

# Markov Chain from Human Feedback

author names withheld

Under Review for NExT-Game 2026

## Abstract

We propose Markov Chain from Human Feedback (MCHF), an alternative approach for aligning generative models from pairwise human preferences. Unlike RLHF, which reduces comparisons to a scalar reward, and NLHF, which preserves pairwise utilities through a KL-regularized minimax objective, MCHF uses pairwise preferences directly to define a transition mechanism over model outputs. Given a pairwise utility  $U(x, y)$  and a reference model  $\mu_{\text{ref}}$ , we define a Markov kernel  $P(x, dy) \propto \exp(U(x, y))\mu_{\text{ref}}(dy)$ , and take its stationary distribution as the aligned model. We show that MCHF converges geometrically fast to the stationary distribution, with a convergence rate governed by the seminorm  $\|U\|_{\oplus} = \inf_{g, f} \|U - g \oplus f\|_{\infty}$ , which quantifies the non-transitive structure of the pairwise utility. We further show that a mirror-descent algorithm for NLHF, with an appropriate step size, satisfies an analogous structure-adaptive convergence guarantee. Finally, through a perturbation analysis, we prove that when  $\|U\|_{\oplus}$  is small, MCHF and NLHF agree up to first order around a RLHF solution, which yields a unified view of reward-based, game-theoretic, and our Markovian alignment.

## 1. Introduction

Alignment from human feedback often begins with pairwise comparisons. RLHF reduces such data to a scalar reward  $R$ , typically through a BTL model [3, 9]:

$$\mathcal{P}(x \prec y) = \sigma_{\text{sigmoid}}(R(y) - R(x)) \quad \text{where} \quad \sigma_{\text{sigmoid}}(t) \equiv (1 + e^{-t})^{-1} \quad (1)$$

and then performs KL-regularized policy optimization [2, 4, 13]:

$$\mu_{\text{RL}} = \operatorname{argmax}_{\mu} \int \mu(dy)R(y) - D_{\text{KL}}(\mu | \mu_{\text{ref}}) \quad (2)$$

More recently, several variants to the standard RLHF have been proposed [1, 6, 8, 10, 16]. However, this reduction to scalar rewards cannot directly represent cyclic or pair-specific preferences. Nash Learning from Human Feedback (NLHF) [11] avoids this reduction by computing the minimax optimization problem:

$$(\mu_{\text{NL}}, \nu_{\text{NL}}) \in \operatorname{argmax}_{\mu} \operatorname{argmin}_{\nu} \iint \nu(dx)\mu(dy)U(x, y) - D_{\text{KL}}(\mu | \mu_{\text{ref}}) + D_{\text{KL}}(\nu | \mu_{\text{ref}}) \quad (3)$$

where  $U(x, y)$  is a pairwise utility, typically chosen as an increasing function of human’s preference model  $\mathcal{P}(x \prec y)$ . More broadly, game-theoretic perspectives on alignment have led to formulations based on other concepts, including Stackelberg games [5, 15] and multiplayer games [20].

The original paper [11] argues that NLHF is attractive because it can yield a more diverse distribution than reward-based alignment methods. More refined analyses have been developed through the lens of social choice theory [7, 18]. Nevertheless, it is not always clear that such a pessimistic objective is the right one for fine-tuning generative models. In addition, existing algorithms (cf. [11, 17, 19]) require repeatedly solving a KL-regularized optimization with an iteration dependent reward  $R_t$  and a local regularization term  $D_{\text{KL}}(\cdot|\mu_t)$ , which is computationally difficult for large generative models, as discussed in [11] (see Section 2.3 for details).

At a structural level, however, pairwise comparisons naturally define a directed graph over the space of possible outputs, where edges encode relative preference information. This perspective is classical in ranking and social choice, and has also motivated spectral methods for rank aggregation, in which the stationary distribution of a Markov chain constructed from pairwise comparisons represents the aggregate ranking [12, 14].

Motivated by this perspective, we propose *Markov Chain from Human Feedback* (MCHF). We construct the Markov kernel  $P$  from the pairwise utility as  $P(x, dy) \propto \exp(U(x, y))\mu_{\text{ref}}(dy)$ . We then iteratively align  $\mu_{\text{ref}}$  by  $\mu_{\text{ref}} \rightarrow \mu_{\text{ref}}P \rightarrow \mu_{\text{ref}}P^2 \dots$  and take the stationary distribution as  $\mu_{\text{MC}}$ . Our results can be summarized as follows:

- In Section 2.2, we prove the geometric convergence of MCHF to the stationary distribution  $\mu_{\text{MC}}$ . In particular, the convergence rate is governed by the seminorm  $\|U\|_{\oplus} = \inf_{g, f \in L^\infty(\mu_{\text{ref}})} \|U - g \oplus f\|_\infty$  which measures the intransitive part of  $U$ . We show that the MCHF implicitly adapts to the additive structure of  $U$ .
- In Section 2.3, we provide a refined convergence rate analysis for NLHF, showing that the NLHF also adapts to this additive structure of  $U$  if step-size is appropriately chosen.
- In Section 2.4, we compare MCHF and NLHF through (1) computational perspective, (2) coupling perspective, and (3) alignment dynamics. Especially, as for (3), we develop a perturbation analysis for MCHF and NLHF, showing that when  $\|U\|_{\oplus}$  is small, MCHF and NLHF agrees up to the first order around the RLHF with reward  $f(y) = \int \mu_{\text{ref}}(dx)U(x, y)$ .

## 2. Main result

### 2.1. Definition of seminorm $\|\cdot\|_{\oplus}$ : Measuring non-additive structure of Utility

Let  $(\mathcal{X}, \mathcal{F}, \mu_{\text{ref}})$  be a probability space. We write  $\|X\|_p = \|X\|_{L^p(\mu_{\text{ref}})}$  for the  $L^p$  norm. Let  $U(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  be a utility function that encodes preference information. Throughout this paper, we assume  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ .

We now define a seminorm on  $L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , which serves as a fundamental complexity measure for the alignment dynamics of MCHF and NLHF throughout this paper.

**Definition 1** Define the seminorm  $\|\cdot\|_{\oplus}$  on  $L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  by

$$\|U\|_{\oplus} = \inf_{g, f \in L^\infty(\mu_{\text{ref}})} \|U - g \oplus f\|_\infty \quad \text{where} \quad g \oplus f(x, y) = g(x) + f(y).$$

This seminorm  $\|U\|_{\oplus}$  captures *additive defect* of the utility; the larger  $\|U\|_{\oplus}$  is, the more  $U$  has a non-additive structure. Furthermore, when  $U$  is antisymmetric,  $\|U\|_{\oplus}$  quantifies the *non-transitivity* of  $U$ :  $\|U\|_{\oplus} \asymp \text{esssup}_{x, y, z} |U(x, y) + U(y, z) + U(z, x)|$  (see Proposition 13).

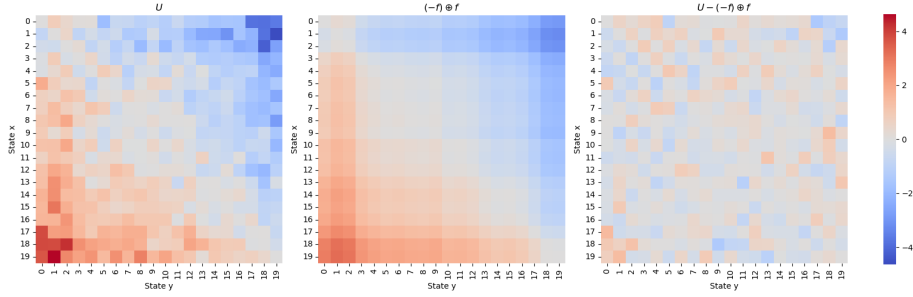


Figure 1: Comparison of antisymmetric utility  $U$  generated by (8),  $(-\hat{f}) \oplus \hat{f}$  with  $\hat{f}(y) = \int \mu_{\text{ref}}(dx)U(x, y)$ , and residual  $U - (-\hat{f}) \oplus \hat{f}$ .

Note that  $\|U\|_{\oplus} \leq \|U\|_{\infty}$  by the definition. In practice, the utility is often taken to be a possibly nonlinear transformation of the preference model  $\mathcal{P}(x \prec y)$ . Thus, the value of  $\|U\|_{\oplus}$  depends both on the transformation and on the structure of  $\mathcal{P}(x \prec y)$ . Importantly, there are natural situations in which  $\|U\|_{\oplus} \ll \|U\|_{\infty}$ .

- Suppose  $U(x, y) = \sigma_{\text{sigmoid}}^{-1} \circ \mathcal{P}(x \prec y)$  where  $\sigma_{\text{sigmoid}}^{-1}(t) = \log t / (1 - t)$  and  $\mathcal{P}(x \prec y)$  takes the form  $\mathcal{P}(x \prec y) = \Phi(R(y) - R(x) + E(x, y))$ , where  $\Phi$  is a nonlinear link function and  $E(x, y)$  is a pairwise interaction term. In this case,  $U(x, y)$  is close to  $R(y) - R(x)$  (and hence  $\|U\|_{\oplus}$  is small), whenever  $\Phi$  is close to  $\sigma_{\text{sigmoid}}$  and the pairwise interaction term  $E$  is small. In contrast,  $\|U\|_{\infty}$  can be large, depending on  $R$  and  $E$  (see Figure 1). In particular, if  $\mathcal{P}(x \prec y)$  follows the BTL model (i.e.,  $\Phi = \sigma_{\text{sigmoid}}$  and  $E = 0$ ), we have  $U(x, y) = R(y) - R(x)$  and hence  $\|U\|_{\oplus} = 0$ .
- Suppose  $U(x, y) = \mathcal{P}(x \prec y) - 1/2$  and  $\mathcal{P}(x \prec y)$  follows the BTL model. If the reward  $R$  is small, using  $\sigma_{\text{sigmoid}}(x) = \frac{1}{2} + \frac{x}{4} + O(x^3)$  as  $x \rightarrow 0$ , one can show  $\|U\|_{\oplus} = O(\|R\|_{\infty}^3)$ . Meanwhile, we can always construct  $R$  such that  $\|U\|_{\infty} \asymp \|R\|_{\infty}$ ; for example, consider the discrete case  $|\mathcal{X}| = 2$  with  $R = (\gamma, -\gamma)$  and let  $\gamma \rightarrow 0$ . In this case, we have  $\|U\|_{\oplus} / \|U\|_{\infty} = O(\|R\|_{\infty}^2) \rightarrow 0$ .

The exact computation of the seminorm  $\|U\|_{\oplus}$  can be challenging, since it is essentially an  $L^{\infty}$ -Linear programming. However, we can estimate  $\|U\|_{\oplus}$  up to a multiplicative order:  $\|U\|_{\oplus} \asymp \|U - \hat{g} \oplus \hat{f}\|_{\infty}$ , where  $(\hat{g}, \hat{f})$  is a  $L^2$ -solution  $(\hat{g}, \hat{f}) \in \inf_{g, f} \|U - g \oplus f\|_2$  (see Proposition 12).

## 2.2. Definition of MCHF, iterative convergence, and equilibrium

Based on the utility function  $U \in L^{\infty}(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , we define the following Markov kernel  $P$ :

$$P(x, dy) = \frac{1}{Z(x)} \exp(U(x, y)) \mu_{\text{ref}}(dy), \quad Z(x) \equiv \int \exp(U(x, y)) \mu_{\text{ref}}(dy) \quad (4)$$

where  $Z(x)$  is a normalizing constant. Here, the conditional distribution  $Y|X \sim P(X, \cdot)$  will move mass towards preferred directions (where  $U(x, \cdot)$  takes large values) while preserving  $\mu_{\text{ref}}$ .

Given the Markov kernel  $P$  defined by (4), we iteratively update  $\mu_{\text{ref}}$  by the Markov chain:

$$\mu_{\text{ref}} \mapsto \mu_{\text{ref}} P \rightarrow \mu_{\text{ref}} P^2 \cdots,$$

where  $\mu P(dy) = \int \mu(dx)P(x, dy)$ . Now we claim that the map  $\mu \mapsto \mu P$  is a contraction mapping on the complete metric space  $(\mathcal{P}_{\mu_{\text{ref}}}, d_{\text{TV}})$ , where  $\mathcal{P}_{\mu_{\text{ref}}}$  is the set of probability measure that is absolutely continuous with respect to  $\mu_{\text{ref}}$ .

**Theorem 2** *For any  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , let  $P$  be the Markov Kernel defined by (4). Then,*

$$\forall \mu, \nu \in \mathcal{P}_{\mu_{\text{ref}}}, \quad d_{\text{TV}}(\mu P, \nu P) \leq c \cdot d_{\text{TV}}(\mu, \nu),$$

where the contraction rate  $c \in [0, 1)$  is given by

$$c = c(\|U\|_{\oplus}) = \min(1 - e^{-2\|U\|_{\oplus}}, \|U\|_{\oplus}) \in [0, 1). \quad (5)$$

Note that  $c(\|U\|_{\oplus})$  is an increasing function of  $\|U\|_{\oplus}$  and strictly less than 1. Thus, by the Banach fixed-point theorem, we obtain the following result.

**Corollary 3** *There exists a unique stationary distribution  $\mu_{\text{MC}} = \mu_{\text{MC}} P$ , and we have*

$$d_{\text{TV}}(\mu_{\text{ref}} P^t, \mu_{\text{MC}}) \leq d_{\text{TV}}(\mu_{\text{ref}}, \mu_{\text{MC}}) \cdot c(\|U\|_{\oplus})^t.$$

We emphasize that the contraction rate  $c(\|U\|_{\oplus})$  is bounded in terms of the semi-norm  $\|U\|_{\oplus}$ , rather than the naive  $L^\infty$ -norm  $\|U\|_\infty$ . Thus, the MCHF implicitly adapts to the additive structure of  $U$ .

### 2.3. Refined convergence rate analysis for NLHF

In this section, we propose a mirror descent algorithm for NLHF, which also adapts to the additive structure of  $U$  if the step size is chosen carefully.

**Theorem 4** *For any step size  $\eta > 0$ , define the iteration  $(\nu_t, \mu_t)$  by*

$$\begin{aligned} \nu_t &= \underset{\nu}{\operatorname{argmin}} \iint \nu(dx) \mu_t(dy) U(x, y) + D_{\text{KL}}(\nu \mid \mu_{\text{ref}}), \\ \mu_{t+1} &= \underset{\mu}{\operatorname{argmax}} \iint \nu_t(dx) \mu(dy) U(x, y) - D_{\text{KL}}(\mu \mid \mu_{\text{ref}}) - \eta^{-1} D_{\text{KL}}(\mu \mid \mu_t), \end{aligned} \quad (6)$$

initialized at  $\mu_0 = \mu_{\text{ref}}$ . Then, for any  $0 < \eta \leq \|U\|_{\oplus}^{-2}$ , it holds that

$$D_{\text{KL}}(\mu_{\text{NL}} \mid \mu_t) \leq (1 + \eta)^{-t} \cdot D_{\text{KL}}(\mu_{\text{NL}} \mid \mu_{\text{ref}}).$$

Taking  $\eta = \|U\|_{\oplus}^{-2}$  gives  $D_{\text{KL}}(\mu_{\text{NL}} \mid \mu_t) \leq (1 + \|U\|_{\oplus}^{-2})^{-t} D_{\text{KL}}(\mu_{\text{NL}} \mid \mu_{\text{ref}})$ . Related last-iterate convergence guarantees were obtained in, e.g., [19], for the special utility  $U(x, y) = \lambda^{-1}(\mathcal{P}(x \prec y) - 1/2)$ . However, their update also regularizes  $\nu_t$  by  $\eta^{-1} D_{\text{KL}}(\nu \mid \nu_{t-1})$ , and their rate depends on  $\|U\|_\infty$  rather than  $\|U\|_{\oplus}$ . Since  $\|U\|_{\oplus} \leq \|U\|_\infty$ , and  $\|U\|_{\oplus} \ll \|U\|_\infty$  can occur when  $U$  has a strong additive component, Theorem 4 yields a sharper, structure-adaptive rate.

Note that our refined rate still requires choosing  $\eta = \|U\|_{\oplus}^{-2}$ , whose exact value may be hard to compute in practice, even though it can be estimated up to multiplicative constants by Proposition 12-13. A further computational issue is the local regularization  $D_{\text{KL}}(\mu \mid \mu_t)$ . Indeed, the update (6) can be written as

$$\nu_t(dx) \propto \exp\left(-\int \mu_t(dy) U(x, y)\right) \mu_{\text{ref}}(dx), \quad \mu_{t+1}(dy) \propto \exp\left(\frac{\eta}{\eta+1} \int \nu_t(dx) U(x, y)\right) \tilde{\mu}_t^\eta(dy),$$

where  $\tilde{\mu}_t^\eta(dy) = (d\mu_t/d\mu_{\text{ref}})^{1/(1+\eta)} \mu_{\text{ref}}(dy)$ . Thus, while  $\nu_t$  is obtained by an exponential tilt of  $\mu_{\text{ref}}$ , updating  $\mu_{t+1}$  requires sampling from the geometric mixture  $\tilde{\mu}_t^\eta$  of  $\mu_t$  and  $\mu_{\text{ref}}$ , which is difficult in practice [11, Section F.1].

## 2.4. Comparison between MCHF and NLHF

**Computational perspective** As discussed above, for general  $U$ , last-iterate convergence of NLHF typically requires the local regularizer  $\eta^{-1}D_{\text{KL}}(\mu \mid \mu_{t-1})$ . Consequently, each iteration must update the reward and sample from a geometric mixture of  $\mu_{t-1}$  and  $\mu_{\text{ref}}$ , which can be computationally challenging.

By contrast, MCHF requires implementing the conditional sampler  $P(x, dy)$ , but we just need to implement it only once. Once such a sampling block is available, the alignment procedure simply consists of repeatedly applying the same Markov kernel. In this sense, MCHF can be viewed as attaching a preference-guided conditional sampling module on top of the existing architecture that generates samples from  $\mu_{\text{ref}}$ .

**Coupling perspective** The key difference between NLHF and MCHF lies in their update dynamics. For simplicity, consider an antisymmetric utility  $U$  and the iterations

$$\mu_{\text{MC}}^t = \mu_{\text{MC}}^{t-1}P, \quad \mu_{\text{NL}}^t = \operatorname{argmax}_{\mu} \int \mu_{\text{NL}}^{t-1}(dx)\mu(dy)U(x, y) - D_{\text{KL}}(\mu \mid \mu_{\text{ref}}), \quad (7)$$

initialized at  $\mu_{\text{MC}}^0 = \mu_{\text{NL}}^0 = \mu_{\text{ref}}$ . Note that when  $\|U\|_{\oplus} < 1$  and  $U$  is antisymmetric, the NLHF iteration  $\mu_{\text{NL}}^t$  also converges to  $\mu_{\text{NL}}$  (see Theorem 18).

To compare the updates, define the induced couplings

$$\pi_{\text{MC}}(dx, dy) = \mu_{t-1}(dx)P(x, dy), \quad \pi_{\text{NL}}(dx, dy) = \mu_{t-1}(dx)\mu_{\text{NL}}^t(dy),$$

whose second marginals are  $\mu_{\text{MC}}^t$  and  $\mu_{\text{NL}}^t$ , respectively. By the definition of the iterates and  $P$ , one can show that both couplings maximize the same objective  $F(\pi)$  but over different feasible sets:

$$\pi_{\text{MC}} = \operatorname{argmax}_{\int \pi(\cdot, dy) = \mu_{t-1}(\cdot)} F(\pi), \quad \pi_{\text{NL}} = \operatorname{argmax}_{\exists \mu: \pi = \mu_{t-1} \otimes \mu} F(\pi), \quad F(\pi) = \iint U(x, y)\pi(dx, dy) - D_{\text{KL}}(\pi \mid \mu_{t-1} \otimes \mu_{\text{ref}}),$$

Thus, NLHF restricts the coupling to be independent, whereas MCHF optimizes over all couplings with first marginal  $\mu_{t-1}$ . In this sense, MCHF can exploit more of the pairwise preference structure than the product coupling induced by NLHF.

