

---

# Variational Inference with Mixtures of Isotropic Gaussians

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Variational inference (VI) is a popular approach in Bayesian inference, that looks  
2 for the best approximation of the posterior distribution within a parametric family,  
3 minimizing a loss that is typically the (reverse) Kullback-Leibler (KL) diver-  
4 gence. In this paper, we focus on the following parametric family: mixtures of  
5 isotropic Gaussians (i.e., with diagonal covariance matrices proportional to the  
6 identity) and uniform weights. We develop a variational framework and provide  
7 efficient algorithms suited for this family. In contrast with mixtures of Gaussian  
8 with generic covariance matrices, this choice presents a balance between accurate  
9 approximations of multimodal Bayesian posteriors, while being memory and  
10 computationally efficient. Our algorithms implement gradient descent on the  
11 location of the mixture components (the modes of the Gaussians), and either (an  
12 entropic) Mirror or Bures descent on their variance parameters. We illustrate the  
13 performance of our algorithms on numerical experiments.

## 14 1 Introduction

15 The core problem of Bayesian inference is to sample from a posterior distribution  $\pi$  over model  
16 parameters, combining prior knowledge with observed data. Unfortunately, the posterior distribution  
17 is generally difficult to compute due to the presence of an intractable integral (the normalization  
18 constant), and its density with respect to the Lebesgue measure on  $\mathbb{R}^d$ , also denoted  $\pi$ , is only known  
19 in unnormalized form as  $\pi \propto e^{-V}$ . Variational inference (VI, [Blei et al., 2017]) is a prominent  
20 alternative to standard Markov chain Monte Carlo (MCMC, [Roberts and Rosenthal, 2004]) that  
21 approximates the posterior by optimizing over a family of tractable distributions. While this restriction  
22 can introduce bias, VI is typically much faster than MCMC, since it reframes the sampling problem  
23 as a (generally finite-dimensional) optimization one over the parameters of the variational family.  
24 Specifically, VI seeks a distribution in a parametric family  $\mathcal{C}$  that minimizes a discrepancy to the  
25 target posterior, typically the reverse Kullback-Leibler (KL) divergence:

$$\min_{\mu \in \mathcal{C}} \text{KL}(\mu|\pi), \quad (1)$$

26 where  $\text{KL}(\mu|\pi) = \int \log(d\mu/d\pi) d\mu$  if  $\mu$  is absolutely continuous with respect to  $\pi$  denoting  $d\mu/d\pi$  its  
27 Radon-Nikodym density, and  $+\infty$  else.

28 There exist various choices of variational families and suited algorithms in the literature of VI.  
29 For instance, Mean-field variational inference (MFVI) aims to find an approximate posterior that  
30 factors as a product of distributions [Lacker, 2023]. Recent studies by Arnese and Lacker [2024] and  
31 Lavenant and Zanella [2024] have investigated solving this problem through the lens of coordinate  
32 ascent variational inference (CAVI, Bishop [2006]), or over polyhedral sets with first-order algorithms  
33 [Jiang et al., 2024]. Traditionally, the variational family is parameterized by a finite-dimensional set  
34 of parameters, yet, Yao and Yang [2022], Yao et al. [2024] have extended this framework to MFVI

over infinite-dimensional families. Alternatively, Gaussian variational inference [Barber and Bishop, 1997, Oppor and Archambeau, 2009] has garnered significant attention due to its computational tractability and theoretical appeal [Challis and Barber, 2013]. More recently, a new class of Gaussian VI algorithms has emerged, based on the discretization of the gradient flow of the KL divergence on the space of Gaussian measures equipped with the Wasserstein-2 metric: Lambert et al. [2022] proposed a forward (explicit) time discretization approach, while Diao et al. [2023], Domke et al. [2023] introduced a forward-backward splitting scheme inspired by [Salim et al., 2020]. Other algorithms were recently introduced for Gaussian VI by optimizing a weighted score-based (Fisher) divergence [Modi et al., 2025, Cai et al., 2024b]. These methods bring valuable insights for Gaussian VI, but may provide a too crude approximation of the target distribution if the latter is multimodal. In this setting, deep generative approaches such as normalizing flows [Tabak and Vanden-Eijnden, 2010, Papamakarios et al., 2021, Kobyzev et al., 2020] offering tractable, normalized densities and efficient sampling, have emerged as flexible and powerful alternatives to traditional Gaussian or factorized variational families. While it offers a flexible and widely applicable approach to VI, it is not tailored to a specific variational family. Moreover, despite their ability to represent multimodal distributions using expressive variational families with neural networks, it can suffer from mode collapse [Soletskyi et al., 2024], especially as the dimension increases, as shown recently in large scale empirical evaluation from [Blessing et al., 2024].

A particularly compelling variational family is the one of mixtures of Gaussians. Indeed, the latter can capture complex, multimodal distributions. Also, they lead to tractable approximate posteriors, since sampling mixtures is computationally cheap, and also marginalizations, or expectations of linear and quadratic functions can be computed in closed-form. For this variational family, several algorithms have been proposed [Lin et al., 2019, Arenz et al., 2018, 2023] based on natural gradient descent over the natural parameters of the Gaussians. An extension of the Wasserstein gradient flow approach for Gaussian mixtures has also been proposed Lambert et al. [2022]. However, their practical utility is often hindered by the computational challenges of inference, particularly when handling high-dimensional distributions. Indeed, a key challenge in deploying Gaussian mixture models lies in parameterizing the covariance matrices. While full covariance matrices provide maximum flexibility, their quadratic scaling with dimensionality leads to high computational demands, even if some solutions have been proposed to store them efficiently [Challis and Barber, 2013, Bonnabel et al., 2024]. To address this issue, we restrict the variational family of mixtures of Gaussians to diagonal, more precisely isotropic covariance matrices (i.e., assuming equal variance across all dimensions) with uniform weights. As mode collapse may result from mean alignment and vanishing weights (see [Soletskyi et al., 2024]), using fixed weights, as we do, prevents the latter. With this structure, each Gaussian in the mixture is parameterized by  $d + 1$  parameters in dimension  $d$ . As a result, a mixture of  $N$  isotropic Gaussians requires a memory cost of  $N(d + 1)$ , compared to  $N(d^2 + d)$  for a full covariance mixture. Then, optimizing mixtures of isotropic Gaussians incurs a memory cost roughly equivalent to that of optimizing the means only. We show in this study that this choice of variational family balances between accuracy of the variational approximation, i.e. its ability to model multimodal target distributions, and the computational efficiency of the associated algorithms.

Our contributions include the development of a variational framework and algorithms tailored to isotropic Gaussian mixtures, as well as an empirical evaluation across synthetic and real-world datasets, demonstrating the trade-offs between accuracy and computational cost. This paper is organized as follows. Section 2 provides the relevant background on the geometry of the space of isotropic Gaussians, and on optimization schemes based on the Bures and entropic mirror descent geometries. In Section 3, we introduce the general setting of optimizing over mixtures of isotropic Gaussians with uniform weights. Section 4 presents our algorithms to efficiently optimize over this variational family. In Section 5 we discuss related work in the Variational Inference literature. Our numerical results are to be found in Section 6.

**Notation.** We write a Gaussian distribution on  $\mathbb{R}^d$  with mean  $m$  and variance  $\epsilon$  as  $\mathcal{N}(m, \epsilon I_d)$ , where  $I_d$  denotes the  $d$ -dimensional identity matrix; and  $\mathcal{N}(x; m, \epsilon I_d)$  its density evaluated at  $x$ . We denote by  $\mathcal{P}_2(\mathbb{R}^d)$  the set of probability distributions on  $\mathbb{R}^d$  with bounded second moments. Consider  $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$ , the Wasserstein-2 ( $W_2$ ) distance is defined as  $W_2^2(\mu, \nu) = \inf_{s \in \mathcal{S}(\mu, \nu)} \int \|x - y\|^2 ds(x, y)$ , where  $\mathcal{S}(\mu, \nu)$  is the set of couplings between  $\mu$  and  $\nu$ . We will denote  $k_\epsilon$  the normalized Gaussian kernel on  $\mathbb{R}^d$  with variance  $\epsilon$ , i.e.  $k_\epsilon(x) = (2\pi\epsilon)^{-d/2} \exp(-\|x\|^2/(2\epsilon))$ . For  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ , we denote by  $k_\epsilon \star \mu$  its convolution with the Gaussian kernel, that writes  $k_\epsilon \star \mu = \int k_\epsilon(\cdot - x) d\mu(x)$ . We denote  $\text{Tr}$  the trace function.

## 92 2 Preliminaries (Gaussian VI)

93 This section introduces key concepts on the space of isotropic Gaussians, as well as different (time-  
94 discretized) gradient flows one can consider through the Bures-Wasserstein or entropic mirror descent  
95 geometries.

### 96 2.1 Isotropic Gaussians (IG)

97 The space of isotropic Gaussians is defined as  $\text{IG} = \{\mathcal{N}(m, \epsilon \text{Id}), m \in \mathbb{R}^d, \epsilon \in \mathbb{R}^{+*}\}$  and is a  
98 subspace of  $\mathcal{P}_2(\mathbb{R}^d)$ . When equipped with the  $W_2$  distance, this space has a particularly tractable ge-  
99 ometric structure. Indeed, the  $W_2$  distance between two isotropic Gaussians  $\mathcal{N}(m, \epsilon \text{Id}), \mathcal{N}(m', \tau \text{Id})$   
100 takes the form of a Bures-Wasserstein (BW) distance:

$$\text{BW}^2(\mathcal{N}(m, \epsilon \text{Id}), \mathcal{N}(m', \tau \text{Id})) = \|m - m'\|^2 + \text{B}^2(\epsilon \text{Id}, \tau \text{Id}), \quad (2)$$

101 where  $\text{B}$  denotes the Bures [Bhatia et al., 2019] metric between positive definite matrices and  
102  $\text{B}^2(\epsilon \text{Id}, \tau \text{Id}) = d(\epsilon + \tau - 2\sqrt{\epsilon\tau})$ . This formula reflects the separable nature (with respect to the  
103 means and variances) of the Wasserstein metric on the space of isotropic Gaussians  $\text{IG}$ . Interestingly,  
104 the metric space  $(\text{IG}, \text{BW})$  of isotropic Gaussians equipped with the BW distance can be seen as a  
105 submanifold of the space of (all) Gaussian distributions equipped with the same metric, which can  
106 itself be seen as a submanifold of the Wasserstein space  $(\mathcal{P}_2(\mathbb{R}^d), W_2)$ . Indeed, the BW geodesic  
107 between  $\mu, \nu \in \text{IG}$  also lies in  $\text{IG}$ , see Appendix A.1.

### 108 2.2 Bures-Wasserstein gradient descent on IG

109 Recall that our goal is to minimize a functional objective  $\text{KL}(\cdot|\pi)$  as defined in Eq (1), where  
110  $\pi \propto e^{-V}$ , firstly on  $\text{IG}$  in this section, before being able to tackle mixtures of isotropic Gaussians.  
111 In this subsection we explain how to derive a gradient flow with respect to the Bures-Wasserstein  
112 geometry, and provide a discrete optimization scheme. To this goal, we first define a minimizing  
113 movement scheme on  $\text{IG}$ . For  $p_0 \in \text{IG}$  and  $\gamma > 0$  a step-size, define:

$$p_{k+1} = \arg \min_{p \in \text{IG}} \left\{ \text{KL}(p|\pi) + \frac{1}{2\gamma} \text{BW}^2(p, p_k) \right\}, \quad (3)$$

114 which corresponds to a JKO scheme Jordan et al. [1998], but where the solution is constrained to lie  
115 in  $\text{IG}$ . In the limit  $\gamma \rightarrow 0$ , we obtain a Wasserstein gradient flow of measures projected on  $\text{IG}$ , i.e. a  
116 continuous curve  $(p_t)_t \in \text{IG}$  decreasing the KL, and which is governed by differential equations for  
117 the mean  $(m_t)_t$  and variance  $(\epsilon_t)_t$  (see Appendix A.3). Such a flow can exhibit a favorable dynamical  
118 behavior under a strong log-concavity assumption on the target distribution, as demonstrated in the  
119 following Proposition.

120 **Proposition 2.1.** *Suppose that  $\nabla^2 V \succeq \alpha \text{Id}$  for some  $\alpha \in \mathbb{R}$ . Then, for any  $p_0 \in \text{IG}$ , there is a unique*  
121 *solution  $(p_t)_t$  to the flow obtained as a limit of Eq (3) as  $\gamma \rightarrow 0$ . Then, for all  $t \geq 0$  and  $p^* \in \text{IG}$ ,*

$$\text{KL}(p_t|\pi) - \text{KL}(p^*|\pi) \leq e^{-2\alpha t} \{ \text{KL}(p_0|\pi) - \text{KL}(p^*|\pi) \},$$

122 *implying that the flow converges linearly when  $\alpha > 0$ .*

123 The full proof of Proposition 2.1 is provided in Appendix A.2, and is a direct application of the one  
124 of [Lambert et al., 2022, Corollary 3]. It relies on the fact that  $\text{KL}(\cdot|\pi)$  is an  $\alpha$ -convex objective  
125 functional along  $W_2$  geodesics when the target potential  $V$  is  $\alpha$ -convex, see e.g. [Villani, 2009,  
126 Theorem 17.15]; and since  $W_2$  (equivalently BW) geodesics between isotropic Gaussians lie in  
127  $\text{IG}$  (Appendix A.1), the objective is also convex on  $(\text{IG}, \text{BW})$ . This result indicates that strongly  
128 log-concave targets can be efficiently approximated using isotropic Gaussians.

129 Here, we propose to evaluate an explicit time-discretization of this gradient flow, as it is computationally  
130 less expensive than an implicit scheme such as Eq (3). To this goal, let  $F : \mathbb{R}^d \times \mathbb{R}^{+*}$  defined as  
131  $F(m, \epsilon) := \text{KL}(\mathcal{N}(m, \epsilon \text{Id})|\pi)$ , and starting from some  $\eta_0 = (m_0, \epsilon_0) \in \mathbb{R}^d \times \mathbb{R}^{+*}$  the scheme:

$$\eta_{k+1} = \underset{\eta \in \mathbb{R}^d \times \mathbb{R}^{+*}}{\text{argmin}} \left\{ \langle \nabla F(\eta_k), \eta - \eta_k \rangle + \frac{1}{2\gamma} \text{BW}^2(\mathcal{N}(\eta), \mathcal{N}(\eta_k)) \right\}, \quad (4)$$

denoting  $\eta = (m, \epsilon)$  and  $\mathcal{N}(\eta) = \mathcal{N}(m, \epsilon \mathbf{I}_d)$ . Note that the scheme above is similar to Eq (3) but where the objective has been linearized. Thanks to the decomposition of the BW distance given in Eq (2), it leads to the following updates on the mean and variances:

$$m_{k+1} = m_k - \gamma \nabla_m F(m_k, \epsilon_k) \quad (5)$$

$$\epsilon_{k+1} = \left(1 - \frac{2\gamma}{d} \nabla_\epsilon F(m_k, \epsilon_k)\right)^2 \epsilon_k, \quad (6)$$

where  $\nabla_m F(m_k, \epsilon_k) = \mathbb{E}_{p_k}[\nabla V]$  and  $\nabla_\epsilon F(m_k, \epsilon_k) = \frac{1}{2} \left( \frac{1}{d\epsilon_k} \mathbb{E}_{p_k}[(\cdot - m_k)^\top \nabla V] - \frac{1}{\epsilon_k} \right)$ . We observe that while the first update on the mean is a simple gradient descent, the latter update ensures that the variance remains positive and differs from a simple Euler discretization of the associated differential equation (see Appendix A.3). We provide the details of the computation as well as an interpretation of these updates as a Riemannian gradient descent in Appendix A.4.

### 2.3 Entropic Mirror Descent on IG

We now turn to an alternative descent scheme on IG, namely a mirror descent one, relying on a different geometry than the  $W_2$  one detailed in the previous subsection. Mirror descent is an optimization algorithm that was introduced to solve constrained convex problems [Nemirovskij and Yudin, 1983], and that uses in the optimization updates a cost (or “geometry”) that is a Bregman divergence [Bregman, 1967], whose definition is given below.

**Definition 2.2.** Let  $\phi : \mathcal{X} \rightarrow \mathbb{R}$  a strictly convex and differentiable functional on a convex set  $\mathcal{X}$ , referred to as a Bregman potential. The  $\phi$ -Bregman divergence is defined for any  $x, y \in \mathcal{X}$  by:

$$B_\phi(y|x) = \phi(y) - \phi(x) - \langle \nabla \phi(x), y - x \rangle.$$

The reader may refer to Appendix A.5 for more background on mirror descent and its relationship with more standard algorithms such as gradient descent.

Hence, we propose to choose an appropriate Bregman divergence on the space of covariance matrices, namely a generalized Kullback-Leibler divergence between positive definite matrices:  $\overline{\text{KL}}(A|B) = \text{Tr}(A(\log A - \log B)) - \text{Tr}(A) + \text{Tr}(B)$ . The latter object, also called Von Neumann relative entropy, is a Bregman divergence whose Bregman potential is the Von Neumann entropy  $\phi : A \mapsto \text{Tr}(A \log A)$  (and where  $\langle A, B \rangle = \text{Tr}(AB)$ ). Note that  $\overline{\text{KL}}(\epsilon \mathbf{I}_d | \tau \mathbf{I}_d) = d \left( \epsilon \log \frac{\epsilon}{\tau} - \epsilon + \tau \right)$ . Then, we can define a descent scheme on IG as follows, starting from some  $\eta_0 = (m_0, \epsilon_0) \in \mathbb{R}^d \times \mathbb{R}^{+*}$ :

$$\eta_{k+1} = \underset{\eta \in \mathbb{R}^d \times \mathbb{R}^{+*}}{\text{argmin}} \left\{ \langle \nabla F(\eta_k), \eta - \eta_k \rangle + \frac{1}{2\gamma} \|m - m_k\|^2 + \frac{1}{2\gamma} \overline{\text{KL}}(\epsilon \mathbf{I}_d | \epsilon_k \mathbf{I}_d) \right\}, \quad (7)$$

denoting again  $\eta = (m, \epsilon)$ . Note that compared to the scheme Eq (4), only the update on the variance differs, and is given by:

$$\epsilon_{k+1} = \epsilon_k \exp \left( -\frac{2\gamma}{d} \nabla_\epsilon F(m_k, \epsilon_k) \right), \quad (8)$$

see Appendix A.6 for the computations. This update, as the one in Eq (6), also guarantees that the variance parameter  $\epsilon$  remains strictly positive; and is known as entropic mirror descent [Beck and Teboulle, 2003]<sup>1</sup>.

