
Mixture weights optimisation for Alpha-Divergence Variational Inference

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This paper focuses on α -divergence minimisation methods for Variational Inference.
2 More precisely, we are interested in algorithms optimising the mixture weights of
3 any given mixture model, without any information on the underlying distribution
4 of its mixture components parameters. The Power Descent, defined for all $\alpha \neq 1$,
5 is one such algorithm and we establish in our work the full proof of its convergence
6 towards the optimal mixture weights when $\alpha < 1$. Since the α -divergence recovers
7 the widely-used forward Kullback-Leibler when $\alpha \rightarrow 1$, we then extend the Power
8 Descent to the case $\alpha = 1$ and show that we obtain an Entropic Mirror Descent.
9 This leads us to investigate the link between Power Descent and Entropic Mirror
10 Descent: first-order approximations allow us to introduce the Renyi Descent, a
11 novel algorithm for which we prove an $O(1/N)$ convergence rate. Lastly, we
12 compare numerically the behavior of the unbiased Power Descent and of the biased
13 Renyi Descent and we discuss the potential advantages of one algorithm over the
14 other.

15 1 Introduction

16 Bayesian Inference involves being able to compute or sample from the posterior density. For many
17 useful models, the posterior density can only be evaluated up to a normalisation constant and we
18 must resort to approximation methods.

19 One major category of approximation methods is Variational Inference, a wide class of optimisation
20 methods which introduce a simpler density family \mathcal{Q} and use it to approximate the posterior density
21 (see for example Variational Bayes [1, 2] and Stochastic Variational Inference [3]). The crux of
22 these methods consists in being able to find the best approximation of the posterior density among
23 the family \mathcal{Q} in the sense of a certain divergence, most typically the Kullback-Leibler divergence.
24 However, The Kullback-Leibler divergence is known to have some undesirable properties (e.g
25 posterior overestimation/underestimation [4]) and as a consequence, the α -divergence [5, 6] and
26 Renyi's α -divergence [7, 8] have gained a lot of attention recently as a more general alternative
27 [9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19].

28 Noticeably, [17] introduced the (α, Γ) -descent, a general family of gradient-based algorithms that
29 are able to optimise the *mixture weights* of mixture models by α -divergence minimisation, without
30 any information on the underlying distribution of its mixture components parameters. The benefit
31 of these types of algorithms is that they allow, in an Sequential Monte Carlo fashion [20], to select
32 the mixture components according to their overall importance in the set of component parameters.
33 From there, one is able to optimise the weights and the components parameters alternatively [17].
34 The (α, Γ) -descent framework recovers the Entropic Mirror Descent algorithm (corresponding to
35 $\Gamma(v) = e^{-\eta v}$ with $\eta > 0$) and includes the Power Descent, an algorithm defined for all $\alpha \in \mathbb{R} \setminus \{1\}$
36 and all $\eta > 0$ that sets $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$. Although these two algorithms are linked to one

37 another from a theoretical perspective through the (α, Γ) -descent framework, numerical experiments
 38 in [17] showed that the Power Descent outperforms the Entropic Mirror Descent when $\alpha < 1$ as the
 39 dimension increases.

40 Nonetheless, the global convergence of the Power Descent algorithm when $\alpha < 1$, as stated in [17],
 41 is subjected to the condition that the limit exists. Furthermore, even though the convergence towards
 42 the global optimum is derived, there is no convergence rate available for the Power Descent when
 43 $\alpha < 1$. While there is no general rule yet on how to select the value of α in practice, the case $\alpha < 1$
 44 has the advantage that it enforces a *mass-covering* property, as opposed to the *mode-seeking* property
 45 exhibited when $\alpha \geq 1$ ([4] and [17]) and which often may lead to posterior variance underestimation.
 46 We are thus interested in studying Variational Inference methods for optimising the mixture weights of
 47 mixture models when $\alpha < 1$. After recalling the basics of the Power Descent algorithm in Section 2,
 48 we make the following contributions in the paper:

- 49 • In Section 3, we derive the full convergence proof of the Power Descent algorithm towards the
 50 optimal mixture weights when $\alpha < 1$ (Theorem 2).
- 51 • Since the α -divergence becomes the traditional forward Kullback-Leibler when $\alpha \rightarrow 1$, we first
 52 bridge in Section 4 the gap between the cases $\alpha < 1$ and $\alpha > 1$ of the Power Descent: we obtain
 53 that the Power Descent recovers an Entropic Mirror Descent performing forward Kullback-Leibler
 54 minimisation (Proposition 1). We then keep on investigating the connections between the Power
 55 Descent and the Entropic Mirror Descent by considering first-order approximations. In doing so, we
 56 are able to go beyond the (α, Γ) -descent framework and to introduce an algorithm closely-related to
 57 the Power Descent that we call the *Renyi Descent* and that is proved in Theorem 3 to converge at an
 58 $O(1/N)$ rate towards its optimum for all $\alpha \in \mathbb{R}$.
- 59 • Finally, we run some numerical experiments in Section 5 to compare the behavior of the Power
 60 Descent and the Renyi Descent altogether, before discussing the potential benefits of one approach
 61 over the other.

62 2 Background

We start by introducing some notation. Let (Y, \mathcal{Y}, ν) be a measured space, where ν is a σ -finite
 measure on (Y, \mathcal{Y}) . Assume that we have access to some observed variables \mathcal{D} generated from a
 probabilistic model $p(\mathcal{D}|y)$ parameterised by a hidden random variable $y \in Y$ that is drawn from a
 certain prior $p_0(y)$. The posterior density of the latent variable y given the data \mathcal{D} is then given by:

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

63 where the normalisation constant $p(\mathcal{D}) = \int_Y p_0(y)p(\mathcal{D}|y)\nu(dy)$ is called the *marginal likelihood* or
 64 *model evidence* and is oftentimes unknown.

To approximate the posterior density, the Power Descent considers a variational family \mathcal{Q} that is large
 enough to contain mixture models and that we redefine now: letting (T, \mathcal{T}) be a measurable space,
 $K : (\theta, A) \mapsto \int_A k(\theta, y)\nu(dy)$ be a Markov transition kernel on $T \times \mathcal{Y}$ with kernel density k defined
 on $T \times Y$, the Power Descent considers the following approximating family

$$\left\{ y \mapsto \int_T \mu(d\theta)k(\theta, y) : \mu \in \mathbf{M} \right\},$$

65 where \mathbf{M} is a convenient subset of $\mathbf{M}_1(T)$, the set of probability measures on (T, \mathcal{T}) . This choice of
 66 approximating family extends the typical parametric family commonly-used in Variational Inference
 67 since it amounts to putting a prior over the parameter θ (in the form of a measure) and does describe
 68 the class of mixture models when μ is a weighted sum of Dirac measures.

69 **Problem statement** Denote by \mathbb{P} the probability measure on (Y, \mathcal{Y}) with corresponding density
 70 $p(\cdot|\mathcal{D})$ with respect to ν and for all $\mu \in \mathbf{M}_1(T)$, for all $y \in Y$, denote $\mu k(y) = \int_T \mu(d\theta)k(\theta, y)$.
 71 Furthermore, given $\alpha \in \mathbb{R}$, let f_α be the convex function on $(0, +\infty)$ defined by $f_0(u) = u - 1 -$
 72 $\log(u)$, $f_1(u) = 1 - u + u \log(u)$ and $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$ for all $\alpha \in \mathbb{R} \setminus \{0, 1\}$.
 73 Then, the α -divergence between μK and \mathbb{P} (extended by continuity to the cases $\alpha = 0$ and $\alpha = 1$ as

74 for example done in [21]) is given by

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

75 and the goal of the Power Descent is to find

$$\operatorname{arginf}_{\mu \in \mathcal{M}} D_\alpha(\mu K || \mathbb{P}) . \quad (1)$$

76 More generally, letting p be any measurable positive function on (Y, \mathcal{Y}) , the Power Descent aims at
77 solving

$$\operatorname{arginf}_{\mu \in \mathcal{M}} \Psi_\alpha(\mu; p) , \quad (2)$$

78 where for all $\mu \in \mathcal{M}_1(\mathbb{T})$, $\Psi_\alpha(\mu; p) = \int_Y f_\alpha(\mu k(y)/p(y)) p(y) \nu(dy)$. The Variational Inference
79 optimisation problem (1) can then be seen as an instance of (2) that is equivalent to optimising
80 $\Psi_\alpha(\mu; p)$ with $p(y) = p(y, \mathcal{D})$ (see Appendix A.1). In the following, the dependency on p in Ψ_α
81 may be dropped throughout the paper for notational ease when no ambiguity occurs and we now
82 present the Power Descent algorithm.

83 **The Power Descent algorithm.** The optimisation problem (2) can be solved for all $\alpha \in \mathbb{R} \setminus \{1\}$ by
84 using the Power Descent algorithm introduced in [17] : given an initial measure $\mu_1 \in \mathcal{M}_1(\mathbb{T})$ such
85 that $\Psi_\alpha(\mu_1) < \infty$, $\alpha \in \mathbb{R} \setminus \{1\}$, $\eta > 0$ and κ such that $(\alpha - 1)\kappa \geq 0$, the Power descent algorithm
86 is an iterative scheme which builds the sequence of probability measures $(\mu_n)_{n \in \mathbb{N}^*}$

$$\mu_{n+1} = \mathcal{I}_\alpha(\mu_n) , \quad n \in \mathbb{N}^* , \quad (3)$$

87 where for all $\mu \in \mathcal{M}_1(\mathbb{T})$, the one-step transition $\mu \mapsto \mathcal{I}_\alpha(\mu)$ is given by Algorithm 1 and where for
88 all $v \in \operatorname{Dom}_\alpha$, $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$ [and $\operatorname{Dom}_\alpha$ denotes an interval of \mathbb{R} such that for all
89 $\theta \in \mathbb{T}$, all $\mu \in \mathcal{M}_1(\mathbb{T})$, $b_{\mu, \alpha}(\theta) + \kappa$ and $\mu(b_{\mu, \alpha}) + \kappa \in \operatorname{Dom}_\alpha$].

Algorithm 1: Power descent one-step transition ($\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$)

- 90 1. Expectation step : $b_{\mu, \alpha}(\theta) = \int_Y k(\theta, y) f'_\alpha \left(\frac{\mu k(y)}{p(y)} \right) \nu(dy)$
91 2. Iteration step : $\mathcal{I}_\alpha(\mu)(d\theta) = \frac{\mu(d\theta) \cdot \Gamma(b_{\mu, \alpha}(\theta) + \kappa)}{\mu(\Gamma(b_{\mu, \alpha} + \kappa))}$
-

91 In this algorithm, $b_{\mu, \alpha}$ can be understood as the gradient of Ψ_α . Algorithm 1 then consists in applying
92 the transform function Γ to the translated gradient $b_{\mu, \alpha} + \kappa$ and projecting back onto the space of
93 probability measures.

94 A remarkable property of the Power Descent algorithm, which has been proven in [17] (it is a special
95 case of [17, Theorem 1] with $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$), is that under (A1) as defined below

96 (A1) The density kernel k on $\mathbb{T} \times Y$, the function p on Y and the σ -finite measure ν on
97 (Y, \mathcal{Y}) satisfy, for all $(\theta, y) \in \mathbb{T} \times Y$, $k(\theta, y) > 0$, $p(y) > 0$ and $\int_Y p(y) \nu(dy) < \infty$.

98 the Power Descent ensures a monotonic decrease in the α -divergence at each step for all $\eta \in (0, 1]$
99 (this result is recalled in Theorem 4 of Appendix A.2 for the sake of completeness). Under the
100 additional assumptions that $\kappa > 0$ and

$$\sup_{\theta \in \mathbb{T}, \mu \in \mathcal{M}_1(\mathbb{T})} |b_{\mu, \alpha}| < \infty \quad \text{and} \quad \Psi_\alpha(\mu_1) < \infty , \quad (4)$$

101 the Power Descent is also known to converge towards its optimal value at an $O(1/N)$ rate when
102 $\alpha > 1$ [17, Theorem 3]. On the other hand, when $\alpha < 1$, the convergence towards the optimum as
103 written in [17] holds under different assumptions including

- 104 (A2) (i) \mathbb{T} is a compact metric space and \mathcal{T} is the associated Borel σ -field;
105 (ii) for all $y \in Y$, $\theta \mapsto k(\theta, y)$ is continuous;

106 (iii) we have $\int_{\mathcal{Y}} \sup_{\theta \in \mathcal{T}} k(\theta, y) \times \sup_{\theta' \in \mathcal{T}} \left(\frac{k(\theta', y)}{p(y)} \right)^{\alpha-1} \nu(dy) < \infty$.

107 If $\alpha = 0$, assume in addition that $\int_{\mathcal{Y}} \sup_{\theta \in \mathcal{T}} \left| \log \left(\frac{k(\theta, y)}{p(y)} \right) \right| p(y) \nu(dy) < \infty$.

108 so that [17, Theorem 4], that is recalled below under the form of Theorem 1, states the convergence
109 of the Power Descent algorithm towards the global optimum.

110 **Theorem 1** ([17, Theorem 4]). *Assume (A1) and (A2). Let $\alpha < 1$ and let $\kappa \leq 0$. Then, for all
111 $\mu \in \mathcal{M}_1(\mathcal{T})$, $\Psi_\alpha(\mu) < \infty$ and any $\eta > 0$ satisfies $0 < \mu(\Gamma(b_{\mu, \alpha} + \kappa)) < \infty$. Further assume
112 that $\eta \in (0, 1]$ and that there exist $\mu_1, \mu^* \in \mathcal{M}_1(\mathcal{T})$ such that the (well-defined) sequence $(\mu_n)_{n \in \mathbb{N}^*}$
113 defined by (3) weakly converges to μ^* as $n \rightarrow \infty$. Finally, denote by $\mathcal{M}_{1, \mu_1}(\mathcal{T})$ the set of probability
114 measures dominated by μ_1 . Then the following assertions hold*

115 (i) $(\Psi_\alpha(\mu_n))_{n \in \mathbb{N}^*}$ is nonincreasing,

116 (ii) μ^* is a fixed point of \mathcal{I}_α ,

117 (iii) $\Psi_\alpha(\mu^*) = \inf_{\zeta \in \mathcal{M}_{1, \mu_1}(\mathcal{T})} \Psi_\alpha(\zeta)$.

