2} \} + c_{-y}$$

→ Plugging this in the optimal rule,

$$q_1(y_1) \propto \exp \left(-\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2} \} \right)$$

$$\propto \exp \left(-\frac{1}{2} \{ y_1^2 \Lambda_{1,1} - 2y_1 [\mu_1 \Lambda_{1,1} - (\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}] \} \right)$$

so that : $q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}, \Lambda_{1,1}^{-1})$

Toy example : Bayesian Linear Regression (BLR) - 2

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

- Mean-field assumption : $q(y) = q_1(y_1)q_2(y_2)$
- Optimal rules : for all $\ell = \{1, 2\}$, $q_\ell(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$
that is,

$$q_1(y_1) \propto \exp(\mathbb{E}_{y_2 \sim q_2}[\log p(y|\mathcal{D})])$$

$$q_2(y_2) \propto \exp(\mathbb{E}_{y_1 \sim q_1}[\log p(y|\mathcal{D})])$$

Notation : $\mu = (\mu_1 \ \mu_2)$, $\Lambda = (\Lambda_{\ell,k})_{1 \leq \ell, k \leq 2}$ with $\Lambda_{1,2} = \Lambda_{2,1}$

$$\log p(y|\mathcal{D}) = -\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(y_2 - \mu_2) \Lambda_{1,2} + (y_2 - \mu_2)^2 \Lambda_{2,2} \} + c_{-y}$$

→ Plugging this in the optimal rule,

$$q_1(y_1) \propto \exp \left(-\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2} \} \right)$$

$$\propto \exp \left(-\frac{1}{2} \{ y_1^2 \Lambda_{1,1} - 2y_1 [\mu_1 \Lambda_{1,1} - (\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}] \} \right)$$

so that : $q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}, \Lambda_{1,1}^{-1})$

Toy example : Bayesian Linear Regression (BLR) - 2

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

- Mean-field assumption : $q(y) = q_1(y_1)q_2(y_2)$
- Optimal rules : for all $\ell = \{1, 2\}$, $q_\ell(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$
that is,

$$q_1(y_1) \propto \exp(\mathbb{E}_{y_2 \sim q_2}[\log p(y|\mathcal{D})])$$

$$q_2(y_2) \propto \exp(\mathbb{E}_{y_1 \sim q_1}[\log p(y|\mathcal{D})])$$

Notation : $\mu = (\mu_1 \ \mu_2)$, $\Lambda = (\Lambda_{\ell,k})_{1 \leq \ell, k \leq 2}$ with $\Lambda_{1,2} = \Lambda_{2,1}$

$$\log p(y|\mathcal{D}) = -\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(y_2 - \mu_2) \Lambda_{1,2} + (y_2 - \mu_2)^2 \Lambda_{2,2} \} + c_{-y}$$

→ Plugging this in the optimal rule,

$$q_1(y_1) \propto \exp \left(-\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2} \} \right)$$

$$\propto \exp \left(-\frac{1}{2} \{ y_1^2 \Lambda_{1,1} - 2y_1 [\mu_1 \Lambda_{1,1} - (\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}] \} \right)$$

so that : $q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}, \Lambda_{1,1}^{-1})$

Toy example : Bayesian Linear Regression (BLR) - 2

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

- Mean-field assumption : $q(y) = q_1(y_1)q_2(y_2)$
- Optimal rules : for all $\ell = \{1, 2\}$, $q_\ell(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$
that is,

$$q_1(y_1) \propto \exp(\mathbb{E}_{y_2 \sim q_2}[\log p(y|\mathcal{D})])$$

$$q_2(y_2) \propto \exp(\mathbb{E}_{y_1 \sim q_1}[\log p(y|\mathcal{D})])$$

Notation : $\mu = (\mu_1 \ \mu_2)$, $\Lambda = (\Lambda_{\ell,k})_{1 \leq \ell, k \leq 2}$ with $\Lambda_{1,2} = \Lambda_{2,1}$

$$\log p(y|\mathcal{D}) = -\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(y_2 - \mu_2) \Lambda_{1,2} + (y_2 - \mu_2)^2 \Lambda_{2,2} \} + c_{-y}$$

→ Plugging this in the optimal rule,

$$q_1(y_1) \propto \exp \left(-\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2} \} \right)$$

$$\propto \exp \left(-\frac{1}{2} \{ y_1^2 \Lambda_{1,1} - 2y_1 [\mu_1 \Lambda_{1,1} - (\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}] \} \right)$$

so that : $q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}, \Lambda_{1,1}^{-1})$

Toy example : Bayesian Linear Regression (BLR) - 2

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

- Mean-field assumption : $q(y) = q_1(y_1)q_2(y_2)$
- Optimal rules : for all $\ell = \{1, 2\}$, $q_\ell(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$
that is,

$$q_1(y_1) \propto \exp(\mathbb{E}_{y_2 \sim q_2}[\log p(y|\mathcal{D})])$$

$$q_2(y_2) \propto \exp(\mathbb{E}_{y_1 \sim q_1}[\log p(y|\mathcal{D})])$$

Notation : $\mu = (\mu_1 \ \mu_2)$, $\Lambda = (\Lambda_{\ell,k})_{1 \leq \ell, k \leq 2}$ with $\Lambda_{1,2} = \Lambda_{2,1}$

$$\log p(y|\mathcal{D}) = -\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(y_2 - \mu_2) \Lambda_{1,2} + (y_2 - \mu_2)^2 \Lambda_{2,2} \} + c_{-y}$$

→ Plugging this in the optimal rule,

$$q_1(y_1) \propto \exp \left(-\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2} \} \right)$$

$$\propto \exp \left(-\frac{1}{2} \{ y_1^2 \Lambda_{1,1} - 2y_1 [\mu_1 \Lambda_{1,1} - (\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}] \} \right)$$

so that : $q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}, \Lambda_{1,1}^{-1})$

Toy example : Bayesian Linear Regression (BLR) - 2

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

- Mean-field assumption : $q(y) = q_1(y_1)q_2(y_2)$
- Optimal rules : for all $\ell = \{1, 2\}$, $q_\ell(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$
that is,

$$q_1(y_1) \propto \exp(\mathbb{E}_{y_2 \sim q_2}[\log p(y|\mathcal{D})])$$

$$q_2(y_2) \propto \exp(\mathbb{E}_{y_1 \sim q_1}[\log p(y|\mathcal{D})])$$

Notation : $\mu = (\mu_1 \ \mu_2)$, $\Lambda = (\Lambda_{\ell,k})_{1 \leq \ell, k \leq 2}$ with $\Lambda_{1,2} = \Lambda_{2,1}$

$$\log p(y|\mathcal{D}) = -\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(y_2 - \mu_2) \Lambda_{1,2} + (y_2 - \mu_2)^2 \Lambda_{2,2} \} + c_{-y}$$

→ Plugging this in the optimal rule,

$$q_1(y_1) \propto \exp \left(-\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2} \} \right)$$

$$\propto \exp \left(-\frac{1}{2} \{ y_1^2 \Lambda_{1,1} - 2y_1 [\mu_1 \Lambda_{1,1} - (\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}] \} \right)$$

so that : $q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}, \Lambda_{1,1}^{-1})$

Toy example : Bayesian Linear Regression (BLR) - 2

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

- Mean-field assumption : $q(y) = q_1(y_1)q_2(y_2)$
- Optimal rules : for all $\ell = \{1, 2\}$, $q_\ell(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$
that is,

$$q_1(y_1) \propto \exp(\mathbb{E}_{y_2 \sim q_2}[\log p(y|\mathcal{D})])$$

$$q_2(y_2) \propto \exp(\mathbb{E}_{y_1 \sim q_1}[\log p(y|\mathcal{D})])$$

Notation : $\mu = (\mu_1 \ \mu_2)$, $\Lambda = (\Lambda_{\ell,k})_{1 \leq \ell, k \leq 2}$ with $\Lambda_{1,2} = \Lambda_{2,1}$

$$\log p(y|\mathcal{D}) = -\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(y_2 - \mu_2) \Lambda_{1,2} + (y_2 - \mu_2)^2 \Lambda_{2,2} \} + c_{-y}$$

→ Plugging this in the optimal rule,

$$\begin{aligned} q_1(y_1) &\propto \exp \left(-\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2} \} \right) \\ &\propto \exp \left(-\frac{1}{2} \{ y_1^2 \Lambda_{1,1} - 2y_1 [\mu_1 \Lambda_{1,1} - (\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}] \} \right) \end{aligned}$$

so that : $q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}, \Lambda_{1,1}^{-1})$

Toy example : Bayesian Linear Regression (BLR) - 2

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

- Mean-field assumption : $q(y) = q_1(y_1)q_2(y_2)$
- Optimal rules : for all $\ell = \{1, 2\}$, $q_\ell(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$
that is,

$$q_1(y_1) \propto \exp(\mathbb{E}_{y_2 \sim q_2}[\log p(y|\mathcal{D})])$$

$$q_2(y_2) \propto \exp(\mathbb{E}_{y_1 \sim q_1}[\log p(y|\mathcal{D})])$$

Notation : $\mu = (\mu_1 \ \mu_2)$, $\Lambda = (\Lambda_{\ell,k})_{1 \leq \ell, k \leq 2}$ with $\Lambda_{1,2} = \Lambda_{2,1}$

$$\log p(y|\mathcal{D}) = -\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(y_2 - \mu_2) \Lambda_{1,2} + (y_2 - \mu_2)^2 \Lambda_{2,2} \} + c_{-y}$$

→ Plugging this in the optimal rule,

$$\begin{aligned} q_1(y_1) &\propto \exp \left(-\frac{1}{2} \{ (y_1 - \mu_1)^2 \Lambda_{1,1} + 2(y_1 - \mu_1)(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2} \} \right) \\ &\propto \exp \left(-\frac{1}{2} \{ y_1^2 \Lambda_{1,1} - 2y_1 [\mu_1 \Lambda_{1,1} - (\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}] \} \right) \end{aligned}$$

so that : $q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2) \Lambda_{1,2}, \Lambda_{1,1}^{-1})$

Toy example : Bayesian Linear Regression (BLR) - 3

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

Optimal updates :

$$q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2)\Lambda_{1,2}, \Lambda_{1,1}^{-1})$$

$$q_2(y_2) = \mathcal{N}(y_2; \mu_2 - \Lambda_{2,2}^{-1}(\mathbb{E}_{y_1 \sim q_1}[y_1] - \mu_1)\Lambda_{1,2}, \Lambda_{2,2}^{-1})$$

Setting $m_1 = \mathbb{E}_{y_1 \sim q_1}[y_1]$ and $m_2 = \mathbb{E}_{y_2 \sim q_2}[y_2]$, the CAVI algorithm alternates between :

$$m_1 \leftarrow \mu_1 - \Lambda_{1,1}^{-1}(m_2 - \mu_2)\Lambda_{1,2}$$

$$m_2 \leftarrow \mu_2 - \Lambda_{2,2}^{-1}(m_1 - \mu_1)\Lambda_{1,2}$$

One stable fixed point : $(m_1, m_2) = (\mu_1, \mu_2)$

$\mu = (0 \ 0)$, $\Lambda_{1,1} = \Lambda_{2,2} = 3$ and $\Lambda_{1,2} = -2$.

Toy example : Bayesian Linear Regression (BLR) - 3

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

Optimal updates :

$$q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2)\Lambda_{1,2}, \Lambda_{1,1}^{-1})$$

$$q_2(y_2) = \mathcal{N}(y_2; \mu_2 - \Lambda_{2,2}^{-1}(\mathbb{E}_{y_1 \sim q_1}[y_1] - \mu_1)\Lambda_{1,2}, \Lambda_{2,2}^{-1})$$

Setting $m_1 = \mathbb{E}_{y_1 \sim q_1}[y_1]$ and $m_2 = \mathbb{E}_{y_2 \sim q_2}[y_2]$, the CAVI algorithm alternates between :

$$m_1 \leftarrow \mu_1 - \Lambda_{1,1}^{-1}(m_2 - \mu_2)\Lambda_{1,2}$$

$$m_2 \leftarrow \mu_2 - \Lambda_{2,2}^{-1}(m_1 - \mu_1)\Lambda_{1,2}$$

One stable fixed point : $(m_1, m_2) = (\mu_1, \mu_2)$

$\mu = (0 \ 0)$, $\Lambda_{1,1} = \Lambda_{2,2} = 3$ and $\Lambda_{1,2} = -2$.

Toy example : Bayesian Linear Regression (BLR) - 3

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

Optimal updates :

$$q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2)\Lambda_{1,2}, \Lambda_{1,1}^{-1})$$

$$q_2(y_2) = \mathcal{N}(y_2; \mu_2 - \Lambda_{2,2}^{-1}(\mathbb{E}_{y_1 \sim q_1}[y_1] - \mu_1)\Lambda_{1,2}, \Lambda_{2,2}^{-1})$$

Setting $m_1 = \mathbb{E}_{y_1 \sim q_1}[y_1]$ and $m_2 = \mathbb{E}_{y_2 \sim q_2}[y_2]$, the CAVI algorithm alternates between :

$$m_1 \leftarrow \mu_1 - \Lambda_{1,1}^{-1}(m_2 - \mu_2)\Lambda_{1,2}$$

$$m_2 \leftarrow \mu_2 - \Lambda_{2,2}^{-1}(m_1 - \mu_1)\Lambda_{1,2}$$

One stable fixed point : $(m_1, m_2) = (\mu_1, \mu_2)$

$\mu = (0 \ 0)$, $\Lambda_{1,1} = \Lambda_{2,2} = 3$ and $\Lambda_{1,2} = -2$.

Toy example : Bayesian Linear Regression (BLR) - 3

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

Optimal updates :

$$q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2)\Lambda_{1,2}, \Lambda_{1,1}^{-1})$$

$$q_2(y_2) = \mathcal{N}(y_2; \mu_2 - \Lambda_{2,2}^{-1}(\mathbb{E}_{y_1 \sim q_1}[y_1] - \mu_1)\Lambda_{1,2}, \Lambda_{2,2}^{-1})$$

Setting $m_1 = \mathbb{E}_{y_1 \sim q_1}[y_1]$ and $m_2 = \mathbb{E}_{y_2 \sim q_2}[y_2]$, the CAVI algorithm alternates between :

$$m_1 \leftarrow \mu_1 - \Lambda_{1,1}^{-1}(m_2 - \mu_2)\Lambda_{1,2}$$

$$m_2 \leftarrow \mu_2 - \Lambda_{2,2}^{-1}(m_1 - \mu_1)\Lambda_{1,2}$$

One stable fixed point : $(m_1, m_2) = (\mu_1, \mu_2)$

$\mu = (0 \ 0)$, $\Lambda_{1,1} = \Lambda_{2,2} = 3$ and $\Lambda_{1,2} = -2$.

Toy example : Bayesian Linear Regression (BLR) - 3

$$p(y|\mathcal{D}) = \mathcal{N}(y; \mu, \Lambda^{-1})$$

Optimal updates :

$$q_1(y_1) = \mathcal{N}(y_1; \mu_1 - \Lambda_{1,1}^{-1}(\mathbb{E}_{y_2 \sim q_2}[y_2] - \mu_2)\Lambda_{1,2}, \Lambda_{1,1}^{-1})$$

$$q_2(y_2) = \mathcal{N}(y_2; \mu_2 - \Lambda_{2,2}^{-1}(\mathbb{E}_{y_1 \sim q_1}[y_1] - \mu_1)\Lambda_{1,2}, \Lambda_{2,2}^{-1})$$

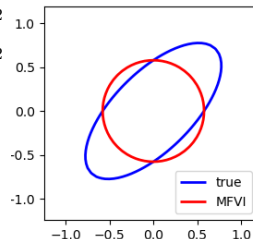
Setting $m_1 = \mathbb{E}_{y_1 \sim q_1}[y_1]$ and $m_2 = \mathbb{E}_{y_2 \sim q_2}[y_2]$, the CAVI algorithm alternates between :

$$m_1 \leftarrow \mu_1 - \Lambda_{1,1}^{-1}(m_2 - \mu_2)\Lambda_{1,2}$$

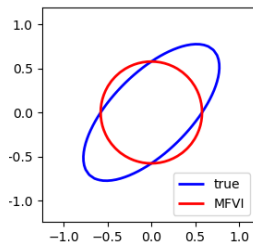
$$m_2 \leftarrow \mu_2 - \Lambda_{2,2}^{-1}(m_1 - \mu_1)\Lambda_{1,2}$$

One stable fixed point : $(m_1, m_2) = (\mu_1, \mu_2)$

$\mu = (0 \ 0)$, $\Lambda_{1,1} = \Lambda_{2,2} = 3$ and $\Lambda_{1,2} = -2$.

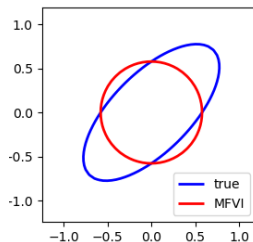


Limitations of MFVI



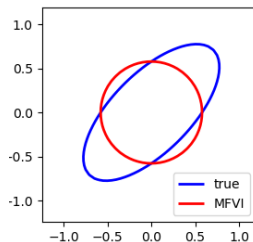
- ❶ The approximative family Q can be **too restrictive** / the updates are **model-specific**.
- ❷ The ELBO tends to **underestimate the posterior variance**.

Limitations of MFVI



- ❶ The approximative family Q can be **too restrictive** / the updates are **model-specific**.
- ❷ The ELBO tends to **underestimate the posterior variance**.

Limitations of MFVI



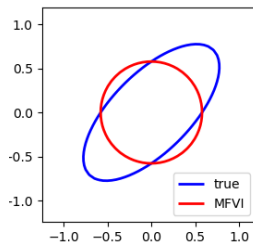
- 1 The approximative family Q can be **too restrictive** / the updates are **model-specific**.

→ **Black-box** Variational Inference

Black Box Variational Inference. R. Ranganath et al. (2014). PMLR.

- 2 The ELBO tends to **underestimate the posterior variance**.

Limitations of MFVI



- 1 The approximative family \mathcal{Q} can be **too restrictive** / the updates are **model-specific**.

→ **Black-box** Variational Inference

Black Box Variational Inference. R. Ranganath et al. (2014). PMLR.

- 2 The ELBO tends to **underestimate the posterior variance**.

→ **Alpha-divergence** Variational Inference

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 Mean-field Variational Inference
- 3 Black-box Variational Inference**
- 4 Alpha-divergence Variational Inference
- 5 Conclusion of Part 1

Black-box Variational Inference (BBVI) - 1

Idea of BBVI : Consider a **parametric** family

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

and perform Gradient **Ascent** on the ELBO with a learning policy $(r_n)_{n \geq 1}$

$$\theta_{n+1} = \theta_n + r_n \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n}$$

We have that :

$$\begin{aligned} \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} &= \nabla_{\theta} \left(\int_{\mathcal{Y}} k(\theta, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta, y)} \right) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= - \int_{\mathcal{Y}} \frac{\partial}{\partial \theta} \left(\frac{k(\theta, y)}{p(y, \mathcal{D})} \log \left(\frac{k(\theta, y)}{p(y, \mathcal{D})} \right) \right) \Big|_{(\theta, y)=(\theta_n, y)} p(y, \mathcal{D}) \nu(dy) \\ &= - \int_{\mathcal{Y}} \left(\log \left(\frac{k(\theta_n, y)}{p(y, \mathcal{D})} \right) + 1 \right) \frac{\partial k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &= \int_{\mathcal{Y}} k(\theta_n, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta_n, y)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \quad (\text{REINFORCE}) \\ &\quad - \int_{\mathcal{Y}} \frac{\partial k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \end{aligned}$$

Black-box Variational Inference (BBVI) - 1

Idea of BBVI : Consider a **parametric** family

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

and perform Gradient **Ascent** on the ELBO with a learning policy $(r_n)_{n \geq 1}$

$$\theta_{n+1} = \theta_n + r_n \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n}$$

We have that :

$$\begin{aligned} \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} &= \nabla_{\theta} \left(\int_{\mathsf{Y}} k(\theta, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta, y)} \right) \nu(\mathrm{d}y) \right) \Big|_{\theta=\theta_n} \\ &= - \int_{\mathsf{Y}} \frac{\partial}{\partial \theta} \left(\frac{k(\theta, y)}{p(y, \mathcal{D})} \log \left(\frac{k(\theta, y)}{p(y, \mathcal{D})} \right) \right) \Big|_{(\theta, y)=(\theta_n, y)} p(y, \mathcal{D}) \nu(\mathrm{d}y) \\ &= - \int_{\mathsf{Y}} \left(\log \left(\frac{k(\theta_n, y)}{p(y, \mathcal{D})} \right) + 1 \right) \frac{\partial k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(\mathrm{d}y) \\ &= \int_{\mathsf{Y}} k(\theta_n, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta_n, y)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(\mathrm{d}y) \quad (\text{REINFORCE}) \\ &\quad - \int_{\mathsf{Y}} \frac{\partial k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(\mathrm{d}y) \end{aligned}$$

Black-box Variational Inference (BBVI) - 1

Idea of BBVI : Consider a **parametric** family

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

and perform Gradient **Ascent** on the ELBO with a learning policy $(r_n)_{n \geq 1}$

$$\theta_{n+1} = \theta_n + r_n \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n}$$

We have that :

$$\begin{aligned} \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} &= \nabla_{\theta} \left(\int_{\mathsf{Y}} k(\theta, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta, y)} \right) \nu(\mathrm{d}y) \right) \Big|_{\theta=\theta_n} \\ &= - \int_{\mathsf{Y}} \frac{\partial}{\partial \theta} \left(\frac{k(\theta, y)}{p(y, \mathcal{D})} \log \left(\frac{k(\theta, y)}{p(y, \mathcal{D})} \right) \right) \Big|_{(\theta, y)=(\theta_n, y)} p(y, \mathcal{D}) \nu(\mathrm{d}y) \\ &= - \int_{\mathsf{Y}} \left(\log \left(\frac{k(\theta_n, y)}{p(y, \mathcal{D})} \right) + 1 \right) \frac{\partial k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(\mathrm{d}y) \\ &= \int_{\mathsf{Y}} k(\theta_n, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta_n, y)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(\mathrm{d}y) \quad (\text{REINFORCE}) \\ &\quad - \int_{\mathsf{Y}} \frac{\partial k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(\mathrm{d}y) \end{aligned}$$

Black-box Variational Inference (BBVI) - 1

Idea of BBVI : Consider a **parametric** family

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

and perform Gradient **Ascent** on the ELBO with a learning policy $(r_n)_{n \geq 1}$

$$\theta_{n+1} = \theta_n + r_n \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n}$$

We have that :

$$\begin{aligned} \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} &= \nabla_{\theta} \left(\int_{\mathcal{Y}} k(\theta, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta, y)} \right) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= - \int_{\mathcal{Y}} \frac{\partial}{\partial \theta} \left(\frac{k(\theta, y)}{p(y, \mathcal{D})} \log \left(\frac{k(\theta, y)}{p(y, \mathcal{D})} \right) \right) \Big|_{(\theta, y)=(\theta_n, y)} p(y, \mathcal{D}) \nu(dy) \\ &= - \int_{\mathcal{Y}} \left(\log \left(\frac{k(\theta_n, y)}{p(y, \mathcal{D})} \right) + 1 \right) \frac{\partial k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &= \int_{\mathcal{Y}} k(\theta_n, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta_n, y)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \quad (\text{REINFORCE}) \\ &\quad - \int_{\mathcal{Y}} \frac{\partial k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \end{aligned}$$

Black-box Variational Inference (BBVI) - 1

Idea of BBVI : Consider a **parametric** family

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

and perform Gradient **Ascent** on the ELBO with a learning policy $(r_n)_{n \geq 1}$

$$\theta_{n+1} = \theta_n + r_n \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n}$$

We have that :

$$\begin{aligned} \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} &= \nabla_{\theta} \left(\int_{\mathcal{Y}} k(\theta, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta, y)} \right) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= - \int_{\mathcal{Y}} \frac{\partial}{\partial \theta} \left(\frac{k(\theta, y)}{p(y, \mathcal{D})} \log \left(\frac{k(\theta, y)}{p(y, \mathcal{D})} \right) \right) \Big|_{(\theta, y)=(\theta_n, y)} p(y, \mathcal{D}) \nu(dy) \\ &= - \int_{\mathcal{Y}} \left(\log \left(\frac{k(\theta_n, y)}{p(y, \mathcal{D})} \right) + 1 \right) \frac{\partial k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &= \int_{\mathcal{Y}} k(\theta_n, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta_n, y)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \quad (\text{REINFORCE}) \\ &\quad - \int_{\mathcal{Y}} \frac{\partial k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \end{aligned}$$

Black-box Variational Inference (BBVI) - 2

Idea of BBVI : Consider a **parametric** family

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

and perform Gradient **Ascent** on the ELBO with a learning policy $(r_n)_{n \geq 1}$

$$\theta_{n+1} = \theta_n + r_n \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n}$$

with

$$\nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \int_{\mathcal{Y}} k(\theta_n, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta_n, y)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y) = (\theta_n, y)} \nu(dy)$$

In practice...

① **Stochastic Gradient Ascent** using the **unbiased** estimate

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y) = (\theta_n, Y_m)}$$

where Y_1, \dots, Y_M : M i.i.d. samples generated from $k(\theta_n, \cdot)$

② Large-scale learning using **mini-batching**

→ convergence towards a **local** optimum of the ELBO ($(r_n)_{n \geq 1}$ follows Robbins-Monro)

Black-box Variational Inference (BBVI) - 2

Idea of BBVI : Consider a **parametric** family

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

and perform Gradient **Ascent** on the ELBO with a learning policy $(r_n)_{n \geq 1}$

$$\theta_{n+1} = \theta_n + r_n \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n}$$

with

$$\nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \int_{\mathcal{Y}} k(\theta_n, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta_n, y)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y) = (\theta_n, y)} \nu(dy)$$

In practice...

① **Stochastic Gradient Ascent** using the **unbiased** estimate

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y) = (\theta_n, Y_m)}$$

where Y_1, \dots, Y_M : M i.i.d. samples generated from $k(\theta_n, \cdot)$

② Large-scale learning using **mini-batching**

→ convergence towards a **local** optimum of the ELBO ($(r_n)_{n \geq 1}$ follows Robbins-Monro)

Black-box Variational Inference (BBVI) - 2

Idea of BBVI : Consider a **parametric** family

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

and perform Gradient **Ascent** on the ELBO with a learning policy $(r_n)_{n \geq 1}$

$$\theta_{n+1} = \theta_n + r_n \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n}$$

with

$$\nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \int_{\mathcal{Y}} k(\theta_n, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta_n, y)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y) = (\theta_n, y)} \nu(dy)$$

In practice...