## 3 Mixtures of Isotropic Gaussians (MIG)

We now turn to the problem of optimizing the KL objective to the target distribution  $\pi$  over the family of mixtures of isotropic Gaussians. We will consider the VI problem Eq (1) for a specific setting where the variational family is the set of mixtures of  $N$  isotropic Gaussians, for some  $N \in \mathbb{N}^*$ , with equally weighted components:

$$\mathcal{C}^N = \left\{ \frac{1}{N} \sum_{j=1}^N \mathcal{N}(m^j, \epsilon^j \mathbf{I}_d), [m^j, \epsilon^j]_{j=1}^N \in (\mathbb{R}^d \times \mathbb{R}^{+*})^N \right\}.$$

<sup>1</sup>[Beck and Teboulle, 2003] used this exponential update followed by a renormalization to optimize over the simplex.

166 Note that any distribution  $\nu \in \mathcal{C}^N$  writes  $\nu = \frac{1}{N} \sum_{j=1}^N k_{\epsilon^j} \star \delta_{m^j}$ , for some  $[m^j, \epsilon^j]_{j=1}^N \in (\mathbb{R}^d \times$   
 167  $\mathbb{R}^{+*})^N$ , where  $k_{\epsilon^j}$  is the Gaussian kernel with variance  $\epsilon^j$ . Then, we define our loss function  
 168  $F : (\mathbb{R}^d)^N \times (\mathbb{R}^+)^N \rightarrow \mathbb{R}^{+*}$  as:

$$F([m^j, \epsilon^j]_{j=1}^N) := \text{KL} \left( \frac{1}{N} \sum_{j=1}^N \mathcal{N}(m^j, \epsilon^j \text{I}_d) \middle| \pi \right). \quad (9)$$

169 The following proposition provides useful formulas regarding the gradients of this objective.

170 **Proposition 3.1.** *Let  $\mu = \frac{1}{N} \sum_{j=1}^N \delta_{m^j}$  and denote  $k_{\epsilon} \otimes \mu = \frac{1}{N} \sum_{j=1}^N k_{\epsilon^j} \star \delta_{m^j}$ . Assume  $\pi \in \mathcal{C}^1(\mathbb{R}^d)$ ,*  
 171 *the gradients of  $F$  with respect to  $m^j, \epsilon^j \in \mathbb{R}^d \times \mathbb{R}^{+*}$  write:*

$$\begin{aligned} \nabla_{\epsilon^j} F([m^j, \epsilon^j]_{j=1}^N) &= \frac{1}{2N\epsilon^j} \mathbb{E}_{k_{\epsilon^j} \star \delta_{m^j}} \left[ (\cdot - m^j)^T \nabla \ln \left( \frac{k_{\epsilon} \otimes \mu}{\pi} \right) (\cdot) \right], \\ \nabla_{m^j} F([m^j, \epsilon^j]_{j=1}^N) &= \frac{1}{N} \mathbb{E}_{k_{\epsilon^j} \star \delta_{m^j}} \left[ \nabla \ln \left( \frac{k_{\epsilon} \otimes \mu}{\pi} \right) (\cdot) \right]. \end{aligned}$$

172 The proof of Proposition 3.1 can be found in Appendix A.7. Note that the means and variances in  
 173 the mixture interact through the terms  $\nabla \ln(k_{\epsilon} \otimes \mu)$  in the gradients. Remarkably, our computations  
 174 provide an expression of the gradient on the variance that only involves a scalar product  $(\cdot -$   
 175  $m^j)^T \nabla \ln \left( \frac{k_{\epsilon} \otimes \mu}{\pi} \right)$ , which can be computed efficiently with a computational cost in  $\mathcal{O}(d)$ . In practice,  
 176 the expectations over Gaussian distributions will be estimated with Monte Carlo integration.

## 177 4 Algorithms for VI on MIG

### 178 4.1 General optimization framework

179 We propose, for the optimization of the objective  $F$  on  $(\mathbb{R}^d)^N \times (\mathbb{R}^{+*})^N$  defined in Eq. (9), or  
 180 equivalently  $\text{KL}(\cdot | \pi)$  on  $\mathcal{C}^N$ , to perform joint optimization on the means and variances of the mixture.  
 181 This joint optimization involves a gradient descent update on the means, and either a Bures or entropic  
 182 mirror descent update on the variances. Our approach is summarized in Algorithm 1.

---

#### Algorithm 1 MIG optimization with IBW or MD

---

**Input:** initial means and variances  $(m_0^j, \epsilon_0^j)_{j=1}^N$ , step-size  $\gamma$ , number of iterations  $T$ .  
**for**  $k = 1$  **to**  $T$  **do**  
  **for**  $i = 1$  **to**  $N$  **do**  
    Update  $m_{k+1}^i = m_k^i - \gamma N \nabla_{m_k^i} F([m_k^j, \epsilon_k^j]_{j=1}^N)$  (GD)  
    Update  $\epsilon_k^i$  with IBW (Eq. 14) or MD (Eq. 16)  
  **end for**  
**end for**

---

183 We now turn to the details on the optimization of the variances with either a Bures or entropic mirror  
 184 descent update in the next subsections. In a nutshell, they rely on careful adaptations of the schemes  
 185 detailed in Section 2, stated for unique isotropic Gaussians, to a mixture of these distributions.

### 186 4.2 Bures (IBW) update

187 This section extends the results of Section 2.2 to Gaussian mixtures and presents a new formulation  
 188 of the JKO scheme adapted to this setting. The  $W_2$  distance between two Gaussian mixtures is  
 189 intractable and does not admit a closed form, in contrast with the BW distance on IG. Then, we  
 190 cannot obtain direct updates on the mean and variances for the projected JKO on  $\mathcal{C}^N$ . To address  
 191 this issue, we will represent a mixture with its associated mixing measure. Lambert et al. [2022]  
 192 proposed a similar approach to derive a Wasserstein gradient flow on Gaussian mixtures. However,  
 193 they considered a mixing measure with an infinite number of components and did not provide a  
 194 formal derivation of a fully explicit discrete (in time and space) scheme.

Any uniform-weight Gaussian mixture  $\nu \in \mathcal{C}^N$  can be identified to a mixing measure  $\hat{p} = \frac{1}{N} \sum_{j=1}^N \delta_{(m^j, \epsilon^j)}$  on  $(\mathbb{R}^d \times \mathbb{R}^{+*})^N$  where for any  $x \in \mathbb{R}^d$  we have:

$$\nu(x) = \int \mathcal{N}(x; m, \epsilon \mathbf{I}_d) d\hat{p}(m, \epsilon) \quad (10)$$

Note that the space  $\mathcal{C}^N$  allows for the full identification of a mixture, up to a reordering of the indices, since the corresponding mixing measure contains only  $N$  particles with equal weights. This avoids the identifiability issues that arise when dealing with a mixing measure supported on a continuous (infinite) set of components, which constitutes an overparameterized model, see [Chewi et al., 2024, Section 5.6]. See also Appendix B.1 for details. Following Chen et al. [2019], we first consider a Wasserstein distance between mixing measures denoted  $W_{bw}$ , where the cost is a squared Bures-Wasserstein distance, i.e.  $c((m, \epsilon), (m', \tau)) = \text{BW}^2(\mathcal{N}(m, \epsilon \mathbf{I}_d), \mathcal{N}(m', \tau \mathbf{I}_d))$ . We then construct the JKO scheme on mixing measures at each step  $k$  as:

$$\hat{p}_{k+1} = \arg \min_{(m^j, \epsilon^j)_{j=1}^N} \left\{ \text{KL}(\nu | \pi) + \frac{1}{2\gamma} W_{bw}^2(\hat{p}, \hat{p}_k) \right\}, \quad (11)$$

where  $\hat{p} = \frac{1}{N} \sum_{j=1}^N \delta_{(m^j, \epsilon^j)}$ ,  $\hat{p}_k = \frac{1}{N} \sum_{j=1}^N \delta_{(m_k^j, \epsilon_k^j)}$  and  $\nu$  is defined as in Eq (10). The resulting Gaussian mixture is  $\nu_{k+1} = \int \mathcal{N}(m, \epsilon \mathbf{I}_d) d\hat{p}_{k+1}$ . Since  $\hat{p}$  and  $\hat{p}_k$  are two discrete measures with an equal number of components, the above Wasserstein distance simplifies as:

$$W_{bw}^2(\hat{p}, \hat{p}_k) = \min_{\sigma} \frac{1}{N} \sum_{j=1}^N \text{BW}^2(\mathcal{N}(m^j, \epsilon^j \mathbf{I}_d), \mathcal{N}(m_k^{\sigma(j)}, \epsilon_k^{\sigma(j)} \mathbf{I}_d)), \quad (12)$$

where  $\sigma$  is a permutation of indices in the mixture. Solving the JKO scheme Eq (11) is now tractable and we can compute the limiting flow as  $\gamma \rightarrow 0$ , since at the limit  $\sigma(i) = i$ . The continuous-time equations of the flow in the isotropic case are given in Appendix B.2. They match the continuous-time equations for the means and covariances derived in [Lambert et al., 2022, Section 5.2] and recalled in Appendix E.1.

Similarly to Section 2.2, we consider an explicit time-discretization of this flow, using a linearization of the objective in Eq (11). This leads us to the scheme:

$$[\eta_{k+1}]_{j=1}^N = \arg \min_{(\eta^j)_{j=1}^N} \left\{ \langle \nabla F([\eta_k^j]_{j=1}^N), [\eta^j]_{j=1}^N - [\eta_k^j]_{j=1}^N \rangle + \frac{1}{2\gamma N} \sum_{j=1}^N \text{BW}^2(\mathcal{N}(\eta^j), \mathcal{N}(\eta_k^j)) \right\}, \quad (13)$$

assuming that  $\sigma(i) = i$  in Eq (12) for  $\gamma$  small enough. Finally, the variance updates for the Gaussian components are:

**IBW update** For  $j = 1, \dots, N$ :  $\epsilon_{k+1}^j = \left( 1 - \frac{2N\gamma}{d} \nabla_{\epsilon^j} F([m_k^j, \epsilon_k^j]) \right)^2 \epsilon_k^j$ .

(14)

The update on the means takes the form of Eq (GD). Details of the computations are deferred to Appendix B.3. Ultimately, we obtain a system of Gaussian particles  $(m^j, \epsilon^j)_{j=1}^N$  that interact through the gradient of the objective.

*Remark 4.1.* Note that we can quantify the discrepancy between our JKO scheme on the mixing measure and the original JKO scheme projected onto Gaussian mixtures. Indeed, the Wasserstein distance on the mixing measure is related to the  $W_2$  distance on  $\mathcal{P}_2(\mathbb{R}^d)$  as follows:  $0 \leq W_{bw}^2(\hat{p}, \hat{p}_k) - W_2^2(\nu, \nu_k) \leq 2\sqrt{2d}\epsilon^*$  where  $\epsilon^*$  is the maximal variance of the mixtures  $\nu, \nu_k$ . This result is a direct consequence of [Delon and Desolneux, 2020, Proposition 6], see Appendix B.4 for further details.

### 4.3 Entropic mirror descent (MD) update

In this section we provide an alternative way to optimize the variances of the mixture, based on mirror descent ideas introduced in Section 2.3. In particular, generalizing to  $N$  components what we have

done for Eq (8), and by analogy with the scheme Eq (13), we consider:

$$[\eta_{k+1}^j]_{j=1}^N = \arg \min_{(\eta^j)_{j=1}^N} \left\{ \langle \nabla F([\eta_k^j]), [\eta^j] - [\eta_k^j] \rangle + \frac{1}{2\gamma N} \sum_{j=1}^N \|m^j - m_k^j\|^2 + \overline{\text{KL}}(\epsilon^j \text{I}_d | \epsilon_k^j \text{I}_d) \right\}. \quad (15)$$

Then, at step  $k \geq 0$ , the update on the variances takes the form:

**MD update** For  $j = 1, \dots, N$  :  $\epsilon_{k+1}^j = \epsilon_k^j \cdot \exp\left(-\frac{2N\gamma}{d} \nabla_{\epsilon^j} F([m_k^j, \epsilon_k^j])\right)$ ,

(16)

while the update on the means remains Eq (GD), see Appendix A.6 for the detailed computations.

## 5 Related Work

In this section, we provide an overview of relevant work on VI with mixtures of Gaussians.

Several works have addressed VI for mixture models, emphasizing computational aspects. Gershman et al. [2012] optimize an approximate ELBO using L-BFGS (a quasi-Newton method), relying on successive approximations of ELBO terms for mixtures of Gaussians. However, their optimization objective diverges significantly from the original KL objective in VI, which is a valid divergence between probability distributions.

Lin et al. [2019], Arenz et al. [2018] propose natural gradient (NGD) updates on the natural parameters of the Gaussians for each component of the mixture, and the categorical distribution over weights. These methods are unified and extended in [Arenz et al., 2023], which introduces computational improvements. Natural gradient descent differs from standard gradient descent by performing steepest descent with respect to changes in the underlying distribution, measured using the Fisher information metric. In other words, the natural gradient is the standard gradient preconditioned by the inverse of the Fisher information matrix. For exponential families, such as Gaussians, the natural gradient of the objective with respect to natural parameters coincides with the standard gradient of the (reparametrized) objective when expressed in terms of expectation parameters (i.e., the moments of the Gaussians). This has some pleasant consequences, including closed-form updates on means and covariances, since the natural parameter admits a simple expression in terms of means and variances for Gaussians. The NGD updates (fixing the weights of the mixture) on means and variances write:

$$\frac{1}{\epsilon_{k+1}} - \frac{1}{\epsilon_k} = \frac{2N\gamma}{d} \nabla_{\epsilon_k} F(m_k, \epsilon_k), \quad m_{k+1} - m_k = -\epsilon_k N \gamma \nabla_{m_k} F(m_k, \epsilon_k). \quad (17)$$

We refer to Appendix C for more details and the computations. *In particular, the latter update does not guarantee that the variances remain positive, and the update on the mean is multiplied by the current covariance.* In contrast, our updates on the means and covariances are decoupled, and our updates on variances enforce positivity.

[Huix et al., 2024] considered the setting where the variances of the mixture are shared, equal to  $\epsilon \text{I}_d$  with  $\epsilon \in \mathbb{R}^{+*}$  that is kept fixed, and only the means  $(m^1, \dots, m^N)$  are optimized with gradient descent. In that setting, they proved a descent lemma showing that the KL objective functional decreases along the gradient descent iterations, under some conditions including a maximal step-size, a boundedness conditions on the second moment of the (distribution) iterations and a finite number of components  $N$ . They also showed a bound on the approximation error, i.e. the KL objective between a  $N$ -component mixture of Gaussians with constant weight and covariance is upper bounded as  $\mathcal{O}(\frac{\log(N)}{N})$ , when  $\pi$  writes as an infinite mixture of these components. Interestingly, this bound is valid for mixtures of isotropic Gaussians with different variances, as we show in Appendix D. Yet, fixing the covariances to a constant factor  $\epsilon$  of the identity limits a lot the expressiveness of the variational family. Hence, we focus on the more general family defined by  $\mathcal{C}^N$ .

## 6 Experiments

In this section, we present some experiments to evaluate our proposed methods summarized in Algorithm 1, namely IBW and MD. In practice, in all our experiments, the gradients in Proposition 3.1

are computed with a Monte Carlo approximation involving  $B$  samples from Gaussian distributions in the components. Our other hyperparameters are the following:  $N$  the number of components in the mixture;  $\gamma$  the step-size; and we also vary the initial values of the means and covariances of the candidate mixture. We will specify these specific values for each experiment in Section 6. Our code will be publicly available. We compare our methods with different algorithms. The algorithm presented in Lambert et al. [2022], based on a Bures geometry but updating full covariances matrix in the mixture, is abbreviated as BW (see Appendix E.1). Note that the latter is of complexity  $\mathcal{O}N(d + d^2)$  instead of  $N(d + 1)$  for our schemes IBW and MD. The algorithm updating only the means with Eq (GD), and not the variances, as proposed in [Huix et al., 2024] is denoted GD. The natural gradient descent algorithm of Lin et al. [2019], adapted to isotropic Gaussians with fixed weights and given in Eq (17), is referred to as NGD. We also evaluate Normalizing Flow (NF). We also used Hamiltonian Monte Carlo (HMC), an MCMC scheme (hence non parametric) and Automatic Differentiation Variational Inference (ADVI) Kucukelbir et al. [2016] relying on mean-field VI, using `stan` framework on specific experiments.

**Gaussian-mixture target in two dimensions.** We first illustrate the behavior of our methods for a two-dimensional target  $\pi$  that is a Gaussians mixture with 5 components. In Figure 1 (top), we evaluate our methods for  $N = 1, 5, 10, 20$ . We observe that  $N = 20$  leads to the lower KL divergence - in accordance with the fact that the approximation error with a mixture of isotropic Gaussians goes to zero as  $N$  tends to infinity, see Section 5. We plot the approximated density with  $N = 10$  for BW, IBW, GD and NF methods together with the target density in Figure 1 (bottom). In this experiment we also evaluate normalizing flow with  $b = 5$  blocks and  $h = 124$  hidden units. The NF method appears slower than ours without improving the approximation. We also evaluated these VI methods on alternative target distributions in two dimensions (other mixtures of Gaussians, Funnel and heavier tails distribution). The results and the details of experiments are deferred to Appendix E.3.