118 The above result assumes there must exist $\mu_1, \mu^* \in \mathcal{M}_1(\mathcal{T})$ such that the sequence $(\mu_n)_{n \in \mathbb{N}^*}$ defined
119 by (3) weakly converges to μ^* as $n \rightarrow \infty$, that is it assumes the limit already exists. Our first
120 contribution consists in showing that this assumption can be alleviated when μ is chosen a weighted
121 sum of Dirac measures, that is when we seek to perform mixture weights optimisation by α -divergence
122 minimisation.

123 3 Convergence of the Power Descent algorithm in the mixture case

124 Before we state our convergence result, let us first make two comments on the assumptions from
125 Theorem 1 that shall be retained in our upcoming convergence result.

126 A first comment is that (A1) is mild since the assumption that $p(y) > 0$ for all $y \in \mathcal{Y}$ can be discarded
127 and is kept for convenience [17, Remark 4]. A second comment is that (A2) is also mild and covers
128 (4) as it amounts to assuming that $b_{\mu, \alpha}(\theta)$ and $\Psi_\alpha(\mu)$ are uniformly bounded with respect to μ and θ .
129 To see this, we give below an example for which (A2) is satisfied.

130 **Example 1.** *Consider the case $\mathcal{Y} = \mathbb{R}^d$ with $\alpha \in [0, 1)$. Let $r > 0$ and let $\mathcal{T} = \mathcal{B}(0, r) \subset \mathbb{R}^d$.
131 Furthermore, let K_h be a Gaussian transition kernel with bandwidth h and denote by k_h its associ-
132 ated kernel density. Finally, let p be a mixture density of two d -dimensional Gaussian distributions
133 multiplied by a positive constant c such that $p(y) = c \times [0.5\mathcal{N}(y; \theta_1^*, \mathbf{I}_d) + 0.5\mathcal{N}(y; \theta_2^*, \mathbf{I}_d)]$ for all
134 $y \in \mathcal{Y}$ where $\theta_1^*, \theta_2^* \in \mathcal{T}$ and \mathbf{I}_d is the identity matrix. Then, (A2) holds (see Appendix B.1).*

Next, we introduce some notation that are specific to the case of mixture models we aim at studying
in this section. Given $J \in \mathbb{N}^*$, we introduce the simplex of \mathbb{R}^J :

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

135 and we also define $\mathcal{S}_J^+ = \{\boldsymbol{\lambda} \in \mathcal{S}_J : \forall j \in \{1, \dots, J\}, \lambda_j > 0\}$. We let $\Theta = (\theta_1, \dots, \theta_J) \in \mathcal{T}^J$
136 be fixed and for all $\boldsymbol{\lambda} \in \mathcal{S}_J$, we define $\mu_{\boldsymbol{\lambda}, \Theta} \in \mathcal{M}_1(\mathcal{T})$ by $\mu_{\boldsymbol{\lambda}, \Theta} = \sum_{j=1}^J \lambda_j \delta_{\theta_j}$.

137 Consequently, $\mu_{\boldsymbol{\lambda}, \Theta} k(y) = \sum_{j=1}^J \lambda_j k(\theta_j, y)$ corresponds to a mixture model and if we let $(\mu_n)_{n \in \mathbb{N}^*}$
138 be defined by $\mu_1 = \mu_{\boldsymbol{\lambda}, \Theta}$ and (3), an immediate induction yields that for every $n \in \mathbb{N}^*$, μ_n can
139 be expressed as $\mu_n = \sum_{j=1}^J \lambda_{j,n} \delta_{\theta_j}$ where $\boldsymbol{\lambda}_n = (\lambda_{1,n}, \dots, \lambda_{J,n}) \in \mathcal{S}_J$ satisfies the initialisation
140 $\boldsymbol{\lambda}_1 = \boldsymbol{\lambda}$ and the update formula:

$$\boldsymbol{\lambda}_{n+1} = \mathcal{I}_\alpha^{\text{mixt}}(\boldsymbol{\lambda}_n), \quad n \in \mathbb{N}^*, \quad (5)$$

where for all $\boldsymbol{\lambda} \in \mathcal{S}_J$,

$$\mathcal{I}_\alpha^{\text{mixt}}(\boldsymbol{\lambda}) := \left(\frac{\lambda_j \Gamma(b_{\mu_{\boldsymbol{\lambda}, \Theta}, \alpha}(\theta_j) + \kappa)}{\sum_{\ell=1}^J \lambda_\ell \Gamma(b_{\mu_{\boldsymbol{\lambda}, \Theta}, \alpha}(\theta_\ell) + \kappa)} \right)_{1 \leq j \leq J}$$

141 with $\Gamma(v) = [(\alpha - 1)v + 1]^{\frac{\eta}{1-\alpha}}$ for all $v \in \text{Dom}_\alpha$. Finally, let us rewrite (A2) in the simplified case
 142 where μ is a sum of Dirac measures, which gives (A3) below.

143 (A3) (i) For all $y \in Y$, $\theta \mapsto k(\theta, y)$ is continuous;

144 (ii) we have $\int_Y \max_{1 \leq j \leq J} k(\theta_j, y) \times \max_{1 \leq j' \leq J} \left(\frac{k(\theta_{j'}, y)}{p(y)} \right)^{\alpha-1} \nu(dy) < \infty$.

145 If $\alpha = 0$, we assume in addition that $\int_Y \max_{1 \leq j \leq J} \left| \log \left(\frac{k(\theta_j, y)}{p(y)} \right) \right| p(y) \nu(dy) < \infty$.

146 We then have the following theorem, which establishes the full proof of the global convergence
 147 towards the optimum for the mixture weights when $\alpha < 1$.

148 **Theorem 2.** Assume (A1) and (A3). Let $\alpha < 1$, let $\Theta = (\theta_1, \dots, \theta_J) \in \mathbb{T}^J$ be fixed and let κ be such
 149 that $\kappa \leq 0$. Then for all $\lambda \in \mathcal{S}_J$, $\Psi_\alpha(\mu_{\lambda, \Theta}) < \infty$ and for any $\eta > 0$ the sequence $(\lambda_n)_{n \in \mathbb{N}^*}$ defined
 150 by $\lambda_1 \in \mathcal{S}_J$ and (5) is well-defined. If in addition $(\lambda_1, \eta) \in \mathcal{S}_J^+ \times (0, 1]$ and $\{K(\theta_1, \cdot), \dots, K(\theta_J, \cdot)\}$
 151 are linearly independent, then

152 (i) $(\Psi_\alpha(\mu_{\lambda_n, \Theta}))_{n \in \mathbb{N}^*}$ is nonincreasing,

153 (ii) the sequence $(\lambda_n)_{n \in \mathbb{N}^*}$ converges to some $\lambda_* \in \mathcal{S}_J$ which is a fixed point of $\mathcal{I}_\alpha^{\text{mixt}}$,

154 (iii) $\Psi_\alpha(\mu_{\lambda_*, \Theta}) = \inf_{\lambda' \in \mathcal{S}_J} \Psi_\alpha(\mu_{\lambda', \Theta})$.

155 The proof of this result builds on Theorem 1 and Theorem 4 and is deferred to Appendix B.2. Notice
 156 that since Ψ_α depends on λ through $\mu_{\lambda, \Theta} K$ in Theorem 2, an identifiability condition was to be
 157 expected in order to achieve the convergence of the sequence $(\lambda_n)_{n \in \mathbb{N}^*}$. Following Example 1, this
 158 identifiability condition notably holds for $J \leq d$ under the assumption that the $\theta_1, \dots, \theta_J$ are full-rank.

159 We thus have the convergence of the Power Descent under less stringent conditions when $\alpha < 1$
 160 and when we consider the particular case of mixture models. This algorithm can easily become
 161 feasible for any choice of kernel K by resorting to an unbiased estimator of $(b_{\mu_{\lambda_n, \Theta}, \alpha}(\theta_j))_{1 \leq j \leq J}$ in
 162 the update formula (5) (see Algorithm 3 of Appendix B.3).

163 Nevertheless, contrary to the case $\alpha > 1$ we still do not have a convergence rate for the Power Descent
 164 when $\alpha < 1$. Furthermore, the important case $\alpha \rightarrow 1$, which corresponds to performing forward
 165 Kullback-Leibler minimisation, is not covered by the Power Descent algorithm. In the next section,
 166 we extend the Power Descent to the case $\alpha = 1$. As we shall see, this will lead us to investigate the
 167 connections between the Power Descent and the Entropic Mirror Descent beyond the (α, Γ) -descent
 168 framework. As a result, we will introduce a novel algorithm closely-related to the Power Descent
 169 that yields an $O(1/N)$ convergence rate when $\mu = \mu_{\lambda, \Theta}$ and $\alpha < 1$ (and more generally when
 170 $\mu \in \mathbb{M}_1(\mathbb{T})$ and $\alpha \in \mathbb{R}$).

171 4 Power Descent and Entropic Mirror Descent

172 Recall from Section 2 that the Power Descent is defined for all $\alpha \in \mathbb{R} \setminus \{1\}$. In this section, we
 173 first establish in Proposition 1 that the Power Descent can be extended to the case $\alpha = 1$ and
 174 that we recover an Entropic Mirror Descent, showing that a deeper connection runs between the
 175 two approaches beyond the one identified by the (α, Γ) -descent framework. This result relies on
 176 typical convergence and differentiability assumptions summarised in (D1) and which are deferred to
 177 Appendix C.1, alongside with the proof of Proposition 1.

Proposition 1 (Limiting case $\alpha \rightarrow 1$). Assume (A1) and (D1). Then, for all continuous and bounded
 real-valued functions h on \mathbb{T} , we have that

$$\lim_{\alpha \rightarrow 1} [\mathcal{I}_\alpha(\mu)](h) = [\mathcal{I}_1(\mu)](h),$$

178 where for all $\mu \in \mathbb{M}_1(\mathbb{T})$ and all $\theta \in \mathbb{T}$, we have set

$$\mathcal{I}_1(\mu)(d\theta) = \frac{\mu(d\theta) e^{-\eta b_{\mu, 1}(\theta)}}{\mu(e^{-\eta b_{\mu, 1}})} \quad \text{and} \quad b_{\mu, 1}(\theta) = \int_Y k(\theta, y) \log \left(\frac{\mu k(y)}{p(y)} \right) \nu(dy). \quad (6)$$

179 Here, we recognise the one-step transition associated to the Entropic Mirror Descent applied to Ψ_1 .
 180 This algorithm is a special case of [17] with $\Gamma(v) = e^{-\eta v}$ and $\alpha = 1$ and as such, it is known to

181 lead to a systematic decrease in the forward Kullback-Leibler divergence and to enjoy an $O(1/N)$
 182 convergence rate under the assumptions that (4) holds and $\eta \in (0, 1)$ [17, Theorem 3].

183 We have thus obtained that the Power Descent coincides exactly with the Entropic Mirror Descent
 184 applied to Ψ_1 when $\alpha = 1$ and we now focus on understanding the links between Power Descent and
 185 Entropic Mirror Descent when $\alpha \in \mathbb{R} \setminus \{1\}$. For this purpose, let κ be such that $(\alpha - 1)\kappa \geq 0$ and
 186 let us study first-order approximations of the Power Descent and the Entropic Mirror Descent applied
 187 to Ψ_α when $b_{\mu_n, \alpha}(\theta) \approx \mu_n(b_{\mu_n, \alpha})$ for all $\theta \in \mathbb{T}$.

