

The Crucial Role of Samplers in Online Direct Preference Optimization

Ruizhe Shi*[†] and Runlong Zhou*[‡] and Simon S. Du[§]

Abstract

In this paper, we provide a rigorous analysis of DPO’s *convergence rates* with different sampling strategies under the exact gradient setting, revealing a separation: uniform sampling achieves *linear* convergence, while our proposed online sampler achieves *quadratic* convergence. We further adapt the sampler to practical settings by incorporating posterior distributions and *logit mixing*, demonstrating significant improvements over previous approaches. Our results not only offer insights into the theoretical standing of DPO but also pave the way for potential algorithm designs in the future.

1. Introduction

Aligning language models (LMs) to human preferences is a critical pursuit due to its great potentials to push forward artificial intelligence (AI) development, and to enable AI to serve humanity better [12]. Reinforcement learning from human feedback (RLHF) [2, 34] has been a widely-used approach, gaining tremendous successes in aligning LMs [19]. However, the multi-stage pipeline of RLHF, including reward model training and RL tuning, is sensitive to hyperparameters and costly to train. DPO [21] directly combines these stages and tunes LMs in an offline way, gaining popularity due to its stability and efficiency.

The empirical success of DPO has recently sparked a significant increase in interest for understanding its theoretical properties. Through modeling RLHF as a KL-regularized contextual bandit problem or Markov decision process, many works [13, 15, 24, 28, 29] obtain strong theoretical results and highlight the role of samplers in DPO. Specifically, they point out drawbacks of the offline sampler in vanilla DPO, and propose on-policy sampler or other samplers as better choices, as validated empirically [8, 10, 26].

However, these theoretical explanations are largely built upon traditional RL and analyze the impact of samplers from the view of data, namely sample complexity, thus involving some *impractical* assumptions, such as the access to an oracle for maximum likelihood estimation (MLE). Meanwhile, from the *optimization* perspective, the *convergence rates* of gradient descent in DPO within different sampling regimes remain an underexplored question. A particular setting of our interest is to give provable guarantees for an *online* sampler depending on the current policy.

* Equal contribution

[†] IIIS, Tsinghua University. Email: srz21@mails.tsinghua.edu.cn. Part of the work was done while Ruizhe Shi was visiting the University of Washington.

[‡] University of Washington. Email: vectorzh@cs.washington.edu

[§] University of Washington. Email: ssdu@cs.washington.edu

1.1. Contributions

To fill this research gap, we focus on analyzing the crucial role of samplers in DPO, from the view of optimization. Based on our theoretical findings, we can further derive a new effective approach, demonstrating advantages in empirical experiments over previous approaches. We summarize our contributions as follows:

- **Theoretical separations.** We analyze the convergence rates of DPO with various samplers under *tabular softmax parametrization*, and demonstrate theoretical advantages brought by specific samplers. Specifically, we show a separation that our proposed samplers, DPO-MIX-R and DPO-MIX-P, achieve quadratic convergence rates, while the commonly used one, DPO-Unif, can only achieve a linear convergence rate. Numerical simulations support our results.
- **Practical improvements.** We design a new sampler for practical DPO. LM alignment experiments show that under the same computation budget, our method demonstrates significant advantages over baselines. [Deferred to Appendix C](#).
- **Explainability and generalizability.** We show that our theoretical framework can explain many existing DPO variants and thus provides a new perspective on their theoretical advantages. [Deferred to Appendix C](#).

2. Preliminaries

We provide a thorough review of related literature in [Appendix A](#), and some basic preliminaries of RLHF in [Appendix B](#). In this paper, we look into the role of *samplers* in the performance of DPO. Now we formally define DPO with samplers, from the perspective of bandit algorithms. We first consider the scenario where we know the exact loss function and its gradient with respect to the model parameter θ .

Definition 1 (Exact DPO) Given an action set \mathcal{Y} , two samplers $\pi^{s1}, \pi^{s2} \in \Pi$ for sampling the first and second action respectively, a human preference oracle $p^* : \mathcal{Y} \times \mathcal{Y} \rightarrow \Delta(\{0, 1\})$, and hyperparameters $\beta, \eta \in \mathbb{R}_+$, the sampling probability and DPO loss function are defined as

$$\begin{aligned} \pi^s(y, y') &:= \text{sg}(\pi^{s1}(y)\pi^{s2}(y') + \pi^{s1}(y')\pi^{s2}(y)) , \\ \mathcal{L}_{\text{DPO}}(\theta) &:= - \sum_{y, y' \in \mathcal{Y}} \pi^s(y, y') p^*(y > y') \log \sigma \left(\beta \log \frac{\pi_\theta(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi_\theta(y')} \right) , \end{aligned} \quad (1)$$

and the parameter is updated by

$$\theta^{(t+1)} = \theta^{(t)} - \eta \alpha(\pi^{s1}, \pi^{s2}) \nabla_\theta \mathcal{L}_{\text{DPO}}(\theta^{(t)}) , \quad (2)$$

where $\alpha(\pi^{s1}, \pi^{s2})$ is a sampling coefficient determined by the samplers.

Remark 1 If the sampling regime is a mixture of ①: loss function \mathcal{L}_1 with sampling coefficient α_1 and ②: loss function \mathcal{L}_2 with sampling coefficient α_2 , the gradient update rule follows

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla_\theta \left(\alpha_1 \mathcal{L}_1(\theta^{(t)}) + \alpha_2 \mathcal{L}_2(\theta^{(t)}) \right) .$$

Note that ① and ② can have different sets of π^{s1} and π^{s2} .

In empirical studies, we do not have access to the exact gradients. Thus, we define the scenario of empirical DPO and make mild assumptions on the gradient estimation.

Definition 2 (Empirical DPO) Given noise scale $\sigma \in \mathbb{R}_+$, DPO (σ) is defined as DPO with the gradient update in Equation (2) as

$$\theta^{(t+1)} = \theta^{(t)} - \eta G^{(t)},$$

where $G_y^{(t)}$ is a random variable s.t. for $\forall y \in \mathcal{Y}$,

$$\frac{1}{\beta A} \left(G_y^{(t)} - \alpha(\pi^{s1}, \pi^{s2}) \nabla_{\theta_y} \mathcal{L}(\theta^{(t)}) \right) \sim \text{sub-Gaussian}(\sigma^2).$$

Remark 2 If the samplers are mixed, e.g., ① and ② in Remark 1, then we assume

$$\frac{1}{\beta A} \left(G_y^{(t)} - \nabla_{\theta_y} \left(\alpha_1 \mathcal{L}_1(\theta^{(t)}) + \alpha_2 \mathcal{L}_2(\theta^{(t)}) \right) \right) \sim \text{sub-Gaussian}(\sigma^2).$$

The closed form solution π^* in Equation (6) satisfies $r(y) - r(y') - \beta \log \frac{\pi^*(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi^*(y')} = 0$, which thus motivates us to study the convergence rate. With the update rule formally defined, now we ask:

How fast can $r(y) - r(y') - \beta \log \frac{\pi_{\theta^{(t)}}(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi_{\theta^{(t)}}(y')}$ **converge to 0, for** $\forall y, y' \in \mathcal{Y}$?

We will study the convergence rates for three sampling regimes: one sampling uniformly on the action space \mathcal{Y} and two with mixtures of samplers. They are defined in Definitions 3 to 5.

Definition 3 (Uniform sampler) DPO-UNIF is defined as DPO with π^{s1}, π^{s2} as

$$\pi^{s1}(\cdot) = \pi^{s2}(\cdot) = \text{Uniform}(\mathcal{Y}),$$

and $\alpha(\pi^{s1}, \pi^{s2}) = 2|\mathcal{Y}|^2$.

Definition 4 (Reward-guided mixed sampler) DPO-MIX-R is defined as DPO with π^{s1}, π^{s2} as

$$\textcircled{1} \begin{cases} \pi^{s1}(\cdot) = \text{Uniform}(\mathcal{Y}), \\ \pi^{s2}(\cdot) = \text{Uniform}(\mathcal{Y}), \end{cases} \quad \textcircled{2} \begin{cases} \pi^{s1}(\cdot) \propto \text{Uniform}(\mathcal{Y}) \cdot \exp(r(\cdot)), \\ \pi^{s2}(\cdot) \propto \text{Uniform}(\mathcal{Y}) \cdot \exp(-r(\cdot)), \end{cases}$$

and $\alpha_1 = |\mathcal{Y}|^2$, $\alpha_2 = \sum_{y, y' \in \mathcal{Y}} \exp(r(y) - r(y'))$.

Definition 5 (Policy-difference-guided mixed sampler) DPO-MIX-P is defined as DPO with π^{s1}, π^{s2} as

$$\textcircled{1} \begin{cases} \pi^{s1}(\cdot) = \text{Uniform}(\mathcal{Y}), \\ \pi^{s2}(\cdot) = \text{Uniform}(\mathcal{Y}), \end{cases} \quad \textcircled{2} \begin{cases} \pi^{s1}(\cdot) \propto \text{Uniform}(\mathcal{Y}) \cdot (\pi(\cdot)/\pi_{\text{ref}}(\cdot))^\beta, \\ \pi^{s2}(\cdot) \propto \text{Uniform}(\mathcal{Y}) \cdot (\pi_{\text{ref}}(\cdot)/\pi(\cdot))^\beta, \end{cases}$$

and $\alpha_1 = |\mathcal{Y}|^2$, $\alpha_2 = \sum_{y, y' \in \mathcal{Y}} \left(\frac{\pi(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi(y')} \right)^\beta$.

Remark 3 Definition 4 does not define a practical sampler as r is unknown, but it is important to display our idea of using a mixture of sampling policies. In Definition 5, ② can also be written as $\pi^{s1} \propto \exp(\beta(\theta - \theta_{\text{ref}}))$, $\pi^{s2} \propto \exp(\beta(\theta_{\text{ref}} - \theta))$. Uniform(\mathcal{Y}) in Definitions 4 and 5 is for consistency with Appendix C, where we use a posterior distribution over \mathcal{Y} .

3. Main Results

We show our main results on convergence rates in this section. In summary, our proposed mixed samplers can provably achieve: 1) exponentially faster convergence rates (*quadratic v.s. linear*) compared with the uniform sampler in the exact gradient setting, and 2) linear convergence rates to the noise scale when we have only unbiased estimations of the gradient. Numerical simulations corroborate these theories.

3.1. Theoretical Findings

We present theories regarding convergence rates of different sampling regimes for exact DPO and empirical DPO in this subsection, along with their proof sketches. We first define important notations:

$$\begin{aligned}\Delta(y, y'; \theta) &:= \sigma(r(y) - r(y')) - \sigma\left(\beta \log \frac{\pi_\theta(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi_\theta(y')}\right), \\ \delta(y, y'; \theta) &:= r(y) - r(y') - \beta \log \frac{\pi_\theta(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi_\theta(y')}.\end{aligned}$$

Then we can obtain

$$\nabla_\theta \mathcal{L}(\theta) = -\beta \sum_{y, y'} \pi^s(y, y') \Delta(y, y'; \theta) \mathbb{1}_y$$

by plugging $p^\star(y, y') = \sigma(r(y) - r(y'))$ and $\sigma(-x) = 1 - \sigma(x)$ into the derivative of Equation (1). Hence, we can derive the iteration equation for δ :

$$\begin{aligned}\delta(y, y'; \theta^{(t+1)}) &= \delta(y, y'; \theta^{(t)}) \\ &\quad - \eta \beta \alpha (\pi^{s1}, \pi^{s2}) \sum_{y''} \left(\pi^s(y, y'') \Delta(y, y''; \theta^{(t)}) - \pi^s(y', y'') \Delta(y', y''; \theta^{(t)}) \right).\end{aligned}\quad (3)$$

We state the common condition for the upper bounds for simplicity:

Condition 1 *Given an action set \mathcal{Y} , it satisfies $r(y) \in [0, 1], \forall y \in \mathcal{Y}$. $\pi_{\theta(0)}$ is initialized as π_{ref} , and the regularization coefficient is $\beta \in \mathbb{R}_+$. Use the learning rate $\eta = \frac{1}{\beta^2 |\mathcal{Y}|}$.*

3.1.1. FOR EXACT DPO

For DPO-`Unif`, we have that $\pi^s(y, y') = 2/|\mathcal{Y}|^2$, making the coefficients of each Δ on the RHS of Equation (3) identical by absolute values. To proceed, we claim a lower bound as $\sigma' \left(\log \frac{\pi_\theta(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi_\theta(y')} \right) \geq \sigma'_{\min}$, and use Lagrange interpolation, namely $\sigma'_{\min} \leq (\sigma(x) - \sigma(y))/(x - y) \leq 1/4$, to transform Δ into δ . By carefully computing the coefficients of each δ and picking learning rate, we arrive at a linear convergence. Using this linear convergence, we can turn back to bound σ'_{\min} , completing the proof. See detailed proof in Appendix D.1.

Theorem 1 (Upper bound of DPO-`Unif`) *Under Condition 1, DPO-`Unif` satisfies*

$$\left| \delta(y, y'; \theta^{(T)}) \right| \leq 0.588^T, \quad \forall y, y' \in \mathcal{Y},$$

where $T \in \mathbb{N}$ is the number of iterations.

The construction of the lower bound is based on a simple 3-armed bandit setting. We use Taylor expansion to transform Δ into δ , and note that the quadratic remainders can be negligible when θ is close to the optimal point. And thus the linear transformation can only achieve linear convergence. See detailed proof in Appendix D.1.

Theorem 2 (Lower bound of DPO-Unif) *Let $|\mathcal{Y}| = 3$, $r(y_1) = 0, r(y_2) = 1/3, r(y_3) = 1$, and $\pi_{\text{ref}} = \text{Uniform}(\mathcal{Y})$. For any $\beta \in \mathbb{R}_+$ and learning rate $\eta \in (0, \frac{2}{\beta^2|\mathcal{Y}|}]$, there always exists small enough $\epsilon \in \mathbb{R}_+$, for any initialization $\pi_{\theta^{(0)}}$ satisfying $\max_{y,y' \in \mathcal{Y}} |\delta(y, y'; \theta^{(0)})| \leq \epsilon$ and $\min_{y,y' \in \mathcal{Y}} |\delta(y, y'; \theta^{(0)})| > 0$, DPO-Unif satisfies*

$$\max_{y,y' \in \mathcal{Y}} |\delta(y, y'; \theta^{(T)})| \geq \gamma^T,$$

where $T \in \mathbb{N}$ is the number of iterations and γ is a constant depending on $\theta^{(0)}$.