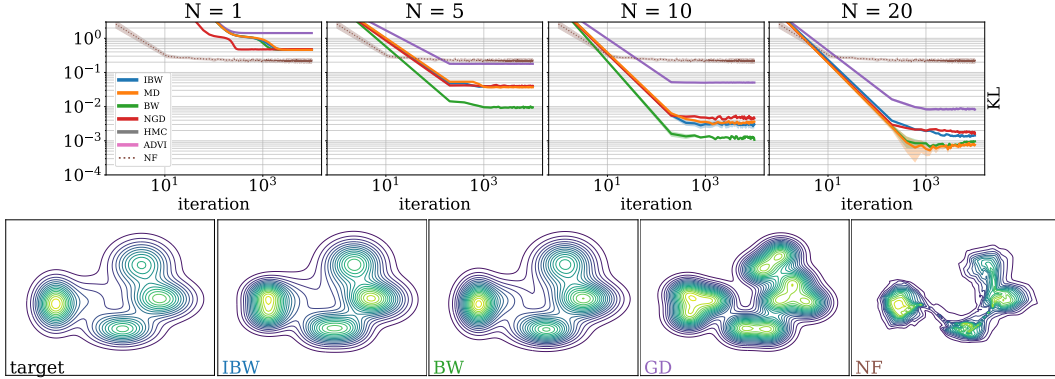


Figure 1: Illustration of convergence of Algorithm 1 for a two-dimensional target distribution.

**Gaussian-mixture target in high dimension.** We then consider a target Gaussian mixture with 10 components in dimension  $d = 20$ , and a variational mixture model with  $N = 20$ . In Figure 3, we plotted the marginals along each dimension together with the KL evolution for the schemes (BW, IBW, MD, NGD). We observe that both our schemes IBW and MD provide a good approximation, along with BW. In the mean update of NGD, the gradient is rescaled by  $\epsilon$ , which leads to a large step for a spread-out Gaussian. While this rescaling allows faster convergence, we observe that it makes the algorithm more sensitive to the initial conditions. Then, we compare in Figure 2 the time per iteration and the KL objective value for BW, that updates full covariances matrices, and our isotropic version (IBW) for a similar target over several dimensions, for  $N = 15$ . We note that IBW performs comparably to BW, while enjoying a faster time execution, and still with a cost in memory linear in the dimension instead of quadratic.

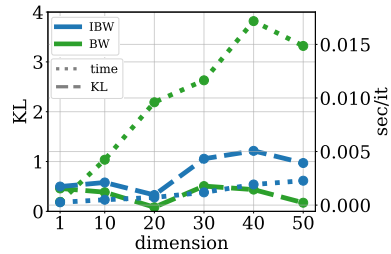


Figure 2: Time vs KL evolution per dimension.

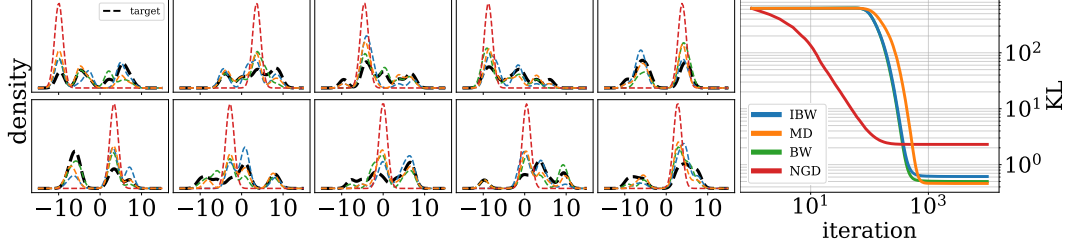


Figure 3: (First 10) Marginals (left) and KL objective (right) for MD, IBW, BW and NGD ( $d = 20$ ).

**Bayesian posteriors.** We evaluate our methods on two probabilistic inference tasks using classical UCI datasets. The first one is Bayesian logistic regression (BLR) for `breast_cancer` (2 labels,  $d = 30$ ) and `wine` (3 labels,  $d = 39$ ). The second one aims to compute a Bayesian neural network (BNN) posterior on a regression task on the `boston` dataset using a single hidden layer neural network of 50 units ( $d = 601$ ), and on the MNIST with a one layer neural network with 256 units ( $d = 203530$ ). In each case, we assume a Gaussian prior, and compute the posterior distribution given observations, more details are given in Appendix E.5. Note that the first task leads to log-concave posteriors (in contrast to the previous mixture of Gaussians targets which are typically non log-concave) while the second leads to a multimodal one.

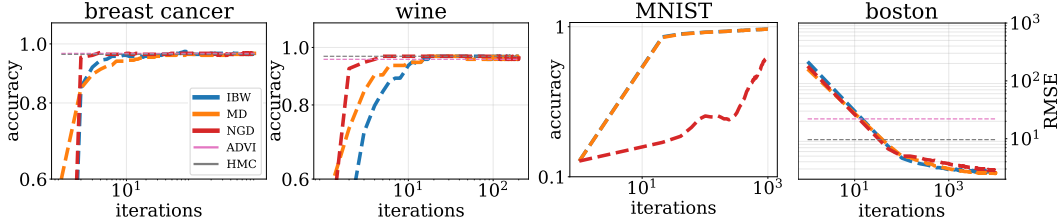


Figure 4: Bayesian logistic regression and BNN regression approximated by mixtures of Gaussians.

We performed optimization with IBW and MD algorithms and plot the accuracy or Root Mean Squared Error (RMSE) score on the test set over iterations, our results are presented Figure 4. The evolution of log-likelihood and unnormalized KL are deferred in E.5. As additional baselines, we compared our results with HMC and ADVI methods, and only plotted the final samples results provided by `stan`. Regarding MNIST, for scalability issues, in analogy with Blundell et al. [2015], we coded a mean-field version of our algorithm where each weight marginal is fitted with our Algorithm 1. More details are provided in Appendix E.5. We observed that we achieved the same order of performance for all algorithms. We also performed the optimization with NGD updates but found out that if the means were not initialized within a very small ball (increasing the chances of missing some modes in the target), the variance  $\varepsilon$  estimated by NGD can become negative, whereas our scheme guarantees the positivity of  $\varepsilon$  by construction.

**Conclusion.** Mixtures of isotropic Gaussians provide a simplified yet powerful tool in variational inference, balancing expressivity for multimodal target distributions with computational and memory efficiency. We presented two optimization schemes, that implement joint optimization on the means through gradient descent, and on the variances through adapted geometries for the space of variance matrices, such as the Bures or Von Neumann entropy ones, that guarantee that they remain positive. Our numerical experiments validate their relevance for different types of target posterior distributions. Future work include establishing more theoretical guarantees regarding our schemes and mixtures of isotropic Gaussians. For instance, comparing the approximation error of full covariance mixtures versus isotropic ones, would be helpful to understand why we observe empirically a great computational cost gain for a very modest increase of the KL loss (eg IBW vs BW). Then, studying optimization guarantees for these schemes is of interest. For instance, when the variance of the gradients is controlled, it is possible to establish descent results guaranteeing that the objective sufficiently decreases at each step. Establishing this, along with practical ways to control the variance of the stochastic gradients such as advanced variance reductions techniques, is of interest especially in high dimensions.

## References

- Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient flows: in metric spaces and in the space of probability measures*. Springer Science & Business Media, 2008.
- Oleg Arenz, Gerhard Neumann, and Mingjun Zhong. Efficient gradient-free variational inference using policy search. In *International conference on machine learning*, pages 234–243. PMLR, 2018.
- Oleg Arenz, Philipp Dahlinger, Zihan Ye, Michael Volpp, and Gerhard Neumann. A unified perspective on natural gradient variational inference with Gaussian mixture models. *Transactions of Machine Learning Research*, 2023.
- Manuel Arnese and Daniel Lacker. Convergence of coordinate ascent variational inference for log-concave measures via optimal transport. *arXiv preprint arXiv:2404.08792*, 2024.
- Pierre-Cyril Aubin-Frankowski, Anna Korba, and Flavien Léger. Mirror descent with relative smoothness in measure spaces, with application to Sinkhorn and EM. *Advances in Neural Information Processing Systems*, 35:17263–17275, 2022.
- David Barber and Christopher Bishop. Ensemble learning for multi-layer networks. In M. Jordan, M. Kearns, and S. Solla, editors, *Advances in Neural Information Processing Systems*, volume 10. MIT Press, 1997.
- Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- Rajendra Bhatia, Tanvi Jain, and Yongdo Lim. On the Bures-Wasserstein distance between positive definite matrices. *Expositiones Mathematicae*, 37(2):165–191, 2019.
- Christopher M. Bishop. *Pattern recognition and machine learning*. Information Science and Statistics. Springer, New York, 2006.
- David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518):859–877, 2017.
- Denis Blessing, Xiaogang Jia, Johannes Esslinger, Francisco Vargas, and Gerhard Neumann. Beyond ELBOs: A large-scale evaluation of variational methods for sampling. *arXiv preprint arXiv:2406.07423*, 2024.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International conference on machine learning*, pages 1613–1622. PMLR, 2015.
- Silvère Bonnabel, Marc Lambert, and Francis Bach. Low-rank plus diagonal approximations for Riccati-like matrix differential equations. *SIAM Journal on Matrix Analysis and Applications*, pages 1669–1688, 2024.
- Lev M Bregman. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR computational mathematics and mathematical physics*, 7(3):200–217, 1967.
- Diana Cai, Chirag Modi, Charles C. Margossian, Robert M. Gower, David M. Blei, and Lawrence K. Saul. Eigenvi: score-based variational inference with orthogonal function expansions, 2024a. URL <https://arxiv.org/abs/2410.24054>.
- Diana Cai, Chirag Modi, Loucas Pillaud-Vivien, Charles C Margossian, Robert M Gower, David M Blei, and Lawrence K Saul. Batch and match: black-box variational inference with a score-based divergence. *International Conference on Machine Learning*, 2024b.
- Edward Challis and David Barber. Gaussian Kullback-Leibler approximate inference. *Journal of Machine Learning Research*, 14(68):2239–2286, 2013.
- Yongxin Chen, Tryphon T. Georgiou, and Allen Tannenbaum. Optimal transport for Gaussian mixture models. *IEEE Access*, 7:6269–6278, 2019.

390 Sinho Chewi, Jonathan Niles-Weed, and Philippe Rigollet. Statistical optimal transport, 2024.

391 Julie Delon and Agnès Desolneux. A Wasserstein-type distance in the space of Gaussian mixture  
392 models. *SIAM Journal on Imaging Sciences*, 13(2):936–970, 2020.

393 Michael Ziyang Diao, Krishna Balasubramanian, Sinho Chewi, and Adil Salim. Forward-backward  
394 Gaussian variational inference via JKO in the Bures-Wasserstein space. In *International Conference  
395 on Machine Learning*, pages 7960–7991. PMLR, 2023.

396 Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp, 2017.  
397 URL <https://arxiv.org/abs/1605.08803>.

398 Justin Domke, Guillaume Garrigos, and Robert Gower. Provable convergence guarantees for black-  
399 box variational inference. *Advances in neural information processing systems*, 2023.

400 Samuel Gershman, Matt Hoffman, and David Blei. Nonparametric variational inference. *International  
401 Conference of Machine Learning*, 2012.

402 Tom Huix, Anna Korba, Alain Durmus, and Eric Moulines. Theoretical guarantees for variational  
403 inference with fixed-variance mixture of Gaussians. *International Conference of Machine Learning*,  
404 2024.

405 Yiheng Jiang, Sinho Chewi, and Aram-Alexandre Pooladian. Algorithms for mean-field variational  
406 inference via polyhedral optimization in the wasserstein space. In *The Thirty Seventh Annual  
407 Conference on Learning Theory*, pages 2720–2721. PMLR, 2024.

408 Richard Jordan, David Kinderlehrer, and Felix Otto. The variational formulation of the Fokker–Planck  
409 equation. *SIAM Journal on Mathematical Analysis*, 29(1):1–17, 1998.

410 Ivan Kobyzev, Simon JD Prince, and Marcus A Brubaker. Normalizing flows: An introduction and  
411 review of current methods. *IEEE transactions on pattern analysis and machine intelligence*, 43  
412 (11):3964–3979, 2020.

413 Alp Kucukelbir, Dustin Tran, Rajesh Ranganath, Andrew Gelman, and David M. Blei. Automatic  
414 differentiation variational inference, 2016. URL <https://arxiv.org/abs/1603.00788>.

415 Daniel Lacker. Independent projections of diffusions: Gradient flows for variational inference and  
416 optimal mean field approximations. *arXiv preprint arXiv:2309.13332*, 2023.

417 Marc Lambert, Sinho Chewi, Francis Bach, Silvére Bonnabel, and Philippe Rigollet. Variational  
418 inference via Wasserstein gradient flows. *Advances in Neural Information Processing Systems*, 35:  
419 14434–14447, 2022.

420 Hugo Lavenant and Giacomo Zanella. Convergence rate of random scan coordinate ascent variational  
421 inference under log-concavity. *SIAM Journal on Optimization*, 34(4):3750–3761, 2024.

422 Wu Lin, Mohammad Emtiyaz Khan, and Mark Schmidt. Fast and simple natural-gradient variational  
423 inference with mixture of exponential-family approximations. In *International Conference on  
424 Machine Learning*, pages 3992–4002. PMLR, 2019.

425 Chirag Modi, Diana Cai, and Lawrence K Saul. Batch, match, and patch: low-rank approximations  
426 for score-based variational inference. *Artificial intelligence and statistics*, 2025.

427 Radford M. Neal. Slice sampling, 2000. URL <https://arxiv.org/abs/physics/0009028>.

428 Arkadij Semenovič Nemirovskij and David Borisovich Yudin. Problem complexity and method  
429 efficiency in optimization. 1983.

430 Manfred Opper and Cédric Archambeau. The variational Gaussian approximation revisited. *Neural  
431 Comput.*, 21(3):786–792, March 2009. ISSN 0899-7667.

432 Felix Otto. The geometry of dissipative evolution equations: the porous medium equation. *Communi-  
433 cations in Partial Differential Equations*, 26(1-2):101–174, 2001.

434 George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji  
435 Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of*  
436 *Machine Learning Research*, 22(57):1–64, 2021.

437 Arthur Pewsey. Sinh-arcsinh distributions. *Biometrika*, 96, 11 2009. doi: 10.1093/biomet/asp053.

438 Garvesh Raskutti and Sayan Mukherjee. The information geometry of mirror descent. *IEEE*  
439 *Transactions on Information Theory*, 61(3):1451–1457, 2015.

440 Gareth O Roberts and Jeffrey S Rosenthal. General state space Markov chains and MCMC algorithms.  
441 *Probability Surveys*, 1, 2004.

442 Adil Salim, Anna Korba, and Giulia Luise. The Wasserstein proximal gradient algorithm. *Advances*  
443 *in Neural Information Processing Systems*, 33:12356–12366, 2020.

444 Roman Soletskyi, Marylou Gabri  , and Bruno Loureiro. A theoretical perspective on mode col-  
445 lapse in variational inference, October 2024. URL <http://arxiv.org/abs/2410.13300>.  
446 arXiv:2410.13300 [stat].

447 Charles M. Stein. Estimation of the mean of a multivariate normal distribution. *The Annals of*  
448 *Statistics*, 9(6):1135 – 1151, 1981. doi: 10.1214/aos/1176345632.

449 Esteban G Tabak and Eric Vanden-Eijnden. Density estimation by dual ascent of the log-likelihood.  
450 2010.

451 C  dric Villani. *Optimal transport: old and new*, volume 338. Springer, 2009.

452 Rentian Yao and Yun Yang. Mean field variational inference via Wasserstein gradient flow. *arXiv*  
453 *preprint arXiv:2207.08074*, 2022.

454 Rentian Yao, Xiaohui Chen, and Yun Yang. Wasserstein proximal coordinate gradient algorithms.  
455 *arXiv preprint arXiv:2405.04628*, 2024.

## A Appendix

### A.1 Geometric structure of Isotropic Gaussians space

We prove here the stability of the isotropic Gaussian model along the Bures-Wasserstein geodesics. We recall that the Bures-Wasserstein space is the space of non-degenerate Gaussian distributions equipped with the Wasserstein-2 metric.