Letting $\eta > 0$, we have that the update formula for the Power Descent is given by

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) [(\alpha - 1)(b_{\mu_n, \alpha}(\theta) + \kappa) + 1]^{\frac{\eta}{1-\alpha}}}{\mu_n([(\alpha - 1)(b_{\mu_n, \alpha} + \kappa) + 1]^{\frac{\eta}{1-\alpha}})}, \quad n \in \mathbb{N}^*.$$

Now using the first order approximation $u^{\frac{\eta}{1-\alpha}} \approx v^{\frac{\eta}{1-\alpha}} - \frac{\eta}{\alpha-1} v^{\frac{\eta}{1-\alpha}-1}(u - v)$ with $u = \frac{(\alpha-1)(b_{\mu_n, \alpha}(\theta)+\kappa)+1}{(\alpha-1)(\mu(b_{\mu_n, \alpha})+\kappa)+1}$ and $v = 1$, we can deduce the following approximated update formula

$$\mu_{n+1}(d\theta) = \mu_n(d\theta) \left[1 - \frac{\eta}{\alpha - 1} \frac{b_{\mu_n, \alpha}(\theta) - \mu_n(b_{\mu_n, \alpha})}{\mu_n(b_{\mu_n, \alpha}) + \kappa + 1/(\alpha - 1)} \right], \quad n \in \mathbb{N}^*.$$

188 Letting $\eta' > 0$, the update formula for the Entropic Mirror Descent applied to Ψ_α can be written as

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \exp[-\eta'(b_{\mu_n, \alpha}(\theta) + \kappa)]}{\mu_n(\exp[-\eta'(b_{\mu_n, \alpha} + \kappa)])}, \quad n \in \mathbb{N}^*, \quad (7)$$

and we obtain in a similar fashion that an approximated version of this iterative scheme is

$$\mu_{n+1}(d\theta) = \mu_n(d\theta) [1 - \eta' (b_{\mu_n, \alpha}(\theta) - \mu_n(b_{\mu_n, \alpha}))], \quad n \in \mathbb{N}^*.$$

189 Thus, for the two approximated formulas above to coincide, we need to set $\eta' =$
 190 $\eta [(\alpha - 1)(\mu_n(b_{\mu_n, \alpha}) + \kappa) + 1]^{-1}$. Now coming back to (7), we see that this leads us to consider
 191 the update formula given by

$$\mu_{n+1}(d\theta) = \frac{\mu_n(d\theta) \exp\left[-\eta \frac{b_{\mu_n, \alpha}(\theta)}{(\alpha-1)(\mu_n(b_{\mu_n, \alpha})+\kappa)+1}\right]}{\mu_n\left(\exp\left[-\eta \frac{b_{\mu_n, \alpha}}{(\alpha-1)(\mu_n(b_{\mu_n, \alpha})+\kappa)+1}\right]\right)}, \quad n \in \mathbb{N}^*. \quad (8)$$

192 Observe then that (8) can again be seen as an Entropic Mirror Descent, but applied this time to the
 193 objective function defined for all $\alpha \in \mathbb{R} \setminus \{0, 1\}$ by

$$\Psi_\alpha^{AR}(\mu) := \frac{1}{\alpha(\alpha - 1)} \log \left(\int_{\mathcal{Y}} \mu k(y)^\alpha p(y)^{1-\alpha} \nu(dy) + (\alpha - 1)\kappa \right),$$

meaning we have applied the monotonic transformation

$$u \mapsto \frac{1}{\alpha(\alpha - 1)} \log \left(\alpha(\alpha - 1)u + \alpha + (1 - \alpha) \int_{\mathcal{Y}} p(y) \nu(dy) + (\alpha - 1)\kappa \right)$$

194 to the initial objective function Ψ_α (see Appendix C.2 for the derivation of (8) based on the objective
 195 function Ψ_α^{AR}). Hence, in the spirit of Renyi's α -divergence gradient-based methods for Variational
 196 Inference (e.g [9, 10]), we can motivate the iterative scheme (8) by observing that we recover the
 197 Variational Renyi bound introduced in [10] up to a constant $-\alpha^{-1}$ when we let $p = p(\cdot, \mathcal{D})$, $\kappa = 0$
 198 and $\alpha > 0$ in Ψ_α^{AR} . For this reason we call the algorithm given by (8) the *Renyi Descent* thereafter.

199 Contrary to the Entropic Mirror Descent applied to Ψ_α , the Renyi Descent now shares the same
 200 first-order approximation as the Power Descent. This might explain why the behavior of the Entropic
 201 Mirror Descent applied to Ψ_α and of the Power Descent differed greatly when $\alpha < 1$ in the numerical
 202 experiments from [17] despite their theoretical connection through the (α, Γ) -descent framework (the
 203 former performing poorly numerically compared to the later as the dimension increased).

Strikingly, we can prove an $O(1/N)$ convergence rate towards the global optimum for the Renyi
 Descent. Letting $\kappa' \in \mathbb{R}$, denoting by Dom_α^{AR} an interval of \mathbb{R} such that for all $\theta \in \mathbb{T}$ and all
 $\mu \in \mathcal{M}_1(\mathbb{T})$,

$$\frac{b_{\mu, \alpha}(\theta) + 1/(\alpha - 1)}{(\alpha - 1)(\mu(b_{\mu, \alpha}) + \kappa) + 1} + \kappa' \quad \text{and} \quad \frac{\mu(b_{\mu, \alpha}) + 1/(\alpha - 1)}{(\alpha - 1)(\mu(b_{\mu, \alpha}) + \kappa) + 1} + \kappa' \in \text{Dom}_\alpha^{AR}$$

204 and introducing the assumption on η

Table 1: Summary of the theoretical results obtained in this paper compared to [17]

	Power Descent	Renyi Descent
[17]	$\alpha < 1$: convergence under restrictive assumptions; $\alpha > 1$: $O(1/N)$ convergence rate	not covered
This paper	$\alpha < 1$: full proof of convergence for mixture weights; extension to $\alpha = 1$ with $O(1/N)$ convergence rate	$O(1/N)$ convergence rate

205 (A4) For all $v \in \text{Dom}_\alpha^{AR}$, $1 - \eta(\alpha - 1)(v - \kappa') \geq 0$.

206 we indeed have the following convergence result.

207 **Theorem 3.** Assume (A1) and (A4). Let $\alpha \in \mathbb{R} \setminus \{1\}$ and let κ be such that $(\alpha - 1)\kappa > 0$. Define
208 $|B|_{\infty, \alpha} := \sup_{\theta \in \mathbb{T}, \mu \in \mathbb{M}_1(\mathbb{T})} |b_{\mu, \alpha}(\theta) + 1/(\alpha - 1)|$ and assume that $|B|_{\infty, \alpha} < \infty$. Moreover, let
209 $\mu_1 \in \mathbb{M}_1(\mathbb{T})$ be such that $\Psi_\alpha(\mu_1) < \infty$. Then, the following assertions hold.

210 (i) The sequence $(\mu_n)_{n \in \mathbb{N}^*}$ defined by (8) is well-defined and the sequence $(\Psi_\alpha(\mu_n))_{n \in \mathbb{N}^*}$ is
211 non-increasing.

212 (ii) For all $N \in \mathbb{N}^*$, we have

$$\Psi_\alpha(\mu_N) - \Psi_\alpha(\mu^*) \leq \frac{L_{\alpha,2}}{N} \left[KL(\mu^* || \mu_1) + L \frac{L_{\alpha,3}}{L_{\alpha,1}(\alpha - 1)\kappa} \Delta_1 \right], \quad (9)$$

213 where μ^* is such that $\Psi_\alpha(\mu^*) = \inf_{\zeta \in \mathbb{M}_1, \mu_1(\mathbb{T})} \Psi_\alpha(\zeta)$, $\mathbb{M}_1, \mu_1(\mathbb{T})$ denotes the set of
214 probability measures dominated by μ_1 , $KL(\mu^* || \mu_1) = \int_{\mathbb{T}} \log(d\mu^*/d\mu_1) d\mu^*$, $\Delta_1 =$
215 $\Psi_\alpha(\mu_1) - \Psi_\alpha(\mu^*)$ and $L_{\alpha,2}$, L , $L_{\alpha,3}$, $L_{\alpha,1}$ are finite constants defined in (20).

216 The proof of this result is deferred to Appendix C.3 and we present in the next example an application
217 of this theorem to the particular case of mixture models.

Example 2. Let $\alpha \in \mathbb{R} \setminus \{1\}$, let $J \in \mathbb{N}^*$, let $\Theta = (\theta_1, \dots, \theta_J) \in \mathbb{T}^J$, let $\mu_1 = J^{-1} \sum_{j=1}^J \delta_{\theta_j}$ and let
218 $\text{Dom}_\alpha^{AR} = [-\frac{|B|_{\infty, \alpha}}{(\alpha-1)\kappa} + \kappa', \frac{|B|_{\infty, \alpha}}{(\alpha-1)\kappa} + \kappa']$ with $\kappa' \in \mathbb{R}$. In addition, assume that $1 - \eta|\kappa|^{-1}|B|_{\infty, \alpha} > 0$.
219 Then, taking $\kappa' = -3\frac{|B|_{\infty, \alpha}}{(\alpha-1)\kappa}$, we obtain

$$\Psi_\alpha(\mu_N) - \Psi_\alpha(\mu^*) \leq \frac{|\alpha - 1|(|B|_{\infty, \alpha} + |\kappa|)}{N} \left[\frac{\log J}{\eta} + \frac{\sqrt{2 \log(J)} |B|_{\infty, \alpha}}{(\alpha - 1)\kappa(1 - \eta|\kappa|^{-1}|B|_{\infty, \alpha})} \right],$$

218 where we have used that $KL(\mu^* || \mu_1) \leq \log J$, $\Delta_1 \leq \sqrt{2 \log J} |B|_{\infty, \alpha}$ and that the constants defined
219 in (20) satisfy $L_{\alpha,2} = \eta^{-1}|\alpha - 1|(|B|_{\infty, \alpha} + |\kappa|)$, $L = \eta^2 e^{\eta \frac{|B|_{\infty, \alpha}}{(\alpha-1)\kappa} - \eta\kappa'}$, $L_{\alpha,3} = e^{\eta \frac{|B|_{\infty, \alpha}}{(\alpha-1)\kappa} + \eta\kappa'}$ and
220 $L_{\alpha,1} = (1 - \eta|\kappa|^{-1}|B|_{\infty, \alpha})\eta e^{-\eta \frac{|B|_{\infty, \alpha}}{(\alpha-1)\kappa} - \eta\kappa'}$.

221 To put things into perspective, notice that the Renyi Descent enjoys an $O(1/\sqrt{N})$ convergence
222 rate as a Entropic Mirror Descent algorithm for the sequence $(\Psi_\alpha(N^{-1} \sum_{n=1}^N \mu_n))_{N \in \mathbb{N}^*}$ under our
223 assumptions when η is proportional to $1/\sqrt{N}$, N being fixed (see [22] or [23, Theorem 4.2.]).

224 The improvement thus lies in the fact that deriving an $O(1/N)$ convergence rate usually requires
225 stronger smoothness assumptions on Ψ_α [23, Theorem 6.2] that we do not assume in Theorem 3.
226 Furthermore, due to the monotonicity property, our result only involves the measure μ_N at time N
227 while typical Entropic Mirror Result are expressed in terms of the average $N^{-1} \sum_{n=1}^N \mu_n$.

228 Finally, observe that the Renyi Descent becomes feasible in practice for any choice of kernel K by
229 letting μ be a weighted sum of Dirac measures i.e $\mu = \mu_{\lambda, \Theta}$ and by resorting to an unbiased estimate
230 of $(b_{\mu, \alpha}(\theta_j))_{1 \leq j \leq J}$ (see Algorithm 4 of Appendix C.4).

231 The theoretical results we have obtained are summarised in Table 1 and we next move on to numerical
232 experiments.

233 **5 Simulation study**

Let the target p be a mixture density of two d -dimensional Gaussian distributions multiplied by a positive constant c such that $p(y) = c \times [0.5\mathcal{N}(y; -s\mathbf{u}_d, \mathbf{I}_d) + 0.5\mathcal{N}(y; s\mathbf{u}_d, \mathbf{I}_d)]$, where \mathbf{u}_d is the d -dimensional vector whose coordinates are all equal to 1, $s = 2$, $c = 2$ and \mathbf{I}_d is the identity matrix. Given $J \in \mathbb{N}^*$, the approximating family is described by

$$\left\{ y \mapsto \mu_\lambda k_h(y) = \sum_{j=1}^J \lambda_j k_h(y - \theta_j) : \lambda \in \mathcal{S}_J, \theta_1, \dots, \theta_J \in \mathbb{T} \right\},$$

234 where K_h is a Gaussian transition kernel with bandwidth h and k_h denotes its associated kernel
235 density.

236 Since the Power Descent and the Renyi Descent operate only on the mixture weights λ of $\mu_\lambda k_h$
237 during the optimisation, a fully adaptive algorithm can be obtained by alternating T times between
238 an *Exploitation step* where the mixture weights are optimised and an *Exploration step* where the
239 $\theta_1, \dots, \theta_J$ are updated, as written in Algorithm 2.

Algorithm 2: Complete Exploitation-Exploration Algorithm

Input: p : measurable positive function, α : α -divergence parameter, q_0 : initial sampler, K_h :
Gaussian transition kernel, T : total number of iterations, J : dimension of the parameter set.

Output: Optimised weights λ and parameter set Θ .

Draw $\theta_{1,1}, \dots, \theta_{J,1}$ from q_0 .

240 **for** $t = 1 \dots T$ **do**

Exploitation step : Set $\Theta = \{\theta_{1,t}, \dots, \theta_{J,t}\}$. Perform the Power Descent or Renyi Descent
and obtain the optimised mixture weights λ .

Exploration step : Perform any exploration step of our choice and obtain
 $\theta_{1,t+1}, \dots, \theta_{J,t+1}$.

241 Many choices of Exploration step can be envisioned in Algorithm 2 since there is no constraint on
242 $\{\theta_1, \dots, \theta_J\}$. Here, we consider the same Exploration step as the one they used in [17]: h is set to be
243 proportional to $J^{-1/(4+d)}$ and the particles are updated by i.i.d sampling according to $\mu_{\lambda, \Theta} k_h$ (and
244 we refer to Appendix C.5 for some details about alternative possible choices of Exploration step).

245 As for the Power Descent and Renyi Descent, we perform N transitions of these algorithms at
246 each time $t = 1 \dots T$ according to Algorithm 3 and 4, in which the initial weights are set to
247 be $[1/J, \dots, 1/J]$, $\eta = \eta_0/\sqrt{N}$ with $\eta_0 > 0$ and M samples are used in the estimation of
248 $(b_{\mu_{\lambda, \Theta, \alpha}}(\theta_{j,t}))_{1 \leq j \leq J}$ at each iteration $n = 1 \dots N$. We take $J = 100$, $M \in \{100, 1000, 2000\}$,
249 $\alpha = 0.5$, $\kappa = 0$, $\eta_0 = 0.3$ and the initial particles $\theta_1, \dots, \theta_J$ are sampled from a centered normal
250 distribution q_0 with covariance matrix $5\mathbf{I}_d$. We let $T = 10$, $N = 20$ and we replicate the experiment
251 100 times independently in dimension $d = 16$ for each algorithm. The convergence is assessed using
252 a Monte Carlo estimate of the Variational Renyi bound introduced in [10] (which requires next to
253 none additional computations).