Next we elaborate the idea of transforming Δ into δ using Taylor expansion, and show how to eliminate the linear term using appropriate samplers and learning rate. For Theorem 3, we can apply Taylor expansion at $r(y_1) - r(y_2)$ (while for Theorem 4 we apply at $\beta \log \frac{\pi_{\theta^{(t)}}(y_1)\pi_{\text{ref}}(y_2)}{\pi_{\text{ref}}(y_1)\pi_{\theta^{(t)}}(y_2)}$), and get

$$\Delta(y_1, y_2; \theta^{(t)}) = \sigma'(r(y_1) - r(y_2))\delta(y_1, y_2; \theta^{(t)}) + \frac{\sigma''(\xi_R)}{2}\delta(y_1, y_2; \theta^{(t)})^2.$$

If we let $\pi^s(y_1, y_2) \propto 1/\sigma'(r(y_1) - r(y_2))$ as in Definition 4, then

$$\pi^s(y, y'')\Delta(y, y''; \theta^{(t)}) - \pi^s(y', y'')\Delta(y', y''; \theta^{(t)}) = \text{constant} \cdot \delta(y, y'; \theta^{(t)}) + \text{quadratic term}.$$

Finally we pick an appropriate η to eliminate the initial linear term in Equation (3) and thus establish a quadratic convergence. This observation motivates our design of samplers and proofs. The detailed proofs of Theorems 3 and 4 can be found in Appendices D.2 and D.3.

Theorem 3 (Upper bound of DPO-Mix-R) *Under Condition 1, DPO-Mix-R satisfies*

$$|\delta(y, y'; \theta^{(T)})| \leq 0.5^{2T-1}, \quad \forall y, y' \in \mathcal{Y},$$

where $T \in \mathbb{N}$ is the number of iterations.

Theorem 4 (Upper bound of DPO-Mix-P) *Under Condition 1, DPO-Mix-P satisfies*

$$|\delta(y, y'; \theta^{(T)})| \leq 0.611^{2T-1}, \quad \forall y, y' \in \mathcal{Y},$$

where $T \in \mathbb{N}$ is the number of iterations.

3.1.2. FOR EMPIRICAL DPO

As in Definition 2, exact gradients are inaccessible in practice. Here we show the guarantees of DPO-Mix-R and DPO-Mix-P with only unbiased estimation of gradients, that they can achieve linear convergence rates to the noise scale. The proofs of Theorems 5 and 6 can be found in Appendix E.

Theorem 5 Under Condition 1 with the noise scale $\sigma \in (0, 1/576)$, $\text{DPO-Mix-R}(\sigma)$ satisfies

$$\sqrt{\mathbb{E}[\delta(y, y'; \theta^{(T)})^2]} \leq 14\sigma, \quad \forall y, y' \in \mathcal{Y},$$

where $T = \lceil \log \frac{1}{\sigma} \rceil$ is the number of iterations.

Theorem 6 Under Condition 1 with the noise scale $\sigma \in (0, 1/576)$, $\text{DPO-Mix-P}^*(\sigma)$ satisfies

$$\sqrt{\mathbb{E}[\delta(y, y'; \theta^{(T)})^2]} \leq 14\sigma, \quad \forall y, y' \in \mathcal{Y},$$

where $T = \lceil \log \frac{1}{\sigma} \rceil$ is the number of iterations, and $\text{DPO-Mix-P}^*(\sigma)$ is $\text{DPO-Mix-P}(\sigma)$ with a rejection sampling process: each time we get $y, y' \in \mathcal{Y}$ sampled from ②, if $\psi(y, y'; \theta^{(t)}) := \left| \beta \log \frac{\pi_{\theta^{(t)}}(y)\pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y)\pi_{\theta^{(t)}}(y')} \right| > 1$, then reject this data pair with probability $1 - \frac{e+e^{-1}}{e^\psi+e^{-\psi}}$; and α_2 needs to be changed to $\frac{1}{2} \sum_{y, y' \in \mathcal{Y}} \min\{e^{\psi(y, y'; \theta^{(t)})} + e^{-\psi(y, y'; \theta^{(t)})}, e + e^{-1}\}$.

3.2. Numerical Simulations

We verify our theoretical findings with numerical simulations in *contextual bandits*. As shown in Figure 1, the two proposed samplers DPO-Mix-P and DPO-Mix-R show great improvements over DPO-Unif . The detailed configurations and more results can be found in Appendix G.1.

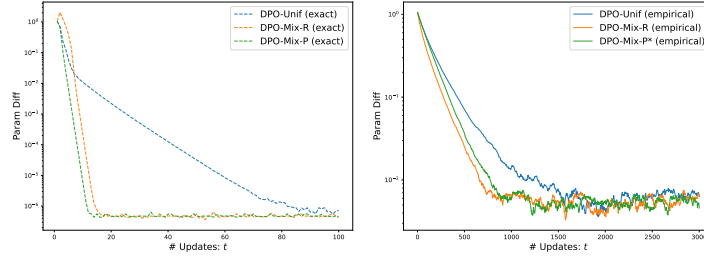


Figure 1: **Contextual bandit experiments for exact DPO and empirical DPO.** The x -axis is the number of gradient updates, and the y -axis is the total parameter difference $\sum_{y, y'} \delta(y, y'; \theta^{(t)})^2$. The left figure illustrates exact DPO, and the right figure illustrates empirical DPO. The separation is clear in exact DPO, and still exists in empirical DPO.

4. Conclusion

This paper studies the convergence rates of DPO with different samplers. We demonstrate that DPO-Mix-R and DPO-Mix-P offer quadratic convergence rates, outperforming the linear rate of DPO-Unif . Our theoretical findings are supported by numerical simulations and LM alignment experiments.

It is also important to acknowledge our limitations. 1) The selection of the posterior distribution is not unique, and thus many useful samplers have yet to be developed from our framework and need further experiments. 2) The convergence analysis is based on *tabular softmax parametrization*, and a future direction would be exploring more practical settings such as *log-linear parametrization* and *function approximation*.

REFERENCES

- [1] Mohammad Gheshlaghi Azar, Mark Rowland, Bilal Piot, Daniel Guo, Daniele Calandriello, Michal Valko, and Rémi Munos. A general theoretical paradigm to understand learning from human preferences. *ArXiv*, abs/2310.12036, 2023.
- [2] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova Dassarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, John Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom B. Brown, Jack Clark, Sam McCandlish, Christopher Olah, Benjamin Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback. *ArXiv*, abs/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. ISSN 00063444.
- [4] Angelica Chen, Sadhika Malladi, Lily H. Zhang, Xinyi Chen, Qiuyi Zhang, Rajesh Ranganath, and Kyunghyun Cho. Preference learning algorithms do not learn preference rankings. *ArXiv*, abs/2405.19534, 2024.
- [5] Paul Francis Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *ArXiv*, abs/1706.03741, 2017.
- [6] Mucong Ding, Souradip Chakraborty, Vibhu Agrawal, Zora Che, Alec Koppel, Mengdi Wang, A. S. Bedi, and Furong Huang. Sail: Self-improving efficient online alignment of large language models. *ArXiv*, abs/2406.15567, 2024.
- [7] Hanze Dong, Wei Xiong, Deepanshu Goyal, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. Raft: Reward ranked finetuning for generative foundation model alignment. *arXiv preprint arXiv:2304.06767*, 2023.
- [8] Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf, 2024.
- [9] Yann Dubois, Bal’azs Galambosi, Percy Liang, and Tatsunori Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators. *ArXiv*, abs/2404.04475, 2024.
- [10] Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, et al. Direct language model alignment from online ai feedback. *arXiv preprint arXiv:2402.04792*, 2024.
- [11] Jiaming Ji, Mickel Liu, Juntao Dai, Xuehai Pan, Chi Zhang, Ce Bian, Ruiyang Sun, Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of llm via a human-preference dataset. *ArXiv*, abs/2307.04657, 2023.
- [12] Jiaming Ji, Tianyi Qiu, Boyuan Chen, Borong Zhang, Hantao Lou, Kaile Wang, Yawen Duan, Zhonghao He, Jiayi Zhou, Zhaowei Zhang, Fanzhi Zeng, Kwan Yee Ng, Juntao Dai, Xuehai Pan, Aidan O’Gara, Yingshan Lei, Hua Xu, Brian Tse, Jie Fu, Stephen Marcus McAleer,

- Yaodong Yang, Yizhou Wang, Song-Chun Zhu, Yike Guo, and Wen Gao. Ai alignment: A comprehensive survey. *ArXiv*, abs/2310.19852, 2023.
- [13] Saeed Khaki, JinJin Li, Lan Ma, Liu Yang, and Prathap Ramachandra. Rs-dpo: A hybrid rejection sampling and direct preference optimization method for alignment of large language models. *ArXiv*, abs/2402.10038, 2024.
- [14] Tianlin Liu, Shangmin Guo, Leonardo Bianco, Daniele Calandriello, Quentin Berthet, Felipe Llinares-López, Jessica Hoffmann, Lucas Dixon, Michal Valko, and Mathieu Blondel. Decoding-time realignment of language models. In *Forty-first International Conference on Machine Learning*, 2024.
- [15] Tianqi Liu, Yao Zhao, Rishabh Joshi, Misha Khalman, Mohammad Saleh, Peter J Liu, and Jialu Liu. Statistical rejection sampling improves preference optimization. In *The Twelfth International Conference on Learning Representations*, 2024.
- [16] Zhihan Liu, Miao Lu, Shenao Zhang, Boyi Liu, Hongyi Guo, Yingxiang Yang, Jose Blanchet, and Zhaoran Wang. Provably mitigating overoptimization in rlhf: Your sft loss is implicitly an adversarial regularizer. *ArXiv*, abs/2405.16436, 2024.
- [17] Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward. *ArXiv*, abs/2405.14734, 2024.
- [18] Rémi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland, Zhaohan Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Andrea Michi, et al. Nash learning from human feedback. *arXiv preprint arXiv:2312.00886*, 2023.
- [19] OpenAI. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- [20] Philippe Rigollet. Lecture Notes for High-Dimensional Statistics - 18.S997, Spring 2015. https://ocw.mit.edu/courses/18-s997-high-dimensional-statistics-spring-2015/619e4ae252f1b26cbe0f7a29d5932978_MIT18_S997S15_CourseNotes.pdf, 2015. Accessed: 2024-08-28.
- [21] 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. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [22] 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*, abs/2404.03715, 2024.
- [23] Ruizhe Shi, Yifang Chen, Yushi Hu, Alisa Liu, Hanna Hajishirzi, Noah A. Smith, and Simon S. Du. Decoding-time language model alignment with multiple objectives. *ArXiv*, abs/2406.18853, 2024.
- [24] Yuda Song, Gokul Swamy, Aarti Singh, J. Andrew Bagnell, and Wen Sun. The importance of online data: Understanding preference fine-tuning via coverage, 2024.

- [25] Gokul Swamy, Christoph Dann, Rahul Kidambi, Zhiwei Steven Wu, and Alekh Agarwal. A minimaximalist approach to reinforcement learning from human feedback. *arXiv preprint arXiv:2401.04056*, 2024.
- [26] Fahim Tajwar, Anika Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal, on-policy data. *ArXiv*, abs/2404.14367, 2024.
- [27] Christian Wirth, Riad Akrouf, Gerhard Neumann, and Johannes Fürnkranz. A survey of preference-based reinforcement learning methods. *J. Mach. Learn. Res.*, 18:136:1–136:46, 2017.
- [28] Tengyang Xie, Dylan J Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Awadallah, and Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit q*-approximation for sample-efficient rlhf. *arXiv preprint arXiv:2405.21046*, 2024.
- [29] Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for RLHF under KL-constraint. In *Forty-first International Conference on Machine Learning*, 2024.
- [30] Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan, Lingfeng Shen, Benjamin Van Durme, Kenton Murray, and Young Jin Kim. Contrastive preference optimization: Pushing the boundaries of llm performance in machine translation. *ArXiv*, abs/2401.08417, 2024.
- [31] Natalia Zhang, Xinqi Wang, Qiwen Cui, Runlong Zhou, Sham M. Kakade, and Simon S. Du. Multi-agent reinforcement learning from human feedback: Data coverage and algorithmic techniques. 2024.
- [32] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Haoteng Zhang, Joseph Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. *ArXiv*, abs/2306.05685, 2023.
- [33] Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human feedback from pairwise or k-wise comparisons. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 43037–43067. PMLR, 23–29 Jul 2023.
- [34] Daniel M. Ziegler, Nisan Stiennon, Jeff Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *ArXiv*, abs/1909.08593, 2019.

Appendix

Appendix A. Related Work

Theoretical study of RLHF/DPO. Zhu et al. [33] formulate RLHF as contextual bandits, and prove the convergence of the maximum likelihood estimator. Xiong et al. [29] further consider KL-regularization and show the benefits in sample complexity of online exploration in DPO. Xie et al. [28] study the online exploration problem from the perspective of KL-regularized Markov decision processes, and show provable guarantees in sample complexity of a exploration bonus. Liu et al. [16] investigate the overoptimization issue, and prove a finite-sample suboptimality gap. Song et al. [24] show a separation of coverage conditions for offline DPO and online RLHF. These works primarily focus on the perspective of data, which is widely adopted in RL literature. For Xie et al. [28], Xiong et al. [29], their policy update iteration is to directly solve MLE instead of doing *gradient descent* as ours. Song et al. [24] focus on data coverage, and have not studied the convergence rates. In contrast, this paper analyzes DPO from the perspective of *optimization*, offering a complementary while more practical viewpoint.

Variants of DPO. There are two line of works exploring the variants of DPO. 1) *Objective function*. Ψ -PO [1] changes the reward term to alternate mappings from preference pairs. RPO [16] adds an imitation loss to mitigate the overoptimization issue. CPO [30] removes the π_{ref} term and adds an imitation loss to ensure that the policy does not deviate too much. SimPO [17] also removes the π_{ref} term for efficiency, while using length normalization for better length control. 2) *Sampler*. Khaki et al. [13], Liu et al. [15] utilize rejection sampling to adjust the data distribution to the theoretically-optimal policy before training. On-policy DPO [6, 10, 26] emphasize the importance of the on-policy sampler. Iterative DPO [8, 29] introduces an iterative training scheme, where an online policy is used to generate data pairs, annotated by a gold reward model, and the DPO training is subsequently applied to update the policy. XPO [28] follows the setting of iterative DPO, and adds an optimistic term to the DPO objective. In this paper, we focus on the latter direction, and only study the original objective.

Other RLHF approaches. There is also a line of works [18, 22, 25, 31] studying RLHF from a game-theoretic perspective. Nash-MD-PG in Munos et al. [18] uses a geometric mixture of online policy and reference policy without specifying the mixing weight. Rosset et al. [22] re-formulates the DPO pipeline and shows theoretical guarantees for the on-policy sampler with an MLE oracle.