❶ **Stochastic Gradient Ascent** using the **unbiased** estimate

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y) = (\theta_n, Y_m)}$$

where Y_1, \dots, Y_M : M i.i.d. samples generated from $k(\theta_n, \cdot)$

❷ Large-scale learning using **mini-batching**

→ convergence towards a **local** optimum of the ELBO ($(r_n)_{n \geq 1}$ follows Robbins-Monro)

Black-box Variational Inference (BBVI) - 2

Idea of BBVI : Consider a **parametric** family

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

and perform Gradient **Ascent** on the ELBO with a learning policy $(r_n)_{n \geq 1}$

$$\theta_{n+1} = \theta_n + r_n \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n}$$

with

$$\nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \int_{\mathcal{Y}} k(\theta_n, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta_n, y)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y) = (\theta_n, y)} \nu(dy)$$

In practice...

❶ **Stochastic Gradient Ascent** using the **unbiased** estimate

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y) = (\theta_n, Y_m)}$$

where Y_1, \dots, Y_M : M i.i.d. samples generated from $k(\theta_n, \cdot)$

❷ Large-scale learning using **mini-batching**

→ convergence towards a **local** optimum of the ELBO ($(r_n)_{n \geq 1}$ follows Robbins-Monro)

Black-box Variational Inference (BBVI) - 2

Idea of BBVI : Consider a **parametric** family

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

and perform Gradient **Ascent** on the ELBO with a learning policy $(r_n)_{n \geq 1}$

$$\theta_{n+1} = \theta_n + r_n \nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n}$$

with

$$\nabla_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \int_{\mathcal{Y}} k(\theta_n, y) \log \left(\frac{p(y, \mathcal{D})}{k(\theta_n, y)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y) = (\theta_n, y)} \nu(dy)$$

In practice...

❶ **Stochastic** Gradient Ascent using the **unbiased** estimate

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y) = (\theta_n, Y_m)}$$

where Y_1, \dots, Y_M : M i.i.d. samples generated from $k(\theta_n, \cdot)$

❷ Large-scale learning using **mini-batching**

→ convergence towards a **local** optimum of the ELBO $((r_n)_{n \geq 1})$ follows Robbins-Monro)

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

① The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

② $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces $\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$ “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

① The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

② $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces $\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$ “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

① The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

② $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces **$\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$** “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

① The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

② $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces $\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$ “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

- ❶ The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

- ❷ $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces $\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$ “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

❶ The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

❷ $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces $\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$ “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

❶ The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
 - Control variates
 - Quasi-Monte Carlo methods
 - ... This is an active area of research!

❷ $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces $\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$ “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

❶ The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

❷ $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces $\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$ “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

- ❶ The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

- ❷ $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces $\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$ “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

- ❶ The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

- ❷ $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces $\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$ “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

- ❶ The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

- ❷ $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces $\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$ “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

- 1 The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

- 2 $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces **$\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$** “zero-forcing”

Remarks on BBVI

In short, BBVI resorts to **Stochastic Gradient Ascent** on the ELBO

- The updates are **not** model-specific (“Black-box”)
- \mathcal{Q} : **parametric** family

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

- Uses the **unbiased** estimator

$$\hat{\nabla}_{\theta} \text{ELBO}(k(\theta, \cdot); \mathcal{D})|_{\theta=\theta_n} = \frac{1}{M} \sum_{m=1}^M \log \left(\frac{p(Y_m, \mathcal{D})}{k(\theta_n, Y_m)} \right) \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}$$

where $Y_1, \dots, Y_M : M$ i.i.d. samples generated from $k(\theta_n, \cdot)$

- 1 The variance of the gradient estimators is an **issue** :

- Rao-blackwellisation
- Reparameterisation (used in VAEs)
- Control variates
- Quasi-Monte Carlo methods
- ... This is an active area of research!

- 2 $Y_1 \sim k(\theta_n, \cdot)$ with $p(Y_1, \mathcal{D}) = 0$ makes the gradient blow up...

→ the ELBO enforces **$\text{supp}(k(\theta_n, \cdot)) \subseteq \text{supp}(p(\cdot|\mathcal{D}))$** “zero-forcing”

Outline

- 1 Introduction
- 2 Mean-field Variational Inference
- 3 Black-box Variational Inference
- 4 Alpha-divergence Variational Inference**
- 5 Conclusion of Part 1

The alpha-divergence family

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

\mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$: $\mathbb{Q} \preceq \nu$, $\mathbb{P}_{|\mathcal{D}} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}_{|\mathcal{D}}}{d\nu} = p(\cdot|\mathcal{D})$

Alpha-divergence between \mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$

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

where

$$f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)] , \quad \alpha \in \mathbb{R} \setminus \{0, 1\}$$

The alpha-divergence family

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

\mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$: $\mathbb{Q} \preceq \nu$, $\mathbb{P}_{|\mathcal{D}} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}_{|\mathcal{D}}}{d\nu} = p(\cdot|\mathcal{D})$

Alpha-divergence between \mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$

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

where

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

The alpha-divergence family

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

\mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$: $\mathbb{Q} \preceq \nu$, $\mathbb{P}_{|\mathcal{D}} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}_{|\mathcal{D}}}{d\nu} = p(\cdot|\mathcal{D})$

Alpha-divergence between \mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$

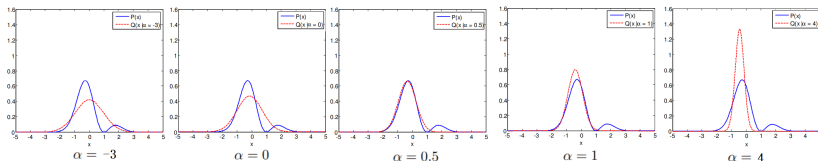
$$D_\alpha(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) = \int_Y f_\alpha \left(\frac{q(y)}{p(y|\mathcal{D})} \right) p(y|\mathcal{D}) \nu(dy) ,$$

where

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

1 A flexible family of divergences...

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



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

The alpha-divergence family

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

\mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$: $\mathbb{Q} \preceq \nu$, $\mathbb{P}_{|\mathcal{D}} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}_{|\mathcal{D}}}{d\nu} = p(\cdot|\mathcal{D})$

Alpha-divergence between \mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$

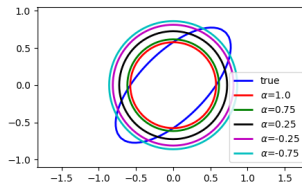
$$D_\alpha(\mathbb{Q}||\mathbb{P}_{|\mathcal{D}}) = \int_Y f_\alpha \left(\frac{q(y)}{p(y|\mathcal{D})} \right) p(y|\mathcal{D}) \nu(dy) ,$$

where

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

- 1 A flexible family of divergences...

Figure: Optimal mean-field approximation for a varying α (BLR revisited)



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

The alpha-divergence family

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

\mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$: $\mathbb{Q} \preceq \nu$, $\mathbb{P}_{|\mathcal{D}} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}_{|\mathcal{D}}}{d\nu} = p(\cdot|\mathcal{D})$

Alpha-divergence between \mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$

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

where

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

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

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

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

- ❸ ...with good **convexity** properties : f_α is convex!

The alpha-divergence family

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

\mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$: $\mathbb{Q} \preceq \nu$, $\mathbb{P}_{|\mathcal{D}} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}_{|\mathcal{D}}}{d\nu} = p(\cdot|\mathcal{D})$

Alpha-divergence between \mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$

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

where

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

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

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

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

- ❸ ...with good **convexity** properties : f_α is convex!

A first approach

Consider a **parametric** family

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

and perform **Gradient Descent** on $\Psi_\alpha(k(\theta, \cdot); p)$

We have that : for all $\alpha \in \mathbb{R} \setminus \{1\}$, $f'_\alpha(u) = \frac{1}{\alpha-1} [u^{\alpha-1} - 1]$ and

$$\begin{aligned}\nabla_\theta \Psi_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \nabla_\theta \left(\int_{\mathcal{Y}} f_\alpha \left(\frac{k(\theta, y)}{p(y)} \right) p(y) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= \dots \\ &= \int_{\mathcal{Y}} \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_n, Y_m)^{\alpha-1} p(Y_m)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}\end{aligned}$$

In practice : **Stochastic** Gradient Descent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation.

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

A first approach

Consider a **parametric** family

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

and perform **Gradient Descent** on $\Psi_\alpha(k(\theta, \cdot); p)$

We have that : for all $\alpha \in \mathbb{R} \setminus \{1\}$, $f'_\alpha(u) = \frac{1}{\alpha-1} [u^{\alpha-1} - 1]$ and

$$\begin{aligned}\nabla_\theta \Psi_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \nabla_\theta \left(\int_{\mathcal{Y}} f_\alpha \left(\frac{k(\theta, y)}{p(y)} \right) p(y) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= \dots \\ &= \int_{\mathcal{Y}} \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_n, Y_m)^{\alpha-1} p(Y_m)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}\end{aligned}$$

In practice : **Stochastic** Gradient Descent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation.

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

A first approach

Consider a **parametric** family

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

and perform **Gradient Descent** on $\Psi_\alpha(k(\theta, \cdot); p)$

We have that : for all $\alpha \in \mathbb{R} \setminus \{1\}$, $f'_\alpha(u) = \frac{1}{\alpha-1} [u^{\alpha-1} - 1]$ and

$$\begin{aligned}\nabla_\theta \Psi_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \nabla_\theta \left(\int_{\mathcal{Y}} f_\alpha \left(\frac{k(\theta, y)}{p(y)} \right) p(y) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= \dots \\ &= \int_{\mathcal{Y}} \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y) = (\theta_n, y)} \nu(dy) \\ &\approx \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_n, Y_m)^{\alpha-1} p(Y_m)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y) = (\theta_n, Y_m)}\end{aligned}$$

In practice : **Stochastic** Gradient Descent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation.

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

A first approach

Consider a **parametric** family

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

and perform **Gradient Descent** on $\Psi_\alpha(k(\theta, \cdot); p)$

We have that : for all $\alpha \in \mathbb{R} \setminus \{1\}$, $f'_\alpha(u) = \frac{1}{\alpha-1} [u^{\alpha-1} - 1]$ and

$$\begin{aligned}\nabla_\theta \Psi_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \nabla_\theta \left(\int_{\mathcal{Y}} f_\alpha \left(\frac{k(\theta, y)}{p(y)} \right) p(y) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= \dots \\ &= \int_{\mathcal{Y}} \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_n, Y_m)^{\alpha-1} p(Y_m)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}\end{aligned}$$

In practice : **Stochastic** Gradient Descent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation.

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

A first approach

Consider a **parametric** family

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

and perform **Gradient Descent** on $\Psi_\alpha(k(\theta, \cdot); p)$

We have that : for all $\alpha \in \mathbb{R} \setminus \{1\}$, $f'_\alpha(u) = \frac{1}{\alpha-1} [u^{\alpha-1} - 1]$ and

$$\begin{aligned}\nabla_\theta \Psi_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \nabla_\theta \left(\int_{\mathcal{Y}} f_\alpha \left(\frac{k(\theta, y)}{p(y)} \right) p(y) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= \dots \\ &= \int_{\mathcal{Y}} \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_n, Y_m)^{\alpha-1} p(Y_m)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}\end{aligned}$$

In practice : **Stochastic** Gradient Descent using $k(\theta_n, \cdot)$ as a sampler

+ Mini-batching + Reparameterisation.

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

A first approach

Consider a **parametric** family

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

and perform **Gradient Descent** on $\Psi_\alpha(k(\theta, \cdot); p)$

We have that : for all $\alpha \in \mathbb{R} \setminus \{1\}$, $f'_\alpha(u) = \frac{1}{\alpha-1} [u^{\alpha-1} - 1]$ and

$$\begin{aligned}\nabla_\theta \Psi_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \nabla_\theta \left(\int_Y f_\alpha \left(\frac{k(\theta, y)}{p(y)} \right) p(y) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= \dots \\ &= \int_Y \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_n, Y_m)^{\alpha-1} p(Y_m)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}\end{aligned}$$

In practice : **Stochastic** Gradient Descent using $k(\theta_n, \cdot)$ as a sampler

+ Mini-batching + Reparameterisation.

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

A first approach

Consider a **parametric** family

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

and perform **Gradient Descent** on $\Psi_\alpha(k(\theta, \cdot); p)$

We have that : for all $\alpha \in \mathbb{R} \setminus \{1\}$, $f'_\alpha(u) = \frac{1}{\alpha-1} [u^{\alpha-1} - 1]$ and

$$\begin{aligned}\nabla_\theta \Psi_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \nabla_\theta \left(\int_Y f_\alpha \left(\frac{k(\theta, y)}{p(y)} \right) p(y) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= \dots \\ &= \int_Y \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_n, Y_m)^{\alpha-1} p(Y_m)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)}\end{aligned}$$

In practice : **Stochastic** Gradient Descent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation.

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

A first approach

Consider a **parametric** family

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

and perform **Gradient Descent** on $\Psi_\alpha(k(\theta, \cdot); p)$

We have that : for all $\alpha \in \mathbb{R} \setminus \{1\}$, $f'_\alpha(u) = \frac{1}{\alpha-1} [u^{\alpha-1} - 1]$ and

$$\begin{aligned}\nabla_\theta \Psi_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \nabla_\theta \left(\int_{\mathcal{Y}} f_\alpha \left(\frac{k(\theta, y)}{p(y)} \right) p(y) \nu(dy) \right) \Big|_{\theta=\theta_n} \\ &= \dots \\ &= \int_{\mathcal{Y}} \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y) = (\theta_n, y)} \nu(dy) \\ &\approx \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_n, Y_m)^{\alpha-1} p(Y_m)^{1-\alpha}}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y) = (\theta_n, Y_m)}\end{aligned}$$

In practice : **Stochastic** Gradient Descent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation.

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

A second approach

Consider a **parametric** family

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

and perform **Gradient Ascent** on the **VR bound** : for all $\alpha \in \mathbb{R} \setminus \{1\}$

$$\mathcal{L}_\alpha(k(\theta, \cdot); p) = \frac{1}{1-\alpha} \log \left(\int_Y k(\theta, y)^\alpha p(y)^{1-\alpha} \nu(dy) \right)$$

→ derived from Rényi's α -divergence, linked to the α -divergence.

$$\begin{aligned} \nabla_\theta \mathcal{L}_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \frac{\alpha}{1-\alpha} \int_Y \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\int_Y k(\theta_n, y')^\alpha p(y')^{1-\alpha} \nu(dy')} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{\alpha}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, Y_m)} \end{aligned}$$

with $w_{n,m} = k(\theta_n, Y_m)^\alpha p(Y_m)^{1-\alpha}$

In practice : **Stochastic** Gradient Ascent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation

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

A second approach

Consider a **parametric** family

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

and perform **Gradient Ascent** on the **VR bound** : for all $\alpha \in \mathbb{R} \setminus \{1\}$

$$\mathcal{L}_\alpha(k(\theta, \cdot); p) = \frac{1}{1-\alpha} \log \left(\int_{\mathcal{Y}} k(\theta, y)^\alpha p(y)^{1-\alpha} \nu(dy) \right)$$

→ derived from Rényi's α -divergence, linked to the α -divergence.

$$\begin{aligned} \nabla_\theta \mathcal{L}_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \frac{\alpha}{1-\alpha} \int_{\mathcal{Y}} \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\int_{\mathcal{Y}} k(\theta_n, y')^\alpha p(y')^{1-\alpha} \nu(dy')} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{\alpha}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, Y_m)} \end{aligned}$$

with $w_{n,m} = k(\theta_n, Y_m)^\alpha p(Y_m)^{1-\alpha}$

In practice : **Stochastic** Gradient Ascent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation

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

A second approach

Consider a **parametric** family

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

and perform **Gradient Ascent** on the **VR bound** : for all $\alpha \in \mathbb{R} \setminus \{1\}$

$$\mathcal{L}_\alpha(k(\theta, \cdot); p) = \frac{1}{1-\alpha} \log \left(\int_{\mathcal{Y}} k(\theta, y)^\alpha p(y)^{1-\alpha} \nu(dy) \right)$$

→ derived from Rényi's α -divergence, linked to the α -divergence.

$$\begin{aligned} \nabla_\theta \mathcal{L}_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \frac{\alpha}{1-\alpha} \int_{\mathcal{Y}} \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\int_{\mathcal{Y}} k(\theta_n, y')^\alpha p(y')^{1-\alpha} \nu(dy')} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{\alpha}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, Y_m)} \end{aligned}$$

with $w_{n,m} = k(\theta_n, Y_m)^\alpha p(Y_m)^{1-\alpha}$

In practice : **Stochastic** Gradient Ascent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation

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

A second approach

Consider a **parametric** family

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

and perform **Gradient Ascent** on the **VR bound** : for all $\alpha \in \mathbb{R} \setminus \{1\}$

$$\mathcal{L}_\alpha(k(\theta, \cdot); p) = \frac{1}{1-\alpha} \log \left(\int_Y k(\theta, y)^\alpha p(y)^{1-\alpha} \nu(dy) \right)$$

→ derived from Rényi's α -divergence, linked to the α -divergence.

$$\begin{aligned} \nabla_\theta \mathcal{L}_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \frac{\alpha}{1-\alpha} \int_Y \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\int_Y k(\theta_n, y')^\alpha p(y')^{1-\alpha} \nu(dy')} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{\alpha}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, Y_m)} \end{aligned}$$

with $w_{n,m} = k(\theta_n, Y_m)^\alpha p(Y_m)^{1-\alpha}$

In practice : **Stochastic** Gradient Ascent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation

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

A second approach

Consider a **parametric** family

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

and perform **Gradient Ascent** on the **VR bound** : for all $\alpha \in \mathbb{R} \setminus \{1\}$

$$\mathcal{L}_\alpha(k(\theta, \cdot); p) = \frac{1}{1-\alpha} \log \left(\int_Y k(\theta, y)^\alpha p(y)^{1-\alpha} \nu(dy) \right)$$

→ derived from Rényi's α -divergence, linked to the α -divergence.

$$\begin{aligned} \nabla_\theta \mathcal{L}_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \frac{\alpha}{1-\alpha} \int_Y \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\int_Y k(\theta_n, y')^\alpha p(y')^{1-\alpha} \nu(dy')} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{\alpha}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, Y_m)} \end{aligned}$$

with $w_{n,m} = k(\theta_n, Y_m)^\alpha p(Y_m)^{1-\alpha}$

In practice : **Stochastic** Gradient Ascent using $k(\theta_n, \cdot)$ as a sampler

+ Mini-batching + Reparameterisation

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

A second approach

Consider a **parametric** family

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

and perform **Gradient Ascent** on the **VR bound** : for all $\alpha \in \mathbb{R} \setminus \{1\}$

$$\mathcal{L}_\alpha(k(\theta, \cdot); p) = \frac{1}{1-\alpha} \log \left(\int_Y k(\theta, y)^\alpha p(y)^{1-\alpha} \nu(dy) \right)$$

→ derived from Rényi's α -divergence, linked to the α -divergence.

$$\begin{aligned} \nabla_\theta \mathcal{L}_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \frac{\alpha}{1-\alpha} \int_Y \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\int_Y k(\theta_n, y')^\alpha p(y')^{1-\alpha} \nu(dy')} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{\alpha}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, Y_m)} \end{aligned}$$

with $w_{n,m} = k(\theta_n, Y_m)^\alpha p(Y_m)^{1-\alpha}$

In practice : **Stochastic** Gradient Ascent using $k(\theta_n, \cdot)$ as a sampler

+ Mini-batching + Reparameterisation

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

A second approach

Consider a **parametric** family

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

and perform **Gradient Ascent** on the **VR bound** : for all $\alpha \in \mathbb{R} \setminus \{1\}$

$$\mathcal{L}_\alpha(k(\theta, \cdot); p) = \frac{1}{1-\alpha} \log \left(\int_Y k(\theta, y)^\alpha p(y)^{1-\alpha} \nu(dy) \right)$$

→ derived from Rényi's α -divergence, linked to the α -divergence.

$$\begin{aligned} \nabla_\theta \mathcal{L}_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \frac{\alpha}{1-\alpha} \int_Y \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\int_Y k(\theta_n, y')^\alpha p(y')^{1-\alpha} \nu(dy')} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{\alpha}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \frac{\partial \log k(\theta, y)}{\partial \theta} \bigg|_{(\theta, y)=(\theta_n, Y_m)} \end{aligned}$$

with $w_{n,m} = k(\theta_n, Y_m)^\alpha p(Y_m)^{1-\alpha}$

In practice : **Stochastic** Gradient Ascent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation

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

A second approach

Consider a **parametric** family

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

and perform **Gradient Ascent** on the **VR bound** : for all $\alpha \in \mathbb{R} \setminus \{1\}$

$$\mathcal{L}_\alpha(k(\theta, \cdot); p) = \frac{1}{1-\alpha} \log \left(\int_{\mathcal{Y}} k(\theta, y)^\alpha p(y)^{1-\alpha} \nu(dy) \right)$$

→ derived from Rényi's α -divergence, linked to the α -divergence.

$$\begin{aligned} \nabla_\theta \mathcal{L}_\alpha(k(\theta, \cdot); p)|_{\theta=\theta_n} &= \frac{\alpha}{1-\alpha} \int_{\mathcal{Y}} \frac{k(\theta_n, y)^\alpha p(y)^{1-\alpha}}{\int_{\mathcal{Y}} k(\theta_n, y')^\alpha p(y')^{1-\alpha} \nu(dy')} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, y)} \nu(dy) \\ &\approx \frac{\alpha}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta, y)=(\theta_n, Y_m)} \end{aligned}$$

with $w_{n,m} = k(\theta_n, Y_m)^\alpha p(Y_m)^{1-\alpha}$

In practice : **Stochastic** Gradient Ascent using $k(\theta_n, \cdot)$ as a sampler
+ Mini-batching + Reparameterisation

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

Alpha-divergence Variational Inference : summary

Alpha-Divergence approach	Rényi's Alpha-Divergence approach
$\inf_{\theta \in \mathcal{T}} \Psi_{\alpha}(k(\theta, \cdot); p)$	$\sup_{\theta \in \mathcal{T}} \mathcal{L}_{\alpha}(k(\theta, \cdot); p)$
$\frac{1}{\alpha-1} \frac{1}{M} \sum_{m=1}^M w_{n,m} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right _{(\theta, y) = (\theta_n, Y_m)}$	$\frac{\alpha}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right _{(\theta, y) = (\theta_n, Y_m)}$

$$\text{Alpha : } \theta_{n+1} = \theta_n + \frac{r_n}{1-\alpha} \frac{1}{M} \sum_{m=1}^M w_{n,m} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right|_{(\theta, y) = (\theta_n, Y_m)}$$

$$\text{Rényi's Alpha : } \theta_{n+1} = \theta_n + \frac{r_n}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right|_{(\theta, y) = (\theta_n, Y_m)}$$

Alpha-divergence Variational Inference : summary

Alpha-Divergence approach	Rényi's Alpha-Divergence approach
$\inf_{\theta \in \mathcal{T}} \Psi_{\alpha}(k(\theta, \cdot); p)$	$\inf_{\theta \in \mathcal{T}} -\alpha^{-1} \mathcal{L}_{\alpha}(k(\theta, \cdot); p)$
$\frac{1}{\alpha-1} \frac{1}{M} \sum_{m=1}^M w_{n,m} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right _{(\theta, y) = (\theta_n, Y_m)}$	$\frac{1}{\alpha-1} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right _{(\theta, y) = (\theta_n, Y_m)}$

$$\text{Alpha : } \theta_{n+1} = \theta_n + \frac{r_n}{1-\alpha} \frac{1}{M} \sum_{m=1}^M w_{n,m} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right|_{(\theta, y) = (\theta_n, Y_m)}$$

$$\text{Rényi's Alpha : } \theta_{n+1} = \theta_n + \frac{r_n}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right|_{(\theta, y) = (\theta_n, Y_m)}$$

Alpha-divergence Variational Inference : summary

Alpha-Divergence approach	Rényi's Alpha-Divergence approach
$\inf_{\theta \in \mathcal{T}} \Psi_{\alpha}(k(\theta, \cdot); p)$	$\inf_{\theta \in \mathcal{T}} -\alpha^{-1} \mathcal{L}_{\alpha}(k(\theta, \cdot); p)$
$\frac{1}{\alpha-1} \frac{1}{M} \sum_{m=1}^M w_{n,m} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right _{(\theta, y) = (\theta_n, Y_m)}$	$\frac{1}{\alpha-1} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right _{(\theta, y) = (\theta_n, Y_m)}$

$$\text{Alpha : } \theta_{n+1} = \theta_n + \frac{r_n}{1-\alpha} \frac{1}{M} \sum_{m=1}^M w_{n,m} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right|_{(\theta, y) = (\theta_n, Y_m)}$$

$$\text{Rényi's Alpha : } \theta_{n+1} = \theta_n + \frac{r_n}{1-\alpha} \sum_{m=1}^M \frac{w_{n,m}}{\sum_{m'=1}^M w_{n,m'}} \left. \frac{\partial \log k(\theta, y)}{\partial \theta} \right|_{(\theta, y) = (\theta_n, Y_m)}$$

Outline

- 1 Introduction
- 2 Mean-field Variational Inference
- 3 Black-box Variational Inference
- 4 Alpha-divergence Variational Inference
- 5 Conclusion of Part 1**

Conclusion of Part 1

Variational Inference : **optimisation-based** methods for Bayesian Inference

Core questions in Variational Inference :

- choice of the **variational family** Q
- choice of the **measure of dissimilarity** D

- ① MFVI : mean-field family, model-specific updates using the ELBO
- ② SVI : scales MFVI to large datasets
- ③ BBVI : parametric family, Stochastic Gradient Ascent on the ELBO
- ④ **Alpha-divergence Variational Inference** : parametric family, extends BBVI to more general objective functions derived from the Alpha-divergence

Conclusion of Part 1

Variational Inference : **optimisation-based** methods for Bayesian Inference

Core questions in Variational Inference :

- choice of the **variational family** Q
- choice of the **measure of dissimilarity** D

- ① MFVI : mean-field family, model-specific updates using the ELBO
- ② SVI : scales MFVI to large datasets
- ③ BBVI : parametric family, Stochastic Gradient Ascent on the ELBO
- ④ **Alpha-divergence Variational Inference** : parametric family, extends BBVI to more general objective functions derived from the Alpha-divergence

