Unintentional Unalignment: Likelihood Displacement in Direct Preference Optimization

Noam Razin, Sadhika Malladi, Adithya Bhaskar, Dangi Chen, Sanjeev Arora, Boris Hanin

Princeton Language & Intelligence **Princeton ORFE**



Likelihood displacement is prevalent, yet not well understood



C³ Question #1

Why does likelihood displacement occur?

C³ Question #2

Benign

(e.g. Pal et al. 2024, Yuan et al. 2024, Rafailov et al. 2024, Tajwar et al. 2024, Liu et al. 2024, Pang et al. 2024)

What are its implications?

Likelihood Displacement Can Cause Unintentional Unalignment

Setting: Train a language model to refuse unsafe prompts via DPO



() Results

Probability shifts from preferred refusals to harmful responses!

E.g., the refusal rate of Llama-3-8B-Instruct drops from 74.4% to 33.4%

Theory: Likelihood Displacement is Driven by the Embedding Geometry

Approach: Characterize evolution of log probabilities

Ω Main Takeaway

Preferences with similar hidden embeddings lead to likelihood displacement!

Similarity measured by the Centered Hidden



3

4

Identifying Sources of Likelihood Displacement

Ω Main Takeaway

CHES score identifies samples causing likelihood displacement, while alternative measures do not



Setting: Llama-3-8B trained via DPO on UltraFeedback subsets

Paper includes similar results for the OLMo-1B and Gemma-2B models and AlpacaFarm dataset

CHES Score Edit Distance Hidden Embedding Similarity

Data Filtering via CHES Score

Ω Main Takeaway

Removing samples with high CHES scores mitigates unintentional unalignment



DPO + SFT Initial DPO DPO (gold data) DPO over samples with lowest length-normalized CHES score

Setting: The same as in

model to refuse unsafe

1 – train a language

prompts via DPO