

# A Fisher-Rao gradient flow for entropic mean-field min-max games

Anonymous authors

Paper under double-blind review

## Abstract

Gradient flows play a substantial role in addressing many machine learning problems. We examine the convergence in continuous-time of a *Fisher-Rao* (Mean-Field Birth-Death) gradient flow in the context of solving convex-concave min-max games with entropy regularization. We propose appropriate Lyapunov functions to demonstrate convergence with explicit rates to the unique mixed Nash equilibrium.

## 1 Introduction

The rapid progress of machine learning (ML) techniques such as Generative Adversarial Networks (GANs) (Goodfellow et al., 2014), adversarial learning (Madry et al., 2018), multi-agent reinforcement learning (Zhang et al., 2021) has propelled a surge of interest in the study of optimization problems on the space of probability measures in recent years. Particularly noteworthy are the numerous works, e.g., (Hsieh et al., 2019; Domingo-Enrich et al., 2020; Wang & Chizat, 2023; Lu, 2023; Trillos & Trillos, 2023; Kim et al., 2024), illustrating how training GANs and addressing adversarial robustness can be cast as min-max games over probability measures. In this setting, understanding the time evolution of the players’ initial strategies to the equilibrium of the game leverages the use of gradient flows on the space of probability measures. A discussion about the applications of gradient flows in optimization and sampling can be found in a recent survey (Trillos et al., 2023).

The Fisher-Rao (FR) gradient flow has been recently studied in the context of mean-field optimization problems (Liu et al., 2023), accelerating Langevin-based sampling from multi-modal distributions (Lu et al., 2019; 2023), and training shallow neural networks in the mean-field regime (Rotskoff et al., 2019). The present paper aims to extend these results and focuses on the continuous-time convergence of the *Fisher-Rao* (FR) gradient flow to the unique mixed Nash equilibrium of an entropy-regularized min-max game.

### 1.1 Notation and setup

For any  $\mathcal{Z} \subseteq \mathbb{R}^d$ , we denote by  $\mathcal{P}_{\text{ac}}(\mathcal{Z})$  the space of probability measures on  $\mathcal{Z}$  which are absolutely continuous with respect to the Lebesgue measure. Following a standard convention, elements in  $\mathcal{P}_{\text{ac}}(\mathcal{Z})$  denote probability measures as well as their densities. Let  $\mathcal{X}, \mathcal{Y} \subseteq \mathbb{R}^d$  and fix two reference probability measures  $\pi(dx) \propto e^{-U^\pi(x)}dx \in \mathcal{P}_{\text{ac}}(\mathcal{X})$  and  $\rho(dy) \propto e^{-U^\rho(y)}dy \in \mathcal{P}_{\text{ac}}(\mathcal{Y})$ , where  $U^\pi : \mathcal{X} \rightarrow \mathbb{R}$  and  $U^\rho : \mathcal{Y} \rightarrow \mathbb{R}$  are two measurable functions. The relative entropy  $D_{\text{KL}}(\cdot|\pi) : \mathcal{P}(\mathcal{X}) \rightarrow [0, \infty]$  with respect to  $\pi$  is given by

$$D_{\text{KL}}(\nu|\pi) = \begin{cases} \int_{\mathcal{X}} \log \left( \frac{\nu(x)}{\pi(x)} \right) \nu(x) dx, & \text{if } \nu \text{ is absolutely continuous with respect to } \pi, \\ +\infty, & \text{otherwise,} \end{cases}$$

and we define  $D_{\text{KL}}(\mu|\rho)$  analogously for any  $\mu \in \mathcal{P}(\mathcal{Y})$ . Let  $F : \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{Y}) \rightarrow \mathbb{R}$  be a convex-concave (possibly non-linear) function and  $\sigma > 0$  be a regularization parameter. We consider the min-max problem

$$\min_{\nu \in \mathcal{P}(\mathcal{X})} \max_{\mu \in \mathcal{P}(\mathcal{Y})} V^\sigma(\nu, \mu), \text{ with } V^\sigma(\nu, \mu) := F(\nu, \mu) + \frac{\sigma^2}{2} (D_{\text{KL}}(\nu|\pi) - D_{\text{KL}}(\mu|\rho)). \quad (1)$$

In order to solve (1), one typically aims to identify *mixed Nash equilibria* (MNEs) (von Neumann et al., 1944; Nash, 1951), characterized by pairs of measures  $(\nu_\sigma^*, \mu_\sigma^*) \in \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{Y})$  that satisfy

$$V^\sigma(\nu_\sigma^*, \mu) \leq V^\sigma(\nu_\sigma^*, \mu_\sigma^*) \leq V^\sigma(\nu, \mu_\sigma^*), \quad \text{for all } (\nu, \mu) \in \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{Y}). \quad (2)$$

It is important to highlight that when  $F$  is bilinear and  $\sigma = 0$ , i.e.,  $V^0(\nu, \mu) = \int_{\mathcal{Y}} \int_{\mathcal{X}} f(x, y) \nu(dx) \mu(dy)$ , for some  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ , measures characterized by (2) represent MNEs in the classical sense of a two-player zero-sum game. Results concerned with the existence and uniqueness of an MNE for (1) are presented in Appendix A in Lascu et al. (2023). Therein, it is also proved that  $V^\sigma$   $\Gamma$ -converges to  $F$  as  $\sigma \downarrow 0$ , under mild assumptions on  $F, \pi$  and  $\rho$ .

In optimization, the monotonic decrease of the objective function along the gradient flow is key to proving convergence. However, for min-max games the monotonic decrease no longer holds due to the conflicting actions of the players. Hence, a suitable Lyapunov function is needed. A common choice is the so-called Nikaidô-Isoda (NI) error (Nikaidô & Isoda, 1955), which, for all  $(\nu, \mu) \in \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{Y})$ , can be defined as

$$\text{NI}(\nu, \mu) := \max_{\mu' \in \mathcal{P}(\mathcal{Y})} V^\sigma(\nu, \mu') - \min_{\nu' \in \mathcal{P}(\mathcal{X})} V^\sigma(\nu', \mu).$$

From the saddle point condition (2), it follows that  $\text{NI}(\nu, \mu) \geq 0$  and  $\text{NI}(\nu, \mu) = 0$  if and only if  $(\nu, \mu)$  is a MNE. We will also consider an alternative Lyapunov function given by (6).

In what follows, we will introduce the FR gradient flow on the space  $(\mathcal{P}_{\text{ac}}(\mathcal{X}) \times \mathcal{P}_{\text{ac}}(\mathcal{Y}), \text{FR})$ , where FR is the Fisher-Rao distance defined by (3).

## 1.2 Fisher-Rao (mean-field birth-death) gradient flow

In this work, we are mainly focusing on the dynamical formulation of the Fisher-Rao metric (see e.g. Gallouët & Monsaingeon (2017); Liu et al. (2023)). For any  $\lambda_0, \lambda_1 \in \mathcal{P}_{\text{ac}}(\mathcal{M})$ , with  $\mathcal{M} \subseteq \mathbb{R}^d$ , the variational representation of the FR distance is given by

$$\text{FR}^2(\lambda_0, \lambda_1) := \inf \left\{ \int_0^1 \int_{\mathcal{M}} \left| r_s(x) - \int_{\mathcal{M}} r_s(y) \lambda_s(dy) \right|^2 \lambda_s(dx) ds : \partial_t \lambda_t = \lambda_t \left( r_t - \int_{\mathcal{M}} r_t(y) \lambda_t(dy) \right) \right\}, \quad (3)$$

where the infimum is taken over all curves  $[0, 1] \ni t \mapsto (\lambda_t, r_t) \in \mathcal{P}_{\text{ac}}(\mathcal{M}) \times L^2(\mathcal{M}; \lambda_t)$  solving

$$\partial_t \lambda_t = \lambda_t \left( r_t - \int_{\mathcal{M}} r_t(y) \lambda_t(dy) \right) \quad (4)$$

in the distributional sense, such that  $t \mapsto \lambda_t$  is weakly continuous with endpoints  $\lambda_0$  and  $\lambda_1$ . The reaction term  $r_t(x) \in \mathbb{R}$  is a scalar that dictates how much mass is created/destroyed at time  $t > 0$  and position  $x \in \mathbb{R}$ . The integral in (4) guarantees that  $\lambda_t$  is a probability measure for all  $t \geq 0$ .

Inspired by Liu et al. (2023), we consider the Fisher-Rao (mean-field birth-death) gradient flow on the space  $(\mathcal{P}_{\text{ac}}(\mathcal{X}) \times \mathcal{P}_{\text{ac}}(\mathcal{Y}), \text{FR})$  in the setting of (1). As opposed to the mean-field Best Response dynamics studied in Lascu et al. (2023), which is another flow that can be used to solve (1), and which relies on introducing a fixed point perspective on min-max games, the FR dynamics utilize a gradient flow  $(\nu_t, \mu_t)_{t \geq 0}$  in the Fisher-Rao geometry. As a first attempt at defining a Fisher-Rao gradient flow for solving (1), consider

$$\begin{cases} \partial_t \nu_t(x) = -\frac{\delta V^\sigma}{\delta \nu}(\nu_t, \mu_t, x) \nu_t(x), \\ \partial_t \mu_t(y) = \frac{\delta V^\sigma}{\delta \mu}(\nu_t, \mu_t, y) \mu_t(y), \end{cases} \quad (5)$$

with initial condition  $(\nu_0, \mu_0) \in \mathcal{P}_{\text{ac}}(\mathcal{X}) \times \mathcal{P}_{\text{ac}}(\mathcal{Y})$ . We adopt the convention that

$$\int_{\mathcal{X}} \frac{\delta V^\sigma}{\delta \nu}(\nu_t, \mu_t, x) \nu_t(dx) = \int_{\mathcal{Y}} \frac{\delta V^\sigma}{\delta \mu}(\nu_t, \mu_t, y) \mu_t(dy) = 0$$

since the flat derivatives of  $V^\sigma$  are uniquely defined up to an additive shift (see Definition B.1 in Section B). Thus, the total mass 1 of probability measures is still preserved along the gradient flow.