Appendix B. Basic Preliminaries

Notations. Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be the sigmoid function, where $\sigma(x) = 1/(1 + \exp(-x))$. For any set X , $\Delta(X)$ represents the set of probability distributions over X . $\text{sg}()$ is the stopping-gradient operator. Let $\mathbb{1}_k$ be a vector with 1 on the dimension corresponding to k and 0 on others (the dimension of this vector is implicitly defined from the context).

B.1. Standard Bandit Learning

Firstly, we give basic concepts of standard bandit learning, which found the basis for RLHF.

Multi-armed bandits and contextual bandits. A multi-armed bandit has an arm (action) space \mathcal{Y} and a reward function $r : \mathcal{Y} \rightarrow [0, 1]$. A contextual bandit has a context space \mathcal{X} , an arm space \mathcal{Y} , and a reward function $r : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$. In this work, the user prompt is viewed as a context, and the agent response is viewed as an arm. To simplify notations, our results are stated in *multi-armed bandits* versions. The statements and proofs can be easily extended to *contextual bandits*. Thus, we will omit the prompts (contexts) and slightly abuse the notations throughout Sections 2 and 3.

Policies. A policy $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ maps each context to a probability simplex over the arm space. For multi-armed bandits, a policy is instead a probability distribution over the arm space. We denote Π as the set of policies we study. Under *tabular softmax parametrization* which is common in previous works [1, 18, 21, 25], the policy π is parameterized by $\theta \in \mathbb{R}^{|\mathcal{Y}|}$: for any $y \in \mathcal{Y}$,

$$\pi_{\theta}(y) = \frac{\exp(\theta_y)}{\sum_{y' \in \mathcal{Y}} \exp(\theta_{y'})} .$$

The goal is to find the optimal policy maximizing the expected reward (with regularization).

B.2. Reinforcement Learning from Human Feedback (RLHF)

Secondly, we introduce RLHF / preference-based reinforcement learning (PBRL) problem [5, 25, 27] and current approaches.

Bradley-Terry (BT) model. Given an implicit reward oracle $r : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$, [3] assume that human preference distribution $p^* : \mathcal{X} \times \mathcal{Y} \times \mathcal{Y} \rightarrow \Delta(\{0, 1\})$ satisfies:

$$p^*(y_1 > y_2 | x) = \sigma(r(x, y_1) - r(x, y_2)) .$$

This means that conditioned on prompt x , response y_1 is favored over y_2 with probability $p^*(y_1 > y_2 | x)$ by human annotators.

RLHF [2, 34]. A human preference dataset $\mathcal{D} = \{(x^{(i)}, y_w^{(i)}, y_l^{(i)})\}_{i=1}^N$ means that in the i^{th} sample, $y_w^{(i)} > y_l^{(i)}$ conditioned on $x^{(i)}$. The reward function $r : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is learned with parameter ϕ using a negative log-likelihood loss:

$$\mathcal{L}_r(\phi) = -\frac{1}{N} \sum_{i=1}^N \left[\log \sigma \left(r_{\phi}(x^{(i)}, y_w^{(i)}) - r_{\phi}(x^{(i)}, y_l^{(i)}) \right) \right] . \quad (4)$$

Given $\pi_1, \pi_2 \in \Pi$, $\mathbb{E}_{x \in \mathcal{X}} \text{KL}(\pi_1(\cdot | x) \| \pi_2(\cdot | x))$ is abbreviated as $\text{KL}(\pi_1 \| \pi_2)$. Based on a reference policy π_{ref} , the goal of RLHF is to maximize the obtained rewards with a KL-divergence penalty:

$$\pi^* = \operatorname{argmax}_{\pi \in \Pi} \mathbb{E}_{x \in \mathcal{X}, y \in \pi(\cdot | x)} r_{\phi}(x, y) - \beta \text{KL}(\pi \| \pi_{\text{ref}}) , \quad (5)$$

where $\beta \in \mathbb{R}_+$ is the regularization coefficient. Additionally, under *tabular softmax parametrization*, we can directly write out the closed-form solution (Equation (4) in [21]):

$$\pi^*(y | x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y | x) \exp \left(\frac{1}{\beta} r_{\phi}(x, y) \right) , \quad \forall x \in \mathcal{X}, y \in \mathcal{Y} , \quad (6)$$

where $Z(x) = \sum_{y \in \mathcal{Y}} \pi_{\text{ref}}(y | x) \exp \left(\frac{1}{\beta} r_{\phi}(x, y) \right)$ is the partition function.

Direct Preference Optimization (DPO, [21]). DPO integrates reward learning with policy learning. Given the human preference dataset $\mathcal{D} = \{(x^{(i)}, y_w^{(i)}, y_l^{(i)})\}_{i=1}^N$, the DPO policy π is learned with parameter θ using a negative log-likelihood loss:

$$\mathcal{L}_\pi(\theta) = -\frac{1}{N} \sum_{i=1}^N \log \sigma \left(\beta \log \frac{\pi_\theta(y_w^{(i)}|x^{(i)})}{\pi_{\text{ref}}(y_w^{(i)}|x^{(i)})} - \beta \log \frac{\pi_\theta(y_l^{(i)}|x^{(i)})}{\pi_{\text{ref}}(y_l^{(i)}|x^{(i)})} \right),$$

which can be directly derived by combining Equations (4) and (6).

Appendix C. Implications for Practical DPO

In this section, we show the implications of theoretical results in Section 3 for practical DPO design.

C.1. Aligning Theory to Practice

Rethinking DPO. We can rewrite a policy $\pi \in \Pi$ as $\pi(y|x) \propto \pi_{\text{ref}}(y|x) e^{\varphi(x,y)/\beta}$, where $\varphi(x, y) \in \mathbb{R}_+$. Then the training objective of DPO can be rewritten as:

$$\varphi^*(x, \cdot) = \underset{\varphi(x, \cdot)}{\operatorname{argmin}} \sum_{y_1, y_2 \in \mathcal{Y}} \pi^s(y_1, y_2|x) \cdot \underbrace{(-\sigma(r(x, y_1) - r(x, y_2)) \log \sigma(\varphi(x, y_1) - \varphi(x, y_2)))}_{\text{cross entropy loss}},$$

which is learning a reward model $\varphi(x, y)$ towards $r(x, y) + C(x)$, where $C(x) \in \mathbb{R}$ is a constant. In Section 3, we have discussed the role of samplers in this implicit reward learning stage. Here we introduce a lemma (for multi-armed bandits) to connect it with the final performance.

Lemma 1 (Performance difference lemma) *For any θ , define its value as*

$$V^\theta := \mathbb{E}_{y \sim \pi_\theta} r(y) - \beta \text{KL}(\pi_\theta \| \pi_{\text{ref}}),$$

and let V^* be the value of the optimal policy π^* in Equation (6), then we have

$$\begin{aligned} V^* - V^\theta &= \sum_{y, y' \in \mathcal{Y}} \pi^*(y) \pi_\theta(y') \left(r(y) - r(y') - \beta \log \frac{\pi_\theta(y) \pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y) \pi_\theta(y')} \right) - \beta \text{KL}(\pi^* \| \pi_\theta) \\ &\leq \sum_{y, y' \in \mathcal{Y}} \pi^*(y) \pi_\theta(y') \left(r(y) - r(y') - \beta \log \frac{\pi_\theta(y) \pi_{\text{ref}}(y')}{\pi_{\text{ref}}(y) \pi_\theta(y')} \right). \end{aligned} \quad (7)$$

Setting the posterior. Lemma 1 indicates that reward learning should concentrate on responses with high probabilities for π^* and π_θ , and thus motivates us to change the distribution over \mathcal{Y} to a posterior distribution close to π^* or π_θ in practical implementation. This perspective provides an alternate explanation for [15], which uses rejection sampling to align the sampling distribution to π^* . Considering the fact that π^* is usually inaccessible, we propose to let $\pi_\theta^{2\beta}$ be the posterior distribution. Setting the sampling temperature as 2β , we can thus derive our new practical algorithm following Definition 5:

$$\textcircled{1} \begin{cases} \pi^{s1}(\cdot|x) = \pi_\theta(\cdot|x), \\ \pi^{s2}(\cdot|x) = \pi_\theta(\cdot|x), \end{cases} \quad \textcircled{2} \begin{cases} \pi^{s1}(\cdot|x) \propto \pi_\theta^{3/2}(\cdot|x) \pi_{\text{ref}}^{-1/2}(\cdot|x), \\ \pi^{s2}(\cdot|x) \propto \pi_\theta^{1/2}(\cdot|x) \pi_{\text{ref}}^{1/2}(\cdot|x), \end{cases}$$

and with a reward margin $r_{\max} \in \mathbb{R}_+$ the mixing ratio can be roughly approximated as

$$\textcircled{1} : \textcircled{2} = 2 : (\exp(r_{\max}) + \exp(-r_{\max})) . \quad (8)$$

Logit mixing. The proposed samplers involve a hybridization between two policies, and a common approach to approximate hybrid distributions is *logit mixing* [14, 23]. Here we show how to understand this point in a theoretically sound way. Given $\pi_1, \pi_2 \in \Pi$, $w_1, w_2 \in \mathbb{R}$, we consider a new logit as $\zeta := w_1\zeta_1 + w_2\zeta_2$, where ζ_1, ζ_2 represent the per-token logits of policies π_1, π_2 , namely $\zeta_k(y_t|x, y_{<t}) = \log \pi_k(y_t|x, y_{<t})$. Note that

$$\begin{aligned} \operatorname{argmax}_{y \in \mathcal{Y}} \pi_1^{w_1}(y|x)\pi_2^{w_2}(y|x) &= \operatorname{argmax}_{y \in \mathcal{Y}} w_1 \log \pi_1(y|x) + w_2 \log \pi_2(y|x) \\ &= \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{t=0}^{|y|} w_1 \zeta_1(y_t|x, y_{<t}) + w_2 \zeta_2(y_t|x, y_{<t}) \\ &= \operatorname{argmax}_{y \in \mathcal{Y}} \sum_{t=0}^{|y|} \zeta(y_t|x, y_{<t}) . \end{aligned}$$

This indicates that, greedy decoding from $\pi \propto \pi_1^{w_1} \pi_2^{w_2}$ is equivalent to greedy decoding from $w_1\zeta_1 + w_2\zeta_2$. Thus, our proposed samplers can be implemented through mixing the logits of π_{ref} and π_θ .

Understanding existing approaches. Vanilla DPO [21] and its online variant [29] can be incorporated into our theoretical framework. As shown in Table 1, vanilla DPO, which assumes that pair-comparison data are sampled from π_{ref} (see Section 4 of Rafailov et al. [21]), can be viewed as DPO-UNIF; On-policy DPO [10, 26] proposes to sample response pairs using π_θ , and is thus equivalent to DPO-UNIF; Hybrid GSHF (Option I in [29]) sets $\pi^{s1} = \pi_\theta$ and $\pi^{s2} = \pi_{\text{ref}}$, equivalent to DPO-MIX-P (2); and Online GSHF (Option II in [29]) adopts the best/worst-of- K response generated by π_θ , which can be approximately viewed as generating from $\pi_\theta(\cdot) \exp(r(\cdot)/\beta)$ and $\pi_\theta(\cdot) \exp(-r(\cdot)/\beta)$, *i.e.* DPO-MIX-R (2). Notably, the 1 part is often omitted in DPO variants, and it can be attributed to the infinitely large reward margin in the implementation [8, 29], making the mixing ratio $\rightarrow 0 : 1$ in Equation (8) (see more details in Section 4 of Rosset et al. [22] and Appendix F).

Table 1: **Comparison with existing approaches.** We find that many baselines can be mapped to components of our proposed samplers, offering an alternative explanation for their advantages.

Algorithm	Practical π^{s1}	Practical π^{s2}	Equivalent Sampler	Posterior Distribution
Vanilla DPO	π_{ref}	π_{ref}	DPO-UNIF	$\pi_{\text{ref}}^{2\beta}$
On-policy DPO	π_θ	π_θ	DPO-UNIF	$\pi_\theta^{2\beta}$
Hybrid GSHF	π_θ	π_{ref}	DPO-MIX-P (2)	$\pi_\theta^\beta \pi_{\text{ref}}^\beta$
Online GSHF	π_θ (best-of- K)	π_θ (worst-of- K)	DPO-MIX-R (2)	$\pi_\theta^{2\beta}$
Ours	$\frac{\pi_\theta^{3/2}}{\pi_\theta} \frac{\pi_{\text{ref}}^{-1/2}}{\pi_{\text{ref}}}$	$\frac{\pi_\theta^{1/2}}{\pi_\theta} \frac{\pi_{\text{ref}}^{1/2}}{\pi_{\text{ref}}}$	DPO-MIX-P	$\pi_\theta^{2\beta}$

Table 2: **Results on Safe-RLHF.** The average reward is scored by the gold reward model on train set and test set, and win-rate is against the reference model. Each algorithm is trained for 3 iterations, and in the final iteration, ours shows advantages over baselines across all metrics.

Algorithm	Iters	Average reward (train)	Win-rate (train)	Average reward (test)	Win-rate (test)
Vanilla DPO	2	-1.486	67.6%	-1.423	68.7%
	3	-1.144	72.5%	-1.203	71.7%
On-policy DPO	2	-1.478	67.6%	-1.510	65.8%
	3	-1.082	73.2%	-1.094	73.2%
Hybrid GSHF	2	-1.517	68.5%	-1.505	66.9%
	3	-1.079	74.8%	-1.002	75.9%
Ours	2	-1.457	68.1%	-1.436	67.6%
	3	-0.908	75.6%	-0.945	76.2%

Table 3: **Results on Iterative-Prompt.** The average reward is scored by the gold reward model on train set and test set, and win-rate is against the reference model. Each algorithm is trained for 3 iterations, and in the final iteration, ours shows advantages over baselines across all metrics.

Algorithm	Iters	Average reward (train)	Win-rate (train)	Average reward (test)	Win-rate (test)
Vanilla DPO	2	1.427	71.4%	1.375	70.0%
	3	2.023	78.4%	2.133	78.8%
On-policy DPO	2	2.106	79.2%	2.157	78.7%
	3	3.131	82.4%	3.327	82.9%
Hybrid GSHF	2	2.116	79.6%	2.224	80.0%
	3	2.386	81.9%	2.500	82.8%
Ours	2	2.026	78.3%	2.068	77.3%
	3	4.149	86.6%	4.221	87.1%

C.2. Alignment Experiments

Experiment setup. We conduct experiments on two datasets, Safe-RLHF [11] and Iterative-Prompt [8, 29]. Our pipeline is mainly borrowed from Dong et al. [8]. For each iteration, responses are generated for a fixed set of prompts. Specifically, given prompt x , we generate $y_1 \sim \pi^{s1}(\cdot|x)$ and $y_2 \sim \pi^{s2}(\cdot|x)$. Each generated pair is annotated by a gold reward model [7] as (r_1, r_2) , and the corresponding loss is

$$\begin{aligned} \mathcal{L}_{(y_1, y_2)}(\theta) = & -\sigma(r_{\max} \cdot (r_1 - r_2)) \log \sigma \left(\beta \log \frac{\pi_{\theta}(y_1|x)\pi_{\text{ref}}(y_2|x)}{\pi_{\text{ref}}(y_1|x)\pi_{\theta}(y_2|x)} \right) \\ & - \sigma(r_{\max} \cdot (r_2 - r_1)) \log \sigma \left(\beta \log \frac{\pi_{\theta}(y_2|x)\pi_{\text{ref}}(y_1|x)}{\pi_{\text{ref}}(y_2|x)\pi_{\theta}(y_1|x)} \right), \end{aligned}$$

where $r_{\max} \in \mathbb{R}_+$ is the reward margin. See more details in Appendix F.