This perspective also reveals a practical distinction. Since MCHF defines a conditional distribution  $P(x, \cdot)$ , it naturally supports inference-time refinement: starting from an output  $x$ , one can repeatedly run the Markov chain by  $P$  until the user is satisfied. This resembles the inference-time alignment mechanism in Stackelberg game-based alignment [15]. In contrast, NLHF directly updates a marginal distribution and does not provide such a conditional refinement mechanism.

**Alignment dynamics** Let  $\mu_{\text{MC}}(U)$  be the stationary distribution of the Markov kernel  $P(x, dy) \propto \exp(U(x, y))\mu_{\text{ref}}(dy)$ , and let  $\mu_{\text{NL}}(U)$  be the NLHF solution (3) based on the utility function  $U$ .

Suppose  $\|U\|_{\oplus} = 0$ . Then by Proposition 12, we have  $U(x, y) = \hat{g}(x) + \hat{f}(y)$  with  $\hat{f}(y) = \int \mu_{\text{ref}}(dx)U(x, y)$ . Moreover, by the definitions of MCHF and NLHF, if  $U = g \oplus f$ , then both  $\mu_{\text{MC}}(U)$  and  $\mu_{\text{NL}}(U)$  collapse to the RLHF solution (2) with reward  $f$ . In summary,

$$\|U\|_{\oplus} = 0 \quad \Rightarrow \quad U = \hat{g} \oplus \hat{f} \quad \Rightarrow \quad \mu_{\text{MC}}(U) = \mu_{\text{NL}}(U) = \mu_{\text{RL}}(\hat{f})$$

where  $\mu_{\text{RL}}(\hat{f})(dz) \propto \exp(\hat{f}(z))\mu_{\text{ref}}(dz)$ .

We aim to provide a quantitative version of this claim: if  $\|U\|_{\oplus}$  is small but not exactly zero, how far are  $\mu_{\text{MC}}(U)$  and  $\mu_{\text{NL}}(U)$  from  $\mu_{\text{RL}}(\hat{f})$ ? For simplicity, we focus on antisymmetric utilities satisfying  $U(x, y) = -U(y, x)$ . See Section E for the general utility case.

**Theorem 5** Suppose  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  is antisymmetric, and let  $\hat{f}(y) = \int \mu_{\text{ref}}(dx) U(x, y)$ . Let  $p_*(U) = d\mu_*(U)/d\mu_{\text{ref}}$  be the density for  $* \in \{\text{MC}, \text{NL}\}$ , and let  $p_{\text{RL}}(\hat{f}) = d\mu_{\text{RL}}(\hat{f})/d\mu_{\text{ref}}$ . Then,

$$\forall * \in \{\text{MC}, \text{NL}\}, \quad \|p_*(U) - p_{\text{RL}}(\hat{f}) - p_{\text{RL}}(\hat{f}) \odot (U - (-\hat{f}) \oplus \hat{f})^* p_{\text{RL}}(\hat{f})\|_{L^1(\mu_{\text{ref}})} = o(\|U\|_\oplus),$$

where  $\odot$  denotes entrywise multiplication and  $K^*h(\cdot) = \int \mu_{\text{ref}}(dx) K(x, \cdot)h(x)$ .

We can view  $p_{\text{RL}}(\hat{f}) + p_{\text{RL}}(\hat{f}) \odot (U - (-\hat{f}) \oplus \hat{f})^* p_{\text{RL}}(\hat{f})$  as a first-order approximation of  $p_*(U)$  around  $U = (-\hat{f}) \oplus \hat{f}$ . In particular, the second term captures the non-additive structure of the utility  $U$ . Interestingly, MCHF and NLHF agree up to the first-order term. Thus, by the triangle inequality,

$$d_{\text{TV}}(\mu_{\text{NL}}(U), \mu_{\text{MC}}(U)) = 2^{-1} \|p_{\text{NL}}(U) - p_{\text{MC}}(U)\|_1 = o(\|U\|_\oplus).$$

We emphasize that the right-hand side is  $o(\|U\|_\oplus)$ , not merely  $O(\|U\|_\oplus)$ .

The next theorem claims that MCHF and NLHF coincide at the iteration level.

**Theorem 6** Let  $\mu_{\text{MC}}^t$  and  $\mu_{\text{NL}}^t$  be the iterations defined in (25), and let  $p_*^t(U) = d\mu_*^t(U)/d\mu_{\text{ref}}$  be the density for  $* \in \{\text{MC}, \text{NL}\}$ . Then, for each  $* \in \{\text{MC}, \text{NL}\}$ ,

$$\begin{aligned} \|p_*^1(U) - p_{\text{RL}}(\hat{f}) - (U - (-\hat{f}) \oplus \hat{f})^* \mathbf{1}\|_{L^1(\mu_{\text{ref}})} &= o(\|U\|_\oplus), \\ \forall t \geq 2, \quad \|p_*^t(U) - p_{\text{RL}}(\hat{f}) - p_{\text{RL}}(\hat{f}) \odot (U - (-\hat{f}) \oplus \hat{f})^* p_{\text{RL}}(\hat{f})\|_{L^1(\mu_{\text{ref}})} &= o(\|U\|_\oplus). \end{aligned}$$

Thus, the MCHF and NLHF iterations also agree up to first order. Moreover, since the first-order approximation is the same for all  $t \geq 2$ , this implies that most of the update is completed within the first two iterations.

Finally, we verify Theorem 5 by numerical simulation. We set  $\mu_{\text{ref}}$  to be the uniform distribution on a discrete space with  $|\mathcal{X}| = 20$ , and generate utility  $U$  as

$$U(x, y) = 2^{-1} \log \frac{\mathcal{P}(x \prec y)}{1 - \mathcal{P}(x \prec y)}, \quad \mathcal{P}(x \prec y) = \Phi(R(y) - R(x) + 0.5E(x, y)), \quad (8)$$

where  $\Phi$  is the CDF of the standard normal distribution. We generate  $R \sim N(0, I_n)$ , and take  $E \in \mathbb{R}^{n \times n}$  as an antisymmetric matrix whose off-diagonal entries are drawn from  $N(0, 1)$ . We visualize  $U$  in Figure 1, which suggests that  $\|U\|_\oplus \ll \|U\|_\infty$ . Figure 2 plots MCHF, NLHF, RLHF, and the first-order approximation given by Theorem 5. We observe that MCHF and NLHF nearly coincide, and that their deviation from RLHF is accurately captured by the first-order approximation.

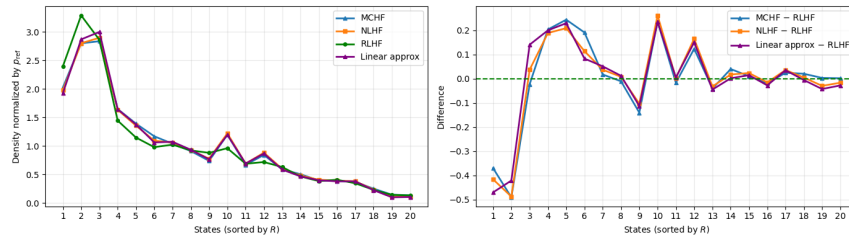


Figure 2: Comparison of MCHF, NLHF, RLHF, and the first-order approximation in Theorem 5.

## References

- [1] Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*, pages 4447–4455. PMLR, 2024.
- [2] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- [3] Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- [4] Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- [5] Xu Chu, Zhixin Zhang, Tianyu Jia, and Yujie Jin. Stackelberg self-annotation: A robust approach to data-efficient llm alignment. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- [6] Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024.
- [7] Paul Gözl, Nika Haghtalab, and Kunhe Yang. Distortion of ai alignment: Does preference optimization optimize for preferences? *arXiv preprint arXiv:2505.23749*, 2025.
- [8] Audrey Huang, Wenhao Zhan, Tengyang Xie, Jason D Lee, Wen Sun, Akshay Krishnamurthy, and Dylan J Foster. Correcting the mythos of kl-regularization: Direct alignment without overoptimization via chi-squared preference optimization. *arXiv preprint arXiv:2407.13399*, 2024.
- [9] R. Duncan Luce. *Individual Choice Behavior: A Theoretical Analysis*. Wiley, 1959.
- [10] Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward. *Advances in Neural Information Processing Systems*, 37:124198–124235, 2024.
- [11] Rémi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland, Zhaohan Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Côme Fiegel, et al. Nash learning from human feedback. In *Forty-first International Conference on Machine Learning*, 2024.
- [12] Sahand Negahban, Sewoong Oh, and Devavrat Shah. Rank centrality: Ranking from pairwise comparisons. *Operations Research*, 65(1):266–287, 2017.

- [13] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744, 2022.
- [14] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.
- [15] Barna Pásztor, Thomas Kleine Buening, and Andreas Krause. Stackelberg learning from human feedback: Preference optimization as a sequential game. In *The Fourteenth International Conference on Learning Representations*, 2026.
- [16] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741, 2023.
- [17] Corby Rosset, Ching-An Cheng, Arindam Mitra, Michael Santacrose, Ahmed Awadallah, and Tengyang Xie. Direct nash optimization: Teaching language models to self-improve with general preferences. *arXiv preprint arXiv:2404.03715*, 2024.
- [18] Zhekun Shi, Kaizhao Liu, Qi Long, Weijie J Su, and Jiancong Xiao. Fundamental limits of game-theoretic llm alignment: Smith consistency and preference matching. *arXiv preprint arXiv:2505.20627*, 2025.
- [19] Daniil Tiapkin, Daniele Calandriello, Denis Belomestny, Eric Moulines, Alexey Naumov, Kashif Rasul, Michal Valko, and Pierre Menard. Accelerating nash learning from human feedback via mirror prox. *arXiv preprint arXiv:2505.19731*, 2025.
- [20] Fang Wu, Xu Huang, Weihao Xuan, Zhiwei Zhang, Yijia Xiao, Guancheng Wan, Xiaomin Li, Bing Hu, Peng Xia, Jure Leskovec, and Yejin Choi. Multiplayer nash preference optimization. In *The Fourteenth International Conference on Learning Representations*, 2026.

## Appendix A. Perturbation analysis

We analyze the sensitivity of the equilibrium distribution to perturbations in  $U$ . This is important because, in practice, the utility function  $U$  may be updated as new preference data become available. We may also consider a setting in which  $U$  is taken to be an increasing transformation of the preference model  $\mathcal{P}(x \succ y)$ , which is itself estimated, for example, by logistic regression as in the original paper [11]. In this case, estimation error in  $\mathcal{P}$  induces a perturbation in  $U$ . It is therefore important to understand how such perturbations affect the resulting equilibrium.

We show that both MCHF and NLHF are locally Lipschitz with respect to perturbations of  $U$  under the  $L^\infty$  metric. The proofs of the theorem stated in this section are given in the subsequent sections.

### A.1. MCHF

First, we discuss the sensitivity of MCHF. Let

$$c(\|U\|_\oplus) = \min(1 - e^{-2\|U\|_\oplus}, \|U\|_\oplus) < 1$$

be the constant given by Theorem 2, which upper bounds the Lipschitz constant of the map  $\mu \mapsto \mu P$  under the TV distance.

**Theorem 7** *For any  $\hat{U}, U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , let  $P$  and  $\hat{P}$  be the Markov kernels (4) constructed from  $U$  and  $\hat{U}$ , respectively. Then, for all  $t \geq 1$ ,*

$$d_{\text{TV}}(\mu_{\text{ref}} P^t, \mu_{\text{ref}} \hat{P}^t) \leq \frac{1}{2} \|U - \hat{U}\|_\infty \cdot \frac{1 - (c(\|U\|_\oplus) \wedge c(\|\hat{U}\|_\oplus))^t}{1 - c(\|U\|_\oplus) \wedge c(\|\hat{U}\|_\oplus)}.$$

See Section C.2 for the proof. Combining this result with Corollary 3, and letting  $\mu_{\text{MC}}(U)$  and  $\mu_{\text{MC}}(\hat{U})$  denote the stationary distributions of  $P$  and  $\hat{P}$ , respectively, we obtain the following corollary.

**Corollary 8** *For any  $U, \hat{U} \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ ,*

$$d_{\text{TV}}(\mu_{\text{MC}}(U), \mu_{\text{MC}}(\hat{U})) \leq \frac{1}{2} \cdot \frac{\|U - \hat{U}\|_\infty}{1 - c(\|U\|_\oplus) \wedge c(\|\hat{U}\|_\oplus)}.$$

Since  $c(\|U\|_\oplus)$  is always strictly less than 1, we conclude that the map  $U \mapsto \mu_{\text{MC}}(U)$  is locally Lipschitz.

### A.2. NLHF

We now turn to NLHF and first establish the following stability result.

**Theorem 9** *For any  $U, \hat{U} \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ ,*

$$D_{\text{KL}}(\mu_{\text{NL}}(U) \parallel \mu_{\text{NL}}(\hat{U})) \leq 2\|U - \hat{U}\|_\infty.$$

This theorem follows from the strong convexity induced by the KL regularization term. Combined with Pinsker's inequality,  $d_{\text{TV}} \leq \sqrt{D_{\text{KL}}/2}$ , it yields

$$d_{\text{TV}}(\mu_{\text{NL}}(U), \mu_{\text{NL}}(\hat{U})) \leq \|U - \hat{U}\|_{\infty}^{1/2}.$$

Thus, the map  $U \mapsto \mu_{\text{NL}}(U)$  is globally Hölder-1/2 continuous. The next theorem shows that this map becomes Lipschitz when the domain is suitably restricted.

**Theorem 10** *If  $\|U\|_{\oplus} \wedge \|\hat{U}\|_{\oplus} < 1$ , then*

$$d_{\text{TV}}(\mu_{\text{NL}}(U), \mu_{\text{NL}}(\hat{U})) \leq \frac{1}{2} \cdot \frac{\|U - \hat{U}\|_{\infty}}{1 - (\|U\|_{\oplus} \wedge \|\hat{U}\|_{\oplus})}.$$

Compared with Corollary 8, this result shows that  $U \mapsto \mu_{\text{NL}}(U)$  is also locally Lipschitz as a map from  $\|\cdot\|_{\infty}$  to  $d_{\text{TV}}$ , but the domain is restricted to the smaller set  $\{U \in L^{\infty}(\mu_{\text{ref}} \otimes \mu_{\text{ref}}) : \|U\|_{\oplus} < 1\}$ . We next show that  $\mu_{\text{NL}}$  satisfies a global pseudo-Lipschitz bound.

**Theorem 11** *There exists an absolute constant  $C$  such that, for any  $U, \hat{U} \in L^{\infty}(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ ,*

$$d_{\text{TV}}(\mu_{\text{NL}}(U), \mu_{\text{NL}}(\hat{U})) \leq C(1 + \|U\|_{\oplus}^6 + \|\hat{U}\|_{\oplus}^6)\|U - \hat{U}\|_{\infty}.$$

The proof of this theorem requires a different technique from the preceding arguments. We first establish the differentiability of the map  $\epsilon \mapsto \mu_{\text{NL}}(U + \epsilon E)$  for fixed  $U, E \in L^{\infty}(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  and then bound the  $L^1$  norm of its derivative.

## Appendix B. Characterization of seminorm $\|\cdot\|_{\oplus}$

**Proposition 12** *Fix  $U \in L^{\infty}(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ . Define  $\square(U)$  as the rectangle defect:*

$$\square(U) = \text{esssup}_{x, x', y, y'} |U(x', y') - U(x, y') - U(x', y) + U(x, y)|.$$

*Let  $(\hat{g}, \hat{f})$  be a solution to  $\min_{g, f} \|U - g \oplus f\|_2$ . Then*

$$\square(U) \asymp \|U\|_{\oplus} \asymp \|U - \hat{g} \oplus \hat{f}\|_{L^{\infty}}.$$

*More precisely,  $\frac{1}{4}\square(U) \leq \|U\|_{\oplus} \leq \|U - \hat{g} \oplus \hat{f}\|_{\infty} \leq \square(U)$ .*

The minimizer of the  $L^2$  projection problem  $\min_{g, f} \|U - g \oplus f\|_2$  is not unique, but one solution is given explicitly by

$$\hat{g}(x) = \int U(x, y) \mu_{\text{ref}}(dy) - m, \quad \hat{f}(y) = \int U(x, y) \mu_{\text{ref}}(dx),$$

where  $m = \iint U(x, y) \mu_{\text{ref}}(dx) \mu_{\text{ref}}(dy)$ .

When  $U$  is antisymmetric, we obtain a sharper characterization.

**Proposition 13** *Suppose  $U \in L^{\infty}(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  is antisymmetric, i.e.,  $U(x, y) = -U(y, x)$ . Define  $\Delta(U)$  as the triangle defect:*

$$\Delta(U) = \text{esssup}_{x, y, z} |U(x, y) + U(y, z) + U(z, x)|.$$

*Let  $\hat{f}(y) = \int \mu_{\text{ref}}(dx) U(x, y)$ . Then*

$$\|U\|_{\oplus} \asymp \Delta(U) \asymp \|U - (-\hat{f}) \oplus \hat{f}\|_{\infty}.$$

*More precisely,  $\frac{1}{3}\Delta(U) \leq \|U\|_{\oplus} \leq \|U - (-\hat{f}) \oplus \hat{f}\|_{\infty} \leq \Delta(U)$ .*

**B.1. Proof of Proposition 12**

Take  $(\hat{g}, \hat{f}) \in \operatorname{argmin}_{g, f \in L^2} \|U - g \oplus f\|_2$  as

$$\hat{g}(x) = \int U(x, y) \mu_{\text{ref}}(dy) - m, \quad \hat{f}(y) = \int U(x, y) \mu_{\text{ref}}(dx), \quad m = \iint U(x, y) \mu_{\text{ref}}(dx) \mu_{\text{ref}}(dy)$$

Here, we have  $\hat{g}, \hat{f} \in L^\infty$  by  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , and hence  $\|U\|_\oplus \leq \|U - \hat{g} \oplus \hat{f}\|_\infty$  holds by the definition of  $\|U\|_\oplus$ . Thus, it suffices to show  $4^{-1}\square(U) \leq \|U\|_\oplus$  and  $\|U - \hat{g} \oplus \hat{f}\|_\infty \leq \square(U)$ .

**Proof of  $4^{-1}\square(U) \leq \|U\|_\oplus$**  Fix arbitrary  $g, f \in L^\infty$ , and set  $U_{g,f} = U - g \oplus f$ . Then for every  $x, x', y, y'$ , noting that the  $g$ - and  $f$ -terms cancel,

$$U(x', y') - U(x, y') - U(x', y) + U(x, y) = U_{g,f}(x', y') - U_{g,f}(x, y') - U_{g,f}(x', y) + U_{g,f}(x, y),$$

Thus, taking  $\operatorname{esssup}_{x, y, x', y'}$  on both sides, the LHS is  $\square(U)$ , while using the triangle inequality for the RHS, we get

$$\square(U) \leq 4\|U_{g,f}\|_\infty$$

Since this holds for every  $f, g$ , taking  $\inf_{g, f \in L^\infty}$  on the RHS, we complete the proof.