**Proposition A.1.** *The space of isotropic Gaussians equipped with the Bures-Wasserstein distance is a geodesically convex subset of the Bures-Wasserstein space, which is itself a geodesically convex subspace of the Wasserstein space  $(\mathcal{P}_2(\mathbb{R}^d), W_2)$ .*

*Proof.* Let  $p = \mathcal{N}(m_p, \epsilon_p \text{Id})$  and  $q = \mathcal{N}(m_q, \epsilon_q \text{Id})$  be two isotropic Gaussian distributions. Since  $p$  is absolutely continuous, the Wasserstein-2 geodesic, (which is also a Bures-Wasserstein geodesic) between  $p$  and  $q$  is given by the pushforward measure:

$$\mu_t = ((1-t)\text{Id} + tT)_{\#} p, \quad t \in [0, 1],$$

where  $\text{Id}$  is the identity map and  $T$  is the optimal transport map between  $p$  and  $q$ . Denoting  $\Sigma_p = \epsilon_p \text{Id}$  and  $\Sigma_q = \epsilon_q \text{Id}$ , this optimal transport map  $T$  is the affine map:

$$T(x) = m_q + A(x - m_p), \text{ where } A = \Sigma_p^{-1/2} \left( \Sigma_p^{1/2} \Sigma_q \Sigma_p^{1/2} \right)^{1/2} \Sigma_p^{-1/2}.$$

$A$  is called the Bures map and satisfies  $\Sigma_q = A \Sigma_p A$ , which can be easily verified from the definition above. Since the map is linear, it preserves densities, ensuring that the transported measure  $\mu_t$  remains Gaussian. The mean and covariance parameters evolve along the Bures-Wasserstein geodesic between  $p$  and  $q$  according to the equations:

$$\begin{aligned} m_t &= (1-t)m_p + tm_q \\ \Sigma_t &= ((1-t)\text{Id} + tA)\Sigma_p((1-t)\text{Id} + tA). \end{aligned}$$

Since both  $p$  and  $q$  are isotropic, the transport map is:

$$T(x) = m_q + a(x - m_p), \text{ where } a = \left( \frac{\epsilon_q}{\epsilon_p} \right)^{1/2}.$$

Then, the covariance  $\Sigma_t$  evolves according to:

$$\Sigma_t = ((1-t)\text{Id} + ta\text{Id})\epsilon_p((1-t)\text{Id} + ta\text{Id}) = ((1-t) + ta)^2 \epsilon_p \text{Id},$$

which is clearly isotropic. Hence, the interpolated distribution  $\mu_t$  remains an isotropic Gaussian for all  $t \in [0, 1]$ . Thus, the space of isotropic Gaussian distributions is geodesically convex in the Bures-Wasserstein space. Since the Bures-Wasserstein space is itself a geodesically convex subspace of the Wasserstein space, we complete the proof.  $\square$

### A.2 Proof of Theorem 2.1

#### A.2.1 Background on Wasserstein gradient flows

Let  $\mathcal{F} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R} \cup \{+\infty\}$  be a functional. We say that  $\mathcal{F}$  is  $\alpha$ -convex along Wasserstein-2 geodesics if for any two probability measures  $\mu_0, \mu_1 \in \mathcal{P}_2(\mathbb{R}^d)$  and any geodesic  $\{\mu_t\}_{t \in [0, 1]}$  in the Wasserstein-2 space connecting  $\mu_0$  and  $\mu_1$ , we have

$$\mathcal{F}(\mu_t) \leq (1-t)\mathcal{F}(\mu_0) + t\mathcal{F}(\mu_1) - \frac{\alpha}{2}t(1-t)W_2^2(\mu_0, \mu_1), \quad \forall t \in [0, 1].$$

A Wasserstein gradient flow of  $\mathcal{F}$  is a solution  $(\mu_t)_{t \in (0, T)}$ ,  $T > 0$ , of the continuity equation

$$\frac{\partial \mu_t}{\partial t} + \nabla \cdot (\mu_t v_t) = 0$$

that holds in the distributional sense, where  $v_t$  is a subgradient of  $\mathcal{F}$  at  $\mu_t$  [Ambrosio et al., 2008, Definition 10.1.1]. Among the possible processes  $(v_t)_t$ , one has a minimal  $L^2(\mu_t)$  norm and is called the velocity field of  $(\mu_t)_t$ . In a Riemannian interpretation of the Wasserstein space Otto

[2001], this minimality condition can be characterized by  $v_t$  belonging to the tangent space to  $\mathcal{P}_2(\mathbb{R}^d)$  at  $\mu_t$  denoted  $T_{\mu_t}\mathcal{P}_2(\mathbb{R}^d)$ , which is a subset of  $L^2(\mu_t)$ , the Hilbert space of square integrable functions with respect to  $\mu_t$ , whose inner product is denoted  $\langle \cdot, \cdot \rangle_{\mu_t}$ . The Wasserstein gradient is defined as this unique element, and is denoted  $\nabla_{W_2}\mathcal{F}(\mu_t)$ . In particular, if  $\mu \in \mathcal{P}_2(\mathbb{R}^d)$  is absolutely continuous with respect to the Lebesgue measure, with density in  $C^1(\mathbb{R}^d)$  and such that  $\mathcal{F}(\mu) < \infty$ ,  $\nabla_{W_2}\mathcal{F}(\mu)(x) = \nabla\mathcal{F}'(\mu)(x)$  for  $\mu$ -a.e.  $x \in \mathbb{R}^d$  [Ambrosio et al., 2008, Lemma 10.4.1], where  $\mathcal{F}'(\mu)$  denotes the first variation of  $\mathcal{F}$  evaluated at  $\mu$ , i.e. (if it exists) the unique function  $\mathcal{F}'(\mu) : \mathbb{R}^d \rightarrow \mathbb{R}$  s.t.

$$\lim_{h \rightarrow 0} \frac{1}{h} (\mathcal{F}(\mu + h\xi) - \mathcal{F}(\mu)) = \int \mathcal{F}'(\mu)(x) d\xi(x) \quad (18)$$

for all  $\xi = \nu - \mu$ , where  $\nu \in \mathcal{P}_2(\mathbb{R}^d)$ .

## A.2.2 Proof

The proof relies on tools on the Wasserstein geometry and calculus introduced above, and is a direct application of the one of [Lambert et al., 2022, Corollary 3], yet we state it for completeness. When  $\nabla^2 V \succeq \alpha \text{Id}$ ,  $\mathcal{F} : \mu \mapsto \text{KL}(\mu|\pi)$  is  $\alpha$ -convex on  $\mathcal{P}_2(\mathbb{R}^d)$  [Villani, 2009, Theorem 17.15]. Let us denote  $p^*$  the minimum of this function and recall we denote BW the Bures-Wasserstein metric on the manifold (IG, BW). We consider the following gradient flow:

$$\frac{\partial p_t}{\partial t} = \text{div}(p_t \nabla_{W_2}\mathcal{F}(p_t)) \text{ with the initial condition } p_0 = p_0.$$

We first want to show that the solution of this problem is unique. Let  $(p_t)_t$  and  $(q_t)_t$  be two solutions of the above gradient flow. Then, using differential calculus in the Wasserstein space and the chain rule, we have

$$\partial_t \text{BW}^2(p_t, q_t) = 2 \langle \log_{p_t}(q_t), \nabla_{W_2}\mathcal{F}(p_t) \rangle_{p_t} + 2 \langle \log_{q_t}(p_t), \nabla_{W_2}\mathcal{F}(q_t) \rangle_{q_t},$$

where  $\nabla_{W_2}\mathcal{F}(p)$  denotes the Wasserstein gradient at  $p \in \mathcal{P}_2(\mathbb{R}^d)$  and  $\log_{p_t}(q_t) = T - \text{Id} \in L^2(p_t)$ , where  $T$  is the optimal transport map from  $p_t$  to  $q_t$ . Moreover since  $\mathcal{F}$  is  $\alpha$  convex, we can write  $\forall p, q \in \mathcal{P}_2(\mathbb{R}^d)$  :

$$\mathcal{F}(q) \geq \mathcal{F}(p) + \langle \nabla_{W_2}\mathcal{F}(p), \log_p(q) \rangle - \frac{\alpha}{2} \text{BW}^2(p, q).$$

Thus we can write

$$\partial_t \text{BW}^2(p_t, q_t) \leq -2\alpha \text{BW}^2(p_t, q_t).$$

Hence, by Grönwall's lemma, we obtain

$$\text{BW}^2(p_t, q_t) \leq \exp(-2\alpha t) \text{BW}^2(p_0, q_0).$$

Since both  $(p_t)_t$  and  $(q_t)_t$  are solution of the gradient flow,  $p_0 = q_0$  and it implies that  $\forall t \in [0, 1], p_t = q_t$ . This proves the uniqueness of the solution.

Moreover, if  $\alpha > 0$ , we can set  $q_t = p^*$  for all  $t \geq 0$  to deduce exponential contraction of the gradient flow to the minimizer  $p^*$ . Observe that by definition of the gradient flow, we have on the one hand that

$$\partial_t \mathcal{F}(p_t) = \langle \nabla_{W_2}\mathcal{F}(p_t), -\nabla_{W_2}\mathcal{F}(p_t) \rangle_{p_t} = -\|\nabla\mathcal{F}(p_t)\|_{p_t}^2. \quad (19)$$

On the other hand, if  $\alpha > 0$ , the convexity inequality and Young's inequality respectively, yield

$$\begin{aligned} 0 = \mathcal{F}(p^*) &\geq \mathcal{F}(p) + \langle \nabla_{W_2}\mathcal{F}(p), \log_p(p^*) \rangle_p + \frac{\alpha}{2} \text{BW}^2(p, p^*) \\ &\geq \mathcal{F}(p) - \frac{1}{2\alpha} \|\nabla_{W_2}\mathcal{F}(p)\|_p^2 - \frac{\alpha}{2} \underbrace{\|\log_p(p^*)\|_p^2}_{=\text{BW}^2(p, p^*)} + \frac{\alpha}{2} \text{BW}^2(p, p^*) \end{aligned}$$

and hence  $\|\nabla_{W_2}\mathcal{F}(p)\|^2 \geq 2\alpha \mathcal{F}(p)$ . Substituting this into Eq (19) and applying Grönwall's inequality again, we deduce

$$\mathcal{F}(p_t) \leq \exp(-2\alpha t) \mathcal{F}(p_0).$$

### 519 A.3 JKO scheme for isotropic Gaussians

520 In this section, we derive the continuous equations for the isotropic Bures-Wasserstein gradient flow.  
 521 Starting from the JKO scheme of Eq (3), we constrain the solution to lie in the space of isotropic  
 522 Gaussians IG. We follow the same method than [Lambert et al., 2022, Appendix A].

523 Let us recall the JKO scheme Jordan et al. [1998]. Starting from  $p_0 = \mathcal{N}(m_0, \epsilon_0 \mathbf{I}_d)$  at time  $k = 0$   
 524 we look for the solution  $p = \mathcal{N}(m, \epsilon \mathbf{I}_d)$  of:

$$p_{k+1} = \arg \min_{p \in \text{IG}} \left\{ \text{KL}(p|\pi) + \frac{1}{2\gamma} \text{BW}^2(p, p_k) \right\}. \quad (20)$$

The gradient of the KL with respect to the variance parameters is given by Eq (31):

$$\nabla_\epsilon \text{KL}(p|\pi) = -\frac{d}{2\epsilon} - \frac{1}{2} \text{Tr} \mathbb{E}_p[\nabla^2 \log \pi] = -\frac{d}{2\epsilon} + \frac{1}{2} \text{Tr} \mathbb{E}_p[\nabla^2 V].$$

525 Then, for two isotropic Gaussian distributions  $p = \mathcal{N}(m, \epsilon \mathbf{I}_d)$  and  $p_k = \mathcal{N}(m_k, \epsilon_k \mathbf{I}_d)$  we have<sup>2</sup>

$$\text{BW}^2(p, p_k) = \|m - m_k\|_2^2 + d(\epsilon + \epsilon_k - 2\sqrt{\epsilon\epsilon_k}), \quad (21)$$

526 and

$$\nabla_\epsilon \text{BW}^2(p, p_k) = d \left( 1 - \sqrt{\frac{\epsilon_k}{\epsilon}} \right). \quad (22)$$

527 Hence, the first order condition on  $\epsilon$  of (20) yield:

$$\begin{aligned} d \left( 1 - \sqrt{\frac{\epsilon_k}{\epsilon}} \right) = \frac{d\gamma}{\epsilon} - \gamma \text{Tr} \mathbb{E}_p[\nabla^2 V] &\iff \epsilon_k = \epsilon \left( 1 - \frac{\gamma}{\epsilon} + \frac{\gamma}{d} \text{Tr} \mathbb{E}_p[\nabla^2 V] \right)^2 \\ &\iff \epsilon_k = \epsilon \left( 1 - \gamma \left( \frac{1}{\epsilon} - \frac{1}{d\epsilon} \mathbb{E}_p[(\cdot - m)^T \nabla V] \right) \right)^2. \end{aligned} \quad (23)$$

528 Developing equation Eq (23) at first order in  $\gamma$  we obtain:

$$\begin{aligned} \epsilon_k &= \epsilon \left( 1 - \frac{2\gamma}{\epsilon} \left( 1 - \frac{1}{d} \mathbb{E}_p[(\cdot - m)^T \nabla V] \right) + O(\gamma^2) \right) \\ &\iff \frac{\epsilon - \epsilon_k}{\gamma} = 2 - \frac{2}{d} \mathbb{E}_p[(\cdot - m)^T \nabla V] + O(\gamma). \end{aligned}$$

529 Taking the limit  $\gamma \rightarrow 0$  yields the differential equation:

$$\dot{\epsilon} = 2 - \frac{2}{d} \mathbb{E}_p[(\cdot - m)^T \nabla V].$$

530 For the first order condition on the mean parameter  $m$  of (20), using the gradient of the KL w.r.t. the  
 531 mean given Eq (30), we obtain:

$$\mathbb{E}_p \left[ \nabla \log \left( \frac{p}{\pi} \right) \right] + \frac{1}{\gamma} (m - m_k) = 0. \quad (24)$$

Since  $\mathbb{E}_p[\nabla \log p] = 0$  we obtain, at the limit  $\gamma \rightarrow 0$ :

$$\dot{m} = \mathbb{E}_p[\nabla \log \pi] = -\mathbb{E}_p[\nabla V].$$

532 Inspecting Eq (23–24), we observe that solving the JKO scheme yields an implicit discrete-time  
 533 update, where the expectations are evaluated under the unknown distribution  $p$ . We now derive the  
 534 explicit form of this update, starting from the formulation in Eq (4).

---

<sup>2</sup>Recall that for two general Gaussians  $p = \mathcal{N}(m, \Sigma)$ ,  $p_k = \mathcal{N}(m_k, \Sigma_k)$ , we have  $\text{BW}^2(p, p_k) = \|m - m_k\|_2^2 + \text{Tr} \left( \Sigma + \Sigma_k - 2 \left( \Sigma^{1/2} \Sigma_k \Sigma^{1/2} \right)^{1/2} \right)$ .

#### 535 A.4 Forward Euler scheme for isotropic Gaussians

536 **Derivations of the updates.** We now derive the updates given by the scheme Eq (4). The first order  
537 condition on  $\epsilon$  is given by:

$$\frac{1}{2}d \left(1 - \sqrt{\frac{\epsilon_k}{\epsilon}}\right) = -\gamma \nabla_{\epsilon} F(m_k, \epsilon_k) \Leftrightarrow \epsilon = \epsilon_k \left(1 + \frac{2\gamma}{d} \nabla_{\epsilon} F(m_k, \epsilon_k)\right)^{-2}.$$

538 Using the Taylor expansion for small step  $\gamma$  given by  $(1+x)^{-1} = 1 - x + o(x)$  we obtain at first  
539 order the explicit update:

$$\epsilon = \left(1 - \frac{2\gamma}{d} \nabla_{\epsilon} F(m_k, \epsilon_k)\right)^2 \epsilon_k.$$

540 The first-order condition on  $m$  gives the explicit update:

$$m = m_k - \gamma \nabla_m F(m_k, \epsilon_k). \quad (25)$$

541 **Riemannian interpretation.** The latter scheme can be identified to Riemannian gradient descent, on  
542 the isotropic Bures-Wasserstein space, i.e. the space IG of isotropic Gaussians equipped with the  
543 Bures-Wasserstein metric, that we will denote iBW = (IG, BW). To achieve this, we first identify  
544 the local tangent space of the isotropic Bures-Wasserstein space. The direction of the tangent vector  
545 is computed by projecting the Wasserstein-2 gradient of the KL objective onto this tangent space (see  
546 Appendix A.2.1 for the definition). We then follow this projected gradient with a step size  $\gamma$ . We can  
547 then retract back to the isotropic Bures-Wasserstein manifold using an exponential map. We now  
548 detail this approach.

549 Let  $\mathcal{F} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$  a functional. The isotropic Bures-Wasserstein Gradient of  $\mathcal{F}$  at  $p \in \text{IG}$ ,  
550 denoted  $\nabla_{\text{iBW}} \mathcal{F}(p)$ , is the projection of its Wasserstein gradient onto the tangent space to iBW at  $p$ .  
551 If  $p = \mathcal{N}(m_p, \epsilon_p I_d)$ , this tangent space writes:

$$T_p \text{iBW}(\mathbb{R}^d) = \{x \mapsto a + s(x - m_p) \mid a \in \mathbb{R}^d, s \in \mathbb{R}\},$$

552 which can be identified with the pair  $(a, s) \in \mathbb{R}^d \times \mathbb{R}$ . Thus,

$$\nabla_{\text{iBW}} \mathcal{F}(p) = \text{proj}_{T_p \text{iBW}(\mathbb{R}^d)} \nabla_{W_2} \mathcal{F}(p) = \underset{w \in T_p \text{iBW}(\mathbb{R}^d)}{\text{argmin}} \|w - \nabla_{W_2} \mathcal{F}(p)\|_p^2.$$

553 The first conditions in  $a \in \mathbb{R}^d$  and  $s \in \mathbb{R}$  of this problem yield:

$$a = \mathbb{E}_p [\nabla_{W_2} \mathcal{F}(p)] \quad \text{and} \quad s = \frac{1}{d\epsilon_p} \mathbb{E}_p [(x - m_p)^T \nabla_{W_2} \mathcal{F}(p)].$$

554 Indeed,

$$\begin{aligned} \nabla_a \int \|a + s(x - m_p) - \nabla_{W_2} \mathcal{F}(p)(x)\|^2 dp(x) &= 0 \\ \Leftrightarrow \int 2(a + s(x - m_p) - \nabla_{W_2} \mathcal{F}(p)(x)) dp(x) &= 0 \\ \Leftrightarrow a = \mathbb{E}_p [\nabla_{W_2} \mathcal{F}(p)], \end{aligned}$$

555 and

$$\begin{aligned} \nabla_s \int \|a + s(x - m_p) - \nabla_{W_2} \mathcal{F}(p)(x)\|^2 dp(x) &= 0 \\ \Leftrightarrow \int 2(x - m_p)^T (a + s(x - m_p) - \nabla_{W_2} \mathcal{F}(p)(x)) dp(x) &= 0 \\ \Leftrightarrow s \epsilon_p d - \int (x - m_p)^T \nabla_{W_2} \mathcal{F}(p)(x) dp(x) &= 0 \\ \Leftrightarrow s = \frac{1}{d\epsilon_p} \mathbb{E}_p [(x - m_p)^T \nabla_{W_2} \mathcal{F}(p)]. \end{aligned}$$

Back to our problem with  $\mathcal{F}(p) = \text{KL}(p|\pi)$ , using Eq (30) and (31) and that  $\nabla_{\mathbf{W}_2} \mathcal{F}(p) = \nabla \log \left( \frac{p}{\pi} \right)$  (see Appendix A.2.1) we have that

$$a = \nabla_m \text{KL}(p|\pi) = \nabla_m F(m_k, \epsilon_k) \quad \text{and} \quad s = \frac{2}{d} \nabla_\epsilon \text{KL}(p|\pi) = \frac{2}{d} \nabla_\epsilon F(m_k, \epsilon_k).$$

We can follow the direction of the gradient from the tangent space back to the isotropic Bures-Wasserstein manifold using an exponential map which is available in closed form

$$\exp_p(a, s) = \mathcal{N}(m_p + a, (1 + s)^2 \epsilon_p \text{Id}), \quad (26)$$

where we have adapted the exponential map formula for the Bures-Wasserstein space [Lambert et al., 2022, Appendix B.3], written as  $\exp_p(a, S) = \mathcal{N}(m_p + a, (S + I) \Sigma_p (S + I))$ , to the case of isotropic Gaussians. We can then construct a discrete update based as follow: we compute the iBW gradient of  $\mathcal{F}(p) = \text{KL}(p|\pi)$  multiplied a step size  $\gamma$ , and map it back onto the manifold using this exponential map. Starting from  $p_k = \mathcal{N}(m_k, \epsilon_k \text{Id})$ , this discrete time update is:

$$p_{k+1} = \exp_{p_k}(-\gamma \nabla_{\text{iBW}} \mathcal{F}(p_k)) = \exp_{p_k}(-\gamma \nabla_m F(m_k, \epsilon_k), -\frac{2\gamma}{d} \nabla_\epsilon F(m_k, \epsilon_k)),$$

which gives the update on the mean and variance parameters:

$$\begin{aligned} m_{k+1} &= m_k - \gamma \nabla_m F(m_k, \epsilon_k), \\ \epsilon_{k+1} &= (1 - \frac{2\gamma}{d} \nabla_\epsilon F(m_k, \epsilon_k))^2 \epsilon_k. \end{aligned}$$

## A.5 Background on Mirror Descent

Mirror descent is an optimization algorithm that was introduced to solve constrained convex problems [Nemirovskij and Yudin, 1983], and that uses in the optimization updates a cost (or “geometry”) that is a Bregman divergence [Bregman, 1967], whose definition is given below.