### 1.2.1 Sketch of convergence proof for the FR gradient flow

For the sake of presenting an intuitive heuristic argument, we ignore here for now that the *flat derivatives*  $(\nu, \mu, x) \mapsto \frac{\delta V^\sigma}{\delta \nu}(\nu, \mu, x)$  and  $(\nu, \mu, y) \mapsto \frac{\delta V^\sigma}{\delta \mu}(\nu, \mu, y)$  may not exist for the  $V^\sigma$  defined in (1), due to the relative entropy term  $D_{\text{KL}}$  being only lower semicontinuous (for this reason, in our analysis in Section 2, we will replace these two derivatives with appropriately defined auxiliary functions  $a$  and  $b$ ). Nevertheless, we will now demonstrate that choosing the flow  $(\nu_t, \mu_t)_{t \geq 0}$  as in (5), makes the function

$$t \mapsto D_{\text{KL}}(\nu_t^* | \nu_t) + D_{\text{KL}}(\mu_t^* | \mu_t) \quad (6)$$

decrease in  $t$ , under the assumption of convexity-concavity of  $F$ . Let  $(\nu, \mu) \in \mathcal{P}_{\text{ac}}(\mathcal{X}) \times \mathcal{P}_{\text{ac}}(\mathcal{Y})$ . Then assuming the existence of the flow  $(\nu_t, \mu_t)_{t \geq 0}$  satisfying (5), and the differentiability of the map  $t \mapsto D_{\text{KL}}(\nu | \nu_t) + D_{\text{KL}}(\mu | \mu_t)$ , we formally have that

$$\begin{aligned} \frac{d}{dt} (D_{\text{KL}}(\nu | \nu_t) + D_{\text{KL}}(\mu | \mu_t)) &= \int_{\mathcal{X}} \partial_t \left( \nu(x) \log \frac{\nu(x)}{\nu_t(x)} \right) dx + \int_{\mathcal{Y}} \partial_t \left( \mu(y) \log \frac{\mu(y)}{\mu_t(y)} \right) dy \\ &= - \int_{\mathcal{X}} (\nu(x) - \nu_t(x)) \frac{\partial_t \nu_t(x)}{\nu_t(x)} dx - \int_{\mathcal{Y}} (\mu(y) - \mu_t(y)) \frac{\partial_t \mu_t(y)}{\mu_t(y)} dy \\ &= \int_{\mathcal{X}} \frac{\delta V^\sigma}{\delta \nu}(\nu_t, \mu_t, x) (\nu - \nu_t)(dx) - \int_{\mathcal{Y}} \frac{\delta V^\sigma}{\delta \mu}(\nu_t, \mu_t, y) (\mu - \mu_t)(dy), \end{aligned}$$

where the second equality follows from the fact that  $\int \partial_t \nu_t(x) dx = \int \partial_t \mu_t(y) dy = 0$ . Then, assuming that  $\nu \mapsto F(\nu, \mu)$  and  $\mu \mapsto F(\nu, \mu)$  are convex and concave, respectively (see Assumption 1), we observe that  $\nu \mapsto V^\sigma(\nu, \mu)$  and  $\mu \mapsto V^\sigma(\nu, \mu)$  are  $\sigma$ -strongly-convex and  $\sigma$ -strongly-concave relative to  $D_{\text{KL}}$ , respectively (see Lemma 3.1 in Section A), that is

$$\begin{aligned} V^\sigma(\nu, \mu_t) - V^\sigma(\nu_t, \mu_t) &\geq \int_{\mathcal{X}} \frac{\delta V^\sigma}{\delta \nu}(\nu_t, \mu_t, x) (\nu - \nu_t)(dx) + \frac{\sigma^2}{2} D_{\text{KL}}(\nu | \nu_t), \\ V^\sigma(\nu_t, \mu) - V^\sigma(\nu_t, \mu_t) &\leq \int_{\mathcal{Y}} \frac{\delta V^\sigma}{\delta \mu}(\nu_t, \mu_t, y) (\mu - \mu_t)(dy) - \frac{\sigma^2}{2} D_{\text{KL}}(\mu | \mu_t). \end{aligned}$$

Therefore, we obtain that

$$\begin{aligned} \frac{d}{dt} (D_{\text{KL}}(\nu | \nu_t) + D_{\text{KL}}(\mu | \mu_t)) &\leq V^\sigma(\nu, \mu_t) - V^\sigma(\nu_t, \mu_t) + V^\sigma(\nu_t, \mu) - V^\sigma(\nu_t, \mu_t) \\ &\quad - \frac{\sigma^2}{2} D_{\text{KL}}(\nu | \nu_t) - \frac{\sigma^2}{2} D_{\text{KL}}(\mu | \mu_t). \end{aligned} \quad (7)$$

Setting  $(\nu, \mu) = (\nu_\sigma^*, \mu_\sigma^*)$  in (7), using the saddle point condition (2), i.e.,  $V^\sigma(\nu_\sigma^*, \mu_t) - V^\sigma(\nu_t, \mu_\sigma^*) \leq 0$ , and applying Gronwall's inequality gives

$$D_{\text{KL}}(\nu_\sigma^* | \nu_t) + D_{\text{KL}}(\mu_\sigma^* | \mu_t) \leq e^{-\frac{\sigma^2}{2} t} (D_{\text{KL}}(\nu_\sigma^* | \nu_0) + D_{\text{KL}}(\mu_\sigma^* | \mu_0)).$$

Since  $D_{\text{KL}}(\nu_\sigma^* | \nu) + D_{\text{KL}}(\mu_\sigma^* | \mu) \geq 0$  with equality if and only if  $\nu = \nu_\sigma^*$  and  $\mu = \mu_\sigma^*$ , it follows that the unique MNE of (1) is achieved with exponential rate  $e^{-\frac{\sigma^2}{2} t}$ . To prove convergence in terms of the NI error, first observe that since

$$\frac{\sigma^2}{2} (D_{\text{KL}}(\nu | \nu_t) + D_{\text{KL}}(\mu | \mu_t)) \geq 0,$$

for all  $(\nu, \mu)$ , (7) becomes

$$\frac{d}{dt} (D_{\text{KL}}(\nu | \nu_t) + D_{\text{KL}}(\mu | \mu_t)) \leq V^\sigma(\nu, \mu_t) - V^\sigma(\nu_t, \mu).$$

We integrate this inequality from 0 to  $t > 0$ , divide by  $t$  and apply Jensen's inequality to  $\nu \mapsto V^\sigma(\nu, \mu)$  and  $\mu \mapsto V^\sigma(\nu, \mu)$ , respectively. Lastly, maximizing over  $(\nu, \mu)$  gives

$$\text{NI} \left( \frac{1}{t} \int_0^t \nu_s ds, \frac{1}{t} \int_0^t \mu_s ds \right) \leq \frac{1}{t} \left( \max_{\nu} D_{\text{KL}}(\nu | \nu_0) + \max_{\mu} D_{\text{KL}}(\mu | \mu_0) \right).$$

### 1.3 Our contribution

We prove the existence of the FR gradient flow  $(\nu_t, \mu_t)_{t \geq 0}$  and show that it converges with rate  $\mathcal{O}\left(e^{-\frac{\sigma^2}{2}t}\right)$  to the unique MNE of (1) with respect to the Lyapunov function  $t \mapsto D_{\text{KL}}(\nu_t^* | \nu_t) + D_{\text{KL}}(\mu_t^* | \mu_t)$ . We also show that  $\text{NI}\left(\frac{1}{t} \int_0^t \nu_s ds, \frac{1}{t} \int_0^t \mu_s ds\right)$  converges to zero with rate  $\mathcal{O}\left(\frac{1}{t}\right)$ .

### 1.4 Related works

Recent intensive research has been dedicated to examining the convergence of the Wasserstein gradient flow to MNEs within the specific formulation of game (1) in which  $F$  is bilinear ( $F(\nu, \mu) = \int_{\mathcal{Y}} \int_{\mathcal{X}} f(x, y) \nu(dx) \mu(dy)$ ) and regularized by entropy rather than relative entropy  $D_{\text{KL}}$ . The spaces  $\mathcal{X}$  and  $\mathcal{Y}$  are assumed to be either compact smooth manifolds without boundary, embedded in the Euclidean space, or Euclidean tori, while  $f$  exhibits sufficient regularity, typically being at least continuously differentiable with Lipschitz conditions satisfied by  $\nabla_x f$  and  $\nabla_y f$ . This line of research is explored in works such as Domingo-Enrich et al. (2020); Ma & Ying (2022); Lu (2023); Wang & Chizat (2023).

In this context, Ma & Ying (2022); Lu (2023) investigate the convergence of the Wasserstein gradient flow, proving exponential convergence to the MNE when the players' flows  $(\nu_t)_{t \geq 0}$  and  $(\mu_t)_{t \geq 0}$  converge at different rates. In Ma & Ying (2022), the analysis considers the scenario where one player's flow reaches equilibrium while the other remains governed by the Wasserstein gradient flow equation. Notably, Ma & Ying (2022, Theorem 5) shows the convergence (without explicit rate) of the flow  $(\nu_t, \mu_t)_{t \geq 0}$  to the unique MNE under these separated dynamics.

On the other hand, Lu (2023) examines the situation where the players' flows evolve at varying speeds but with finite timescale separation, meaning that neither player has reached equilibrium. Here, Lu (2023, Theorem 2.1) proves exponential convergence of the finitely timescale-separated Wasserstein gradient flow to the unique MNE, with the convergence rate depending upon regularization and timescale separation parameters. In contrast to Ma & Ying (2022); Lu (2023), we prove that the Fisher-Rao gradient flow (5) converges exponentially fast to the unique MNE while players' dynamics converge at the same speed.

Other works focused on the convergence of the Wasserstein-Fisher-Rao (WFR) gradient flow, combining both the Wasserstein gradient flow (allowing particles to diffuse in space) and the Fisher-Rao gradient flow (forcing particles to evade low probability regions). Assuming that  $F$  is bilinear and  $\sigma = 0$ , Domingo-Enrich et al. (2020) investigates the Wasserstein-Fisher-Rao gradient flow's convergence (without giving explicit rates). For  $t_0 > 0$  (dependent on the parameters governing the individual contributions of the Wasserstein and Fisher-Rao components in the WFR gradient flow), and under circumstances where the Fisher-Rao component predominates over the Wasserstein component, Domingo-Enrich et al. (2020, Theorem 2) establishes that the pair  $\left(\frac{1}{t_0} \int_0^{t_0} \nu_s ds, \frac{1}{t_0} \int_0^{t_0} \mu_s ds\right)$  is an  $\epsilon$ -approximate MNE of the game, i.e.,  $\text{NI}\left(\frac{1}{t_0} \int_0^{t_0} \nu_s ds, \frac{1}{t_0} \int_0^{t_0} \mu_s ds\right) \leq \epsilon$ , for any arbitrarily chosen  $\epsilon > 0$ .

Lastly, Wang & Chizat (2023) introduces a proximal point method that can be viewed as a discrete-time counterpart of the WFR gradient flow. Working within the framework of bilinear  $F$ , with  $\sigma = 0$ , and unique MNE, Wang & Chizat (2023, Theorem 2.2) establishes the local exponential convergence of the iterates to the unique MNE of the game with respect to the NI error and the WFR distance, provided that the initialization is done in close vicinity to the MNE.

## 2 Main results

As we explained in the introduction, we study the convergence of the FR gradient flow to the unique MNE of the entropy-regularized two-player zero-sum game given by (1), where  $F : \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{Y}) \rightarrow \mathbb{R}$  is a non-linear function and  $\sigma > 0$ . Throughout the paper, we have the following assumptions.