**Proof of  $\|U - \hat{g} \oplus \hat{f}\|_\infty \leq \square(U)$**  By the definition of  $\hat{g}$  and  $\hat{f}$ , for all  $(x, y)$ ,

$$\begin{aligned} (U - \hat{g} \oplus \hat{f})(x, y) &= U(x, y) - \int U(x, y') \mu_{\text{ref}}(dy') - \int U(x', y) \mu_{\text{ref}}(dx') + \iint U(x', y') \mu_{\text{ref}}(dx') \mu_{\text{ref}}(dy') \\ &= \iint (U(x, y) - U(x', y) - U(x, y') + U(x', y')) \mu_{\text{ref}}(dx') \mu_{\text{ref}}(dy'). \end{aligned}$$

Hence, taking  $\operatorname{esssup}_{x, y} |\cdot|$  on both sides,

$$\begin{aligned} \|U - \hat{g} \oplus \hat{f}\|_\infty &\leq \operatorname{esssup}_{x, y} \left| \iint U(x, y) - U(x', y) - U(x, y') + U(x', y') \mu_{\text{ref}}(dx') \mu_{\text{ref}}(dy') \right| \\ &\leq \operatorname{esssup}_{x, y, x', y'} |U(x, y) - U(x', y) - U(x, y') + U(x', y')| \\ &= \square(U), \end{aligned}$$

so the proof is complete.

**B.2. Proof of Proposition 13**

We first claim  $\|U\|_\oplus = \inf_f \|U - (-f) \oplus f\|_\infty$ . To prove this, it suffices to show  $\inf_{g, f} \|U - g \oplus f\|_\infty \geq \inf_f \|U - (-f) \oplus f\|_\infty$ . Let us fix  $g, f \in L^\infty$  and let  $h = 2^{-1}(g + f)$ . Note

$$(-h) \oplus h = \frac{1}{2}g \oplus f - \frac{1}{2}f \oplus g$$

and the  $\|U + f \oplus g\|_\infty = \|-U + g \oplus f\|_\infty$  since  $U$  is antisymmetric. Therefore, by the triangle inequality,

$$\inf_f \|U - (-f) \oplus f\|_\infty \leq \|U - (-h) \oplus h\|_\infty \leq \frac{1}{2}\|U - g \oplus f\|_\infty + \frac{1}{2}\|U + f \oplus g\|_\infty = \|U - g \oplus f\|_\infty.$$

Taking  $\inf_{g, f}$  on the RHS, we obtain the claim.

Therefore, we have  $\|U\|_\oplus \leq \|U - (\hat{f}) \oplus \hat{f}\|_\infty$  for  $\hat{f}(y) = \int \mu_{\text{ref}}(dx) U(x, y)$ . Furthermore, using  $\|U\|_\oplus = \inf_f \|U - (-f) \oplus f\|_\infty$  and the same argument as that of Section B.1, one can show  $3^{-1}\Delta(U) \leq \|U\|_\oplus$  and  $\|U - (\hat{f}) \oplus \hat{f}\|_\infty \leq \Delta(U)$ . Thus, we complete the proof.

## Appendix C. Proof of MCHF

### C.1. Proof of Theorem 2

We prove  $d_{\text{TV}}(\mu\text{P}, \nu\text{P}) \leq (1 - e^{-2\|U\|_{\oplus}})d_{\text{TV}}(\mu, \nu)$  and  $d_{\text{TV}}(\mu\text{P}, \nu\text{P}) \leq \|U\|_{\oplus}d_{\text{TV}}(\mu, \nu)$  separately.

C.1.1. PROOF OF  $d_{\text{TV}}(\mu\text{P}, \nu\text{P}) \leq (1 - e^{-2\|U\|_{\oplus}})d_{\text{TV}}(\mu, \nu)$  VIA MINORIZATION LEMMA

Let  $L^1 = L^1(\mu_{\text{ref}})$  and  $L_0^1 = \{h \in L^1 : \int \mu_{\text{ref}}(dy)h(y) = 1\}$ . We denote the  $L_p$  norm as  $\|f\|_p = \|f\|_{L^p}$ . We also write  $L^\infty = L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  and  $\|U\|_\infty = \|U\|_{L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})}$  for any  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  when the context is clear. For any  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , define the Markov transition density  $P \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  as

$$P(x, y) = \frac{\exp(U(x, y))}{\int \mu_{\text{ref}}(dy') \exp(U(x, y'))}, \quad (9)$$

and we define the bounded linear operator  $P^* : L^1 \rightarrow L^1$  as

$$P^*f(y) = \int \mu_{\text{ref}}(dx) f(x) P(x, y).$$

For all  $\mu, \nu \in \mathcal{P}_{\mu_{\text{ref}}}$ , we can write the TV distance as

$$\begin{aligned} d_{\text{TV}}(\mu, \nu) &= \frac{1}{2} \left\| \frac{d\mu}{d\mu_{\text{ref}}} - \frac{d\nu}{d\mu_{\text{ref}}} \right\|_1 \\ d_{\text{TV}}(\mu\text{P}, \nu\text{P}) &= \frac{1}{2} \left\| P^* \left( \frac{d\mu}{d\mu_{\text{ref}}} - \frac{d\nu}{d\mu_{\text{ref}}} \right) \right\|_1. \end{aligned}$$

Since  $\frac{d\mu}{d\mu_{\text{ref}}} - \frac{d\nu}{d\mu_{\text{ref}}} \in L_0^1 = \{h \in L^1 : \int \mu_{\text{ref}}(dy)h(y) = 0\}$ , it suffices to show the following lemma.

**Lemma 14 (Minorization lemma)** *For any  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , let  $P$  be the Markov density given by (9). Then, the map  $h \mapsto P^*h$  is a contraction mapping on  $L_0^1 = \{h \in L^1 : \int \mu_{\text{ref}}(dy)h(y) = 0\}$  and with its contraction rate given by*

$$\sup_{h \in L_0^1, h \neq 0} \frac{\|P^*h\|_1}{\|h\|_1} \leq 1 - \exp(-2\|U\|_{\oplus})$$

where  $\|U\|_{\oplus} = \inf_{g, f} \|U - g \oplus f\|_\infty$ .

**Proof** Let us fix  $g, f \in L^\infty$  and let  $U_{g, f} = U - g \oplus f$ . Define  $p_f \in L^\infty$  as

$$p_f(y) = \frac{\exp(f(y))}{Z}, \quad Z = \int \mu_{\text{ref}}(dy') \exp(f(y'))$$

so that  $p_f$  is a valid density with respect to  $\mu_{\text{ref}}$ , i.e.,  $p_f(\cdot) \geq 0$  and  $\int \mu_{\text{ref}}(dy)p_f(y) = 1$ . Now, we can rewrite the Markov density  $P$  as

$$\begin{aligned} P(x, y) &= \frac{\exp(U(x, y))}{\int \mu_{\text{ref}}(dy') \exp(U(x, y'))} \\ &= \frac{\exp(U(x, y) - g(x) - f(y)) \exp(f(y))}{\int \mu_{\text{ref}}(dy') \exp(U(x, y') - g(x) - f(y')) \exp(f(y'))} \\ &= \frac{\exp(U_{g, f}(x, y)) p_f(y)}{\int \mu_{\text{ref}}(dy') \exp(U_{g, f}(x, y')) p_f(y')} \end{aligned} \quad (10)$$

Then, for almost every  $(x, y)$  with respect to the product measure  $\mu_{\text{ref}} \otimes \mu_{\text{ref}}$ , noting  $\int \mu_{\text{ref}}(dy') p_f(y') = 1$ ,

$$P(x, y) \geq \frac{\exp(-\|U_{g,f}\|_\infty)}{\exp(\|U_{g,f}\|_\infty)} \cdot \frac{p_f(y)}{\int \mu_{\text{ref}}(dy') p_f(y')} = \delta \cdot p_f(y)$$

where we take  $\delta$  as  $\delta = \exp(-2\|U_{g,f}\|_\infty) \in (0, 1]$ . If we define  $\tilde{P}$  as

$$\tilde{P}(x, y) = \frac{P(x, y) - \delta p_f(y)}{1 - \delta}$$

then  $\tilde{P}$  satisfies  $\tilde{P}(x, \cdot) \geq 0$  and  $\int \mu_{\text{ref}}(dy) \tilde{P}(x, y) = 1$ , that is,  $\tilde{P}$  is a valid Markov transition density with respect to  $\mu_{\text{ref}}$ . Rearranging this,

$$P(x, y) = \delta \cdot p_f(y) + (1 - \delta) \tilde{P}(x, y).$$

For any  $h \in L_0^1$  such that  $h \neq 0$ , noting  $\int \mu_{\text{ref}}(dx) h(x) = 0$ , we get

$$\begin{aligned} P^*h(y) &= \int \mu_{\text{ref}}(dx) P(x, y) h(x) \\ &= \delta p_f(y) \int \mu_{\text{ref}}(dx) h(x) + (1 - \delta) \int \mu_{\text{ref}}(dx) \tilde{P}(x, y) h(x) \\ &= (1 - \delta) \tilde{P}^*h(y). \end{aligned}$$

Taking the  $L_1$ -norm and using the fact (a.k.a. Data Processing Inequality):

$$\|\tilde{P}^*h\|_1 = \int \mu_{\text{ref}}(dy) \left| \int \mu_{\text{ref}}(dx) \tilde{P}(x, y) h(x) \right| \leq \int \mu_{\text{ref}}(dx) \mu_{\text{ref}}(dy) \tilde{P}(x, y) |h(x)| = \int \mu_{\text{ref}}(dx) |h(x)| = \|h\|_1,$$

where we have used the Fubini's theorem and the fact that  $\tilde{P}$  is the Markov transition probability, we get

$$\|P^*h\|_1 \leq (1 - \delta) \|\tilde{P}^*h\|_1 \leq (1 - \delta) \|h\|_1$$

so that

$$\sup_{h \in L_0^1: h \neq 0} \frac{\|P^*h\|_1}{\|h\|_1} \leq 1 - \delta.$$

Since  $\delta = \exp(-2\|U_{g,f}\|_\infty)$  and  $\|U_{g,f}\|_\infty = \|U - g \oplus f\|_\infty$ , taking  $\inf_{g,f}$  on the RHS, we complete the proof of (14).  $\blacksquare$

### C.1.2. PROOF OF $d_{\text{TV}}(\mu\mathbf{P}, \nu\mathbf{P}) \leq \|U\|_\oplus d_{\text{TV}}(\mu, \nu)$

It suffices to show

$$\sup_{h \in L_0^1: h \neq 0} \frac{\|P^*h\|_1}{\|h\|_1} \leq \|U\|_\oplus.$$

We prove this using the following two lemmas.

**Lemma 15 (Doebelin-Dobrushin characterization)** For any  $P \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ ,

$$\sup_{h \in L_0^1: h \neq 0} \frac{\|P^*h\|_1}{\|h\|_1} \leq \text{esssup}_{x, x'} \frac{1}{2} \|P(x, \cdot) - P(x', \cdot)\|_1$$

where  $P^*h(y) = \int \mu_{\text{ref}}(dx) h(x) P(x, y)$ .

**Lemma 16 (Lipschitz continuity of Gibbs measure)** Fix a probability measure  $\mu$  and define the map  $\mathcal{G} : L^\infty(\mu) \rightarrow L^\infty(\mu)$  as

$$\mathcal{G}(g) \equiv \frac{\exp(g(y))}{\int \mu(dy') \exp(g(y'))}.$$

Then, for all  $g, \tilde{g} \in L^\infty(\mu)$ , we have

$$\|\mathcal{G}(g) - \mathcal{G}(\tilde{g})\|_{L^1(\mu)} \leq \|g - \tilde{g}\|_{L^\infty(\mu)}$$

Applying Lemma 15, we have

$$\sup_{h \in L_0^1: h \neq 0} \frac{\|P^*h\|_1}{\|h\|_1} \leq \text{esssup}_{x, x'} \frac{1}{2} \|P(x, \cdot) - P(x', \cdot)\|_1.$$

Recall that for any  $g, f \in L^\infty$ ,  $P(x, y)$  can be written as follows (see (10) for the derivation):

$$P(x, y) = \frac{\exp(U_{g,f}(x, y)) p_f(y)}{\int \mu_{\text{ref}}(dy') \exp(U_{g,f}(x, y')) p_f(y')}, \quad p_f(y) = \frac{\exp(f(y))}{\int \mu_{\text{ref}}(dy') \exp(f(y'))}$$

where  $U_{g,f} = U - g \oplus f$ . Therefore, letting  $\mu_f(dy) = \mu_{\text{ref}}(dy) p_f(y)$  be the exponentially tilted probability measure by  $f$ , using Lemma 16  $\mu = \mu_f$  and  $g(y) = U_{g,f}(x, y)$  and  $\tilde{g}(y) = U_{g,f}(x', y)$ , we have

$$\|P(x, \cdot) - P(x', \cdot)\|_1 \leq \|U_{g,f}(x, \cdot) - U_{g,f}(x', \cdot)\|_{L^\infty(\mu_f)} = \|U_{g,f}(x, \cdot) - U_{g,f}(x', \cdot)\|_{L^\infty(\mu_{\text{ref}})},$$

where the second equation follows from since  $\mu_f$  and  $\mu_{\text{ref}}$  are absolutely continuous with respect to each other. Thus, by the definition:  $\text{esssup}_y |f(y)| = \|f\|_{L^\infty(\mu_{\text{ref}})}$ ,

$$\begin{aligned} \sup_{h \in L_0^1: h \neq 0} \frac{\|P^*h\|_1}{\|h\|_1} &\leq \text{esssup}_{x, x'} \frac{1}{2} \text{esssup}_y |U_{g,f}(x, y) - U_{g,f}(x', y)| \\ &\leq \text{esssup}_{x, y} |U_{g,f}(x, y)| && \text{triangle inequality} \\ &= \|U_{g,f}\|_\infty \end{aligned}$$

Since  $U_{g,f} = U - g \oplus f$  and  $g, f \in L^\infty(\mu_{\text{ref}})$  are arbitrary, taking  $\inf_{g,f}$  on the RHS, we complete the proof of  $\sup_{h \in L_0^1: h \neq 0} \frac{\|P^*h\|_1}{\|h\|_1} \leq \|U\|_\oplus$ . Below, we prove the intermediate lemmas we used.

## C.1.3. PROOF OF LEMMA 15

Below we denote  $\|X\|_p = \|X\|_{L^p(\mu_{\text{ref}})}$  for simplicity. Note that  $\|P^*h\|_1$  can be represented as

$$\|P^*h\|_1 = \sup_{\|f\|_\infty \leq 1} \left| \int \mu_{\text{ref}}(dy) f(y) P^*h(y) \right|.$$

Using Fubini's theorem, for all  $f \in L^\infty$ ,

$$\int \mu_{\text{ref}}(dy) f(y) P^*h(y) = \int \mu_{\text{ref}}(dy) f(y) \int \mu_{\text{ref}}(dx) P(x, y) h(x) = \int \mu(dx) h(x) Pf(x)$$

where  $Pf(x) \equiv \int \mu_{\text{ref}}(dy) P(x, y) f(y)$ . Since  $\int \mu_{\text{ref}}(dx) h(x) = 0$  by  $h \in L_0^1$ , for any constant  $c \in \mathbb{R}$ , we have

$$\int \mu_{\text{ref}}(dx) h(x) Pf(x) = \int \mu_{\text{ref}}(dx) h(x) (Pf(x) - c).$$

Now we let  $c = 2^{-1}(\text{esssup}_x Pf(x) + \text{essinf}_x Pf(x)) \in \mathbb{R}$  so that

$$\text{esssup}_x |Pf(x) - c| \leq \frac{1}{2} \left( \text{esssup}_x Pf(x) - \text{essinf}_x Pf(x) \right).$$

Combined with Hölder's inequality

$$\left| \int \mu_{\text{ref}}(dx) h(x) (Pf(x) - c) \right| \leq \|h\|_1 \cdot \frac{1}{2} \left( \text{esssup}_x Pf(x) - \text{essinf}_x Pf(x) \right).$$

Notice that for almost every  $x, x'$ ,

$$\begin{aligned} \frac{1}{2} \left( Pf(x) - Pf(x') \right) &= \frac{1}{2} \int \mu_{\text{ref}}(dy) f(y) (P(x, y) - P(x', y)) \\ &\leq \frac{1}{2} \|P(x, \cdot) - P(x', \cdot)\|_1 \cdot \|f\|_\infty \end{aligned}$$

Therefore, taking  $\text{esssup}_{x, x'}$ , we are left with

$$\frac{1}{2} \left( \text{esssup}_x Pf(x) - \text{essinf}_x Pf(x) \right) \leq \|f\|_\infty \cdot \text{esssup}_{x, x'} \frac{1}{2} \|P(x, \cdot) - P(x', \cdot)\|_1$$

Putting all together, we obtain

$$\left| \int \mu_{\text{ref}}(dy) f(y) P^*h(y) \right| \leq \|f\|_\infty \|h\|_1 \text{esssup}_{x, x'} \frac{1}{2} \|P(x, \cdot) - P(x', \cdot)\|_1$$

for all  $f \in L^\infty$  and  $h \in L_0^1$ . Taking supremum of  $f$  such that  $\|f\|_\infty \leq 1$ , we complete the proof.

## C.1.4. PROOF OF LEMMA 16

Let  $h(x) = g(x) - \tilde{g}(x)$  and define the interpolator  $g_t \in L^\infty$  for  $t \in [0, 1]$  as

$$\forall t \in [0, 1], \quad g_t \equiv \tilde{g} + th$$

so that  $g_0 = \tilde{g}$  and  $g_1 = g$ . Now we define the tilted probability measure  $\mu_t$  as

$$\mu_t(dy) \equiv \mathcal{G}(g_t)(y)\mu(dy).$$

Note that  $\mu_t$  and  $\mu$  are absolutely continuous with respect to each other. Then,

$$\|\mathcal{G}(g) - \mathcal{G}(\tilde{g})\|_1 = \sup_{\|f\|_{L^\infty(\mu)} \leq 1} \left| \int f(x)\mu_1(dx) - \int f(x)\mu_0(dx) \right| = \sup_{\|f\|_{L^\infty(\mu)} \leq 1} \left| \phi_f(1) - \phi_f(0) \right|$$

where

$$\forall t \in [0, 1], \quad \phi_f(t) \equiv \int f(x)\mu_t(dx) \equiv \int \mu(dx)f(x) \frac{\exp(g_t(x))}{\int \mu(dx') \exp(g_t(x'))}.$$

Taking the derivative of  $\phi_f(t)$ , noting  $\frac{d}{dt}g_t(x) = h(x)$ , and  $h$  is bounded,

$$\begin{aligned} \frac{d}{dt}\phi_f(t) &= \int \mu(dx)f(x) \frac{\exp(g_t(x))h(x)}{\int \mu(dx') \exp(g_t(x'))} - \int \mu(dx)f(x) \exp(g_t(x)) \cdot \frac{\int \mu(dx') \exp(g_t(x'))h(x')}{\left(\int \mu(dx') \exp(g_t(x'))\right)^2} \\ &= \int \mu_t(dx)f(x)h(x) - \left(\int \mu_t(dx)f(x)\right) \cdot \left(\int \mu_t(dx)h(x)\right) \\ &= \text{Cov}_{X \sim \mu_t}(f(X), h(X)) \end{aligned}$$

Using the Cauchy–Schwarz inequality  $|\text{Cov}(X, Y)| \leq \sqrt{\text{Var}(X)\text{Var}(Y)} \leq \|X\|_\infty \|Y\|_\infty$ , noting that  $\mu_t$  and  $\mu$  are absolutely continuous with respect to each other, we obtain the uniform upper bound:

$$\left| \frac{d}{dt}\phi_f(t) \right| \leq \|f\|_{L^\infty(\mu_t)} \cdot \|h\|_{L^\infty(\mu_t)} = \|f\|_{L^\infty(\mu)} \|h\|_{L^\infty(\mu)}.$$

Therefore, we get

$$\|\mathcal{G}(g) - \mathcal{G}(\tilde{g})\|_1 = \sup_{\|f\|_{L^\infty} \leq 1} \left| \phi_f(1) - \phi_f(0) \right| \leq \sup_{\|f\|_{L^\infty(\mu)} \leq 1} \int_0^1 \left| \frac{d}{dt}\phi_f(t) \right| dt \leq \|h\|_{L^\infty(\mu)}.$$

This completes the proof.