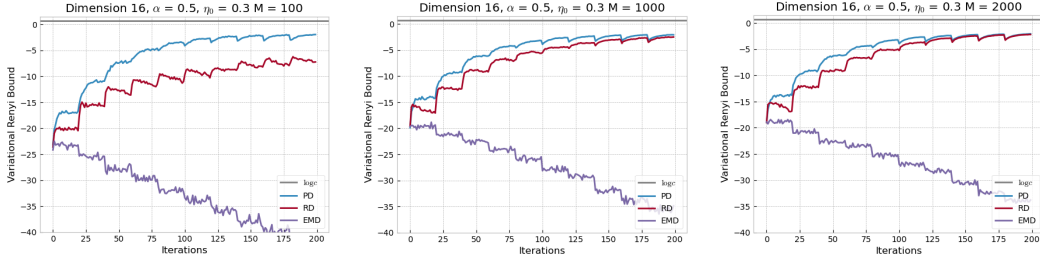
254 The results for the Power Descent and the Renyi Descent are displayed on Figure 1 below and we
255 add the Entropic Mirror Descent applied to Ψ_α as a reference.

256 We then observe that the Renyi Descent is indeed better-behaved compared to the Entropic Mirror
257 Descent applied to Ψ_α , which fails in dimension 16. Furthermore, it matches the performances of the
258 Power Descent as M increases in our numerical experiment, which illustrates the link between the
259 two algorithms we have established in the previous section.

260 **Discussion** From a theoretical standpoint, no convergence rate is yet available for the Power Descent
261 algorithm when $\alpha < 1$. An advantage of the novel Renyi Descent algorithm is then that while being
262 close to the Power Descent, it also benefits from the Entropic Mirror Descent optimisation literature
263 and as such $O(1/\sqrt{N})$ convergence rates hold, which we have been able to improve to $O(1/N)$
264 convergence rates.

265 A practical use of the Power Descent and of the Renyi Descent algorithms requires approximations to
266 handle intractable integrals appearing in the update formulas so that the Power Descent applies the

Figure 1: Plotted is the average Variational Renyi bound for the Power Descent (PD), the Renyi Descent (RD) and the Entropic Mirror Descent applied to Ψ_α (EMD) in dimension $d = 16$ computed over 100 replicates with $\eta_0 = 0.3$ and $\alpha = 0.5$ and an increasing number of samples M .



267 function $\Gamma(v) = [(\alpha - 1)v + 1]^{\eta/(1-\alpha)}$ to an unbiased estimator of the translated gradient $b_{\mu,\alpha}(\theta) + \kappa$
 268 before renormalising, while the the Renyi Descent applies the Entropic Mirror Descent function
 269 $\Gamma(v) = e^{-\eta v}$ to a biased estimator of $b_{\mu_n,\alpha}(\theta)/(\mu_n(b_{\mu_n,\alpha}) + \kappa + 1/(\alpha - 1))$ before renormalising.

270 Finding which approach is most suitable between biased and unbiased α -divergence minimisation
 271 is still an open issue in the literature, both theoretically and empirically [15, 16, 19]. Due to the
 272 exponentiation, considering the α -divergence instead of Renyi's α -divergence has for example been
 273 said to lead to high-variance gradients [11, 10] and low Signal-to-Noise ratio when $\alpha \neq 0$ [16] during
 274 the stochastic gradient descent optimization.

275 In that regard, our work sheds light on additional links between unbiased and biased α -divergence
 276 methods beyond the framework of stochastic gradient descent algorithms, as both the unbiased Power
 277 Descent and the biased Renyi Descent share the same first order approximation.

278 6 Conclusion

279 We investigated algorithms that can be used to perform mixture weights optimisation for α -divergence
 280 minimisation regardless of how the mixture parameters are obtained. We have established the full
 281 proof of the convergence of the Power Descent algorithm in the case $\alpha < 1$ when we consider mixture
 282 models and bridged the gap with the case $\alpha = 1$. We also introduced a closely-related algorithm
 283 called the Renyi Descent. We proved it enjoys an $O(1/N)$ convergence rate and illustrated in practice
 284 the proximity between these two algorithms when the number of samples M increases.

285 Further work could include establishing theoretical results regarding the stochastic version of these
 286 two algorithms, as well as providing complementary empirical results comparing the performances
 287 of the unbiased α -divergence-based Power Descent algorithm to those of the biased Renyi's α -
 288 divergence-based Renyi Descent. Since our contributions are mainly theoretical, we believe these
 289 will not result in any negative societal impacts.

290 References

- 291 [1] Michael I. Jordan, Zoubin Ghahramani, Tommi S. Jaakkola, and Lawrence K. Saul. An introduction to
 292 variational methods for graphical models. *Machine Learning*, 37(2):183–233, 1999.
- 293 [2] Matthew James. Beal. Variational algorithms for approximate bayesian inference. *PhD thesis*, 01 2003.
- 294 [3] Matthew D. Hoffman, David M. Blei, Chong Wang, and John Paisley. Stochastic variational inference.
 295 *Journal of Machine Learning Research*, 14(4):1303–1347, 2013.
- 296 [4] Tom Minka. Divergence measures and message passing. Technical Report MSR-TR-2005-173, January
 297 2005.
- 298 [5] Huaiyu Zhu and Richard Rohwer. Information geometric measurements of generalisation. Technical
 299 Report NCRG/4350, Aug 1995.
- 300 [6] Huaiyu Zhu and Richard Rohwer. Bayesian invariant measurements of generalization. *Neural Processing*
 301 *Letters*, 2:28–31, December 1995.

- 302 [7] Alfréd Rényi. On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium*
303 *on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, pages
304 547–561, Berkeley, Calif., 1961. University of California Press.
- 305 [8] Tim van Erven and Peter Harremoës. Rényi divergence and kullback-leibler divergence. *IEEE Transactions*
306 *on Information Theory*, 60(7):3797–3820, Jul 2014.
- 307 [9] Jose Hernandez-Lobato, Yingzhen Li, Mark Rowland, Thang Bui, Daniel Hernandez-Lobato, and Richard
308 Turner. Black-box alpha divergence minimization. In Maria Florina Balcan and Kilian Q. Weinberger,
309 editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings*
310 *of Machine Learning Research*, pages 1511–1520, New York, New York, USA, 20–22 Jun 2016. PMLR.
- 311 [10] Yingzhen Li and Richard E Turner. Rényi divergence variational inference. In D. D. Lee, M. Sugiyama,
312 U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*,
313 pages 1073–1081. Curran Associates, Inc., 2016.
- 314 [11] Adji Bousso Dieng, Dustin Tran, Rajesh Ranganath, John Paisley, and David Blei. Variational inference
315 via lchi upper bound minimization. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus,
316 S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages
317 2732–2741. Curran Associates, Inc., 2017.
- 318 [12] Volodymyr Kuleshov and Stefano Ermon. Neural variational inference and learning in undirected graphical
319 models. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett,
320 editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 321 [13] Robert Bamler, Cheng Zhang, Manfred Opper, and Stephan Mandt. Perturbative black box variational
322 inference. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett,
323 editors, *Advances in Neural Information Processing Systems 30*, pages 5079–5088. Curran Associates,
324 Inc., 2017.
- 325 [14] Dilin Wang, Hao Liu, and Qiang Liu. Variational inference with tail-adaptive f-divergence. In S. Bengio,
326 H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural*
327 *Information Processing Systems 31*, pages 5737–5747. Curran Associates, Inc., 2018.
- 328 [15] Tomas Geffner and Justin Domke. Empirical evaluation of biased methods for alpha divergence minimiza-
329 tion. In *3rd Symposium on Advances in Approximate Bayesian Inference*, pages 1–12, 2020.
- 330 [16] Tomas Geffner and Justin Domke. On the difficulty of unbiased alpha divergence minimization. *arXiv*
331 *preprint arXiv:2010.09541*, 2020.
- 332 [17] Kamélia Daudel, Randal Douc, and François Portier. Infinite-dimensional gradient-based descent for
333 alpha-divergence minimisation. *To appear in the Annals of Statistics*, 2021.
- 334 [18] Kamélia Daudel, Randal Douc, and François Roueff. Monotonic alpha-divergence minimisation. *arXiv*
335 *preprint arxiv:2103.05684*, 2021.
- 336 [19] Akash Kumar Dhaka, Alejandro Catalina, Manushi Welandawe, Michael Riis Andersen, Jonathan Huggins,
337 and Aki Vehtari. Challenges and opportunities in high-dimensional variational inference. *arxiv preprint*
338 *arxiv:2103.01085*, 2021.
- 339 [20] Arnaud Doucet, Nando Freitas, Kevin Murphy, and Stuart Russell. Sequential monte carlo methods in
340 practice. 01 2013.
- 341 [21] Andrzej Cichocki and Shun-ichi Amari. Families of alpha- beta- and gamma- divergences: Flexible and
342 robust measures of similarities. *Entropy*, 12(6):1532–1568, Jun 2010.
- 343 [22] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex
344 optimization. *Operations Research Letters*, 31(3):167 – 175, 2003.
- 345 [23] Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in*
346 *Machine Learning*, 8(3-4):231–357, 01 2015.

347 **Checklist**

- 348 1. For all authors...
- 349 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
350 contributions and scope? [Yes]
- 351 (b) Did you describe the limitations of your work? [Yes] End of Section 5 and Section 6.
- 352 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Section 6.
- 353 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
354 them? [Yes]
- 355 2. If you are including theoretical results...
- 356 (a) Did you state the full set of assumptions of all theoretical results? [Yes] Section 3 and
357 Section 4.
- 358 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendices.
- 359 3. If you ran experiments...
- 360 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
361 mental results (either in the supplemental material or as a URL)? [Yes] Section 5 and
362 supplementary.
- 363 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
364 were chosen)? [Yes] Section 5.
- 365 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
366 ments multiple times)? [N/A]
- 367 (d) Did you include the total amount of compute and the type of resources used (e.g., type
368 of GPUs, internal cluster, or cloud provider)? [N/A]
- 369 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 370 (a) If your work uses existing assets, did you cite the creators? [Yes] Supplementary.
- 371 (b) Did you mention the license of the assets? [N/A]
- 372 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 373 (d) Did you discuss whether and how consent was obtained from people whose data you're
374 using/curating? [N/A]
- 375 (e) Did you discuss whether the data you are using/curating contains personally identifiable
376 information or offensive content? [N/A]
- 377 5. If you used crowdsourcing or conducted research with human subjects...
- 378 (a) Did you include the full text of instructions given to participants and screenshots, if
379 applicable? [N/A]
- 380 (b) Did you describe any potential participant risks, with links to Institutional Review
381 Board (IRB) approvals, if applicable? [N/A]
- 382 (c) Did you include the estimated hourly wage paid to participants and the total amount
383 spent on participant compensation? [N/A]

384 **A**

385 **A.1 Equivalence between (1) and (2) with $p(y) = p(y, \mathcal{D})$**

386 • Case $\alpha = 1$ with $f_1(u) = 1 - u + u \log(u)$ for all $u > 0$. Then,

$$\begin{aligned} D_1(\mu K || \mathbb{P}) &= \int_{\mathcal{Y}} f_1\left(\frac{\mu k(y)}{p(y|\mathcal{D})}\right) p(y|\mathcal{D}) \nu(dy) \\ &= \int_{\mathcal{Y}} \mu k(y) \log\left(\frac{\mu k(y)}{p(y|\mathcal{D})}\right) \nu(dy) + 0 \\ &= \int_{\mathcal{Y}} \mu k(y) \log\left(\frac{\mu k(y)}{p(y, \mathcal{D})}\right) \nu(dy) + \log p(\mathcal{D}) \\ &= \int_{\mathcal{Y}} f_1\left(\frac{\mu k(y)}{p(y, \mathcal{D})}\right) p(y, \mathcal{D}) \nu(dy) + 1 - p(\mathcal{D}) + \log p(\mathcal{D}) \end{aligned}$$

Thus,

$$\operatorname{arginf}_{\mu \in \mathcal{M}} D_1(\mu K || \mathbb{P}) = \operatorname{arginf}_{\mu \in \mathcal{M}} \Psi_1(\mu; p) \quad \text{with } p(y) = p(y, \mathcal{D})$$

387 • Case $\alpha = 0$ with $f_0(u) = u - 1 - \log(u)$ for all $u > 0$.

$$\begin{aligned} D_0(\mu K || \mathbb{P}) &= \int_{\mathcal{Y}} f_0\left(\frac{\mu k(y)}{p(y|\mathcal{D})}\right) p(y|\mathcal{D}) \nu(dy) \\ &= \int_{\mathcal{Y}} -\log\left(\frac{\mu k(y)}{p(y|\mathcal{D})}\right) p(y|\mathcal{D}) \nu(dy) \\ &= \int_{\mathcal{Y}} -\log\left(\frac{\mu k(y)}{p(y, \mathcal{D})}\right) p(y|\mathcal{D}) \nu(dy) - \log p(\mathcal{D}) \\ &= \frac{1}{p(\mathcal{D})} \left[\int_{\mathcal{Y}} f_1\left(\frac{\mu k(y)}{p(y, \mathcal{D})}\right) p(y, \mathcal{D}) \nu(dy) + p(\mathcal{D}) - 1 - p(\mathcal{D}) \log p(\mathcal{D}) \right] \end{aligned}$$

Thus

$$\operatorname{arginf}_{\mu \in \mathcal{M}} D_0(\mu K || \mathbb{P}) = \operatorname{arginf}_{\mu \in \mathcal{M}} \Psi_0(\mu; p) \quad \text{with } p(y) = p(y, \mathcal{D})$$