**Assumption 1** (Convexity-concavity of  $F$ ). *Suppose  $F$  admits first order flat derivatives with respect to both  $\nu$  and  $\mu$  as stated in Definition B.1. Furthermore, suppose that  $F$  is convex in  $\nu$  and concave in  $\mu$ , i.e.,*

for any  $\nu, \nu' \in \mathcal{P}(\mathcal{X})$  and any  $\mu, \mu' \in \mathcal{P}(\mathcal{Y})$ , we have

$$F(\nu', \mu) - F(\nu, \mu) \geq \int_{\mathcal{X}} \frac{\delta F}{\delta \nu}(\nu, \mu, x)(\nu' - \nu)(dx), \quad (8)$$

$$F(\nu, \mu') - F(\nu, \mu) \leq \int_{\mathcal{Y}} \frac{\delta F}{\delta \mu}(\nu, \mu, y)(\mu' - \mu)(dy). \quad (9)$$

**Assumption 2** (Boundedness of first order flat derivatives). *There exist constants  $C_\nu, C_\mu > 0$  such that for all  $(\nu, \mu) \in \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{Y})$  and for all  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ , we have*

$$\left| \frac{\delta F}{\delta \nu}(\nu, \mu, x) \right| \leq C_\nu, \quad \left| \frac{\delta F}{\delta \mu}(\nu, \mu, y) \right| \leq C_\mu.$$

**Assumption 3** (Boundedness of second order flat derivatives). *Suppose  $F$  admits second order flat derivatives and that there exist constants  $C_{\nu, \nu}, C_{\mu, \mu}, C_{\nu, \mu}, C_{\mu, \nu} > 0$  such that for all  $(\nu, \mu) \in \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{Y})$  and for all  $(x, y), (x', y') \in \mathcal{X} \times \mathcal{Y}$ , we have*

$$\left| \frac{\delta^2 F}{\delta \nu^2}(\nu, \mu, x, x') \right| \leq C_{\nu, \nu}, \quad \left| \frac{\delta^2 F}{\delta \mu^2}(\nu, \mu, y, y') \right| \leq C_{\mu, \mu},$$

$$\left| \frac{\delta^2 F}{\delta \nu \delta \mu}(\nu, \mu, y, x) \right| \leq C_{\nu, \mu}, \quad \left| \frac{\delta^2 F}{\delta \mu \delta \nu}(\nu, \mu, x, y) \right| \leq C_{\mu, \nu}.$$

Note that the order of the flat derivatives in  $\nu$  and  $\mu$  can be interchanged due to Lascu et al. (2023, Lemma B.2). Using Assumption 3, it is straightforward to check that there exist constants  $C'_\nu, C'_\mu > 0$  such that for all  $(\nu, \mu) \in \mathcal{P}_{\text{ac}}(\mathcal{X}) \times \mathcal{P}_{\text{ac}}(\mathcal{Y})$ ,  $(\nu', \mu') \in \mathcal{P}_{\text{ac}}(\mathcal{X}) \times \mathcal{P}_{\text{ac}}(\mathcal{Y})$  and all  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ , we have that

$$\left| \frac{\delta F}{\delta \nu}(\nu, \mu, x) - \frac{\delta F}{\delta \nu}(\nu', \mu', x) \right| \leq C'_\nu (\text{TV}(\nu, \nu') + \text{TV}(\mu, \mu')), \quad (10)$$

$$\left| \frac{\delta F}{\delta \mu}(\nu, \mu, y) - \frac{\delta F}{\delta \mu}(\nu', \mu', y) \right| \leq C'_\mu (\text{TV}(\nu, \nu') + \text{TV}(\mu, \mu')). \quad (11)$$

**Remark 2.1.** *An example of a function  $F$  which satisfies Assumptions 1, 2, 3 is  $F(\nu, \mu) = \int_{\mathcal{Y}} \int_{\mathcal{X}} f(x, y) \nu(dx) \mu(dy)$ , for a function  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  which is bounded but possibly non-convex-non-concave. Indeed, Assumption 1 is trivially satisfied by such  $F$ , while Assumptions 2 and 3 hold due to the boundedness of  $f$ . This type of objective function  $F$  is prototypical in applications including the training of GANs (see, e.g., Arjovsky et al. (2017); Hsieh et al. (2019)) and distributionally robust optimization (see, e.g, Madry et al. (2018); Sinha et al. (2018)).*

**Assumption 4** (Ratio condition). *Suppose  $(\nu_0, \mu_0) \in \mathcal{P}(\mathcal{X}) \times \mathcal{P}(\mathcal{Y})$  are absolutely continuous and comparable with  $\pi$  and  $\rho$ , respectively, in the sense that*

1. *There exist constants  $r_\nu, r_\mu > 0$  such that*

$$\inf_{x \in \mathcal{X}} \frac{\nu_0(x)}{\pi(x)} \geq r_\nu, \quad \inf_{y \in \mathcal{Y}} \frac{\mu_0(y)}{\rho(y)} \geq r_\mu. \quad (12)$$

2. *There exist constants  $R_\nu, R_\mu > 1$  such that*

$$\sup_{x \in \mathcal{X}} \frac{\nu_0(x)}{\pi(x)} \leq R_\nu, \quad \sup_{y \in \mathcal{Y}} \frac{\mu_0(y)}{\rho(y)} \leq R_\mu. \quad (13)$$

It can be proved (see, e.g., Lascu et al. (2023, Proposition A.1)) that the MNE  $(\nu_\sigma^*, \mu_\sigma^*)$  of 1 is given implicitly by the equations

$$\nu_\sigma^*(x) = \frac{1}{Z(\nu_\sigma^*, \mu_\sigma^*)} \exp \left( -\frac{2}{\sigma^2} \frac{\delta F}{\delta \nu}(\nu_\sigma^*, \mu_\sigma^*, x) - U^\pi(x) \right), \quad (14)$$

$$\mu_\sigma^*(y) = \frac{1}{Z'(\nu_\sigma^*, \mu_\sigma^*)} \exp \left( \frac{2}{\sigma^2} \frac{\delta F}{\delta \mu}(\nu_\sigma^*, \mu_\sigma^*, y) - U^\rho(y) \right), \quad (15)$$

where  $Z(\nu_\sigma^*, \mu_\sigma^*)$  and  $Z'(\nu_\sigma^*, \mu_\sigma^*)$  are normalizing constants.

Combining (14) and (15) and Assumption 2, we deduce that Assumption 4 is equivalent to assuming that there exist constants  $\bar{r}_\nu, \bar{r}_\mu > 0$  and  $\bar{R}_\nu, \bar{R}_\mu > 1$  such that for all  $(x, y) \in \mathcal{X} \times \mathcal{Y}$ ,

$$\bar{r}_\nu \leq \frac{\nu_0(x)}{\nu_\sigma^*(x)} \leq \bar{R}_\nu, \quad \bar{r}_\mu \leq \frac{\mu_0(y)}{\mu_\sigma^*(y)} \leq \bar{R}_\mu. \quad (16)$$

We emphasize that Assumption 4 is natural in the context of Fisher-Rao flows. From (5), we observe that the support of the measures along the gradient flow does not increase. Thus, it is essential that the initialization is comparable with the MNE (cf. (16)) as reflected in Assumption 4. It is observed in Liu et al. (2023), that Assumption 4 can be understood as a “warm start” type of condition.

Returning to the question of flat differentiability of  $V^\sigma$ , which was raised in Subsection 1.2, if we assume that  $F$  is flat differentiable with respect to both  $\nu$  and  $\mu$  (see Assumption 1), then the maps  $(\nu, \mu, x) \mapsto a(\nu, \mu, x) := \frac{\delta F}{\delta \nu}(\nu, \mu, x) + \frac{\sigma^2}{2} \log \left( \frac{\nu(x)}{\pi(x)} \right) - \frac{\sigma^2}{2} \text{D}_{\text{KL}}(\nu|\pi)$  and  $(\nu, \mu, y) \mapsto b(\nu, \mu, y) := \frac{\delta F}{\delta \mu}(\nu, \mu, y) - \frac{\sigma^2}{2} \log \left( \frac{\mu(y)}{\rho(y)} \right) + \frac{\sigma^2}{2} \text{D}_{\text{KL}}(\mu|\rho)$  are well-defined and formally correspond to the flat derivatives  $\frac{\delta V^\sigma}{\delta \nu}(\nu, \mu, \cdot)$  and  $\frac{\delta V^\sigma}{\delta \mu}(\nu, \mu, \cdot)$ , respectively, for those measures  $\nu$  and  $\mu$  for which such derivatives exist (note that we will only need to consider  $a$  and  $b$  along our gradient flow  $(\nu_t, \mu_t)_{t \geq 0}$ , so our argument can be interpreted as stating that, while  $V^\sigma$  is not flat differentiable everywhere, it is indeed flat differentiable along our gradient flow). The relative entropy terms  $\text{D}_{\text{KL}}$  appear in the definition of  $a$  and  $b$  as normalizing constants to ensure that  $\int_{\mathcal{X}} a(\nu, \mu, x) \nu(dx) = 0$  and  $\int_{\mathcal{Y}} b(\nu, \mu, y) \mu(dy) = 0$ , since we adopt the convention that the flat derivatives of  $F$  are uniquely defined up to an additive shift (see Definition B.1 in Section B). Motivated by this discussion, we define the Fisher-Rao gradient flow  $(\nu_t, \mu_t)_{t \geq 0}$  on the space  $(\mathcal{P}_{\text{ac}}(\mathcal{X}) \times \mathcal{P}_{\text{ac}}(\mathcal{Y}), \text{FR})$  by

$$\begin{cases} \partial_t \nu_t(x) = -a(\nu_t, \mu_t, x) \nu_t(x), \\ \partial_t \mu_t(y) = b(\nu_t, \mu_t, y) \mu_t(y), \end{cases} \quad (17)$$

with initial condition  $(\nu_0, \mu_0) \in \mathcal{P}_{\text{ac}}(\mathcal{X}) \times \mathcal{P}_{\text{ac}}(\mathcal{Y})$ .

The following result extends Liu et al. (2023, Theorem 2.1) to the case of two-player zero-sum games by showing that the Fisher-Rao gradient flow (17) admits a unique solution  $(\nu_t, \mu_t)_{t \geq 0}$ .

**Theorem 2.2.** *Suppose that Assumption 2, 3 and condition (13) from Assumption 4 hold. Then the system of equations (17) has a unique solution  $(\nu_t, \mu_t)_{t \geq 0}$ . Moreover, for  $t \geq 0$ ,*

$$\text{D}_{\text{KL}}(\nu_t|\pi) \leq 2 \log R_\nu + \frac{4}{\sigma^2} C_\nu, \quad \text{D}_{\text{KL}}(\mu_t|\rho) \leq 2 \log R_\mu + \frac{4}{\sigma^2} C_\mu \quad (18)$$

and there exist constants  $R_{1,\nu}, R_{1,\mu} > 1$  such that for all  $t \geq 0$ ,

$$\sup_{x \in \mathcal{X}} \frac{\nu_t(x)}{\pi(x)} \leq R_{1,\nu}, \quad \sup_{y \in \mathcal{Y}} \frac{\mu_t(y)}{\rho(y)} \leq R_{1,\mu}. \quad (19)$$

Additionally, if condition (12) from Assumption 4 holds, then there exist constants  $r_{1,\nu}, r_{1,\mu} > 0$  such that for all  $t \geq 0$ ,

$$\inf_{x \in \mathcal{X}} \frac{\nu_t(x)}{\pi(x)} \geq r_{1,\nu}, \quad \inf_{y \in \mathcal{Y}} \frac{\mu_t(y)}{\rho(y)} \geq r_{1,\mu}. \quad (20)$$

The bounds obtained in this theorem allow us to prove that the map  $t \mapsto \text{D}_{\text{KL}}(\nu_\sigma^*|\nu_t) + \text{D}_{\text{KL}}(\mu_\sigma^*|\mu_t)$  is differentiable along the FR dynamics (17).

Next, we state the other main result of this paper. We prove two different types of convergence results for the FR dynamics (17) in min-max games given by (1), one in terms of the players’ strategies and the other in terms of the payoff function. Whereas the proof of the existence of the flow (Theorem 2.2) follows a

route similar to Liu et al. (2023), the proof of convergence in Theorem 2.3 is significantly different from Liu et al. (2023) since, as indicated in the introduction, in the present paper, convergence has to be studied with respect to appropriately chosen Lyapunov functions and, moreover, we do not rely on the Polyak-Łojasiewicz inequality.