## C.2. Proof of Theorem 7

Let  $c$  and  $\hat{c}$  be the Lipschitz constants (given by Theorem 2) of  $\mu \mapsto P\mu$  and  $\mu \mapsto \hat{P}\mu$ , respectively. Consider the iterations

$$\mu_t = \mu_{t-1}P, \quad \hat{\mu}_t = \hat{\mu}_{t-1}\hat{P},$$

initialized with  $\mu_0 = \hat{\mu}_0 \in \mathcal{P}_{\mu_{\text{ref}}}$ . Using the triangle inequality, and the Lipschitz continuity of  $\mu \mapsto \mu P$ ,

$$\begin{aligned} d_{\text{TV}}(\mu_t, \hat{\mu}_t) &\leq d_{\text{TV}}(\mu_{t-1}P, \hat{\mu}_{t-1}P) + d_{\text{TV}}(\hat{\mu}_{t-1}P, \hat{\mu}_{t-1}\hat{P}) \\ &\leq c \cdot d_{\text{TV}}(\mu_{t-1}, \hat{\mu}_{t-1}) + d_{\text{TV}}(\hat{\mu}_{t-1}P, \hat{\mu}_{t-1}\hat{P}). \end{aligned}$$

By the definition of TV distance,

$$d_{\text{TV}}(\hat{\mu}_{t-1}\mathbf{P}, \hat{\mu}_{t-1}\hat{\mathbf{P}}) = \frac{1}{2} \sup_{\|f\|_{\infty} \leq 1} \left| \int f(y) [(\hat{\mu}_{t-1}\mathbf{P})(dy) - (\hat{\mu}_{t-1}\hat{\mathbf{P}})(dy)] \right|.$$

By Fubini's lemma, for any measurable function  $f$  with  $\|f\|_{\infty} \leq 1$ ,

$$\begin{aligned} & \frac{1}{2} \left| \int f(y) [(\hat{\mu}_{t-1}\mathbf{P})(dy) - (\hat{\mu}_{t-1}\hat{\mathbf{P}})(dy)] \right| \\ &= \frac{1}{2} \left| \int \mu_{\text{ref}}(dy) f(y) \int \hat{\mu}_{t-1}(dx) (P(x, y) - \hat{P}(x, y)) \right| \\ &= \frac{1}{2} \left| \int \hat{\mu}_{t-1}(dx) \int \mu_{\text{ref}}(dy) f(y) (P(x, y) - \hat{P}(x, y)) \right| \\ &\leq \frac{1}{2} \int \hat{\mu}_{t-1}(dx) \left| \int \mu_{\text{ref}}(dy) f(y) (P(x, y) - \hat{P}(x, y)) \right| \\ &\leq \frac{1}{2} \int \hat{\mu}_{t-1}(dx) \|P(x, \cdot) - \hat{P}(x, \cdot)\|_1 \quad \|f\|_{\infty} \leq 1 \\ &\leq \frac{1}{2} \int \hat{\mu}_{t-1}(dx) \|U(x, \cdot) - \hat{U}(x, \cdot)\|_{\infty} \quad \text{by Lemma 16 with } \mu = \mu_{\text{ref}}, g(\cdot) = U(x, \cdot) \text{ and } \tilde{g}(\cdot) = \hat{U}(x, \cdot) \\ &\leq \frac{1}{2} \|U - \hat{U}\|_{\infty} \quad \hat{\mu}_{t-1} \ll \mu_{\text{ref}} \end{aligned}$$

Therefore,

$$d_{\text{TV}}(\mu_t, \hat{\mu}_t) \leq c \cdot d_{\text{TV}}(\mu_{t-1}, \hat{\mu}_{t-1}) + \frac{1}{2} \|U - \hat{U}\|_{\infty}.$$

Iterating this inequality and using the initial condition  $\mu_0 = \hat{\mu}_0$ , we get

$$d_{\text{TV}}(\mu_t, \hat{\mu}_t) \leq \frac{1}{2} \|U - \hat{U}\|_{\infty} \sum_{s=0}^{t-1} c^s.$$

By symmetry, the same upper bound also holds with  $c$  replaced by  $\hat{c}$ . Thus, we complete the proof.

### C.3. Differentiability of MCHF

Let us fix  $U, E \in L^{\infty}(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ . For each  $\epsilon \in \mathbb{R}$ , let  $p_{\epsilon}$  be the unique stationary distribution of the Markov transition density  $P_{\epsilon}$  defined by

$$P_{\epsilon}(x, y) = \frac{\exp(U_{\epsilon}(x, y))}{\int \mu_{\text{ref}}(dy') \exp(U_{\epsilon}(x, y'))}, \quad U_{\epsilon}(x, y) \equiv U(x, y) + \epsilon E(x, y),$$

Now we define the linear operator  $P_{\epsilon}^* : L_1 \rightarrow L_1$  as

$$P_{\epsilon}^* f(y) = \int \mu_{\text{ref}}(dx) f(x) P_{\epsilon}(x, y).$$

Let  $\mu_{\epsilon}$  be the stationary distribution of the Markov kernel  $P_{\epsilon}(x, dy) = P_{\epsilon}(x, y) \mu_{\text{ref}}(y)$  and let  $p_{\epsilon} = \frac{d\mu_{\epsilon}}{d\mu_{\text{ref}}}$  be its density. Note that this  $p_{\epsilon}$  satisfies

$$P_{\epsilon}^* p_{\epsilon} = p_{\epsilon}, \quad \int \mu_{\text{ref}}(dy) p_{\epsilon}(y) = 1.$$

The next theorem calculates the derivative of  $p_{\epsilon}$  in  $\epsilon$  using the Implicit function theorem.

**Theorem 17** *The map  $\epsilon \mapsto p_\epsilon$  is differentiable at  $\epsilon = 0$  in the sense that*

$$\lim_{\epsilon \rightarrow 0} \left\| \frac{p_\epsilon - p_0}{\epsilon} - \dot{p}_0 \right\|_1 = 0$$

where the derivative  $\dot{p}_0 \in L_0^1$  is given by

$$\dot{p}_0 = \left( (I - P_0^*) \big|_{L_0^1} \right)^{-1} \left( y \mapsto \int \mu_{\text{ref}}(dx) p_0(x) P_0(x, y) \left( E(x, y) - \int \mu_{\text{ref}}(dy') P_0(x, y') E(x, y') \right) \right)$$

where  $L_0^1 = \{f \in L^1 : \int \mu_{\text{ref}}(dy) f(y) = 0\}$ .

**Proof** Let  $f_\epsilon = p_\epsilon - 1$ , which lies in  $L_0^1$  by  $\int \mu_{\text{ref}}(dy) p_\epsilon(y) = 1$ . Define  $F : L_0^1 \times \mathbb{R} \rightarrow L_0^1$  as

$$F(f, \epsilon) = (I - P_\epsilon^*)(f + 1)$$

so that  $F(f_\epsilon, \epsilon) = 0$ . Note that the image of  $F$  is included in  $L_0^1$  by the fact that  $\int \mu_{\text{ref}}(dy) P_\epsilon(x, y) = 1$ .

Below, we derive the derivative of  $f_\epsilon$  in the Banach space  $(L_0^1, \|\cdot\|_1)$ . Notice that  $F$  is linear in  $f$ , especially the Jacobian  $D_f F(f, \epsilon) : L_0^1 \rightarrow L_0^1$  with respect to  $f$  is given by

$$D_f F(f, \epsilon) = I - P_\epsilon^*.$$

Now we claim that  $I - P_0^*$  has a bounded inverse on  $L_0^1$ . Noting that for all  $h \in L_0^1$ , Lemma 14 implies that there exists a constant  $c = c(U) \in [0, 1)$  such that

$$\|P_0^* h\|_1 \leq c \cdot \|h\|_1.$$

Since  $0 \leq c < 1$ ,  $(I - P_0^*)$  has the inverse  $(I - P_0^*)^{-1} = \sum_{t=0}^{\infty} (P_0^*)^t$  and its operator norm is bounded by  $1/(1-c)$  since for all  $h \in L_0^1$ ,

$$\|(I - P_0^*)^{-1} h\|_1 = \left\| \sum_{t=0}^{\infty} (P_0^*)^t h \right\|_1 \leq \sum_{t=0}^{\infty} \|(P_0^*)^t h\|_1 = \sum_{t=0}^{\infty} c^t \|h\|_1 = \frac{1}{1-c} \|h\|_1.$$

Thus,  $\sum_{t=0}^{\infty} (P_0^*)^t$  is the bounded inverse of  $(I - P_0^*)$  on  $L_0^1$ .

On the other hand, for the derivative with respect to  $\epsilon$ , one can show that

$$D_\epsilon F(f, \epsilon)(y) = - \int \mu_{\text{ref}}(dx) (f(x) + 1) P_\epsilon(x, y) \left( E(x, y) - \int \mu_{\text{ref}}(dy') P_\epsilon(x, y') E(x, y') \right).$$

Therefore, by the Implicit function theorem for the Banach space  $(L_0^1, \|\cdot\|_1)$ ,  $f_\epsilon$  is differentiable at  $\epsilon = 0$  in the sense that

$$\lim_{\epsilon \rightarrow 0} \left\| \frac{f_\epsilon - f_0}{\epsilon} - \dot{f}_0 \right\|_1 = 0$$

where the derivative  $\dot{f}_0$  is given by

$$\dot{f}_0 = \left( (I - P_0^*) \big|_{L_0^1} \right)^{-1} \left( y \mapsto \int \mu(dx) (f_0 + 1)(x) P_0(x, y) \left( E(x, y) - \int \mu(dy') P_0(x, y') E(x, y') \right) \right).$$

Substituting  $f_\epsilon = p_\epsilon - 1$ , we complete the proof. ■

## Appendix D. Proof for NLHF

### D.1. Geometric convergence under TV when $\|U\|_{\oplus} < 1$

**Theorem 18** *The algorithm (6) with  $\eta = \infty$  (i.e., no local regularization term) satisfies*

$$d_{\text{TV}}(\nu_t, \nu_{\text{NL}}) \leq \|U\|_{\oplus} d_{\text{TV}}(\mu_t, \mu_{\text{NL}}), \quad d_{\text{TV}}(\mu_{t+1}, \mu_{\text{NL}}) \leq \|U\|_{\oplus} d_{\text{TV}}(\nu_t, \nu_{\text{NL}}).$$

Thus, if  $\|U\|_{\oplus} < 1$ , we have  $d_{\text{TV}}(\mu_t, \mu_{\text{NL}}) \leq \|U\|_{\oplus}^{2t} \cdot d_{\text{TV}}(\mu_{\text{ref}}, \mu_{\text{NL}})$ .

Here, the convergence rate is bounded by  $\|U\|_{\oplus}$ , rather than  $\|U\|_{\infty}$ . Thus, similarly to the MCHF update in Corollary 3, this algorithm implicitly adapts to the additive structure of  $U$ .

**Proof** Let  $p_t = \frac{d\mu_{\text{NL}}^t}{d\mu_{\text{ref}}}$ ,  $q_t = \frac{d\nu_{\text{NL}}^t}{d\mu_{\text{ref}}}$ ,  $p_{\text{NL}} = \frac{d\mu_{\text{NL}}}{d\mu_{\text{ref}}}$ , and  $q_{\text{NL}} = \frac{d\nu_{\text{NL}}}{d\mu_{\text{ref}}}$  so that

$$\begin{aligned} q_t &= \operatorname{argmin}_{q \in \Delta^{\infty}} \langle q, U p_t \rangle + w(q), & q_{\text{NL}} &= \operatorname{argmin}_{q \in \Delta^{\infty}} \langle q, U p_{\text{NL}} \rangle + w(q), \\ p_{t+1} &= \operatorname{argmax}_{p \in \Delta^{\infty}} \langle q_t, U p \rangle - w(p), & p_{\text{NL}} &= \operatorname{argmax}_{p \in \Delta^{\infty}} \langle q_{\text{NL}}, U p \rangle - w(p) \end{aligned}$$

By Lemma 20-(1) and  $\|U^*\|_{\oplus} = \|U\|_{\oplus}$ , we have

$$\|q_t - q_{\text{NL}}\|_1 \leq \|U\|_{\oplus} \|p_t - p_{\text{NL}}\|_1, \quad \|p_{t+1} - p_{\text{NL}}\|_1 \leq \|U\|_{\oplus} \|q_t - q_{\text{NL}}\|_1,$$

which completes the proof. ■

### D.2. Proof of Theorem 4

**Lemma 19** *Let  $X$  be a convex subset of a Hilbert space  $(\mathcal{H}, \langle \cdot, \cdot \rangle)$  and consider the minimization problem*

$$\min_{p \in X} \ell(p) + w(p),$$

where both  $\ell$  and  $w$  are differentiable and convex. We further assume that  $\ell$  is relatively  $L$ -smooth with respect to  $w$  in the sense that

$$\forall p, q \in X, \quad D_{\ell}(p||q) \leq L D_w(p||q),$$

where  $D_h(p||q) = h(p) - h(q) - \langle \nabla h(q), p - q \rangle$  denotes the Bregman divergence associated with a convex function  $h \in \{\ell, w\}$ . Now we consider the (proximal) mirror descent iteration

$$p_{t+1} = \operatorname{argmin}_{p \in X} \left\{ \langle \nabla \ell(p_t), p \rangle + w(p) + \frac{1}{\eta} D_w(p||p_t) \right\}.$$

Then for any  $\eta \leq 1/L$ , letting  $p_* \in \operatorname{argmin} \ell(p) + w(p)$  be the minimizer, we have

$$D_w(p_*||p_t) \leq \left( \frac{1}{1 + \eta} \right)^t D_w(p_*||p_0).$$

We will prove Lemma 19 in Section D.2.1. Let us show Theorem 4 using Lemma 19. Define  $p_{\text{NL}}$  and  $q_{\text{NL}}$  as

$$p_{\text{NL}} = \frac{d\mu_{\text{NL}}}{d\mu_{\text{ref}}}, \quad q_{\text{NL}} = \frac{d\nu_{\text{NL}}}{d\mu_{\text{ref}}}.$$

Consider the density simplex

$$\Delta^\infty = \{f \in L^\infty : f(\cdot) \geq 0, \int \mu_{\text{ref}}(dy)f(y) = 1\},$$

which is a convex subset of the Hilbert space  $L^2$  with the usual inner product  $\langle f, g \rangle = \int \mu_{\text{ref}}(dx)g(x)f(x)$ , and define the linear operator  $h \mapsto Uh$  by  $Uh(x) = \int \mu_{\text{ref}}(dy)U(x, y)h(y)$ , and let  $U^*$  be its adjoint operator. Then,  $(p_{\text{NL}}, q_{\text{NL}})$  solves the following minimax optimization problem:

$$\operatorname{argmax}_{p \in \Delta^\infty} \min_{q \in \Delta^\infty} \langle q, Up \rangle + w(q) - w(p) \quad \text{where} \quad w(p) = \int \mu_{\text{ref}}(dy)p(y) \log p(y)$$

If we define the convex function  $\ell : \Delta^\infty \rightarrow \mathbb{R}$  as

$$\ell(p) = -\left(\min_{q \in \Delta^\infty} \langle q, Up \rangle + w(q)\right)$$

then  $p_{\text{NL}}$  is the solution to the following convex optimization problem:

$$p_{\text{NL}} \in \operatorname{argmin}_{p \in \Delta^\infty} \ell(p) + w(p).$$

Now we claim that  $\ell$  is relatively  $\|U\|_\oplus^2$ -smooth with respect to  $w$ .

**Lemma 20** Fix  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , and for each  $p \in \Delta^\infty$ , let  $\ell(p) = -(\min_{q \in \Delta^\infty} \langle q, Up \rangle + w(q))$  and  $q_p \in \operatorname{argmin}_{q \in \Delta^\infty} \langle q, Up \rangle + w(q)$ . Then,

1.  $\|q_p - q_{p'}\|_1 \leq \|U\|_\oplus \|p - p'\|_1$  for all  $p, p' \in \Delta^\infty$ .
2.  $\ell$  is  $\|U\|_\oplus^2$ -smooth with respect to  $w$ , i.e.,  $D_\ell(p||q) \leq \|U\|_\oplus^2 D_w(p||q)$

where  $\|U\|_\oplus = \inf_{g, f \in L^\infty} \|U - g \oplus f\|_\infty$ .

Thus, combined with Lemma 20-(2), we may apply Lemma 19 with

$$L = \|U\|_\oplus^2, \quad w(p) = \int \mu_{\text{ref}}(dy)p(y) \log p(y), \quad \ell(p) = -\left(\min_{q \in \Delta^\infty} \langle q, Up \rangle + w(q)\right)$$

and conclude that for any  $\eta \leq 1/\|U\|_\oplus^2$ , the iteration  $(q_t, p_t)$  defined by

$$\begin{aligned} q_t &= \operatorname{argmin}_q \left\{ \langle q, Up_t \rangle + w(q) \right\} \\ p_{t+1} &= \operatorname{argmin}_p \left\{ \langle \nabla \ell(p_t), p \rangle + w(p) + \frac{1}{\eta} D_w(p||p_t) \right\} \\ &= \operatorname{argmin}_p \left\{ -\langle q_t, Up \rangle + w(p) + \frac{1}{\eta} D_w(p||p_t) \right\} \end{aligned}$$

satisfies that

$$D_w(p_*||p_t) \leq \left(\frac{1}{1+\eta}\right)^t D_w(p_*||p_0).$$

Finally, letting  $\mu(dy) = p(y)d\mu_{\text{ref}}(dy)$  and  $\nu(dx) = q(x)d\mu_{\text{ref}}(dx)$ , noting that  $w(p) = D_{\text{KL}}(\mu||\mu_{\text{ref}})$  and  $D_w(p||q) = D_{\text{KL}}(\mu||\nu)$ , we obtain Theorem 4.