Results. Experimental results on LM alignment are provided in Tables 2 and 3. On Safe-RLHF dataset, our method is 4.5% better than vanilla DPO and 3.0% better than on-policy DPO. On Iterative-Prompt dataset, ours improves by 8.3% compared to vanilla DPO and by 4.2% compared to on-policy DPO. We also show the reward-KL curves in Figure 2, to indicate that the tuned models do not deviate much from π_{ref} .

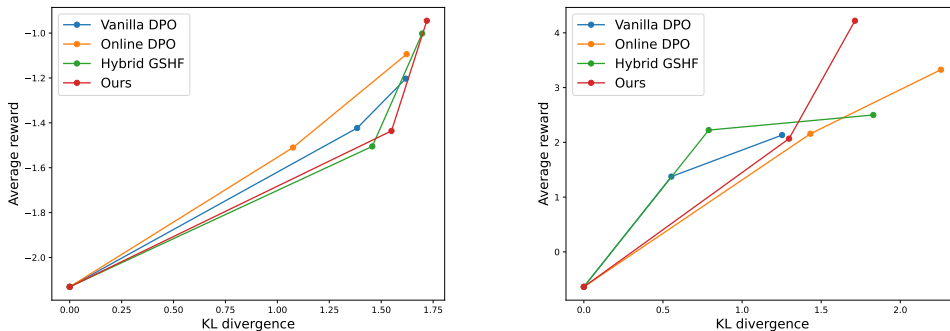


Figure 2: **The Reward-KL curves.** The left figure illustrates results on Safe-RLHF, and the right one illustrates results on Iterative-Prompt. The KL-divergence is measured on a subset of prompts in the test set. The results indicate that the KL-divergence of trained models does not deviate much from the reference model, and our method performs best in balancing reward and KL-divergence.

Clarification on evaluations. It is not enough to only show the results scored by reward models, since DPO algorithm is not explicitly learning the reward rankings [4, 17]. Due to restricted resources, we have not evaluated on open-benchmarks [9, 32]. Our work has demonstrated the potential to train models more effectively with minimal changes to the existing DPO pipeline. We hope this will inspire the community, especially those with rich computational resources, to conduct more systematic experiments.

Appendix D. Proofs of Convergence Rates of Exact DPO

Without loss of generality, we assume π_{ref} to be uniform distribution throughout this section. In the main text, we use \mathcal{Y} to represent the action space and y to represent an action for compatibility with other LM papers. From here, we turn back to \mathcal{A} for action space, a for an action, and A for the size of \mathcal{A} since all the proofs are conducted in bandit environments. And for notational ease, we make the following definitions:

$$\begin{aligned}\Delta(a, a'; \theta) &:= \sigma(r(a) - r(a')) - \sigma(\beta(\theta_a - \theta_{a'})) , \\ \delta(a, a'; \theta) &:= r(a) - r(a') - \beta(\theta_a - \theta_{a'}) .\end{aligned}$$

D.1. Theorems 1 and 2: Linear Convergence of Exact DPO–Unif

D.1.1. PROOF OF UPPER BOUND

For DPO with uniform sampler on action pairs, we first claim that for any θ appearing in the optimization process,

$$\max_{a, a'} \{\beta(\theta_a - \theta_{a'})\} \leq R_{\max} ,$$

where R_{\max} will be bounded later, and let $\sigma'_{\min} := \sigma'(R_{\max}) = \sigma(R_{\max})\sigma(-R_{\max})$. Then we have

$$\sigma'_{\min} \leq \frac{\sigma(x) - \sigma(y)}{x - y} \leq \frac{1}{4} \text{ when } |x|, |y| \leq R_{\max} \text{ and } x \neq y , \quad (9)$$

$$\mathcal{L}(\theta) = -\frac{2}{A^2} \sum_{a, a'} p^*(a > a') \log \sigma \left(\beta \log \frac{\pi_\theta(a)}{\pi_\theta(a')} \right) , \quad (10)$$

$$\nabla_{\theta} \mathcal{L}(\theta) = -\frac{2\beta}{A^2} \sum_{a, a'} \Delta(a, a'; \theta) \mathbb{1}_a . \quad (11)$$

Equation (11) reduces to

$$\nabla_{\theta_a} \mathcal{L}(\theta) = -\frac{2\beta}{A^2} \sum_{a'} \Delta(a, a'; \theta) .$$

Thus for any action pair (a, a') ,

$$\begin{aligned}(\theta_a - \theta_{a'})^{(t+1)} &= (\theta_a - \theta_{a'})^{(t)} + \frac{2\eta\beta\alpha(\pi^{s1}, \pi^{s2})}{A^2} \sum_{a''} \left(\Delta(a, a''; \theta^{(t)}) - \Delta(a', a''; \theta^{(t)}) \right) \\ &= (\theta_a - \theta_{a'})^{(t)} + 4\eta\beta \sum_{a''} \left(\Delta(a, a''; \theta^{(t)}) - \Delta(a', a''; \theta^{(t)}) \right) .\end{aligned}$$

At time t , sort the actions in the order that $r(a_i) - \beta\theta_{a_i}^{(t)} \leq r(a_{i+1}) - \beta\theta_{a_{i+1}}^{(t)}$. Then we have $\Delta(a_i, a_j; \theta^{(t)}) \geq 0$ if $i > j$. Note that it is possible that the order of actions at time $t+1$ is different, and in the following proof for any index i , a_i is from the order at time t . Let $l < r$, then

$$\delta(a_r, a_l; \theta^{(t+1)})$$

$$\begin{aligned}
 &= \delta(a_r, a_l; \theta^{(t)}) - 4\eta\beta^2 \sum_{i=1}^A \left(\Delta(a_r, a_i; \theta^{(t)}) - \Delta(a_l, a_i; \theta^{(t)}) \right) \\
 &\stackrel{(i)}{\leq} \delta(a_r, a_l; \theta^{(t)}) - 4\eta\beta^2 \sum_{i=1}^{l-1} \left(\sigma'_{\min} \delta(a_r, a_i; \theta^{(t)}) - \frac{1}{4} \delta(a_l, a_i; \theta^{(t)}) \right) \\
 &\quad - 4\eta\beta^2 \sum_{i=l}^r \left(\sigma'_{\min} \delta(a_r, a_i; \theta^{(t)}) - \sigma'_{\min} \delta(a_l, a_i; \theta^{(t)}) \right) - 4\eta\beta^2 \sum_{i=r+1}^A \left(\frac{1}{4} \delta(a_r, a_i; \theta^{(t)}) - \sigma'_{\min} \delta(a_l, a_i; \theta^{(t)}) \right) \\
 &= \delta(a_r, a_l; \theta^{(t)}) - 4\eta\beta^2 \left[\sigma'_{\min} (l-1) \delta(a_r, a_l; \theta^{(t)}) - \left(\frac{1}{4} - \sigma'_{\min} \right) \sum_{i=1}^{l-1} \delta(a_l, a_i; \theta^{(t)}) \right] \\
 &\quad - 4\eta\beta^2 \sigma'_{\min} (r-l+1) \delta(a_r, a_l; \theta^{(t)}) - 4\eta\beta^2 \left[\sigma'_{\min} (A-r) \delta(a_r, a_l; \theta^{(t)}) - \left(\frac{1}{4} - \sigma'_{\min} \right) \sum_{i=r+1}^A \delta(a_i, a_r; \theta^{(t)}) \right] \\
 &= (1 - 4\eta\beta^2 A \sigma'_{\min}) \delta(a_r, a_l; \theta^{(t)}) + 4\eta\beta^2 \left(\frac{1}{4} - \sigma'_{\min} \right) \left(\sum_{i=1}^{l-1} \delta(a_l, a_i; \theta^{(t)}) + \sum_{i=r+1}^A \delta(a_i, a_r; \theta^{(t)}) \right),
 \end{aligned}$$

where (i) is by using Equation (9) for different cases of x and y and whether $x - y > 0$. Similarly, for the lower bound:

$$\begin{aligned}
 &- \delta(a_r, a_l; \theta^{(t+1)}) \\
 &= 4\eta\beta^2 \sum_{i=1}^A \left(\Delta(a_r, a_i; \theta^{(t)}) - \Delta(a_l, a_i; \theta^{(t)}) \right) - \delta(a_r, a_l; \theta^{(t)}) \\
 &\leq 4\eta\beta^2 \sum_{i=1}^{l-1} \left(\frac{1}{4} \delta(a_r, a_i; \theta^{(t)}) - \sigma'_{\min} \delta(a_l, a_i; \theta^{(t)}) \right) + 4\eta\beta^2 \sum_{i=l}^r \left(\frac{1}{4} \delta(a_r, a_i; \theta^{(t)}) - \frac{1}{4} \delta(a_l, a_i; \theta^{(t)}) \right) \\
 &\quad + 4\eta\beta^2 \sum_{i=r+1}^A \left(\sigma'_{\min} \delta(a_r, a_i; \theta^{(t)}) - \frac{1}{4} \delta(a_l, a_i; \theta^{(t)}) \right) - \delta(a_r, a_l; \theta^{(t)}) \\
 &= 4\eta\beta^2 \left[\frac{1}{4} (l-1) \delta(a_r, a_l; \theta^{(t)}) + \left(\frac{1}{4} - \sigma'_{\min} \right) \sum_{i=1}^{l-1} \delta(a_l, a_i; \theta^{(t)}) \right] + 4\eta\beta^2 \cdot \frac{1}{4} (r-l+1) \delta(a_r, a_l; \theta^{(t)}) \\
 &\quad + 4\eta\beta^2 \left[\frac{1}{4} (A-r) \delta(a_r, a_l; \theta^{(t)}) + \left(\frac{1}{4} - \sigma'_{\min} \right) \sum_{i=r+1}^A \delta(a_i, a_r; \theta^{(t)}) \right] - \delta(a_r, a_l; \theta^{(t)}) \\
 &= (\eta\beta^2 A - 1) \delta(a_r, a_l; \theta^{(t)}) + 4\eta\beta^2 \left(\frac{1}{4} - \sigma'_{\min} \right) \left(\sum_{i=1}^{l-1} \delta(a_l, a_i; \theta^{(t)}) + \sum_{i=r+1}^A \delta(a_i, a_r; \theta^{(t)}) \right).
 \end{aligned}$$

Now taking $\eta = \frac{1}{\beta^2 A}$, then we have

$$\begin{aligned}
 \delta(a_r, a_l; \theta^{(t+1)}) &\leq (2 - 8\sigma'_{\min}) \max_{a, a'} \delta(a, a'; \theta^{(t)}), \\
 -\delta(a_r, a_l; \theta^{(t+1)}) &\leq (1 - 4\sigma'_{\min}) \max_{a, a'} \delta(a, a'; \theta^{(t)}).
 \end{aligned}$$

Define

$$\gamma := 2 - 8\sigma'_{\min}$$

as the contraction factor, then

$$\left| \delta(a_r, a_l; \theta^{(t+1)}) \right| \leq \gamma \max_{a, a'} \left| \delta(a, a'; \theta^{(t)}) \right|. \quad (12)$$

Recall that we initialize $\theta^{(0)} = \vec{0}$. Next we use induction to verify that throughout the process ($t \geq 0$),

$$\left| \delta(a_r, a_l; \theta^{(t+1)}) \right| \leq 0.214\gamma^t, \quad \text{and} \quad \left| \beta(\theta_a - \theta_{a'})^{(t+1)} \right| < 1.214. \quad (13)$$

For time $t = 0$, we have special versions: $r(a_1) \leq r(a_2) \leq \dots \leq r(a_A)$.

$$\begin{aligned} \delta(a_r, a_l; \theta^{(1)}) &= r(a_r) - r(a_l) - 4\eta\beta^2 \sum_{i=1}^A (\Delta(a_r, a_i; \theta^{(0)}) - \Delta(a_l, a_i; \theta^{(0)})) \\ &\stackrel{(i)}{=} r(a_r) - r(a_l) - 4\eta\beta^2 \sum_{i=1}^A (\sigma(r(a_r) - r(a_i)) - \sigma(r(a_l) - r(a_i))) \\ &\stackrel{(ii)}{\leq} r(a_r) - r(a_l) - 4\eta\beta^2 \sum_{i=1}^A \sigma'(1) [r(a_r) - r(a_i) - (r(a_l) - r(a_i))] \\ &= (1 - 4\eta\beta^2 A \sigma'(1)) (r(a_r) - r(a_l)) \\ &\stackrel{(iii)}{\leq} 0.214; \\ -\delta(a_r, a_l; \theta^{(1)}) &= 4\eta\beta^2 \sum_{i=1}^A (\Delta(a_r, a_i; \theta^{(0)}) - \Delta(a_l, a_i; \theta^{(0)})) - (r(a_r) - r(a_l)) \\ &= 4\eta\beta^2 \sum_{i=1}^A (\sigma(r(a_r) - r(a_i)) - \sigma(r(a_l) - r(a_i))) - (r(a_r) - r(a_l)) \\ &\leq 4\eta\beta^2 \sum_{i=1}^A \frac{1}{4} [r(a_r) - r(a_i) - (r(a_l) - r(a_i))] - (r(a_r) - r(a_l)) \\ &= (\eta\beta^2 A - 1) (r(a_r) - r(a_l)) \\ &= 0, \end{aligned}$$

where (i) is by $\theta^{(0)} = \vec{0}$; (ii) is by Equation (9) and $r(a_r) - r(a_i) \geq r(a_l) - r(a_i)$; (iii) is by $r(a_r) - r(a_l) \leq 1$. So $|\delta(a_r, a_l; \theta^{(1)})| \leq 0.214$, and $|\beta(\theta_a - \theta_{a'})_1| \leq |r(a) - r(a')| + |\delta(a_r, a_l; \theta^{(1)})| \leq 1.214$. Suppose for time $t - 1$, Equation (13) holds, then Equation (12) holds. So for time t ,

$$\left| \delta(a_r, a_l; \theta^{(t+1)}) \right| \leq \gamma \max_{a, a'} \left| \delta(a_r, a_l; \theta^{(t)}) \right| \leq 0.214\gamma^t \leq 0.214,$$

and

$$\left| \beta(\theta_a - \theta_{a'})^{(t+1)} \right| \leq |r(a) - r(a')| + \left| \delta(a_r, a_l; \theta^{(t+1)}) \right| \leq 1.214.$$