**Theorem 2.3.** *Suppose that Assumption 2, 3 and 4 hold and let  $(\nu, \mu) \in \mathcal{P}_{ac}(\mathcal{X}) \times \mathcal{P}_{ac}(\mathcal{Y})$ . Then the map  $t \mapsto D_{KL}(\nu|\nu_t) + D_{KL}(\mu|\mu_t)$  is differentiable along the FR dynamics (17). Suppose furthermore that Assumption 1 holds and that  $(\nu_0, \mu_0)$  are chosen such that  $\max_{\nu \in \mathcal{P}_{ac}(\mathcal{X})} D_{KL}(\nu|\nu_0) + \max_{\mu \in \mathcal{P}_{ac}(\mathcal{Y})} D_{KL}(\mu|\mu_0) < \infty$ . Then, for all  $t \geq 0$ , we have*

$$D_{KL}(\nu_\sigma^*|\nu_t) + D_{KL}(\mu_\sigma^*|\mu_t) \leq e^{-\frac{\sigma^2}{2}t} (D_{KL}(\nu_\sigma^*|\nu_0) + D_{KL}(\mu_\sigma^*|\mu_0)),$$

$$NI \left( \frac{1}{t} \int_0^t \nu_s ds, \frac{1}{t} \int_0^t \mu_s ds \right) \leq \frac{1}{t} \left( \max_{\nu \in \mathcal{P}_{ac}(\mathcal{X})} D_{KL}(\nu|\nu_0) + \max_{\mu \in \mathcal{P}_{ac}(\mathcal{Y})} D_{KL}(\mu|\mu_0) \right). \quad (21)$$

In game theoretic language, the first result of Theorem 2.3 says that convergence of the FR dynamics (17) to the unique MNE  $(\nu_\sigma^*, \mu_\sigma^*)$  in terms of the strategies  $\nu_t$  and  $\mu_t$  is achieved with exponential rate depending on  $\sigma$ . On the other hand, the second result gives a convergence estimate in terms of the NI error, which in turn is expressed as a function of the payoff  $V^\sigma$ . More precisely, the NI error of the time-averaged flow converges to 0 with rate  $\mathcal{O}(1/t)$  regardless of the value of  $\sigma$ .

**Remark 2.4.** *Exponential convergence of the single mean-field birth-death flow  $(\nu_t)_{t \geq 0}$  with respect to  $D_{KL}(\nu_\sigma^*|\nu_t)$  can be shown to also hold in the setting of Liu et al. (2023), however it was not studied in Liu et al. (2023), which considered only convergence of  $V^\sigma(\nu_t)$  for a convex energy function  $V^\sigma : \mathcal{P}(\mathbb{R}^d) \rightarrow \mathbb{R}$ .*

**Remark 2.5.** *The FR gradient flow can also be interpreted as the continuous-time limit of mirror descent (MD) algorithm with  $D_{KL}$  divergence over the space of probability measures (see, e.g., Liu et al. (2023)). In the context of a minimization problem, assuming convexity of the objective function, Aubin-Frankowski et al. (2022, Theorem 4) proves that the MD algorithm converges with rate  $\mathcal{O}(\frac{1}{N})$  to the minimizer, where  $N$  is the algorithm's total number of iterations.*

*It is crucial to emphasize that the rate  $\mathcal{O}(\frac{1}{t})$  obtained in Theorem 2.3 is not immediately implicit from the optimization result. Although one might expect that the mirror descent-ascent (MDA) algorithm for a min-max game with convex-concave payoff function converges with rate  $\mathcal{O}(\frac{1}{N})$ , this is not the case because, as stressed in Lascu et al. (2024), the dynamics of the game are greatly influenced by the order the players move: either simultaneously (players move at the same time) or sequentially (each player moves upon observing the opponents' moves). The continuous-time analysis does not capture these two situations. Consequently, Lascu et al. (2024, Theorem 2.6, 2.16) show that the convergence rates of the simultaneous and sequential MDA schemes to mixed Nash equilibria of convex-concave min-max games, measured in the NI error, are  $\mathcal{O}(N^{-1/2})$  and  $\mathcal{O}(N^{-2/3})$ , respectively.*

**Remark 2.6.** *It is worth commenting that the second result (21) of Theorem 2.3 can be obtained by only requiring  $F$  to be convex-concave (see Assumption 1) and setting  $\sigma = 0$ , i.e., the proof of (21) does not use the fact that  $V^\sigma$  is regularized by  $D_{KL}$  with a positive  $\sigma$ . Although we are presently unaware of a way to prove the existence of the FR flow (17) for  $\sigma = 0$  using the Picard iteration technique we employed in Theorem 2.2, we believe that a possible way to prove the existence of the FR flow for  $\sigma = 0$  could be via a minimizing movement scheme approach in the spirit of Gallouët & Monsaingeon (2017). Indeed, the FR gradient flow can be viewed as the zero-step-size limit of a proximal gradient algorithm under the Fisher-Rao distance (see equation (4.5) in Gallouët & Monsaingeon (2017)). Furthermore, for absolutely continuous measures, one could explicitly write the geodesics minimizing (3) (see Subsection 2.2 in Gallouët & Monsaingeon (2017)) along which the iterates generated by the proximal gradient algorithm can be interpolated.*

### 3 Proof of Theorem 2.3

Before we present the proof of Theorem 2.3, we begin with a technical lemma, which extends Liu et al. (2023, Lemma 2.5) to the min-max games setup, and which is proved in Section A. The proof of Theorem 2.3 is done in two steps:

- First, we show that, for fixed  $(\nu, \mu)$ , the map  $t \mapsto D_{\text{KL}}(\nu|\nu_t) + D_{\text{KL}}(\mu|\mu_t)$  is differentiable when  $(\nu_t, \mu_t)_{t \geq 0}$  is defined as the FR flow (17).
- Second, we show that  $\frac{d}{dt} (D_{\text{KL}}(\nu_\sigma^*|\nu_t) + D_{\text{KL}}(\mu_\sigma^*|\mu_t))$  is bounded from above by  $-\frac{\sigma^2}{2} (D_{\text{KL}}(\nu_\sigma^*|\nu_t) + D_{\text{KL}}(\mu_\sigma^*|\mu_t))$ , and then we apply Gronwall's inequality to obtain exponential convergence. Subsequently, we establish linear convergence for  $t \mapsto \text{NI} \left( \frac{1}{t} \int_0^t \nu_s ds, \frac{1}{t} \int_0^t \mu_s ds \right)$ .

**Lemma 3.1** (Relative  $\sigma$ -strong-convexity-concavity to  $D_{\text{KL}}$ ). *For  $V^\sigma$  given by (1), if Assumption 1 holds, then  $V^\sigma$  satisfies the following inequalities for all  $(\nu, \mu), (\nu', \mu') \in \mathcal{P}_{\text{ac}}(\mathcal{X}) \times \mathcal{P}_{\text{ac}}(\mathcal{Y})$ :*

$$\begin{aligned} V^\sigma(\nu', \mu) - V^\sigma(\nu, \mu) &\geq \int_{\mathcal{X}} a(\nu, \mu, x)(\nu' - \nu)(dx) + \frac{\sigma^2}{2} D_{\text{KL}}(\nu'|\nu), \\ V^\sigma(\nu, \mu') - V^\sigma(\nu, \mu) &\leq \int_{\mathcal{Y}} b(\nu, \mu, y)(\mu' - \mu)(dy) - \frac{\sigma^2}{2} D_{\text{KL}}(\mu'|\mu). \end{aligned}$$

*Proof of Theorem 2.3. Step 1: Differentiability of  $D_{\text{KL}}$  with respect to the FR flow (17):* Suppose that Assumption 2, 3 and 4 hold. In order to show the differentiability of  $t \mapsto D_{\text{KL}}(\nu|\nu_t)$  for fixed  $\nu \in \mathcal{P}_{\text{ac}}(\mathcal{X})$  with respect to (17), it suffices to show that there exists an integrable function  $f : \mathcal{X} \rightarrow \mathbb{R}$  such that

$$\left| \partial_t \left( \nu(x) \log \frac{\nu(x)}{\nu_t(x)} \right) \right| \leq f(x),$$

for all  $t \geq 0$ . Indeed, using (17) and (33), we have that

$$\begin{aligned} \left| \partial_t \left( \nu(x) \log \frac{\nu(x)}{\nu_t(x)} \right) \right| &= \left| \nu(x) \frac{\partial_t \nu_t(x)}{\nu_t(x)} \right| = \left| \nu(x) \left( \frac{\delta F}{\delta \nu}(\nu_t, \mu_t, x) + \frac{\sigma^2}{2} \log \left( \frac{\nu_t(x)}{\pi(x)} \right) - \frac{\sigma^2}{2} D_{\text{KL}}(\nu_t|\pi) \right) \right| \\ &\leq \left( 3C_\nu + \frac{\sigma^2}{2} (\max\{|\log r_{1,\nu}|, \log R_{1,\nu}\} + 2 \log R_\nu) \right) \nu(x) := f(x). \end{aligned}$$

An identical argument gives the differentiability of  $t \mapsto D_{\text{KL}}(\mu|\mu_t)$  for fixed  $\mu \in \mathcal{P}_{\text{ac}}(\mathcal{Y})$ . Then, we have that the map  $t \mapsto D_{\text{KL}}(\nu|\nu_t) + D_{\text{KL}}(\mu|\mu_t)$  is differentiable.

*Step 2: Convergence of the FR flow:* Since  $t \mapsto D_{\text{KL}}(\nu|\nu_t) + D_{\text{KL}}(\mu|\mu_t)$  is differentiable, we have that

$$\begin{aligned} \frac{d}{dt} (D_{\text{KL}}(\nu|\nu_t) + D_{\text{KL}}(\mu|\mu_t)) &= \int_{\mathcal{X}} \partial_t \left( \nu(x) \log \frac{\nu(x)}{\nu_t(x)} \right) dx + \int_{\mathcal{Y}} \partial_t \left( \mu(y) \log \frac{\mu(y)}{\mu_t(y)} \right) dy \\ &= - \int_{\mathcal{X}} (\nu(x) - \nu_t(x)) \frac{\partial_t \nu_t(x)}{\nu_t(x)} dx - \int_{\mathcal{Y}} (\mu(y) - \mu_t(y)) \frac{\partial_t \mu_t(y)}{\mu_t(y)} dy \\ &= \int_{\mathcal{X}} a(\nu_t, \mu_t, x)(\nu - \nu_t)(dx) - \int_{\mathcal{Y}} b(\nu_t, \mu_t, y)(\mu - \mu_t)(dy), \end{aligned}$$

where in the second equality we used the fact that  $\int_{\mathcal{X}} \partial_t \nu_t(x) dx = \int_{\mathcal{Y}} \partial_t \mu_t(y) dy = 0$ .