388 • Case $\alpha \in \mathbb{R} \setminus \{1\}$ with $f_\alpha(u) = \frac{1}{\alpha(\alpha-1)} [u^\alpha - 1 - \alpha(u-1)]$ for all $u > 0$.

$$\begin{aligned} D_\alpha(\mu K || \mathbb{P}) &= \int_{\mathcal{Y}} f_\alpha\left(\frac{\mu k(y)}{p(y|\mathcal{D})}\right) p(y|\mathcal{D}) \nu(dy) \\ &= \int_{\mathcal{Y}} \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{\mu k(y)}{p(y|\mathcal{D})}\right)^\alpha - 1 \right] p(y|\mathcal{D}) \nu(dy) \\ &= p(\mathcal{D})^{\alpha-1} \int_{\mathcal{Y}} \frac{1}{\alpha(\alpha-1)} \left[\left(\frac{\mu k(y)}{p(y, \mathcal{D})}\right)^\alpha - 1 \right] p(y, \mathcal{D}) \nu(dy) + \frac{p(\mathcal{D})^\alpha - 1}{\alpha(\alpha-1)} \\ &= p(\mathcal{D})^{\alpha-1} \int_{\mathcal{Y}} f_\alpha\left(\frac{\mu k(y)}{p(y, \mathcal{D})}\right) p(y, \mathcal{D}) \nu(dy) + \frac{\alpha p(\mathcal{D})^{\alpha-1} + (1-\alpha)p(\mathcal{D})^\alpha - 1}{\alpha(\alpha-1)} \quad (10) \end{aligned}$$

Thus,

$$\operatorname{arginf}_{\mu \in \mathcal{M}} D_\alpha(\mu K || \mathbb{P}) = \operatorname{arginf}_{\mu \in \mathcal{M}} \Psi_\alpha(\mu; p) \quad \text{with } p(y) = p(y, \mathcal{D})$$

389 **A.2 [17, Theorem 1] with $\Gamma(v) = [(\alpha-1)v + 1]^{\eta/(1-\alpha)}$**

390 **Theorem 4** ([17, Theorem 1] with $\Gamma(v) = [(\alpha-1)v + 1]^{\eta/(1-\alpha)}$). Assume that p and k are as in
391 (A1). Let $\alpha \in \mathbb{R} \setminus \{1\}$, let κ be such that $(\alpha-1)\kappa \geq 0$, let $\mu \in \mathcal{M}_1(\mathcal{T})$ and let $\eta \in (0, 1]$ be such
392 that

$$0 < \mu(\Gamma(b_{\mu, \alpha} + \kappa)) < \infty \quad (11)$$

393 holds and $\Psi_\alpha(\mu) < \infty$. Then, the two following assertions hold.

394 (i) We have $\Psi_\alpha \circ \mathcal{I}_\alpha(\mu) \leq \Psi_\alpha(\mu)$.

395 (ii) We have $\Psi_\alpha \circ \mathcal{I}_\alpha(\mu) = \Psi_\alpha(\mu)$ if and only if $\mu = \mathcal{I}_\alpha(\mu)$.

396 **B**

397 **B.1 Proof that (A2) is satisfied in Example 1**

398 *Proof that (A2) is satisfied in Example 1.*

399

400 We have $k_h(\theta, y) = \frac{e^{-\|y-\theta\|^2/(2h^2)}}{(2\pi h^2)^{d/2}}$ and $p(y) = c \times \left[0.5 \frac{e^{-\|y-\theta_1^*\|^2/2}}{(2\pi)^{d/2}} + 0.5 \frac{e^{-\|y-\theta_2^*\|^2/2}}{(2\pi)^{d/2}} \right]$ for all $\theta \in \mathbb{T}$
 401 and all $y \in \mathbb{Y}$. Recall that by assumption $\mathbb{T} = \mathcal{B}(0, r) \subset \mathbb{R}^d$ with $r > 0$. Then, for all $\alpha \in [0, 1)$, we
 402 are interested in proving

$$\int_{\mathbb{Y}} \sup_{\theta \in \mathbb{T}} k(\theta, y) \times \sup_{\theta' \in \mathbb{T}} \left(\frac{k(\theta', y)}{p(y)} \right)^{\alpha-1} \nu(dy) < \infty \quad (12)$$

403 and

$$\int_{\mathbb{Y}} \sup_{\theta \in \mathbb{T}} \left| \log \left(\frac{k_h(\theta, y)}{p(y)} \right) \right| p(y) \nu(dy) < \infty. \quad (13)$$

404 (i) We start by proving (12). First note that for all $\theta, \theta' \in \mathbb{T}$ and for all $y \in \mathbb{Y}$ we can write

$$\begin{aligned} \frac{k_h(\theta, y)}{k_h(\theta', y)} &= e^{\frac{-\|y-\theta\|^2 + \|y-\theta'\|^2}{2h^2}} = e^{\frac{2\langle y, \theta-\theta' \rangle - \|\theta\|^2 + \|\theta'\|^2}{2h^2}} \\ &\leq e^{\frac{2|\langle y, \theta-\theta' \rangle| + \|\theta\|^2 + \|\theta'\|^2}{2h^2}} \leq e^{\frac{\|y\| \|\theta-\theta'\| + r^2}{h^2}}. \end{aligned}$$

405 from which we deduce that for all $\theta, \theta' \in \mathbb{T}$ and for all $y \in \mathbb{Y}$,

$$\frac{k_h(\theta, y)}{k_h(\theta', y)} \leq e^{\frac{\|y\|2r+r^2}{h^2}} \quad (14)$$

and that

$$\int_{\mathbb{Y}} \sup_{\theta \in \mathbb{T}} k(\theta, y) \times \sup_{\theta' \in \mathbb{T}} \left(\frac{k(\theta', y)}{p(y)} \right)^{\alpha-1} \nu(dy) \leq \int_{\mathbb{Y}} k(\theta, y) e^{\frac{\|y\|2r+r^2}{h^2}} \sup_{\theta' \in \mathbb{T}} \left(\frac{k(\theta', y)}{p(y)} \right)^{\alpha-1} \nu(dy).$$

406 Additionally, Jensen's inequality applied to the concave function $u \mapsto u^{1-\alpha}$ implies

$$\begin{aligned} \int_{\mathbb{Y}} k(\theta, y) e^{\frac{\|y\|2r+r^2}{h^2}} \sup_{\theta' \in \mathbb{T}} \left(\frac{k(\theta', y)}{p(y)} \right)^{\alpha-1} \nu(dy) &\leq \left(\int_{\mathbb{Y}} k(\theta, y) e^{\frac{\|y\|2r+r^2}{(1-\alpha)h^2}} \sup_{\theta' \in \mathbb{T}} \frac{p(y)}{k(\theta', y)} \nu(dy) \right)^{1-\alpha} \\ &\leq \left(\int_{\mathbb{Y}} \sup_{\theta, \theta' \in \mathbb{T}} \frac{k_h(\theta, y)}{k_h(\theta', y)} e^{\frac{\|y\|2r+r^2}{(1-\alpha)h^2}} p(y) \nu(dy) \right)^{1-\alpha} \end{aligned}$$

Now using (14), we can deduce

$$\int_{\mathbb{Y}} \sup_{\theta, \theta' \in \mathbb{T}} \frac{k_h(\theta, y)}{k_h(\theta', y)} e^{\frac{\|y\|2r+r^2}{(1-\alpha)h^2}} p(y) \nu(dy) \leq \int_{\mathbb{Y}} e^{\frac{\|y\|2r+r^2}{h^2} (1+\frac{1}{1-\alpha})} p(y) \nu(dy) < \infty,$$

407 which yields the desired result.

408 (ii) We now prove (13). For all $y \in \mathbb{Y}$ and all $\theta \in \mathbb{T}$, we have

$$\begin{aligned} e^{-\sup_{\theta \in \mathbb{T}} \frac{\|y-\theta\|^2}{2h^2}} &\leq (2\pi h^2)^{d/2} k_h(\theta, y) \leq 1 \\ e^{-\max_{i \in \{1,2\}} \frac{\|y-\theta_i^*\|^2}{2}} &\leq c^{-1} (2\pi)^{d/2} p(y) \leq 1 \end{aligned}$$

409 and we can deduce for all $y \in \mathbb{Y}$ and all $\theta \in \mathbb{T}$

$$\begin{aligned} \left| \log \left(\frac{k_h(\theta, y)}{p(y)} \right) \right| &\leq \sup_{\theta \in \mathbb{T}} \frac{\|y-\theta\|^2}{2h^2} + \max_{i \in \{1,2\}} \frac{\|y-\theta_i^*\|^2}{2} + d|\log h| + |\log c| \\ &\leq \frac{(\|y\| + r)^2}{2} \left[\frac{1}{h^2} + 1 \right] + d|\log h| + |\log c|. \end{aligned} \quad (15)$$

Since we have

$$\int_{\mathbb{Y}} \left(\frac{(\|y\| + r)^2}{2} \left[\frac{1}{h^2} + 1 \right] + d|\log h| + |\log c| \right) p(y) \nu(dy) < \infty$$

410 we deduce that (13) holds.

412 B.2 Proof of Theorem 2

413 We start with some preliminary results. Let $\zeta, \zeta' \in M_1(\mathbb{T})$. Recall that we say that $\zeta \mathcal{R} \zeta'$ if and only
414 if $\zeta K = \zeta' K$ and that $M_{1,\zeta}(\mathbb{T})$ denotes the set of probability measures dominated by ζ .

Lemma 2. *Assume (A1). Let M be a convex subset of $M_1(\mathbb{T})$ and let $\zeta_1, \zeta_2 \in M_1(\mathbb{T})$ be such that*

$$\Psi_\alpha(\zeta_1) = \Psi_\alpha(\zeta_2) = \inf_{\zeta \in M} \Psi_\alpha(\zeta).$$

415 *Then, we have $\zeta_1 \mathcal{R} \zeta_2$.*

416 *Proof.* For all $y \in Y$, set $u_y = \zeta_1 k(y)/p(y)$ and $v_y = \zeta_2 k(y)/p(y)$. Then, for all $y \in Y$ and for all
417 $t \in (0, 1)$, $f_\alpha(tu_y + (1-t)v_y) \leq t f_\alpha(u_y) + (1-t)f_\alpha(v_y)$ by convexity of f_α and we obtain

$$\Psi_\alpha(t\zeta_1 + (1-t)\zeta_2) \leq t\Psi_\alpha(\zeta_1) + (1-t)\Psi_\alpha(\zeta_2) = \inf_{\zeta \in M} \Psi_\alpha(\zeta). \quad (16)$$

418 Furthermore, $t\zeta_1 + (1-t)\zeta_2 \in M$ which implies that we have equality in (16).

Consequently, for all $t \in (0, 1)$:

$$\int_Y \underbrace{[t f_\alpha(u_y) + (1-t)f_\alpha(v_y) - f_\alpha(tu_y + (1-t)v_y)]}_{\geq 0} p(y) \nu(dy) = 0.$$

419 Now using that f_α is strictly convex, we deduce that for p -almost all $y \in Y$, $\zeta_1 k(y) = \zeta_2 k(y)$ that is
420 $\zeta_1 \mathcal{R} \zeta_2$. □

421 **Lemma 3.** *Assume (A1). Let $\alpha \in \mathbb{R} \setminus \{1\}$, let κ be such that $(\alpha - 1)\kappa \geq 0$ and let $\mu^* \in M_1(\mathbb{T})$ be
422 a fixed point of \mathcal{I}_α . Then,*

$$\Psi_\alpha(\mu^*) = \inf_{\zeta \in M_{1,\mu^*}(\mathbb{T})} \Psi_\alpha(\zeta). \quad (17)$$

423 *Furthermore, for all $\zeta \in M_{1,\mu^*}(\mathbb{T})$, $\Psi_\alpha(\mu^*) = \Psi_\alpha(\zeta)$ implies that $\mu^* \mathcal{R} \zeta$.*

424 *Proof.* Let $\zeta \in M_{1,\mu^*}(\mathbb{T})$ be such that $\Psi_\alpha(\zeta) \leq \Psi_\alpha(\mu^*)$. We have that

$$\zeta(b_{\mu^*,\alpha} - \mu^*(b_{\mu^*,\alpha})) \leq \Psi_\alpha(\zeta) - \Psi_\alpha(\mu^*) \leq 0. \quad (18)$$

425 Furthermore, since μ^* is a fixed point of \mathcal{I}_α , $\Gamma(b_{\mu^*,\alpha} + \kappa)$, hence $|b_{\mu^*,\alpha} + \kappa + 1/(\alpha - 1)|$ is μ^* -almost
426 all constant. In addition, $b_{\mu^*,\alpha} + \kappa + 1/(\alpha - 1)$ is of constant sign by assumption on κ . Since $\zeta \preceq \mu^*$,
427 we thus deduce that

$$\zeta(b_{\mu^*,\alpha} - \mu^*(b_{\mu^*,\alpha})) = 0.$$

428 Combining this result with (18) yields $\Psi_\alpha(\zeta) = \Psi_\alpha(\mu^*)$ and we recover (17).

429 Finally, assume there exists $\zeta \in M_{1,\mu^*}(\mathbb{T})$ such that $\Psi_\alpha(\mu^*) = \Psi_\alpha(\zeta)$. Then, since $M_{1,\mu^*}(\mathbb{T})$ is a
430 convex set, we have by Lemma 2 that $\mu^* \mathcal{R} \zeta$. □

431 We now move on to the proof of Theorem 2.

432 *Proof of Theorem 2.* For convenience, we define the notation $\Psi_{\alpha,\theta}(\lambda) := \Psi_\alpha(\mu_{\lambda,\theta})$ for all $\lambda \in \mathcal{S}_J$.
433 In this proof, we will use the equivalence relation \mathcal{R} defined by: $\zeta \mathcal{R} \zeta'$ if and only if $\zeta K = \zeta' K$ and
434 we write $M_{1,\zeta}(\mathbb{T})$ the set of probability measures dominated by ζ .