**Definition A.2.** Let  $\phi : \mathcal{X} \rightarrow \mathbb{R}$  a strictly convex and differentiable functional on a convex set  $\mathcal{X}$ , referred to as a Bregman potential. The  $\phi$ -Bregman divergence is defined for any  $x, y \in \mathcal{X}$  by:

$$B_\phi(y|x) = \phi(y) - \phi(x) - \langle \nabla \phi(x), y - x \rangle. \quad (27)$$

Let  $\gamma$  a fixed step-size and  $G$  an objective function on  $\mathcal{X}$ . Mirror descent with  $\phi$ -Bregman divergence writes at each step  $k \geq 0$  as:

$$x_{k+1} = \underset{x \in \mathcal{X}}{\text{argmin}} \langle \nabla G(x_k), x - x_k \rangle + \frac{1}{\gamma} B_\phi(x|x_k). \quad (28)$$

Writing the first order conditions of the problem above, mirror descent writes

$$x_{k+1} = \nabla \phi^*(\nabla \phi(x_k) - \gamma \nabla G(x_k)), \quad (29)$$

where  $\phi^*(x) = \sup_{y \in \mathcal{X}} \langle y, x \rangle - \phi(y)$  is the convex conjugate of  $\phi$ , and  $\nabla \phi^* = (\nabla \phi)^{-1}$ . Note that, if  $\mathcal{X}$  is a subset of  $\mathbb{R}^d$  and that  $\phi(x) = \frac{1}{2} \|x\|^2$ , that  $B_\phi(x, y) = \frac{1}{2} \|x - y\|^2$  and mirror descent coincides with gradient descent. Yet, this scheme is more general and is useful for constrained optimization, as if one chooses  $\phi$  wisely, the inverse  $\nabla \phi^*$  of the so-called “mirror map”  $\nabla \phi$  maps the iterates into the domain of  $\phi$ .

## A.6 Entropic mirror descent updates

**Derivations of the updates (Gaussian case).** We now derive the updates given by the scheme Eq (7). The first order condition on  $\epsilon$  is given by:

$$\frac{1}{2} d \log \frac{\epsilon}{\epsilon_k} = -\gamma \nabla_\epsilon F(m_k, \epsilon_k) \Leftrightarrow \epsilon = \epsilon_k \exp \left( -\frac{2\gamma}{d} \nabla_\epsilon F(m_k, \epsilon_k) \right).$$

The first order conditions on the means gives the same explicit update as Eq (4):

$$m = m_k - \gamma \nabla_m F(m_k, \epsilon_k)$$

579 *Remark A.3.* Note that Eq (7) involves a factor  $1/2$  in front of the KL penalization term. Our  
 580 motivation is that in Eq (4), the Bures distance  $\text{BW}^2(\epsilon \text{Id}, \epsilon_k \text{Id}) = d(\sqrt{\epsilon} - \sqrt{\epsilon_k})^2$  is a distance  
 581 between the square root of the matrices and thus its derivative w.r.t  $\epsilon$  leads to a term implying only the  
 582 square root of the matrices, or in our isotropic Gaussian setting the square root of the scale isotropic  
 583 value:  $\nabla_{\epsilon} \text{BW}^2(\epsilon \text{Id}, \epsilon_k \text{Id}) = d(1 + \sqrt{\frac{\epsilon_k}{\epsilon}})$ . We find out it was judicious to have a derivative for the  
 584  $\overline{\text{KL}}$  implying also the square roots of the isotropic coefficient, i.e.,  $\nabla_{\epsilon} \overline{\text{KL}}(\epsilon \text{Id} | \epsilon_k \text{Id}) = \frac{d}{2} \log\left(\frac{\epsilon}{\epsilon_k}\right) =$   
 585  $d \log\left(\sqrt{\frac{\epsilon}{\epsilon_k}}\right)$ .

586 **Mixture Case.** In the mixture of Gaussians case, we derive the following updates from Eq (15). The  
 587 first order condition on  $\epsilon^i$  for  $i = 1, \dots, N$  gives:

$$\frac{1}{2N} d \log \frac{\epsilon^j}{\epsilon_k^i} = -\gamma \nabla_{\epsilon^i} F([m_k^j, \epsilon_k^j]_{j=1}^N) \Leftrightarrow \epsilon^i = \epsilon_k^i \exp\left(-\frac{2N\gamma}{d} \nabla_{\epsilon^i} F([m_k^j, \epsilon_k^j]_{j=1}^N)\right),$$

588 and for the means:

$$m^i = m_k^i - N\gamma \nabla_{m^i} F([m_k^j, \epsilon_k^j]_{j=1}^N).$$

## 589 A.7 Proof of Theorem 3.1

590 In this section, we aim to compute the gradients of the  $\text{KL}(\cdot | \pi)$  objective function defined Eq (9). To  
 591 achieve this, a fundamental tool are Stein's identities Stein [1981] which relate the derivatives with  
 592 respect to the parameters to the derivatives of the integrand.

### 593 A.7.1 Stein's identities for isotropic Gaussians

594 **Lemma A.4.** Let  $p$  be an isotropic Gaussian  $\mathcal{N}(m, \epsilon \text{Id})$  and  $p(x)$  its density for  $x \in \mathbb{R}^d$ . Assume  
 595 that  $\pi \in (C^1(\mathbb{R}^d))$  and  $\lim_{\|x\| \rightarrow \infty} p(x) \log \pi(x) = 0$ . We have

$$\nabla_m \text{KL}(p | \pi) = \mathbb{E}_p \left[ \nabla \log \left( \frac{p}{\pi} \right) \right] \quad (30)$$

$$\nabla_{\epsilon} \text{KL}(p | \pi) = \frac{1}{2\epsilon} \mathbb{E}_p[(\cdot - m)^T \nabla \log \left( \frac{p}{\pi} \right)]. \quad (31)$$

596 Note that the latter equation does not require to compute the Hessian (of dimension  $d \times d$ ), which  
 597 can be computationally more efficient than an alternative formula given in Eq (31).

598 *Proof.* Recall that  $p(x) = (2\pi\epsilon)^{-d/2} \exp\left(-\frac{\|x-m\|^2}{2\epsilon}\right)$  and  $\text{KL}(p | \pi) = \mathbb{E}_p \left[ \log \left( \frac{p}{\pi} \right) \right]$ . First, note  
 599 that for  $\theta = (m, \epsilon)$ :

$$\nabla_{\theta} \int \log p(x | \theta) p(x | \theta) dx = \int \log p(x | \theta) \nabla_{\theta} p(x | \theta) dx, \quad (32)$$

600 which comes from the fact that the expectation of the score function is null, namely  
 601  $\int \nabla_{\theta} \log p(x | \theta) p(x | \theta) dx = 0$ . Hence,

$$\nabla_m \mathbb{E}_p \left[ \log \left( \frac{p}{\pi} \right) \right] = \int \log \left( \frac{p}{\pi}(x) \right) \nabla_m p(x) dx = \int \nabla_x \log \left( \frac{p}{\pi}(x) \right) p(x) dx,$$

602 where we used  $\nabla_m p(x) = -\nabla_x p(x)$  and an integration by parts. We now compute the gradient of  $p$   
 603 with respect to  $\epsilon$ :

$$\begin{aligned} \nabla_{\epsilon} p(x) &= \frac{1}{2} p(x) \left( \frac{1}{\epsilon^2} \|x - m\|^2 - \frac{d}{\epsilon} \right) \\ &= \frac{1}{2} \text{Tr} \left[ p(x) \left( \frac{1}{\epsilon^2} (x - m)(x - m)^T - \frac{1}{\epsilon} \text{Id} \right) \right] = \frac{1}{2} \text{Tr} \nabla_x^2 p(x). \end{aligned} \quad (33)$$

604 Then, we have, using an integration by parts:

$$\begin{aligned} \nabla_{\epsilon} \mathbb{E}_p \left[ \log \left( \frac{p}{\pi} \right) \right] &= \frac{1}{2} \text{Tr} \int \log \left( \frac{p}{\pi} \right) (x) \nabla_x^2 p(x) dx \\ &= -\frac{1}{2} \text{Tr} \int \nabla_x p(x) \nabla_x \log \left( \frac{p}{\pi}(x) \right)^T dx = \frac{1}{2\epsilon} \mathbb{E}_p[(\cdot - m)^T \nabla \log \left( \frac{p}{\pi} \right)]. \end{aligned}$$

605 Note that applying twice an integration by parts would yield another formula:

$$\nabla_{\epsilon} \mathbb{E}_p \left[ \log \left( \frac{p}{\pi} \right) \right] = \frac{1}{2} \text{Tr} \mathbb{E}_p \left[ \nabla^2 \log \left( \frac{p}{\pi} \right) \right] = -\frac{d}{2\epsilon} - \frac{1}{2} \text{Tr} \mathbb{E}_p \left[ \nabla^2 \log \pi \right]. \quad \square$$

## 606 A.7.2 Stein’s identities for mixtures of isotropic Gaussians

607 The application of the previous results for a mixture of isotropic Gaussians is straightforward. Let  
 608  $\mu = \frac{1}{N} \sum_{i=1}^N \delta_{m^i}$ , and recall we denote the associated isotropic Gaussian mixture as  $k_{\epsilon} \otimes \mu =$   
 609  $\frac{1}{N} \sum_{i=1}^N k_{\epsilon^i} \star \delta_{m^i} = \frac{1}{N} \sum_{i=1}^N \mathcal{N}(m^i, \epsilon^i \text{Id})$  where  $k_{\epsilon}(x) = (2\pi\epsilon)^{-d/2} \exp(-\|x\|^2/(2\epsilon))$  is the  
 610 Gaussian kernel.

611 Using the fact that the expectation of the score function is null and that  $\nabla_{m^i} k_{\epsilon} \otimes \mu = \frac{1}{N} \nabla_{m^i} k_{\epsilon^i} \star \delta_{m^i}$ ,  
 612 the gradient of the KL with respect to the mean parameter  $m^i$  is then given by:

$$\nabla_{m^i} \text{KL} \left( k_{\epsilon} \otimes \mu \middle| \pi \right) = \nabla_{m^i} \mathbb{E}_{k_{\epsilon} \otimes \mu} \left[ \ln \left( \frac{k_{\epsilon} \otimes \mu}{\pi} \right) \right] = \frac{1}{N} \mathbb{E}_{k_{\epsilon^i} \star \delta_{m^i}} \left[ \nabla \ln \left( \frac{k_{\epsilon} \otimes \mu}{\pi} \right) \right].$$

613 With the same arguments and using Eq (33), the gradient with respect to the variance parameter  $\epsilon^i$  is:

$$\nabla_{\epsilon^i} \text{KL} \left( k_{\epsilon} \otimes \mu \middle| \pi \right) = \nabla_{\epsilon^i} \mathbb{E}_{k_{\epsilon} \otimes \mu} \left[ \ln \left( \frac{k_{\epsilon} \otimes \mu}{\pi} \right) \right] = \frac{1}{2N\epsilon^i} \mathbb{E}_{k_{\epsilon^i} \star \delta_{m^i}} \left[ (\cdot - m^i)^T \nabla \ln \left( \frac{k_{\epsilon} \otimes \mu}{\pi} \right) \right].$$

614 Equivalently, note that for the latter gradient, using twice an integration by parts would yield the  
 615 formula

$$\nabla_{\epsilon^i} \text{KL} \left( k_{\epsilon} \otimes \mu \middle| \pi \right) = \frac{1}{2N} \text{Tr} \mathbb{E}_{k_{\epsilon^i} \star \delta_{m^i}} \left[ \nabla^2 \ln \left( \frac{k_{\epsilon} \otimes \mu}{\pi} \right) \right].$$

## 616 B Bures-Wasserstein gradient flow for mixtures

617 In this section we consider the Bures-Wasserstein gradient flow for mixtures proposed in Lambert  
 618 et al. [2022] in the particular case of isotropic covariance matrices.

### 619 B.1 Additional discussion

620 Starting from a mixture of isotropic Gaussians  $\nu_0 \in \mathcal{C}^N$  at  $k = 0$ , the JKO scheme where we  
 621 constrain the distribution to be such a mixture writes recursively at subsequent step  $k + 1$ :

$$\hat{\nu}_{k+1} = \arg \min_{\nu \in \mathcal{C}^N} \left\{ \text{KL}(\nu | \pi) + \frac{1}{2\gamma} W_2^2(\nu, \nu_k) \right\}. \quad (34)$$

622 Unfortunately, a closed-form solution for this scheme is not available, as the Wasserstein distance  
 623 between mixtures is not tractable. Regarding the geometry of the space of mixtures of Gaussians,  
 624 Lambert et al. [2022] considered a mixing measure with infinitely many components  $\mu \in \mathcal{P}(\Theta)$ ,  
 625 which can be identified with a Gaussian mixture of the form  $\int \mathcal{N}_{\theta} dp(\theta)$ , where  $\mathcal{N}_{\theta}$  is a Gaussian  
 626 distribution with parameters  $\theta \in \Theta$ . Transport maps can be defined for this model, and gradient flows  
 627 can subsequently be computed [Lambert et al., 2022, Theorem 5]. However, this model is highly  
 628 *overparameterized*. For instance, as illustrated in [Chewi et al., 2024], a single standard Gaussian in  $\mathbb{R}$   
 629 can be represented by infinitely many mixing measures of the form  $\mathcal{N}(0, \text{Id}) = \int \mathcal{N}(x, \tau \text{Id}) dp(x)$ ,  
 630 where  $p = \mathcal{N}(0, (1 - \tau)I)$  for any  $\tau \in [0, 1]$ .

631 An alternative solution is to constrain the mixing model by considering a fixed and finite number of  
 632 components with uniform weights as we do in this work. This resolves the identifiability issue: in  
 633 particular, two co-localized components with weight  $w$  cannot be confused with a single component  
 634 of weight  $2w$ , since all weights are equal. See [Delon and Desolneux, 2020, Proposition 2] for further  
 635 details on identifiability for the mixture model. Moreover discrete measures with the same number  
 636 of atoms are stable along Wasserstein-2 geodesics. This property still holds if we consider discrete  
 637 measures on the isotropic Bures-Wasserstein space  $\hat{p} = \frac{1}{N} \sum_{i=1}^N \delta_{(m^i, \epsilon^i)}$ . We can therefore consider  
 638 discrete mixing measures with uniform weights as stable and identifiable representatives of Gaussian  
 639 mixtures.

Following Section 4.2 we can reconsider the JKO scheme 34 where we replace the  $W_2$  distance with the  $W_{bw}$  on mixing measures:

$$\min_{m^i, \epsilon^i} \left\{ \text{KL}(\nu|\pi) + \frac{1}{2\gamma} W_{bw}^2(\hat{p}, \hat{p}_k) \right\}, \quad (35)$$

which, at the limit  $\gamma \rightarrow 0$ , is equivalent to:

$$\min_{m^i, \epsilon^i} \left\{ \text{KL}(\nu|\pi) + \frac{1}{2N\gamma} \sum_{i=1}^N \text{BW}^2(\mathcal{N}(m^i, \epsilon^i \text{Id}), \mathcal{N}(m_k^i, \epsilon_k^i \text{Id})) \right\}. \quad (36)$$

The problem in Eq. (36) corresponds to the discrete scheme introduced for the full covariance case  $\Sigma^i = \epsilon^i I$  in [Lambert et al., 2022, Appendix B], without further justification. This scheme can therefore be reinterpreted as a gradient flow on the space of discrete mixing measures.