If  $\sigma > 0$  and Assumption 1 holds then, using Lemma 3.1 in the equation above, we obtain

$$\frac{d}{dt} (D_{\text{KL}}(\nu|\nu_t) + D_{\text{KL}}(\mu|\mu_t)) \leq V^\sigma(\nu, \mu_t) - V^\sigma(\nu_t, \mu) - \frac{\sigma^2}{2} D_{\text{KL}}(\nu|\nu_t) - \frac{\sigma^2}{2} D_{\text{KL}}(\mu|\mu_t). \quad (22)$$

Setting  $(\nu, \mu) = (\nu_\sigma^*, \mu_\sigma^*)$  in (22) and using the saddle point condition (2) gives

$$\frac{d}{dt} (D_{\text{KL}}(\nu_\sigma^*|\nu_t) + D_{\text{KL}}(\mu_\sigma^*|\mu_t)) \leq -\frac{\sigma^2}{2} D_{\text{KL}}(\nu_\sigma^*|\nu_t) - \frac{\sigma^2}{2} D_{\text{KL}}(\mu_\sigma^*|\mu_t).$$

Hence, by Gronwall's inequality, we obtain that

$$D_{\text{KL}}(\nu_\sigma^*|\nu_t) + D_{\text{KL}}(\mu_\sigma^*|\mu_t) \leq e^{-\frac{\sigma^2}{2}t} (D_{\text{KL}}(\nu_\sigma^*|\nu_0) + D_{\text{KL}}(\mu_\sigma^*|\mu_0)).$$



On the other hand, since  $\frac{\sigma^2}{2} (\mathrm{D}_{\mathrm{KL}}(\nu|\nu_t) + \mathrm{D}_{\mathrm{KL}}(\mu|\mu_t)) \geq 0$ , for all  $(\nu, \mu)$ , (22) becomes

$$\frac{d}{dt} (\mathrm{D}_{\mathrm{KL}}(\nu|\nu_t) + \mathrm{D}_{\mathrm{KL}}(\mu|\mu_t)) \leq V^\sigma(\nu, \mu_t) - V^\sigma(\nu_t, \mu).$$

Hence, integrating this inequality from 0 to  $t > 0$  and dividing by  $t$ , it follows that

$$\begin{aligned} \frac{1}{t} \int_0^t (V^\sigma(\nu_s, \mu) - V^\sigma(\nu, \mu_s)) ds &\leq \frac{1}{t} (\mathrm{D}_{\mathrm{KL}}(\nu|\nu_0) + \mathrm{D}_{\mathrm{KL}}(\mu|\mu_0) - \mathrm{D}_{\mathrm{KL}}(\nu|\nu_t) - \mathrm{D}_{\mathrm{KL}}(\mu|\mu_t)) \\ &\leq \frac{1}{t} (\mathrm{D}_{\mathrm{KL}}(\nu|\nu_0) + \mathrm{D}_{\mathrm{KL}}(\mu|\mu_0)). \end{aligned} \quad (23)$$

Then, using the fact that  $\nu \mapsto V^\sigma(\nu, \mu)$  is convex and  $\mu \mapsto V^\sigma(\nu, \mu)$  is concave, it follows from Jensen's inequality that

$$\frac{1}{t} \int_0^t V^\sigma(\nu_s, \mu) ds - \frac{1}{t} \int_0^t V^\sigma(\nu, \mu_s) ds \geq V^\sigma\left(\frac{1}{t} \int_0^t \nu_s ds, \mu\right) - V^\sigma\left(\nu, \frac{1}{t} \int_0^t \mu_s ds\right). \quad (24)$$

Combining (23) with (24) and taking maximum over  $(\nu, \mu)$ , we obtain

$$\mathrm{NI}\left(\frac{1}{t} \int_0^t \nu_s ds, \frac{1}{t} \int_0^t \mu_s ds\right) \leq \frac{1}{t} \left( \max_{\nu} \mathrm{D}_{\mathrm{KL}}(\nu|\nu_0) + \max_{\mu} \mathrm{D}_{\mathrm{KL}}(\mu|\mu_0) \right).$$

□

## References

- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 214–223. PMLR, 06–11 Aug 2017.
- Pierre-Cyril Aubin-Frankowski, Anna Korba, and Flavien Léger. Mirror descent with relative smoothness in measure spaces, with application to sinkhorn and em. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 17263–17275. Curran Associates, Inc., 2022.
- René A. Carmona and François Delarue. *Probabilistic Theory of Mean Field Games with Applications I: Mean Field FBSDEs, Control, and Games*. Springer International Publishing, 2018.
- Carles Domingo-Enrich, Samy Jelassi, Arthur Mensch, Grant Rotskoff, and Joan Bruna. A mean-field analysis of two-player zero-sum games. In *Advances in Neural Information Processing Systems*, volume 33, pp. 20215–20226. Curran Associates, Inc., 2020.
- Thomas O. Gallouët and Léonard Monsaingeon. A JKO splitting scheme for Kantorovich-Fisher-Rao gradient flows. *SIAM J. Math. Anal.*, 49(2):1100–1130, 2017. ISSN 0036-1410. doi: 10.1137/16M106666X.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- Ya-Ping Hsieh, Chen Liu, and Volkan Cevher. Finding mixed Nash equilibria of generative adversarial networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2810–2819. PMLR, 09–15 Jun 2019.
- Juno Kim, Kakei Yamamoto, Kazusato Oko, Zhuoran Yang, and Taiji Suzuki. Symmetric mean-field langevin dynamics for distributional minimax problems. In *The Twelfth International Conference on Learning Representations*, 2024.

- Razvan-Andrei Lascu, Mateusz B. Majka, and Łukasz Szpruch. Entropic mean-field min-max problems via best response flow, 2023. arXiv:2306.03033.
- Razvan-Andrei Lascu, Mateusz B. Majka, and Łukasz Szpruch. Mirror descent-ascent for mean-field min-max problems, 2024. arXiv:2402.08106.
- Linshan Liu, Mateusz B. Majka, and Łukasz Szpruch. Polyak–Łojasiewicz inequality on the space of measures and convergence of mean-field birth-death processes. *Applied Mathematics and Optimization*, 87(3), March 2023. ISSN 0095-4616. doi: 10.1007/s00245-022-09962-0.
- Yulong Lu. Two-scale gradient descent ascent dynamics finds mixed nash equilibria of continuous games: A mean-field perspective. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- Yulong Lu, Jianfeng Lu, and James Nolen. Accelerating Langevin Sampling with Birth-death, May 2019. arXiv:1905.09863.
- Yulong Lu, Dejan Slepčev, and Lihan Wang. Birth-death dynamics for sampling: global convergence, approximations and their asymptotics. *Nonlinearity*, 36:5731–5772, 09 2023.
- Chao Ma and Lexing Ying. Provably convergent quasistatic dynamics for mean-field two-player zero-sum games. In *International Conference on Learning Representations*, 2022.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.
- John Nash. Non-cooperative games. *Annals of Mathematics*, 54(2):286–295, 1951. ISSN 0003486X.
- Hukukane Nikaidô and Kazuo Isoda. Note on non-cooperative convex game. *Pacific Journal of Mathematics*, 5:807–815, 1955.
- Grant Rotskoff, Samy Jelassi, Joan Bruna, and Eric Vanden-Eijnden. Neuron birth-death dynamics accelerates gradient descent and converges asymptotically. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 5508–5517. PMLR, 09–15 Jun 2019.
- Aman Sinha, Hongseok Namkoong, and John Duchi. Certifiable distributional robustness with principled adversarial training. In *International Conference on Learning Representations*, 2018.
- Camilo Garcia Trillos and Nicolas Garcia Trillos. On adversarial robustness and the use of wasserstein ascent-descent dynamics to enforce it, 2023. arXiv:2301.03662.
- Nicolás García Trillos, Bahram Hosseini, and Daniel Sanz-Alonso. From optimization to sampling through gradient flows. *Notices of the American Mathematical Society*, 2023.
- A.B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Series in Statistics. Springer New York, 2008. ISBN 9780387790527.
- John von Neumann, Oskar Morgenstern, and Ariel Rubinstein. *Theory of Games and Economic Behavior (60th Anniversary Commemorative Edition)*. Princeton University Press, 1944. ISBN 9780691130613.
- David Šiška and Łukasz Szpruch. Gradient Flows for Regularized Stochastic Control Problems. *arXiv e-prints*, art. arXiv:2006.05956, June 2020.
- Guillaume Wang and Lénaïc Chizat. An exponentially converging particle method for the mixed nash equilibrium of continuous games, 2023. arXiv:2211.01280.
- Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms, 2021. arXiv:1911.10635.

## A Technical results and proofs

In this section, we present the proofs of the remaining results formulated in Section 2 of the paper.

### A.1 Proof of Lemma 3.1.

In this subsection, we present the proof of Lemma 3.1.

*Proof of Lemma 3.1.* Using (8) from Assumption 1, it follows that

$$\begin{aligned}
V^\sigma(\nu', \mu) - V^\sigma(\nu, \mu) &\geq \int_{\mathcal{X}} \frac{\delta F}{\delta \nu}(\nu, \mu, x)(\nu' - \nu)(dx) + \frac{\sigma^2}{2} D_{\text{KL}}(\nu'|\pi) - \frac{\sigma^2}{2} D_{\text{KL}}(\nu|\pi) \\
&= \int_{\mathcal{X}} \left( \frac{\delta F}{\delta \nu}(\nu, \mu, x) + \frac{\sigma^2}{2} \log \left( \frac{\nu(x)}{\pi(x)} \right) \right) (\nu' - \nu)(dx) - \frac{\sigma^2}{2} \int_{\mathcal{X}} \log \left( \frac{\nu(x)}{\pi(x)} \right) (\nu' - \nu)(dx) \\
&\quad + \frac{\sigma^2}{2} \int_{\mathcal{X}} \log \left( \frac{\nu'(x)}{\pi(x)} \right) \nu'(dx) - \frac{\sigma^2}{2} \int_{\mathcal{X}} \log \left( \frac{\nu(x)}{\pi(x)} \right) \nu(dx) \\
&= \int_{\mathcal{X}} \left( \frac{\delta F}{\delta \nu}(\nu, \mu, x) + \frac{\sigma^2}{2} \log \left( \frac{\nu(x)}{\pi(x)} \right) \right) (\nu' - \nu)(dx) + \frac{\sigma^2}{2} D_{\text{KL}}(\nu'|\nu) \\
&= \int_{\mathcal{X}} a(\nu, \mu, x)(\nu' - \nu)(dx) + \frac{\sigma^2}{2} D_{\text{KL}}(\nu'|\nu).
\end{aligned}$$

Similarly, using (9) from Assumption 1, it follows that

$$\begin{aligned}
V^\sigma(\nu, \mu') - V^\sigma(\nu, \mu) &\leq \int_{\mathcal{Y}} \frac{\delta F}{\delta \mu}(\nu, \mu, y)(\mu' - \mu)(dy) - \frac{\sigma^2}{2} D_{\text{KL}}(\mu'|\rho) + \frac{\sigma^2}{2} D_{\text{KL}}(\mu|\rho) \\
&= \int_{\mathcal{Y}} \left( \frac{\delta F}{\delta \mu}(\nu, \mu, y) - \frac{\sigma^2}{2} \log \left( \frac{\mu(y)}{\rho(y)} \right) \right) (\mu' - \mu)(dy) + \frac{\sigma^2}{2} \int_{\mathcal{Y}} \log \left( \frac{\mu(y)}{\rho(y)} \right) (\mu' - \mu)(dy) \\
&\quad - \frac{\sigma^2}{2} \int_{\mathcal{Y}} \log \left( \frac{\mu'(y)}{\rho(y)} \right) \mu'(dy) + \frac{\sigma^2}{2} \int_{\mathcal{Y}} \log \left( \frac{\mu(y)}{\rho(y)} \right) \mu(dy) \\
&= \int_{\mathcal{Y}} \left( \frac{\delta F}{\delta \mu}(\nu, \mu, y) - \frac{\sigma^2}{2} \log \left( \frac{\mu(y)}{\rho(y)} \right) \right) (\mu' - \mu)(dy) - \frac{\sigma^2}{2} D_{\text{KL}}(\mu'|\mu) \\
&= \int_{\mathcal{Y}} b(\nu, \mu, y)(\mu' - \mu)(dy) - \frac{\sigma^2}{2} D_{\text{KL}}(\mu'|\mu).
\end{aligned}$$

□

### A.2 Existence and uniqueness of the FR flow

In this subsection, we present the proof of our main result concerning the existence and uniqueness of the Fisher-Rao (FR) flow, i.e., Theorem 2.2. We construct a Picard iteration which is proved to be well-defined in Lemma A.1. Lemma A.2 shows that the Picard iteration is contractive in an appropriate metric. Then in order to conclude the proof of Theorem 2.2 we show the ratio condition (19).