## D.2.1. PROOF OF LEMMA 19

Let  $G(p) = \ell(p) + w(p)$  and define the surrogate objective

$$G_t(p) = \ell(p_t) + \langle \nabla \ell(p_t), p - p_t \rangle + \eta^{-1} D_w(p \| p_t) + w(p).$$

Then  $p_{t+1} \in \operatorname{argmin}_{p \in X} G_t(p)$ . By the definition of Bregman divergence,

$$\begin{aligned} G(p) &= \ell(p) + w(p) \\ &= \ell(p_t) + \langle \nabla \ell(p_t), p - p_t \rangle + D_\ell(p \| p_t) + w(p) \quad \text{by } D_\ell(p \| p_t) = \ell(p) - \ell(p_t) - \langle \nabla \ell(p_t), p - p_t \rangle \\ &= G_t(p) + D_\ell(p \| p_t) - \eta^{-1} D_w(p \| p_t). \end{aligned}$$

Using the assumption  $D_\ell(p \| p_t) \leq L D_w(p \| p_t)$  and the non-negativity of the Bregman divergence  $D_\ell(p \| p_t) \geq 0$ , the approximation error  $G(p) - G_t(p)$  can be controlled as

$$-\eta^{-1} D_w(p \| p_t) \leq G(p) - G_t(p) \leq -(\eta^{-1} - L) D_w(p \| p_t). \quad (11)$$

Now, since  $p_{t+1} \in \operatorname{argmin} G_t(p)$  and  $G_t$  is convex and differentiable, we have

$$\langle \nabla G_t(p_{t+1}), q - p_{t+1} \rangle \geq 0 \quad \forall q \in X.$$

Equivalently, by the definition of Bregman divergence,

$$D_{G_t}(q \| p_{t+1}) \leq G_t(q) - G_t(p_{t+1}), \quad \forall q \in X.$$

Applying this with  $q = p_*$ , we get

$$\begin{aligned} (\eta^{-1} + 1) D_w(p_* \| p_{t+1}) &= D_{G_t}(p_* \| p_{t+1}) \quad \text{since } G_t(p) = (\eta^{-1} + 1)w(p) + \text{linear terms} \\ &\leq G_t(p_*) - G_t(p_{t+1}) \\ &\leq \left( F(p_*) + \eta^{-1} D_w(p_* \| p_t) \right) - \left( F(p_{t+1}) + (\eta^{-1} - L) D_w(p_{t+1} \| p_t) \right) \quad \text{by (11)} \\ &\leq \eta^{-1} D_w(p_* \| p_t) - (\eta^{-1} - L) D_w(p_{t+1} \| p_t) \quad \text{since } F(p_*) \leq F(p_{t+1}) \\ &\leq \eta^{-1} D_w(p_* \| p_t), \end{aligned}$$

where the last step uses  $\eta^{-1} \geq L$  and  $D_w(\cdot \| \cdot) \geq 0$ . Dividing both sides by  $\eta^{-1} + 1$  and iterating completes the proof.

## D.2.2. PROOF OF LEMMA 20

Fix  $g, f \in L^\infty$  and let  $U_{g,f} = U - g \oplus f$ . Then, we can rewrite  $\ell(p)$  as

$$\ell(p) = - \left( \min_{q \in \Delta^\infty} \langle q, U_{g,f} p \rangle + \langle q, g \rangle + \langle p, f \rangle + w(q) \right).$$

By the envelope theorem,  $\ell$  is differentiable with its derivative given by

$$\nabla \ell(p) = -U_{g,f}^* q_p + f, \quad (12)$$

Here,  $q_p$  is the solution of  $\min_q \langle q, U_{g,f}p \rangle + \langle q, g \rangle + \langle p, f \rangle + w(q)$ , which can be written explicitly as

$$q_p(x) = \frac{\exp(-U_{g,f}p(x))e^{-g(x)}}{\int \mu_{\text{ref}}(dx') \exp(-U_{g,f}p(x'))e^{-g(x')}}. \quad (13)$$

Then, for any  $p, p' \in \Delta^\infty$ , using  $\|Uf\|_\infty \leq \|U\|_\infty \|f\|_1$  by the Cauchy–Schwarz inequality and Lemma 16 with  $\mu_{-g}(dx) \propto e^{-g(x)}\mu_{\text{ref}}(dx)$ , noting that  $\mu_{-g}$  and  $\mu_{\text{ref}}$  are absolutely continuous with respect to each other,

$$\begin{aligned} \|q_p - q_{p'}\|_1 &\leq \| -U_{g,f}p + U_{g,f}p' \|_{L^\infty(\mu_{-g})} \quad \text{by (13) and Lemma 16 with } \mu_{-g}(dx) \propto e^{-g(x)}\mu_{\text{ref}}(dx) \\ &= \| -U_{g,f}p + U_{g,f}p' \|_\infty \quad \| \cdot \|_{L^\infty(\mu_{-g})} = \| \cdot \|_{L^\infty(\mu_{\text{ref}})} \\ &\leq \|U_{g,f}\|_\infty \|p - p'\|_1 \quad \text{Cauchy–Schwarz.} \end{aligned}$$

Since  $g, f \in L^\infty$  are arbitrary, taking  $\inf_{g,f \in L^\infty}$  on the RHS, we complete the proof of  $\|q_p - q_{p'}\|_1 \leq \|U\|_{\oplus} \|p - p'\|_1$ .

Now, using the derivative form (12),

$$\begin{aligned} \|\nabla \ell(p) - \nabla \ell(p')\|_\infty &= \| -U_{g,f}^*q_p + U_{g,f}^*q_{p'} \|_\infty \quad \text{by (12)} \\ &\leq \|U_{g,f}^*\|_\infty \|q_p - q_{p'}\|_1 \quad \text{Cauchy–Schwarz} \\ &\leq \|U_{g,f}\|_\infty \|U\|_{\oplus} \|p - p'\|_1 \quad \text{Cauchy–Schwarz.} \end{aligned}$$

Again, since  $g, f \in L^\infty$  are arbitrary, taking  $\inf_{g,f \in L^\infty}$  on the RHS, we get

$$\forall p, p' \in \Delta^\infty, \quad \|\nabla \ell(p) - \nabla \ell(p')\|_\infty \leq \|U\|_{\oplus}^2 \|p - p'\|_1. \quad (14)$$

Therefore, letting  $\phi(t) = \ell(tp + (1-t)q)$ , we have

$$\begin{aligned} D_\ell(p||q) &= \ell(p) - \ell(q) - \langle \nabla \ell(q), p - q \rangle \\ &= \int_0^1 \phi'(t) dt - \langle \nabla \ell(q), p - q \rangle \\ &= \int_0^1 \langle \nabla \ell(tp + (1-t)q) - \nabla \ell(q), p - q \rangle dt \\ &\leq \int_0^1 \|\nabla \ell(tp + (1-t)q) - \nabla \ell(q)\|_\infty \|p - q\|_1 dt \\ &\leq \int_0^1 \|U\|_{\oplus}^2 \|t(p - q)\|_1 \|p - q\|_1 dt \quad \text{by (14)} \\ &= \frac{1}{2} \|U\|_{\oplus}^2 \|p - q\|_1^2 \\ &\leq \|U\|_{\oplus}^2 D_w(p||q) \quad \text{by Pinsker's inequality} \end{aligned}$$

where the last inequality follows from Pinsker's inequality  $D_w(p||q) \geq \frac{1}{2} \|p - q\|_1^2$  for  $w(p) = \int \mu_{\text{ref}}(dy) p(y) \log p(y)$ . Therefore,  $\ell$  is relatively  $\|U\|_{\oplus}^2$ -smooth with respect to  $w$ .

### D.3. Proof of Theorem 9

Let  $(\mu_{\text{NL}}(U), \nu_{\text{NL}}(U))$  be the solution to the NLHF with utility  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  and let  $p(U) = \frac{d\mu_{\text{NL}}(U)}{d\mu_{\text{ref}}}$ . We showed in Section D.3 that  $p(U)$  is the solution to the strongly convex optimization:

$$p(U) = \operatorname{argmin}_{p \in \Delta^\infty} G_U(p) \equiv \ell_U(p) + w(p), \quad \ell_U(p) = -\left( \min_{q \in \Delta^\infty} \langle q, Up \rangle + w(q) \right)$$

where  $w(p) = \int \mu_{\text{ref}}(dz) p(z) \log p(z)$ . Since  $G_U$  is strongly convex with respect to  $w$ ,  $p(U)$  also solves

$$p(U) \in \operatorname{argmin}_{p \in \Delta^\infty} G_U(p) - D_w(p \| p(U)).$$

Then, for any  $U, \hat{U} \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , using this optimality of  $p(U)$  against  $p(\hat{U})$ , noting  $D_w(p(U) \| p(U)) = 0$ , we have

$$G_U(p(U)) \leq G_U(p(\hat{U})) - D_w(p(\hat{U}) \| p(U)).$$

Rearranging this inequality,

$$\begin{aligned} D_w(p(\hat{U}) \| p(U)) &\leq G_U(p(\hat{U})) - G_U(p(U)) \\ &= G_U(p(\hat{U})) - G_{\hat{U}}(p(\hat{U})) + G_{\hat{U}}(p(\hat{U})) - G_U(p(U)) \\ &\leq G_U(p(\hat{U})) - G_{\hat{U}}(p(\hat{U})) + G_{\hat{U}}(p(U)) - G_U(p(U)) \quad p(\hat{U}) \in \operatorname{argmin}_p G_{\hat{U}}(p) \\ &\leq 2 \sup_{p \in \Delta^\infty} |G_U(p) - G_{\hat{U}}(p)| \end{aligned} \quad (15)$$

Here, letting  $q_p \in \operatorname{argmin}_q \langle q, Up \rangle + w(q)$ ,

$$\begin{aligned} G_U(p) &= -\langle q_p, Up \rangle - w(q_p) + w(p) \\ &= -\langle q_p, (U - \hat{U})p \rangle - \langle q_p, \hat{U}p \rangle - w(q_p) + w(p) \\ &\leq -\langle q_p, (U - \hat{U})p \rangle - \min_q \left( \langle q, \hat{U}p \rangle + w(q) \right) + w(p) \\ &= -\langle q_p, (U - \hat{U})p \rangle + G_{\hat{U}}(p), \end{aligned}$$

so that

$$G_U(p) - G_{\hat{U}}(p) \leq |\langle q_p, (U - \hat{U})p \rangle| \leq \|q_p\|_1 \|p\|_1 \|U - \hat{U}\|_\infty = \|U - \hat{U}\|_\infty$$

for all  $p \in \Delta^\infty$ . By symmetry, the above inequality also holds with  $U$  and  $\hat{U}$  swapped. Thus, we get

$$\sup_{p \in \Delta^\infty} |G_U(p) - G_{\hat{U}}(p)| \leq \|U - \hat{U}\|_\infty.$$

Combined with (15), we get

$$D_w(p(\hat{U}) \| p(U)) \leq 2\|U - \hat{U}\|_\infty.$$

Substituting  $p(U) = d\mu_{\text{NL}}(U)/d\mu_{\text{ref}}$  and  $w(p) = \int \mu_{\text{ref}}(dz) p(z) \log p(z)$ , we have  $D_w(p(\hat{U}) \| p(U)) = D_{\text{KL}}(\mu_{\text{NL}}(\hat{U}) \| \mu_{\text{NL}}(U))$  and hence complete the proof.

**D.4. Proof of Theorem 10**

Let  $(\mu_{\text{NL}}(U), \nu_{\text{NL}}(U))$  be the solution to the NLHF with utility  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  and let  $p(U) = \frac{d\mu_{\text{NL}}(U)}{d\mu_{\text{ref}}}$  and  $q(U) = \frac{d\nu_{\text{NL}}(U)}{d\mu_{\text{ref}}}$ . Note in passing that for any  $U, \hat{U}$ ,

$$q(U) = \operatorname{argmin}_{q \in \Delta^\infty} \langle q, Up(U) \rangle + w(q), \quad q(\hat{U}) = \operatorname{argmin}_{q \in \Delta^\infty} \langle q, \hat{U}p(\hat{U}) \rangle + w(q).$$

Now we define an intermediate  $\tilde{q}$  as

$$\tilde{q} = \operatorname{argmin}_{q \in \Delta^\infty} \langle q, \hat{U}p(U) \rangle + w(q)$$

and decompose  $\|q(U) - q(\hat{U})\|_1$  as

$$\|q(U) - q(\hat{U})\|_1 \leq \|q(U) - \tilde{q}\|_1 + \|\tilde{q} - q(\hat{U})\|_1.$$

For the first term, we may apply Lemma 16 with  $g = Up(U)$  and  $\tilde{g} = \hat{U}p(U)$  and obtain

$$\|q(U) - \tilde{q}\|_1 \leq \|Up(U) - \hat{U}p(U)\|_\infty \leq \|U - \hat{U}\|_\infty \|p(U)\|_1 = \|U - \hat{U}\|_\infty.$$

For the second term, we apply Lemma 20-(1) and obtain

$$\|\tilde{q} - q(\hat{U})\|_1 \leq \|\hat{U}\|_\oplus \|p(U) - p(\hat{U})\|_1.$$

Combining the above displays, we get

$$\|q(U) - q(\hat{U})\|_1 \leq \|U - \hat{U}\|_\infty + \|\hat{U}\|_\oplus \|p(U) - p(\hat{U})\|_1.$$

Here, swapping  $U$  and  $\hat{U}$ , the above inequality also holds with  $\|\hat{U}\|_\oplus$  replaced by  $\|U\|_\oplus$ . Thus, we get

$$\|q(U) - q(\hat{U})\|_1 \leq \|U - \hat{U}\|_\infty + (\|U\|_\oplus \wedge \|\hat{U}\|_\oplus) \|p(U) - p(\hat{U})\|_1.$$

Applying the same argument to the pair  $(p(U), p(\hat{U}))$  which satisfy

$$p(U) = \operatorname{argmin}_{p \in \Delta^\infty} \langle p, -U^*q(U) \rangle + w(p), \quad p(\hat{U}) = \operatorname{argmin}_{p \in \Delta^\infty} \langle p, -\hat{U}^*q(\hat{U}) \rangle + w(p),$$

noting  $\| -U^* \|_\oplus = \|U\|_\oplus$ , it also holds that

$$\|p(U) - p(\hat{U})\|_1 \leq \|U - \hat{U}\|_\infty + (\|U\|_\oplus \wedge \|\hat{U}\|_\oplus) \|q(U) - q(\hat{U})\|_1.$$

Combining the above two displays,

$$\|p(U) - p(\hat{U})\|_1 \leq \|U - \hat{U}\|_\infty + (\|U\|_\oplus \wedge \|\hat{U}\|_\oplus) \left( \|U - \hat{U}\|_\infty + (\|U\|_\oplus \wedge \|\hat{U}\|_\oplus) \|p(U) - p(\hat{U})\|_1 \right).$$

Rearranging this,

$$\left( 1 - (\|U\|_\oplus \wedge \|\hat{U}\|_\oplus)^2 \right) \|p(U) - p(\hat{U})\|_1 \leq (1 + \|U\|_\oplus \wedge \|\hat{U}\|_\oplus) \|U - \hat{U}\|_\infty.$$

By the assumption  $\|U\|_\oplus \wedge \|\hat{U}\|_\oplus < 1$ , we may divide by  $(1 - (\|U\|_\oplus \wedge \|\hat{U}\|_\oplus)^2) = (1 - (\|U\|_\oplus \wedge \|\hat{U}\|_\oplus))(1 + (\|U\|_\oplus \wedge \|\hat{U}\|_\oplus))$  and complete the proof.

### D.5. Proof of Theorem 11

We denote  $\|\cdot\|_p = \|\cdot\|_{L^p(\mu_{\text{ref}})}$  and  $\|U\|_\infty = \|U\|_{L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})}$  for simplicity. Now for each  $t \in [0, 1]$ , let  $\mu_t$  be the solution to the following minimax problem

$$\mu_t \in \operatorname{argmax}_\mu \min_\nu \left( \mathbb{E}_{y \sim \mu, x \sim \nu} [tU(x, y) + (1-t)\hat{U}(x, y)] - D_{\text{KL}}(\mu|\mu_{\text{ref}}) + D_{\text{KL}}(\nu|\mu_{\text{ref}}) \right).$$

Our goal is to bound  $d_{\text{TV}}(\mu_1, \mu_0)$ . Now we claim that the density  $p_t(y) \equiv \frac{\mu_t(dy)}{\mu_{\text{ref}}(dy)}$  is differentiable in  $t$  in the sense that there exists  $\dot{p}_t \in L^1$  such that

$$\lim_{\epsilon \rightarrow 0} \left\| \frac{p_{t+\epsilon} - p_t}{\epsilon} - \dot{p}_t \right\|_1 = 0,$$

and there exists an absolute constant  $C$  such that the L1 norm of the derivative is bounded as

$$\|\dot{p}_t\|_1 = \int \mu_{\text{ref}}(dy) |\dot{p}_t(y)| \leq C \cdot (1 + \|tU + (1-t)\hat{U}\|_\oplus^6) \cdot \|U - \hat{U}\|_\infty \quad (16)$$

for all  $t \in [0, 1]$ . Then, by Fubini's theorem,

$$\begin{aligned} d_{\text{TV}}(\mu_1, \mu_0) &= \frac{1}{2} \int \mu_{\text{ref}}(dy) |p_1(y) - p_0(y)| \\ &= \frac{1}{2} \int \mu_{\text{ref}}(dy) \left| \int_0^1 \dot{p}_t(y) dt \right| \\ &\leq \frac{1}{2} \int \mu_{\text{ref}}(dy) \int_0^1 |\dot{p}_t(y)| dt \\ &\leq \frac{1}{2} \int_0^1 dt \int \mu_{\text{ref}}(dy) |\dot{p}_t(y)| \\ &\leq \frac{C}{2} \|U - \hat{U}\|_\infty \int_0^1 dt (1 + \|tU + (1-t)\hat{U}\|_\oplus^6) \quad \text{by (16)} \end{aligned}$$

Here, by the triangle inequality for the semi-norm  $\|\cdot\|_\oplus$ , one can show  $\int_0^1 dt \|tU + (1-t)\hat{U}\|_\oplus^6 \leq C'(\|U\|_\oplus^6 + \|\hat{U}\|_\oplus^6)$  for an absolute constant  $C'$ . Then we get

$$d_{\text{TV}}(p_{\text{NL}}(U), p_{\text{NL}}(\hat{U})) \leq 2^{-1} C C' \cdot \|U - \hat{U}\|_\infty \cdot (1 + \|\hat{U}\|_\oplus^6 + \|U\|_\oplus^6).$$

Since  $2^{-1} C C'$  is an absolute constant independent of all other quantities, the proof of Theorem 11 is complete. Thus, the rest of the goal is to show (16). Notice that it suffices to prove the following lemma:

**Lemma 21** Fix  $U, E \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ . For any  $\epsilon \in \mathbb{R}$ , let  $\mu_\epsilon$  be the solution to the following minimax problem

$$\mu_\epsilon \in \operatorname{argmax}_\mu \min_\nu \left( \mathbb{E}_{y \sim \mu, x \sim \nu} [U(x, y) + \epsilon E(x, y)] - D_{\text{KL}}(\mu|\mu_{\text{ref}}) + D_{\text{KL}}(\nu|\mu_{\text{ref}}) \right). \quad (17)$$

Then, the density  $p_\epsilon(y) = \frac{p_\epsilon}{\mu_{\text{ref}}}(dy)$  is differentiable at  $\epsilon = 0$  in the sense of

$$\lim_{\epsilon \rightarrow 0} \left\| \frac{p_\epsilon - p_0}{\epsilon} - \dot{p}_0 \right\|_1 = 0,$$

where the norm of its derivative bounded as

$$\int |\dot{p}_0(y)| \mu_{\text{ref}}(dy) \leq C \cdot (1 + \|U\|_\infty^6) \cdot \|E\|_\infty$$

where  $C$  is an absolute constant.

Indeed, if the above lemma holds, then replacing  $U$  by  $tU + (1-t)U$  and  $E$  by  $U - \hat{U}$ , we obtain (16). Thus, it suffices to show Lemma 21.