Thus we have

$$\begin{aligned} \gamma &= 2 - 8\sigma'_{\min} \leq 2 - 8\sigma'(1.214) < 0.588, \\ \left| \delta(a_r, a_l; \theta^{(T)}) \right| &\leq 0.588^T. \end{aligned}$$

D.1.2. CONSTRUCTION OF LOWER BOUND

Consider a three-armed bandit setting with rewards $r(a_1) = 0, r(a_2) = 1/3, r(a_3) = 1$ and any regularization parameter $\beta \in \mathbb{R}_+$. The update rule satisfies:

$$\delta(a_2, a_1; \theta^{(t+1)}) = \delta(a_2, a_1; \theta^{(t)}) - 4\eta\beta^2 \left(2\Delta(a_2, a_1; \theta^{(t)}) + \Delta(a_3, a_1; \theta^{(t)}) - \Delta(a_3, a_2; \theta^{(t)}) \right), \quad (14)$$

$$\delta(a_3, a_2; \theta^{(t+1)}) = \delta(a_3, a_2; \theta^{(t)}) - 4\eta\beta^2 \left(2\Delta(a_3, a_2; \theta^{(t)}) + \Delta(a_3, a_1; \theta^{(t)}) - \Delta(a_2, a_1; \theta^{(t)}) \right), \quad (15)$$

$$\delta(a_3, a_1; \theta^{(t+1)}) = \delta(a_3, a_1; \theta^{(t)}) - 4\eta\beta^2 \left(2\Delta(a_3, a_1; \theta^{(t)}) + \Delta(a_3, a_2; \theta^{(t)}) + \Delta(a_2, a_1; \theta^{(t)}) \right).$$

Define $x_t := \delta(a_2, a_1; \theta^{(t)})$, and $y_t := \delta(a_3, a_2; \theta^{(t)})$. Clearly we have $\delta(a_3, a_1; \theta^{(t)}) = x_t + y_t$. We can perform Taylor expansion on Equations (14) and (15) and get

$$\begin{pmatrix} x_{t+1} \\ y_{t+1} \end{pmatrix} = \underbrace{\begin{pmatrix} 1 - 4\eta\beta^2(2\sigma'(1/3) + \sigma'(1)) & 4\eta(\sigma'(2/3) - \sigma'(1)) \\ 4\eta\beta^2(\sigma'(1/3) - \sigma'(1)) & 1 - 4\eta\beta^2(2\sigma'(2/3) + \sigma'(1)) \end{pmatrix}}_{:=B} \begin{pmatrix} x_t \\ y_t \end{pmatrix} + \eta\beta^2 \begin{pmatrix} u_t \\ v_t \end{pmatrix}, \quad (16)$$

where

$$|u_t| \leq \frac{4x_t^2 + 3y_t^2}{3\sqrt{3}} \leq x_t^2 + y_t^2, \quad |v_t| \leq \frac{3x_t^2 + 4y_t^2}{3\sqrt{3}} \leq x_t^2 + y_t^2. \quad (17)$$

Now we analyze the eigenvalues of B under three scenarios.

1. If

$$0 < \eta\beta^2 < \frac{1}{4(2\sigma'(1/3) + \sigma'(1))} \approx 0.366,$$

then we have

$$\det(\lambda I - B) = \lambda^2 - (B_{11} + B_{22})\lambda + \underbrace{B_{11}B_{22}}_{\leq (B_{11}+B_{22})^2/4} - \underbrace{B_{12}B_{21}}_{>0}.$$

2. If

$$\frac{1}{4(2\sigma'(1/3) + \sigma'(1))} \leq \eta\beta^2 \leq \frac{1}{4(2\sigma'(2/3) + \sigma'(1))} \approx 0.388,$$

then we have

$$\det(\lambda I - B) = \lambda^2 - (B_{11} + B_{22})\lambda + \underbrace{B_{11}B_{22}}_{\leq 0} - \underbrace{B_{12}B_{21}}_{>0}.$$

3. If

$$\frac{1}{4(2\sigma'(2/3) + \sigma'(1))} < \eta\beta^2 < \frac{1}{2(2\sigma'(1/3) + \sigma'(2/3))} \approx 0.704 ,$$

then we have

$$\det(\lambda I - B) = \lambda^2 - (B_{11} + B_{22})\lambda + \underbrace{B_{11}B_{22}}_{\leq (B_{11}+B_{22})^2/4} - \underbrace{B_{12}B_{21}}_{>0} .$$

Therefore B has two different eigenvalues $\lambda_1, \lambda_2 \in (-1, 0) \cup (0, 1)$, with normalized eigenvectors w_1, w_2 . Clearly $w_{ij} \in (-1, 0) \cup (0, 1)$, $\forall i, j \in \{1, 2\}$. Then we define $\lambda_{\max} := \max(|\lambda_1|, |\lambda_2|)$, $\lambda_{\min} := \min(|\lambda_1|, |\lambda_2|)$, Now perform basis transformation with new basis (w_1, w_2) . Thus Equation (16) can be rewritten as

$$\begin{pmatrix} p_{t+1} \\ q_{t+1} \end{pmatrix} = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \begin{pmatrix} p_t \\ q_t \end{pmatrix} + \begin{pmatrix} u'_t \\ v'_t \end{pmatrix} ,$$

Let w'_1, w'_2 be the inverse basis, and define $\alpha := \max_{i,j \in \{1,2\}} |w'_{ij}|$, and $\epsilon := \min(\lambda_{\min}, 1 - \lambda_{\max}) / (64\alpha^2)$.

Now initialize $|x_0|, |y_0| \in (0, \epsilon)$. Then we have $\max_{i,j \in \{1,2,3\}} |\delta(a_i, a_j; \theta^{(0)})| \leq 2\epsilon$. Therefore

$$|p_0|, |q_0| \stackrel{(i)}{\leq} 2\alpha\epsilon ,$$

and

$$\begin{aligned} |u'_t|, |v'_t| &\stackrel{(ii)}{\leq} 2\alpha(x_t^2 + y_t^2) \\ &\stackrel{(iii)}{\leq} 4\alpha(p_t^2 + q_t^2) . \end{aligned}$$

(i) and (ii) comes from the fact that $p_t = w'_{11}x_t + w'_{21}y_t$ and $q_t = w'_{12}x_t + w'_{22}y_t$, and Equation (17); (iii) is from the fact that $x_t = w_{11}p_t + w_{21}q_t$ and $y_t = w_{12}p_t + w_{22}q_t$, and Cauchy-Schwarz inequality. Now we have

$$\begin{aligned} |p_{t+1}| + |q_{t+1}| &\leq [\lambda_{\max} + 8\alpha(|p_t| + |q_t|)] (|p_t| + |q_t|) \\ &\stackrel{(iv)}{\leq} (\lambda_{\max} + 32\alpha^2\epsilon) (|p_t| + |q_t|) \\ &\leq \frac{1 + \lambda_{\max}}{2} (|p_t| + |q_t|) . \\ |p_{t+1}| + |q_{t+1}| &\geq [\lambda_{\min} - 8\alpha(|p_t| + |q_t|)] (|p_t| + |q_t|) \\ &\stackrel{(v)}{\geq} (\lambda_{\min} - 32\alpha^2\epsilon) (|p_t| + |q_t|) \\ &\geq \frac{\lambda_{\min}}{2} (|p_t| + |q_t|) , \end{aligned}$$

where (iv) and (v) are based on simple induction that $|p_t| + |q_t|$ will not increase. And it thus indicates that $\max(|x_t|, |y_t|)$ can at most be linear convergence when $\eta\beta^2 \leq \frac{2}{A} \approx 0.667$.

D.2. Theorem 3: Quadratic Convergence of Exact DPO-Mix-R

We study DPO with a mixture of fixed samplers: $Z^+ Z^- \cdot \pi^{s1} \times \pi^{s2} + A^2 \cdot \text{Uniform}(\mathcal{A}) \times \text{Uniform}(\mathcal{A})$, where $Z^+ = \sum_a \exp(r(a))$, $\pi^{s1}(a) = \exp(r(a))/Z^+$ and $Z^- = \sum_a \exp(-r(a))$, $\pi^{s2}(A) = \exp(-r(a))/Z^-$. We have

$$\begin{aligned}
 & \alpha_1 \mathcal{L}_1(\theta) + \alpha_2 \mathcal{L}_2(\theta) \\
 &= - \sum_{a,a'} \left(A^2 \cdot \frac{1}{A^2} + Z^+ Z^- \cdot \pi^{s1}(a) \pi^{s2}(a') \right) \left[p^*(a > a') \log \sigma \left(\beta \log \frac{\pi_\theta(a)}{\pi_\theta(a')} \right) + p^*(a' > a) \log \sigma \left(\beta \log \frac{\pi_\theta(a')}{\pi_\theta(a)} \right) \right] \\
 &= - \sum_{a,a'} (\exp(r(a) - r(a')) + 1) \left[p^*(a > a') \log \sigma \left(\beta \log \frac{\pi_\theta(a)}{\pi_\theta(a')} \right) + p^*(a' > a) \log \sigma \left(\beta \log \frac{\pi_\theta(a')}{\pi_\theta(a)} \right) \right], \\
 & \nabla_\theta (\alpha_1 \mathcal{L}_1(\theta) + \alpha_2 \mathcal{L}_2(\theta)) \\
 &= -\beta \sum_{a,a'} (\exp(r(a) - r(a')) + 1) \Delta(a, a'; \theta) (\mathbb{1}_a - \mathbb{1}_{a'}) \\
 &= -\beta \sum_{a,a'} (\exp(r(a) - r(a')) + \exp(r(a') - r(a)) + 2) \Delta(a, a'; \theta) \mathbb{1}_a \\
 &= -\beta \sum_{a,a'} \frac{\Delta(a, a'; \theta)}{\sigma'(r(a) - r(a'))} \mathbb{1}_a. \tag{18}
 \end{aligned}$$

Equation (18) reduces to

$$\nabla_{\theta_a} (\alpha_1 \mathcal{L}_1(\theta) + \alpha_2 \mathcal{L}_2(\theta)) = -\beta \sum_{a'} \frac{\Delta(a, a'; \theta)}{\sigma'(r(a) - r(a'))}.$$

Fix parameter θ . For any action pair a, a' , through Taylor expansion we have that

$$\Delta(a, a'; \theta) = \sigma'(r(a) - r(a')) \delta(a, a'; \theta) - \frac{\sigma''(\xi_R(a, a'; \theta))}{2} \delta(a, a'; \theta)^2,$$

where $\xi_R(a, a'; \theta)$ is between $r(a) - r(a')$ and $\beta(\theta_a - \theta_{a'})$. We have that at time step t , for any action pair (a, a') ,

$$\begin{aligned}
 \delta(a, a'; \theta^{(t+1)}) &= \delta(a, a'; \theta^{(t)}) - \eta \beta^2 \sum_{a''} \left(\frac{\Delta(a, a''; \theta^{(t)})}{\sigma'(r(a) - r(a''))} - \frac{\Delta(a', a''; \theta^{(t)})}{\sigma'(r(a') - r(a''))} \right) \\
 &= \delta(a, a'; \theta^{(t)}) - \eta \beta^2 \sum_{a''} (\delta(a, a''; \theta^{(t)}) - \delta(a', a''; \theta^{(t)})) \\
 &\quad + \frac{\eta \beta^2}{2} \sum_{a''} \left(\frac{\sigma''(\xi_R(a, a''; \theta^{(t)}))}{\sigma'(r(a) - r(a''))} \delta(a, a''; \theta^{(t)})^2 - \frac{\sigma''(\xi_R(a', a''; \theta^{(t)}))}{\sigma'(r(a') - r(a''))} \delta(a', a''; \theta^{(t)})^2 \right) \\
 &= (1 - \eta \beta^2 A) \delta(a, a'; \theta^{(t)}) \\
 &\quad + \frac{\eta \beta^2}{2} \sum_{a''} \left(\frac{\sigma''(\xi_R(a, a''; \theta^{(t)}))}{\sigma'(r(a) - r(a''))} \delta(a, a''; \theta^{(t)})^2 - \frac{\sigma''(\xi_R(a', a''; \theta^{(t)}))}{\sigma'(r(a') - r(a''))} \delta(a', a''; \theta^{(t)})^2 \right).
 \end{aligned}$$

From the range of r , we know that $\sigma'(r(a) - r(a')) \geq \sigma'(1) > 0.196$. We have $|\sigma''(\xi_R(a, a'; \theta^{(t)}))| \leq \sigma''_{\max} := \sup_{0 \leq x \leq 1} x(1-x)(1-2x) = 1/(6\sqrt{3}) < 0.097$. Set

$$\eta = \frac{1}{\beta^2 A},$$

then

$$\begin{aligned} |\delta(a, a'; \theta^{(t+1)})| &\leq \frac{1}{2A} \sum_{a''} \left(\frac{\sigma''_{\max}}{\sigma'(1)} \delta(a, a''; \theta^{(t)})^2 + \frac{\sigma''_{\max}}{\sigma'(1)} \delta(a', a''; \theta^{(t)})^2 \right) \\ &\leq \frac{\sigma''_{\max}}{\sigma'(1)} \max_{a, a'} \delta(a, a'; \theta^{(t)})^2 \\ &< \frac{1}{2} \max_{a, a'} \delta(a, a'; \theta^{(t)})^2. \end{aligned}$$

Since $\max_{a, a'} |\delta(a, a'; \theta^{(0)})| \leq 1$, we can show a quadratic convergence for this regime:

$$|\delta(a, a'; \theta^t)| \leq 0.5^{2^t - 1}.$$

D.3. Theorem 4: Quadratic Convergence of Exact DPO-Mix-P

We study DPO with a mixture of on-policy samplers (with gradient stopped) and uniform samplers: $Z^+ Z^- \cdot \pi^{s1} \times \pi^{s2} + A^2 \cdot \text{Uniform}(\mathcal{A}) \times \text{Uniform}(\mathcal{A})$, where $Z^+ = \sum_a \exp(\beta\theta_a)$, $\pi^{s1}(a) = \exp(\beta\theta_a)/Z^+$ and $Z^- = \sum_a \exp(-\beta\theta_a)$, $\pi^{s2}(a) = \exp(-\beta\theta_a)/Z^-$. Samely we have

$$\nabla_{\theta_a} (\alpha_1 \mathcal{L}_1(\theta) + \alpha_2 \mathcal{L}_2(\theta)) = -\beta \sum_{a'} \frac{\Delta(a, a'; \theta)}{\sigma'(\beta(\theta_a - \theta_{a'}))}.$$