435 (i) *Any possible limit of convergent subsequence of $(\lambda_n)_{n \in \mathbb{N}^*}$ is a fixed point of $\mathcal{I}_\alpha^{\text{mixt}}$.*

First note that by (A3), we have that $|\Psi_{\alpha,\theta}(\lambda)| < \infty$ and that (11) is satisfied for all $\mu_{\lambda,\theta}$ such that
 $\lambda \in \mathcal{S}_J$. This means that the sequence $(\lambda_n)_{n \in \mathbb{N}^*}$ defined by (5) is well-defined, that the sequence
 $(\Psi_{\alpha,\theta}(\lambda_n))_{n \in \mathbb{N}^*}$ is lower-bounded and that $\Psi_{\alpha,\theta}(\lambda_n)$ is finite for all $n \in \mathbb{N}^*$. As $(\Psi_{\alpha,\theta}(\lambda_n))_{n \in \mathbb{N}^*}$
is nonincreasing by Theorem 4-(i), it converges in \mathbb{R} and in particular we have

$$\lim_{n \rightarrow \infty} \Psi_{\alpha,\theta} \circ \mathcal{I}_\alpha^{\text{mixt}}(\lambda_n) - \Psi_{\alpha,\theta}(\lambda_n) = 0.$$

436 Let $(\lambda_{\varphi(n)})_{n \in \mathbb{N}^*}$ be a convergent subsequence of $(\lambda_n)_{n \in \mathbb{N}^*}$ and denote by $\bar{\lambda}$ its limit. Since the
437 function $\lambda \mapsto \Psi_{\alpha,\theta} \circ \mathcal{I}_\alpha^{\text{mixt}}(\lambda) - \Psi_{\alpha,\theta}(\lambda)$ is continuous we obtain that $\Psi_{\alpha,\theta} \circ \mathcal{I}_\alpha^{\text{mixt}}(\bar{\lambda}) = \Psi_{\alpha,\theta}(\bar{\lambda})$
438 and hence by Theorem 4-(ii), $\bar{\lambda}$ is a fixed point of $\mathcal{I}_\alpha^{\text{mixt}}$.

439 (ii) The set $F = \{\lambda \in \mathcal{S}_J : \lambda = \mathcal{I}_\alpha^{\text{mixt}}(\lambda)\}$ of fixed points of $\mathcal{I}_\alpha^{\text{mixt}}$ is finite.

440 For any subset $R \subset \{1, \dots, J\}$, define

$$\begin{aligned} S_{J,R} &= \{\lambda \in \mathcal{S}_J : \forall i \in R^c, \lambda_i = 0, \forall j \in R, \lambda_j \neq 0\}, \\ \tilde{S}_{J,R} &= \{\lambda \in \mathcal{S}_J : \forall i \in R^c, \lambda_i = 0\}, \end{aligned}$$

and write

$$F = \bigcup_{R \subset \{1, \dots, J\}} (S_{J,R} \cap F).$$

441 In order to show that F is finite, we prove by contradiction that for any $R \subset \{1, \dots, J\}$, $S_{J,R} \cap F$
 442 contains at most one element. Assume indeed the existence of two distinct elements $\lambda \neq \lambda'$
 443 belonging to $S_{J,R} \cap F$. Since $M_{1, \mu_{\lambda, \Theta}}(\mathbb{T}) = M_{1, \mu_{\lambda', \Theta}}(\mathbb{T}) = \{\mu_{\lambda'', \Theta} : \lambda'' \in \tilde{S}_{J,R}\}$, Lemma 3
 444 implies that

$$\Psi_{\alpha, \Theta}(\lambda) = \inf_{\lambda'' \in \tilde{S}_{J,R}} \Psi_{\alpha, \Theta}(\lambda'') = \Psi_{\alpha, \Theta}(\lambda').$$

445 Applying again Lemma 3, we get $\mu_{\lambda, \Theta} \mathcal{R} \mu_{\lambda', \Theta}$, that is, $\mu_{\lambda, \Theta} K = \mu_{\lambda', \Theta} K$. This means that
 446 $\sum_{j=1}^J (\lambda_j - \lambda'_j) K(\theta_j, \cdot)$ is the null measure, which in turns implies the identity $\lambda = \lambda'$ since the
 447 family of measures $\{K(\theta_1, \cdot), \dots, K(\theta_J, \cdot)\}$ is assumed to be linearly independent.

448 (iii) *Conclusion.*

According to Lemma 2 applied to the convex subset of measures $M = \mathcal{S}_J$, the function $\Psi_{\alpha, \Theta}$ attains
 its global infimum at a unique $\lambda_* \in \mathcal{S}_J$. The uniqueness of λ_* actually follows from the fact that, as
 shown above, $\mu_{\lambda, \Theta} \mathcal{R} \mu_{\lambda', \Theta}$ if and only if $\lambda = \lambda'$. Then, by Theorem 4-(i) and by definition of λ_*

$$\Psi_{\alpha, \Theta} \circ \mathcal{I}_\alpha^{\text{mixt}}(\lambda_*) \leq \Psi_{\alpha, \Theta}(\lambda_*) = \inf_{\lambda' \in \mathcal{S}_J} \Psi_{\alpha, \Theta}(\lambda') \leq \Psi_{\alpha, \Theta} \circ \mathcal{I}_\alpha^{\text{mixt}}(\lambda_*),$$

449 and hence, $\Psi_{\alpha, \Theta} \circ \mathcal{I}_\alpha^{\text{mixt}}(\lambda_*) = \Psi_{\alpha, \Theta}(\lambda_*)$, showing that $\lambda_* \in F$ by Theorem 4-(ii). Since by (ii), F
 450 is finite, there exists $L \geq 1$ such that $F = \{\lambda^\ell : 1 \leq \ell \leq L\}$, where for $i \neq j$, $\lambda^i \neq \lambda^j$. Without
 451 any loss of generality, we set $\lambda^1 = \lambda_*$ to simplify the notation.

452 We now introduce a sequence $(W_\ell)_{1 \leq \ell \leq L}$ of disjoint open neighborhoods of $(\lambda^\ell)_{1 \leq \ell \leq L}$ such that
 453 for any $\ell \in \{1, \dots, L\}$,

$$\mathcal{I}_\alpha^{\text{mixt}}(W_\ell) \cap \left(\bigcup_{j \neq \ell} W_j \right) = \emptyset \quad (19)$$

454 This is possible since $\mathcal{I}_\alpha^{\text{mixt}}(\lambda^\ell) = \lambda^\ell$ and $\lambda \mapsto \mathcal{I}_\alpha^{\text{mixt}}(\lambda)$ is continuous.

455 By (i), the set F contains all the possible limits of any subsequence of $(\lambda_n)_{n \in \mathbb{N}^*}$. As a consequence,
 456 there exists $N > 0$ such that for all $n \geq N$, $\lambda_n \in \bigcup_{1 \leq \ell \leq L} W_\ell$. Combining with (19), there exists
 457 $\ell \in \{1, \dots, L\}$ such that for all $n \geq N$, $\lambda_n \in W_\ell$. Therefore λ^ℓ is the only possible limit of any
 458 convergent subsequence of $(\lambda_n)_{n \in \mathbb{N}^*}$ and as a consequence, $\lim_{n \rightarrow \infty} \lambda_n = \lambda^\ell$.

Thus, the sequence $(\mu_{\lambda_n, \Theta})_{n \in \mathbb{N}^*}$ weakly converges to $\mu_{\lambda^\ell, \Theta}$ as $n \rightarrow \infty$ and Theorem 1 can be
 applied. Since $\lambda_1 \in \mathcal{S}_J^+$, we have $M_{1, \mu_{\lambda_1, \Theta}}(\mathbb{T}) = \{\mu_{\lambda', \Theta} : \lambda' \in \mathcal{S}_J\}$ and Theorem 1-(iii) then
 shows that $\mu_{\lambda^\ell, \Theta}$ is the global arginf of Ψ_α over all $\{\mu_{\lambda', \Theta} : \lambda' \in \mathcal{S}_J\}$. Therefore, $\ell = 1$, i.e.,
 $\lambda^\ell = \lambda^1 = \lambda_*$ and

$$\Psi_{\alpha, \Theta}(\lambda_*) = \inf_{\lambda' \in \mathcal{S}_J} \Psi_{\alpha, \Theta}(\lambda').$$

459 □

460 B.3 The Power Descent for mixture models: practical version

461 The algorithm below provides one possible approximated version of the Power Descent algorithm,
 462 where we have set $\Gamma(v) = [(\alpha - 1)v + 1]^{-\frac{\eta}{1-\alpha}}$ with $\eta \in (0, 1]$.

Algorithm 3: Practical version of the Power Descent for mixture models

Input: p : measurable positive function, K : Markov transition kernel, M : number of samples, $\Theta = \{\theta_1, \dots, \theta_J\} \subset \mathbb{T}$: parameter set, $\Gamma(v) = [(\alpha - 1)v + 1]^{\frac{\eta}{1-\alpha}}$ with $\eta \in (0, 1]$, N : total number of iterations.

Output: Optimised weights λ .

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

for $n = 1 \dots N$ **do**

Sampling step : Draw independently M samples Y_1, \dots, Y_M from $\mu_{\lambda, \Theta} k$.

Expectation step : Compute $B_{\lambda} = (b_j)_{1 \leq j \leq J}$ where for all $j = 1 \dots J$

$$b_j = \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_j, Y_m)}{\mu_{\lambda, \Theta} k(Y_m)} f'_{\alpha} \left(\frac{\mu_{\lambda, \Theta} k(Y_m)}{p(Y_m)} \right)$$

and deduce $W_{\lambda} = (\lambda_j \Gamma(b_j + \kappa))_{1 \leq j \leq J}$ and $w_{\lambda} = \sum_{j=1}^J \lambda_j \Gamma(b_j + \kappa)$.

Iteration step : Set

$$\lambda \leftarrow \frac{1}{w_{\lambda}} W_{\lambda}$$

463 **C**

464 **C.1 Proof of Proposition 1**

465 We first state (D1), which summarises the necessary convergence and differentiability assumptions
466 needed in the proof of proposition 1.

467 (D1) (i) we have $\int_{\mathbb{Y}} \sup_{\theta \in \mathbb{T}} k(\theta, y) \times \sup_{\theta' \in \mathbb{T}} \left(\frac{k(\theta', y)}{p(y)} \right)^{\alpha-1} \nu(dy) < \infty$;

468 (ii) we have $\int_{\mathbb{Y}} \sup_{\theta \in \mathbb{T}} k(\theta, y) \times \sup_{\theta' \in \mathbb{T}} \left| \log \left(\frac{k(\theta', y)}{p(y)} \right) \right| \times \sup_{\theta'' \in \mathbb{T}} \left(\frac{k(\theta'', y)}{p(y)} \right)^{\alpha-1} \nu(dy) < \infty$;

469 (iii) we have $\int_{\mathbb{Y}} \inf_{\theta \in \mathbb{T}} k(\theta, y) \times \inf_{\theta' \in \mathbb{T}} \left(\frac{k(\theta', y)}{p(y)} \right)^{\alpha-1} \nu(dy) > 0$.

470 Note that these assumptions are mild if we assume that \mathbb{T} is a compact metric space, which is
471 generally the case. Assumption (D1)-(iii) is only required when $\alpha > 1$ to ensure that the quantity
472 $[(\alpha - 1)(b_{\mu, \alpha} + \kappa) + 1]^{\frac{\eta}{1-\alpha}}$ is bounded from above. This assumption could also be replaced by the
473 assumption that κ is such that $(\alpha - 1)\kappa > 0$.

Proof of proposition 1. For all $\theta \in \mathbb{T}$, the Dominated Convergence Theorem and (D1)-(i) yield

$$\lim_{\alpha \rightarrow 1} (\alpha - 1)(b_{\mu, \alpha}(\theta) + \kappa) + 1 = \lim_{\alpha \rightarrow 1} \int_{\mathbb{Y}} k(\theta, y) \left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} \nu(dy) + 0 = 1.$$

474 Then, using (D1)-(ii) we have that for all $\theta \in \mathbb{T}$,

$$\begin{aligned} \lim_{\alpha \rightarrow 1} [(\alpha - 1)(b_{\mu, \alpha}(\theta) + \kappa) + 1]^{\frac{\eta}{1-\alpha}} &= \exp \left(\lim_{\alpha \rightarrow 1} -\eta \frac{\log [(\alpha - 1)(b_{\mu, \alpha}(\theta) + \kappa) + 1]}{\alpha - 1} \right) \\ &= \exp \left(\lim_{\alpha \rightarrow 1} -\eta \frac{\int_{\mathbb{Y}} k(\theta, y) \left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} \log \left(\frac{\mu k(y)}{p(y)} \right) \nu(dy) + \kappa}{\int_{\mathbb{Y}} k(\theta, y) \left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} \nu(dy) + (\alpha - 1)\kappa} \right) \\ &= \exp \left[-\eta \int_{\mathbb{Y}} k(\theta, y) \log \left(\frac{\mu k(y)}{p(y)} \right) \nu(dy) \right] \exp(-\eta\kappa) \end{aligned}$$

In addition, by the Dominated Convergence Theorem (and (D1)-(iii) when $\alpha > 1$), we have

$$\lim_{\alpha \rightarrow 1} \mu \left([(\alpha - 1)(b_{\mu, \alpha} + \kappa) + 1]^{\frac{\eta}{1-\alpha}} \right) = \mu \left(\exp \left[-\eta \int_{\mathcal{Y}} k(\cdot, y) \log \left(\frac{\mu k(y)}{p(y)} \right) \nu(dy) \right] \right) \exp(-\eta \kappa) .$$

Thus,

$$\lim_{\alpha \rightarrow 1} [\mathcal{I}_\alpha(\mu)](h) = \int_{\mathcal{T}} \frac{\mu(d\theta) h(\theta) e^{-\eta \int_{\mathcal{Y}} k(\theta, y) \log \left(\frac{\mu k(y)}{p(y)} \right) \nu(dy)}}{\mu \left(e^{-\eta \int_{\mathcal{Y}} k(\cdot, y) \log \left(\frac{\mu k(y)}{p(y)} \right) \nu(dy)} \right)} = [\mathcal{I}_1(\mu)](h) .$$

475

□

476 C.2 Derivation of the update formula for the Renyi Descent

477 For all $\alpha \in \mathbb{R} \setminus \{0, 1\}$ and κ such that $(\alpha - 1)\kappa \geq 0$, we are interested applying the Entropic Mirror
478 Descent algorithm to the following objective function

$$\Psi_\alpha^{AR}(\mu) := \frac{1}{\alpha(\alpha - 1)} \log \left(\int_{\mathcal{Y}} \mu k(y)^\alpha p(y)^{1-\alpha} \nu(dy) + (\alpha - 1)\kappa \right)$$

479 **Lemma 4.** Assume (A1). The gradient of $\Psi_\alpha^{AR}(\mu)$ is given by $\theta \mapsto \frac{b_{\mu, \alpha}(\theta) + 1/(\alpha - 1)}{(\alpha - 1)(\mu(b_{\mu, \alpha}) + \kappa) + 1}$.