## B.2 Flow of mixtures of isotropic Gaussian

We now solve the problem (36) when  $\nu$  is a mixture of isotropic Gaussians. We recall our notation for an isotropic Gaussian mixture:  $\nu = k_\epsilon \otimes \mu := \frac{1}{N} \sum_{i=1}^N \mathcal{N}(m^i, \epsilon^i \text{Id})$  where  $k_{\epsilon^i} \star \delta_{m^i} := \mathcal{N}(m^i, \epsilon^i \text{Id})$ , as well as our loss function:

$$F([m^i, \epsilon^i]_{i=1}^N) := \text{KL} \left( \frac{1}{N} \sum_{i=1}^N \mathcal{N}(m^i, \epsilon^i \text{Id}) \middle| \pi \right) = \text{KL}(k_\epsilon \otimes \mu | \pi). \quad (37)$$

The derivative of this loss w.r.t. the parameters of the Gaussian mixture are given in Appendix A.7.2 and we report them here:

$$\nabla_{m^i} F = \frac{1}{N} \mathbb{E}_{k_{\epsilon^i} \star \delta_{m^i}} \left[ \nabla \ln \left( \frac{k_\epsilon \otimes \mu}{\pi} \right) \right], \quad (38)$$

$$\nabla_{\epsilon^i} F = \frac{1}{2N\epsilon^i} \mathbb{E}_{k_{\epsilon^i} \star \delta_{m^i}} \left[ (\cdot - m^i)^T \nabla \ln \left( \frac{k_\epsilon \otimes \mu}{\pi} \right) \right]. \quad (39)$$

Using the derivative of the Bures-Wasserstein distance computed in Appendix A.3, we obtain the first order condition w.r.t.  $\epsilon^i$  for the JKO-like scheme (36):

$$d \left( 1 - \sqrt{\frac{\epsilon_k^i}{\epsilon^i}} \right) = -2N\gamma \nabla_{\epsilon^i} F \Rightarrow \epsilon_k^i = \left( 1 + \frac{2N\gamma}{d} \nabla_{\epsilon^i} F \right)^2 \epsilon^i, \quad (40)$$

and following Appendix A.3, as  $\gamma \rightarrow 0$  we obtain the flow:

$$\dot{\epsilon}^i = -\frac{2}{d} \mathbb{E}_{k_{\epsilon^i} \star \delta_{m^i}} \left[ (\cdot - m^i)^T \nabla \ln \left( \frac{k_\epsilon \otimes \mu}{\pi} \right) \right].$$

On the other hand, first order condition on the mean gives the implicit update:

$$m^i = m_k^i - N\gamma \nabla_{m^i} F,$$

and at the limit  $\gamma \rightarrow 0$  we obtain the flow:

$$\dot{m}^i = -\mathbb{E}_{k_{\epsilon^i} \star \delta_{m^i}} \left[ \nabla \ln \left( \frac{k_\epsilon \otimes \mu}{\pi} \right) \right].$$

## B.3 Discrete update and Gaussian particles

**Derivations of the updates.** For the scheme given in Eq (13), we obtain the following first-order conditions on the parameters for  $i = 1, \dots, N$ :

$$\frac{1}{2N} d \left( 1 - \sqrt{\frac{\epsilon_k^i}{\epsilon^i}} \right) = -\gamma \nabla_{\epsilon^i} F([m_k^j, \epsilon_k^j]_{j=1}^N) \Leftrightarrow \epsilon^i = \epsilon_k^i \left( 1 + \frac{2N\gamma}{d} \nabla_{\epsilon^i} F([m_k^j, \epsilon_k^j]_{j=1}^N) \right)^{-2},$$

and using a Taylor expansion for small step size  $\gamma$ ,  $(1+x)^{-1} = 1-x+o(x)$  we obtain:

$$\epsilon^i = \epsilon_k^i \left( 1 - \frac{2N\gamma}{d} \nabla_{\epsilon^i} F([m_k^j, \epsilon_k^j]_{j=1}^N) \right)^2. \quad (41)$$

The first-order condition on  $m^i$  gives the explicit update:

$$m^i = m_k^i - N\gamma \nabla_{m^i} F([m_k^j, \epsilon_k^j]_{j=1}^N). \quad (42)$$

**Riemannian interpretation.** We now show that this update has also a Riemannian interpretation, extending our geometric analysis from Appendix A.4 to the case of mixtures. We can compute the isotropic Bures-Wasserstein gradient for each Gaussian component, by projecting the Wasserstein-2 gradient of the KL objective to the iBW-tangent space of each component. Each Gaussian particle (component) follows its own trajectory ruled by  $\nabla_{\text{iBW}^i} \mathcal{F}$ . Adopting this point of view, we obtain the following system of updates for  $i = 1, \dots, N$ :

$$p_{k+1}^i = \exp_{p_k^i}(-\gamma \nabla_{\text{iBW}^i} \mathcal{F}(\nu_k)), \quad (43)$$

where  $\exp_{p_k^i}$  is the exponential map from the iBW tangent space at  $p_k^i$  defined in Appendix A.4. We then get:

$$\nabla_{\text{iBW}^i} \mathcal{F}(\nu) = \text{proj}_{T_{p^i} \text{iBW}(\mathbb{R}^d)} \nabla_{W_2} \mathcal{F}(\nu) = \underset{w \in T_{p^i} \text{iBW}(\mathbb{R}^d)}{\text{argmin}} \|w - \nabla_{W_2} \mathcal{F}(\nu)\|_{p^i}^2,$$

with  $w = (a, s) \in \mathbb{R}^d \times \mathbb{R}$ . Together with (38,39) it gives:

$$\begin{aligned} a &= \mathbb{E}_{p^i} [\nabla_{W_2} \mathcal{F}(\nu)] = N \nabla_{m^i} F([m^j, \epsilon^j]_{j=1}^N), \\ s &= \frac{1}{d\epsilon^i} \mathbb{E}_{p^i} [(\cdot - m^i)^T \nabla_{W_2} \mathcal{F}(\nu)] = \frac{2N}{d} \nabla_{\epsilon^i} F([m^j, \epsilon^j]_{j=1}^N). \end{aligned}$$

Using Eq (26), we obtain the discrete updates Eq (41)-(42).

#### B.4 Background on Wasserstein distances for Gaussian mixtures [Delon and Desolneux, 2020]

In Delon and Desolneux [2020], the authors introduce the  $MW_2$  distance as a Wasserstein distance between Gaussian mixtures, where the transport plan is itself constrained to be a Gaussian mixture (with any number of components). We denote by  $GMM$  the latter space.

Namely, let  $p_0, p_1$  being two general Gaussians mixtures  $p_0 = \sum_{i=1}^{K_0} \pi_0^i p_0^i$  and  $p_1 = \sum_{i=1}^{K_1} \pi_1^i p_1^i$ , the  $MW_2$  distance is defined as:

$$MW_2^2(p_0, p_1) = \min_{\gamma \in \Pi(p_0, p_1) \cap GMM} \int \|x_1 - x_2\|^2 d\gamma(x_1, x_2) \geq W_2^2(p_0, p_1), \quad (44)$$

which is an upper bound to the true Wasserstein distance since the transport plan is constrained.

It has been shown in [Delon and Desolneux, 2020, Proposition 4] that this distance is also equal to:

$$MW_2^2(p_0, p_1) = \min_{W \in \Pi(\pi_0, \pi_1)} \sum_{i=1}^{K_0} \sum_{j=1}^{K_1} W_{i,j} BW^2(\mathcal{N}(m_i, \Sigma_i), \mathcal{N}(m_j, \Sigma_j)), \quad (45)$$

where  $\Pi(\pi_0, \pi_1)$  is the set of coupling matrices between the vector of weights  $\pi_0$  and  $\pi_1$  of the two mixtures defined by  $\Pi(\pi_0, \pi_1) = \{W \in \mathcal{M}_{K_0 \times K_1}(\mathbb{R}^+) | \forall i, \sum_j W_{i,j} = \pi_0^i; \forall j, \sum_i W_{i,j} = \pi_1^j\}$ , where we note  $\mathcal{M}_{n \times p}(\mathbb{R}^+)$  the set of matrices of size  $n \times p$  with positive values. Moreover, Delon and Desolneux [2020] showed that the optimal transport plan takes the form:

$$\gamma(x, y) = \sum_{i=1}^{K_0} \sum_{j=1}^{K_1} W_{i,j}^* p_0^i(x) \delta_{y=T_{i,j}^{BW}(x)}(y), \quad (46)$$

where  $W^*$  is the optimal coupling matrix and  $T_{i,j}^{BW}$  is the BW transport map from Gaussian component  $i$  to Gaussian component  $j$ . The transport plan  $\gamma$  is a GMM with at most  $K_0 K_1$  Gaussian components which are degenerated.

**Mixture model with fixed number of components.** We now consider the case of mixture of exactly  $N$  Gaussians with equal, fixed weights:

$$p = \frac{1}{N} \sum_{i=1}^N \mathcal{N}(m^i, \Sigma^i) = \frac{1}{N} \sum_{i=1}^N p^i,$$

where  $p^i$  is the  $i^{th}$  Gaussian component of the mixture.

688 The  $MW_2$  distance takes an even simpler expression when considering a distance between two mix-  
 689 tures  $p_0, p_1$  of exactly  $N$  Gaussians with equal, fixed weights. If we note  $\mathfrak{S}_N$  the set of permutations  
 690 over  $\{1, \dots, N\}$ , we have:

$$MW_2^2(p_0, p_1) = \min_{\sigma \in \mathfrak{S}_N} \frac{1}{N} \sum_{i=1}^N BW^2(\mathcal{N}(m_i, \Sigma_i), \mathcal{N}(m_{\sigma(i)}, \Sigma_{\sigma(i)})) = W_{bw}^2(\hat{p}_0, \hat{p}_1),$$

691 where  $W_{bw}$  is the Wasserstein distance between mixing measures defined in Section 4.2. The lower  
 692 and upper bounds for  $MW_2^2$  given in [Delon and Desolneux, 2020, Proposition 6] then transfer to  
 693  $W_{bw}$  and give:

$$W_2(p_0, p_1) \leq W_{bw}(\hat{p}, \hat{p}_k) \leq W_2(p_0, p_1) + \sqrt{\frac{2}{N} \sum_{k=1}^N \text{Tr} \Sigma_0^k} + \sqrt{\frac{2}{N} \sum_{k=1}^N \text{Tr} \Sigma_1^k}. \quad (47)$$

694 The last term simplifies for isotropic Gaussians and we finally obtain:

$$0 \leq W_{bw}^2(\hat{p}, \hat{p}_k) - W_2^2(\nu, \nu_k) \leq 2\sqrt{2d}\epsilon^*, \quad (48)$$

695 where  $\epsilon^*$  is the maximal variance of the mixtures  $\nu, \nu_k$ . When  $\epsilon^* \rightarrow 0$ , such that the Gaussian mixture  
 696 degenerates into an empirical measure, the two distance matches.

**Geodesics on mixtures** Finally, when considering mixtures with  $N$  equal weights, the transport  
 plan has exactly  $N$  components and can be written:

$$\gamma(x, y) = \frac{1}{N} \sum_{i=1}^N p_0^i(x) \delta_{y=T_{i, \sigma^*(i)}^{BW}(x)}(y),$$

697 such that the mixture model with exactly  $N$  components is stable along the geodesics transported by  
 698 this plan. Indeed, the intermediate measure between two GMM  $p_0$  and  $p_1$  is given by the formula for  
 699  $t \in [0, 1]$ :

$$\mu_t = P_t \# \gamma \text{ where } P_t(x) = (1-t)x + ty. \quad (49)$$

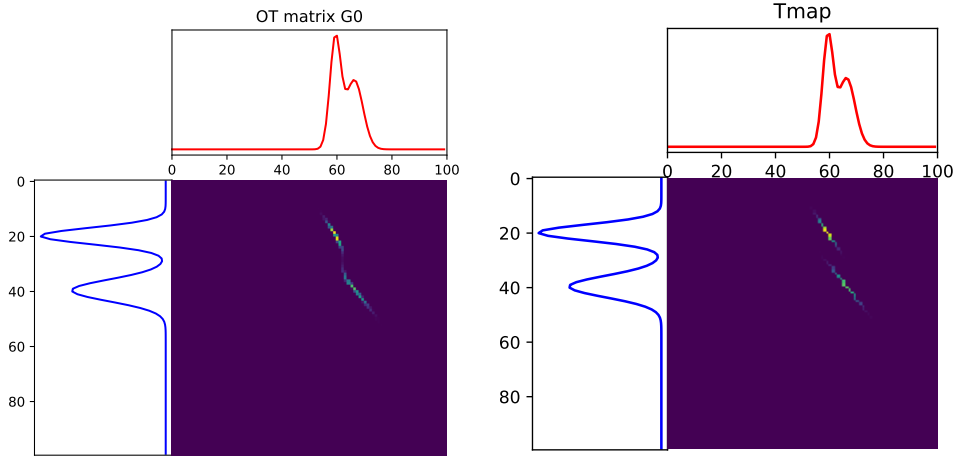


Figure 5: Optimal transport plan between two mixtures of two Gaussians with equal weights  $\frac{1}{2}$ . On the left is the Wasserstein distance  $W_2$  case, where the optimal transport plan is not constrained. On the right is the  $MW_2$  case, where the optimal transport plan is constrained to be a mixture of Gaussians. In the  $W_2$  case, there exists a bijective map. We don't have such a bijective map for the  $MW_2$  case. Indeed, in the right figure, some points have two images. These figures are generated using the Python Optimal Transport library (<https://pythonot.github.io/>).

700 Applying this to our specific case, we obtain:

$$p_t = \frac{1}{N} \sum_{i=1}^N ((1-t)\text{Id} + tT_{i, \sigma^*(i)}^{BW}) \# p_0^i, \quad (50)$$

701 where  $p_t$  has exactly  $N$  components, proving that our GMM structure with  $N$  components is stable  
 702 along the geodesics.

703 We may wonder if a transport map exists in our simpler framework of mixtures with a fixed number  
 704 of components and equal weights. Unfortunately, this is not the case, as illustrated in Figure 5, where  
 705 we observe that the map between the two mixtures is not bijective and cannot be represented by a  
 706 function  $T(x)$ .

## 707 C Natural gradient descent updates

708 In this section, we provide more details on Natural gradient descent on IG, which corresponds to the  
 709 algorithm proposed by [Lin et al., 2019].

710 Recall that  $F : \mathbb{R}^d \times \mathbb{R}^{+*}$ , the objective function, is defined as :  $F(m, \epsilon) = \text{KL}(\mathcal{N}(m, \epsilon \text{Id}) | \pi)$ . Let  
 711  $\gamma > 0$  a step-size, consider the scheme at each step  $k \geq 0$ :

$$(m_{k+1}, \epsilon_{k+1}) = \underset{m, \epsilon \in \mathbb{R}^d \times \mathbb{R}^{+*}}{\text{argmin}} \langle \nabla F(m_k, \epsilon_k), (m, \epsilon) - (m_k, \epsilon_k) \rangle + \frac{1}{\gamma} \text{KL}(\mathcal{N}(m, \epsilon \text{Id}) | \mathcal{N}(m_k, \epsilon_k \text{Id})). \quad (51)$$

712 Since

$$\text{KL}(\mathcal{N}(m, \epsilon \text{Id}) | \mathcal{N}(m_k, \epsilon_k \text{Id})) = \frac{1}{2} \left( d \cdot \frac{\epsilon}{\epsilon_k} + \frac{\|m - m_k\|^2}{\epsilon_k} - d + d \log \frac{\epsilon_k}{\epsilon} \right), \quad (52)$$

713 first order conditions of Eq (51) yield

$$\frac{1}{\epsilon_{k+1}} - \frac{1}{\epsilon_k} = \frac{2\gamma}{d} \nabla_{\epsilon_k} F(m_k, \epsilon_k), \quad (53)$$

$$\frac{m_{k+1} - m_k}{\epsilon_k} = -\gamma \nabla_{m_k} F(m_k, \epsilon_k). \quad (54)$$

714 Thus, we find similar updates as in [Lin et al., 2019] (see for instance Eq (16) therein, with the  
 715 difference that we do not update the weights of the components in the mixture as in their paper).  
 716 This is expected as the latter paper implements a natural gradient descent on the natural parameter  
 717  $\lambda = [\frac{m}{\epsilon}, -\frac{1}{2\epsilon}]$ . The former is known to be equivalent [Raskutti and Mukherjee, 2015] to mirror  
 718 descent on the mean parameters, where the geometry is induced by a Bregman divergence generated  
 719 by the Legendre transform of the partition function.

720 **Mixture case.** In the mixture case, we consider the scheme at each step  $k \geq 0$ :

$$[\eta_{k+1}]_{j=1}^N = \arg \min_{(\eta^j)_{j=1}^N} \left\{ \langle \nabla F([\eta_k^j]_{j=1}^N), [\eta^j]_{j=1}^N - [\eta_k^j]_{j=1}^N \rangle + \frac{1}{\gamma N} \sum_{j=1}^N \text{KL}(\mathcal{N}(\eta^j), \mathcal{N}(\eta_k^j)) \right\}, \quad (55)$$

721 and have the following updates for  $i = 1, \dots, N$ :

$$\begin{aligned} \frac{1}{\epsilon_{k+1}^i} - \frac{1}{\epsilon_k^i} &= \frac{2N\gamma}{d} \nabla_{\epsilon_k^i} F([m_k^j, \epsilon_k^j]_{j=1}^N), \\ \frac{m_{k+1}^i - m_k^i}{\epsilon_k^i} &= -\gamma N \nabla_{m_k^i} F([m_k^j, \epsilon_k^j]_{j=1}^N). \end{aligned}$$

722 *Remark C.1.* Note that while  $\text{KL}(\cdot | \cdot)$  is known to be a Bregman divergence on the space of probability  
 723 distributions over  $\mathbb{R}^d$  [Aubin-Frankowski et al., 2022], it is not a Bregman divergence on  $\mathbb{R}^d \times \mathbb{R}^{+*}$ .  
 724 Indeed, note that Eq (52) does not decouple the mean and variance terms, resulting in coupled updates  
 725 in Eq (53)-(54).