Conclusion of Part 1

Variational Inference : **optimisation-based** methods for Bayesian Inference

Core questions in Variational Inference :

- choice of the **variational family** Q
- choice of the **measure of dissimilarity** D

- ① MFVI : mean-field family, model-specific updates using the ELBO
- ② SVI : scales MFVI to large datasets
- ③ BBVI : parametric family, Stochastic Gradient Ascent on the ELBO
- ④ **Alpha-divergence Variational Inference** : parametric family, extends BBVI to more general objective functions derived from the Alpha-divergence

Conclusion of Part 1

Variational Inference : **optimisation-based** methods for Bayesian Inference

Core questions in Variational Inference :

- choice of the **variational family** Q
- choice of the **measure of dissimilarity** D

- ① MFVI : mean-field family, model-specific updates using the ELBO
- ② SVI : scales MFVI to large datasets
- ③ BBVI : parametric family, Stochastic Gradient Ascent on the ELBO
- ④ **Alpha-divergence Variational Inference** : parametric family, extends BBVI to more general objective functions derived from the Alpha-divergence

Conclusion of Part 1

Variational Inference : **optimisation-based** methods for Bayesian Inference

Core questions in Variational Inference :

- choice of the **variational family** Q
- choice of the **measure of dissimilarity** D

- ① MFVI : mean-field family, model-specific updates using the ELBO
- ② SVI : scales MFVI to large datasets
- ③ BBVI : parametric family, Stochastic Gradient Ascent on the ELBO
- ④ **Alpha-divergence Variational Inference** : parametric family, extends BBVI to more general objective functions derived from the Alpha-divergence

Conclusion of Part 1

Variational Inference : **optimisation-based** methods for Bayesian Inference

Core questions in Variational Inference :

- choice of the **variational family** Q
- choice of the **measure of dissimilarity** D

- ➊ MFVI : mean-field family, model-specific updates using the ELBO
- ➋ SVI : scales MFVI to large datasets
- ➌ BBVI : parametric family, Stochastic Gradient Ascent on the ELBO
- ➍ **Alpha-divergence Variational Inference** : parametric family, extends BBVI to more general objective functions derived from the Alpha-divergence

Conclusion of Part 1

Variational Inference : **optimisation-based** methods for Bayesian Inference

Core questions in Variational Inference :

- choice of the **variational family** Q
- choice of the **measure of dissimilarity** D

- ➊ MFVI : mean-field family, model-specific updates using the ELBO
- ➋ SVI : scales MFVI to large datasets
- ➌ BBVI : parametric family, Stochastic Gradient Ascent on the ELBO
- ➍ **Alpha-divergence Variational Inference** : parametric family, extends BBVI to more general objective functions derived from the Alpha-divergence

Conclusion of Part 1

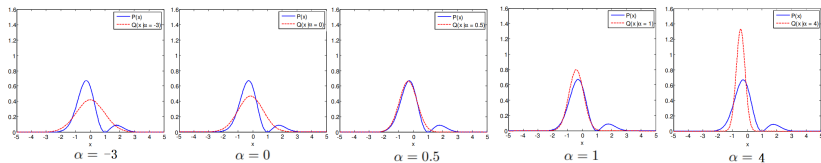
Variational Inference : **optimisation-based** methods for Bayesian Inference

Core questions in Variational Inference :

- choice of the **variational family** Q
- choice of the **measure of dissimilarity** D

- ➊ MFVI : mean-field family, model-specific updates using the ELBO
- ➋ SVI : scales MFVI to large datasets
- ➌ BBVI : parametric family, Stochastic Gradient Ascent on the ELBO
- ➍ **Alpha-divergence Variational Inference** : parametric family, extends BBVI to more general objective functions derived from the Alpha-divergence

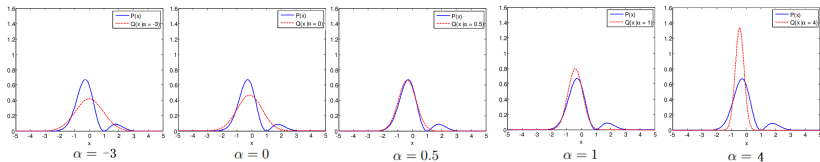
Food for thoughts



Question : Can we further extend the approximating family \mathcal{Q} in the context of Alpha-divergence Variational Inference?

Some answers in Part 2 and 3!

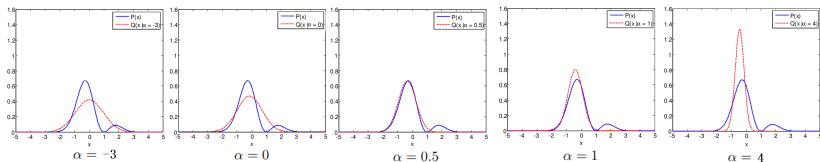
Food for thoughts



Question : Can we further extend the approximating family \mathcal{Q} in the context of Alpha-divergence Variational Inference?

Some answers in Part 2 and 3!

Food for thoughts



Question : Can we further extend the approximating family \mathcal{Q} in the context of Alpha-divergence Variational Inference?

Some answers in Part 2 and 3!

Proof of the optimal rule $q_\ell^*(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$

For notational convenience : $\mathbf{Y} = \mathbf{Y}_\ell \times \mathbf{Y}_{-\ell}$, $q(y) = q_\ell(y_\ell)q_{-\ell}(y_{-\ell})$
and $\nu(dy) = \nu_\ell(dy_\ell)\nu_{-\ell}(dy_{-\ell})$.

$$\begin{aligned}\text{ELBO}(q; \mathcal{D}) &= \int_{\mathbf{Y}} q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) \\&= \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log \left(\frac{p(y, \mathcal{D})}{q_\ell(y_\ell) q_{-\ell}(y_{-\ell})} \right) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \left(\int_{\mathbf{Y}_{-\ell}} q_{-\ell}(y_{-\ell}) \log p(y, \mathcal{D}) \nu_{-\ell}(dy_{-\ell}) \right) \nu_\ell(dy_\ell) \\&\quad - \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log (q_\ell(y_\ell) q_{-\ell}(y_{-\ell})) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&:= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \mathbb{E}_{-\ell} [\log p(y, \mathcal{D})] \nu_\ell(dy_\ell) - \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \log (q_\ell(y_\ell)) \nu_\ell(dy_\ell) + c_{-\ell}\end{aligned}$$

$$\text{ELBO}(q; \mathcal{D}) = \int_{\mathbf{Y}} q_\ell(y_\ell) \log \left(\frac{\exp(\mathbb{E}_{-\ell} [\log p(y, \mathcal{D})])}{q_\ell(y_\ell)} \right) \nu_\ell(dy_\ell) + c_{-\ell}.$$

Proof of the optimal rule $q_\ell^*(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$

For notational convenience : $\mathbf{Y} = \mathbf{Y}_\ell \times \mathbf{Y}_{-\ell}$, $q(y) = q_\ell(y_\ell)q_{-\ell}(y_{-\ell})$
and $\nu(dy) = \nu_\ell(dy_\ell)\nu_{-\ell}(dy_{-\ell})$.

$$\begin{aligned}\text{ELBO}(q; \mathcal{D}) &= \int_{\mathbf{Y}} q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) \\&= \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log \left(\frac{p(y, \mathcal{D})}{q_\ell(y_\ell) q_{-\ell}(y_{-\ell})} \right) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \left(\int_{\mathbf{Y}_{-\ell}} q_{-\ell}(y_{-\ell}) \log p(y, \mathcal{D}) \nu_{-\ell}(dy_{-\ell}) \right) \nu_\ell(dy_\ell) \\&\quad - \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log (q_\ell(y_\ell) q_{-\ell}(y_{-\ell})) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&:= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \mathbb{E}_{-\ell} [\log p(y, \mathcal{D})] \nu_\ell(dy_\ell) - \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \log (q_\ell(y_\ell)) \nu_\ell(dy_\ell) + c_{-\ell}\end{aligned}$$

$$\text{ELBO}(q; \mathcal{D}) = \int_{\mathbf{Y}} q_\ell(y_\ell) \log \left(\frac{\exp(\mathbb{E}_{-\ell} [\log p(y, \mathcal{D})])}{q_\ell(y_\ell)} \right) \nu_\ell(dy_\ell) + c_{-\ell}.$$

Proof of the optimal rule $q_\ell^*(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$

For notational convenience : $\mathbf{Y} = \mathbf{Y}_\ell \times \mathbf{Y}_{-\ell}$, $q(y) = q_\ell(y_\ell)q_{-\ell}(y_{-\ell})$
and $\nu(dy) = \nu_\ell(dy_\ell)\nu_{-\ell}(dy_{-\ell})$.

$$\begin{aligned}\text{ELBO}(q; \mathcal{D}) &= \int_{\mathbf{Y}} q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) \\&= \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log \left(\frac{p(y, \mathcal{D})}{q_\ell(y_\ell) q_{-\ell}(y_{-\ell})} \right) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \left(\int_{\mathbf{Y}_{-\ell}} q_{-\ell}(y_{-\ell}) \log p(y, \mathcal{D}) \nu_{-\ell}(dy_{-\ell}) \right) \nu_\ell(dy_\ell) \\&\quad - \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log (q_\ell(y_\ell) q_{-\ell}(y_{-\ell})) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&:= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \mathbb{E}_{-\ell} [\log p(y, \mathcal{D})] \nu_\ell(dy_\ell) - \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \log (q_\ell(y_\ell)) \nu_\ell(dy_\ell) + c_{-\ell}\end{aligned}$$

$$\text{ELBO}(q; \mathcal{D}) = \int_{\mathbf{Y}} q_\ell(y_\ell) \log \left(\frac{\exp(\mathbb{E}_{-\ell} [\log p(y, \mathcal{D})])}{q_\ell(y_\ell)} \right) \nu_\ell(dy_\ell) + c_{-\ell}.$$

Proof of the optimal rule $q_\ell^*(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$

For notational convenience : $\mathbf{Y} = \mathbf{Y}_\ell \times \mathbf{Y}_{-\ell}$, $q(y) = q_\ell(y_\ell)q_{-\ell}(y_{-\ell})$
and $\nu(dy) = \nu_\ell(dy_\ell)\nu_{-\ell}(dy_{-\ell})$.

$$\begin{aligned}\text{ELBO}(q; \mathcal{D}) &= \int_{\mathbf{Y}} q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) \\&= \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log \left(\frac{p(y, \mathcal{D})}{q_\ell(y_\ell) q_{-\ell}(y_{-\ell})} \right) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \left(\int_{\mathbf{Y}_{-\ell}} q_{-\ell}(y_{-\ell}) \log p(y, \mathcal{D}) \nu_{-\ell}(dy_{-\ell}) \right) \nu_\ell(dy_\ell) \\&\quad - \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log (q_\ell(y_\ell) q_{-\ell}(y_{-\ell})) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&:= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \mathbb{E}_{-\ell} [\log p(y, \mathcal{D})] \nu_\ell(dy_\ell) - \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \log (q_\ell(y_\ell)) \nu_\ell(dy_\ell) + c_{-\ell}\end{aligned}$$

$$\text{ELBO}(q; \mathcal{D}) = \int_{\mathbf{Y}} q_\ell(y_\ell) \log \left(\frac{\exp(\mathbb{E}_{-\ell} [\log p(y, \mathcal{D})])}{q_\ell(y_\ell)} \right) \nu_\ell(dy_\ell) + c_{-\ell}.$$

Proof of the optimal rule $q_\ell^*(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$

For notational convenience : $\mathbf{Y} = \mathbf{Y}_\ell \times \mathbf{Y}_{-\ell}$, $q(y) = q_\ell(y_\ell)q_{-\ell}(y_{-\ell})$
and $\nu(dy) = \nu_\ell(dy_\ell)\nu_{-\ell}(dy_{-\ell})$.

$$\begin{aligned}\text{ELBO}(q; \mathcal{D}) &= \int_{\mathbf{Y}} q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) \\&= \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log \left(\frac{p(y, \mathcal{D})}{q_\ell(y_\ell) q_{-\ell}(y_{-\ell})} \right) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \left(\int_{\mathbf{Y}_{-\ell}} q_{-\ell}(y_{-\ell}) \log p(y, \mathcal{D}) \nu_{-\ell}(dy_{-\ell}) \right) \nu_\ell(dy_\ell) \\&\quad - \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log (q_\ell(y_\ell) q_{-\ell}(y_{-\ell})) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&:= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \mathbb{E}_{-\ell} [\log p(y, \mathcal{D})] \nu_\ell(dy_\ell) - \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \log (q_\ell(y_\ell)) \nu_\ell(dy_\ell) + c_{-\ell}\end{aligned}$$

$$\text{ELBO}(q; \mathcal{D}) = \int_{\mathbf{Y}} q_\ell(y_\ell) \log \left(\frac{\exp(\mathbb{E}_{-\ell} [\log p(y, \mathcal{D})])}{q_\ell(y_\ell)} \right) \nu_\ell(dy_\ell) + c_{-\ell}.$$

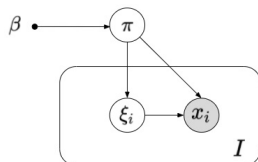
Proof of the optimal rule $q_\ell^*(y_\ell) \propto \exp(\mathbb{E}_{-\ell}[\log p(y, \mathcal{D})])$

For notational convenience : $\mathbf{Y} = \mathbf{Y}_\ell \times \mathbf{Y}_{-\ell}$, $q(y) = q_\ell(y_\ell)q_{-\ell}(y_{-\ell})$
and $\nu(dy) = \nu_\ell(dy_\ell)\nu_{-\ell}(dy_{-\ell})$.

$$\begin{aligned}\text{ELBO}(q; \mathcal{D}) &= \int_{\mathbf{Y}} q(y) \log \left(\frac{p(y, \mathcal{D})}{q(y)} \right) \nu(dy) \\&= \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log \left(\frac{p(y, \mathcal{D})}{q_\ell(y_\ell) q_{-\ell}(y_{-\ell})} \right) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \left(\int_{\mathbf{Y}_{-\ell}} q_{-\ell}(y_{-\ell}) \log p(y, \mathcal{D}) \nu_{-\ell}(dy_{-\ell}) \right) \nu_\ell(dy_\ell) \\&\quad - \int_{\mathbf{Y}_\ell} \int_{\mathbf{Y}_{-\ell}} q_\ell(y_\ell) q_{-\ell}(y_{-\ell}) \log (q_\ell(y_\ell) q_{-\ell}(y_{-\ell})) \nu_{-\ell}(dy_{-\ell}) \nu_\ell(dy_\ell) \\&:= \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \mathbb{E}_{-\ell} [\log p(y, \mathcal{D})] \nu_\ell(dy_\ell) - \int_{\mathbf{Y}_\ell} q_\ell(y_\ell) \log (q_\ell(y_\ell)) \nu_\ell(dy_\ell) + c_{-\ell}\end{aligned}$$

$$\text{ELBO}(q; \mathcal{D}) = \int_{\mathbf{Y}} q_\ell(y_\ell) \log \left(\frac{\exp(\mathbb{E}_{-\ell} [\log p(y, \mathcal{D})])}{q_\ell(y_\ell)} \right) \nu_\ell(dy_\ell) + c_{-\ell} .$$

CAVI for large datasets : Stochastic Variational Inference



- $\mathcal{D} = \{x_1, \dots, x_I\}$, x_1, \dots, x_I : i.i.d. observations
- $y = \{\pi, \xi_1, \dots, \xi_I\}$, π : global latent variable, ξ_1, \dots, ξ_I : local latent variables (β : hyperparameter)

In that case, $p(y, \mathcal{D}) = p(\pi|\beta) \prod_{i=1}^I p(\xi_i|\pi)p(x_i|\xi_i, \pi)$

Mean-field approximation :

$$q(y) = q(\pi|\gamma) \prod_{i=1}^I q(\xi_i|\phi_i)$$

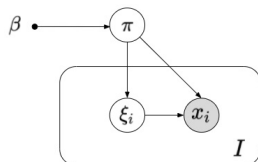
γ : global variational parameter, ϕ_1, \dots, ϕ_I : local variational parameters

Problem : I is often very large (e.g. 1.8M articles from the New York Times)

→ The use of **stochastic** optimisation enabled large scale learning

Stochastic Variational Inference. M. D. Hoffman et al. (2013). JMLR.

CAVI for large datasets : Stochastic Variational Inference



- $\mathcal{D} = \{x_1, \dots, x_I\}$, x_1, \dots, x_I : i.i.d. observations
- $y = \{\pi, \xi_1, \dots, \xi_I\}$, π : global latent variable, ξ_1, \dots, ξ_I : local latent variables (β : hyperparameter)

In that case, $p(y, \mathcal{D}) = p(\pi|\beta) \prod_{i=1}^I p(\xi_i|\pi)p(x_i|\xi_i, \pi)$

Mean-field approximation :

$$q(y) = q(\pi|\gamma) \prod_{i=1}^I q(\xi_i|\phi_i)$$

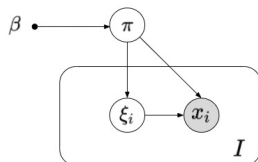
γ : global variational parameter, ϕ_1, \dots, ϕ_I : local variational parameters

Problem : I is often very large (e.g. 1.8M articles from the New York Times)

→ The use of **stochastic** optimisation enabled large scale learning

Stochastic Variational Inference. M. D. Hoffman et al. (2013). JMLR.

CAVI for large datasets : Stochastic Variational Inference



- $\mathcal{D} = \{x_1, \dots, x_I\}$, x_1, \dots, x_I : i.i.d. observations
- $y = \{\pi, \xi_1, \dots, \xi_I\}$, π : global latent variable, ξ_1, \dots, ξ_I : local latent variables (β : hyperparameter)

In that case, $p(y, \mathcal{D}) = p(\pi|\beta) \prod_{i=1}^I p(\xi_i|\pi)p(x_i|\xi_i, \pi)$

Mean-field approximation :

$$q(y) = q(\pi|\gamma) \prod_{i=1}^I q(\xi_i|\phi_i)$$

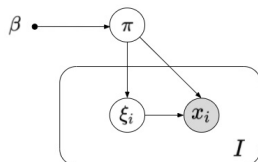
γ : global variational parameter, ϕ_1, \dots, ϕ_I : local variational parameters

Problem : I is often very large (e.g. 1.8M articles from the New York Times)

→ The use of **stochastic** optimisation enabled large scale learning

Stochastic Variational Inference. M. D. Hoffman et al. (2013). JMLR.

CAVI for large datasets : Stochastic Variational Inference



- $\mathcal{D} = \{x_1, \dots, x_I\}$, x_1, \dots, x_I : i.i.d. observations
- $y = \{\pi, \xi_1, \dots, \xi_I\}$, π : global latent variable, ξ_1, \dots, ξ_I : local latent variables (β : hyperparameter)

In that case, $p(y, \mathcal{D}) = p(\pi|\beta) \prod_{i=1}^I p(\xi_i|\pi)p(x_i|\xi_i, \pi)$

Mean-field approximation :

$$q(y) = q(\pi|\gamma) \prod_{i=1}^I q(\xi_i|\phi_i)$$

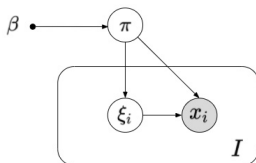
γ : global variational parameter, ϕ_1, \dots, ϕ_I : local variational parameters

Problem : I is often very large (e.g. 1.8M articles from the New York Times)

→ The use of **stochastic** optimisation enabled large scale learning

Stochastic Variational Inference. M. D. Hoffman et al. (2013). JMLR.

CAVI for large datasets : Stochastic Variational Inference



- $\mathcal{D} = \{x_1, \dots, x_I\}$, x_1, \dots, x_I : i.i.d. observations
- $y = \{\pi, \xi_1, \dots, \xi_I\}$, π : global latent variable, ξ_1, \dots, ξ_I : local latent variables (β : hyperparameter)

In that case, $p(y, \mathcal{D}) = p(\pi|\beta) \prod_{i=1}^I p(\xi_i|\pi)p(x_i|\xi_i, \pi)$

Mean-field approximation :

$$q(y) = q(\pi|\gamma) \prod_{i=1}^I q(\xi_i|\phi_i)$$

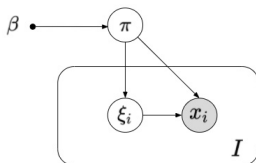
γ : global variational parameter, ϕ_1, \dots, ϕ_I : local variational parameters

Problem : I is often very large (e.g. 1.8M articles from the New York Times)

→ The use of **stochastic** optimisation enabled large scale learning

Stochastic Variational Inference. M. D. Hoffman et al. (2013). JMRL.

CAVI for large datasets : Stochastic Variational Inference



- $\mathcal{D} = \{x_1, \dots, x_I\}$, x_1, \dots, x_I : i.i.d. observations
- $y = \{\pi, \xi_1, \dots, \xi_I\}$, π : global latent variable, ξ_1, \dots, ξ_I : local latent variables (β : hyperparameter)

In that case, $p(y, \mathcal{D}) = p(\pi|\beta) \prod_{i=1}^I p(\xi_i|\pi)p(x_i|\xi_i, \pi)$

Mean-field approximation :

$$q(y) = q(\pi|\gamma) \prod_{i=1}^I q(\xi_i|\phi_i)$$

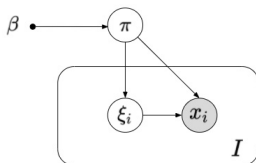
γ : global variational parameter, ϕ_1, \dots, ϕ_I : local variational parameters

Problem : I is often very large (e.g. 1.8M articles from the New York Times)

→ The use of **stochastic** optimisation enabled large scale learning

Stochastic Variational Inference. M. D. Hoffman et al. (2013). JMRL.

CAVI for large datasets : Stochastic Variational Inference



- $\mathcal{D} = \{x_1, \dots, x_I\}$, x_1, \dots, x_I : i.i.d. observations
- $y = \{\pi, \xi_1, \dots, \xi_I\}$, π : global latent variable, ξ_1, \dots, ξ_I : local latent variables (β : hyperparameter)

In that case, $p(y, \mathcal{D}) = p(\pi|\beta) \prod_{i=1}^I p(\xi_i|\pi)p(x_i|\xi_i, \pi)$

Mean-field approximation :

$$q(y) = q(\pi|\gamma) \prod_{i=1}^I q(\xi_i|\phi_i)$$

γ : global variational parameter, ϕ_1, \dots, ϕ_I : local variational parameters

Problem : I is often very large (e.g. 1.8M articles from the New York Times)

→ The use of **stochastic** optimisation enabled large scale learning

Stochastic Variational Inference. M. D. Hoffman et al. (2013). JMLR.

Variational Inference

Foundations and recent advances

(Part 2)

Kamélia Daudel



University of Bristol – 09/03/2022

Reminder - 1

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

\mathbb{Q} and $\mathbb{P}_{|\mathcal{D}} : \mathbb{Q} \preceq \nu, \mathbb{P}_{|\mathcal{D}} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q, \frac{d\mathbb{P}_{|\mathcal{D}}}{d\nu} = p(\cdot|\mathcal{D}) = \frac{p(\cdot, \mathcal{D})}{p(\mathcal{D})}$

- Variational Inference optimisation problem :

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

where \mathcal{Q} is the variational family and D is the measure of dissimilarity

- **Alpha-Divergence Variational Inference** : Two possible objective functions

$$\begin{aligned}\Psi_{\alpha}(q; p) &= \int_Y f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) \\ \mathcal{L}_{\alpha}(q; p) &= \frac{1}{1-\alpha} \log \left(\int_Y q(y)^{\alpha} p(y)^{1-\alpha} \nu(dy) \right)\end{aligned}$$

with $p = p(\cdot, \mathcal{D})$ and

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

Reminder - 1

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

\mathbb{Q} and $\mathbb{P}_{|\mathcal{D}} : \mathbb{Q} \preceq \nu, \mathbb{P}_{|\mathcal{D}} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q, \frac{d\mathbb{P}_{|\mathcal{D}}}{d\nu} = p(\cdot|\mathcal{D}) = \frac{p(\cdot, \mathcal{D})}{p(\mathcal{D})}$

- Variational Inference optimisation problem :

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

where \mathcal{Q} is the variational family and D is the measure of dissimilarity

- **Alpha-Divergence Variational Inference** : Two possible objective functions

$$\Psi_{\alpha}(q; p) = \int_Y f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$$

$$\mathcal{L}_{\alpha}(q; p) = \frac{1}{1 - \alpha} \log \left(\int_Y q(y)^{\alpha} p(y)^{1-\alpha} \nu(dy) \right)$$

with $p = p(\cdot, \mathcal{D})$ and

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

Reminder - 1

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

\mathbb{Q} and $\mathbb{P}_{|\mathcal{D}}$: $\mathbb{Q} \preceq \nu$, $\mathbb{P}_{|\mathcal{D}} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q$, $\frac{d\mathbb{P}_{|\mathcal{D}}}{d\nu} = p(\cdot|\mathcal{D}) = \frac{p(\cdot, \mathcal{D})}{p(\mathcal{D})}$

- Variational Inference optimisation problem :

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

where \mathcal{Q} is the variational family and D is the measure of dissimilarity

- **Alpha-Divergence Variational Inference** : Two possible objective functions

$$\Psi_{\alpha}(q; p) = \int_Y f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$$

$$\mathcal{L}_{\alpha}(q; p) = \frac{1}{1 - \alpha} \log \left(\int_Y q(y)^{\alpha} p(y)^{1-\alpha} \nu(dy) \right)$$

with $p = p(\cdot, \mathcal{D})$ and

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

Reminder - 1

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

\mathbb{Q} and $\mathbb{P}_{|\mathcal{D}} : \mathbb{Q} \preceq \nu, \mathbb{P}_{|\mathcal{D}} \preceq \nu$ with $\frac{d\mathbb{Q}}{d\nu} = q, \frac{d\mathbb{P}_{|\mathcal{D}}}{d\nu} = p(\cdot|\mathcal{D}) = \frac{p(\cdot, \mathcal{D})}{p(\mathcal{D})}$

- Variational Inference optimisation problem :

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

where \mathcal{Q} is the variational family and D is the measure of dissimilarity

- **Alpha-Divergence Variational Inference** : Two possible objective functions

$$\begin{aligned}\Psi_{\alpha}(q; p) &= \int_Y f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy) \\ -\alpha^{-1} \mathcal{L}_{\alpha}(q; p) &= \frac{1}{\alpha(\alpha - 1)} \log \left(\int_Y q(y)^{\alpha} p(y)^{1-\alpha} \nu(dy) \right)\end{aligned}$$

with $p = p(\cdot, \mathcal{D})$ and

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

Reminder - 2

When \mathcal{Q} is parametric,

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

we can perform Stochastic Gradient Descent w.r.t θ on $\Psi_\alpha(q; p)$
(resp. $-\alpha^{-1}\mathcal{L}_\alpha(q; p)$)

Question : Can we further extend the approximating family \mathcal{Q} in the context of Alpha-divergence Variational Inference?