480 *Proof.* Let $\varepsilon > 0$ be small and let $\mu, \mu' \in M_1(\mathcal{T})$. Then,

$$\begin{aligned} \Psi_\alpha^{AR}(\mu + \varepsilon \mu') &= \frac{1}{\alpha(\alpha - 1)} \log \left(\int_{\mathcal{Y}} [(\mu + \varepsilon \mu') k(y)]^\alpha p(y)^{1-\alpha} \nu(dy) + (\alpha - 1)\kappa \right) \\ &= \frac{1}{\alpha(\alpha - 1)} \log \left(\int_{\mathcal{Y}} \mu k(y)^\alpha \left[1 + \alpha \varepsilon \frac{\mu' k(y)}{\mu k(y)} \right] p(y)^{1-\alpha} \nu(dy) + (\alpha - 1)\kappa + o(\varepsilon) \right) \end{aligned}$$

481 where we used that $(1 + u)^\alpha = 1 + \alpha u + o(u)$ as $u \rightarrow 0$. Thus,

$$\begin{aligned} \Psi_\alpha^{AR}(\mu + \varepsilon \mu') &= \Psi_\alpha^{AR}(\mu) + \frac{1}{\alpha(\alpha - 1)} \log \left(1 + \alpha \varepsilon \frac{\int_{\mathcal{Y}} \mu' k(y) \left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} \nu(dy)}{\int_{\mathcal{Y}} \mu k(y)^\alpha p(y)^{1-\alpha} \nu(dy) + (\alpha - 1)\kappa} + o(\varepsilon) \right) \\ &= \Psi_\alpha^{AR}(\mu) + \varepsilon \frac{1}{\alpha - 1} \frac{\int_{\mathcal{Y}} \mu' k(y) \left(\frac{\mu k(y)}{p(y)} \right)^{\alpha-1} \nu(dy)}{\int_{\mathcal{Y}} \mu k(y)^\alpha p(y)^{1-\alpha} \nu(dy) + (\alpha - 1)\kappa} + o(\varepsilon) \\ &= \Psi_\alpha^{AR}(\mu) + \varepsilon \int_{\mathcal{T}} \mu'(d\theta) \frac{1}{\alpha - 1} \frac{b_{\mu, \alpha}(\theta) + 1/(\alpha - 1)}{\mu(b_{\mu, \alpha}) + \kappa + 1/(\alpha - 1)} + o(\varepsilon) \end{aligned}$$

482 using that $\log(1 + u) = u + o(u)$ as $u \rightarrow 0$. □

483 Consequently, the iterative update formula for the Entropic Mirror Descent applied to the objective
484 function Ψ_α^{AR} is given by

$$\mu_{n+1}(d\theta) = \mu_n(d\theta) \frac{e^{-\frac{\eta}{\alpha-1} \frac{b_{\mu_n, \alpha}(\theta)}{\mu_n(b_{\mu_n, \alpha}) + \kappa + 1/(\alpha-1)}}}{\mu_n \left(e^{-\frac{\eta}{\alpha-1} \frac{b_{\mu_n, \alpha}}{\mu_n(b_{\mu_n, \alpha}) + \kappa + 1/(\alpha-1)}} \right)} , \quad n \in \mathbb{N}^* .$$

485 C.3 Proof of Theorem 3

As we shall see, the proof can be adapted from the proof of [17, Theorem 2]. For all $\mu \in M_1(\mathcal{T})$, we will use the notation

$$\mathcal{I}_\alpha^{AR}(\mu)(d\theta) = \frac{\mu(d\theta) \exp \left[-\eta \frac{b_{\mu, \alpha}(\theta)}{(\alpha-1)(\mu(b_{\mu, \alpha}) + \kappa) + 1} \right]}{\mu \left(\exp \left[-\eta \frac{b_{\mu, \alpha}}{(\alpha-1)(\mu(b_{\mu, \alpha}) + \kappa) + 1} \right] \right)}$$

to designate the one-step transition of the Renyi Descent algorithm. Note in passing that for all $\kappa' \in \mathbb{R}$, this definition can also be rewritten under the form

$$\mathcal{I}_\alpha^{AR}(\mu)(d\theta) = \frac{\mu(d\theta) \exp \left[-\eta \frac{b_{\mu,\alpha}(\theta)}{(\alpha-1)(\mu(b_{\mu,\alpha}) + \kappa) + 1} + \kappa' \right]}{\mu \left(\exp \left[-\eta \frac{b_{\mu,\alpha}}{(\alpha-1)(\mu(b_{\mu,\alpha}) + \kappa) + 1} + \kappa' \right] \right)}.$$

486 We also define

$$\begin{aligned} L_{\alpha,2} &= \eta^{-1} \sup_{\theta \in \mathbb{T}, \mu \in \mathbb{M}_1(\mathbb{T})} [(\alpha-1)(b_{\mu,\alpha}(\theta) + \kappa) + 1] \\ L &= \eta^2 \sup_{v \in \text{Dom}_\alpha^{AR}} e^{-\eta v} \\ L_{\alpha,3} &= \sup_{v \in \text{Dom}_\alpha^{AR}} e^{\eta v} \\ L_{\alpha,1} &= \inf_{v \in \text{Dom}_\alpha^{AR}} \{1 - \eta(\alpha-1)(v - \kappa')\} \times \eta \inf_{v \in \text{Dom}_\alpha^{AR}} e^{-\eta v}. \end{aligned} \quad (20)$$

487 C.3.1 Recalling [17, Lemma 5]

488 Let (ζ, μ) be a couple of probability measures where ζ is dominated by μ which we denote by $\zeta \preceq \mu$
489 and define

$$A_\alpha := \int_{\mathbb{Y}} \nu(dy) \int_{\mathbb{T}} \mu(d\theta) k(\theta, y) f'_\alpha \left(\frac{g(\theta)\mu k(y)}{p(y)} \right) [1 - g(\theta)], \quad (21)$$

490 where g is the density of ζ w.r.t μ , i.e. $\zeta(d\theta) = \mu(d\theta)g(\theta)$. We recall [17, Lemma 5] in Lemma 5
491 below.

492 **Lemma 5.** [17, Lemma 5] Assume (A1). Then, for all $\mu, \zeta \in \mathbb{M}_1(\mathbb{T})$ such that $\zeta \preceq \mu$ and
493 $\Psi_\alpha(\mu) < \infty$, we have

$$A_\alpha \leq \Psi_\alpha(\mu) - \Psi_\alpha(\zeta). \quad (22)$$

494 Moreover, equality holds in (22) if and only if $\zeta = \mu$.

495 C.3.2 Adaptation of [17, Theorem 1]

496 **Lemma 6.** Assume (A1) and (A4). Let $\alpha \in \mathbb{R} \setminus \{1\}$, let κ be such that $(\alpha-1)\kappa \geq 0$ and let
497 $\mu \in \mathbb{M}_1(\mathbb{T})$ be such that

$$0 < \mu \left\{ \exp \left(-\eta \frac{b_{\mu,\alpha} + 1/(\alpha-1)}{(\alpha-1)(\mu(b_{\mu,\alpha}) + \kappa) + 1} \right) \right\} < \infty \quad (23)$$

498 holds and $\Psi_\alpha(\mu) < \infty$. Then, the two following assertions hold.

499 (i) We have $\Psi_\alpha \circ \mathcal{I}_\alpha^{AR}(\mu) \leq \Psi_\alpha(\mu)$.

500 (ii) We have $\Psi_\alpha \circ \mathcal{I}_\alpha^{AR}(\mu) = \Psi_\alpha(\mu)$ if and only if $\mu = \mathcal{I}_\alpha^{AR}(\mu)$.

501 *Proof.* The proof builds on the proof of [17, Theorem 1] in the particular case $\alpha \in \mathbb{R} \setminus \{1\}$. Indeed,
502 in this case,

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

503 so that

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

where $(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1 > 0$ under **(A1)**. Set

$$g = \tilde{\Gamma} \circ \left(\frac{b_{\mu,\alpha} + 1/(\alpha - 1)}{(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1} \right)$$

where for all $v \in \text{Dom}_{\alpha}^{AR}$,

$$\tilde{\Gamma}(v) = \frac{e^{-\eta v}}{\mu \left\{ \exp \left(-\eta \frac{b_{\mu,\alpha} + 1/(\alpha - 1)}{(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1} - \eta \kappa' \right) \right\}}.$$

Finally, let us consider the probability space $(\mathbb{T}, \mathcal{T}, \mu)$ and let V be the random variable

$$V(\theta) = \frac{b_{\mu,\alpha}(\theta) + 1/(\alpha - 1)}{(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1} + \kappa'.$$

504 Then, we have $\mathbb{E}[1 - \tilde{\Gamma}(V)] = 0$ and we can write

$$\begin{aligned} A_{\alpha} &= [(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1] \times \mathbb{E}[(V - \kappa')\tilde{\Gamma}^{\alpha-1}(V)(1 - \tilde{\Gamma}(V))] \\ &= [(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1] \times \mathbb{Cov}((V - \kappa')\tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V)). \end{aligned} \quad (24)$$

505 Under **(A4)** with $\alpha \in \mathbb{R} \setminus \{1\}$, $v \mapsto (v - \kappa')\tilde{\Gamma}^{\alpha-1}(v)$ and $v \mapsto 1 - \tilde{\Gamma}(v)$ are increasing on Dom_{α}^{AR}
506 which implies $\mathbb{Cov}(V\tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V)) \geq 0$ and thus $A_{\alpha} \geq 0$ since $(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1 > 0$.
507 \square

508 C.3.3 Adaptation of [17, Lemma 6]

509 Consider the probability space $(\mathbb{T}, \mathcal{T}, \mu)$ and denote by $\mathbb{V}\text{ar}_{\mu}$ the associated variance operator.

510 **Lemma 7.** Assume **(A1)** and **(A4)**. Let $\alpha \in \mathbb{R} \setminus \{1\}$, let κ be such that $(\alpha - 1)\kappa > 0$, and let
511 $\mu \in \mathbb{M}_1(\mathbb{T})$ be such that (23) holds and $\Psi_{\alpha}(\mu) < \infty$. Then,

$$\frac{(\alpha - 1)\kappa L_{\alpha,1}}{2} \mathbb{V}\text{ar}_{\mu} \left(\frac{b_{\mu,\alpha} + 1/(\alpha - 1)}{(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1} \right) \leq \Psi_{\alpha}(\mu) - \Psi_{\alpha} \circ \mathcal{I}_{\alpha}^{AR}(\mu), \quad (25)$$

where

$$L_{\alpha,1} := \inf_{v \in \text{Dom}_{\alpha}^{AR}} \{1 - \eta(\alpha - 1)(v - \kappa')\} \times \inf_{v \in \text{Dom}_{\alpha}^{AR}} \eta e^{-\eta v}.$$

512 *Proof.* The proof of Lemma 7 builds on the proof of [17, Lemma 6], which can be found in the
513 supplementary material of [17]. Using (24) combined with the fact that under **(A1)**, $(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1 > 0$
514 $(\alpha - 1)\kappa > 0$

$$\begin{aligned} A_{\alpha} &= [(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1] \times \mathbb{Cov}((V - \kappa')\tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V)) \\ &> (\alpha - 1)\kappa \times \mathbb{Cov}((V - \kappa')\tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V)) \end{aligned}$$

515 Furthermore,

$$\begin{aligned} &\mathbb{Cov}((V - \kappa')\tilde{\Gamma}^{\alpha-1}(V), 1 - \tilde{\Gamma}(V)) \\ &= \frac{1}{2} \mathbb{E} \left[((U - \kappa')\tilde{\Gamma}^{\alpha-1}(U) - (V - \kappa')\tilde{\Gamma}^{\alpha-1}(V))(-\tilde{\Gamma}(U) + \tilde{\Gamma}(V)) \right] \\ &= \frac{1}{2} \mathbb{E} \left[\frac{(U - \kappa')\tilde{\Gamma}^{\alpha-1}(U) - (V - \kappa')\tilde{\Gamma}^{\alpha-1}(V)}{U - V} \frac{-\tilde{\Gamma}(U) + \tilde{\Gamma}(V)}{U - V} (U - V)^2 \right] \\ &\geq \frac{L_{\alpha,1}}{2} \mathbb{V}\text{ar}_{\mu} \left(\frac{b_{\mu,\alpha} + 1/(\alpha - 1)}{(\alpha - 1)(\mu(b_{\mu,\alpha}) + \kappa) + 1} \right) \end{aligned}$$

516 and we thus obtain (25). \square

517 **C.3.4 Adaptation of the proof of [17, Theorem 2] to obtain Theorem 3**

518 *Proof of Theorem 3.* The proof of Theorem 3 builds on the proof of [17, Theorem 2], which can be
 519 found in the supplementary material of [17]. We prove the assertions successively.