Fix parameter θ . For any action pair a, a' , through Taylor expansion we have that

$$\Delta(a, a'; \theta) = \sigma'(\beta(\theta_a - \theta_{a'})) \delta(a, a'; \theta) + \frac{\sigma''(\xi_P(a, a'; \theta))}{2} \delta(a, a'; \theta)^2,$$

where $\xi_P(a, a'; \theta)$ is between $r(a) - r(a')$ and $\beta(\theta_a - \theta_{a'})$. We have that at time step t , for any action pair (a, a') ,

$$\begin{aligned} \delta(a, a'; \theta^{(t+1)}) &= \delta(a, a'; \theta^{(t)}) - \eta \beta^2 \sum_{a''} \left(\frac{\Delta(a, a''; \theta^{(t)})}{\sigma'(\beta(\theta_a - \theta_{a''})^{(t)})} - \frac{\Delta(a', a''; \theta^{(t)})}{\sigma'(\beta(\theta_{a'} - \theta_{a''})^{(t)})} \right) \\ &= \delta(a, a'; \theta^{(t)}) - \eta \beta^2 \sum_{a''} (\delta(a, a''; \theta^{(t)}) - \delta(a', a''; \theta^{(t)})) \\ &\quad - \frac{\eta \beta^2}{2} \sum_{a''} \left(\frac{\sigma''(\xi_P(a, a''; \theta^{(t)}))}{\sigma'(\beta(\theta_a - \theta_{a''})^{(t)})} \delta(a, a''; \theta^{(t)})^2 - \frac{\sigma''(\xi_P(a', a''; \theta^{(t)}))}{\sigma'(\beta(\theta_{a'} - \theta_{a''})^{(t)})} \delta(a', a''; \theta^{(t)})^2 \right) \\ &= (1 - \eta \beta^2 A) \delta(a, a'; \theta^{(t)}) \\ &\quad - \frac{\eta \beta^2}{2} \sum_{a''} \left(\frac{\sigma''(\xi_P(a, a''; \theta^{(t)}))}{\sigma'(\beta(\theta_a - \theta_{a''})^{(t)})} \delta(a, a''; \theta^{(t)})^2 - \frac{\sigma''(\xi_P(a', a''; \theta^{(t)}))}{\sigma'(\beta(\theta_{a'} - \theta_{a''})^{(t)})} \delta(a', a''; \theta^{(t)})^2 \right). \end{aligned}$$

We still first claim that $\sigma'(\beta(\theta_a - \theta_{a'})_t) \geq \sigma'_{\min}$, and will bound it later. We have $|\sigma''(\xi_{\mathcal{P}}(a, a''; \theta^{(t)}))| \leq \sigma''_{\max} < 0.097$. Set

$$\eta = \frac{1}{\beta^2 A},$$

then

$$\begin{aligned} |\delta(a, a'; \theta^{(t+1)})| &\leq \frac{\sigma''_{\max}}{2A\sigma'_{\min}} \sum_{a''} (\delta(a, a''; \theta^{(t)})^2 + \delta(a', a''; \theta^{(t)})^2) \\ &\leq \frac{\sigma''_{\max}}{\sigma'_{\min}} \max_{a, a'} \delta(a, a'; \theta^{(t)})^2. \end{aligned}$$

At time step $t = 0$ we have $\sigma'(\beta(\theta_a - \theta_{a'})^{(0)}) = \sigma'(0) = 0.25$ and $\max_{a, a'} |\delta(a, a'; \theta^{(0)})| \leq 1$, so

$$\max_{a, a'} |\delta(a, a'; \theta^{(1)})| < 0.388.$$

By simple induction, we have that

$$\begin{aligned} \sigma'_{\min} &\geq \sigma'(1 + \max_{a, a'} |\delta(a, a'; \theta^{(t)})|) \geq \sigma'(1.388) > 0.159, \\ \max_{a, a'} |\delta(a, a'; \theta^{(t+1)})| &\leq \frac{0.097}{0.159} \max_{a, a'} \delta(a, a'; \theta^{(t)})^2 < 0.611 \max_{a, a'} \delta(a, a'; \theta^{(t)})^2. \end{aligned}$$

which is a quadratic convergence:

$$|\delta(a, a'; \theta^t)| \leq 0.611^{2^t - 1}.$$

Appendix E. Proof of Convergence Rates of Empirical DPO

For notational ease, we make the following definitions throughout this section:

$$\begin{aligned} \Delta(a, a'; \theta) &:= \sigma(r(a) - r(a')) - \sigma(\beta(\theta_a - \theta_{a'})), \\ \delta(a, a'; \theta) &:= r(a) - r(a') - \beta(\theta_a - \theta_{a'}). \end{aligned}$$

This section conforms to Definition 2. Denote the filtration \mathcal{F}_t as all the samples on and before time step t .

E.1. Technical Lemma

Lemma 2 (Lemma 1.4 in Philippe Rigollet [20]) *Let X be a random variable such that*

$$\mathbb{P}[|X| > t] \leq 2 \exp\left(-\frac{t^2}{2\sigma^2}\right),$$

then for any positive integer $k \geq 2$,

$$\mathbb{E}[|X|^k] \leq (\sigma e^{1/e} \sqrt{k})^k,$$

and

$$\mathbb{E}[|X|] \leq \sigma \sqrt{2\pi}.$$

E.2. Theorem 5: Convergence of Empirical DPO-Mix-R

Similar to Appendix D.2, at time step t , conditioned on \mathcal{F}_t , we have that for any action pair (a, a') ,

$$\begin{aligned} \mathbb{E}[(G_a - G_{a'})^{(t)}] &= -\beta A \delta(a, a'; \theta^{(t)}) \\ &\quad - \underbrace{\frac{\beta}{2} \sum_{a''} \left(\frac{\sigma''(\xi_R(a, a''; \theta^{(t)}))}{\sigma'(r(a) - r(a''))} \delta(a, a''; \theta^{(t)})^2 - \frac{\sigma''(\xi_R(a', a''; \theta^{(t)}))}{\sigma'(r(a') - r(a''))} \delta(a', a''; \theta^{(t)})^2 \right)}_{=: N_t(a, a')}, \\ |N_t(a, a')| &< \frac{1}{2} \sum_{a''} (\delta(a, a''; \theta^{(t)})^2 + \delta(a', a''; \theta^{(t)})^2). \end{aligned} \quad (19)$$

From Definition 2 and Lemma 2, we have that

$$\mathbb{E} \left[\left| \frac{G_a^{(t)} - \mathbb{E}[G_a^{(t)}]}{\beta A} \right|^k \right] \leq (3\sigma\sqrt{k})^k.$$

Therefore, from Minkowski inequality,

$$\mathbb{E} \left[\left| \frac{(G_a - G_{a'})^{(t)} - \mathbb{E}[(G_a - G_{a'})^{(t)}]}{\beta A} \right|^k \right] \leq (6\sigma\sqrt{k})^k.$$

Now we take $\eta = 1/(\beta^2 A)$, then by taking expectation conditioning on \mathcal{F}_t we obtain

$$\begin{aligned} \mathbb{E}[\delta(a, a'; \theta^{(t+1)})^{2n}] &= \mathbb{E}[(\delta(a, a'; \theta^{(t)}) + \eta\beta(G_a - G_{a'}))^{2n}] \\ &= \mathbb{E}[(\delta(a, a'; \theta^{(t)}) + \eta\beta\mathbb{E}[G_a - G_{a'}] + \eta\beta(G_a - G_{a'} - \mathbb{E}[G_a - G_{a'}]))^{2n}] \\ &= \sum_{k=0}^{2n} \binom{2n}{k} (\delta(a, a'; \theta^{(t)}) + \eta\beta\mathbb{E}[G_a - G_{a'}])^{2n-k} \cdot (\eta\beta)^k \mathbb{E}[(G_a - G_{a'} - \mathbb{E}[G_a - G_{a'}])^k] \\ &\stackrel{(i)}{=} \sum_{k=0}^{2n} \binom{2n}{k} \left(-\frac{1}{2A} N_t(a, a') \right)^{2n-k} \cdot \frac{1}{(\beta A)^k} \mathbb{E}[(G_a - G_{a'} - \mathbb{E}[G_a - G_{a'}])^k] \\ &\leq \sum_{k=0}^{2n} \binom{2n}{k} \left(\frac{1}{2A} |N_t(a, a')| \right)^{2n-k} (6\sigma\sqrt{k})^k, \end{aligned}$$

where (i) is by substituting $\eta = 1/(\beta^2 A)$.

Further taking expectation over \mathcal{F}_t , we have

$$\begin{aligned} \mathbb{E}[\delta(a, a'; \theta^{(t+1)})^{2n}] &\leq \sum_{k=0}^{2n} \frac{\binom{2n}{k}}{(2A)^{2n-k}} (6\sigma\sqrt{k})^k \mathbb{E}[|N_t|^{2n-k}(a, a')] \\ &\stackrel{(i)}{\leq} \sum_{k=0}^{2n} \frac{\binom{2n}{k}}{(2A)^{2n-k}} (6\sigma\sqrt{k})^k \cdot \frac{1}{2^{2n-k}} \mathbb{E} \left[\left[\sum_{a''} (\delta(a, a''; \theta^{(t)})^2 + \delta(a', a''; \theta^{(t)})^2) \right]^{2n-k} \right] \end{aligned}$$

$$\begin{aligned}
 &\stackrel{(ii)}{\leq} \sum_{k=0}^{2n} \frac{\binom{2n}{k}}{(2A)^{2n-k}} (6\sigma\sqrt{k})^k \cdot \frac{1}{2^{2n-k}} \cdot (2A)^{2n-k-1} \sum_{a''} (\mathbb{E}[\delta(a, a''; \theta^{(t)})^{4n-2k}] + \mathbb{E}[\delta(a', a''; \theta^{(t)})^{4n-2k}]) \\
 &\leq \sum_{k=0}^{2n} \binom{2n}{k} (6\sigma\sqrt{k})^k \cdot \frac{1}{2^{2n-k}} \max_{a_1, a_2} \mathbb{E}[\delta(a_1, a_2; \theta^{(t)})^{4n-2k}] \\
 &\leq \sum_{k=0}^{2n} \binom{2n}{k} (6\sigma\sqrt{n})^k \cdot \frac{1}{2^{2n-k}} \max_{a_1, a_2} \mathbb{E}[\delta(a_1, a_2; \theta^{(t)})^{4n-2k}],
 \end{aligned}$$

where (i) is by Equation (19); (ii) is by Hölder inequality.

Take $T = \lceil \log(1/\sigma) \rceil$. When $\sigma \leq 1/576 < 0.00174$, we will show that $\forall n, t \in \mathbb{N}$ such that $n \cdot 2^t \leq 1/\sigma$,

$$\mathbb{E}[\delta(a, a'; \theta^{(t)})^{2n}] \leq \left(12\sqrt{n}\sigma + \frac{1}{2^t}\right)^{2n}.$$

This can be proved using induction on t . For $t \leq 1$, we have that for any n :

$$\begin{aligned}
 \mathbb{E}[\delta(a, a'; \theta^{(0)})^{2n}] &\leq 1, \\
 \mathbb{E}[\delta(a, a'; \theta^{(1)})^{2n}] &\leq \left(6\sqrt{n}\sigma + \frac{1}{2}\right)^{2n}.
 \end{aligned}$$

For $t = 2$ and $n \leq 1/(4\sigma)$,

$$\begin{aligned}
 \mathbb{E}[\delta(a, a'; \theta^{(2)})^{2n}] &\leq \sum_{k=0}^{2n} \binom{2n}{k} (6\sigma\sqrt{n})^k \cdot \frac{1}{2^{2n-k}} \left(6\sqrt{2n}\sigma + \frac{1}{2}\right)^{4n-2k} \\
 &\leq \left(6\sqrt{n}\sigma + \frac{(6\sqrt{2n}\sigma + \frac{1}{2})^2}{2}\right)^{2n} \\
 &= \left(36n\sigma^2 + (6 + 3\sqrt{2})\sqrt{n}\sigma + \frac{1}{8}\right)^{2n} \\
 &\stackrel{(i)}{\leq} \left(12\sqrt{n}\sigma + \frac{1}{2^2}\right)^{2n},
 \end{aligned}$$

where (i) is by plugging in the range of n and σ . Suppose the arguments holds for $t \geq 2$, then

$$\begin{aligned}
 \mathbb{E}[\delta(a, a'; \theta^{(t+1)})^{2n}] &\leq \sum_{k=0}^{2n} \binom{2n}{k} (6\sigma\sqrt{n})^k \cdot \frac{1}{2^{2n-k}} \left(12\sqrt{2n}\sigma + \frac{1}{2^t}\right)^{4n-2k} \\
 &= \left[6\sqrt{n}\sigma + \frac{(12\sqrt{2n}\sigma + \frac{1}{2^t})^2}{2}\right]^{2n} \\
 &\leq \left[\left(6 + \frac{12\sqrt{2}}{2^t}\right)\sqrt{n}\sigma + 144n\sigma^2 + \frac{1}{2^{2t+1}}\right]^{2n} \\
 &\stackrel{(i)}{\leq} \left[\left(6 + \frac{12\sqrt{2}}{2^t}\right)\sqrt{n}\sigma + \frac{288\sigma + \frac{1}{2^t}}{2^{t+1}}\right]^{2n}
 \end{aligned}$$

$$\stackrel{\text{(ii)}}{\leq} \left(12\sqrt{n}\sigma + \frac{1}{2^{t+1}} \right)^{2n},$$

where (i) is by $n \leq 1/(\sigma \cdot 2^t)$; (ii) is by $t \geq 2$ and range of σ .