The proof of Theorem 2.2 is an adaptation of the proof of Liu et al. (2023, Theorem 2.1) to the min-max setting (1).

*Proof of Theorem 2.2. Step 1: Existence of the gradient flow and bound (18) on  $[0, T]$ .* In order to prove the existence of a solution  $(\nu_t, \mu_t)_{t \geq 0}$  to

$$\begin{cases} \partial_t \nu_t(x) = - \left( \frac{\delta F}{\delta \nu}(\nu_t, \mu_t, x) + \frac{\sigma^2}{2} \log \left( \frac{\nu_t(x)}{\pi(x)} \right) - \frac{\sigma^2}{2} D_{\text{KL}}(\nu_t|\pi) \right) \nu_t(x), \\ \partial_t \mu_t(y) = \left( \frac{\delta F}{\delta \mu}(\nu_t, \mu_t, y) - \frac{\sigma^2}{2} \log \left( \frac{\mu_t(y)}{\rho(y)} \right) + \frac{\sigma^2}{2} D_{\text{KL}}(\mu_t|\rho) \right) \mu_t(y), \end{cases} \quad (25)$$

we first notice that (25) is equivalent to

$$\begin{cases} \partial_t \log \nu_t(x) = - \left( \frac{\delta F}{\delta \nu}(\nu_t, \mu_t, x) + \frac{\sigma^2}{2} \log \left( \frac{\nu_t(x)}{\pi(x)} \right) - \frac{\sigma^2}{2} D_{\text{KL}}(\nu_t | \pi) \right), \\ \partial_t \log \mu_t(y) = \left( \frac{\delta F}{\delta \mu}(\nu_t, \mu_t, y) - \frac{\sigma^2}{2} \log \left( \frac{\mu_t(y)}{\rho(y)} \right) + \frac{\sigma^2}{2} D_{\text{KL}}(\mu_t | \rho) \right). \end{cases} \quad (26)$$

By Duhamel's formula, (26) is equivalent to

$$\begin{cases} \log \nu_t(x) = e^{-\frac{\sigma^2}{2}t} \log \nu_0(x) - \int_0^t \frac{\sigma^2}{2} e^{-\frac{\sigma^2}{2}(t-s)} \left( \frac{2}{\sigma^2} \frac{\delta F}{\delta \nu}(\nu_s, \mu_s, x) - \log \pi(x) - D_{\text{KL}}(\nu_s | \pi) \right) ds, \\ \log \mu_t(y) = e^{-\frac{\sigma^2}{2}t} \log \mu_0(y) + \int_0^t \frac{\sigma^2}{2} e^{-\frac{\sigma^2}{2}(t-s)} \left( \frac{2}{\sigma^2} \frac{\delta F}{\delta \mu}(\nu_s, \mu_s, y) + \log \rho(y) + D_{\text{KL}}(\mu_s | \rho) \right) ds. \end{cases}$$

Based on these formulas, we will define a Picard iteration scheme. To this end, let us first fix  $T > 0$  and choose a pair of flows of probability measures  $(\nu_t^{(0)}, \mu_t^{(0)})_{t \in [0, T]}$  such that

$$\int_0^T D_{\text{KL}}(\nu_s^{(0)} | \pi) ds < \infty, \quad \int_0^T D_{\text{KL}}(\mu_s^{(0)} | \rho) ds < \infty.$$

For each  $n \geq 1$ , we fix  $\nu_0^{(n)} = \nu_0^{(0)} = \nu_0$  and  $\mu_0^{(n)} = \mu_0^{(0)} = \mu_0$  (with  $\nu_0$  and  $\mu_0$  satisfying condition (13) from Assumption 4) and define  $(\nu_t^{(n)}, \mu_t^{(n)})_{t \in [0, T]}$  by

$$\begin{aligned} \log \nu_t^{(n)}(x) &= e^{-\frac{\sigma^2}{2}t} \log \nu_0(x) \\ &\quad - \int_0^t \frac{\sigma^2}{2} e^{-\frac{\sigma^2}{2}(t-s)} \left( \frac{2}{\sigma^2} \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, x) - \log \pi(x) - D_{\text{KL}}(\nu_s^{(n-1)} | \pi) \right) ds, \end{aligned} \quad (27)$$

$$\begin{aligned} \log \mu_t^{(n)}(y) &= e^{-\frac{\sigma^2}{2}t} \log \mu_0(y) \\ &\quad + \int_0^t \frac{\sigma^2}{2} e^{-\frac{\sigma^2}{2}(t-s)} \left( \frac{2}{\sigma^2} \frac{\delta F}{\delta \mu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, y) + \log \rho(y) + D_{\text{KL}}(\mu_s^{(n-1)} | \rho) \right) ds. \end{aligned} \quad (28)$$

We have the following result.

**Lemma A.1.** *The sequence of flows  $((\nu_t^{(n)}, \mu_t^{(n)})_{t \in [0, T]})_{n=0}^\infty$  given by (27) and (28) is well-defined and such that for all  $n \geq 1$  and all  $t \in [0, T]$  we have*

$$D_{\text{KL}}(\nu_t^{(n)} | \pi) \leq 2 \log R_\nu + \frac{4}{\sigma^2} C_\nu, \quad D_{\text{KL}}(\mu_t^{(n)} | \rho) \leq 2 \log R_\mu + \frac{4}{\sigma^2} C_\mu.$$

*Proof of Lemma A.1.* The proof follows from the same induction argument used to prove Liu et al. (2023, Lemma 3.1).  $\square$

For fixed  $T > 0$ , we consider the sequence of flows  $((\nu_t^{(n)}, \mu_t^{(n)})_{t \in [0, T]})_{n=0}^\infty$  in  $(\mathcal{P}(\mathcal{X})^{[0, T]} \times \mathcal{P}(\mathcal{Y})^{[0, T]}, \mathcal{TV}^{[0, T]})$ , where, for any  $(\nu_t, \mu_t)_{t \in [0, T]} \in \mathcal{P}(\mathcal{X})^{[0, T]} \times \mathcal{P}(\mathcal{Y})^{[0, T]}$ , the distance  $\mathcal{TV}^{[0, T]}$  is defined by

$$\mathcal{TV}^{[0, T]}((\nu_t, \mu_t)_{t \in [0, T]}, (\nu'_t, \mu'_t)_{t \in [0, T]}) := \int_0^T \text{TV}(\nu_t, \nu'_t) dt + \int_0^T \text{TV}(\mu_t, \mu'_t) dt.$$

Since  $(\mathcal{P}(\mathcal{X}), \text{TV})$  is complete, we can apply the argument from Šiška & Szpruch (2020, Lemma A.5) with  $p = 1$  to conclude that  $(\mathcal{P}(\mathcal{X})^{[0, T]}, \int_0^T \text{TV}(\nu_t, \nu'_t) dt)$  and  $(\mathcal{P}(\mathcal{Y})^{[0, T]}, \int_0^T \text{TV}(\mu_t, \mu'_t) dt)$  are complete. Therefore, one can deduce that  $(\mathcal{P}(\mathcal{X})^{[0, T]} \times \mathcal{P}(\mathcal{Y})^{[0, T]}, \mathcal{TV}^{[0, T]})$  is also complete. We consider the Picard iteration mapping  $\phi((\nu_t^{(n-1)}, \mu_t^{(n-1)})_{t \in [0, T]}) := (\nu_t^{(n)}, \mu_t^{(n)})_{t \in [0, T]}$  defined via (27) and (28) and show that  $\phi$  is contractive in  $(\mathcal{P}(\mathcal{X})^{[0, T]} \times \mathcal{P}(\mathcal{Y})^{[0, T]}, \mathcal{TV}^{[0, T]})$ . Then the Banach fixed point theorem will give us the existence and uniqueness of the solution to (17).

**Lemma A.2.** *The mapping  $\phi\left((\nu_t^{(n-1)}, \mu_t^{(n-1)})_{t \in [0, T]}\right) := (\nu_t^{(n)}, \mu_t^{(n)})_{t \in [0, T]}$  defined via (27) and (28) is contractive in  $(\mathcal{P}(\mathcal{X})^{[0, T]} \times \mathcal{P}(\mathcal{Y})^{[0, T]}, \mathcal{TV}^{[0, T]})$ .*

*Proof of Lemma A.2.* From (27), we have

$$\begin{aligned} \log \nu_t^{(n)}(x) - \log \nu_t^{(n-1)}(x) &= - \int_0^t \frac{\sigma^2}{2} e^{-\frac{\sigma^2}{2}(t-s)} \times \\ &\times \left[ \frac{2}{\sigma^2} \left( \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, x) - \frac{\delta F}{\delta \nu}(\nu_s^{(n-2)}, \mu_s^{(n-2)}, x) \right) - \text{D}_{\text{KL}}(\nu_s^{(n-1)}|\pi) + \text{D}_{\text{KL}}(\nu_s^{(n-2)}|\pi) \right] ds. \end{aligned}$$

Multiplying both sides by  $\nu_t^{(n)}(x)$  and integrating with respect to  $x$ , we obtain

$$\begin{aligned} \text{D}_{\text{KL}}(\nu_t^{(n)}|\nu_t^{(n-1)}) &= - \int_0^t \frac{\sigma^2}{2} e^{-\frac{\sigma^2}{2}(t-s)} \left[ \frac{2}{\sigma^2} \int_{\mathcal{X}} \left( \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, x) - \frac{\delta F}{\delta \nu}(\nu_s^{(n-2)}, \mu_s^{(n-2)}, x) \right) \nu_t^{(n)}(dx) \right. \\ &\quad \left. - \text{D}_{\text{KL}}(\nu_s^{(n-1)}|\pi) + \text{D}_{\text{KL}}(\nu_s^{(n-2)}|\pi) \right] ds. \quad (29) \end{aligned}$$