#### D.5.1. PROOF OF LEMMA 21

**Notation** We first fix notation. Let  $L^2 = L^2(\mu_{\text{ref}})$  be the Hilbert space with usual inner product  $\langle f, g \rangle = \int f(z)g(z)\mu_{\text{ref}}(dz)$ . For any  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , let us write  $\|U\|_\infty = \|U\|_{L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})}$ . Now we define the linear operators  $U$  as

$$U : L^2 \rightarrow L^\infty, \quad f \mapsto Uf(\cdot) = \int U(\cdot, y)f(y)\mu_{\text{ref}}(dy)$$

and let  $U^*$  be its adjoint operator such that  $\langle g, Uf \rangle = \langle U^*g, f \rangle$  for all  $f, g \in L^2$ . Note that  $U$  is a bounded operator with operator norm bounded as

$$\|U\|_{L^2 \rightarrow L^\infty} = \sup_{\|f\|_2=1} \|Uf\|_\infty \leq \sup_{\|f\|_2=1} \|U\|_\infty \|f\|_1 \leq \sup_{\|f\|_2=1} \|U\|_\infty \|f\|_2 = \|U\|_\infty.$$

Next, let  $\Delta^\infty \subset L^\infty$  be the density simplex defined as

$$\Delta^\infty = \{p \in L^\infty : p(\cdot) \geq 0, \int \mu_{\text{ref}}(dz)p(z) = 1\}.$$

Notice that for any  $p \in \Delta^\infty$ , it induces a tilted probability measure  $\mu(dy) = p(y)\mu_{\text{ref}}(dy)$ . We define the operator  $\mathcal{G}$  taking value in  $\Delta^\infty$  as

$$\mathcal{G} : L^\infty \rightarrow \Delta^\infty, \quad f \mapsto \mathcal{G}(f)(x) = \frac{\exp(f(x))}{\int \mu_{\text{ref}}(dx') \exp(f(x'))}$$

For each  $p \in \Delta^\infty$ , define the linear operator  $J_p : L^2 \rightarrow L^2_0$  as  $J_p = \text{Diag}(p) - pp^*$  or more precisely

$$J_p : L^2 \rightarrow L^2_0, \quad f \mapsto J_p f(z) = p(z)f(z) - p(z) \int \mu_{\text{ref}}(dz')p(z')f(z')$$

where  $L^2_0 = \{h \in L^2 : \int \mu_{\text{ref}}(dz)h(z) = 0\} \subset L^2$ . Notice that  $J_p$  is bounded, self-adjoint, and positive semidefinite operator on  $L^2$ . Indeed, for any  $f, g \in L^2$ , letting  $\mu_p(dz) = p(z)\mu_{\text{ref}}(dz)$  be the tilted probability measure by the density  $p \in \Delta^\infty$ , one can show

$$\langle f, J_p g \rangle = \text{Cov}_{Z \sim \mu_p}(f(Z), g(Z))$$

from which we conclude that  $J_p$  is self-adjoint and positive semidefinite operator. Moreover, the operator norm  $\|J_p\|_{L^2 \rightarrow L^2}$  is bounded by  $\|p\|_\infty$  because the Cauchy–Schwarz inequality yields

$$|\langle f, J_p g \rangle| \leq \sqrt{\text{Var}_{Z \sim \mu_p} f(Z)} \cdot \sqrt{\text{Var}_{Z \sim \mu_p} g(Z)} \leq \|f\|_{L^2(\mu_p)} \cdot \|g\|_{L^2(\mu_p)} \leq \|p\|_\infty \cdot \|f\|_{L^2(\mu_{\text{ref}})} \cdot \|g\|_{L^2(\mu_{\text{ref}})}$$

where the last inequality follows from  $\|f\|_{L^2(\mu_p)} \leq \|\frac{d\mu_p}{d\mu_{\text{ref}}}\|_{L^\infty(\mu_{\text{ref}})}^{1/2} \cdot \|f\|_{L^2(\mu_{\text{ref}})} = \|p\|_\infty^{1/2} \|f\|_\infty$  for any  $f \in L^2$ . Therefore, by the spectral theorem, there exists a bounded self-adjoint positive semidefinite operator  $T$  such that  $T \cdot T = J_p$ . We denote such  $T$  by  $J_p^{1/2}$ .

Finally, the two operators  $\mathcal{G}, J_p$  are related via Fréchet derivative:  $D_f \mathcal{G} = J_{\mathcal{G}(f)}$ . That is, for any  $f, \delta f \in L^\infty$ , by the quotient rule, it holds that

$$\mathcal{G}(f + \delta f) = \mathcal{G}(f) + J_{\mathcal{G}(f)} \delta f + o(\|\delta f\|_\infty).$$

We now claim that the density of NLHF,  $p(U) = d\mu_{\text{NL}}(U)/d\mu_{\text{ref}}$ , is Fréchet differentiable.

**Theorem 22** *Fix  $U, E \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ . For any  $\epsilon \in \mathbb{R}$ , let  $(\mu_\epsilon, \nu_\epsilon)$  be the solution to the following minimax problem*

$$\max_{\mu} \min_{\nu} \left( \mathbb{E}_{y \sim \mu, x \sim \nu} [U(x, y) + \epsilon E(x, y)] - D_{\text{KL}}(\mu | \mu_{\text{ref}}) + D_{\text{KL}}(\nu | \mu_{\text{ref}}) \right).$$

Let  $p_\epsilon(y) = \frac{\mu_\epsilon(dy)}{\mu_{\text{ref}}(dy)}$  and  $q_\epsilon(x) = \frac{\nu_\epsilon(dx)}{\mu_{\text{ref}}(dx)}$  be their densities. The map  $\epsilon \mapsto (p_\epsilon, q_\epsilon)$  is differentiable at  $\epsilon = 0$  in the sense that

$$\lim_{\epsilon \rightarrow 0} \|\epsilon^{-1}(p_\epsilon - p_0) - \dot{p}_0\|_2 = \lim_{\epsilon \rightarrow 0} \|\epsilon^{-1}(q_\epsilon - q_0) - \dot{q}_0\|_2 = 0$$

where the derivatives are given by

$$\begin{aligned} \dot{p}_0 &= (I + J_{p_0} U^* J_{q_0} U)^{-1} (J_{p_0} E^* q_0 - J_{p_0} U^* J_{q_0} E p_0) \\ \dot{q}_0 &= -(I + J_{q_0} U J_{p_0} U^*)^{-1} (J_{q_0} E p_0 + J_{q_0} U J_{p_0} E^* q_0) \end{aligned}$$

where  $I + J_{p_0} U^* J_{q_0} U$  and  $I + J_{q_0} U J_{p_0} U^*$  have bounded inverses on  $L_0^2 = \{h \in L^2 : \int \mu_{\text{ref}}(dz) h(z) = 0\} \subset L^2$ .

**Proof** Let  $x_\epsilon = p_\epsilon - 1$  and  $y_\epsilon = q_\epsilon - 1$  so that  $x_\epsilon, y_\epsilon \in L_0^2$ . We then define  $F : \mathbb{R} \times L_0^2 \times L_0^2 \rightarrow L_0^2 \times L_0^2$  as the residual of the fixed point equation:

$$F(\epsilon, x, y) = \begin{pmatrix} F_1(\epsilon, x, y) \\ F_2(\epsilon, x, y) \end{pmatrix} \equiv \begin{pmatrix} (x+1) - \mathcal{G}((U + \epsilon E)^*(y+1)) \\ (y+1) - \mathcal{G}(-(U + \epsilon E)(x+1)) \end{pmatrix}$$

so that  $F(\epsilon, x_\epsilon, y_\epsilon) = 0$  for all  $\epsilon$ . Our goal is to show the differentiability of  $x_\epsilon$  (as a map from  $\mathbb{R}$  to  $L^2$ ) via the Implicit function theorem on the subset  $L_0^2$  of Hilbert spaces  $L^2$ . First, note that  $F$  indeed takes values in  $L_0^2 \times L_0^2$  since  $\mathcal{G}$  maps from  $L^\infty$  to the density simplex  $\Delta^\infty$ . Thus  $F : \mathbb{R} \times L_0^2 \times L_0^2 \rightarrow L_0^2 \times L_0^2$  is well-defined.

Let  $U_\epsilon = U + \epsilon E$ . We next compute the derivative of  $F$  with respect to  $(x, y)$ . For perturbations  $(\delta x, \delta y) \in L_0^2 \times L_0^2$ , using

$$\mathcal{G}(f + \delta f) = \mathcal{G}(f) + J_{\mathcal{G}(f)} \delta f + o(\|\delta f\|_\infty)$$

and the bound  $\|U_\epsilon \delta x\|_\infty \leq \|U_\epsilon\|_\infty \|\delta x\|_1 \leq \|U_\epsilon\|_\infty \|\delta x\|_2$ , we obtain

$$\begin{aligned} F_1(\epsilon, x + \delta x, y + \delta y) &= F_1(\epsilon, x, y) + \delta x - J_{\mathcal{G}(U_\epsilon^*(y+1))} U_\epsilon^* \delta y + o(\|\delta x\|_2 + \|\delta y\|_2), \\ F_2(\epsilon, x + \delta x, y + \delta y) &= F_2(\epsilon, x, y) + \delta y + J_{\mathcal{G}(-U_\epsilon(x+1))} U_\epsilon \delta x + o(\|\delta x\|_2 + \|\delta y\|_2). \end{aligned}$$

In particular, at the solution  $(x_\epsilon, y_\epsilon)$ , where  $p_\epsilon = x_\epsilon + 1$  and  $q_\epsilon = y_\epsilon + 1$ , we have  $p_\epsilon = \mathcal{G}(U_\epsilon^* q_\epsilon)$  and  $q_\epsilon = \mathcal{G}(-U_\epsilon p_\epsilon)$ , and hence

$$D_{x,y}F(\epsilon, x_\epsilon, y_\epsilon) = \begin{bmatrix} I & -J_{p_\epsilon} U_\epsilon^* \\ J_{q_\epsilon} U_\epsilon & I \end{bmatrix}.$$

Here the derivative is understood as a bounded linear operator from  $L_0^2 \times L_0^2$  to  $L_0^2 \times L_0^2$  since  $J_{p_0}, J_{q_0}$  map  $L^2$  to  $L_0^2$ .

We now show that  $D_{x,y}F(0, x_0, y_0)$  is an isomorphism on  $L_0^2 \times L_0^2$ . Let  $(r_1, r_2) \in L_0^2 \times L_0^2$ . Solving

$$\begin{bmatrix} I & -J_{p_0} U^* \\ J_{q_0} U & I \end{bmatrix} \begin{pmatrix} \delta x \\ \delta y \end{pmatrix} = \begin{pmatrix} r_1 \\ r_2 \end{pmatrix}$$

gives

$$\delta y = r_2 - J_{q_0} U \delta x.$$

Substituting this into the first equation yields the Schur complement equation

$$(I + J_{p_0} U^* J_{q_0} U) \delta x = r_1 + J_{p_0} U^* r_2.$$

Thus it is enough to show that  $(I + J_{p_0} U^* J_{q_0} U)$  has a bounded inverse on  $L_0^2$ . Indeed, once  $\delta x$  is obtained from the Schur complement equation,  $\delta y$  is given by  $\delta y = r_2 - J_{q_0} U \delta x$ , and the boundedness of the inverse follows from the boundedness of  $(I + J_{p_0} U^* J_{q_0} U)^{-1}$ ,  $J_{p_0} U^*$ , and  $J_{q_0} U$ .

On  $L_0^2$ , write

$$I + J_{p_0} U^* J_{q_0} U = I + AB, \quad A = J_{p_0}^{1/2}, \quad B = J_{p_0}^{1/2} U^* J_{q_0} U.$$

The operators  $A$  and  $B$  are bounded on  $L_0^2$ . Moreover,

$$BA = J_{p_0}^{1/2} U^* J_{q_0} U J_{p_0}^{1/2}$$

is self-adjoint and positive semidefinite on  $L_0^2$  since  $J_{q_0}$  and  $J_{p_0}$  are self-adjoint positive semidefinite. Hence

$$\langle h, (I + BA)h \rangle_{L^2} = \|h\|_2^2 + \langle h, BA h \rangle_{L^2} \geq \|h\|_2^2.$$

By the Lax–Milgram theorem,  $I + BA$  has a bounded inverse on  $L_0^2$ , with

$$\|(I + BA)^{-1}\|_{L_0^2 \rightarrow L_0^2} \leq 1.$$

Therefore  $I + AB$  is also invertible on  $L_0^2$ , with inverse

$$(I + AB)^{-1} = I - A(I + BA)^{-1}B.$$

Consequently,  $I + J_{p_0} U^* J_{q_0} U$  has a bounded inverse on  $L_0^2$ , and hence

$$D_{x,y}F(0, x_0, y_0) = \begin{bmatrix} I & -J_{p_0} U^* \\ J_{q_0} U & I \end{bmatrix}$$

has a bounded inverse on  $L_0^2 \times L_0^2$ .

We may therefore apply the Implicit Function Theorem on the Hilbert space  $L_0^2 \times L_0^2$ . Since  $F(\epsilon, x_\epsilon, y_\epsilon) = 0$ , the map  $\epsilon \mapsto (x_\epsilon, y_\epsilon)$  is differentiable at  $\epsilon = 0$  as a map into  $L_0^2 \times L_0^2$ . Since  $p_\epsilon = 1 + x_\epsilon$  and  $q_\epsilon = 1 + y_\epsilon$ , this is equivalent to the differentiability of  $\epsilon \mapsto (p_\epsilon, q_\epsilon)$  in  $L^2$ .

It remains to compute the derivative. Differentiating  $F(\epsilon, x_\epsilon, y_\epsilon) = 0$  at  $\epsilon = 0$ , we obtain

$$D_{x,y}F(0, x_0, y_0) \begin{pmatrix} \dot{x}_0 \\ \dot{y}_0 \end{pmatrix} + \partial_\epsilon F(0, x_0, y_0) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

Since  $\dot{p}_0 = \dot{x}_0$  and  $\dot{q}_0 = \dot{y}_0$ , it is enough to compute the system with  $(\dot{x}_0, \dot{y}_0)$  replaced by  $(\dot{p}_0, \dot{q}_0)$ . The partial derivative with respect to  $\epsilon$  is

$$\partial_\epsilon F_1(0, x_0, y_0) = -J_{p_0} E^* q_0, \quad \partial_\epsilon F_2(0, x_0, y_0) = J_{q_0} E p_0.$$

Therefore

$$\begin{bmatrix} I & -J_{p_0} U^* \\ J_{q_0} U & I \end{bmatrix} \begin{pmatrix} \dot{p}_0 \\ \dot{q}_0 \end{pmatrix} = \begin{pmatrix} J_{p_0} E^* q_0 \\ -J_{q_0} E p_0 \end{pmatrix}.$$

Equivalently,

$$\begin{aligned} \dot{p}_0 - J_{p_0} U^* \dot{q}_0 &= J_{p_0} E^* q_0, \\ J_{q_0} U \dot{p}_0 + \dot{q}_0 &= -J_{q_0} E p_0. \end{aligned}$$

Substituting  $\dot{q}_0 = -J_{q_0} U \dot{p}_0 - J_{q_0} E p_0$  into the first equation gives

$$(I + J_{p_0} U^* J_{q_0} U) \dot{p}_0 = J_{p_0} E^* q_0 - J_{p_0} U^* J_{q_0} E p_0.$$

Hence

$$\dot{p}_0 = (I + J_{p_0} U^* J_{q_0} U)^{-1} (J_{p_0} E^* q_0 - J_{p_0} U^* J_{q_0} E p_0).$$

Similarly, eliminating  $\dot{p}_0$  instead gives

$$\dot{q}_0 = -(I + J_{q_0} U J_{p_0} U^*)^{-1} (J_{q_0} E p_0 + J_{q_0} U J_{p_0} E^* q_0).$$

Here, the bounded invertibility of  $I + J_{q_0} U J_{p_0} U^*$  on  $L_0^2$  follows from the same argument as  $I + J_{p_0} U^* J_{q_0} U$  above. This completes the proof.  $\blacksquare$

We have shown that the derivative is given by

$$\dot{p}_0 = (I + J_{p_0} U^* J_{q_0} U)^{-1} (J_{p_0} E^* q_0 - J_{p_0} U^* J_{q_0} E p_0)$$

in the sense of  $\lim_{\epsilon \rightarrow 0} \|\epsilon^{-1}(p_\epsilon - p_0) - \dot{p}_0\|_2 = 0$ . The rest of our goal is to bound the  $L_1$  norm of  $\dot{p}_0$ .

**Lemma 23** *For any  $U \in L^\infty(\mu_{\text{ref}} \times \mu_{\text{ref}})$ ,  $p \in \Delta^\infty$ , and  $f \in L^2$ , it holds that*

$$\forall f \in L^2, \quad \|J_p U f\|_1 \leq \|U\|_\infty \|f\|_1$$

**Proof** Letting  $\mu(dz) = p(z)\mu_{\text{ref}}(dz)$  be the tilted measure induced by  $p \in \Delta^\infty$ , for any  $f \in L^2$ ,

$$\begin{aligned} \|J_p U f\|_1 &= \int \mu_{\text{ref}}(dz) \left| p(z) U f(z) - p(z) \int \mu_{\text{ref}}(dz') p(z') U f(z') \right| \\ &= \int \mu(dz) \left| U f(z) - \int \mu(dz') U f(z') \right| \\ &\leq \|U f\|_{L^\infty(\mu)} && \mathbb{E}[|X - \mathbb{E}[X]|] \leq \|X\|_\infty \\ &\leq \|U f\|_{L^\infty(\mu_{\text{ref}})} && \mu \ll \mu_{\text{ref}} \\ &\leq \|U\|_\infty \|f\|_1, \end{aligned}$$

where the last inequality follows from the Cauchy–Schwarz inequality.  $\blacksquare$

**Lemma 24** For any  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ ,  $g \in L^2$ , and  $p, q \in \Delta^\infty$ ,

$$\|(I + J_p U^* J_q U)^{-1} g\|_1 \leq \left( \frac{21}{8} (1 + \|U\|_\infty^2)^2 \|U\|_\infty + 1 + \frac{1}{2} \|U\|_\infty \right) \|g\|_1.$$

**Proof** See Section D.5.2.  $\blacksquare$

Let us finish the proof of Lemma 21. We fix  $f, g \in L^\infty$  and let  $U_{g,f} = U - g \oplus f$ . For any  $h \in L_0^2$ , using  $J_{q_0} \mathbf{1} = 0$  and  $\langle \mathbf{1}, h \rangle = 0$ , we have

$$J_{q_0} U h = J_{q_0} (U - g \oplus f) h + J_{q_0} \mathbf{1} \langle f, h \rangle + J_{q_0} g \langle \mathbf{1}, h \rangle = J_{q_0} U_{g,f} h$$

and hence  $J_{q_0} U$  operates on  $L_0^2$  as

$$J_{q_0} U|_{L_0^2} = J_{q_0} U_{g,f}$$

By the same argument, we have

$$J_{p_0} U^*|_{L_0^2} = J_{p_0} U_{g,f}^*$$

Thus, we can rewrite  $\dot{p}_0$  as

$$\begin{aligned} \dot{p}_0 &= (I + J_{p_0} U^* J_{q_0} U)^{-1} (J_{p_0} E^* q_0 - J_{p_0} U^* J_{q_0} E p_0) \\ &= (I + J_{p_0} U_{g,f}^* J_{q_0} U_{g,f})^{-1} (J_{p_0} E^* q_0 - J_{p_0} U_{g,f}^* J_{q_0} E p_0). \end{aligned}$$

From Lemma 20, noting  $\|p_0\|_1 = \|q_0\|_1 = 1$ , we get

$$\begin{aligned} \|J_{p_0} E^* q_0 - J_{p_0} U_{g,f}^* J_{q_0} E p_0\|_1 &\leq \|J_{p_0} E^* q_0\|_1 + \|J_{p_0} U_{g,f}^* J_{q_0} E p_0\|_1 \\ &\leq \|E\|_\infty \|q_0\|_1 + \|U_{g,f}\|_\infty \|E\|_\infty \|p_0\|_1 \\ &= (1 + \|U_{g,f}\|_\infty) \|E\|_\infty. \end{aligned}$$

Combined with Lemma 24, we have

$$\|\dot{p}_0\|_1 \leq \left( \frac{21}{8} (1 + \|U_{g,f}\|_\infty^2)^2 \|U_{g,f}\|_\infty + 1 + \frac{1}{2} \|U_{g,f}\|_\infty \right) (1 + \|U_{g,f}\|_\infty) \|E\|_\infty.$$

Here, expanding the RHS, we can find an absolute constant  $C$  such that

$$\|\dot{p}_0\|_1 \leq C (1 + \|U_{g,f}\|_\infty^6) \|E\|_\infty.$$

Since  $g, f \in L^\infty$  are arbitrary, taking  $\inf_{g,f}$  on the RHS, with  $\|U\|_\oplus = \inf_{g,f} \|U - g \oplus f\|_\infty$ , we complete the proof of Lemma 21.