Reminder - 2

When \mathcal{Q} is parametric,

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

we can perform Stochastic Gradient Descent w.r.t θ on $\Psi_\alpha(q; p)$
(resp. $-\alpha^{-1}\mathcal{L}_\alpha(q; p)$)

Question : Can we further extend the approximating family \mathcal{Q} in the context of Alpha-divergence Variational Inference?

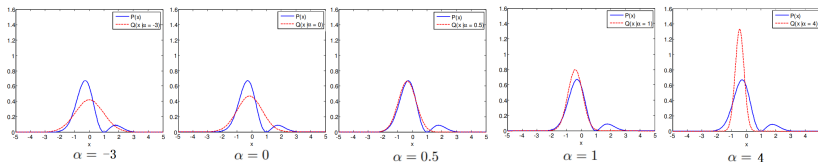
Reminder - 2

When \mathcal{Q} is parametric,

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

we can perform Stochastic Gradient Descent w.r.t θ on $\Psi_\alpha(q; p)$
(resp. $-\alpha^{-1} \mathcal{L}_\alpha(q; p)$)

Question : Can we further extend the approximating family \mathcal{Q} in the context of Alpha-divergence Variational Inference?



Outline

- ① Infinite-dimensional Alpha-divergence minimisation
- ② Numerical experiments
- ③ Conclusion of Part 2

Outline

- 1 Infinite-dimensional Alpha-divergence minimisation
- 2 Numerical experiments
- 3 Conclusion of Part 2

Infinite-dimensional Alpha-divergence minimisation

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

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

<https://doi.org/10.1214/20-AOS2035>.

Mixture weights optimisation for Alpha-Divergence Variational Inference.

K. Daudel and R. Douc (2021). To appear in NeurIPS2021

Idea : Extend the traditional variational parametric family

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

by putting a prior on the variational parameter θ

$$\mathcal{Q} = \left\{ q : y \mapsto \mu k(y) := \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) : \mu \in \mathcal{M} \right\}$$

and propose an update formula for μ that ensures a systematic decrease in $\mu \mapsto \Psi_{\alpha}(\mu k; p)$ at each step

→ Finite Mixture Models : $\mu = \sum_{j=1}^J \lambda_j \delta_{\theta_j}$

Infinite-dimensional Alpha-divergence minimisation

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

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

<https://doi.org/10.1214/20-AOS2035>.

Mixture weights optimisation for Alpha-Divergence Variational Inference.

K. Daudel and R. Douc (2021). To appear in NeurIPS2021

Idea : Extend the traditional variational parametric family

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

by putting a prior on the variational parameter θ

$$\mathcal{Q} = \left\{ q : y \mapsto \mu k(y) := \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) : \mu \in \mathcal{M} \right\}$$

and propose an update formula for μ that ensures a **systematic decrease** in $\mu \mapsto \Psi_{\alpha}(\mu k; p)$ at each step

$$\rightarrow \text{Finite Mixture Models : } \mu = \sum_{j=1}^J \lambda_j \delta_{\theta_j}$$

Infinite-dimensional Alpha-divergence minimisation

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

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

<https://doi.org/10.1214/20-AOS2035>.

Mixture weights optimisation for Alpha-Divergence Variational Inference.

K. Daudel and R. Douc (2021). To appear in NeurIPS2021

Idea : Extend the traditional variational parametric family

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

by putting a prior on the variational parameter θ

$$\mathcal{Q} = \left\{ q : y \mapsto \mu k(y) := \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) : \mu \in \mathbb{M} \right\}$$

and propose an update formula for μ that ensures a **systematic decrease** in $\mu \mapsto \Psi_{\alpha}(\mu k; p)$ at each step

$$\rightarrow \text{Finite Mixture Models : } \mu = \sum_{j=1}^J \lambda_j \delta_{\theta_j}$$

Infinite-dimensional Alpha-divergence minimisation

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

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

<https://doi.org/10.1214/20-AOS2035>.

Mixture weights optimisation for Alpha-Divergence Variational Inference.

K. Daudel and R. Douc (2021). To appear in NeurIPS2021

Idea : Extend the traditional variational parametric family

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

by putting a prior on the variational parameter θ

$$\mathcal{Q} = \left\{ q : y \mapsto \mu k(y) := \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) : \mu \in \mathbb{M} \right\}$$

and propose an update formula for μ that ensures a **systematic decrease** in $\mu \mapsto \Psi_{\alpha}(\mu k; p)$ at each step

$$\rightarrow \text{Finite Mixture Models} : \mu = \sum_{j=1}^J \lambda_j \delta_{\theta_j}$$

The (α, Γ) -descent algorithm

Optimisation problem

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

- p is a **nonnegative measurable function** defined on $(\mathcal{Y}, \mathcal{Y})$
- \mathcal{M} is a subset of $\mathcal{M}_1(\mathcal{T})$, the space of probability measures on \mathcal{T}
- $K : (\theta, A) \mapsto \int_A k(\theta, y) \nu(dy)$ is a Markov transition kernel defined on $\mathcal{T} \times \mathcal{Y}$ with density k

Algorithm

Let $\mu_1 \in \mathcal{M}_1(\mathcal{T})$ be such that $\Psi_{\alpha}(\mu_1 k) < \infty$. The sequence of probability measures $(\mu_n)_{n \geq 1}$ is defined iteratively by

$$\mu_{n+1} = \mathcal{I}_{\alpha}(\mu_n), \quad n \geq 1$$

where for all $\mu \in \mathcal{M}_1(\mathcal{T})$ and all $\theta \in \mathcal{T}$,

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

The (α, Γ) -descent algorithm

Optimisation problem

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

- p is a **nonnegative measurable function** defined on $(\mathcal{Y}, \mathcal{Y})$
- \mathcal{M} is a subset of $\mathcal{M}_1(\mathcal{T})$, the space of probability measures on \mathcal{T}
- $K : (\theta, A) \mapsto \int_A k(\theta, y) \nu(dy)$ is a Markov transition kernel defined on $\mathcal{T} \times \mathcal{Y}$ with density k

Algorithm

Let $\mu_1 \in \mathcal{M}_1(\mathcal{T})$ be such that $\Psi_{\alpha}(\mu_1 k) < \infty$. The sequence of probability measures $(\mu_n)_{n \geq 1}$ is defined iteratively by

$$\mu_{n+1} = \mathcal{I}_{\alpha}(\mu_n), \quad n \geq 1$$

where for all $\mu \in \mathcal{M}_1(\mathcal{T})$ and all $\theta \in \mathcal{T}$,

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

The (α, Γ) -descent algorithm

Optimisation problem

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

- p is a **nonnegative measurable function** defined on $(\mathcal{Y}, \mathcal{Y})$
- \mathcal{M} is a subset of $\mathcal{M}_1(\mathcal{T})$, the space of probability measures on \mathcal{T}
- $K : (\theta, A) \mapsto \int_A k(\theta, y) \nu(dy)$ is a Markov transition kernel defined on $\mathcal{T} \times \mathcal{Y}$ with density k

Algorithm

Let $\mu_1 \in \mathcal{M}_1(\mathcal{T})$ be such that $\Psi_{\alpha}(\mu_1 k) < \infty$. The sequence of probability measures $(\mu_n)_{n \geq 1}$ is defined iteratively by

$$\mu_{n+1} = \mathcal{I}_{\alpha}(\mu_n), \quad n \geq 1$$

where for all $\mu \in \mathcal{M}_1(\mathcal{T})$ and all $\theta \in \mathcal{T}$,

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

The (α, Γ) -descent algorithm

Optimisation problem

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

- p is a **nonnegative measurable function** defined on $(\mathcal{Y}, \mathcal{Y})$
- \mathcal{M} is a subset of $\mathcal{M}_1(\mathcal{T})$, the space of probability measures on \mathcal{T}
- $K : (\theta, A) \mapsto \int_A k(\theta, y) \nu(dy)$ is a Markov transition kernel defined on $\mathcal{T} \times \mathcal{Y}$ with density k

Algorithm

Let $\mu_1 \in \mathcal{M}_1(\mathcal{T})$ be such that $\Psi_{\alpha}(\mu_1 k) < \infty$. The sequence of probability measures $(\mu_n)_{n \geq 1}$ is defined iteratively by

$$\mu_{n+1} = \mathcal{I}_{\alpha}(\mu_n), \quad n \geq 1$$

where for all $\mu \in \mathcal{M}_1(\mathcal{T})$ and all $\theta \in \mathcal{T}$,

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

The (α, Γ) -descent algorithm

Optimisation problem

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

- p is a **nonnegative measurable function** defined on $(\mathcal{Y}, \mathcal{Y})$
- \mathcal{M} is a subset of $\mathcal{M}_1(\mathcal{T})$, the space of probability measures on \mathcal{T}
- $K : (\theta, A) \mapsto \int_A k(\theta, y) \nu(dy)$ is a Markov transition kernel defined on $\mathcal{T} \times \mathcal{Y}$ with density k

Algorithm

Let $\mu_1 \in \mathcal{M}_1(\mathcal{T})$ be such that $\Psi_{\alpha}(\mu_1 k) < \infty$. The sequence of probability measures $(\mu_n)_{n \geq 1}$ is defined iteratively by

$$\mu_{n+1} = \mathcal{I}_{\alpha}(\mu_n), \quad n \geq 1$$

where for all $\mu \in \mathcal{M}_1(\mathcal{T})$ and all $\theta \in \mathcal{T}$,

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

The (α, Γ) -descent algorithm

Optimisation problem

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

- p is a **nonnegative measurable function** defined on $(\mathcal{Y}, \mathcal{Y})$
- \mathcal{M} is a subset of $\mathcal{M}_1(\mathcal{T})$, the space of probability measures on \mathcal{T}
- $K : (\theta, A) \mapsto \int_A k(\theta, y) \nu(dy)$ is a Markov transition kernel defined on $\mathcal{T} \times \mathcal{Y}$ with density k

Algorithm

Let $\mu_1 \in \mathcal{M}_1(\mathcal{T})$ be such that $\Psi_{\alpha}(\mu_1 k) < \infty$. The sequence of probability measures $(\mu_n)_{n \geq 1}$ is defined iteratively by

$$\mu_{n+1} = \mathcal{I}_{\alpha}(\mu_n), \quad n \geq 1$$

where for all $\mu \in \mathcal{M}_1(\mathcal{T})$ and all $\theta \in \mathcal{T}$,

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

Conditions for a monotonic decrease

(A1) For all $(\theta, y) \in \mathsf{T} \times \mathsf{Y}$, $k(\theta, y) > 0$, $p(y) \geq 0$ 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

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

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

Conditions for a monotonic decrease

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

(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

Assume (A1) and (A2). Let $\mu \in \mathcal{M}_1(\mathcal{T})$ be such that $\Psi_{\alpha}(\mu k) < \infty$ and $\mu(\Gamma(b_{\mu, \alpha} + \kappa)) < \infty$. Then,

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

Conditions for a monotonic decrease

(A1) For all $(\theta, y) \in \mathsf{T} \times \mathsf{Y}$, $k(\theta, y) > 0$, $p(y) \geq 0$ 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

Assume (A1) and (A2). Let $\mu \in M_1(\mathsf{T})$ be such that $\Psi_{\alpha}(\mu k) < \infty$ and $\mu(\Gamma(b_{\mu, \alpha} + \kappa)) < \infty$. Then,

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

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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



First idea

By convexity of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{\zeta k(y)}{p(y)} \right) + f'_\alpha \left(\frac{\zeta k(y)}{p(y)} \right) \frac{\mu k(y) - \zeta k(y)}{p(y)}$$

✗ Not the best idea!

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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



First idea

By **convexity** of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{\zeta k(y)}{p(y)} \right) + f'_\alpha \left(\frac{\zeta k(y)}{p(y)} \right) \frac{\mu k(y) - \zeta k(y)}{p(y)}$$

✗ Not the best idea!

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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



First idea

By convexity of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{\zeta k(y)}{p(y)} \right) + f'_\alpha \left(\frac{\zeta k(y)}{p(y)} \right) \frac{\mu k(y) - \zeta k(y)}{p(y)}$$

✗ Not the best idea!

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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



First idea

By convexity of f_α ,

$$f_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \geq f_\alpha \left(\frac{\zeta k(y)}{p(y)} \right) + f'_\alpha \left(\frac{\zeta k(y)}{p(y)} \right) \frac{\mu k(y) - \zeta k(y)}{p(y)}$$

✗ Not the best idea!

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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



Second idea

By **convexity** of f_α : for all $y \in Y$

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

→ Next, we integrate w.r.t to $\frac{\mu(d\theta)k(\theta, y)}{\mu k(y)}$,

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

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

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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



Second idea

By **convexity** of f_α : for all $y \in Y$

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

→ Next, we integrate w.r.t to $\frac{\mu(d\theta)k(\theta, y)}{\mu k(y)}$,

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

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

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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



Second idea

By **convexity** of f_α : for all $y \in Y$

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

→ Next, we integrate w.r.t to $\frac{\mu(d\theta)k(\theta, y)}{\mu k(y)}$,

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

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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



Second idea

By **convexity** of f_α : for all $y \in Y$

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

→ Next, we integrate w.r.t to $\frac{\mu(d\theta)k(\theta, y)}{\mu k(y)}$,

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

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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

→ At this stage : for all $y \in \mathcal{Y}$,

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

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

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

Choice of A_α

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

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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

→ At this stage : for all $y \in \mathcal{Y}$,

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

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

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

Choice of A_α

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

Proof : 1) Proving a general lower bound

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

We want to find A_α such that

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

and equality holds iff $\zeta = \mu$.

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

→ At this stage : for all $y \in \mathcal{Y}$,

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

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

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

Choice of A_α

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

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

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

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

and thus

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

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

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

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

and thus

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

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

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

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

and thus

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

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

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

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

and thus

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

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

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

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

and thus

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ 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

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

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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$!

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v) / \mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!

□

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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)!$

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v) / \mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!

□

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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$!

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v) / \mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!

□

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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$!

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v) / \mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!

□

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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$!

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v) / \mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!

□

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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$!

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v)/\mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!

□

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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$!

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v) / \mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!

□

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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$!

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v) / \mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!



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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$!

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v) / \mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!



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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$!

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v) / \mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!



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

Proving that $A_\alpha \geq 0 \rightarrow$ We treat the case $\alpha \in \mathbb{R} \setminus \{1\}$.

We have obtained that

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

It's time to use that $g \propto \Gamma(b_{\mu,\alpha} + \kappa)$!

(i) Let V be the random variable $V(\theta) = b_{\mu,\alpha}(\theta) + \kappa$ (probability space $(\mathcal{T}, \mathcal{T}, \mu)$)

(ii) Set $\tilde{\Gamma}(v) = \Gamma(v) / \mu(\Gamma(b_{\mu,\alpha} + \kappa))$ for all $v \in \text{Dom}_\alpha$.

Then,

$$\begin{aligned} A_\alpha &= \mathbb{E} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V) [1 - \tilde{\Gamma}(V)] \right) \\ &= \mathbb{Cov} \left(\left[V - \kappa + \frac{1}{\alpha - 1} \right] \tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V) \right) \quad \text{since } \mathbb{E}[1 - \tilde{\Gamma}(V)] = 0 \end{aligned}$$

(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: $A_\alpha \geq 0$!

□

Reminder : Conditions for a monotonic decrease

(A1) For all $(\theta, y) \in \mathsf{T} \times \mathsf{Y}$, $k(\theta, y) > 0$, $p(y) \geq 0$ 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

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

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

Examples satisfying (A2)

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

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

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

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

$$\mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\eta \int_Y k(\theta, y) \log \left(\frac{\mu_n k(y)}{p(y)} \right) \nu(\mathrm{d}y) \right]$$

→ NB : η corresponds to the **learning rate**

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

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

$$\mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \left[\int_Y k(\theta, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1} \nu(\mathrm{d}y) + (\alpha - 1)\kappa \right]^{\frac{\eta}{1-\alpha}}$$

Examples satisfying (A2)

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

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

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

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

$$\mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\eta \int_Y k(\theta, y) \log \left(\frac{\mu_n k(y)}{p(y)} \right) \nu(\mathrm{d}y) \right]$$

→ NB : η corresponds to the **learning rate**

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

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

$$\mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \left[\int_Y k(\theta, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1} \nu(\mathrm{d}y) + (\alpha - 1)\kappa \right]^{\frac{\eta}{1-\alpha}}$$

Examples satisfying (A2)

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

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

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

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

$$\mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\eta \int_Y k(\theta, y) \log \left(\frac{\mu_n k(y)}{p(y)} \right) \nu(\mathrm{d}y) \right]$$

→ NB : η corresponds to the **learning rate**

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

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

$$\mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \left[\int_Y k(\theta, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1} \nu(\mathrm{d}y) + (\alpha - 1)\kappa \right]^{\frac{\eta}{1-\alpha}}$$

Examples satisfying (A2)

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

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

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

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

$$\mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\eta \int_Y k(\theta, y) \log \left(\frac{\mu_n k(y)}{p(y)} \right) \nu(\mathrm{d}y) \right]$$

→ NB : η corresponds to the **learning rate**

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

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

$$\mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \left[\int_Y k(\theta, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1} \nu(\mathrm{d}y) + (\alpha - 1)\kappa \right]^{\frac{\eta}{1-\alpha}}$$

Convergence results

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L-smooth and $-\log \Gamma$ is concave increasing

- Entropic Mirror Descent : $\Gamma(v) = e^{-\eta v}$ with $\eta \in (0, \frac{1}{|\alpha-1||b|_{\infty, \alpha}+1})$, any α, κ
- Power Descent : $\Gamma(v) = [(\alpha-1)v + 1]^{\eta/(1-\alpha)}$ with $\eta \in (0, 1]$, $\alpha > 1, \kappa > 0$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \leq 0$

Under additional assumptions on Ψ_α and $b_{\mu, \alpha}$, if $(\mu_n)_{n \geq 1}$ weakly converges to μ^* , then :

$$\mu^* \text{ is a fixed point of } \mathcal{I}_\alpha \text{ and } \Psi_\alpha(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_\alpha(\zeta k)$$

Convergence results

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L-smooth and $-\log \Gamma$ is concave increasing

- Entropic Mirror Descent : $\Gamma(v) = e^{-\eta v}$ with $\eta \in (0, \frac{1}{|\alpha-1||b|_{\infty, \alpha}+1})$, any α, κ
- Power Descent : $\Gamma(v) = [(\alpha-1)v+1]^{\eta/(1-\alpha)}$ with $\eta \in (0, 1]$, $\alpha > 1$, $\kappa > 0$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \leq 0$

Under additional assumptions on Ψ_α and $b_{\mu, \alpha}$, if $(\mu_n)_{n \geq 1}$ weakly converges to μ^* , then :

$$\mu^* \text{ is a fixed point of } \mathcal{I}_\alpha \text{ and } \Psi_\alpha(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_\alpha(\zeta k)$$

Convergence results

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L-smooth and $-\log \Gamma$ is concave increasing

- Entropic Mirror Descent : $\Gamma(v) = e^{-\eta v}$ with $\eta \in (0, \frac{1}{|\alpha-1||b|_{\infty, \alpha}+1})$, any α, κ
- Power Descent : $\Gamma(v) = [(\alpha-1)v+1]^{\eta/(1-\alpha)}$ with $\eta \in (0, 1]$, $\alpha > 1$, $\kappa > 0$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \leq 0$

Under additional assumptions on Ψ_α and $b_{\mu, \alpha}$, if $(\mu_n)_{n \geq 1}$ weakly converges to μ^* , then :

$$\mu^* \text{ is a fixed point of } \mathcal{I}_\alpha \text{ and } \Psi_\alpha(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_\alpha(\zeta k)$$

Convergence results

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L-smooth and $-\log \Gamma$ is concave increasing

- Entropic Mirror Descent : $\Gamma(v) = e^{-\eta v}$ with $\eta \in (0, \frac{1}{|\alpha-1||b|_{\infty, \alpha}+1})$, any α, κ
- Power Descent : $\Gamma(v) = [(\alpha-1)v + 1]^{\eta/(1-\alpha)}$ with $\eta \in (0, 1]$, $\alpha > 1$, $\kappa > 0$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \leq 0$

Under additional assumptions on Ψ_α and $b_{\mu, \alpha}$, if $(\mu_n)_{n \geq 1}$ weakly converges to μ^* , then :

$$\mu^* \text{ is a fixed point of } \mathcal{I}_\alpha \text{ and } \Psi_\alpha(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_\alpha(\zeta k)$$

Convergence results

$$\mu_{n+1}(\mathrm{d}\theta) = \frac{\mu_n(\mathrm{d}\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L-smooth and $-\log \Gamma$ is concave increasing

- Entropic Mirror Descent : $\Gamma(v) = e^{-\eta v}$ with $\eta \in (0, \frac{1}{|\alpha-1||b|_{\infty, \alpha}+1})$, any α, κ
- Power Descent : $\Gamma(v) = [(\alpha-1)v + 1]^{\eta/(1-\alpha)}$ with $\eta \in (0, 1]$, $\alpha > 1, \kappa > 0$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \leq 0$

Under additional assumptions on Ψ_α and $b_{\mu, \alpha}$, if $(\mu_n)_{n \geq 1}$ weakly converges to μ^* , then :

$$\mu^* \text{ is a fixed point of } \mathcal{I}_\alpha \text{ and } \Psi_\alpha(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_\alpha(\zeta k)$$

Convergence results

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L-smooth and $-\log \Gamma$ is concave increasing

- Entropic Mirror Descent : $\Gamma(v) = e^{-\eta v}$ with $\eta \in (0, \frac{1}{|\alpha-1||b|_{\infty, \alpha}+1})$, any α, κ
- Power Descent : $\Gamma(v) = [(\alpha-1)v + 1]^{\eta/(1-\alpha)}$ with $\eta \in (0, 1]$, $\alpha > 1, \kappa > 0$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \leq 0$

Under additional assumptions on Ψ_α and $b_{\mu, \alpha}$, if $(\mu_n)_{n \geq 1}$ weakly converges to μ^* , then :

$$\mu^* \text{ is a fixed point of } \mathcal{I}_\alpha \text{ and } \Psi_\alpha(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_\alpha(\zeta k)$$

Convergence results

$$\mu_{n+1}(\mathrm{d}\theta) = \frac{\mu_n(\mathrm{d}\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L-smooth and $-\log \Gamma$ is concave increasing

- Entropic Mirror Descent : $\Gamma(v) = e^{-\eta v}$ with $\eta \in (0, \frac{1}{|\alpha-1||b|_{\infty, \alpha}+1})$, any α, κ
- Power Descent : $\Gamma(v) = [(\alpha-1)v + 1]^{\eta/(1-\alpha)}$ with $\eta \in (0, 1]$, $\alpha > 1, \kappa > 0$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \leq 0$

Under additional assumptions on Ψ_α and $b_{\mu, \alpha}$, if $(\mu_n)_{n \geq 1}$ weakly converges to μ^* , then :

$$\mu^* \text{ is a fixed point of } \mathcal{I}_\alpha \text{ and } \Psi_\alpha(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_\alpha(\zeta k)$$

Convergence results

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L-smooth and $-\log \Gamma$ is concave increasing

- Entropic Mirror Descent : $\Gamma(v) = e^{-\eta v}$ with $\eta \in (0, \frac{1}{|\alpha-1||b|_{\infty, \alpha}+1})$, any α, κ
- Power Descent : $\Gamma(v) = [(\alpha-1)v + 1]^{\eta/(1-\alpha)}$ with $\eta \in (0, 1]$, $\alpha > 1$, $\kappa > 0$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \leq 0$

Under additional assumptions on Ψ_α and $b_{\mu, \alpha}$, if $(\mu_n)_{n \geq 1}$ weakly converges to μ^* , then :

$$\mu^* \text{ is a fixed point of } \mathcal{I}_\alpha \text{ and } \Psi_\alpha(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_\alpha(\zeta k)$$

The special case of finite mixture models

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

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

Let $\Theta = (\theta_1, \dots, \theta_J) \in \mathbb{T}^J$, $\boldsymbol{\lambda}_1 = (\lambda_{1,1}, \dots, \lambda_{J,1}) \in \mathcal{S}_J$ and denote

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

Then, $\mu_n = \underbrace{\mathcal{I}_{\alpha} \circ \dots \circ \mathcal{I}_{\alpha}}_{n \text{ times}}(\mu_{\boldsymbol{\lambda}_1})$ is of the form $\mu_n = \sum_{j=1}^J \lambda_{j,n} \delta_{\theta_j}$ with

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J, \quad n \geq 1$$

$$\text{NB : } \mu_n k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_j, y)$$

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

The special case of finite mixture models

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

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

Let $\Theta = (\theta_1, \dots, \theta_J) \in \mathbb{T}^J$, $\boldsymbol{\lambda}_1 = (\lambda_{1,1}, \dots, \lambda_{J,1}) \in \mathcal{S}_J$ and denote

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

Then, $\mu_n = \underbrace{\mathcal{I}_{\alpha} \circ \dots \circ \mathcal{I}_{\alpha}}_{n \text{ times}}(\mu_{\boldsymbol{\lambda}_1})$ is of the form $\mu_n = \sum_{j=1}^J \lambda_{j,n} \delta_{\theta_j}$ with

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J, \quad n \geq 1$$

$$\text{NB : } \mu_n k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_j, y)$$

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

The special case of finite mixture models

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

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

Let $\Theta = (\theta_1, \dots, \theta_J) \in \mathbb{T}^J$, $\boldsymbol{\lambda}_1 = (\lambda_{1,1}, \dots, \lambda_{J,1}) \in \mathcal{S}_J$ and denote

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

Then, $\mu_n = \underbrace{\mathcal{I}_{\alpha} \circ \dots \circ \mathcal{I}_{\alpha}}_{n \text{ times}}(\mu_{\boldsymbol{\lambda}_1})$ is of the form $\mu_n = \sum_{j=1}^J \lambda_{j,n} \delta_{\theta_j}$ with