Moreover, note that

$$\begin{aligned} &\int_{\mathcal{X}} \left( \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, x) - \frac{\delta F}{\delta \nu}(\nu_s^{(n-2)}, \mu_s^{(n-2)}, x) \right) \nu_t^{(n)}(dx) \\ &= \int_{\mathcal{X}} \left( \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, x) - \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-2)}, x) + \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-2)}, x) - \frac{\delta F}{\delta \nu}(\nu_s^{(n-2)}, \mu_s^{(n-2)}, x) \right) \nu_t^{(n)}(dx) \\ &= \int_{\mathcal{X}} \int_{\mathcal{Y}} \int_0^1 \frac{\delta^2 F}{\delta \mu \delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-2)} + \lambda(\mu_s^{(n-1)} - \mu_s^{(n-2)}), x, w) d\lambda (\mu_s^{(n-1)} - \mu_s^{(n-2)})(dw) \nu_t^{(n)}(dx) \\ &\quad + \int_{\mathcal{X}} \int_{\mathcal{X}} \int_0^1 \frac{\delta^2 F}{\delta \nu^2}(\nu_s^{(n-2)} + \lambda(\nu_s^{(n-1)} - \nu_s^{(n-2)}), \mu_s^{(n-2)}, x, z) d\lambda (\nu_s^{(n-1)} - \nu_s^{(n-2)})(dz) \nu_t^{(n)}(dx). \end{aligned}$$

Similarly, again from (27) we have

$$\begin{aligned} \log \nu_t^{(n-1)}(x) - \log \nu_t^{(n)}(x) &= - \int_0^t \frac{\sigma^2}{2} e^{-\frac{\sigma^2}{2}(t-s)} \times \\ &\times \left[ \frac{2}{\sigma^2} \int_{\mathcal{X}} \left( \frac{\delta F}{\delta \nu}(\nu_s^{(n-2)}, \mu_s^{(n-2)}, x) - \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, x) \right) \nu_t^{(n)}(dx) \right. \\ &\quad \left. - \text{D}_{\text{KL}}(\nu_s^{(n-2)}|\pi) + \text{D}_{\text{KL}}(\nu_s^{(n-1)}|\pi) \right] ds. \end{aligned}$$

Multiplying both sides by  $\nu_t^{(n-1)}(x)$  and integrating with respect to  $x$ , we obtain

$$\begin{aligned} \text{D}_{\text{KL}}(\nu_t^{(n-1)}|\nu_t^{(n)}) &= - \int_0^t \frac{\sigma^2}{2} e^{-\frac{\sigma^2}{2}(t-s)} \left[ \frac{2}{\sigma^2} \int_{\mathcal{X}} \left( \frac{\delta F}{\delta \nu}(\nu_s^{(n-2)}, \mu_s^{(n-2)}, x) - \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, x) \right) \nu_t^{(n-1)}(dx) \right. \\ &\quad \left. - \text{D}_{\text{KL}}(\nu_s^{(n-2)}|\pi) + \text{D}_{\text{KL}}(\nu_s^{(n-1)}|\pi) \right] ds. \quad (30) \end{aligned}$$

Similarly as before, we note that

$$\begin{aligned}
& \int_{\mathcal{X}} \left( \frac{\delta F}{\delta \nu}(\nu_s^{(n-2)}, \mu_s^{(n-2)}, x) - \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, x) \right) \nu_t^{(n-1)}(dx) \\
&= - \int_{\mathcal{X}} \left( \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, x) - \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-2)}, x) + \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-2)}, x) - \frac{\delta F}{\delta \nu}(\nu_s^{(n-2)}, \mu_s^{(n-2)}, x) \right) \nu_t^{(n-1)}(dx) \\
&= - \int_{\mathcal{X}} \int_{\mathcal{Y}} \int_0^1 \frac{\delta^2 F}{\delta \mu \delta \nu} \left( \nu_s^{(n-1)}, \mu_s^{(n-2)} + \lambda \left( \mu_s^{(n-1)} - \mu_s^{(n-2)} \right), x, w \right) \times d\lambda \left( \mu_s^{(n-1)} - \mu_s^{(n-2)} \right) (dw) \nu_t^{(n-1)}(dx) \\
&\quad - \int_{\mathcal{X}} \int_{\mathcal{X}} \int_0^1 \frac{\delta^2 F}{\delta \nu^2} \left( \nu_s^{(n-2)} + \lambda \left( \nu_s^{(n-1)} - \nu_s^{(n-2)} \right), \mu_s^{(n-2)}, x, z \right) \times d\lambda \left( \nu_s^{(n-1)} - \nu_s^{(n-2)} \right) (dz) \nu_t^{(n-1)}(dx).
\end{aligned}$$

Combining (29) and (30), we obtain

$$\begin{aligned}
& \text{D}_{\text{KL}}(\nu_t^{(n)} | \nu_t^{(n-1)}) + \text{D}_{\text{KL}}(\nu_t^{(n-1)} | \nu_t^{(n)}) = - \int_0^t e^{-\frac{\sigma^2}{2}(t-s)} \times \left[ \right. \\
& \int_{\mathcal{X}} \int_{\mathcal{Y}} \int_0^1 \frac{\delta^2 F}{\delta \mu \delta \nu} \left( \nu_s^{(n-1)}, \mu_s^{(n-2)} + \lambda \left( \mu_s^{(n-1)} - \mu_s^{(n-2)} \right), x, w \right) d\lambda \left( \mu_s^{(n-1)} - \mu_s^{(n-2)} \right) (dw) \left( \nu_t^{(n)} - \nu_t^{(n-1)} \right) (dx) \\
& \left. + \int_{\mathcal{X}} \int_{\mathcal{X}} \int_0^1 \frac{\delta^2 F}{\delta \nu^2} \left( \nu_s^{(n-2)} + \lambda \left( \nu_s^{(n-1)} - \nu_s^{(n-2)} \right), \mu_s^{(n-2)}, x, z \right) d\lambda \left( \nu_s^{(n-1)} - \nu_s^{(n-2)} \right) (dz) \left( \nu_t^{(n)} - \nu_t^{(n-1)} \right) (dx) \right] ds.
\end{aligned}$$

Hence, due to Assumption 3, we get

$$\begin{aligned}
& \text{D}_{\text{KL}}(\nu_t^{(n)} | \nu_t^{(n-1)}) + \text{D}_{\text{KL}}(\nu_t^{(n-1)} | \nu_t^{(n)}) \\
& \leq \text{TV}(\nu_t^{(n)}, \nu_t^{(n-1)}) \int_0^t e^{-\frac{\sigma^2}{2}(t-s)} \left( C_{\mu, \nu} \text{TV}(\mu_s^{(n-1)}, \mu_s^{(n-2)}) + C_{\nu, \nu} \text{TV}(\nu_s^{(n-1)}, \nu_s^{(n-2)}) \right) ds \\
& \leq \max\{C_{\mu, \nu}, C_{\nu, \nu}\} \text{TV}(\nu_t^{(n)}, \nu_t^{(n-1)}) \int_0^t e^{-\frac{\sigma^2}{2}(t-s)} \left( \text{TV}(\mu_s^{(n-1)}, \mu_s^{(n-2)}) + \text{TV}(\nu_s^{(n-1)}, \nu_s^{(n-2)}) \right) ds.
\end{aligned}$$

By the Pinsker-Csizsar inequality,  $\text{TV}^2(\nu_t^{(n)}, \nu_t^{(n-1)}) \leq \frac{1}{2} \text{D}_{\text{KL}}(\nu_t^{(n)} | \nu_t^{(n-1)})$ , and hence

$$\begin{aligned}
4 \text{TV}^2(\nu_t^{(n)}, \nu_t^{(n-1)}) & \leq \max\{C_{\mu, \nu}, C_{\nu, \nu}\} \text{TV}(\nu_t^{(n)}, \nu_t^{(n-1)}) \int_0^t e^{-\frac{\sigma^2}{2}(t-s)} \left( \text{TV}(\mu_s^{(n-1)}, \mu_s^{(n-2)}) \right. \\
& \quad \left. + \text{TV}(\nu_s^{(n-1)}, \nu_s^{(n-2)}) \right) ds,
\end{aligned}$$

which gives

$$\text{TV}(\nu_t^{(n)}, \nu_t^{(n-1)}) \leq \frac{1}{4} \max\{C_{\mu, \nu}, C_{\nu, \nu}\} \int_0^t e^{-\frac{\sigma^2}{2}(t-s)} \left( \text{TV}(\mu_s^{(n-1)}, \mu_s^{(n-2)}) + \text{TV}(\nu_s^{(n-1)}, \nu_s^{(n-2)}) \right) ds.$$

An almost identical argument leads to

$$\text{TV}(\mu_t^{(n)}, \mu_t^{(n-1)}) \leq \frac{1}{4} \max\{C_{\nu, \mu}, C_{\mu, \mu}\} \int_0^t e^{-\frac{\sigma^2}{2}(t-s)} \left( \text{TV}(\mu_s^{(n-1)}, \mu_s^{(n-2)}) + \text{TV}(\nu_s^{(n-1)}, \nu_s^{(n-2)}) \right) ds.$$

If we set  $C_{\max} := \max\{C_{\mu,\nu}, C_{\nu,\nu}\} + \max\{C_{\nu,\mu}, C_{\mu,\mu}\}$ , and add the previous two inequalities, we obtain

$$\begin{aligned} \text{TV}(\nu_t^{(n)}, \nu_t^{(n-1)}) + \text{TV}(\mu_t^{(n)}, \mu_t^{(n-1)}) &\leq \frac{C_{\max}}{4} \int_0^t e^{-\frac{\sigma^2}{2}(t-s)} \left( \text{TV}(\mu_s^{(n-1)}, \mu_s^{(n-2)}) + \text{TV}(\nu_s^{(n-1)}, \nu_s^{(n-2)}) \right) ds. \\ &\leq \left( \frac{C_{\max}}{4} \right)^{n-1} e^{-\frac{\sigma^2}{2}t} \int_0^t \int_0^{t_1} \dots \int_0^{t_{n-2}} e^{\frac{\sigma^2}{2}t_{n-1}} \left( \text{TV}(\nu_{t_{n-1}}^{(1)}, \nu_{t_{n-1}}^{(0)}) + \text{TV}(\mu_{t_{n-1}}^{(1)}, \mu_{t_{n-1}}^{(0)}) \right) dt_{n-1} \dots dt_2 dt_1 \\ &\leq \left( \frac{C_{\max}}{4} \right)^{n-1} e^{-\frac{\sigma^2}{2}t} \frac{t^{n-2}}{(n-2)!} \int_0^t e^{\frac{\sigma^2}{2}t_{n-1}} \left( \text{TV}(\nu_{t_{n-1}}^{(1)}, \nu_{t_{n-1}}^{(0)}) + \text{TV}(\mu_{t_{n-1}}^{(1)}, \mu_{t_{n-1}}^{(0)}) \right) dt_{n-1} \\ &\leq \left( \frac{C_{\max}}{4} \right)^{n-1} \frac{t^{n-2}}{(n-2)!} \int_0^t \left( \text{TV}(\nu_{t_{n-1}}^{(1)}, \nu_{t_{n-1}}^{(0)}) + \text{TV}(\mu_{t_{n-1}}^{(1)}, \mu_{t_{n-1}}^{(0)}) \right) dt_{n-1}, \end{aligned}$$

where in the third inequality we bounded  $\int_0^{t_{n-2}} dt_{n-1} \leq \int_0^t dt_{n-1}$  and in the fourth inequality we bounded  $e^{\frac{\sigma^2}{2}t_{n-1}} \leq e^{\frac{\sigma^2}{2}t}$ . Hence, we obtain

$$\begin{aligned} \int_0^T \text{TV}(\nu_t^{(n)}, \nu_t^{(n-1)}) dt + \int_0^T \text{TV}(\mu_t^{(n)}, \mu_t^{(n-1)}) dt \\ \leq \left( \frac{C_{\max}}{4} \right)^{n-1} \frac{T^{n-1}}{(n-1)!} \left( \int_0^T \text{TV}(\nu_{t_{n-1}}^{(1)}, \nu_{t_{n-1}}^{(0)}) dt_{n-1} + \int_0^T \text{TV}(\mu_{t_{n-1}}^{(1)}, \mu_{t_{n-1}}^{(0)}) dt_{n-1} \right). \end{aligned}$$