## D.5.2. PROOF OF LEMMA 24

Let  $f = (I + J_p U^* J_q U)^{-1} g$ . We aim to bound  $\|f\|_1$  by  $\|g\|_1$ . If we define  $h = J_q U f$ , the triple  $(f, h, g)$  satisfy

$$\begin{aligned} f + J_p U^* h &= g \\ -J_q U f + h &= 0. \end{aligned}$$

Now, for any  $p \in \Delta^\infty$ , we rewrite the operator  $J_p$  as  $J_p = D_{\sqrt{p}} \cdot D_{\sqrt{p}} - pp^*$  where  $D_{\sqrt{p}}$  and  $\nu\nu^*$  are bounded linear operators on  $L^2$  defined as

$$D_{\sqrt{p}} f(z) = \sqrt{p(z)} f(z), \quad pp^* f(z) = p(z) \cdot \langle p, f \rangle.$$

Substituting this into the previous system, letting  $u_f = D_{\sqrt{q}} U f$ ,  $u_h = D_{\sqrt{p}} U^* h$ , we are left with

$$\begin{aligned} f + D_{\sqrt{p}} u_h - p \cdot c_h &= g \\ -D_{\sqrt{q}} u_f + q \cdot c_f + h &= 0 \end{aligned} \tag{18}$$

where we defined the scalars  $c_h, c_f$  as

$$\begin{aligned} c_h &= \langle p, U^* h \rangle = \langle \sqrt{p}, D_{\sqrt{p}} U^* h \rangle = \langle \sqrt{p}, u_h \rangle, \\ c_f &= \langle q, U f \rangle = \langle \sqrt{q}, D_{\sqrt{q}} U f \rangle = \langle \sqrt{q}, u_f \rangle \end{aligned}$$

Here, from the first equation, taking L1 norm  $\|\cdot\|_1$  and using triangle inequality, we see

$$\|f\|_1 \leq \|D_{\sqrt{p}} u_h\|_1 + c_h \|p\|_1 + \|g\|_1 \leq \|u_h\|_2 + |c_h| + \|g\|_1 \tag{19}$$

where the second inequality follows from  $\|p\|_1 = 1$  and the Cauchy–Schwarz inequality applied to  $\|D_{\sqrt{p}} u_h\|_1 = \int \mu_{\text{ref}}(dz) |\sqrt{p}(x) u_h(x)| \leq \|p\|_1^{1/2} \|u_h\|_2$ . From this result, we see that it suffices to bound the  $\|u_h\|_2$  and  $|c_h| = |\langle \sqrt{p}, u_h \rangle|$ .

Now, multiplying the first equation in (18) by  $D_{\sqrt{q}} U$  and the second equation by  $D_{\sqrt{p}} U^*$ , letting

$$A = D_{\sqrt{q}} U D_{\sqrt{p}},$$

noting  $D_{\sqrt{p}} \sqrt{p} = p$  and  $D_{\sqrt{q}} \sqrt{q} = q$ , we have

$$\begin{aligned} u_f + A u_h - c_h A \sqrt{p} &= D_{\sqrt{q}} U g \\ -A^* u_f + u_h + c_f A^* \sqrt{q} &= 0. \end{aligned} \tag{20}$$

Here  $A$  is a bounded linear operator on  $L^2$  with its operator norm bounded by  $\|U\|_\infty$ . Indeed, using the Cauchy–Schwarz inequality and  $p, q \in \Delta^\infty$ , for any  $f \in L^\infty$  with  $\|f\|_2 \leq 1$ ,

$$\begin{aligned} \|A f\|_2^2 &= \|D_{\sqrt{q}} U D_{\sqrt{p}} f\|_2^2 \\ &= \int \mu_{\text{ref}}(dx) (\sqrt{q(x)} \int \mu_{\text{ref}}(dy) U(x, y) f(y) \sqrt{p(y)})^2 \\ &\leq \|q\|_1 \cdot \text{esssup}_x \left( \int \mu_{\text{ref}}(dy) U(x, y) f(y) \sqrt{p(y)} \right)^2 \\ &\leq \|q\|_1 \cdot \|U\|_\infty^2 \left( \int \mu_{\text{ref}}(dy) |f(y)| \sqrt{p(y)} \right)^2 \\ &\leq \|q\|_1 \cdot \|U\|_\infty^2 \cdot \|f\|_2^2 \cdot \|p\|_1 \\ &= \|U\|_\infty^2 \end{aligned} \tag{by } \|q\|_1 = \|p\|_1 = 1 \text{ and } \|f\|_2 \leq 1$$

Now we solve the system for  $(u_f, u_h)$  and  $(c_f, c_h)$ . Substituting the second equation  $u_h = A^*u_f - c_f A^* \sqrt{q}$  in (20) to the first equation, rearranging it, we get

$$(I + AA^*)u_f = c_f AA^* \sqrt{q} + c_h A \sqrt{p} + D_{\sqrt{q}} U g.$$

Now we claim that  $(I + AA^*)$  has a bounded inverse on  $L^2$ . Indeed,  $(I + A^*A)$  is coercive since  $AA^*$  is a self adjoint positive semidefinite linear operator on  $L^2$ . Moreover, the operator norm is bounded as

$$\|AA^*\|_{L^2 \rightarrow L^2} \leq \|A\|_{L^2 \rightarrow L^2}^2 \leq \|U\|_{\infty}^2.$$

Thus, by the Lax-Milgram theorem,  $(I + AA^*)$  has a bounded inverse, and letting

$$T = (I + AA^*)^{-1},$$

the spectrum of  $T$  is controlled as

$$\forall f \in L^2, \quad (1 + \|U\|_{\infty}^2)^{-1} \|f\|_2^2 \leq \langle f, Tf \rangle \leq \|f\|_2^2. \quad (21)$$

Then,  $u_f$  can be solved as

$$u_f = T(c_f AA^* \sqrt{q} + c_h A \sqrt{p} + D_{\sqrt{p}} U g) = c_f (I - T) \sqrt{q} + c_h T A \sqrt{p} + T D_{\sqrt{q}} U g$$

where we used  $I - T = T A^*$  and  $T A A^* = I - T$ , which follow from the definition  $T = (I + AA^*)^{-1}$ .

By the same argument,  $(I + A^*A)$  has a bounded inverse on  $L^2$ , and letting

$$S = (I + A^*A)^{-1},$$

the spectrum of  $S$  is controlled as

$$\forall f \in L^2, \quad (1 + \|U\|_{\infty}^2)^{-1} \|f\|_2^2 \leq \langle f, Sf \rangle \leq \|f\|_2^2 \quad (22)$$

and  $u_h$  can be solved explicitly as

$$u_h = c_h (I - S) \sqrt{p} - c_f \cdot S A^* \sqrt{q} + S A^* D_{\sqrt{q}} U g$$

Putting all together, we get

$$\begin{aligned} u_f &= c_f (I - T) \sqrt{q} + c_h T A \sqrt{p} + T D_{\sqrt{q}} U g \\ u_h &= c_h (I - S) \sqrt{p} - c_f S A^* \sqrt{q} + S A^* D_{\sqrt{q}} U g. \end{aligned}$$

Taking the inner products  $\langle \sqrt{q}, \cdot \rangle$  and  $\langle \sqrt{p}, \cdot \rangle$  on the first equation and the second equation respectively, recalling  $c_f = \langle \sqrt{q}, u_f \rangle$  and  $c_h = \langle \sqrt{p}, u_h \rangle$ , with  $\langle \sqrt{p}, \sqrt{p} \rangle = \langle \sqrt{q}, \sqrt{q} \rangle = 1$  since  $p, q \in \Delta^{\infty}$ , we are left with

$$\begin{aligned} c_f &= c_f (1 - \langle \sqrt{q}, T \sqrt{q} \rangle) + c_h \langle \sqrt{q}, T A \sqrt{p} \rangle + \langle \sqrt{q}, T D_{\sqrt{q}} U g \rangle \\ c_h &= c_h (1 - \langle \sqrt{p}, S \sqrt{p} \rangle) - c_f \langle \sqrt{p}, S A^* \sqrt{q} \rangle + \langle \sqrt{p}, S A^* D_{\sqrt{q}} U g \rangle \end{aligned}$$

Noting that  $c_f$  and  $c_h$  are cancelled out, using the identity  $TA = A^*S$  from the definition of  $S$  and  $T$ , we are left with the linear system of  $(c_f, c_h)$ :

$$\begin{bmatrix} \alpha & -\gamma \\ \gamma & \beta \end{bmatrix} \begin{bmatrix} c_f \\ c_h \end{bmatrix} = \begin{bmatrix} \langle \sqrt{q}, TD_{\sqrt{p}}Ug \rangle \\ \langle \sqrt{p}, SA^*D_{\sqrt{q}}Ug \rangle \end{bmatrix} \quad \text{where} \quad \begin{cases} \alpha & = \langle \sqrt{q}, T\sqrt{q} \rangle \\ \beta & = \langle \sqrt{p}, S\sqrt{p} \rangle \\ \gamma & = \langle \sqrt{q}, TA\sqrt{p} \rangle \end{cases}$$

Since  $\|\sqrt{p}\|_2 = \|\sqrt{q}\|_2 = 1$ , using the established bound of the spectrum of  $S$  and  $T$  (21)-(22), we know

$$\frac{1}{1 + \|U\|_\infty^2} \leq \alpha \leq 1, \quad \frac{1}{1 + \|U\|_\infty^2} \leq \beta \leq 1, \quad |\gamma| \leq \frac{1}{2}. \quad (23)$$

Here  $|\gamma| \leq 1/2$  follows from the fact that  $\|TA\|_{L^2 \rightarrow L^2} = \sqrt{\|(TA)(TA)^*\|_{L^2 \rightarrow L^2}}$  where  $(TA)(TA)^* = AA^*(I + AA^*)^{-2}$ , whose spectrum is given by  $\sigma(1 + \sigma)^{-2}$  with  $\sigma \in [0, +\infty)$  where  $\sigma$  is the spectrum of  $AA^*$ . Since  $\inf_{x \geq 0} x/(1+x)^2 = 1/4$ , taking square root, we get  $\|TA\|_{L^2 \rightarrow L^2} \leq 1/2$ . By the same argument,  $\|SA^*\|_{L^2 \rightarrow L^2} \leq 1/2$  holds.

Now, using the lower bound of  $\alpha$  and  $\beta$  in (23), the determinant of the matrix is strictly positive:

$$\det \begin{bmatrix} \alpha & -\gamma \\ \gamma & \beta \end{bmatrix} = \alpha\beta + \gamma^2 \geq \frac{1}{1 + \|U\|_\infty^2}.$$

Therefore, that matrix is invertible, and we obtain

$$\begin{bmatrix} c_f \\ c_h \end{bmatrix} = \frac{1}{\alpha\beta + \gamma^2} \begin{bmatrix} \beta & \gamma \\ -\gamma & \alpha \end{bmatrix} \begin{bmatrix} |\langle \sqrt{q}, TD_{\sqrt{p}}Ug \rangle| \\ |\langle \sqrt{p}, SA^*D_{\sqrt{q}}Ug \rangle| \end{bmatrix}$$

Taking the absolute value and using the estimate of  $|\alpha|, |\beta|, |\gamma|$  from (23), with  $|\langle \sqrt{q}, TD_{\sqrt{p}}Ug \rangle| \leq \|D_{\sqrt{p}}Ug\|_2$  and  $|\langle \sqrt{p}, SA^*D_{\sqrt{q}}Ug \rangle| \leq 2^{-1}\|D_{\sqrt{q}}Ug\|_2$ , the following inequality holds for each coordinate:

$$\begin{bmatrix} |c_f| \\ |c_h| \end{bmatrix} \leq (1 + \|U\|_\infty^2)^2 \begin{bmatrix} 1 & 2^{-1} \\ 2^{-1} & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2^{-1} \end{bmatrix} \cdot \|D_{\sqrt{p}}Ug\|_2 = (1 + \|U\|_\infty^2)^2 \|D_{\sqrt{p}}Ug\|_2 \begin{bmatrix} 5/4 \\ 1 \end{bmatrix}$$

With  $\|p\|_1 = 1$ , we have

$$\begin{aligned} \|D_{\sqrt{p}}Ug\|_2^2 &= \int \mu_{\text{ref}}(dx) (\sqrt{p}(x) \int U(x, y)g(y)\mu_{\text{ref}}(dy))^2 \\ &\leq \|p\|_1 \cdot \|U\|_\infty^2 \cdot \|g\|_1^2 \\ &= (\|U\|_\infty \|g\|_1)^2, \end{aligned}$$

so we get

$$\begin{bmatrix} |c_f| \\ |c_h| \end{bmatrix} \leq (1 + \|U\|_\infty^2)^2 \|U\|_\infty \|g\|_1 \begin{bmatrix} 5/4 \\ 1 \end{bmatrix}. \quad (24)$$

Let us finish the proof. By the second equation of (20),

$$u_h = c_h(I - S)\sqrt{p} - c_f SA^* \sqrt{q} + SA^* D_{\sqrt{q}} U g.$$

Taking  $\|\cdot\|_2$ , using  $\|I - S\|_{L^2 \rightarrow L^2} \leq 1$ ,  $\|SA^*\|_{L^2 \rightarrow L^2} \leq 1/2$  and  $\|D_{\sqrt{p}}Ug\|_2 \leq \|U\|_\infty \|g\|_1$ , we obtain

$$\|u_h\|_2 \leq |c_h| + 2^{-1}|c_f| + 2^{-1}\|U\|_\infty \|g\|_1.$$

Combined with the upper bound  $\|f\|_1 \leq \|u_h\|_2 + |c_h| + \|g\|_1$  from (19),

$$\|f\|_1 \leq |c_h| + 2^{-1}|c_f| + 2^{-1}\|U\|_\infty \|g\|_1 + |c_h| + \|g\|_1 = [2^{-1}, 2] \begin{bmatrix} |c_f| \\ |c_h| \end{bmatrix} + (2^{-1}\|U\|_\infty + 1)\|g\|_1$$

Finally, substituting the upper bounds of  $[|c_f|, |c_h|]$  in (24), with  $[2^{-1}, 2] \cdot [5/4, 1]^\top = 21/8$ , we complete the proof.

## Appendix E. Dynamical analysis for general utility

We provide dynamical analysis of MCHF and NLHF for general  $U$ . Note that Theorem 5 and Theorem 6 in the main document hold as corollaries of the general theorem.

For any  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$ , we define the linear operators  $U, U^* : L^1(\mu_{\text{ref}}) \rightarrow L^1(\mu_{\text{ref}})$  as

$$Uf(x) = \int \mu_{\text{ref}}(dy)U(x, y)f(y), \quad U^*f(y) = \int \mu_{\text{ref}}(dx)f(x)U(x, y).$$

For any  $p \in L^\infty(\mu_{\text{ref}})$  such that  $p(\cdot) \geq 0$  and  $\|p\|_{L^1(\mu_{\text{ref}})} = 1$ , define the linear operator  $J_p : L^1(\mu_{\text{ref}}) \rightarrow L^1(\mu_{\text{ref}})$  as  $J_p = \text{Diag}(p) + pp^*$  or more precisely

$$J_p f(x) = p(x)f(x) - p(x) \int \mu_{\text{ref}}(dx')p(x')f(x').$$

For  $*$  in  $\{\text{MC}, \text{NL}, \text{RL}\}$ , let  $p_* = \frac{d\mu_*}{d\mu_{\text{ref}}}$  be the density with respect to  $\mu_{\text{ref}}$ . The next theorem gives a Taylor expansion of  $p_{\text{MC}}$  and  $p_{\text{NL}}$  around a given additive utility  $g \oplus f$ .

**Theorem 25** *For any  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  and  $f, g \in L^\infty(\mu_{\text{ref}})$ , it holds that*

$$\begin{aligned} \|p_{\text{MC}}(U) - p_{\text{RL}}(f) - J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_{\text{RL}}(f)\|_{L^1(\mu_{\text{ref}})} &= o(\|U - g \oplus f\|_\infty), \\ \|p_{\text{NL}}(U) - p_{\text{RL}}(f) - J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_{\text{RL}}(-g)\|_{L^1(\mu_{\text{ref}})} &= o(\|U - g \oplus f\|_\infty) \end{aligned}$$

where  $\|\cdot\|_\infty = \|\cdot\|_{L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})}$ .

Note that the above theorem holds for any  $f, g \in L^\infty(\mu_{\text{ref}})$  and any  $U \in L^\infty(\mu_{\text{ref}} \otimes \mu_{\text{ref}})$  (possibly non-antisymmetric). If we take the L2 solution  $(\hat{g}, \hat{f}) \in \text{argmin}_{g, f} \|U - g \oplus f\|_2$  for Theorem 25, using  $\|U - \hat{f} \oplus \hat{g}\|_\infty \asymp \|U\|_\oplus$  from Proposition 12, we get

$$\begin{aligned} \|p_{\text{MC}}(U) - p_{\text{RL}}(\hat{f}) - J_{p_{\text{RL}}(\hat{f})}(U - \hat{g} \oplus \hat{f})^* p_{\text{RL}}(\hat{f})\|_{L^1(\mu_{\text{ref}})} &= o(\|U\|_\oplus) \\ \|p_{\text{NL}}(U) - p_{\text{RL}}(\hat{f}) - J_{p_{\text{RL}}(\hat{f})}(U - \hat{g} \oplus \hat{f})^* p_{\text{RL}}(-\hat{g})\|_{L^1(\mu_{\text{ref}})} &= o(\|U\|_\oplus). \end{aligned}$$

We see that the MHCH and the NLHF differs in their first order term by  $p_{\text{ref}}(\hat{f})$  and  $p_{\text{ref}}(-\hat{g})$ .

Next, we characterize the dynamics of the algorithm for MCHF.