Therefore, we have for $\sigma \leq 1/576$ and $T = \lfloor \log(1/\sigma) \rfloor > \log(1/\sigma) - 1$,

$$\sqrt{\mathbb{E}[\delta(a, a'; \theta^{(T)})^2]} \leq 12\sigma + \frac{1}{2^T} < 14\sigma.$$

E.3. Theorem 6: Convergence of Empirical DPO-Mix-P*

Here we use the joint probability weights $\psi(a, a') \propto \exp(z(a, a'))$ such that $z(a, a') = -z(a', a)$ and let $Z := \sum_{a, a'} \exp(z(a, a'))$:

$$\begin{aligned} & \alpha_1 \mathcal{L}_1(\theta) + \alpha_2 \mathcal{L}_2(\theta) \\ &= - \sum_{a, a'} \text{sg} \left(A^2 \cdot \frac{1}{A^2} + Z \cdot \psi(a, a') \right) \left[p^*(a > a') \log \sigma \left(\beta \log \frac{\pi_\theta(a)}{\pi_\theta(a')} \right) + p^*(a' > a) \log \sigma \left(\beta \log \frac{\pi_\theta(a')}{\pi_\theta(a)} \right) \right], \\ & \nabla_\theta (\alpha_1 \mathcal{L}_1(\theta) + \alpha_2 \mathcal{L}_2(\theta)) \\ &= -\beta \sum_{a, a'} (\exp(z(a, a')) + 1) \Delta(a, a'; \theta) (\mathbb{1}_a - \mathbb{1}_{a'}) \\ &= -\beta \sum_{a, a'} (\exp(z(a, a')) + \exp(-z(a, a')) + 2) \Delta(a, a'; \theta) \mathbb{1}_a \\ &= -\beta \sum_{a, a'} \frac{\Delta(a, a'; \theta)}{\sigma'(z(a, a'))} \mathbb{1}_a. \end{aligned} \tag{20}$$

Equation (20) reduces to

$$\nabla_{\theta_a} \mathcal{L}(\theta) = -\beta \sum_{a'} \frac{\Delta(a, a'; \theta)}{\sigma'(z(a, a'))}.$$

Fix parameter θ . For any action pair a, a' , through Taylor expansion we have that

$$\begin{aligned} \Delta(a, a'; \theta) &= (\sigma(r(a) - r(a')) - \sigma(z(a, a'))) - (\sigma(\beta(\theta_a - \theta_{a'})) - \sigma(z(a, a'))) \\ &= [\sigma'(z(a, a'))(r(a) - r(a') - z(a, a')) + \frac{\sigma''(\xi_1(a, a'; \theta))}{2}(r(a) - r(a') - z(a, a'))^2] \\ & \quad - \{\sigma'(z(a, a'))[\beta(\theta_a - \theta_{a'}) - z(a, a')] + \frac{\sigma''(\xi_2(a, a'; \theta))}{2}[\beta(\theta_a - \theta_{a'}) - z(a, a')]^2\} \\ &= \sigma'(z(a, a'))\delta(a, a'; \theta) + \frac{\sigma''(\xi_1(a, a'; \theta))}{2}(r(a) - r(a') - z(a, a'))^2 \\ & \quad - \frac{\sigma''(\xi_2(a, a'; \theta))}{2}[\beta(\theta_a - \theta_{a'}) - z(a, a')]^2, \end{aligned}$$

where $\xi_1(a, a'; \theta)$ is between $r(a) - r(a')$ and $z(a, a')$, and $\xi_2(a, a'; \theta)$ is between $z(a, a')$ and $\beta(\theta_a - \theta_{a'})$.

If we set

$$z(a, a') = \begin{cases} 1, & \text{if } \beta(\theta_a - \theta_{a'}) > 1, \\ -1, & \text{if } \beta(\theta_a - \theta_{a'}) < -1, \\ \beta(\theta_a - \theta_{a'}), & \text{otherwise,} \end{cases}$$

then we can conclude that

$$[r(a) - r(a') - z(a, a')]^2 + [\beta(\theta_a - \theta_{a'}) - z(a, a')]^2 \leq \delta(a, a'; \theta)^2.$$

Note that this construction satisfies $z(a, a') = -z(a', a)$. We have that at time step t , conditioning on \mathcal{F}_t , for any action pair (a, a') ,

$$\begin{aligned} & \mathbb{E}[\delta(a, a'; \theta^{(t+1)})] \\ &= \delta(a, a'; \theta^{(t)}) - \eta\beta^2 \sum_{a''} \left(\frac{\Delta(a, a''; \theta^{(t)})}{\sigma'(z(a, a''))} - \frac{\Delta(a', a''; \theta^{(t)})}{\sigma'(z(a', a''))} \right) \\ &= \delta(a, a'; \theta^{(t)}) - \eta\beta^2 \sum_{a''} (\delta(a, a''; \theta^{(t)}) - \delta(a', a''; \theta^{(t)})) \\ &\quad - \frac{\eta\beta^2}{2} \sum_{a''} \left\{ \frac{\sigma''(\xi_1(a, a''; \theta^{(t)}))}{\sigma'(z(a, a''))} (r(a) - r(a'') - z(a, a''))^2 - \frac{\sigma''(\xi_2(a, a''; \theta^{(t)}))}{\sigma'(z(a, a''))} [\beta(\theta_a - \theta_{a''})^{(t)} - z(a, a'')]^2 \right\} \\ &\quad + \frac{\eta\beta^2}{2} \sum_{a''} \left\{ \frac{\sigma''(\xi_1(a', a''; \theta^{(t)}))}{\sigma'(z(a', a''))} (r(a') - r(a'') - z(a', a''))^2 - \frac{\sigma''(\xi_2(a', a''; \theta^{(t)}))}{\sigma'(z(a', a''))} [\beta(\theta_{a'} - \theta_{a''})^{(t)} - z(a', a'')]^2 \right\} \\ &= (1 - \eta\beta^2 A) \delta(a, a'; \theta^{(t)}) \\ &\quad - \frac{\eta\beta^2}{2} \sum_{a''} \left\{ \frac{\sigma''(\xi_1(a, a''; \theta^{(t)}))}{\sigma'(z(a, a''))} (r(a) - r(a'') - z(a, a''))^2 - \frac{\sigma''(\xi_2(a, a''; \theta^{(t)}))}{\sigma'(z(a, a''))} [\beta(\theta_a - \theta_{a''})^{(t)} - z(a, a'')]^2 \right\} \\ &\quad + \frac{\eta\beta^2}{2} \sum_{a''} \left\{ \frac{\sigma''(\xi_1(a', a''; \theta^{(t)}))}{\sigma'(z(a', a''))} (r(a') - r(a'') - z(a', a''))^2 - \frac{\sigma''(\xi_2(a', a''; \theta^{(t)}))}{\sigma'(z(a', a''))} [\beta(\theta_{a'} - \theta_{a''})^{(t)} - z(a', a'')]^2 \right\}. \end{aligned}$$

Set

$$\eta = \frac{1}{\beta^2 A},$$

then

$$\begin{aligned} \mathbb{E} \left| \delta(a, a'; \theta^{(t+1)}) \right| &\leq \frac{\sigma''_{\max}}{2A\sigma'(1)} \sum_{a''} \{ (r(a) - r(a'') - z(a, a''))^2 + [\beta(\theta_a - \theta_{a''})^{(t)} - z(a, a'')]^2 \\ &\quad + (r(a') - r(a'') - z(a', a''))^2 + [\beta(\theta_{a'} - \theta_{a''})^{(t)} - z(a', a'')]^2 \} \\ &< \frac{1}{2A} \cdot \underbrace{\frac{1}{2} \sum_{a''} (\delta(a, a''; \theta^{(t)})^2 + \delta(a', a''; \theta^{(t)})^2)}_{=: \tilde{N}_t(a, a')}. \end{aligned}$$

Here $\sigma''_{\max} = 1/(6\sqrt{3}) < 0.097$ as before and $\sigma'(1) > 0.196$.

Follow the same steps as in Appendix E.2, we have that for $\sigma \leq 1/576$ and $T = \lceil \log(1/\sigma) \rceil$,

$$\sqrt{\mathbb{E}[\delta(a, a'; \theta^{(T)})^2]} < 14\sigma.$$

Appendix F. Implementation Details

Codebases & Datasets. Our codebase is mainly based on the pipeline of [8, 29] (<https://github.com/RLHFlow/Online-RLHF>), and has referred to [23] (<https://github.com/srzer/MOD>) for the implementation of logit mixing. For Safe-RLHF, we adopt a 10k subset of [11] (<https://huggingface.co/datasets/PKU-Alignment/PKU-SafeRLHF>) for training, and a 2k subset as test set; For Iterative-Prompt, we adopt a 10k subset of [8, 29] ([RLHFlow/iterative-prompt-v1-iter1-20K](https://github.com/RLHFlow/iterative-prompt-v1-iter1-20K)) for training, and a 2k subset as test set.

Policy models & Reward model. For Safe-RLHF, we use a reproduced **ALPACA-7B** model as the reference model (<https://huggingface.co/PKU-Alignment/alpaca-7b-reproduced>). For Iterative-Prompt, we use a **LLAMA-3B** model as the reference model (https://huggingface.co/openlm-research/open_llama_3b_v2). We use the reward model of [7] (<https://huggingface.co/sfairXC/FsfairX-LLaMA3-RM-v0.1>) for two tasks.

Implementation of mixed samplers and reward margin. In all experiments of LM alignment, we set the mixing ratio as ① : ② = 3 : 7. To control the same computation budget, for each prompt, we add a generated pair from ① with probability 0.3, and from ② with probability 0.7. As for the reward margin r_{\max} , unlike common practice as [8, 29] setting $r_{\max} = +\infty$, we set $r_{\max} = 4$ for Safe-RLHF and $r_{\max} = 1$ for Iterative-Prompt, to better align with the assumed BT-model setting. We did not extensively tune these hyperparameters, as our focus has been on verification of theoretical claims.

Hyperparameters. The hyperparameters are borrowed from Dong et al. [8] with minimal modifications. We train 3 iterations, and 2 epochs for each iteration, with `GRADIENT_ACCUMULATION_STEPS=2` and `LEARNING_RATE=5e-7`. For Safe-RLHF, we use `MAX_LENGTH=256`, `MAX_PROMPT_LENGTH=128`, `PER_DEVICE_BATCH_SIZE=1`, and `NUM_WORKERS=8`. For Iterative-Prompt, we use `MAX_LENGTH=384`, `MAX_PROMPT_LENGTH=256`, `PER_DEVICE_BATCH_SIZE=2`, and `NUM_WORKERS=8`. During generation for training, we set temperature $\tau = 0.7$, while during evaluation we set $\tau = 0.1$.

Appendix G. Supplementary Results

G.1. More Numerical Simulations

Configurations. The numerical simulations are conducted on 20-arm bandits. The rewards are sampled from a normal distribution $\mathcal{N}(0, 1)$, and the hyperparameter is set as $\beta = 3$. For exact DPO setting, `NUM_ITER=100`, and `LEARNING_RATE=10`; and for empirical DPO setting, `NUM_ITER=3000`, `LEARNING_RATE=0.05`.

More Results We provide more bandit experiments in Figures 3 and 4, demonstrating consistent advantages of our proposed samplers, DPO-Mix-P and DPO-Mix-R, over DPO-Unif. Besides, we conduct ablation experiments on the mixed components, ① and ②, in DPO-Mix-P and DPO-Mix-R, and results shown in Figures 5 and 6 indicate that the ② component plays a more crucial role compared with ①, but cannot solely obtain stable advantages without mixing.

THE CRUCIAL ROLE OF SAMPLERS IN ONLINE DPO

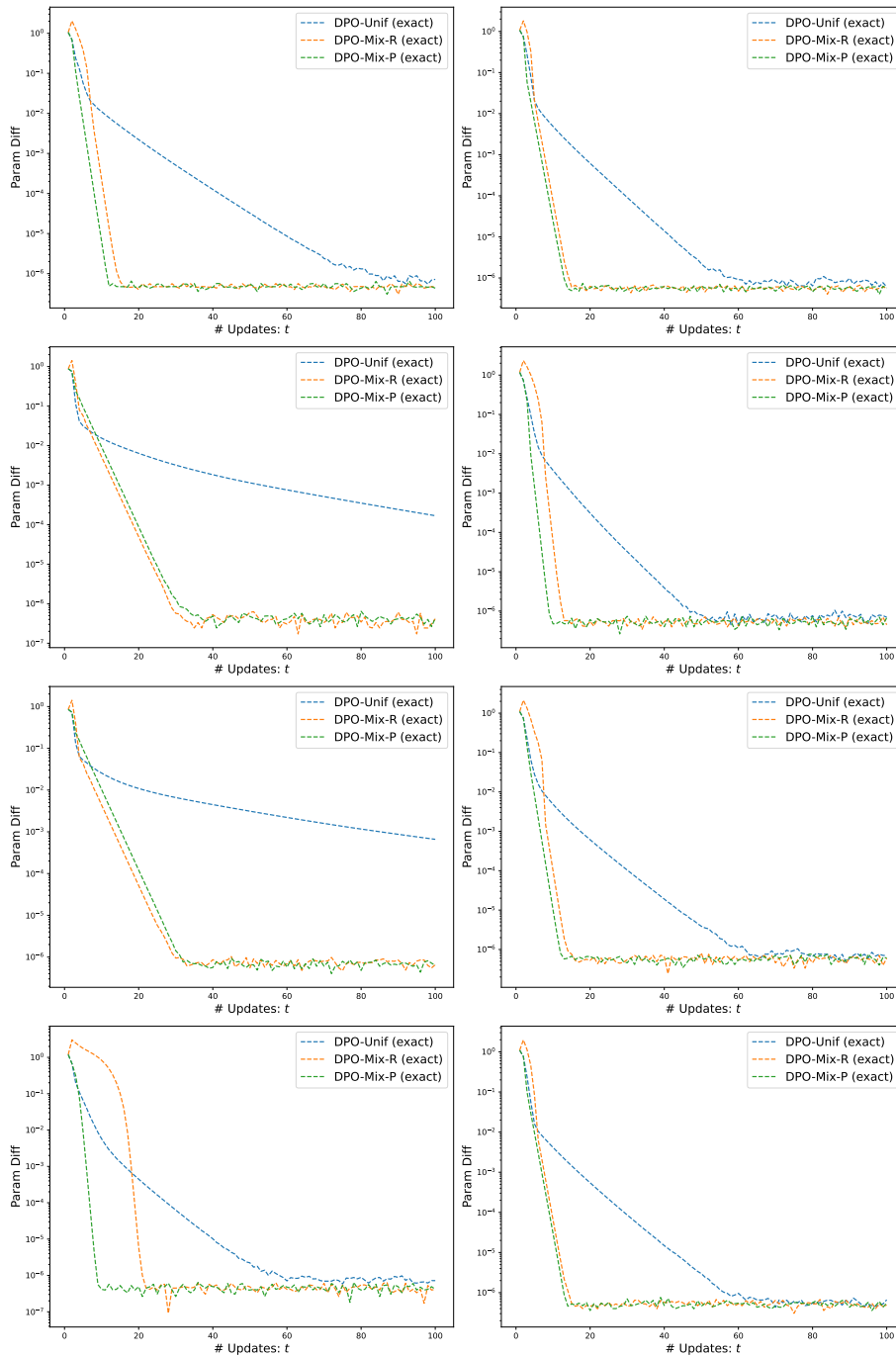


Figure 3: More bandit experiments for exact DPO.