## 726 D Approximation error for mixtures of isotropic Gaussians

727 [Huix et al., 2024, Theorem 7] states that the approximation error of VI within the family of  
 728 mixtures of Gaussian distributions with equal weight and constant isotropic covariance in the (reverse)

729 Kullback-Leibler tends to 0 as  $N$  tends to infinity, under the assumption that the target distribution  
 730 writes as an infinite mixture of these isotropic Gaussian components with same covariance. Here  
 731 we derive a similar result for our richer family, i.e., mixtures of (isotropic) Gaussian distributions  
 732 with equal weights (and possibly different covariances). Note that we are deriving the results for  
 733 mixtures of isotropic Gaussians, but the result and computation are the same for Gaussians with full  
 734 covariance matrix.

735 **Assumption D.1.** There exists  $p^*$  on  $\mathbb{R}^d \times \mathbb{R}^{+*}$  such that the target  $\pi$  writes as:

$$\pi := \int_{\Theta} k_{\epsilon}^m dp^*(m, \epsilon), \quad (56)$$

736 where  $k_{\epsilon}^m := k_{\epsilon}(\cdot - m)$  for any  $x \in \mathbb{R}^d$ .

737 Recall that we use the notation  $\rho_N = k_{\epsilon} \otimes \mu = \int k_{\epsilon}^m d\hat{p}(m, \epsilon)$  with  $\mu = \frac{1}{N} \sum_{i=1}^N \delta_{m^i}$  for an isotropic  
 738 Gaussian mixture with  $N$  components. We now state and prove our generalization of [Huix et al.,  
 739 2024, Theorem 7].

740 **Theorem D.2.** Suppose that assumption 56 holds, then

$$\min_{\rho \in \mathcal{C}^N} \text{KL}(\rho|\pi) \leq C_{\pi}^2 \frac{\log(N) + 1}{N}, \quad \text{where} \quad C_{\pi}^2 = \int \frac{\int k_{\epsilon}^m(x)^2 dp^*(m, \epsilon)}{\int k_{\epsilon}^m(x) dp^*(m, \epsilon)} dx.$$

741 *Proof.* We denote

$$D_N = \min_{\rho \in \mathcal{C}^N} \text{KL}(\rho|\pi), \quad \rho_N = \operatorname{argmin}_{\rho \in \mathcal{C}^N} \text{KL}(\rho|\pi).$$

742 For any  $m \in \mathbb{R}^d$ , we consider  $\rho_{N+1}^{m, \epsilon} \in \mathcal{C}_{N+1}$  defined as

$$\rho_{N+1}^{m, \epsilon} = (1 - \alpha)\rho_N + \alpha k_{\epsilon}^m,$$

743 with  $\alpha = \frac{1}{N+1}$ . By definition of  $D_N$ , we have that,  $D_{N+1} \leq \text{KL}(\rho_{N+1}^{m, \epsilon}|\pi)$ . Denoting  $f(x) =$   
 744  $x \log x$ , we have  $\text{KL}(\rho_{N+1}^{m, \epsilon}|\pi) = \int f(r_{N+1}) d\pi$ , where we define:

$$r_{N+1} := \frac{\rho_{N+1}^{m, \epsilon}}{\pi} = (1 - \alpha) \frac{\rho_N}{\pi} + \alpha \frac{k_{\epsilon}^m}{\pi} := r_0 + \alpha \frac{k_{\epsilon}^m}{\pi}.$$

745 Defining  $B(x) = \frac{x \log x - x + 1}{(x-1)^2}$  for  $x \in \mathbb{R}_+^* \setminus \{1\}$ . By Lemma D.3, this function is decreasing; and  
 746 since  $r_{N+1}(x) \geq r_0(x) \forall x$ , we have  $B(r_{N+1}(x)) \leq B(r_0(x))$ . It follows that

$$\begin{aligned} r_{N+1} \log(r_{N+1}) &\leq r_{N+1} - 1 + B(r_0)(r_{N+1} - 1)^2 \\ &= r_0 + \alpha \frac{k_{\epsilon}^m}{\pi} - 1 + B(r_0)(r_0 + \alpha \frac{k_{\epsilon}^m}{\pi} - 1)^2 \\ &= \alpha \frac{k_{\epsilon}^m}{\pi} + r_0 \log(r_0) + \alpha^2 \left( \frac{k_{\epsilon}^m}{\pi} \right)^2 B(r_0) + 2\alpha \frac{k_{\epsilon}^m}{\pi} B(r_0)(r_0 - 1). \end{aligned} \quad (57)$$

747 Moreover, we have:

$$\begin{aligned}
D_{N+1} &= \int D_{N+1} dp^*(m, \epsilon) \\
&\leq \int \text{KL}(\rho_{N+1}^{m, \epsilon} | \pi) dp^*(m, \epsilon) \\
&= \int \int f(r_{N+1}) d\pi dp^*(m, \epsilon) \\
&= \int \int r_{N+1} \log(r_{N+1}) d\pi dp^*(m, \epsilon) \\
&\leq \int \int \alpha \frac{k_\epsilon^m}{\pi} d\pi dp^*(m, \epsilon) + \int \int r_0 \log(r_0) d\pi dp^*(m, \epsilon) \\
&\quad + \int \int \alpha^2 \left( \frac{k_\epsilon^m}{\pi} \right)^2 B(r_0) d\pi dp^*(m, \epsilon) \\
&\quad + \int \int 2\alpha \frac{k_\epsilon^m}{\pi} B(r_0)(r_0 - 1) d\pi dp^*(m, \epsilon) \quad \text{from (57)} \\
&= \alpha + \int r_0(x) \log(r_0(x)) d\pi(x) + \alpha^2 \int \int \frac{k_\epsilon^m(x)^2}{\pi(x)} B(r_0(x)) dx dp^*(m, \epsilon) \\
&\quad + 2\alpha \int \int k_\epsilon^m B(r_0(x))(r_0(x) - 1) dx dp^*(m, \epsilon) \\
&= \alpha + \int r_0(x) \log(r_0(x)) d\pi(x) \tag{a} \\
&\quad + \alpha^2 \int \int \frac{k_\epsilon^m(x)^2}{\pi(x)} B(r_0(x)) dx dp^*(m, \epsilon) \tag{b} \\
&\quad + 2\alpha \int B(r_0(x))(r_0(x) - 1) \pi(x) dx. \tag{c}
\end{aligned}$$

748 We first observe that we can write (a) in function of  $D_N$ . Indeed,  $r_0(x) = (1 - \alpha) \frac{\rho_N}{\pi}$ , so

$$\begin{aligned}
\text{(a)} &= \int r_0(x) \log(r_0(x)) d\pi(x) \\
&= \int (1 - \alpha) \frac{\rho_N}{\pi} \log \left( (1 - \alpha) \frac{\rho_N}{\pi} \right) d\pi \\
&= (1 - \alpha) \int \rho_N(x) \log \left( \frac{\rho_N(x)}{\pi(x)} \right) dx + (1 - \alpha) \log(1 - \alpha) \\
&= (1 - \alpha) D_N + (1 - \alpha) \log(1 - \alpha).
\end{aligned}$$

749 For the second term (b), we have that  $\lim_{x \rightarrow 0^+} B(x) = 1$  and since  $B$  decrease,  $B(x) \leq 1$ , thus  
750  $B(r_0(x)) \leq 1$ , this implies :

$$\begin{aligned}
\text{(b)} &= \alpha^2 \int \int \frac{k_\epsilon^m(x)^2}{\pi(x)} B(r_0(x)) dx dp^*(m, \epsilon) \leq \alpha^2 \int \int \frac{k_\epsilon^m(x)^2}{\pi(x)} dx dp^*(m, \epsilon) \\
&= \alpha^2 C_\pi^2.
\end{aligned}$$

751 And for the third term (c), we have that  $B(x)(x - 1) \leq \sqrt{x} - 1$ , see Lemma D.4. Thus,

$$\begin{aligned}
\text{(c)} &= 2\alpha \int B(r_0(x))(r_0(x) - 1) \pi(x) dx \leq 2\alpha \int \left( \sqrt{r_0(x)} - 1 \right) \pi(x) dx \\
&= 2\alpha \int \sqrt{(1 - \alpha) \frac{\rho_N(x)}{\pi(x)}} \pi(x) dx - 2\alpha \\
&= 2\alpha \sqrt{1 - \alpha} \int \sqrt{\rho_N(x) \pi(x)} dx - 2\alpha \\
&= 2\alpha \sqrt{1 - \alpha} (1 - H^2(\rho_N, \pi)) - 2\alpha \\
&\leq 2\alpha \sqrt{1 - \alpha} - 2\alpha,
\end{aligned}$$

752 where we have used the definition of the squared Hellinger distance  $H^2(f, g) = \int \sqrt{f(x)g(x)}dx$   
 753 and the property stating that for any densities of probability  $f, g$ ,  $0 \leq H^2(f, g) \leq 1$ .

754 Finally, we have

$$\begin{aligned} D_{N+1} &\leq \alpha + (1 - \alpha)D_N + (1 - \alpha)\log(1 - \alpha) + \alpha^2 C_\pi^2 + 2\alpha(\sqrt{1 - \alpha} - 1) \\ &\leq (1 - \alpha)D_N + \alpha^2 C_\pi^2, \end{aligned} \quad (58)$$

755 using Lemma D.5, stating that  $\alpha + (1 - \alpha)\log(1 - \alpha) + 2\alpha(\sqrt{1 - \alpha} - 1) \leq 0$ .

756 The previous inequality (58) is true for any  $n \geq 0$  and recalling that  $\alpha = \frac{1}{n+1}$  we have,

$$\begin{aligned} D_{n+1} &\leq (1 - \alpha)D_n + \alpha^2 C_\pi^2 \\ (n+1)D_{n+1} - nD_n &\leq \frac{1}{n+1} C_\pi^2 \\ \sum_{n=0}^{N-1} (n+1)D_{n+1} - nD_n &\leq C_\pi^2 \sum_{n=0}^{N-1} \frac{1}{n+1} \\ ND_N &\leq C_\pi^2 (\log(N) + 1) \\ D_N &\leq C_\pi^2 \frac{\log(N) + 1}{N}, \end{aligned}$$

757 where the Harmonic number  $\sum_{n=0}^{N-1} \frac{1}{n+1}$  has been bounded by  $\log(N) + 1$ .  $\square$

758 **Lemma D.3.** The function  $B(x) = \frac{x \log x - x + 1}{(x-1)^2} \quad \forall x \in \mathbb{R}^{+*} \setminus \{1\}$ , is decreasing.

759 *Proof.* For all  $x \in \mathbb{R}^{+*} \setminus \{1\}$  the gradient of  $B$  writes:

$$\nabla B(x) = \frac{(x-1)\log x - 2(x \log x - x + 1)}{(x-1)^3}$$

- 760 • For  $x \in (0, 1)$ , the denominator is strictly negative and the numerator strictly positive, thus  
 761  $\nabla B(x) \leq 0$ .
- 762 • For  $x \in (1, \infty)$ , the denominator is strictly positive and the numerator is strictly negative,  
 763 thus  $\nabla B(x) \leq 0$ .

764 So  $B$  is decreasing on both intervals, and  $\lim_{x \rightarrow 1^-} = \frac{1}{2}$  and  $\lim_{x \rightarrow 1^+} = \frac{1}{2}$  by Hospital's rule.  $\square$

765 **Lemma D.4.** The function  $B$  satisfies:  $B(x)(x-1) \leq \sqrt{x} - 1 \quad \forall x \in \mathbb{R}^{+*} \setminus \{1\}$ .

766 *Proof.* Let  $C(x) := B(x)(x-1) = \frac{x \log x - x + 1}{x-1}$

- 767 • For  $x \in (0, 1)$ ,  $\log x \geq \frac{x-1}{\sqrt{x}}$  implies  $\frac{\log x}{x-1} \leq \frac{1}{\sqrt{x}} \Rightarrow \frac{x \log x}{x-1} \leq \frac{x}{\sqrt{x}} = \sqrt{x}$ ,
- 768 • For  $x \in (1, \infty)$ ,  $\log x \leq \frac{x-1}{\sqrt{x}}$  implies  $\frac{\log x}{x-1} \leq \frac{1}{\sqrt{x}} \Rightarrow \frac{x \log x}{x-1} \leq \frac{x}{\sqrt{x}} = \sqrt{x}$ .

769 and  $C(x) - \sqrt{x} - 1 = \frac{x \log x}{x-1} - \sqrt{x} \leq 0$ .  $\square$

770 **Lemma D.5.** Let consider  $\alpha = \frac{1}{n+1} \quad \forall n \in \mathbb{N}$ , then  $\alpha + (1 - \alpha)\log(1 - \alpha) + 2\alpha(\sqrt{1 - \alpha} - 1) \leq 0$ .

771 *Proof.* We have:

$$\begin{aligned} \alpha + (1 - \alpha)\log(1 - \alpha) + 2\alpha(\sqrt{1 - \alpha} - 1) &= -\alpha + (1 - \alpha)\log(1 - \alpha) + 2\alpha\sqrt{1 - \alpha} \\ &\leq \alpha(\alpha - 2 + 2\sqrt{1 - \alpha}) \\ &\leq \alpha - 2 + 2\sqrt{1 - \alpha} \\ &\leq 0, \end{aligned}$$

772 using  $\log(1 - \alpha) \leq -\alpha$  and the fact that  $\alpha = 1/(n + 1) \leq 1$  for the first and second inequality. For  
 773 the last inequality we have used the fact that the before last expression is decreasing and is equal to 0  
 774 when  $\alpha$  goes to 0.  $\square$

## 775 E Additional experiments and details

### 776 E.1 Updates for the full-covariance scheme of Lambert et al. [2022, Section 5.2]

777 In this section we detail the updates for the full-covariance scheme of Lambert et al. [2022, Section  
 778 5.2]. The parameter space is  $\Theta = \mathbb{R}^d \times \mathbb{S}_{++}^d$  (the space of means and covariance matrices). Consider  
 779 initializing this evolution at a finitely supported distribution  $p_0$ :

$$p_0 = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_0^{(i)}} = \frac{1}{N} \sum_{i=1}^N \delta_{(m_0^{(i)}, \Sigma_0^{(i)})}$$

780 It has been checked in Lambert et al. [2022] that the system of ODEs thus initialized maintains a  
 781 finite mixture distribution:

$$p_t = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_t^{(i)}} = \frac{1}{N} \sum_{i=1}^N \delta_{(m_t^{(i)}, \Sigma_t^{(i)})},$$

782 where the parameters  $\theta_t^{(i)} = (m_t^{(i)}, \Sigma_t^{(i)})$  evolve according to the following interacting particle  
 783 system, for  $i \in [N]$

$$\dot{m}_t^{(i)} = -\mathbb{E} \nabla \ln \frac{\nu_t}{\pi} (Y_t^{(i)}), \quad (11) \quad (59)$$

$$\dot{\Sigma}_t^{(i)} = -\mathbb{E} \nabla^2 \ln \frac{\nu_t}{\pi} (Y_t^{(i)}) \Sigma_t^{(i)} - \Sigma_t^{(i)} \mathbb{E} \nabla^2 \ln \frac{\nu_t}{\pi} (Y_t^{(i)}), \quad (12) \quad (60)$$

784 where  $Y_t^{(i)} \sim p_{\theta_t^{(i)}}$  and  $\nu_t = \int \mathcal{N}_\theta dp_t(\theta)$ . The ODEs are solved with a Runge–Kutta scheme of 4th  
 785 order. We used the code provided at <https://github.com/marc-h-lambert/W-VI>.

### 786 E.2 Details experimentations

787 As mentioned in Section 6, we detail here our experimental setup and hyperparameters whose values  
 788 are provided in Table 1.

789 **Initialization of the variational mixture:** We initialize the variational mixture by sampling the  
 790 means in a ball of size  $[-s, s]^d$ , where  $s \in \mathbb{R}^{+*}$ , setting each covariance matrix to  $r I_d$ , where  
 791  $r \in \mathbb{R}^{+*}$ , and the number of components to  $N \in \mathbb{N}^*$ .

792 **Optimization hyperparameters:** We set the step-size  $\gamma$ , the number of iteration  $n_{\text{iter}}$ , the number of  
 793 Monte Carlo samples for gradient estimation  $B_{\text{grad}} = 10$  and for KL estimation  $B_{\text{KL}} = 1000$ .

794 **Normalizing Flows:** For NF baseline, we used a simplified RealNVP architecture [Dinh et al., 2017]  
 795 with  $b = 2$  coupling layers and hidden dimension  $h = 124$ , which yields a Neural Network (NN)  
 796 with 4976 parameters in dimension  $d = 2$ . Our isotropic mixture framework, has  $N(d + 1) = 3N$   
 797 parameters, even for large  $N$ , the NF model remains more complex and costly to optimize. Therefore,  
 798 for each target distribution, we tuned the learning rate and number of iterations for the NF method  
 799 separately, rather than using the same settings as for the VI mixture methods, since their optimization  
 800 dynamics differ significantly.

801 **MOG targets:** To generate target MoG distributions, as in the initialization of variational MOG, we  
 802 fix  $s_{\text{tg}}$  and  $N_{\text{tg}}$ . Each component covariance matrix is constructed by sampling a random symmetric  
 803 positive-definite matrix (full diagonal or isotropic) and scaling it by  $r_{\text{tg}}$ . We draw raw weights  
 804 uniformly in  $\{1, \dots, 2N_{\text{tg}}\}$  and normalize them to one. In Figure 6(c) all weights are equal except in  
 805 case (c-3), where one component has weight 0.1 and the remaining ones share the remaining mass  
 806 equally and the component means are placed at  $(\pm \text{pt}, \pm \text{pt})$ .

807 **Datasets:** We have used popular datasets from the UCI repository, as well as MNIST. The training  
 808 ratio has been set to 0.5 for UCI datasets and 0.8 for MNIST.