For sufficiently large  $n$ , the constant on the right hand side becomes less than 1 and the proof is complete.  $\square$

By Lemma A.2, for any  $T > 0$  we obtain the existence of a pair of flows  $(\nu_t, \mu_t)_{t \in [0, T]}$  satisfying (25). Moreover, for Lebesgue-almost all  $t \in [0, T]$  we have

$$\text{TV}(\nu_t^{(n)}, \nu_t) \rightarrow 0, \quad \text{TV}(\mu_t^{(n)}, \mu_t) \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

which implies

$$\nu_t^{(n)} \rightarrow \nu_t, \quad \mu_t^{(n)} \rightarrow \mu_t \quad \text{weakly as } n \rightarrow \infty.$$

Hence, using the lower semi-continuity of the entropy, we obtain

$$\text{D}_{\text{KL}}(\nu_t | \pi) \leq \liminf_{n \rightarrow \infty} \text{D}_{\text{KL}}(\nu_t^{(n)} | \pi) \leq 2 \log R_\nu + \frac{4}{\sigma^2} C_\nu, \quad (31)$$

$$\text{D}_{\text{KL}}(\mu_t | \rho) \leq \liminf_{n \rightarrow \infty} \text{D}_{\text{KL}}(\mu_t^{(n)} | \rho) \leq 2 \log R_\mu + \frac{4}{\sigma^2} C_\mu, \quad (32)$$

where both second inequalities follow from Lemma A.1. In order to ensure that the solution  $(\nu_t, \mu_t)_{t \in [0, T]}$  can be extended to all  $t \geq 0$ , we first need to prove the bound on the ratios  $\nu_t/\pi$  and  $\mu_t/\rho$  in (19).

*Step 2: Ratio condition (19).* Using (27) and (28), we see that for any  $t \in [0, T]$  we have

$$\begin{aligned} \log \frac{\nu_t^{(n)}(x)}{\pi(x)} &= e^{-\frac{\sigma^2}{2}t} \log \frac{\nu_0(x)}{\pi(x)} - \int_0^t \frac{\sigma^2}{2} e^{-\frac{\sigma^2}{2}(t-s)} \left( \frac{2}{\sigma^2} \frac{\delta F}{\delta \nu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, x) - \text{D}_{\text{KL}}(\nu_s^{(n-1)} | \pi) \right) ds, \\ \log \frac{\mu_t^{(n)}(y)}{\rho(y)} &= e^{-\frac{\sigma^2}{2}t} \log \frac{\mu_0(y)}{\rho(y)} + \int_0^t \frac{\sigma^2}{2} e^{-\frac{\sigma^2}{2}(t-s)} \left( \frac{2}{\sigma^2} \frac{\delta F}{\delta \mu}(\nu_s^{(n-1)}, \mu_s^{(n-1)}, y) + \text{D}_{\text{KL}}(\mu_s^{(n-1)} | \rho) \right) ds. \end{aligned}$$

Using Assumption 2, (13), (31) and (32) we obtain

$$\begin{aligned} \log \frac{\nu_t(x)}{\pi(x)} &\leq \log R_\nu + C_\nu + \frac{\sigma^2}{2} \left( 2 \log R_\nu + \frac{4}{\sigma^2} C_\nu \right), \\ \log \frac{\mu_t(y)}{\rho(y)} &\leq \log R_\mu + C_\mu + \frac{\sigma^2}{2} \left( 2 \log R_\mu + \frac{4}{\sigma^2} C_\mu \right). \end{aligned}$$

Hence we can choose  $R_{1,\nu} := 1 + \exp\left(\log R_\nu + C_\nu + \frac{\sigma^2}{2} (2 \log R_\nu + \frac{4}{\sigma^2} C_\nu)\right)$  and  $R_{1,\mu} := 1 + \exp\left(\log R_\mu + C_\mu + \frac{\sigma^2}{2} (2 \log R_\mu + \frac{4}{\sigma^2} C_\mu)\right)$ . Note  $R_{1,\nu}, R_{1,\mu} > 1$  are conveniently chosen so that  $\log R_{1,\nu}, \log R_{1,\mu} > 0$  in our subsequent calculations. Obtaining a lower bound on  $\frac{\nu_t(x)}{\pi(x)}$  and  $\frac{\mu_t(y)}{\rho(y)}$  follows similarly, by using (12) instead of (13).

*Step 3: Existence of the gradient flow on  $[0, \infty)$ .* In order to complete our proof, note that the unique solution  $(\nu_t, \mu_t)_{t \in [0, T]}$  to (25) can also be expressed as

$$\begin{aligned}\nu_t(x) &= \nu_0(x) \exp\left(-\int_0^t \left(\frac{\delta F}{\delta \nu}(\nu_s, \mu_s, x) + \frac{\sigma^2}{2} \log\left(\frac{\nu_s(x)}{\pi(x)}\right) - \frac{\sigma^2}{2} D_{\text{KL}}(\nu_s|\pi)\right) ds\right), \\ \mu_t(y) &= \mu_0(y) \exp\left(\int_0^t \left(\frac{\delta F}{\delta \mu}(\nu_s, \mu_s, y) - \frac{\sigma^2}{2} \log\left(\frac{\mu_s(y)}{\rho(y)}\right) + \frac{\sigma^2}{2} D_{\text{KL}}(\mu_s|\rho)\right) ds\right).\end{aligned}$$

From 2, (31), (32) and (19), we obtain for any  $t \in [0, T]$

$$\begin{aligned}\left|\frac{\delta F}{\delta \nu}(\nu_t, \mu_t, x) + \frac{\sigma^2}{2} \log\left(\frac{\nu_t(x)}{\pi(x)}\right) - \frac{\sigma^2}{2} D_{\text{KL}}(\nu_t|\pi)\right| \\ \leq 3C_\nu + \frac{\sigma^2}{2} (\max\{|\log r_{1,\nu}|, \log R_{1,\nu}\} + 2 \log R_\nu) =: C_{V,\nu},\end{aligned}\quad (33)$$

$$\begin{aligned}\left|\frac{\delta F}{\delta \mu}(\nu_t, \mu_t, y) - \frac{\sigma^2}{2} \log\left(\frac{\mu_t(y)}{\rho(y)}\right) + \frac{\sigma^2}{2} D_{\text{KL}}(\mu_t|\rho)\right| \\ \leq 3C_\mu + \frac{\sigma^2}{2} (\max\{|\log r_{1,\mu}|, \log R_{1,\mu}\} + 2 \log R_\mu) =: C_{V,\mu}.\end{aligned}$$

This gives  $\|\nu_t\|_{TV} \leq \|\nu_0\|_{TV} e^{C_{V,\nu} t}$  and  $\|\mu_t\|_{TV} \leq \|\mu_0\|_{TV} e^{C_{V,\mu} t}$ , and shows that  $\nu_t$  and  $\mu_t$  do not explode in any finite time, hence we obtain a global solution  $(\nu_t, \mu_t)_{t \in [0, \infty)}$ . In particular, the bounds in (31), (32), (19) and (20) hold for all  $t \geq 0$ .  $\square$

## B Notation and definitions

In this section we recall some important definitions. Following Carmona & Delarue (2018, Definition 5.43), we start with the notion of differentiability on the space of probability measure that we utilize throughout the paper.

**Definition B.1.** Fix  $p \geq 0$ . For any  $\mathcal{M} \subseteq \mathbb{R}^d$ , let  $\mathcal{P}_p(\mathcal{M})$  be the space of probability measures on  $\mathcal{M}$  with finite  $p$ -th moments. A function  $F : \mathcal{P}_p(\mathcal{M}) \rightarrow \mathbb{R}$  admits first-order flat derivative on  $\mathcal{P}_p(\mathcal{M})$ , if there exists a function  $\frac{\delta F}{\delta \nu} : \mathcal{P}_p(\mathcal{M}) \times \mathcal{M} \rightarrow \mathbb{R}$ , such that

1. the map  $\mathcal{P}_p(\mathcal{M}) \times \mathcal{M} \ni (m, x) \mapsto \frac{\delta F}{\delta m}(m, x)$  is jointly continuous with respect to the product topology, where  $\mathcal{P}_p(\mathcal{M})$  is endowed with the weak topology,
2. For any  $m \in \mathcal{P}_p(\mathcal{M})$ , there exists  $C > 0$  such that, for all  $x \in \mathcal{M}$ , we have

$$\left|\frac{\delta F}{\delta m}(m, x)\right| \leq C (1 + |x|^p),$$

3. For all  $m, m' \in \mathcal{P}_p(\mathcal{M})$ , it holds that

$$F(m') - F(m) = \int_0^1 \int_{\mathcal{M}} \frac{\delta F}{\delta m}(m + \varepsilon(m' - m), x) (m' - m) (dx) d\varepsilon. \quad (34)$$



The functional  $\frac{\delta F}{\delta m}$  is then called the flat derivative of  $F$  on  $\mathcal{P}_p(\mathcal{M})$ . We note that  $\frac{\delta F}{\delta m}$  exists up to an additive constant, and thus we make the normalizing convention  $\int_{\mathcal{M}} \frac{\delta F}{\delta m}(m, x) m(dx) = 0$ .

If, for fixed  $x \in \mathcal{M}$ , the map  $m \mapsto \frac{\delta F}{\delta m}(m, x)$  satisfies Definition B.1, we say that  $F$  admits a second-order flat derivative denoted by  $\frac{\delta^2 F}{\delta m^2}$ . Consequently, by Definition B.1, there exists a functional  $\frac{\delta^2 F}{\delta m^2} : \mathcal{P}_p(\mathcal{M}) \times \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$  such that

$$\frac{\delta F}{\delta m}(m', x) - \frac{\delta F}{\delta m}(m, x) = \int_0^1 \int_{\mathcal{M}} \frac{\delta^2 F}{\delta m^2}(\nu + \varepsilon(m' - m), x, x') (m' - m) (dx') d\varepsilon. \quad (35)$$

**Definition B.2** (TV distance between probability measures; (Tsybakov, 2008), Definition 2.4). *Let  $(\mathcal{M}, \mathcal{A})$  be a measurable space and let  $P$  and  $Q$  be probability measures on  $(\mathcal{M}, \mathcal{A})$ . Assume that  $\mu$  is a  $\sigma$ -finite measure on  $(\mathcal{M}, \mathcal{A})$  such that  $P$  and  $Q$  are absolutely continuous with respect to  $\mu$  and let  $p$  and  $q$  denote their probability density functions, respectively. The total variation distance between  $P$  and  $Q$  is defined as:*

$$\text{TV}(P, Q) := \sup_{A \in \mathcal{A}} |P(A) - Q(A)| = \sup_{A \in \mathcal{A}} \left| \int_A (p - q) d\mu \right|.$$