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J, \quad n \geq 1$$

$$\text{NB : } \mu_n k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_j, y)$$

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

The special case of finite mixture models

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \cdot \Gamma(b_{\mu_n, \alpha}(\theta) + \kappa)}{\mu_n(\Gamma(b_{\mu_n, \alpha} + \kappa))}, \quad n \geq 1$$

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

Let $\Theta = (\theta_1, \dots, \theta_J) \in \mathbb{T}^J$, $\boldsymbol{\lambda}_1 = (\lambda_{1,1}, \dots, \lambda_{J,1}) \in \mathcal{S}_J$ and denote

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

Then, $\mu_n = \underbrace{\mathcal{I}_{\alpha} \circ \dots \circ \mathcal{I}_{\alpha}}_{n \text{ times}}(\mu_{\boldsymbol{\lambda}_1})$ is of the form $\mu_n = \sum_{j=1}^J \lambda_{j,n} \delta_{\theta_j}$ with

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J, \quad n \geq 1$$

$$\text{NB : } \mu_n k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_j, y)$$

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

Convergence results for finite mixture models

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L -smooth and $-\log \Gamma$ is concave increasing
e.g. Entropic Mirror Descent: when $\alpha = 1$, we have for all $\eta \in (0, 1)$

$$\Psi_{\alpha}(\mu_{\lambda_n} k) - \Psi_{\alpha}(\mu^* k) \leq \frac{\log J}{\eta N} + \frac{\sqrt{2 \log J} |b|_{\infty, 1}}{(1 - \eta) N}$$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \geq 0$

Under additional assumptions on Ψ_{α} and $b_{\mu, \alpha}$, if $\{K(\theta_1, \cdot), \dots, K(\theta_J, \cdot)\}$ are linearly independent, then :

- $(\lambda_n)_{n \geq 1}$ converges to some λ^*
- $\mu^* = \mu_{\lambda^*}$ is a fixed point of \mathcal{I}_{α} and $\Psi_{\alpha}(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_{\alpha}(\zeta k)$

Convergence results for finite mixture models

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L -smooth and $-\log \Gamma$ is concave increasing
e.g. Entropic Mirror Descent: when $\alpha = 1$, we have for all $\eta \in (0, 1)$

$$\Psi_{\alpha}(\mu_{\lambda_n} k) - \Psi_{\alpha}(\mu^* k) \leq \frac{\log J}{\eta N} + \frac{\sqrt{2 \log J} |b|_{\infty, 1}}{(1 - \eta) N}$$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \geq 0$

Under additional assumptions on Ψ_{α} and $b_{\mu, \alpha}$, if $\{K(\theta_1, \cdot), \dots, K(\theta_J, \cdot)\}$ are linearly independent, then :

- $(\lambda_n)_{n \geq 1}$ converges to some λ^*
- $\mu^* = \mu_{\lambda^*}$ is a fixed point of \mathcal{I}_{α} and $\Psi_{\alpha}(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_{\alpha}(\zeta k)$

Convergence results for finite mixture models

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L -smooth and $-\log \Gamma$ is concave increasing
e.g. Entropic Mirror Descent: when $\alpha = 1$, we have for all $\eta \in (0, 1)$

$$\Psi_{\alpha}(\mu_{\lambda_n} k) - \Psi_{\alpha}(\mu^* k) \leq \frac{\log J}{\eta N} + \frac{\sqrt{2 \log J} |b|_{\infty, 1}}{(1 - \eta) N}$$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \geq 0$

Under additional assumptions on Ψ_{α} and $b_{\mu, \alpha}$, if $\{K(\theta_1, \cdot), \dots, K(\theta_J, \cdot)\}$ are linearly independent, then :

- $(\lambda_n)_{n \geq 1}$ converges to some λ^*
- $\mu^* = \mu_{\lambda^*}$ is a fixed point of \mathcal{I}_{α} and $\Psi_{\alpha}(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_{\alpha}(\zeta k)$

Convergence results for finite mixture models

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L -smooth and $-\log \Gamma$ is concave increasing
e.g. Entropic Mirror Descent: when $\alpha = 1$, we have for all $\eta \in (0, 1)$

$$\Psi_{\alpha}(\mu_{\lambda_n} k) - \Psi_{\alpha}(\mu^* k) \leq \frac{\log J}{\eta N} + \frac{\sqrt{2 \log J} |b|_{\infty, 1}}{(1 - \eta) N}$$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \geq 0$

Under additional assumptions on Ψ_{α} and $b_{\mu, \alpha}$, if $\{K(\theta_1, \cdot), \dots, K(\theta_J, \cdot)\}$ are linearly independent, then :

- $(\lambda_n)_{n \geq 1}$ converges to some λ^*
- $\mu^* = \mu_{\lambda^*}$ is a fixed point of \mathcal{I}_{α} and $\Psi_{\alpha}(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_{\alpha}(\zeta k)$

Convergence results for finite mixture models

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L -smooth and $-\log \Gamma$ is concave increasing
e.g. Entropic Mirror Descent: when $\alpha = 1$, we have for all $\eta \in (0, 1)$

$$\Psi_{\alpha}(\mu_{\lambda_n} k) - \Psi_{\alpha}(\mu^* k) \leq \frac{\log J}{\eta N} + \frac{\sqrt{2 \log J} |b|_{\infty, 1}}{(1 - \eta) N}$$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \geq 0$

Under additional assumptions on Ψ_{α} and $b_{\mu, \alpha}$, if $\{K(\theta_1, \cdot), \dots, K(\theta_J, \cdot)\}$ are linearly independent, then :

- $(\lambda_n)_{n \geq 1}$ converges to some λ^*
- $\mu^* = \mu_{\lambda^*}$ is a fixed point of \mathcal{I}_{α} and $\Psi_{\alpha}(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_{\alpha}(\zeta k)$

Convergence results for finite mixture models

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J, \quad n \geq 1$$

Assume (A1) and that $|b|_{\infty, \alpha} = \sup_{\theta \in \mathcal{T}, \mu \in \mathcal{M}_1(\mathcal{T})} |b_{\mu, \alpha}(\theta)| < \infty$

→ $O(1/N)$ convergence rates when Γ is L -smooth and $-\log \Gamma$ is concave increasing
e.g. Entropic Mirror Descent: when $\alpha = 1$, we have for all $\eta \in (0, 1)$

$$\Psi_{\alpha}(\mu_{\lambda_n} k) - \Psi_{\alpha}(\mu^* k) \leq \frac{\log J}{\eta N} + \frac{\sqrt{2 \log J} |b|_{\infty, 1}}{(1 - \eta) N}$$

→ The case $\alpha < 1$ for the Power Descent is trickier... $\eta \in (0, 1]$, $\kappa \geq 0$

Under additional assumptions on Ψ_{α} and $b_{\mu, \alpha}$, if $\{K(\theta_1, \cdot), \dots, K(\theta_J, \cdot)\}$ are linearly independent, then :

- $(\lambda_n)_{n \geq 1}$ converges to some λ^*
- $\mu^* = \mu_{\lambda^*}$ is a fixed point of \mathcal{I}_{α} and $\Psi_{\alpha}(\mu^* k) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_{\alpha}(\zeta k)$

Towards a practical implementation

Algorithm

Let $\Theta = (\theta_1, \dots, \theta_J) \in \mathcal{T}^J$ be **fixed** and let $\lambda_1 \in \mathcal{S}_J$. At time $n \geq 1$, define

$$\mu_{n+1}^k = \sum_{j=1}^J \lambda_{j,n+1} k(\theta_j, \cdot)$$

where

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J$$

→ Monte Carlo approximations to estimate $b_{\mu_n, \alpha}(\theta_j)$, e.g.

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

with $Y_{1,n}, \dots, Y_{M,n} \stackrel{\text{i.i.d.}}{\sim} \mu_n^k$.

→ **Exploitation step** not requiring any information on the distribution of $\theta_1, \dots, \theta_J$

→ Idea : combine this step with and *Exploration Step* updating Θ

Towards a practical implementation

Algorithm

Let $\Theta = (\theta_1, \dots, \theta_J) \in \mathcal{T}^J$ be **fixed** and let $\lambda_1 \in \mathcal{S}_J$. At time $n \geq 1$, define

$$\mu_{n+1}^k = \sum_{j=1}^J \lambda_{j,n+1} k(\theta_j, \cdot)$$

where

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J$$

→ Monte Carlo approximations to estimate $b_{\mu_n, \alpha}(\theta_j)$, e.g.

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

with $Y_{1,n}, \dots, Y_{M,n} \stackrel{\text{i.i.d.}}{\sim} \mu_n^k$.

→ **Exploitation step** not requiring any information on the distribution of $\theta_1, \dots, \theta_J$

→ Idea : combine this step with and *Exploration Step* updating Θ

Towards a practical implementation

Algorithm

Let $\Theta = (\theta_1, \dots, \theta_J) \in \mathcal{T}^J$ be **fixed** and let $\lambda_1 \in \mathcal{S}_J$. At time $n \geq 1$, define

$$\mu_{n+1}^k = \sum_{j=1}^J \lambda_{j,n+1} k(\theta_j, \cdot)$$

where

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J$$

→ Monte Carlo approximations to estimate $b_{\mu_n, \alpha}(\theta_j)$, e.g.

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

with $Y_{1,n}, \dots, Y_{M,n} \stackrel{\text{i.i.d.}}{\sim} \mu_n^k$.

→ **Exploitation step** not requiring any information on the distribution of $\theta_1, \dots, \theta_J$

→ Idea : combine this step with and *Exploration Step* updating Θ

Towards a practical implementation

Algorithm

Let $\Theta = (\theta_1, \dots, \theta_J) \in \mathcal{T}^J$ be **fixed** and let $\lambda_1 \in \mathcal{S}_J$. At time $n \geq 1$, define

$$\mu_{n+1}^k = \sum_{j=1}^J \lambda_{j,n+1} k(\theta_j, \cdot)$$

where

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \Gamma(b_{\mu_n, \alpha}(\theta_j) + \kappa)}{\sum_{i=1}^J \lambda_{i,n} \Gamma(b_{\mu_n, \alpha}(\theta_i) + \kappa)}, \quad j = 1 \dots J$$

→ Monte Carlo approximations to estimate $b_{\mu_n, \alpha}(\theta_j)$, e.g.

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

with $Y_{1,n}, \dots, Y_{M,n} \stackrel{\text{i.i.d.}}{\sim} \mu_n^k$.

→ **Exploitation step** not requiring any information on the distribution of $\theta_1, \dots, \theta_J$

→ Idea : combine this step with and *Exploration Step* updating Θ

Outline

- 1 Infinite-dimensional Alpha-divergence minimisation
- 2 Numerical experiments
- 3 Conclusion of Part 2

Numerical experiments

- Gaussian kernel with density k_h and bandwidth h , $\mathsf{T} = \mathbb{R}^d$

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

Algorithm

- ① **Exploitation step** : optimise λ using the (α, Γ) -descent.
 - ② *Exploration step* : update Θ (e.g. by sampling under $\mu_{\lambda, \Theta} k_h$, $h \propto J^{-1/(4+d)}$)
- Toy example
 $p(y) = Z \times [0.5\mathcal{N}(y; -2u_d, I_d) + 0.5\mathcal{N}(y; 2u_d, I_d)]$, $Z = 2$
 - Bayesian Logistic Regression Covertypes dataset (581,012 data points and 54 features)

Numerical experiments

- Gaussian kernel with density k_h and bandwidth h , $\mathsf{T} = \mathbb{R}^d$

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

Algorithm

- ① **Exploitation step** : optimise λ using the (α, Γ) -descent.
 - ② *Exploration step* : update Θ (e.g. by sampling under $\mu_{\lambda, \Theta} k_h$, $h \propto J^{-1/(4+d)}$)
- Toy example
 $p(y) = Z \times [0.5\mathcal{N}(y; -2u_d, I_d) + 0.5\mathcal{N}(y; 2u_d, I_d)], Z = 2$
 - Bayesian Logistic Regression Covertypes dataset (581,012 data points and 54 features)

Numerical experiments

- Gaussian kernel with density k_h and bandwidth h , $\mathsf{T} = \mathbb{R}^d$

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

Algorithm

- ① **Exploitation step** : optimise λ using the (α, Γ) -descent.
- ② *Exploration step* : update Θ (e.g. by sampling under $\mu_{\lambda, \Theta} k_h$, $h \propto J^{-1/(4+d)}$)
- Toy example
 $p(y) = Z \times [0.5\mathcal{N}(y; -2u_d, I_d) + 0.5\mathcal{N}(y; 2u_d, I_d)], Z = 2$
- Bayesian Logistic Regression Covertypes dataset (581,012 data points and 54 features)

Numerical experiments

- Gaussian kernel with density k_h and bandwidth h , $\mathsf{T} = \mathbb{R}^d$

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

Algorithm

- ❶ **Exploitation step** : optimise λ using the (α, Γ) -descent.
 - ❷ *Exploration step* : update Θ (e.g. by sampling under $\mu_{\lambda, \Theta} k_h$, $h \propto J^{-1/(4+d)}$)
- Toy example
 $p(y) = Z \times [0.5\mathcal{N}(y; -2u_d, I_d) + 0.5\mathcal{N}(y; 2u_d, I_d)]$, $Z = 2$
 - Bayesian Logistic Regression Covertypes dataset (581,012 data points and 54 features)

Numerical experiments

- Gaussian kernel with density k_h and bandwidth h , $\mathsf{T} = \mathbb{R}^d$

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

Algorithm

- ❶ **Exploitation step** : optimise $\boldsymbol{\lambda}$ using the (α, Γ) -descent.
 - ❷ *Exploration step* : update Θ (e.g. by sampling under $\mu_{\boldsymbol{\lambda}, \Theta} k_h$, $h \propto J^{-1/(4+d)}$)
- Toy example
 $p(y) = Z \times [0.5\mathcal{N}(\mathbf{y}; -2\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(\mathbf{y}; 2\mathbf{u}_d, \mathbf{I}_d)], Z = 2$
 - Bayesian Logistic Regression Covertypes dataset (581,012 data points and 54 features)

Numerical experiments

- Gaussian kernel with density k_h and bandwidth h , $\mathsf{T} = \mathbb{R}^d$

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

Algorithm

- ❶ **Exploitation step** : optimise $\boldsymbol{\lambda}$ using the (α, Γ) -descent.
- ❷ *Exploration step* : update Θ (e.g. by sampling under $\mu_{\boldsymbol{\lambda}, \Theta} k_h$, $h \propto J^{-1/(4+d)}$)
- Toy example
 $p(y) = Z \times [0.5\mathcal{N}(\mathbf{y}; -2\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(\mathbf{y}; 2\mathbf{u}_d, \mathbf{I}_d)], Z = 2$
- Bayesian Logistic Regression Covertypes dataset (581,012 data points and 54 features)

Toy example : Entropic Mirror Descent vs Power Descent

Comparison between

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

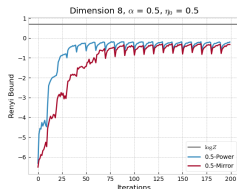
$J = M = 100$, initial mixture weights : $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$
 $\eta_n = \eta_0/\sqrt{n}$, $\eta_0 = 0.5$, cv criterion : VR-Bound averaged over 100 trials

Toy example : Entropic Mirror Descent vs Power Descent

Comparison between

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

$J = M = 100$, initial mixture weights : $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$
 $\eta_n = \eta_0/\sqrt{n}$, $\eta_0 = 0.5$, cv criterion : VR-Bound averaged over 100 trials

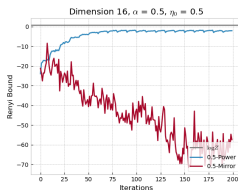
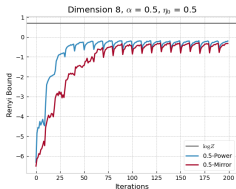


Toy example : Entropic Mirror Descent vs Power Descent

Comparison between

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

$J = M = 100$, initial mixture weights : $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$
 $\eta_n = \eta_0/\sqrt{n}$, $\eta_0 = 0.5$, cv criterion : VR-Bound averaged over 100 trials

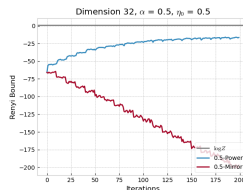
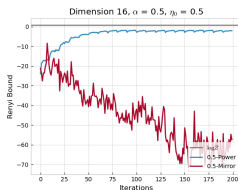
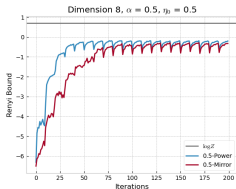


Toy example : Entropic Mirror Descent vs Power Descent

Comparison between

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

$J = M = 100$, initial mixture weights : $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$
 $\eta_n = \eta_0/\sqrt{n}$, $\eta_0 = 0.5$, cv criterion : VR-Bound averaged over 100 trials



Toy example : the case $\alpha = 1$

Comparison between:

- 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 $\alpha = 0.5$.

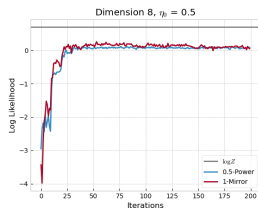
$J = M = 100$, initial mixture weights : $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$
 $\eta_n = \eta_0/\sqrt{n}$, $\eta_0 = 0.5$, cv criterion : llh averaged over 100 trials

Toy example : the case $\alpha = 1$

Comparison between:

- 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 $\alpha = 0.5$.

$J = M = 100$, initial mixture weights : $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$
 $\eta_n = \eta_0/\sqrt{n}$, $\eta_0 = 0.5$, cv criterion : llh averaged over 100 trials

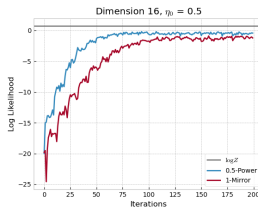
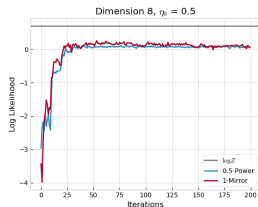


Toy example : the case $\alpha = 1$

Comparison between:

- 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 $\alpha = 0.5$.

$J = M = 100$, initial mixture weights : $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$
 $\eta_n = \eta_0/\sqrt{n}$, $\eta_0 = 0.5$, cv criterion : llh averaged over 100 trials

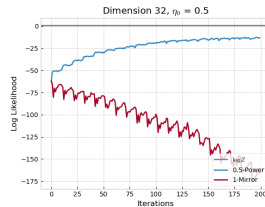
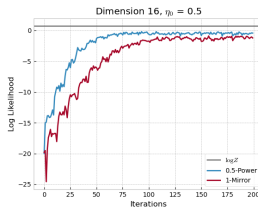
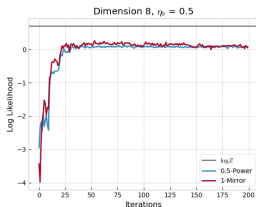


Toy example : the case $\alpha = 1$

Comparison between:

- 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 $\alpha = 0.5$.

$J = M = 100$, initial mixture weights : $[1/J, \dots, 1/J]$, $N = 10$, $T = 20$
 $\eta_n = \eta_0/\sqrt{n}$, $\eta_0 = 0.5$, cv criterion : llh averaged over 100 trials



Bayesian Logistic Regression

→ $\mathcal{D} = \{\mathbf{c}, \mathbf{x}\} : I$ binary class labels, $c_i \in \{-1, 1\}$, L covariates for each datapoint, $\mathbf{x}_i \in \mathbb{R}^L$

→ Model : L regression coefficients $w_l \in \mathbb{R}$, precision parameter $\beta \in \mathbb{R}^+$

$$p_0(\beta) = \text{Gamma}(\beta; a, b) ,$$

$$p_0(w_l | \beta) = \mathcal{N}(w_l; 0, \beta^{-1}) , \quad 1 \leq l \leq L$$

$$p(c_i = 1 | \mathbf{x}_i, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}} , \quad 1 \leq i \leq I$$

where $a = 1$ and $b = 0.01$

Nonparametric variational inference S. Gershman, M. Hoffman, and D. Blei (2012). ICML

→ Quantity of interest : $p(y | \mathcal{D})$ with $y = [\mathbf{w}, \log \beta]$

Comparison between

- 0.5-Power descent
- Typical AIS

$N = 1$, $T = 500$, $J_0 = M_0 = 20$, $J_{t+1} = M_{t+1} = J_t + 1$

initial mixture weights : $[1/J_t, \dots, 1/J_t]$, $\eta_n = \eta_0 / \sqrt{n}$ with $\eta_0 = 0.05$

Bayesian Logistic Regression

→ $\mathcal{D} = \{\mathbf{c}, \mathbf{x}\}$: I binary class labels, $c_i \in \{-1, 1\}$, L covariates for each datapoint, $\mathbf{x}_i \in \mathbb{R}^L$

→ Model : L regression coefficients $w_l \in \mathbb{R}$, precision parameter $\beta \in \mathbb{R}^+$

$$p_0(\beta) = \text{Gamma}(\beta; a, b) ,$$

$$p_0(w_l | \beta) = \mathcal{N}(w_l; 0, \beta^{-1}) , \quad 1 \leq l \leq L$$

$$p(c_i = 1 | \mathbf{x}_i, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}} , \quad 1 \leq i \leq I$$

where $a = 1$ and $b = 0.01$

Nonparametric variational inference S. Gershman, M. Hoffman, and D. Blei (2012). ICML

→ Quantity of interest : $p(y | \mathcal{D})$ with $y = [\mathbf{w}, \log \beta]$

Comparison between

- 0.5-Power descent
- Typical AIS

$N = 1$, $T = 500$, $J_0 = M_0 = 20$, $J_{t+1} = M_{t+1} = J_t + 1$

initial mixture weights : $[1/J_t, \dots, 1/J_t]$, $\eta_n = \eta_0 / \sqrt{n}$ with $\eta_0 = 0.05$

Bayesian Logistic Regression

→ $\mathcal{D} = \{\mathbf{c}, \mathbf{x}\}$: I binary class labels, $c_i \in \{-1, 1\}$, L covariates for each datapoint, $\mathbf{x}_i \in \mathbb{R}^L$

→ Model : L regression coefficients $w_l \in \mathbb{R}$, precision parameter $\beta \in \mathbb{R}^+$

$$p_0(\beta) = \text{Gamma}(\beta; a, b) ,$$

$$p_0(w_l | \beta) = \mathcal{N}(w_l; 0, \beta^{-1}) , \quad 1 \leq l \leq L$$

$$p(c_i = 1 | \mathbf{x}_i, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}} , \quad 1 \leq i \leq I$$

where $a = 1$ and $b = 0.01$

Nonparametric variational inference S. Gershman, M. Hoffman, and D. Blei (2012). ICML

→ Quantity of interest : $p(y | \mathcal{D})$ with $y = [\mathbf{w}, \log \beta]$

Comparison between

- 0.5-Power descent
- Typical AIS

$N = 1$, $T = 500$, $J_0 = M_0 = 20$, $J_{t+1} = M_{t+1} = J_t + 1$

initial mixture weights : $[1/J_t, \dots, 1/J_t]$, $\eta_n = \eta_0 / \sqrt{n}$ with $\eta_0 = 0.05$

Bayesian Logistic Regression

→ $\mathcal{D} = \{\mathbf{c}, \mathbf{x}\} : I$ binary class labels, $c_i \in \{-1, 1\}$, L covariates for each datapoint, $\mathbf{x}_i \in \mathbb{R}^L$

→ Model : L regression coefficients $w_l \in \mathbb{R}$, precision parameter $\beta \in \mathbb{R}^+$

$$p_0(\beta) = \text{Gamma}(\beta; a, b),$$

$$p_0(w_l | \beta) = \mathcal{N}(w_l; 0, \beta^{-1}), \quad 1 \leq l \leq L$$

$$p(c_i = 1 | \mathbf{x}_i, \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}_i}}, \quad 1 \leq i \leq I$$

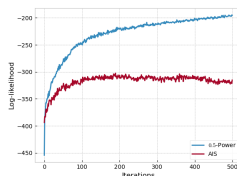
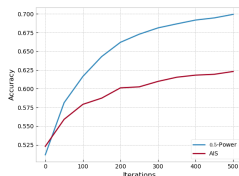
where $a = 1$ and $b = 0.01$

Nonparametric variational inference S. Gershman, M. Hoffman, and D. Blei (2012). ICML

→ Quantity of interest : $p(y | \mathcal{D})$ with $y = [\mathbf{w}, \log \beta]$

Comparison between

- 0.5-Power descent
- Typical AIS



$N = 1, T = 500, J_0 = M_0 = 20, J_{t+1} = M_{t+1} = J_t + 1$

initial mixture weights : $[1/J_t, \dots, 1/J_t], \eta_n = \eta_0 / \sqrt{n}$ with $\eta_0 = 0.05$

Outline

- 1 Infinite-dimensional Alpha-divergence minimisation
- 2 Numerical experiments
- 3 Conclusion of Part 2**

Conclusion of Part 2

General framework for infinite-dimensional α -divergence minimisation over

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

- recovers the Entropic Mirror Descent algorithm
- novel Power Descent algorithm
- conditions for a systematic decrease + convergence results
- applicable to mixture models :

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

→ Exploitation - Exploration algorithm

- ① Update for Θ not specified (e.g. your favorite update for Θ)
- ② Empirical advantages of using the Power Descent algorithm

Conclusion of Part 2

General framework for infinite-dimensional α -divergence minimisation over

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

- recovers the **Entropic Mirror Descent** algorithm
- novel **Power Descent** algorithm
- conditions for a **systematic decrease + convergence** results
- applicable to **mixture models** :

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

→ Exploitation - Exploration algorithm

- ① Update for Θ **not specified** (e.g. **your** favorite update for Θ)
- ② Empirical advantages of using the **Power Descent** algorithm

Conclusion of Part 2

General framework for infinite-dimensional α -divergence minimisation over

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

- recovers the **Entropic Mirror Descent** algorithm
- novel **Power Descent** algorithm
- conditions for a systematic decrease + convergence results
- applicable to mixture models :

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

→ Exploitation - Exploration algorithm

- ① Update for Θ not specified (e.g. your favorite update for Θ)
- ② Empirical advantages of using the **Power Descent** algorithm

Conclusion of Part 2

General framework for infinite-dimensional α -divergence minimisation over

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

- recovers the **Entropic Mirror Descent** algorithm
- novel **Power Descent** algorithm
- conditions for a **systematic decrease** + **convergence** results
- applicable to **mixture models** :

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

→ Exploitation - Exploration algorithm

- ① Update for Θ **not specified** (e.g. **your** favorite update for Θ)
- ② Empirical advantages of using the **Power Descent** algorithm

Conclusion of Part 2

General framework for infinite-dimensional α -divergence minimisation over

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

- recovers the **Entropic Mirror Descent** algorithm
- novel **Power Descent** algorithm
- conditions for a **systematic decrease** + **convergence** results
- applicable to **mixture models** :

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

→ Exploitation - Exploration algorithm

- ① Update for Θ **not specified** (e.g. **your** favorite update for Θ)
- ② Empirical advantages of using the **Power Descent** algorithm

Conclusion of Part 2

General framework for infinite-dimensional α -divergence minimisation over

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

- recovers the **Entropic Mirror Descent** algorithm
- novel **Power Descent** algorithm
- conditions for a **systematic decrease** + **convergence** results
- applicable to **mixture models** :

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

→ Exploitation - Exploration algorithm

- ① Update for Θ **not specified** (e.g. **your** favorite update for Θ)
- ② Empirical advantages of using the **Power Descent** algorithm

Conclusion of Part 2

General framework for infinite-dimensional α -divergence minimisation over

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

- recovers the **Entropic Mirror Descent** algorithm
- novel **Power Descent** algorithm
- conditions for a **systematic decrease** + **convergence** results
- applicable to **mixture models** :

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

→ Exploitation - Exploration algorithm

- ① Update for Θ **not specified** (e.g. **your** favorite update for Θ)
- ② Empirical advantages of using the **Power Descent** algorithm

Conclusion of Part 2

General framework for infinite-dimensional α -divergence minimisation over

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

- recovers the **Entropic Mirror Descent** algorithm
- novel **Power Descent** algorithm
- conditions for a **systematic decrease** + **convergence** results
- applicable to **mixture models** :

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

→ Exploitation - Exploration algorithm

- ① Update for Θ **not specified** (e.g. **your** favorite update for Θ)
- ② Empirical advantages of using the **Power Descent** algorithm

Food for thoughts - 1

- Question : What about Entropic Mirror Descent applied to $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$?