Table 1: Hyperparameters

Exp.	MOG Init.		Optim	Targets
	$s$	$r$	$\gamma$	$n_{\text{iter}}$
Figure 1				
mog	15	2	$10^{-1}$	$10^3$ $\{s_{\text{tg}}, r_{\text{tg}}, N_{\text{tg}}\} = \{8, 5, 5\}$
Figure 2				
mog	$10^2 d^{-1}$	10	$10^{-2} d^{-1}$	$10^3$ $\{s_{\text{tg}}, r_{\text{tg}}, N_{\text{tg}}\} = \{10^2 d^{-1}, 5, 5\}$
Figure 3				
mog	30	100	$10^{-2}$	$10^4$ $\{s_{\text{tg}}, r_{\text{tg}}, N_{\text{tg}}\} = \{10, 1, 10\}$
Figure 4				
breast_cancer	20	10	$10^{-2}$	$10^4$ $\sigma_{\text{prior}}^2 = 100$
wine	20	10	$10^{-2}$	$10^2$ $\sigma_{\text{prior}}^2 = 100$
boston	10	10	$10^{-6}$	$10^4$ $\{\sigma_{\text{prior}}^2, h\} = \{10, 50\}$
Figure 6				
funnel (a)	5	0.5	$10^{-2}$	$10^4$ $\sigma^2 = 1.2$
sinh-arcsinh (b-1)	10	2	$10^{-3}$	$10^4$ $\text{skw} = (-0.2, -0.2)$
sinh-arcsinh (b-2)	10	2	$10^{-3}$	$10^4$ $\text{skw} = (-0.2, -0.5)$
mog (c-1)	10	5	$10^{-2}$	$10^4$ $\{\text{pt}, r_{\text{tg}}, N_{\text{tg}}\} = \{3, 2, 4\}$
mog (c-2)	10	5	$10^{-2}$	$10^4$ $\{\text{pt}, r_{\text{tg}}, N_{\text{tg}}\} = \{4, 1, 4\}$
mog (c-3)	10	5	$10^{-2}$	$10^4$ $\{\text{pt}, r_{\text{tg}}, N_{\text{tg}}\} = \{3, 2, 4\}$
Figure 8				
mog	30	100	$10^{-3}$	$10^4$ $\{s_{\text{tg}}, r_{\text{tg}}, N_{\text{tg}}\} = \{20, 5, 5\}$
Figure 9				
mog	30	100	$10^{-3}$	$10^4$ $\{s_{\text{tg}}, r_{\text{tg}}, N_{\text{tg}}\} = \{10, 10, 10\}$

809 **Computational resources:** All experiments (except MNIST) were conducted on a MacBook Air  
810 (M3, 2024) with an Apple M3 processor and 16 GB of RAM. The MNIST experiments were run on  
811 an NVIDIA 50-90 GPU. Experiment runtimes ranged from a few seconds to up to two hours.

### 812 E.3 Additional 2-D examples

813 We present more experiments on 2D synthetic target distributions on Figure 6.

814 **Funnel distribution:** The funnel distribution [Neal, 2000] in dimension  $d = 2$  has density

$$p(x_1, x_2) = \mathcal{N}(x_1; 0, \sigma^2) \times \mathcal{N}(x_2; 0, e^{x_1}),$$

815 for  $x = (x_1, x_2) \in \mathbb{R}^2$ . Although unimodal, its “funnel” shape is difficult to capture with isotropic  
816 Gaussians. We experimented with  $N = 5, 20, 40$  components, but even for large  $N$ , our isotropic  
817 mixtures struggled, and the BW and NF methods still outperformed them. We follow Cai et al.  
818 [2024a] in setting  $\sigma^2 = 1.2$ .

819 **Sinh-arcsinh normal distribution:** This distribution [Pewsey, 2009] applies a sinh–arcsinh transfor-  
820 mation to a multivariate Gaussian to control skewness  $\text{skw}$  and tail weight  $\tau$ . Let

$$Z_0 \sim \mathcal{N}(m, \Sigma), \quad Z = \sinh(\tau \sinh^{-1}(Z_0) - \text{skw}).$$

821 In our experiments, we use  $\tau = (0.8, 0.8)$ ,  $m = (0, 0)$ ,  $\Sigma = \begin{pmatrix} 1 & 0.4 \\ 0.4 & 1 \end{pmatrix}$ , and vary the skew  
822 parameter  $\text{skw}$  as specified in Table 1.

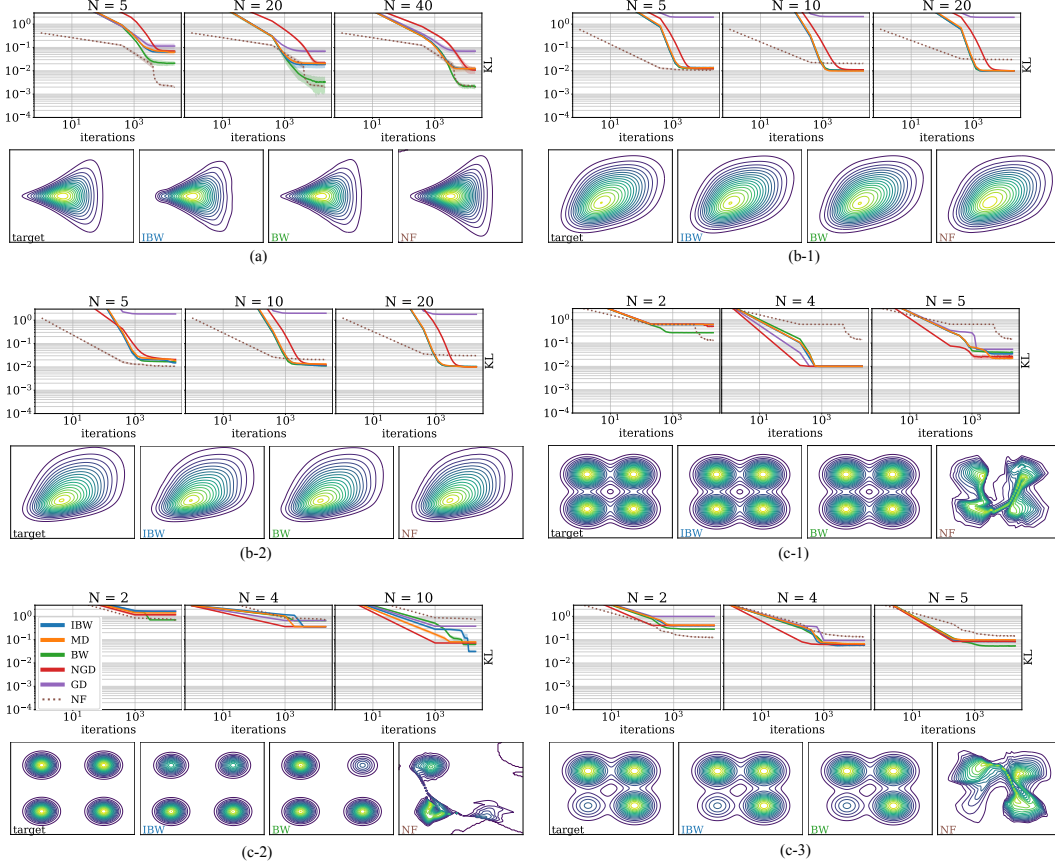


Figure 6: Evolution of the KL divergence over iterations together with optimized variational mixture density on different type of target distribution, (a) Funnel, (b) sinh-arcsinh normal distribution and (c) Gaussian mixture. Optimization performed with different methods (IBW, MD, BW, NGD, GD and NF) and varying  $N$  values.

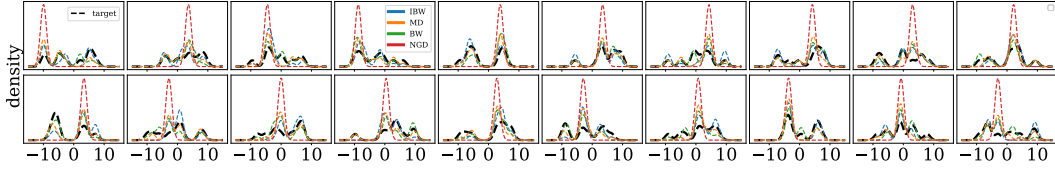


Figure 7: Marginals for MD, IBW, BW and NGD ( $d = 20$ ).

#### 823 E.4 Additional high-dimensional mixtures

824 We first provide the full marginals of the experiment described in Figure 3 in Figure 7.

825 We also performed experiments in  $d = 10$  Figure 8 and  $d = 50$  Figure 9.

#### 826 E.5 More details on Bayesian posteriors experiment

827 In this section we provide some background on logistic and linear regression as well as additional  
828 experiments.

##### 829 E.5.1 Definition of the target distributions

830 Let  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  be a labeled dataset, where  $x_i \in \mathbb{R}^d$  and  $y_i$  is the associated label.

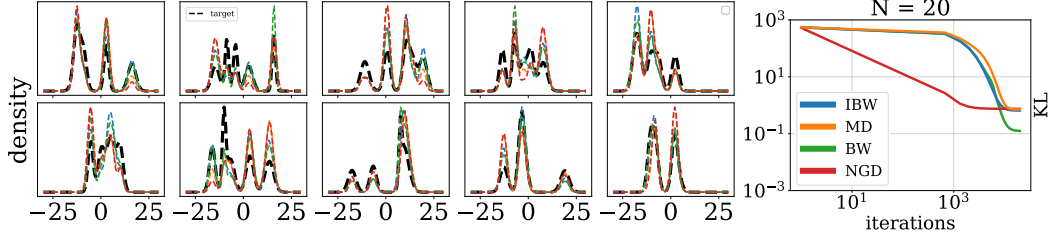


Figure 8: Marginals (left) and KL objective (right) for MD, IBW, BW and NGD ( $d = 10$ ).

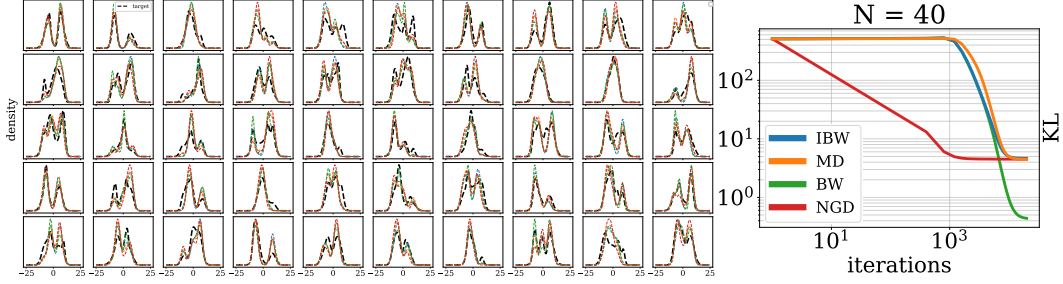


Figure 9: Marginals (left) and KL objective (right) for MD, IBW, BW and NGD ( $d = 50$ ).

**Binary logistic regression:** We model the probability of a binary label  $y_i \in \{0, 1\}$  given  $x_i$  and parameter  $z \in \mathbb{R}^d$  by

$$\pi(y_i | x_i, z) = \sigma(x_i^T z)^{y_i} (1 - \sigma(x_i^T z))^{1-y_i},$$

where  $\sigma(t) = 1/(1 + e^{-t})$  is the logistic function. The likelihood is

$$\mathcal{L}(\mathcal{D} | z) = \prod_{i=1}^n \pi(y_i | x_i, z),$$

and the log-likelihood is

$$\ell(\mathcal{D} | z) = \sum_{i=1}^n \log \pi(y_i | x_i, z) = \sum_{i=1}^n [y_i (x_i^T z) - \log(1 + e^{x_i^T z})].$$

With a Gaussian prior  $\pi(z) = \mathcal{N}(0, \sigma_{prior}^2 I_d)$ , the posterior is

$$\pi(z | \mathcal{D}) \propto \mathcal{L}(\mathcal{D} | z) \pi(z),$$

and its gradient is

$$\nabla_z \log \pi(z | \mathcal{D}) = \sum_{i=1}^n (y_i - \sigma(x_i^T z)) x_i - \frac{z}{\epsilon_z}.$$

**Multi class logistic regression:** For  $L$  classes, let  $z = (z_1, \dots, z_L)$  with each  $z_l \in \mathbb{R}^d$ . Then

$$\pi(y_i = l | x_i, z) = \frac{\exp(x_i^T z_l)}{\sum_{l=1}^L \exp(x_i^T z_l)}, \quad l = 1, \dots, L.$$

**Linear regression:** The classical linear model is

$$y_i = z^T x_i + \xi_i, \quad \xi_i \sim \mathcal{N}(0, \sigma^2),$$

so that

$$y_i \sim \mathcal{N}(x_i^T z, \sigma^2).$$

and the ordinary least squares estimator is

$$\hat{z} = \arg \min_{z \in \mathbb{R}^d} \sum_{i=1}^n (y_i - z^T x_i)^2.$$

841 for which we are able to find a close form when  $X = (x_i)_{i=1}^n$  is invertible.

842 In our Bayesian setting we aim at finding a distribution on  $z$ , put a prior on it, and approximate the  
843 resulting posterior  $\pi(z \mid \mathcal{D})$  via variational inference.

844 **Bayesian Neural Network:** In the Bayesian neural network (BNN) setting, the linear predictor  $z^T x_i$   
845 is replaced with a neural network output  $f(x_i \mid z)$  and model

$$y_i \sim \mathcal{N}(f(x_i \mid z), \sigma^2).$$

846 In our experiments we use a single hidden layer with  $h$  hidden units, ReLU activation function and  
847 output dimension  $c$ . Thus, the dimension of parameters and thus of the problem is

$$d = h(d_{\text{data}} + 1) + c(h + 1).$$

848 For a  $L$ -class classification task,  $c = L$  and the BNN output class probabilities  $\pi(y_i = l \mid x_i, z) =$   
849  $f(x_i \mid z)_l$ .

850 Once the variational approximation to the posterior is optimized, we can make predictions by Bayesian  
851 model averaging:

$$p(y \mid x) = \int \pi(y \mid x, z) \pi_{\text{post}}(z) dz, \quad \text{or} \quad \hat{y} = \int f(x \mid z) \pi_{\text{post}}(z) dz.$$

852 When  $d$  is large (e.g. MNIST, where  $d \approx 10^5$ ), sampling or expectation under a full  $d$ -dimensional  
853 mixture becomes too expensive. To address this, we adopt a mean-field-style approximation: we  
854 model the posterior as a product of identical univariate Gaussian-mixture marginals,

$$z_j \sim \frac{1}{N} \sum_{i=1}^N \mathcal{N}(m^i[j], \epsilon^i), \quad j = 1, \dots, d,$$

855 so that all  $d$  dimensions share the same  $N$ -component mixture. This reduces both memory and  
856 computational cost while retaining multimodality in each coordinate. We updated  $m^i, \epsilon^i$  using the  
857 presented algorithms. In this setting, we follow a classical deep-learning framework. We use a single-  
858 layer neural network with  $h = 256$  hidden units and ReLU activation. The means of the variational  
859 mixture of Gaussians are initialized by sampling from a normal distribution, and a Gaussian prior is  
860 placed on the model parameters.

## 861 E.5.2 Plots

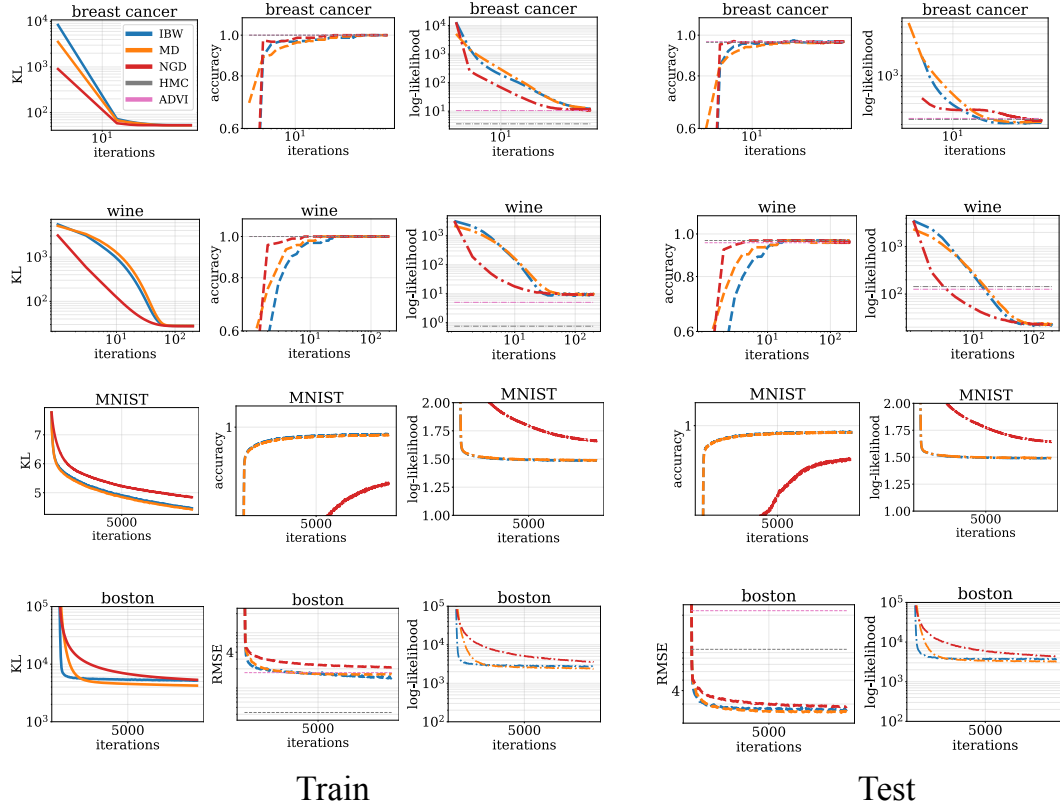


Figure 10: Evolution of the KL, accuracy (or RMSE), and log-likelihood over iterations on the **Train** set (left) and **test** set (Tight) for the `breast_cancer` (upper row), `wine` (middle row), `MNIST` and `boston` (bottom row) datasets for  $N = 5$ .