**Theorem 26** Consider the iteration  $\mu_{\text{MC}}^t = \mu_{\text{MC}}^{t-1}P$  initialize at  $\mu_{\text{MC}}^0 = \mu_{\text{ref}}$ . Then,

$$\begin{aligned} & \|p_{\text{MC}}^1 - p_{\text{RL}}(f) - J_{p_{\text{RL}}(f)}(U - g \oplus f)^* \mathbf{1}\|_1 = o(\|U - g \oplus f\|_\infty) \\ \forall t \geq 2, & \|p_{\text{MC}}^t - p_{\text{RL}}(f) - J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_{\text{RL}}(f)\|_1 = o(\|U - g \oplus f\|_\infty). \end{aligned}$$

Next, we discuss iterative algorithms for NLHF.

**Theorem 27** Consider the following iteration:

$$\begin{aligned} \mu_{\text{NL}}^t & \in \operatorname{argmax}_{\mu} \int \nu_{\text{NL}}^{t-1}(dx) \mu(dy) U(x, y) - D_{\text{KL}}(\mu \mid \mu_{\text{ref}}) \\ \nu_{\text{NL}}^t & \in \operatorname{argmin}_{\nu} \int \nu(dx) \mu_{\text{NL}}^t(dy) U(x, y) + D_{\text{KL}}(\nu \mid \mu_{\text{ref}}) \end{aligned} \quad (25)$$

with  $\nu_{\text{NL}}^0 = \mu_{\text{ref}}$ . Then, the densities  $p_{\text{NL}}^t = d\mu_{\text{NL}}^t/d\mu_{\text{ref}}$  and  $q_{\text{NL}}^t = d\nu_{\text{NL}}^t/d\mu_{\text{ref}}$  satisfy

$$\begin{aligned} & \|p_{\text{NL}}^1 - p_{\text{RL}}(f) - J_{p_{\text{RL}}(f)}(U - g \oplus f)^* \mathbf{1}\|_1 = o(\|U - g \oplus f\|_{L^\infty}) \\ \forall t \geq 1, & \|q_{\text{NL}}^t - p_{\text{RL}}(-g) + J_{p_{\text{RL}}(-g)}(U - g \oplus f) p_{\text{RL}}(f)\|_1 = o(\|U - g \oplus f\|_{L^\infty}) \\ \forall t \geq 2, & \|p_{\text{NL}}^t - p_{\text{RL}}(f) - J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_{\text{RL}}(-g)\|_1 = o(\|U - g \oplus f\|_{L^\infty}). \end{aligned}$$

### E.1. Proof of Theorem 5

If  $U$  is antisymmetric, substituting  $(g, f) = (-\hat{f}, \hat{f})$  for Theorem 25, noting that  $J_p X p = (\text{Diag}(p) - p p^*) X p = \text{Diag}(p) X p$  for any antisymmetric  $X$  due to  $p^* X p = 0$ , and using  $\|U - (-\hat{f}) \oplus \hat{f}\|_\infty \asymp \|U\|_\oplus$  from Proposition 13, we complete the proof.

### E.2. Proof of Theorem 6

When  $U$  is antisymmetric, the update of  $\nu$  from  $\mu$  and the update of  $\mu$  from  $\nu$  in the definition of iterates (25) can be written as the same mapping. Thus, if we define the new iterate  $\tilde{\mu}_{\text{NL}}^t$  as  $\tilde{\mu}_{\text{NL}}^{2t-1} = \mu_{\text{NL}}^t$  and  $\tilde{\mu}_{\text{NL}}^{2t} = \nu_{\text{NL}}^t$ , and substituting  $-g = f = \hat{f}$  to Theorem 27, we complete the proof.

### E.3. Proof of Theorem 25

#### E.3.1. MCHF

We write  $U$  as

$$U = g \oplus f + \epsilon E, \quad E = \frac{U - g \oplus f}{\|U - g \oplus f\|_\infty}, \quad \epsilon = \|U - g \oplus f\|_\infty$$

and let

$$p_\epsilon \equiv p_{\text{MC}}(g \oplus f + \epsilon E).$$

Notice that  $p_0 = p_{\text{MC}}(g \oplus f) = p_{\text{RL}}(f)$ . Applying Theorem 17, we have

$$\|p_\epsilon - p_0 - \epsilon \dot{p}_0\|_1 = o(\epsilon)$$

where  $\dot{p}_0$  is given by

$$\dot{p}_0 = \left( (I - P_0^*) \Big|_{L_0^1} \right)^{-1} \left( y \mapsto \int \mu_{\text{ref}}(dx) p_0(x) P_0(x, y) \left( E(x, y) - \int \mu_{\text{ref}}(dy') P_0(x, y') E(x, y') \right) \right)$$

where

$$P_0(x, y) = \frac{\exp(g(x) + f(y))}{\int \mu_{\text{ref}}(dy) \exp(g(x) + f(y'))} = \frac{\exp(f(y))}{\int \mu_{\text{ref}}(dy) \exp(f(y'))} = p_{\text{RL}}(f)(y) = p_0(y).$$

Note that  $P^*h = 0$  for any  $h \in L_0^1$  since

$$P_0^*h(y) = \int \mu_{\text{ref}}(dx) h(x) P_0(x, y) = \left( \int \mu_{\text{ref}}(dx) h(x) \right) p_0(y) = 0,$$

and hence

$$\left( (I - P_0^*) \big|_{L_0^1} \right)^{-1} = I$$

On the other hand, substituting  $P_0(x, y) = p_0(y)$ , we have

$$\begin{aligned} & \int \mu_{\text{ref}}(dx) p_0(x) P_0(x, y) \left( E(x, y) - \int \mu_{\text{ref}}(dy') P_0(x, y') E(x, y') \right) \\ &= \int \mu_{\text{ref}}(dx) p_0(x) p_0(y) \left( E(x, y) - \int \mu_{\text{ref}}(dy') p_0(y') E(x, y') \right) \\ &= p_0(y) \left( \int \mu_{\text{ref}}(dx) p_0(x) E(x, y) - \int \int \mu_{\text{ref}}(dx) \mu_{\text{ref}}(dy') p_0(x) p_0(y') E(x, y') \right) \\ &= p_0(y) \left( E^* p_0(y) - \int \mu_{\text{ref}}(dy') p_0(y') E^* p_0(y) \right) \\ &= J_{p_0} E^* p_0(y). \end{aligned}$$

Therefore, the derivative  $p_0$  is given by

$$\dot{p}_0 = J_{p_0} E^* p_0$$

and hence

$$\|p_\epsilon - p_0 - \epsilon J_{p_0} E^* p_0\|_1 = o(\epsilon).$$

Note that  $\epsilon J_{p_0} E^* p_0 = J_{p_0} (\epsilon E)^* p_0$  since  $J_{p_0}$  and  $E^*$  are linear operators. Therefore, substituting  $p_\epsilon = p_{\text{MC}}(U)$ ,  $p_0 = p_{\text{RL}}(f)$ ,  $\epsilon E^* = (U - g \oplus f)$ , and  $\epsilon = \|U - g \oplus f\|_\infty$ , we complete the proof for MCHF.

### E.3.2. NLHF

We write  $U$  as

$$U = g \oplus f + \epsilon E, \quad E = \frac{U - g \oplus f}{\|U - g \oplus f\|_\infty}, \quad \epsilon = \|U - g \oplus f\|_\infty$$

and let

$$p_\epsilon = p_{\text{NL}}(g \oplus f + \epsilon E), \quad q_\epsilon = p_{\text{NL}}(g \oplus f - \epsilon E).$$

Note that  $p_0 = p_{\text{RL}}(f)$  and  $q_0 = p_{\text{RL}}(-g)$ . By Theorem 22, we have

$$\|p_\epsilon - p_0 - \epsilon \dot{p}_0\|_2 = o(\epsilon)$$

where the derivative is given by

$$\begin{aligned}\dot{p}_0 &= (I + J_{p_0}(g \oplus f)^* J_{q_0}(g \oplus f))^{-1} (J_{p_0} E^* q_0 - J_{p_0}(g \oplus f)^* J_{q_0} E p_0) \\ &= (I + J_{p_0}(f \oplus g) J_{q_0}(g \oplus f))^{-1} (J_{p_0} E^* q_0 - J_{p_0}(f \oplus g) J_{q_0} E p_0).\end{aligned}$$

Now, for any  $h \in L_0^2 = \{h \in L^2 : \int \mu_{\text{ref}}(dy) h(y) = 0\}$ ,

$$\begin{aligned}((f \oplus g)h)(x) &= \int \mu_{\text{ref}}(dy) (f \oplus g)(x, y) h(y) \\ &= \int \mu_{\text{ref}}(dy) (f(x) + g(y)) h(y) \\ &= \int \mu_{\text{ref}}(dy) g(y) h(y).\end{aligned}$$

Thus,  $(f \oplus g)h$  is a constant. Since  $J_{p_0}$  is a mean-subtracting operator, we get  $J_{p_0}(f \oplus g)h = 0$  and hence

$$J_{p_0}(f \oplus g) \big|_{L_0^2} = 0.$$

By the same argument, we also get

$$J_{q_0}(g \oplus f) \big|_{L_0^2} = 0.$$

Since the image space of  $J_{p_0}$  and  $J_{q_0}$  are both  $L_0^2$ , we get

$$\begin{aligned}\left( (I + J_{p_0}(f \oplus g) J_{q_0}(g \oplus f)) \big|_{L_0^2} \right)^{-1} &= I \\ J_{p_0} E^* q_0 - J_{p_0}(f \oplus g) J_{q_0} E p_0 &= J_{p_0} E^* q_0\end{aligned}$$

and hence

$$\dot{p}_0 = J_{p_0} E^* q_0.$$

Substituting  $p_0 = p_{\text{RL}}(f)$ ,  $q_0 = p_{\text{RL}}(-g)$ ,  $E = (U - g \oplus f)/\epsilon$ , and  $\epsilon = \|U - g \oplus f\|_\infty$ , we get

$$\|p_{\text{NL}}(U) - p_{\text{RL}}(f) - J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_{\text{RL}}(-g)\|_2 = o(\|U - g \oplus f\|_\infty).$$

Since  $\|\cdot\|_{L^1} \leq \|\cdot\|_2$ , this completes the proof.

#### E.4. Proof of Theorem 26

We write  $U$  as

$$U = g \oplus f + \epsilon E, \quad E = \frac{U - g \oplus f}{\|U - g \oplus f\|_\infty}, \quad \epsilon = \|U - g \oplus f\|_\infty$$

and let

$$p_\epsilon^t = P_\epsilon^* p_\epsilon^{t-1}, \quad p_\epsilon^0 = 1$$

where

$$P_\epsilon(x, y) = \frac{\exp((g \oplus f + \epsilon E)(x, y))}{\int \mu_{\text{ref}}(dy') \exp((g \oplus f + \epsilon E)(x, y'))} = \frac{\exp(f(y) + \epsilon E(x, y))}{\int \mu_{\text{ref}}(dy') \exp(f(y') + \epsilon E(x, y'))}.$$

Then it suffices to show

$$\lim_{\epsilon \rightarrow 0} \|\epsilon^{-1}(p_\epsilon^t - p_0^t) - \dot{p}_0^t\|_1 = 0$$

where

$$\dot{p}_0^t = \begin{cases} J_{p_{\text{RL}}(f)}(U - g \oplus f)^* \mathbf{1} & t = 1 \\ J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_{\text{RL}}(f) & t \geq 2 \end{cases}$$

Note that  $p_0^t = p_{\text{RL}}(f)$  for all  $t \geq 1$  and  $P_0(x, y) = p_{\text{RL}}(f)(y)$ . Furthermore, by the chain rule, we have

$$\partial_\epsilon P_\epsilon(x, y) = P_\epsilon(x, y)E(x, y) - P_\epsilon(x, y) \int \mu_{\text{ref}}(dy') P_\epsilon(x, y') E(x, y')$$

Taking  $\epsilon = 0$ ,

$$\dot{P}_0(x, y) = p_{\text{RL}}(f)(y)E(x, y) - p_{\text{RL}}(f)(y) \int \mu_{\text{ref}}(dy') p_{\text{RL}}(f)(y') E(x, y')$$

By the chain rule,  $p_0^1 = P_0^* \mathbf{1}$  is differentiable with derivative given by

$$\begin{aligned} \dot{p}_0^1(y) &= \dot{P}_0^* \mathbf{1}(y) = \int \mu_{\text{ref}}(dx) \left( p_{\text{RL}}(f)(y) E(x, y) - p_{\text{RL}}(f)(y) \int \mu_{\text{ref}}(dy') p_{\text{RL}}(f)(y') E(x, y') \right) \\ &= p_{\text{RL}}(f)(y) \int \mu_{\text{ref}}(dx) E(x, y) - p_{\text{RL}}(f)(y) \iint \mu_{\text{ref}}(dx) \mu_{\text{ref}}(dy') p_{\text{RL}}(f)(y') E(x, y') \\ &= p_{\text{RL}}(f)(y) E^* \mathbf{1}(y) - p_{\text{RL}}(f)(y) \int \mu_{\text{ref}}(dy') p_{\text{RL}}(f)(y') E^* \mathbf{1}(y') \\ &= J_{p_{\text{RL}}(f)}(U - g \oplus f)^* \mathbf{1}. \end{aligned}$$

For  $t = 2$ , by the chain rule,  $p_\epsilon^2$  is differentiable at  $\epsilon = 0$  with derivative given by

$$\dot{p}_0^2 = \dot{P}_0^* p_0^1(y) + P_0^* \dot{p}_0^1$$

Since  $P_0^* |_{L_0^1} = 0$  and  $\dot{p}_0^1 \in L_0^1$ , we get  $\dot{p}_0^2 = \dot{P}_0^* p_0^1(y)$ . By the same calculation for  $\dot{p}_0^1$ , we get

$$\dot{p}_0^2 = J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_0^1.$$

Since  $p_0^1 = p_{\text{RL}}(f)$ , we obtain

$$\dot{p}_0^2 = J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_{\text{RL}}(f).$$

Now we suppose that  $p_\epsilon^t$  is differentiable at  $\epsilon = 0$  with derivative given by  $\dot{p}_0^t = J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_{\text{RL}}(f)$ . Since  $\dot{p}_0^t \in L_0^1$  and  $p_0^t = p_{\text{RL}}(f)$ , by the chain rule,  $p_\epsilon^{t+1}$  is differentiable with derivative given by

$$\dot{p}_0^{t+1} = \dot{P}_0^* p_0^t + P_0^* \dot{p}_0^t = \dot{P}_0^* p_0^t = J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_0^t = J_{p_{\text{RL}}(f)}(U - g \oplus f)^* p_{\text{RL}}(f).$$

By induction, this completes the proof.

**E.5. Proof of Theorem 27**

We write  $U$  as

$$U = g \oplus f + \epsilon E, \quad E = \frac{U - g \oplus f}{\|U - g \oplus f\|_\infty}, \quad \epsilon = \|U - g \oplus f\|_\infty$$

and let  $p_\epsilon = d\mu_{\text{NL}}^t(U)/d\mu_{\text{ref}}$  and  $q_\epsilon = d\nu_{\text{NL}}^t(U)/d\mu_{\text{ref}}$  where

$$\begin{aligned} \mu_{\text{NL}}^t &\in \operatorname{argmax}_\mu \int \nu_{\text{NL}}^{t-1}(dx)\mu(dy)U(x, y) - D_{\text{KL}}(\mu \mid \mu_{\text{ref}}) \\ \nu_{\text{NL}}^t &\in \operatorname{argmin}_\nu \int \nu(dx)\mu_{\text{NL}}^t(dy)U(x, y) + D_{\text{KL}}(\nu \mid \mu_{\text{ref}}) \end{aligned}$$

with  $\nu_{\text{NL}}^0 = \mu_{\text{ref}}$ . Note that  $p_\epsilon^t$  and  $q_\epsilon^t$  satisfies

$$p_\epsilon^t = \mathcal{G}((g \oplus f + \epsilon E)^* q_\epsilon^{t-1}), \quad q_\epsilon^t = \mathcal{G}(-(g \oplus f + \epsilon E)p_\epsilon^t)$$

with  $q_\epsilon^0 = 1$ . Noting  $p_0 = p_{\text{RL}}(f)$  and  $q_0 = p_{\text{ref}}(-g)$ , it suffices to show

$$\dot{p}_0^t = \begin{cases} J_{p_{\text{RL}}(f)} E^* 1 & t = 1 \\ J_{p_{\text{RL}}(f)} E^* p_{\text{RL}}(-g) & t \geq 2 \end{cases}$$

and

$$\forall t \geq 1, \quad \dot{q}_0^t = -J_{q_{\text{RL}}(-g)} E p_{\text{RL}}(f).$$

Now, by the chain rule,  $p_\epsilon^1$  is differentiable at  $\epsilon$  with derivative given by

$$p_\epsilon^1 = J_{p_\epsilon^1} E^* q_\epsilon^0.$$

Taking  $\epsilon = 0$ , noting  $q_0^0 = 1$  and  $p_0^1 = p_{\text{ref}}(f)$ , we get

$$\dot{p}_0^1 = J_{p_{\text{RL}}(f)} E^* 1.$$

By the same argument, applying the chain rule for  $q_\epsilon^1 = \mathcal{G}(-(g \oplus f + \epsilon E)p_\epsilon^1)$ , we get

$$\dot{q}_\epsilon^1 = J_{q_\epsilon^1} (-E p_\epsilon^1 - (g \oplus f + \epsilon E) \dot{p}_\epsilon^1).$$

Taking  $\epsilon = 0$ , recalling that  $J_q(g \oplus f) \mid_{L_0^1} = 0$  and  $p_\epsilon^1 \in L_0^1$ , we get

$$\dot{q}_0^1 = J_{q_0^1} (-E p_0^1) = -J_{q_0^1} E p_0^1 = -J_{p_{\text{RL}}(-g)} E p_{\text{RL}}(f)$$

Now we assume that  $\dot{q}_\epsilon^t$  is differentiable for  $t \geq 1$ . Then, by the chain rule for  $p_\epsilon^{t+1} = \mathcal{G}((g \oplus f + \epsilon E)^* q_\epsilon^t)$ ,  $p_\epsilon^{t+1}$  is also differentiable with derivative given by

$$\dot{p}_\epsilon^{t+1} = J_{p_\epsilon^{t+1}} (E^* q_\epsilon^t + (g \oplus f + \epsilon E)^* \dot{q}_\epsilon^t).$$

Taking  $\epsilon = 0$ , using  $J_p(g \oplus f) \mid_{L_0^1} = 0$  and  $\dot{q}_\epsilon^t \in L_0^1$ , we get

$$\dot{p}_0^{t+1} = J_{p_0^{t+1}} (E^* q_0^t + (g \oplus f)^* \dot{q}_0^t) = J_{p_0^{t+1}} E^* q_0^t = J_{p_{\text{RL}}(f)} E^* p_{\text{RL}}(-g)$$

Now we further assume that  $p_0^{t+1}$  is differentiable. Then applying the chain rule for  $q_\epsilon^{t+1} = \mathcal{G}(-(g \oplus f + \epsilon E)p_\epsilon^{t+1})$ , we know that  $q_\epsilon^{t+1}$  is also differentiable with

$$\dot{q}_\epsilon^{t+1} = -J_{q_\epsilon^{t+1}}(Ep_\epsilon^{t+1} + (g \oplus f + \epsilon)\dot{p}_\epsilon^{t+1}).$$

Taking  $\epsilon = 0$ , using  $J_q(g \oplus f) |_{L_0^1} = 0$  and  $\dot{p}_\epsilon^{t+1} \in L_0^1$ , we have

$$\dot{q}_\epsilon^{t+1} = -J_{q_0^{t+1}}Ep_0^{t+1} = -J_{p_{\text{RL}}(-g)}Ep_{\text{RL}}(f).$$

By induction, this completes the proof.