Mixture weights optimisation for Alpha-Divergence Variational Inference.

K. Daudel and R. Douc (2021). To appear in NeurIPS2021

$$\text{Alpha : } \mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\frac{\eta}{\alpha - 1} \int_Y k(\theta, y) \mu_n k(y)^{\alpha-1} p(y)^{1-\alpha} \nu(\mathrm{d}y) \right]$$

$$\text{Rényi's Alpha : } \mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\frac{\eta}{\alpha - 1} \frac{\int_Y k(\theta, y) \mu_n k(y)^{\alpha-1} p(y)^{1-\alpha} \nu(\mathrm{d}y)}{\int_Y \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(\mathrm{d}y)} \right]$$

→ The Entropic Mirror Descent applied to $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$ is in fact closely-related to the **Power Descent**

Food for thoughts - 1

- Question : What about Entropic Mirror Descent applied to $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$?

Mixture weights optimisation for Alpha-Divergence Variational Inference.

K. Daudel and R. Douc (2021). To appear in NeurIPS2021

$$\text{Alpha : } \mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\frac{\eta}{\alpha - 1} \int_Y k(\theta, y) \mu_n k(y)^{\alpha-1} p(y)^{1-\alpha} \nu(\mathrm{d}y) \right]$$

$$\text{Rényi's Alpha : } \mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\frac{\eta}{\alpha - 1} \frac{\int_Y k(\theta, y) \mu_n k(y)^{\alpha-1} p(y)^{1-\alpha} \nu(\mathrm{d}y)}{\int_Y \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(\mathrm{d}y)} \right]$$

→ The Entropic Mirror Descent applied to $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$ is in fact closely-related to the **Power Descent**

Food for thoughts - 1

- Question : What about Entropic Mirror Descent applied to $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$?

Mixture weights optimisation for Alpha-Divergence Variational Inference.

K. Daudel and R. Douc (2021). To appear in NeurIPS2021

$$\text{Alpha : } \mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\frac{\eta}{\alpha - 1} \int_{\mathcal{Y}} k(\theta, y) \mu_n k(y)^{\alpha-1} p(y)^{1-\alpha} \nu(\mathrm{d}y) \right]$$

$$\text{Rényi's Alpha : } \mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\frac{\eta}{\alpha - 1} \frac{\int_{\mathcal{Y}} k(\theta, y) \mu_n k(y)^{\alpha-1} p(y)^{1-\alpha} \nu(\mathrm{d}y)}{\int_{\mathcal{Y}} \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(\mathrm{d}y)} \right]$$

→ The Entropic Mirror Descent applied to $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$ is in fact closely-related to the **Power Descent**

Food for thoughts - 1

- Question : What about Entropic Mirror Descent applied to $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$?

Mixture weights optimisation for Alpha-Divergence Variational Inference.

K. Daudel and R. Douc (2021). To appear in NeurIPS2021

$$\text{Alpha : } \mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\frac{\eta}{\alpha - 1} \int_Y k(\theta, y) \mu_n k(y)^{\alpha-1} p(y)^{1-\alpha} \nu(\mathrm{d}y) \right]$$

$$\text{Rényi's Alpha : } \mu_{n+1}(\mathrm{d}\theta) \propto \mu_n(\mathrm{d}\theta) \exp \left[-\frac{\eta}{\alpha - 1} \frac{\int_Y k(\theta, y) \mu_n k(y)^{\alpha-1} p(y)^{1-\alpha} \nu(\mathrm{d}y)}{\int_Y \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(\mathrm{d}y)} \right]$$

→ The Entropic Mirror Descent applied to $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$ is in fact closely-related to the Power Descent

Food for thoughts - 1

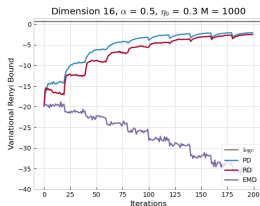
- Question : What about Entropic Mirror Descent applied to $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$?

Mixture weights optimisation for Alpha-Divergence Variational Inference.

K. Daudel and R. Douc (2021). To appear in NeurIPS2021

$$\text{Alpha : } \mu_{n+1}(d\theta) \propto \mu_n(d\theta) \exp \left[-\frac{\eta}{\alpha - 1} \int_Y k(\theta, y) \mu_n k(y)^{\alpha-1} p(y)^{1-\alpha} \nu(dy) \right]$$

$$\text{Rényi's Alpha : } \mu_{n+1}(d\theta) \propto \mu_n(d\theta) \exp \left[-\frac{\eta}{\alpha - 1} \frac{\int_Y k(\theta, y) \mu_n k(y)^{\alpha-1} p(y)^{1-\alpha} \nu(dy)}{\int_Y \mu_n k(y)^\alpha p(y)^{1-\alpha} \nu(dy)} \right]$$



→ The Entropic Mirror Descent applied to $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$ is in fact closely-related to the **Power Descent**

Food for thoughts - 2

- Question : What is a good choice for the exploration step?

Some answers in Part 3!

Food for thoughts - 2

- Question : What is a good choice for the exploration step?

Some answers in Part 3!

Variational Inference

Foundations and recent advances

(Part 3)

Kamélia Daudel



University of Bristol – 09/03/2022

Reminder - 1

Alpha-Divergence Variational Inference : Two possible objective functions

$$\Psi_{\alpha}(q; p) = \int_Y f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$$
$$-\alpha^{-1} \mathcal{L}_{\alpha}(q; p) = \frac{1}{\alpha(\alpha - 1)} \log \left(\int_Y q(y)^{\alpha} p(y)^{1-\alpha} \nu(dy) \right)$$

with $p = p(\cdot, \mathcal{D})$ and

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

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

Stochastic Gradient Descent w.r.t θ on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

- $\mathcal{Q} = \{q : y \mapsto \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) : \mu \in \mathbb{M}\}$

Power Descent, Entropic Mirror Descent on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

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

Reminder - 1

Alpha-Divergence Variational Inference : Two possible objective functions

$$\Psi_{\alpha}(q; p) = \int_Y f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$$
$$-\alpha^{-1} \mathcal{L}_{\alpha}(q; p) = \frac{1}{\alpha(\alpha - 1)} \log \left(\int_Y q(y)^{\alpha} p(y)^{1-\alpha} \nu(dy) \right)$$

with $p = p(\cdot, \mathcal{D})$ and

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

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

Stochastic Gradient Descent w.r.t θ on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

- $\mathcal{Q} = \{q : y \mapsto \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) : \mu \in \mathbb{M}\}$

Power Descent, Entropic Mirror Descent on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

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

Reminder - 1

Alpha-Divergence Variational Inference : Two possible objective functions

$$\Psi_{\alpha}(q; p) = \int_Y f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$$
$$-\alpha^{-1} \mathcal{L}_{\alpha}(q; p) = \frac{1}{\alpha(\alpha - 1)} \log \left(\int_Y q(y)^{\alpha} p(y)^{1-\alpha} \nu(dy) \right)$$

with $p = p(\cdot, \mathcal{D})$ and

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

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

Stochastic Gradient Descent w.r.t θ on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

- $\mathcal{Q} = \{q : y \mapsto \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) : \mu \in \mathbb{M}\}$

Power Descent, Entropic Mirror Descent on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

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

Reminder - 1

Alpha-Divergence Variational Inference : Two possible objective functions

$$\Psi_{\alpha}(q; p) = \int_Y f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$$
$$-\alpha^{-1} \mathcal{L}_{\alpha}(q; p) = \frac{1}{\alpha(\alpha - 1)} \log \left(\int_Y q(y)^{\alpha} p(y)^{1-\alpha} \nu(dy) \right)$$

with $p = p(\cdot, \mathcal{D})$ and

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

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

Stochastic Gradient Descent w.r.t θ on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

- $\mathcal{Q} = \{q : y \mapsto \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) : \mu \in \mathbb{M}\}$

Power Descent, Entropic Mirror Descent on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

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

Reminder - 1

Alpha-Divergence Variational Inference : Two possible objective functions

$$\Psi_{\alpha}(q; p) = \int_Y f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$$
$$-\alpha^{-1} \mathcal{L}_{\alpha}(q; p) = \frac{1}{\alpha(\alpha - 1)} \log \left(\int_Y q(y)^{\alpha} p(y)^{1-\alpha} \nu(dy) \right)$$

with $p = p(\cdot, \mathcal{D})$ and

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

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

Stochastic Gradient Descent w.r.t θ on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

- $\mathcal{Q} = \{q : y \mapsto \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) : \mu \in \mathbb{M}\}$

Power Descent, Entropic Mirror Descent on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

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

Reminder - 1

Alpha-Divergence Variational Inference : Two possible objective functions

$$\Psi_{\alpha}(q; p) = \int_Y f_{\alpha} \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$$
$$-\alpha^{-1} \mathcal{L}_{\alpha}(q; p) = \frac{1}{\alpha(\alpha - 1)} \log \left(\int_Y q(y)^{\alpha} p(y)^{1-\alpha} \nu(dy) \right)$$

with $p = p(\cdot, \mathcal{D})$ and

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

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

Stochastic Gradient Descent w.r.t θ on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

- $\mathcal{Q} = \{q : y \mapsto \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) : \mu \in \mathbb{M}\}$

Power Descent, Entropic Mirror Descent on $\Psi_{\alpha}(q; p)$ (resp. $-\alpha^{-1} \mathcal{L}_{\alpha}(q; p)$)

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

Reminder - 2

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

Stochastic Gradient Descent w.r.t θ on $\Psi_\alpha(q; p)$ (resp. $-\alpha^{-1}\mathcal{L}_\alpha(q; p)$)

- $\mathcal{Q} = \{q : y \mapsto \int_{\mathcal{T}} \mu(d\theta)k(\theta, y) : \mu \in \mathcal{M}\}$

Power Descent, Entropic Mirror Descent on $\Psi_\alpha(q; p)$ (resp. $-\alpha^{-1}\mathcal{L}_\alpha(q; p)$)

→ applies to $\mathcal{Q} = \left\{q : y \mapsto \sum_{j=1}^J \lambda_j k(\theta_j, y) : \boldsymbol{\lambda} \in \mathcal{S}_J\right\}$

Question : Can we propose valid updates for

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

Reminder - 2

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

Stochastic Gradient Descent w.r.t θ on $\Psi_\alpha(q; p)$ (resp. $-\alpha^{-1}\mathcal{L}_\alpha(q; p)$)

- $\mathcal{Q} = \{q : y \mapsto \int_{\mathcal{T}} \mu(d\theta)k(\theta, y) : \mu \in \mathcal{M}\}$

Power Descent, Entropic Mirror Descent on $\Psi_\alpha(q; p)$ (resp. $-\alpha^{-1}\mathcal{L}_\alpha(q; p)$)

→ applies to $\mathcal{Q} = \left\{q : y \mapsto \sum_{j=1}^J \lambda_j k(\theta_j, y) : \boldsymbol{\lambda} \in \mathcal{S}_J\right\}$

Question : Can we propose valid updates for

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

Outline

- ① Monotonic Alpha-Divergence Minimisation
- ② Maximisation approach
- ③ Gradient-based approach
- ④ Numerical Experiments
- ⑤ Conclusion of Part 3

Outline

- 1 Monotonic Alpha-Divergence Minimisation
- 2 Maximisation approach
- 3 Gradient-based approach
- 4 Numerical Experiments
- 5 Conclusion of Part 3

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 alpha-divergence (i.e. Ψ_α) at each step.

→ Optimising w.r.t λ and Θ is the novelty compared to Part 2!

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). Let $\alpha \in [0, 1)$, $J \in \mathbb{N}^*$. Then, choosing $(\lambda_n, \Theta_n)_{n \geq 1}$ so that: $\Psi_\alpha(\mu_{\lambda_1, \Theta_1} k) < \infty$ and $\forall n \geq 1$,

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

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

where $\varphi_{j,n}^{(\alpha)}(y) = k(\theta_{j,n}, y) \left(\frac{\mu_{\lambda_n, \Theta_n} k(y)}{p(y)} \right)^{\alpha-1}$, 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). Let $\alpha \in [0, 1)$, $J \in \mathbb{N}^*$. Then, choosing $(\lambda_n, \Theta_n)_{n \geq 1}$ so that: $\Psi_\alpha(\mu_{\lambda_1, \Theta_1} k) < \infty$ and $\forall n \geq 1$,

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

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

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

Conditions for a monotonic decrease

Optimisation problem

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

(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). Let $\alpha \in [0, 1)$, $J \in \mathbb{N}^*$. Then, choosing $(\lambda_n, \Theta_n)_{n \geq 1}$ so that: $\Psi_\alpha(\mu_{\lambda_1, \Theta_1} k) < \infty$ and $\forall n \geq 1$,

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

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

where $\varphi_{j,n}^{(\alpha)}(y) = k(\theta_{j,n}, y) \left(\frac{\mu_{\lambda_n, \Theta_n} k(y)}{p(y)} \right)^{\alpha-1}$, 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). Let $\alpha \in [0, 1)$, $J \in \mathbb{N}^*$. Then, choosing $(\lambda_n, \Theta_n)_{n \geq 1}$ so that: $\Psi_\alpha(\mu_{\lambda_1, \Theta_1} k) < \infty$ and $\forall n \geq 1$,

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

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

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

Proof : 1) Proving a general lower bound - 1

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha = 0$: $f_0(u) = -\log(u) + u - 1$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \frac{q(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \left(\frac{q(y)}{p(y)} - \frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{q'(y)} \right) - \log \left(\frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \end{aligned}$$

so that

$$\Psi_\alpha(q) = \int_Y -\log \left(\frac{q(y)}{q'(y)} \right) p(y) \nu(dy) + \Psi_\alpha(q')$$

Proof : 1) Proving a general lower bound - 1

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha = 0$: $f_0(u) = -\log(u) + u - 1$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \frac{q(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \left(\frac{q(y)}{p(y)} - \frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{q'(y)} \right) - \log \left(\frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \end{aligned}$$

so that

$$\Psi_\alpha(q) = \int_Y -\log \left(\frac{q(y)}{q'(y)} \right) p(y) \nu(dy) + \Psi_\alpha(q')$$

Proof : 1) Proving a general lower bound - 1

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha = 0$: $f_0(u) = -\log(u) + u - 1$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \frac{q(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \left(\frac{q(y)}{p(y)} - \frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{q'(y)} \right) - \log \left(\frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \end{aligned}$$

so that

$$\Psi_\alpha(q) = \int_Y -\log \left(\frac{q(y)}{q'(y)} \right) p(y) \nu(dy) + \Psi_\alpha(q')$$

Proof : 1) Proving a general lower bound - 1

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha = 0$: $f_0(u) = -\log(u) + u - 1$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \frac{q(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \left(\frac{q(y)}{p(y)} - \frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{q'(y)} \right) - \log \left(\frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \end{aligned}$$

so that

$$\Psi_\alpha(q) = \int_Y -\log \left(\frac{q(y)}{q'(y)} \right) p(y) \nu(dy) + \Psi_\alpha(q')$$

Proof : 1) Proving a general lower bound - 1

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha = 0$: $f_0(u) = -\log(u) + u - 1$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \frac{q(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \left(\frac{q(y)}{p(y)} - \frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{q'(y)} \right) - \log \left(\frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \end{aligned}$$

so that

$$\Psi_\alpha(q) = \int_Y -\log \left(\frac{q(y)}{q'(y)} \right) p(y) \nu(dy) + \Psi_\alpha(q')$$

Proof : 1) Proving a general lower bound - 1

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha = 0$: $f_0(u) = -\log(u) + u - 1$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \frac{q(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \left(\frac{q(y)}{p(y)} - \frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{q'(y)} \right) - \log \left(\frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \end{aligned}$$

so that

$$\Psi_\alpha(q) = \int_Y -\log \left(\frac{q(y)}{q'(y)} \right) p(y) \nu(dy) + \Psi_\alpha(q')$$

Proof : 1) Proving a general lower bound - 1

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha = 0$: $f_0(u) = -\log(u) + u - 1$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \frac{q(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{p(y)} \right) + \left(\frac{q(y)}{p(y)} - \frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \\ &= \int_Y \left(-\log \left(\frac{q(y)}{q'(y)} \right) - \log \left(\frac{q'(y)}{p(y)} \right) + \frac{q'(y)}{p(y)} - 1 \right) p(y) \nu(dy) \end{aligned}$$

so that

$$\Psi_\alpha(q) = \int_Y -\log \left(\frac{q(y)}{q'(y)} \right) p(y) \nu(dy) + \Psi_\alpha(q')$$

Proof : 1) Proving a general lower bound - 2

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha \in (0, 1)$: $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q'(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha + \left(\frac{q'(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q'(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha \right] p(y) \nu(dy) + \Psi_\alpha(q') \end{aligned}$$

Proof : 1) Proving a general lower bound - 2

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha \in (0, 1)$: $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q'(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha + \left(\frac{q'(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q'(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha \right] p(y) \nu(dy) + \Psi_\alpha(q') \end{aligned}$$

Proof : 1) Proving a general lower bound - 2

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha \in (0, 1)$: $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q'(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha + \left(\frac{q'(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q'(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha \right] p(y) \nu(dy) + \Psi_\alpha(q') \end{aligned}$$

Proof : 1) Proving a general lower bound - 2

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha \in (0, 1)$: $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q'(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha + \left(\frac{q'(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q'(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha \right] p(y) \nu(dy) + \Psi_\alpha(q') \end{aligned}$$

Proof : 1) Proving a general lower bound - 2

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha \in (0, 1)$: $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q'(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha + \left(\frac{q'(y)}{p(y)} \right)^\alpha - 1 - \alpha \left(\frac{q'(y)}{p(y)} - 1 \right) \right] p(y) \nu(dy) \\ &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha \right] p(y) \nu(dy) + \Psi_\alpha(q') \end{aligned}$$

Proof : 1) Proving a general lower bound - 2

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha \in (0, 1)$: $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$

$$\Psi_\alpha(q) = \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha \right] p(y) \nu(dy) + \Psi_\alpha(q')$$

Proof : 1) Proving a general lower bound - 2

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1)$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha \in (0, 1)$: $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha \right] p(y) \nu(dy) + \Psi_\alpha(q') \\ &= \int_Y \left(\frac{q'(y)}{p(y)} \right)^\alpha \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{q'(y)} \right)^\alpha - 1 \right] p(y) \nu(dy) + \Psi_\alpha(q') \end{aligned}$$

Proof : 1) Proving a general lower bound - 2

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha \in (0, 1)$: $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha \right] p(y) \nu(dy) + \Psi_\alpha(q') \\ &= \int_Y \left(\frac{q'(y)}{p(y)} \right)^\alpha \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{q'(y)} \right)^\alpha - 1 \right] p(y) \nu(dy) + \Psi_\alpha(q') \end{aligned}$$

Since $\log(u) \leq u - 1$ for all $u > 0$ and $\alpha \in (0, 1)$,

$$\frac{u-1}{\alpha(\alpha-1)} \leq \frac{\log(u)}{\alpha(\alpha-1)}$$

Proof : 1) Proving a general lower bound - 2

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1)$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha \in (0, 1)$: $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha \right] p(y) \nu(dy) + \Psi_\alpha(q') \\ &= \int_Y \left(\frac{q'(y)}{p(y)} \right)^\alpha \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{q'(y)} \right)^\alpha - 1 \right] p(y) \nu(dy) + \Psi_\alpha(q') \end{aligned}$$

Since $\log(u) \leq u - 1$ for all $u > 0$ and $\alpha \in (0, 1)$,

$$\frac{u^\alpha - 1}{\alpha(\alpha-1)} \leq \frac{\log(u^\alpha)}{\alpha(\alpha-1)}$$

Proof : 1) Proving a general lower bound - 2

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1]$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

By definition $\Psi_\alpha(q) = \int_Y f_\alpha \left(\frac{q(y)}{p(y)} \right) p(y) \nu(dy)$.

→ **Case** $\alpha \in (0, 1)$: $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$

$$\begin{aligned} \Psi_\alpha(q) &= \int_Y \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{p(y)} \right)^\alpha - \left(\frac{q'(y)}{p(y)} \right)^\alpha \right] p(y) \nu(dy) + \Psi_\alpha(q') \\ &= \int_Y \left(\frac{q'(y)}{p(y)} \right)^\alpha \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{q(y)}{q'(y)} \right)^\alpha - 1 \right] p(y) \nu(dy) + \Psi_\alpha(q') \end{aligned}$$

Since $\log(u) \leq u - 1$ for all $u > 0$ and $\alpha \in (0, 1)$,

$$\frac{u^\alpha - 1}{\alpha(\alpha-1)} \leq \frac{\log(u^\alpha)}{\alpha(\alpha-1)} = \frac{\log(u)}{\alpha-1}$$

Proof : 2) Derive (Weights) and (Components)

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1)$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

Notation : $\mu_n k(y) := \mu_{\lambda_n, \Theta_n} k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_{j,n}, y)$, for all $n \geq 1$ and all $y \in Y$

Proof : 2) Derive (Weights) and (Components)

Let $q, q' \in \mathcal{Q}$ and assume that $\Psi_\alpha(q') < \infty$. For all $\alpha \in [0, 1)$, it holds that

$$\frac{1}{1-\alpha} \int_Y q'(y)^\alpha p(y)^{1-\alpha} \log \left(\frac{q(y)}{q'(y)} \right) \nu(dy) \leq \Psi_\alpha(q') - \Psi_\alpha(q)$$

Notation : $\mu_n k(y) := \mu_{\lambda_n, \Theta_n} k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_{j,n}, y)$, for all $n \geq 1$ and all $y \in Y$

Proof : 2) Derive (Weights) and (Components)

Assume that $\Psi_\alpha(\mu_n k) < \infty$. For all $\alpha \in [0, 1)$, it holds that

$$\frac{1}{1-\alpha} \int_Y (\mu_n k(y))^\alpha p(y)^{1-\alpha} \log \left(\frac{\mu_{n+1} k(y)}{\mu_n k(y)} \right) \nu(dy) \leq \Psi_\alpha(\mu_n k) - \Psi_\alpha(\mu_{n+1} k)$$


Notation : $\mu_n k(y) := \mu_{\lambda_n, \Theta_n} k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_{j,n}, y)$, for all $n \geq 1$ and all $y \in Y$

Proof : 2) Derive (Weights) and (Components)

Assume that $\Psi_\alpha(\mu_n k) < \infty$. For all $\alpha \in [0, 1)$, it holds that

$$\frac{1}{1-\alpha} \int_Y (\mu_n k(y))^\alpha p(y)^{1-\alpha} \log \left(\frac{\mu_{n+1} k(y)}{\mu_n k(y)} \right) \nu(dy) \leq \Psi_\alpha(\mu_n k) - \Psi_\alpha(\mu_{n+1} k)$$

Notation : $\mu_n k(y) := \mu_{\lambda_n, \Theta_n} k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_{j,n}, y)$, for all $n \geq 1$ and all $y \in Y$


 $u \mapsto \frac{1}{1-\alpha} \log(u)$ is concave

Proof : 2) Derive (Weights) and (Components)

Assume that $\Psi_\alpha(\mu_n k) < \infty$. For all $\alpha \in [0, 1)$, it holds that

$$\frac{1}{1-\alpha} \int_Y (\mu_n k(y))^\alpha p(y)^{1-\alpha} \log \left(\frac{\mu_{n+1} k(y)}{\mu_n k(y)} \right) \nu(dy) \leq \Psi_\alpha(\mu_n k) - \Psi_\alpha(\mu_{n+1} k)$$

Notation : $\mu_n k(y) := \mu_{\lambda_n, \Theta_n} k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_{j,n}, y)$, for all $n \geq 1$ and all $y \in Y$

 $u \mapsto \frac{1}{1-\alpha} \log(u)$ is concave

Jensen's inequality: for all $y \in Y$ and all $n \geq 1$,


$$\begin{aligned} \frac{1}{1-\alpha} \log \left(\frac{\mu_{n+1} k(y)}{\mu_n k(y)} \right) &= \frac{1}{1-\alpha} \log \left(\frac{\sum_{j=1}^J \frac{\lambda_{j,n} k(\theta_{j,n}, y)}{\sum_{\ell=1}^J \lambda_{\ell,n} k(\theta_{\ell,n}, y)} \frac{\lambda_{j,n+1} k(\theta_{j,n+1}, y)}{\lambda_{j,n} k(\theta_{j,n}, y)}}{\sum_{j=1}^J \frac{\lambda_{j,n} k(\theta_{j,n}, y)}{\sum_{\ell=1}^J \lambda_{\ell,n} k(\theta_{\ell,n}, y)}} \right) \\ &\geq \frac{1}{1-\alpha} \sum_{j=1}^J \frac{\lambda_{j,n} k(\theta_{j,n}, y)}{\sum_{\ell=1}^J \lambda_{\ell,n} k(\theta_{\ell,n}, y)} \log \left(\frac{\lambda_{j,n+1} k(\theta_{j,n+1}, y)}{\lambda_{j,n} k(\theta_{j,n}, y)} \right) \end{aligned}$$