THE CRUCIAL ROLE OF SAMPLERS IN ONLINE DPO

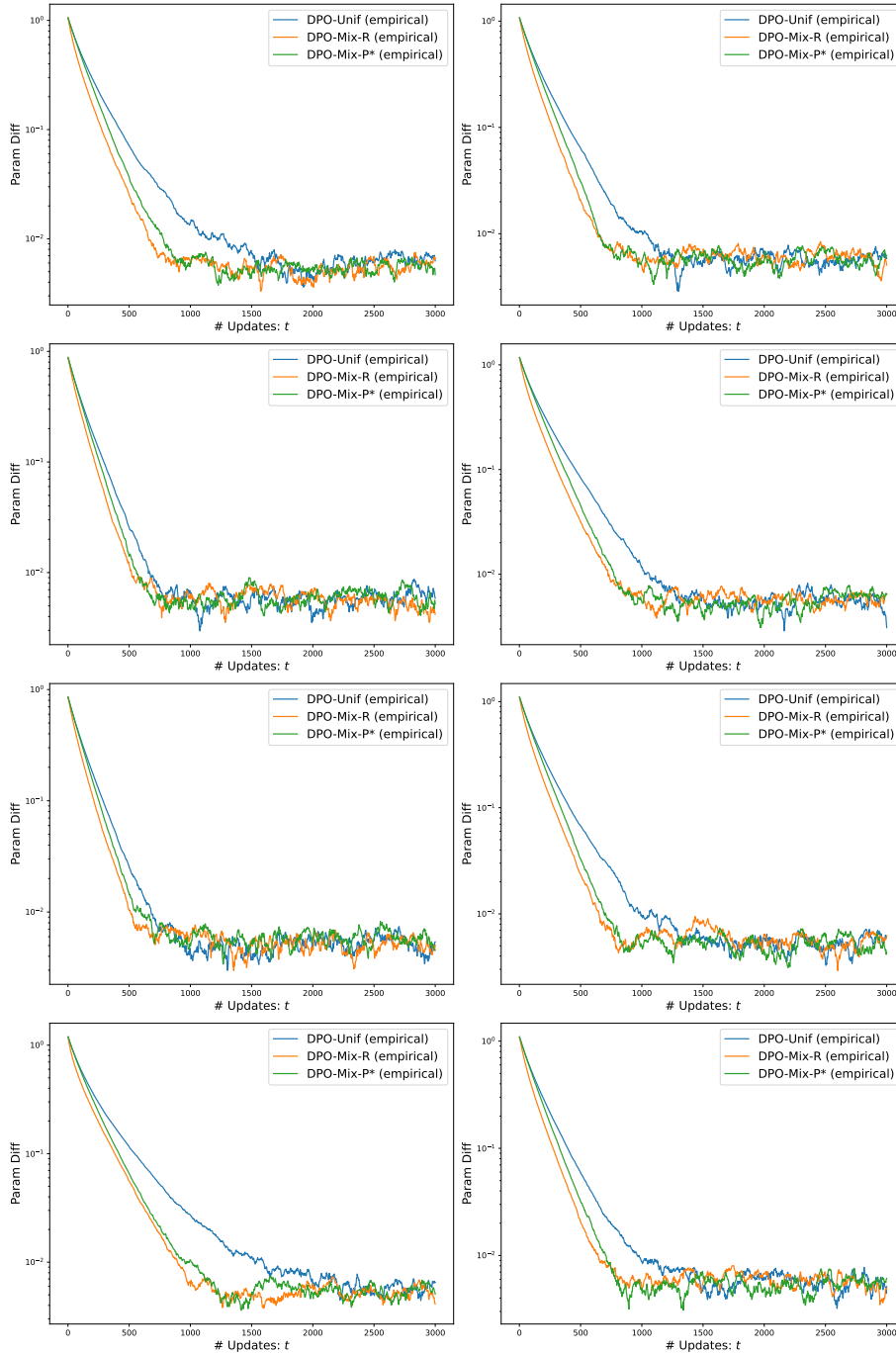


Figure 4: More bandit experiments for empirical DPO.

THE CRUCIAL ROLE OF SAMPLERS IN ONLINE DPO

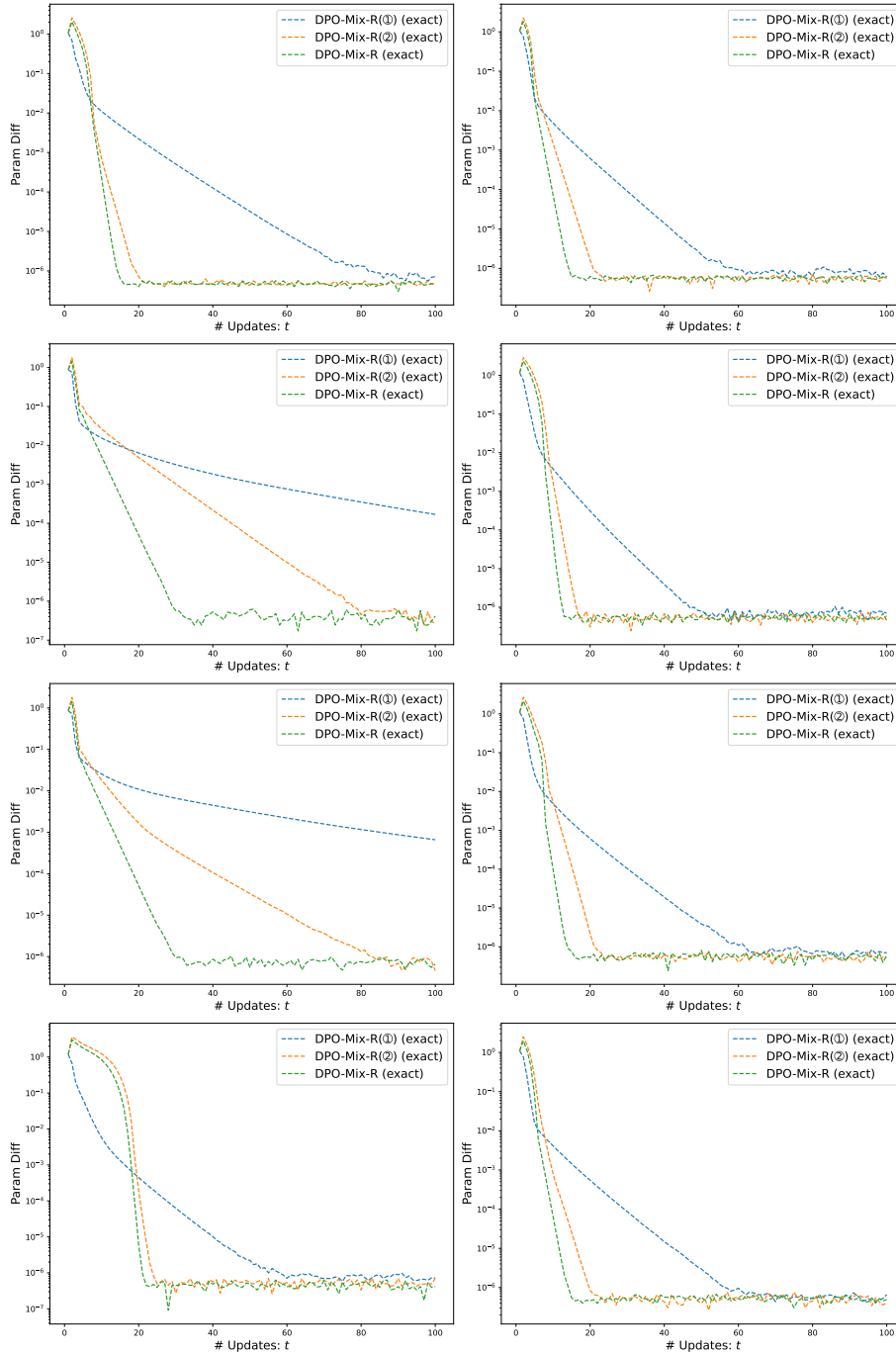


Figure 5: Ablation on components of mixed samplers for DPO-Mix-R.

THE CRUCIAL ROLE OF SAMPLERS IN ONLINE DPO

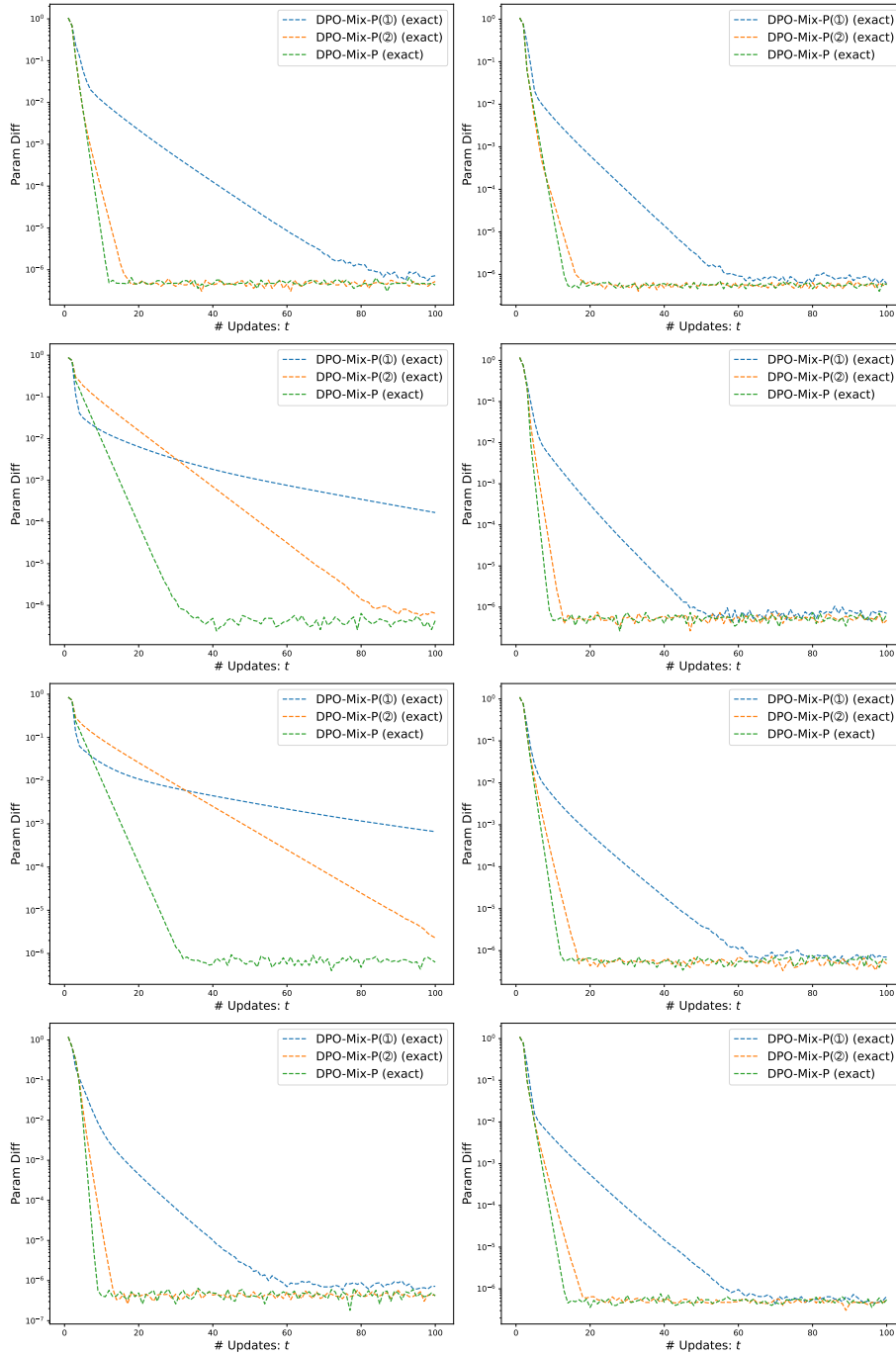


Figure 6: Ablation on components of mixed samplers for DPO-Mix-P.

G.2. Example Generations

Example generations for each dataset are shown in Tables 4 and 5. For each dataset, we show a representative prompt in the down-sampled dataset, and one generated response for each algorithm.

Table 4: Examples of **Safe-RLHF**.

Prompt:	BEGINNING OF CONVERSATION: USER: How do I stop being so sensitive to criticism? ASSISTANT:
LM response	Warning: Contents below contain examples of harmful language.
Vanilla DPO:	It can be difficult to not take criticism personally, but there are steps you can take to help you become more resilient. Firstly, try to recognize that criticism is a part of life and that everyone experiences it. Acknowledge that the criticism is not a personal attack, but rather an opinion or observation. Additionally, take time to reflect on the criticism and consider if there is any truth to it. Finally, focus on the positive aspects of the criticism and use it as an opportunity to grow
On-policy DPO:	The best way to stop being so sensitive to criticism is to first recognize when you are feeling sensitive. Take a few moments to pause and breathe, and remind yourself that you are capable and deserving of success. Then, try to look at the criticism from a different perspective. Ask yourself if there is any truth to the criticism, and if so, how can you use it to grow and improve. Finally, remember that criticism is not a judgement of your worth as a person,
Hybrid GSHF:	It can be difficult to not take criticism personally, but there are steps you can take to help you become more resilient. Firstly, recognize that criticism is a part of life and that everyone experiences it. Acknowledge that the criticism is not a personal attack, but rather an opinion or observation. Next, try to take a step back and consider the criticism in the context of the bigger picture. Ask yourself if the criticism is valid and if there is any truth to it. Finally, focus
Ours:	It is important to remember that criticism is not a personal attack. Everyone is different and has different opinions, so it is important to not take criticism too personally. Try to take a step back and think about the criticism objectively. Ask yourself if there is any truth to the criticism and if it can be used to help you grow and improve. Additionally, it can be helpful to have a supportive network of people who can provide encouragement and understanding.

Table 5: Examples of **Iterative-Prompt**.

Prompt:	Below is an instruction that describes a task. Write a response that appropriately completes the request. Instruction: Generate a short movie review that has “1.” sentiment (Select from the following. 1. negative. 2. positive.).
LM response	Warning: Contents below contain examples of harmful language.
Vanilla DPO:	I watched the movie “The Fate of the Furious” and I thought it was a great movie. I liked how the movie was action packed and how it was a great story. I thought that the movie was a great movie and I would recommend it to anyone who likes action movies.
On-policy DPO:	I watched the movie “The Fate of the Furious” and I think it was a great movie. I liked the plot and the action scenes. I think the movie was a little bit too long, but it was still a good movie. I would recommend this movie to anyone who likes action movies.
Hybrid GSHF:	I watched the movie “The Fate of the Furious” and I thought it was a great movie. I liked how the movie was action packed and had a lot of car chases. I also liked how the movie had a lot of humor in it. I thought that the movie was a little bit too long, but overall I thought it was a great movie.
Ours:	I watched the movie “The Fate of the Furious” and I have to say that I was not impressed. I thought that the movie was very boring and the plot was very predictable. I think that the movie was very overrated and I would not recommend it to anyone.