520 (i) The proof of (i) simply consists in verifying that we can apply Lemma 6. For all $\mu \in M_1(\mathbb{T})$, (23)
 521 with $\mu = \mu_n$ holds for all $n \in \mathbb{N}^*$ by assumption on $|B|_{\infty, \alpha}$ and since at each step $n \in \mathbb{N}^*$, Lemma 6
 522 combined with $\Psi_\alpha(\mu_n) < \infty$ implies that $\Psi_\alpha(\mu_{n+1}) \leq \Psi_\alpha(\mu_n) < \infty$, we obtain by induction that
 523 $(\Psi_\alpha(\mu_n))_{n \in \mathbb{N}^*}$ is non-increasing.

524 (ii) Let $n \in \mathbb{N}^*$, set $\Delta_n = \Psi_\alpha(\mu_n) - \Psi_\alpha(\mu^*)$ and for all $\theta \in \mathbb{T}$, $V_n(\theta) = \frac{b_{\mu_n, \alpha}(\theta) + \frac{1}{\alpha-1}}{(\alpha-1)(\mu_n(b_{\mu_n, \alpha}) + \kappa) + 1} + \kappa'$,
 525 such that $d\mu_{n+1} \propto e^{-\eta V_n} d\mu_n$.

526 We first show that

$$\Delta_n \leq L_{\alpha, 2} \left[\int_{\mathbb{T}} \log \left(\frac{d\mu_{n+1}}{d\mu_n} \right) d\mu^* + \frac{L}{2} \mathbb{V}\text{ar}_{\mu_n}(V_n) L_{\alpha, 3} \right]. \quad (26)$$

527 The convexity of f_α implies that

$$\Delta_n \leq \int_{\mathbb{T}} b_{\mu_n, \alpha} (d\mu_n - d\mu^*) \quad (27)$$

$$\begin{aligned} &= \int_{\mathbb{T}} \left(b_{\mu_n, \alpha} + \frac{1}{\alpha-1} \right) (d\mu_n - d\mu^*) \\ &= \frac{(\alpha-1)(\mu_n(b_{\mu_n, \alpha}) + \kappa) + 1}{\eta} \int_{\mathbb{T}} (\mu_n(\eta V_n) - \eta V_n) d\mu^*. \end{aligned} \quad (28)$$

Then, noting that

$$-\eta V_n = \log \mu_n (e^{-\eta V_n}) + \log \left(\frac{d\mu_{n+1}}{d\mu_n} \right)$$

528 we deduce

$$\Delta_n \leq L_{\alpha, 2} \int_{\mathbb{T}} \left[\mu_n(\eta V_n) + \log \mu_n (e^{-\eta V_n}) + \log \left(\frac{d\mu_{n+1}}{d\mu_n} \right) \right] d\mu^*. \quad (29)$$

529 Since $v \mapsto e^{-\eta v}$ is L -smooth on Dom_α^{AR} , for all $\theta \in \mathbb{T}$ and for all $n \in \mathbb{N}^*$ we can write

$$e^{-\eta V_n(\theta)} \leq e^{-\eta \mu_n(V_n)} + \eta e^{-\eta \mu_n(V_n)} (V_n(\theta) - \mu_n(V_n)) + \frac{L}{2} (V_n(\theta) - \mu_n(V_n))^2$$

which in turn implies

$$\mu_n(e^{-\eta V_n}) \leq e^{-\eta \mu_n(V_n)} + \frac{L}{2} \mathbb{V}\text{ar}_{\mu_n}(V_n).$$

Finally, we obtain

$$\log \mu_n(e^{-\eta V_n}) \leq \log e^{-\eta \mu_n(V_n)} + \log \left(1 + \frac{L \mathbb{V}\text{ar}_{\mu_n}(V_n)}{2 e^{-\eta \mu_n(V_n)}} \right).$$

Using that $\log(1+u) \leq u$ when $u \geq 0$ and by definition of $L_{\alpha, 3}$, we deduce

$$\log \mu_n(e^{-\eta V_n}) \leq -\eta \mu_n(V_n) + \frac{L}{2} \mathbb{V}\text{ar}_{\mu_n}(V_n) L_{\alpha, 3},$$

which combined with (29) implies (26). To conclude, we apply Lemma 7 to $g = \frac{d\mu_{n+1}}{d\mu_n}$ and combining with (26), we obtain

$$\Delta_n \leq L_{\alpha, 2} \left[\int_{\mathbb{T}} \log \left(\frac{d\mu_{n+1}}{d\mu_n} \right) d\mu^* + \frac{L L_{\alpha, 3}}{L_{\alpha, 1}(\alpha-1)\kappa} (\Delta_n - \Delta_{n+1}) \right],$$

530 where by assumption $L_{\alpha, 1}$, $L_{\alpha, 2}$ and $L_{\alpha, 3} > 0$. As the r.h.s involves two telescopic sums, we deduce

$$\frac{1}{N} \sum_{n=1}^N \Psi_\alpha(\mu_n) - \Psi_\alpha(\mu^*) \leq \frac{L_{\alpha, 2}}{N} \left[KL(\mu^* || \mu_1) - KL(\mu^* || \mu_{N+1}) + L \frac{L_{\alpha, 3}}{L_{\alpha, 1}(\alpha-1)\kappa} (\Delta_1 - \Delta_{N+1}) \right]$$

531 and we recover (9) using (i), that $KL(\mu^* || \mu_{N+1}) \geq 0$ and that $\Delta_{N+1} \geq 0$.

532 □

533 **C.4 The Renyi Descent for mixture models: practical version**

534 The algorithm below provides one possible approximated version of the Renyi Descent algorithm,
 535 where we have set $\Gamma(v) = e^{-\eta v}$ with $\eta > 0$.

Algorithm 4: *Practical version of the Renyi Descent for mixture models*

Input: p : measurable positive function, K : Markov transition kernel, M : number of samples,
 $\Theta = \{\theta_1, \dots, \theta_J\} \subset \mathbb{T}$: parameter set, $\Gamma(v) = e^{-\eta v}$ with $\eta > 0$, N : total number of iterations.

Output: Optimised weights λ .

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

for $n = 1 \dots N$ **do**

Sampling step : Draw independently M samples Y_1, \dots, Y_M from $\mu_{\lambda, \Theta} k$.

Expectation step : Compute $B_\lambda = (b'_j)_{1 \leq j \leq J}$ where for all $j = 1 \dots J$

$$b_j = \frac{1}{M} \sum_{m=1}^M \frac{k(\theta_j, Y_m)}{\mu_{\lambda, \Theta} k(Y_m)} f'_\alpha \left(\frac{\mu_{\lambda, \Theta} k(Y_m)}{p(Y_m)} \right)$$

and for all $j = 1 \dots J$

$$b'_j = \frac{b_j}{(\alpha - 1)(\sum_{\ell=1}^J b_\ell + \kappa) + 1}$$

and deduce $W_\lambda = (\lambda_j \Gamma(b'_j + \kappa'))_{1 \leq j \leq J}$ and $w_\lambda = \sum_{j=1}^J \lambda_j \Gamma(b'_j + \kappa')$.

Iteration step : Set

$$\lambda \leftarrow \frac{1}{w_\lambda} W_\lambda$$

536 **C.5 Alternative Exploration step in Algorithm 2**

537 We present here several possible alternative choices of Exploration step in Algorithm 2, beyond the
 538 one we have made in Section 5 and that is based on [18]. Our goal here is not to discriminate between
 539 all of them, but to illustrate the generality of our approach.

540 **Gradient Descent.** One could use a Gradient Descent approach to optimise the mixture components
 541 parameters $\{\theta_{1,t+1}, \dots, \theta_{J,t+1}\}$ in the spirit of Renyi's α -divergence gradient-based methods (e.g
 542 [9, 10]) or α -divergence gradient-based methods (e.g [11, 12]).

543 **The particular case** $\alpha \in [0, 1)$. Following [18], if we consider the specific case $\alpha \in [0, 1)$ another
 544 possibility would be to set at time t : for all $j = 1 \dots J$

$$\theta_{j,t+1} = \operatorname{argmax}_{\theta_j \in \mathbb{T}} \int_{\mathcal{Y}} \gamma_{j,\alpha}^t(y) \log(k(\theta_j, y)) \nu(dy) \quad (30)$$

545 where for all $y \in \mathcal{Y}$,

$$\gamma_{j,\alpha}^t(y) = k(\theta_{j,t}, y) \left(\frac{\mu_{\lambda, \Theta} k(y)}{p(y)} \right)^{\alpha-1}.$$

546 Indeed, [18] showed that the above update formulas for $\{\theta_{1,t+1}, \dots, \theta_{J,t+1}\}$ ensure a systematic
 547 decrease in the α -divergence and they notably explained how these update formulas could even
 548 outperform typical Renyi's α / α -divergence gradient-based approaches (we refer to [18] for details).

549 Furthermore, in the particular case of d -dimensional Gaussian kernels with $k(\theta_{j,t}, y) =$
 550 $\mathcal{N}(y; m_{j,t}, \Sigma_{j,t})$ and where $\theta_{j,t} = (m_{j,t}, \Sigma_{j,t}) \in \mathbb{T}$ denotes the mean and covariance matrix of
 551 the j -th Gaussian component density, they obtained that the maximisation procedure (30) amounts to

552 setting

$$\forall j = 1 \dots J, \quad m_{j,t+1} = \frac{\int_{\mathcal{Y}} \gamma_{j,\alpha}^t(y) y \nu(dy)}{\int_{\mathcal{Y}} \gamma_{j,\alpha}^t(y) \nu(dy)}$$

$$\Sigma_{j,t+1} = \frac{\int_{\mathcal{Y}} \gamma_{j,\alpha}^t(y) (y - m_{j,t})(y - m_{j,t})^T \nu(dy)}{\int_{\mathcal{Y}} \gamma_{j,\alpha}^t(y) \nu(dy)}.$$

553 These update formulas can then always be made feasible by resorting to Monte Carlo approximations
 554 and can be used as a valid Exploration step. If we were to focus on solely updating the means
 555 $(m_{j,t+1})_{1 \leq j \leq J}$, we could for example consider the Exploration step given by:

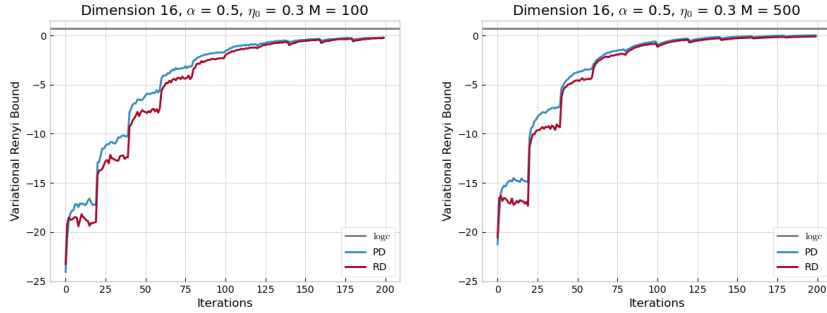
$$\forall j = 1 \dots J, \quad \theta_{j,t+1} = m_{j,t+1} = \frac{\sum_{m=1}^M \hat{\gamma}_j^{(t)}(Y'_m; \boldsymbol{\lambda}) \cdot Y'_m}{\sum_{m=1}^M \hat{\gamma}_j^{(t)}(Y'_m; \boldsymbol{\lambda})}$$

where the M samples $(Y'_m)_{1 \leq m \leq M}$ have been drawn independently from the proposal $\mu_{\boldsymbol{\lambda}, \Theta}$ and where we have set

$$\hat{\gamma}_j^{(t)}(y; \boldsymbol{\lambda}) = \frac{k(\theta_{j,t}, y)}{\mu_{\boldsymbol{\lambda}, \Theta} k(y)} \left(\frac{\mu_{\boldsymbol{\lambda}, \Theta} k(y)}{p(y)} \right)^{\alpha-1}.$$

556 We ran Algorithm 2 over 100 replicates for this choice of Exploration step with $M \in \{100, 500\}$ (and
 557 keeping the same target p , initial sampler q_0 , and hyperparameters $N = 20$, $T = 10$, $\eta = \eta_0/\sqrt{N}$
 558 with $\eta_0 = 0.3$, $\alpha = 0.5$, $J = 100$, $\kappa = 0$ and $d = 16$ as those chosen in Section 5). The results
 559 when using the Power and the Renyi Descent as Exploitation steps can be visualised in the figure
 560 below.

Figure 2: Plotted is the average Variational Renyi bound for the Power Descent (PD) and the Renyi Descent (RD) in dimension $d = 16$ computed over 100 replicates with $\eta_0 = 0.3$ and $\alpha = 0.5$ and an increasing number of samples M .



561 We then observe a similar behavior for the Power and the Renyi Descent, which illustrates the
 562 closeness between both algorithms, irrespective of the choice of the Exploration step.