Proof : 2) Derive (Weights) and (Components)

Assume that $\Psi_\alpha(\mu_n k) < \infty$. For all $\alpha \in [0, 1)$, it holds that

$$\frac{1}{1-\alpha} \int_Y (\mu_n k(y))^\alpha p(y)^{1-\alpha} \log \left(\frac{\mu_{n+1} k(y)}{\mu_n k(y)} \right) \nu(dy) \leq \Psi_\alpha(\mu_n k) - \Psi_\alpha(\mu_{n+1} k)$$

Notation : $\mu_n k(y) := \mu_{\lambda_n, \Theta_n} k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_{j,n}, y)$, for all $n \geq 1$ and all $y \in Y$

 $u \mapsto \frac{1}{1-\alpha} \log(u)$ is concave

Jensen's inequality: for all $y \in Y$ and all $n \geq 1$,

$$\begin{aligned} \frac{1}{1-\alpha} \log \left(\frac{\mu_{n+1} k(y)}{\mu_n k(y)} \right) &= \frac{1}{1-\alpha} \log \left(\frac{\sum_{j=1}^J \frac{\lambda_{j,n} k(\theta_{j,n}, y)}{\sum_{\ell=1}^J \lambda_{\ell,n} k(\theta_{\ell,n}, y)} \frac{\lambda_{j,n+1} k(\theta_{j,n+1}, y)}{\lambda_{j,n} k(\theta_{j,n}, y)} \right) \\ &\geq \frac{1}{1-\alpha} \sum_{j=1}^J \frac{\lambda_{j,n} k(\theta_{j,n}, y)}{\sum_{\ell=1}^J \lambda_{\ell,n} k(\theta_{\ell,n}, y)} \log \left(\frac{\lambda_{j,n+1} k(\theta_{j,n+1})}{\lambda_{j,n} k(\theta_{j,n}, y)} \right) \end{aligned}$$

that is :


$$\frac{1}{1-\alpha} \log \left(\frac{\mu_{n+1} k(y)}{\mu_n k(y)} \right) \geq \frac{1}{1-\alpha} \sum_{j=1}^J \lambda_{j,n} \frac{k(\theta_{j,n}, y)}{\mu_n k(y)} \log \left(\frac{\lambda_{j,n+1} k(\theta_{j,n+1})}{\lambda_{j,n} k(\theta_{j,n}, y)} \right)$$

Proof : 2) Derive (Weights) and (Components)

Assume that $\Psi_\alpha(\mu_n k) < \infty$. For all $\alpha \in [0, 1)$, it holds that

$$\frac{1}{1-\alpha} \int_Y (\mu_n k(y))^\alpha p(y)^{1-\alpha} \log \left(\frac{\mu_{n+1} k(y)}{\mu_n k(y)} \right) \nu(dy) \leq \Psi_\alpha(\mu_n k) - \Psi_\alpha(\mu_{n+1} k)$$

Notation : $\mu_n k(y) := \mu_{\lambda_n, \Theta_n} k(y) = \sum_{j=1}^J \lambda_{j,n} k(\theta_{j,n}, y)$, for all $n \geq 1$ and all $y \in Y$

 $u \mapsto \frac{1}{1-\alpha} \log(u)$ is concave

Jensen's inequality: for all $y \in Y$ and all $n \geq 1$,

$$\begin{aligned} \frac{1}{1-\alpha} \log \left(\frac{\mu_{n+1} k(y)}{\mu_n k(y)} \right) &= \frac{1}{1-\alpha} \log \left(\frac{\sum_{j=1}^J \frac{\lambda_{j,n} k(\theta_{j,n}, y)}{\sum_{\ell=1}^J \lambda_{\ell,n} k(\theta_{\ell,n}, y)} \frac{\lambda_{j,n+1} k(\theta_{j,n+1}, y)}{\lambda_{j,n} k(\theta_{j,n}, y)} \right) \\ &\geq \frac{1}{1-\alpha} \sum_{j=1}^J \frac{\lambda_{j,n} k(\theta_{j,n}, y)}{\sum_{\ell=1}^J \lambda_{\ell,n} k(\theta_{\ell,n}, y)} \log \left(\frac{\lambda_{j,n+1} k(\theta_{j,n+1}, y)}{\lambda_{j,n} k(\theta_{j,n}, y)} \right) \end{aligned}$$

that is :

$$\frac{1}{1-\alpha} \log \left(\frac{\mu_{n+1} k(y)}{\mu_n k(y)} \right) \geq \frac{1}{1-\alpha} \sum_{j=1}^J \lambda_{j,n} \frac{k(\theta_{j,n}, y)}{\mu_n k(y)} \log \left(\frac{\lambda_{j,n+1} k(\theta_{j,n+1}, y)}{\lambda_{j,n} k(\theta_{j,n}, y)} \right)$$

To finish the proof :

- (i) multiply by $(\mu_n k(y))^\alpha p(y)^{1-\alpha}$ on both sides $(\varphi_{j,n}^{(\alpha)}(y) = k(\theta_{j,n}, y) \left(\frac{p(y)}{\mu_n k(y)} \right)^{1-\alpha})$
- (ii) integrate with respect to $\nu(dy)$

Commenting the conditions for a monotonic decrease

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

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

where $\varphi_{j,n}^{(\alpha)}(y) = k(\theta_{j,n}, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1}$

- ① (Weights) and (Components) permit **separate/simultaneous** updates
- ② They are satisfied for $\lambda_{n+1} = \lambda_n$ and $\Theta_{n+1} = \Theta_n$ respectively
- ③ The dependency is **simpler** in (Weights)

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

$$\lambda_{n+1} = \operatorname{argmax}_{\lambda \in S_J} \sum_{j=1}^J \left[\lambda_{j,n} \int_Y \varphi_{j,n}^{(\alpha)}(y) \nu(dy) \right] \log \left(\frac{\lambda_j}{\lambda_{j,n}} \right)$$

Commenting the conditions for a monotonic decrease

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

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

where $\varphi_{j,n}^{(\alpha)}(y) = k(\theta_{j,n}, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1}$

- 1 (Weights) and (Components) permit **separate/simultaneous** updates
- 2 They are satisfied for $\lambda_{n+1} = \lambda_n$ and $\Theta_{n+1} = \Theta_n$ respectively
- 3 The dependency is **simpler** in (Weights)

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

$$\lambda_{n+1} = \operatorname{argmax}_{\lambda \in S_J} \sum_{j=1}^J \left[\lambda_{j,n} \int_Y \varphi_{j,n}^{(\alpha)}(y) \nu(dy) \right] \log \left(\frac{\lambda_j}{\lambda_{j,n}} \right)$$

Commenting the conditions for a monotonic decrease

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

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

where $\varphi_{j,n}^{(\alpha)}(y) = k(\theta_{j,n}, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1}$

- ❶ (Weights) and (Components) permit **separate/simultaneous** updates
- ❷ They are satisfied for $\lambda_{n+1} = \lambda_n$ and $\Theta_{n+1} = \Theta_n$ respectively
- ❸ The dependency is **simpler** in (Weights)

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

$$\lambda_{n+1} = \operatorname{argmax}_{\lambda \in S_J} \sum_{j=1}^J \left[\lambda_{j,n} \int_Y \varphi_{j,n}^{(\alpha)}(y) \nu(dy) \right] \log \left(\frac{\lambda_j}{\lambda_{j,n}} \right)$$

Commenting the conditions for a monotonic decrease

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

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

where $\varphi_{j,n}^{(\alpha)}(y) = k(\theta_{j,n}, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1}$

- ❶ (Weights) and (Components) permit **separate/simultaneous** updates
- ❷ They are satisfied for $\lambda_{n+1} = \lambda_n$ and $\Theta_{n+1} = \Theta_n$ respectively
- ❸ The dependency is **simpler** in (Weights)

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

$$\lambda_{n+1} = \operatorname{argmax}_{\lambda \in S_J} \sum_{j=1}^J \left[\lambda_{j,n} \int_Y \varphi_{j,n}^{(\alpha)}(y) \nu(dy) \right] \log \left(\frac{\lambda_j}{\lambda_{j,n}} \right)$$

Commenting the conditions for a monotonic decrease

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

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

where $\varphi_{j,n}^{(\alpha)}(y) = k(\theta_{j,n}, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1}$

- ❶ (Weights) and (Components) permit **separate/simultaneous** updates
- ❷ They are satisfied for $\lambda_{n+1} = \lambda_n$ and $\Theta_{n+1} = \Theta_n$ respectively
- ❸ The dependency is **simpler** in (Weights)

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

$$\lambda_{n+1} = \operatorname{argmax}_{\lambda \in S_J} \sum_{j=1}^J \left[\lambda_{j,n} \int_Y \varphi_{j,n}^{(\alpha)}(y) \nu(dy) \right] \log \left(\frac{\lambda_j}{\lambda_{j,n}} \right)$$

Commenting the conditions for a monotonic decrease

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

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

where $\varphi_{j,n}^{(\alpha)}(y) = k(\theta_{j,n}, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1}$

- ❶ (Weights) and (Components) permit **separate/simultaneous** updates
- ❷ They are satisfied for $\lambda_{n+1} = \lambda_n$ and $\Theta_{n+1} = \Theta_n$ respectively
- ❸ The dependency is **simpler** in (Weights)

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

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \int_Y \varphi_{j,n}^{(\alpha)}(y) \nu(dy)}{\sum_{\ell=1}^J \lambda_{\ell,n} \int_Y \varphi_{\ell,n}^{(\alpha)}(y) \nu(dy)}, \quad j = 1 \dots J$$

Commenting the conditions for a monotonic decrease

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

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

where $\varphi_{j,n}^{(\alpha)}(y) = k(\theta_{j,n}, y) \left(\frac{\mu_n k(y)}{p(y)} \right)^{\alpha-1}$

- ❶ (Weights) and (Components) permit **separate/simultaneous** updates
- ❷ They are satisfied for $\lambda_{n+1} = \lambda_n$ and $\Theta_{n+1} = \Theta_n$ respectively
- ❸ The dependency is **simpler** in (Weights)

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

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

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:

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

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

→ We recover the **Power Descent** algorithm from Part 2

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:

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

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

→ We recover the **Power Descent** algorithm from Part 2

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:

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

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

→ We recover the **Power Descent** algorithm from Part 2

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:

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

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

→ We recover the **Power Descent** algorithm from Part 2

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:

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

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

→ We recover the **Power Descent** algorithm from Part 2

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

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

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

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

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

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

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

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

Towards simultaneous updates

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

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

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

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

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

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

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

Towards simultaneous updates

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

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

$$\theta_{j,n+1} = \operatorname{argmax}_{\theta \in \mathsf{T}} \int_Y \varphi_{j,n}^{(\alpha)}(y) \log \left(\frac{k(\theta, y)}{k(\theta_{j,n}, y)} \right) \nu(dy)$$

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

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

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

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

Two questions at this stage

- ① Can we derive practical updates from the maximisation / gradient-based approaches?
- ② Do those approaches relate to the existing literature?

Two questions at this stage

- ① Can we derive practical updates from the maximisation / gradient-based approaches?
- ② Do those approaches relate to the existing literature?

Outline

- 1 Monotonic Alpha-Divergence Minimisation
- 2 Maximisation approach
- 3 Gradient-based approach
- 4 Numerical Experiments
- 5 Conclusion of Part 3

Maximisation approach

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

For all $j = 1 \dots J$,

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

Maximisation approach

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

For all $j = 1 \dots J$,

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

Maximisation approach

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

For all $j = 1 \dots J$, $b_{j,n} \geq 0$ and

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

→ We have added a regularisation term!

Maximisation approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

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

$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

where $\tilde{\Sigma}_{j,n} = \Sigma_{j,n} + (m_{j,n+1} - m_{j,n})(m_{j,n+1} - m_{j,n})^T$ and $\gamma_{j,n}$ depends on $b_{j,n}$.

→ Considering all possible values of $b_{j,n}$, we have $\gamma_{j,n} \in (0, 1]$

Interpretation : tradeoff between

- an update close to $\theta_{j,n} = (m_{j,n}, \Sigma_{j,n})$ [$\gamma_{j,n} \rightarrow 0$]
- an update that chooses the Gaussian with the same mean and covariance matrix as $\check{\varphi}_{j,n}^{(\alpha)}$ [$\gamma_{j,n} = 1$]

Why does it matter? In practice, Monte Carlo approximations!

Maximisation approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

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

$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

where $\tilde{\Sigma}_{j,n} = \Sigma_{j,n} + (m_{j,n+1} - m_{j,n})(m_{j,n+1} - m_{j,n})^T$ and $\gamma_{j,n}$ depends on $b_{j,n}$.

→ Considering all possible values of $b_{j,n}$, we have $\gamma_{j,n} \in (0, 1]$

Interpretation : tradeoff between

- an update close to $\theta_{j,n} = (m_{j,n}, \Sigma_{j,n})$ [$\gamma_{j,n} \rightarrow 0$]
- an update that chooses the Gaussian with the same mean and covariance matrix as $\check{\varphi}_{j,n}^{(\alpha)}$ [$\gamma_{j,n} = 1$]

Why does it matter? In practice, Monte Carlo approximations!

Maximisation approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

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

$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

where $\tilde{\Sigma}_{j,n} = \Sigma_{j,n} + (m_{j,n+1} - m_{j,n})(m_{j,n+1} - m_{j,n})^T$ and $\gamma_{j,n}$ depends on $b_{j,n}$.

→ Considering all possible values of $b_{j,n}$, we have $\gamma_{j,n} \in (0, 1]$

Interpretation : tradeoff between

- an update close to $\theta_{j,n} = (m_{j,n}, \Sigma_{j,n})$ [$\gamma_{j,n} \rightarrow 0$]
- an update that chooses the Gaussian with the same mean and covariance matrix as $\check{\varphi}_{j,n}^{(\alpha)}$ [$\gamma_{j,n} = 1$]

Why does it matter? In practice, Monte Carlo approximations!

Maximisation approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

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

$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

where $\tilde{\Sigma}_{j,n} = \Sigma_{j,n} + (m_{j,n+1} - m_{j,n})(m_{j,n+1} - m_{j,n})^T$ and $\gamma_{j,n}$ depends on $b_{j,n}$.

→ Considering all possible values of $b_{j,n}$, we have $\gamma_{j,n} \in (0, 1]$

Interpretation : tradeoff between

- an update close to $\theta_{j,n} = (m_{j,n}, \Sigma_{j,n})$ [$\gamma_{j,n} \rightarrow 0$]
- an update that chooses the Gaussian with the same mean and covariance matrix as $\check{\varphi}_{j,n}^{(\alpha)}$ [$\gamma_{j,n} = 1$]

Why does it matter? In practice, Monte Carlo approximations!

Maximisation approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

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

$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

where $\tilde{\Sigma}_{j,n} = \Sigma_{j,n} + (m_{j,n+1} - m_{j,n})(m_{j,n+1} - m_{j,n})^T$ and $\gamma_{j,n}$ depends on $b_{j,n}$.

→ Considering all possible values of $b_{j,n}$, we have $\gamma_{j,n} \in (0, 1]$

Interpretation : tradeoff between

- an update close to $\theta_{j,n} = (m_{j,n}, \Sigma_{j,n})$ [$\gamma_{j,n} \rightarrow 0$]
- an update that chooses the Gaussian with the same mean and covariance matrix as $\check{\varphi}_{j,n}^{(\alpha)}$ [$\gamma_{j,n} = 1$]

Why does it matter? In practice, Monte Carlo approximations!

Maximisation approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

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

$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

where $\tilde{\Sigma}_{j,n} = \Sigma_{j,n} + (m_{j,n+1} - m_{j,n})(m_{j,n+1} - m_{j,n})^T$ and $\gamma_{j,n}$ depends on $b_{j,n}$.

→ Considering all possible values of $b_{j,n}$, we have $\gamma_{j,n} \in (0, 1]$

Interpretation : tradeoff between

- an update close to $\theta_{j,n} = (m_{j,n}, \Sigma_{j,n})$ [$\gamma_{j,n} \rightarrow 0$]
- an update that chooses the Gaussian with the same mean and covariance matrix as $\check{\varphi}_{j,n}^{(\alpha)}$ [$\gamma_{j,n} = 1$]

Why does it matter? In practice, Monte Carlo approximations!

Maximisation approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

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

$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

where $\tilde{\Sigma}_{j,n} = \Sigma_{j,n} + (m_{j,n+1} - m_{j,n})(m_{j,n+1} - m_{j,n})^T$ and $\gamma_{j,n}$ depends on $b_{j,n}$.

→ Considering all possible values of $b_{j,n}$, we have $\gamma_{j,n} \in (0, 1]$

Interpretation : tradeoff between

- an update close to $\theta_{j,n} = (m_{j,n}, \Sigma_{j,n})$ [$\gamma_{j,n} \rightarrow 0$]
- an update that chooses the Gaussian with the same mean and covariance matrix as $\check{\varphi}_{j,n}^{(\alpha)}$ [$\gamma_{j,n} = 1$]

Why does it matter? In practice, Monte Carlo approximations!

Maximisation approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

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

$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

where $\tilde{\Sigma}_{j,n} = \Sigma_{j,n} + (m_{j,n+1} - m_{j,n})(m_{j,n+1} - m_{j,n})^T$ and $\gamma_{j,n}$ depends on $b_{j,n}$.

→ Considering all possible values of $b_{j,n}$, we have $\gamma_{j,n} \in (0, 1]$

Interpretation : tradeoff between

- an update close to $\theta_{j,n} = (m_{j,n}, \Sigma_{j,n})$ [$\gamma_{j,n} \rightarrow 0$]
- an update that chooses the Gaussian with the same mean and covariance matrix as $\check{\varphi}_{j,n}^{(\alpha)}$ [$\gamma_{j,n} = 1$]

Why does it matter? In practice, Monte Carlo approximations!

Maximisation approach for GMMs : related work

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

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

$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

Consider the case $\alpha = 0$, $\gamma_{j,n} = 1$, $\eta_n = 1$, $\kappa = 0$, set $t_{j,n} = \frac{\lambda_{j,n} k(\theta_{j,n}, \cdot)}{\mu_{\lambda_n, \Theta_n} k}$ and $\tilde{p} = p / \int p d\nu$

→ The M-PMC algorithm a.k.a 'Integrated EM' for GMMs

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

Core insight : We have **generalised** an integrated EM algorithm for mixture models optimisation!

Maximisation approach for GMMs : related work

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

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

$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = (1 - \gamma_{j,n}) \tilde{\Sigma}_{j,n} + \gamma_{j,n} \int_Y \check{\varphi}_{j,n}^{(\alpha)}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

Consider the case $\alpha = 0$, $\gamma_{j,n} = 1$, $\eta_n = 1$, $\kappa = 0$, set $t_{j,n} = \frac{\lambda_{j,n} k(\theta_{j,n}, \cdot)}{\mu_{\lambda_n, \Theta_n} k}$ and $\tilde{p} = p / \int p d\nu$

→ The M-PMC algorithm a.k.a 'Integrated EM' for GMMs

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

Core insight : We have **generalised** an integrated EM algorithm for mixture models optimisation!

Maximisation approach for GMMs : related work

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

$$\lambda_{j,n+1} = \int_Y t_{j,n}(y) \tilde{p}(y) \nu(dy)$$

$$m_{j,n+1} = \int_Y t_{j,n}(y) \tilde{p}(y) y \nu(dy)$$

$$\Sigma_{j,n+1} = \int_Y t_{j,n}(y) \tilde{p}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)$$

Consider the case $\alpha = 0$, $\gamma_{j,n} = 1$, $\eta_n = 1$, $\kappa = 0$, set $t_{j,n} = \frac{\lambda_{j,n} k(\theta_{j,n}, \cdot)}{\mu_{\lambda_n, \Theta_n} k}$ and $\tilde{p} = p / \int p d\nu$

→ The M-PMC algorithm a.k.a 'Integrated EM' for GMMs

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

Core insight : We have generalised an integrated EM algorithm for mixture models optimisation!

Maximisation approach for GMMs : related work

Set $k(\theta, y) = \mathcal{N}(y; m, \Sigma)$ with $\theta = (m, \Sigma)$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$.

For all $j = 1 \dots J$,

$$\begin{aligned}\lambda_{j,n+1} &= \int_{\mathcal{Y}} t_{j,n}(y) \tilde{p}(y) \nu(dy) \\ m_{j,n+1} &= \int_{\mathcal{Y}} t_{j,n}(y) \tilde{p}(y) y \nu(dy) \\ \Sigma_{j,n+1} &= \int_{\mathcal{Y}} t_{j,n}(y) \tilde{p}(y) (y - m_{j,n+1})(y - m_{j,n+1})^T \nu(dy)\end{aligned}$$

Consider the case $\alpha = 0$, $\gamma_{j,n} = 1$, $\eta_n = 1$, $\kappa = 0$, set $t_{j,n} = \frac{\lambda_{j,n} k(\theta_{j,n}, \cdot)}{\mu_{\lambda_n, \Theta_n} k}$ and $\tilde{p} = p / \int p d\nu$

→ The M-PMC algorithm a.k.a 'Integrated EM' for GMMs

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

Core insight : We have **generalised** an integrated EM algorithm for mixture models optimisation!

Outline

- 1 Monotonic Alpha-Divergence Minimisation
- 2 Maximisation approach
- 3 Gradient-based approach**
- 4 Numerical Experiments
- 5 Conclusion of Part 3

Gradient-based approach

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

For all $j = 1 \dots J$, $\gamma_{j,n} \in (0, 1]$

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

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

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

We have that

$$\nabla g_{j,n}(\theta)|_{\theta=\theta_{j,n}} = \int_Y \frac{\varphi_{j,n}^{(\alpha)}(y)}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta,y)=(\theta_{j,n},y)} \nu(dy)$$

→ There might be links with Gradient Descent steps...

Gradient-based approach

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

For all $j = 1 \dots J$, $\gamma_{j,n} \in (0, 1]$

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

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

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

We have that

$$\nabla g_{j,n}(\theta)|_{\theta=\theta_{j,n}} = \int_Y \frac{\varphi_{j,n}^{(\alpha)}(y)}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta,y)=(\theta_{j,n},y)} \nu(dy)$$

→ There might be links with Gradient Descent steps...

Gradient-based approach

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

For all $j = 1 \dots J$, $\gamma_{j,n} \in (0, 1]$

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

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

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

We have that

$$\nabla g_{j,n}(\theta)|_{\theta=\theta_{j,n}} = \int_Y \frac{\varphi_{j,n}^{(\alpha)}(y)}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta,y)=(\theta_{j,n},y)} \nu(dy)$$

→ There might be links with Gradient Descent steps...

Gradient-based approach

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

For all $j = 1 \dots J$, $\gamma_{j,n} \in (0, 1]$

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

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

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

We have that

$$\nabla g_{j,n}(\theta)|_{\theta=\theta_{j,n}} = \int_Y \frac{\varphi_{j,n}^{(\alpha)}(y)}{\alpha - 1} \frac{\partial \log k(\theta, y)}{\partial \theta} \Big|_{(\theta,y)=(\theta_{j,n},y)} \nu(dy)$$

→ There might be links with Gradient Descent steps...

Gradient-based approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \sigma^2 \mathbf{I}_d)$ with $\theta = m$, fixed $\sigma > 0$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$

For all $j = 1 \dots J$,

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

where $\gamma_{j,n} \in (0, 1]$

→ Interpretation :

- 1 Maximisation and gradient-based approach coincide when $\Sigma = \sigma^2 \mathbf{I}_d$ with σ fixed
- 2 We recognise Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_\alpha(\mu k; p)$ by setting

$$\gamma_{j,n} = \gamma'_{j,n} \frac{\lambda_{j,n} \int_Y \varphi_{j,n}^{(\alpha)}(y) \nu(dy)}{\int_Y (\mu_n k(y))^\alpha p(y)^{1-\alpha} \nu(dy)} \quad \text{with} \quad \gamma'_{j,n} \in (0, 1]$$

Compatibility between Gradient Descent steps w.r.t Θ and mixture weights updates (and even covariance matrices updates)!

Gradient-based approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \sigma^2 \mathbf{I}_d)$ with $\theta = m$, fixed $\sigma > 0$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$

For all $j = 1 \dots J$,

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

where $\gamma_{j,n} \in (0, 1]$

→ Interpretation :

- 1 Maximisation and gradient-based approach coincide when $\Sigma = \sigma^2 \mathbf{I}_d$ with σ fixed
- 2 We recognise Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_\alpha(\mu k; p)$ by setting

$$\gamma_{j,n} = \gamma'_{j,n} \frac{\lambda_{j,n} \int_Y \varphi_{j,n}^{(\alpha)}(y) \nu(dy)}{\int_Y (\mu_n k(y))^\alpha p(y)^{1-\alpha} \nu(dy)} \quad \text{with} \quad \gamma'_{j,n} \in (0, 1]$$

Compatibility between Gradient Descent steps w.r.t Θ and mixture weights updates (and even covariance matrices updates)!

Gradient-based approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \sigma^2 \mathbf{I}_d)$ with $\theta = m$, fixed $\sigma > 0$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$

For all $j = 1 \dots J$,

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_{\mathbf{Y}} \varphi_{j,n}^{(\alpha)}(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_{\mathbf{Y}} \varphi_{\ell,n}^{(\alpha)}(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}$$
$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_{\mathbf{Y}} \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

where $\gamma_{j,n} \in (0, 1]$

→ Interpretation :

- 1 Maximisation and gradient-based approach coincide when $\Sigma = \sigma^2 \mathbf{I}_d$ with σ fixed
- 2 We recognise Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_\alpha(\mu k; p)$ by setting

$$\gamma_{j,n} = \gamma'_{j,n} \frac{\lambda_{j,n} \int_{\mathbf{Y}} \varphi_{j,n}^{(\alpha)}(y) \nu(dy)}{\int_{\mathbf{Y}} (\mu_n k(y))^\alpha p(y)^{1-\alpha} \nu(dy)} \quad \text{with} \quad \gamma'_{j,n} \in (0, 1]$$

Compatibility between Gradient Descent steps w.r.t Θ and mixture weights updates (and even covariance matrices updates)!

Gradient-based approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \sigma^2 \mathbf{I}_d)$ with $\theta = m$, fixed $\sigma > 0$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$

For all $j = 1 \dots J$,

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

where $\gamma_{j,n} \in (0, 1]$

→ Interpretation :

- 1 Maximisation and gradient-based approach coincide when $\Sigma = \sigma^2 \mathbf{I}_d$ with σ fixed
- 2 We recognise Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_\alpha(\mu k; p)$ by setting

$$\gamma_{j,n} = \gamma'_{j,n} \frac{\lambda_{j,n} \int_Y \varphi_{j,n}^{(\alpha)}(y) \nu(dy)}{\int_Y (\mu_n k(y))^\alpha p(y)^{1-\alpha} \nu(dy)} \quad \text{with} \quad \gamma'_{j,n} \in (0, 1]$$

Compatibility between Gradient Descent steps w.r.t Θ and mixture weights updates (and even covariance matrices updates)!

Gradient-based approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \sigma^2 \mathbf{I}_d)$ with $\theta = m$, fixed $\sigma > 0$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$

For all $j = 1 \dots J$,

$$\lambda_{j,n+1} = \frac{\lambda_{j,n} \left[\int_{\mathbf{Y}} \varphi_{j,n}^{(\alpha)}(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}{\sum_{\ell=1}^J \lambda_{\ell,n} \left[\int_{\mathbf{Y}} \varphi_{\ell,n}^{(\alpha)}(y) \nu(dy) + (\alpha - 1) \kappa \right]^{\eta_n}}$$
$$m_{j,n+1} = (1 - \gamma_{j,n}) m_{j,n} + \gamma_{j,n} \int_{\mathbf{Y}} \check{\varphi}_{j,n}^{(\alpha)}(y) y \nu(dy)$$

where $\gamma_{j,n} \in (0, 1]$

→ Interpretation :

- 1 Maximisation and gradient-based approach coincide when $\Sigma = \sigma^2 \mathbf{I}_d$ with σ fixed
- 2 We recognise Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_{\alpha}(\mu k; p)$ by setting

$$\gamma_{j,n} = \gamma'_{j,n} \frac{\lambda_{j,n} \int_{\mathbf{Y}} \varphi_{j,n}^{(\alpha)}(y) \nu(dy)}{\int_{\mathbf{Y}} (\mu_n k(y))^{\alpha} p(y)^{1-\alpha} \nu(dy)} \quad \text{with} \quad \gamma'_{j,n} \in (0, 1]$$

Compatibility between Gradient Descent steps w.r.t Θ and mixture weights updates (and even covariance matrices updates)!

Gradient-based approach for GMMs

Set $k(\theta, y) = \mathcal{N}(y; m, \sigma^2 \mathbf{I}_d)$ with $\theta = m$, fixed $\sigma > 0$ and $\check{\varphi}_{j,n}^{(\alpha)} = \varphi_{j,n}^{(\alpha)} / \int \varphi_{j,n}^{(\alpha)} d\nu$

For all $j = 1 \dots J$,

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

where $\gamma_{j,n} \in (0, 1]$

→ Interpretation :

- 1 Maximisation and gradient-based approach coincide when $\Sigma = \sigma^2 \mathbf{I}_d$ with σ fixed
- 2 We recognise Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_\alpha(\mu k; p)$ by setting

$$\gamma_{j,n} = \gamma'_{j,n} \frac{\lambda_{j,n} \int_Y \varphi_{j,n}^{(\alpha)}(y) \nu(dy)}{\int_Y (\mu_n k(y))^\alpha p(y)^{1-\alpha} \nu(dy)} \quad \text{with} \quad \gamma'_{j,n} \in (0, 1]$$

Compatibility between Gradient Descent steps w.r.t Θ and mixture weights updates (and even covariance matrices updates)!

At this stage

We expressed conditions on λ and Θ ensuring a systematic decrease in $\Psi_\alpha(\mu_{\lambda,\Theta})$:

- ① Updates on λ linked to the gradient-based **Power Descent**
- ② Updates on Θ :
 - Maximisation approach : generalises an Integrated EM
 - Gradient-based approach : links with Gradient Descent algorithms

	Improvements of our framework
Gradient Descent w.r.t Θ on $-\alpha^{-1}\mathcal{L}_\alpha(\mu_k; p)$	Simultaneous optimisation w.r.t $(\lambda_n)_{n \geq 1}$ $\lambda_{j,n}$ needs not to be as a factor in the means updates Covariance matrices update formulas
Power Descent	Simultaneous optimisation w.r.t $(\Theta_n)_{n \geq 1}$ Convergence towards a local optimum of the full algorithm
M-PMC algorithm	$\alpha \in [0, 1)$ (prev. $\alpha = 0$) $\eta_n \in (0, 1]$ and $(\alpha - 1)\kappa_n \geq 0$ (prev. $\eta_n = 1, \kappa_n = 0$) $b_{j,n} \geq 0$ (prev. $b_{j,n} = 0$)

At this stage

We expressed conditions on λ and Θ ensuring a systematic decrease in $\Psi_\alpha(\mu_{\lambda,\Theta})$:

- ① Updates on λ linked to the gradient-based **Power Descent**
- ② Updates on Θ :
 - Maximisation approach : generalises an Integrated EM
 - Gradient-based approach : links with Gradient Descent algorithms

	Improvements of our framework
Gradient Descent w.r.t Θ on $-\alpha^{-1}\mathcal{L}_\alpha(\mu_k; p)$	Simultaneous optimisation w.r.t $(\lambda_n)_{n \geq 1}$ $\lambda_{j,n}$ needs not to be as a factor in the means updates Covariance matrices update formulas
Power Descent	Simultaneous optimisation w.r.t $(\Theta_n)_{n \geq 1}$ Convergence towards a local optimum of the full algorithm
M-PMC algorithm	$\alpha \in [0, 1)$ (prev. $\alpha = 0$) $\eta_n \in (0, 1]$ and $(\alpha - 1)\kappa_n \geq 0$ (prev. $\eta_n = 1, \kappa_n = 0$) $b_{j,n} \geq 0$ (prev. $b_{j,n} = 0$)

At this stage

We expressed conditions on λ and Θ ensuring a systematic decrease in $\Psi_\alpha(\mu_{\lambda,\Theta})$:

- ❶ Updates on λ linked to the gradient-based **Power Descent**
- ❷ Updates on Θ :
 - Maximisation approach : generalises an Integrated EM
 - Gradient-based approach : links with Gradient Descent algorithms

	Improvements of our framework
Gradient Descent w.r.t Θ on $-\alpha^{-1}\mathcal{L}_\alpha(\mu_k; p)$	Simultaneous optimisation w.r.t $(\lambda_n)_{n \geq 1}$ $\lambda_{j,n}$ needs not to be as a factor in the means updates Covariance matrices update formulas
Power Descent	Simultaneous optimisation w.r.t $(\Theta_n)_{n \geq 1}$ Convergence towards a local optimum of the full algorithm
M-PMC algorithm	$\alpha \in [0, 1)$ (prev. $\alpha = 0$) $\eta_n \in (0, 1]$ and $(\alpha - 1)\kappa_n \geq 0$ (prev. $\eta_n = 1, \kappa_n = 0$) $b_{j,n} \geq 0$ (prev. $b_{j,n} = 0$)

At this stage

We expressed conditions on λ and Θ ensuring a systematic decrease in $\Psi_\alpha(\mu_{\lambda,\Theta})$:

- ❶ Updates on λ linked to the gradient-based **Power Descent**
- ❷ Updates on Θ :
 - Maximisation approach : generalises an Integrated EM
 - Gradient-based approach : links with Gradient Descent algorithms

	Improvements of our framework
Gradient Descent w.r.t Θ on $-\alpha^{-1}\mathcal{L}_\alpha(\mu_k; p)$	Simultaneous optimisation w.r.t $(\lambda_n)_{n \geq 1}$ $\lambda_{j,n}$ needs not to be as a factor in the means updates Covariance matrices update formulas
Power Descent	Simultaneous optimisation w.r.t $(\Theta_n)_{n \geq 1}$ Convergence towards a local optimum of the full algorithm
M-PMC algorithm	$\alpha \in [0, 1)$ (prev. $\alpha = 0$) $\eta_n \in (0, 1]$ and $(\alpha - 1)\kappa_n \geq 0$ (prev. $\eta_n = 1, \kappa_n = 0$) $b_{j,n} \geq 0$ (prev. $b_{j,n} = 0$)

At this stage

We expressed conditions on λ and Θ ensuring a systematic decrease in $\Psi_\alpha(\mu_{\lambda,\Theta})$:

- ❶ Updates on λ linked to the gradient-based **Power Descent**
- ❷ Updates on Θ :
 - Maximisation approach : generalises an Integrated EM
 - Gradient-based approach : links with Gradient Descent algorithms

	Improvements of our framework
Gradient Descent w.r.t Θ on $-\alpha^{-1}\mathcal{L}_\alpha(\mu_k; p)$	Simultaneous optimisation w.r.t $(\lambda_n)_{n \geq 1}$ $\lambda_{j,n}$ needs not to be as a factor in the means updates Covariance matrices update formulas
Power Descent	Simultaneous optimisation w.r.t $(\Theta_n)_{n \geq 1}$ Convergence towards a local optimum of the full algorithm
M-PMC algorithm	$\alpha \in [0, 1)$ (prev. $\alpha = 0$) $\eta_n \in (0, 1]$ and $(\alpha - 1)\kappa_n \geq 0$ (prev. $\eta_n = 1, \kappa_n = 0$) $b_{j,n} \geq 0$ (prev. $b_{j,n} = 0$)

At this stage

We expressed conditions on λ and Θ ensuring a systematic decrease in $\Psi_\alpha(\mu_{\lambda,\Theta})$:

- ❶ Updates on λ linked to the gradient-based **Power Descent**
- ❷ Updates on Θ :
 - Maximisation approach : generalises an Integrated EM
 - Gradient-based approach : links with Gradient Descent algorithms

	Improvements of our framework
Gradient Descent w.r.t Θ on $-\alpha^{-1}\mathcal{L}_\alpha(\mu k; p)$	Simultaneous optimisation w.r.t $(\lambda_n)_{n \geq 1}$ $\lambda_{j,n}$ needs not to be as a factor in the means updates Covariance matrices update formulas
Power Descent	Simultaneous optimisation w.r.t $(\Theta_n)_{n \geq 1}$ Convergence towards a local optimum of the full algorithm
M-PMC algorithm	$\alpha \in [0, 1)$ (prev. $\alpha = 0$) $\eta_n \in (0, 1]$ and $(\alpha - 1)\kappa_n \geq 0$ (prev. $\eta_n = 1, \kappa_n = 0$) $b_{j,n} \geq 0$ (prev. $b_{j,n} = 0$)

Outline

- 1 Monotonic Alpha-Divergence Minimisation
- 2 Maximisation approach
- 3 Gradient-based approach
- 4 Numerical Experiments**
- 5 Conclusion of Part 3

Monte Carlo approximations

Algorithm 1: Gaussian Mixture Models optimisation

At iteration n ,

- ➊ Draw independently M samples $(Y_{m,n})_{1 \leq m \leq M}$ from the proposal q_n .
- ➋ For all $j = 1 \dots J$, set:

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

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

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

→ Here $\hat{\varphi}_{j,n}^{(\alpha)}(y) = \frac{\varphi_{j,n}^{(\alpha)}(y)}{q_n(y)}$, $\gamma_{j,n} := \gamma_n \in (0, 1]$

→ 2 possible algorithms :

- RGD : updates derived from Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_\alpha(\mu k; p)$
- MG : maximisation approach without $\lambda_{j,n}$ as a factor

→ 2 possible 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 .
- ❷ For all $j = 1 \dots J$, set:

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

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

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

→ Here $\hat{\varphi}_{j,n}^{(\alpha)}(y) = \frac{\varphi_{j,n}^{(\alpha)}(y)}{q_n(y)}$, $\gamma_{j,n} := \gamma_n \in (0, 1]$

→ 2 possible algorithms :

- RGD : updates derived from Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_{\alpha}(\mu k; p)$
- MG : maximisation approach without $\lambda_{j,n}$ as a factor

→ 2 possible 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 .
- ❷ For all $j = 1 \dots J$, set:

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

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

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

→ Here $\hat{\varphi}_{j,n}^{(\alpha)}(y) = \frac{\varphi_{j,n}^{(\alpha)}(y)}{q_n(y)}$, $\gamma_{j,n} := \gamma_n \in (0, 1]$

→ 2 possible algorithms :

- RGD : updates derived from Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_\alpha(\mu k; p)$
- MG : maximisation approach without $\lambda_{j,n}$ as a factor

→ 2 possible 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 .
- ❷ For all $j = 1 \dots J$, set:

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

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

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

→ Here $\hat{\varphi}_{j,n}^{(\alpha)}(y) = \frac{\varphi_{j,n}^{(\alpha)}(y)}{q_n(y)}$, $\gamma_{j,n} := \gamma_n \in (0, 1]$

→ 2 possible algorithms :

- RGD : updates derived from Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_\alpha(\mu k; p)$
- MG : maximisation approach without $\lambda_{j,n}$ as a factor

→ 2 possible 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 .
- ❷ For all $j = 1 \dots J$, set:

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

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

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

→ Here $\hat{\varphi}_{j,n}^{(\alpha)}(y) = \frac{\varphi_{j,n}^{(\alpha)}(y)}{q_n(y)}$, $\gamma_{j,n} := \gamma_n \in (0, 1]$

→ 2 possible algorithms :

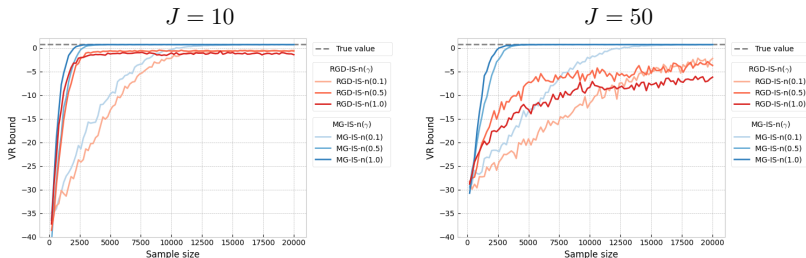
- RGD : updates derived from Gradient Descent steps w.r.t Θ on $-\alpha^{-1} \mathcal{L}_\alpha(\mu k; p)$
- MG : maximisation approach without $\lambda_{j,n}$ as a factor

→ 2 possible 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 $q_n = \mu_n k$.]



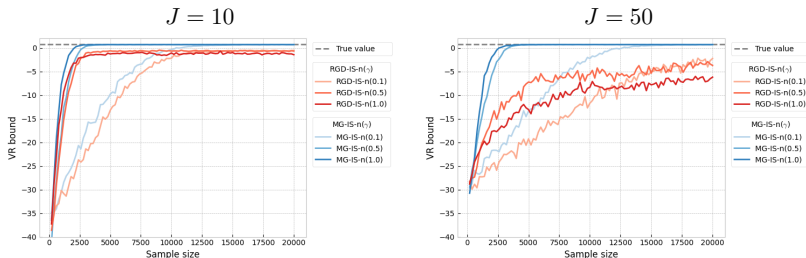
- LogMSE averaged over 30 trials for RGD and MG.

	$J = 10$			$J = 50$		
	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$
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 $q_n = \mu_n k$.]



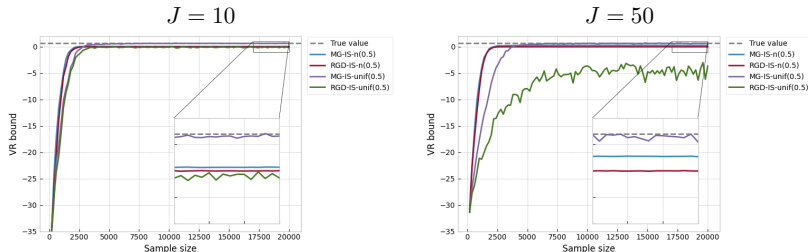
- LogMSE averaged over 30 trials for RGD and MG.

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

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, $\alpha = 0.2$, $d = 16$, $M = 200$, $\eta = 0.1$, $\kappa_n = 0$.]



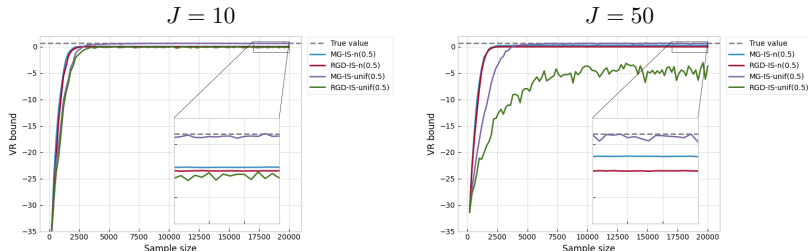
- LogMSE averaged over 30 trials for RGD and MG.

	$J = 10$			$J = 50$		
	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 1.0$
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)

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

- 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) - 2

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

- 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

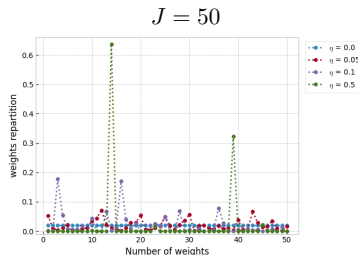
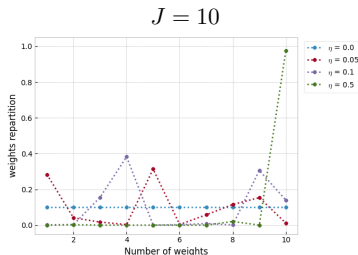
Comparing RGD to MG (varying λ) - 2

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

- 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 Monotonic Alpha-Divergence Minimisation
- 2 Maximisation approach
- 3 Gradient-based approach
- 4 Numerical Experiments
- 5 Conclusion of Part 3

Conclusion of Part 3

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

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **links** with an Integrated EM algorithm and with gradient-based approaches
- **empirical benefits** of our general framework

Some perspectives

- Additional convergence results
- Hyperparameters tuning
- ML applications...

Conclusion of Part 3

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

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **links** with an Integrated EM algorithm and with gradient-based approaches
- **empirical benefits** of our general framework

Some perspectives

- Additional convergence results
- Hyperparameters tuning
- ML applications...

Conclusion of Part 3

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

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **links** with an Integrated EM algorithm and with gradient-based approaches
- **empirical benefits** of our general framework

Some perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications...

Conclusion of Part 3

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

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **links** with an Integrated EM algorithm and with gradient-based approaches
- **empirical benefits** of our general framework

Some perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications...

Conclusion of Part 3

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

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **links** with an Integrated EM algorithm and with gradient-based approaches
- **empirical benefits** of our general framework

Some perspectives

- Additional convergence results
- Hyperparameters tuning
- ML applications...

Conclusion of Part 3

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

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **links** with an Integrated EM algorithm and with gradient-based approaches
- **empirical benefits** of our general framework

Some perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications...

Conclusion of Part 3

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

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **links** with an Integrated EM algorithm and with gradient-based approaches
- **empirical benefits** of our general framework

Some perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications...

Conclusion of Part 3

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

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **links** with an Integrated EM algorithm and with gradient-based approaches
- **empirical benefits** of our general framework

Some perspectives

- Additionnal convergence results
- Hyperparameters tuning
- ML applications...

Conclusion of Part 3

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

- applicable to **mixture models** optimisation
- mixture weights and mixture components parameters can be updated **simultaneously**
- **links** with an Integrated EM algorithm and with gradient-based approaches
- **empirical benefits** of our general framework

Some perspectives

- Additional convergence results
- Hyperparameters tuning
- ML applications...

Overall conclusion

- **Part 1.** *General introduction to Variational Inference :*
MFVI, BBVI, Alpha-divergence VI.
- **Part 2.** *Infinite-dimensional Alpha-divergence minimisation :*
Mixture weights optimisation to **increase expressiveness**
Links with the Entropic Mirror Descent algorithm
- **Part 3.** *Monotonic Alpha-divergence minimisation :*
Mixture models optimisation with **convergence guarantees**
Links with an Integrated EM algorithm and gradient-based approaches.

There is still a lot to do in Variational Inference regarding :

- ① the expressiveness of the variational family
- ② the choice of the measure of dissimilarity
- ③ the theory of Variational Inference
- ④ the interface between Variational Inference and Monte Carlo methods
- ⑤ and so much more!

Overall conclusion

- **Part 1.** *General introduction to Variational Inference :*
MFVI, BBVI, Alpha-divergence VI.
- **Part 2.** *Infinite-dimensional Alpha-divergence minimisation :*
Mixture weights optimisation to **increase expressiveness**
Links with the Entropic Mirror Descent algorithm
- **Part 3.** *Monotonic Alpha-divergence minimisation :*
Mixture models optimisation with **convergence guarantees**
Links with an Integrated EM algorithm and gradient-based approaches.

There is still a lot to do in Variational Inference regarding :

- ① the expressiveness of the variational family
- ② the choice of the measure of dissimilarity
- ③ the theory of Variational Inference
- ④ the interface between Variational Inference and Monte Carlo methods
- ⑤ and so much more!

Overall conclusion

- **Part 1.** *General introduction to Variational Inference* :
MFVI, BBVI, Alpha-divergence VI.
- **Part 2.** *Infinite-dimensional Alpha-divergence minimisation* :
Mixture weights optimisation to **increase expressiveness**
Links with the Entropic Mirror Descent algorithm
- **Part 3.** *Monotonic Alpha-divergence minimisation* :
Mixture models optimisation with **convergence guarantees**
Links with an Integrated EM algorithm and gradient-based approaches.

There is still a lot to do in Variational Inference regarding :

- ① the expressiveness of the variational family
- ② the choice of the measure of dissimilarity
- ③ the theory of Variational Inference
- ④ the interface between Variational Inference and Monte Carlo methods
- ⑤ and so much more!

Overall conclusion

- **Part 1.** *General introduction to Variational Inference* :
MFVI, BBVI, Alpha-divergence VI.
- **Part 2.** *Infinite-dimensional Alpha-divergence minimisation* :
Mixture weights optimisation to **increase expressiveness**
Links with the Entropic Mirror Descent algorithm
- **Part 3.** *Monotonic Alpha-divergence minimisation* :
Mixture models optimisation with **convergence guarantees**
Links with an Integrated EM algorithm and gradient-based approaches.

There is still a lot to do in Variational Inference regarding :

- ① the expressiveness of the variational family
- ② the choice of the measure of dissimilarity
- ③ the theory of Variational Inference
- ④ the interface between Variational Inference and Monte Carlo methods
- ⑤ and so much more!

Overall conclusion

- **Part 1.** *General introduction to Variational Inference* :
MFVI, BBVI, Alpha-divergence VI.
- **Part 2.** *Infinite-dimensional Alpha-divergence minimisation* :
Mixture weights optimisation to **increase expressiveness**
Links with the Entropic Mirror Descent algorithm
- **Part 3.** *Monotonic Alpha-divergence minimisation* :
Mixture models optimisation with **convergence guarantees**
Links with an Integrated EM algorithm and gradient-based approaches.

There is still a lot to do in Variational Inference regarding :

- ❶ the expressiveness of the variational family
- ❷ the choice of the measure of dissimilarity
- ❸ the theory of Variational Inference
- ❹ the interface between Variational Inference and Monte Carlo methods
- ❺ and so much more!

Overall conclusion

- **Part 1.** *General introduction to Variational Inference* :
MFVI, BBVI, Alpha-divergence VI.
- **Part 2.** *Infinite-dimensional Alpha-divergence minimisation* :
Mixture weights optimisation to **increase expressiveness**
Links with the Entropic Mirror Descent algorithm
- **Part 3.** *Monotonic Alpha-divergence minimisation* :
Mixture models optimisation with **convergence guarantees**
Links with an Integrated EM algorithm and gradient-based approaches.

There is still a lot to do in Variational Inference regarding :

- ❶ the expressiveness of the variational family
- ❷ the choice of the measure of dissimilarity
- ❸ the theory of Variational Inference
- ❹ the interface between Variational Inference and Monte Carlo methods
- ❺ and so much more!

Overall conclusion

- **Part 1.** *General introduction to Variational Inference* :
MFVI, BBVI, Alpha-divergence VI.
- **Part 2.** *Infinite-dimensional Alpha-divergence minimisation* :
Mixture weights optimisation to **increase expressiveness**
Links with the Entropic Mirror Descent algorithm
- **Part 3.** *Monotonic Alpha-divergence minimisation* :
Mixture models optimisation with **convergence guarantees**
Links with an Integrated EM algorithm and gradient-based approaches.

There is still a lot to do in Variational Inference regarding :

- ❶ the expressiveness of the variational family
- ❷ the choice of the measure of dissimilarity
- ❸ the theory of Variational Inference
- ❹ the interface between Variational Inference and Monte Carlo methods
- ❺ and so much more!

Overall conclusion

- **Part 1.** *General introduction to Variational Inference* :
MFVI, BBVI, Alpha-divergence VI.
- **Part 2.** *Infinite-dimensional Alpha-divergence minimisation* :
Mixture weights optimisation to **increase expressiveness**
Links with the Entropic Mirror Descent algorithm
- **Part 3.** *Monotonic Alpha-divergence minimisation* :
Mixture models optimisation with **convergence guarantees**
Links with an Integrated EM algorithm and gradient-based approaches.

There is still a lot to do in Variational Inference regarding :

- ① the expressiveness of the variational family
- ② the choice of the measure of dissimilarity
- ③ the theory of Variational Inference
- ④ the interface between Variational Inference and Monte Carlo methods
- ⑤ and so much more!

Overall conclusion

- **Part 1.** *General introduction to Variational Inference* :
MFVI, BBVI, Alpha-divergence VI.
- **Part 2.** *Infinite-dimensional Alpha-divergence minimisation* :
Mixture weights optimisation to **increase expressiveness**
Links with the Entropic Mirror Descent algorithm
- **Part 3.** *Monotonic Alpha-divergence minimisation* :
Mixture models optimisation with **convergence guarantees**
Links with an Integrated EM algorithm and gradient-based approaches.

There is still a lot to do in Variational Inference regarding :

- ① the expressiveness of the variational family
- ② the choice of the measure of dissimilarity
- ③ the theory of Variational Inference
- ④ the interface between Variational Inference and Monte Carlo methods
- ⑤ and so much more!