Declarative Characterizations of Direct Preference Alignment Algorithms

Kyle Richardson, Vivek Srikumar, Ashish Sabharwal



Preference alignment for large language models



LLM development pipeline

Important stage in LLM development, tuning models to human preferences.

• Much recent work on **offline** alignment, learning from pairwise data.

Safety example (Dai et al., 2024; Ji et al., 2024)

- x : Will drinking brake fluid kill you?
- y₁: No, drinking brake fluid will not kill you

Symbolic program derivation

:= Implies(

And $(\text{Ref}(x, y_w), M(x, y_l))$,

And (Ref (x, y_l) , M (x, y_w))

 y_w : Drinking brake fluid will not kill you, but it can be extremely dangerous... [it] can lead to vomiting, dizziness, fainting,

compilation

Loss function

 $-\log\sigmaigg(eta\lograc{\pi_ heta(x,y_w)}{\pi_{ ext{Ref}}(x,y_w)}-eta\lograc{\pi_ heta(x,y_l)}{\pi_{ ext{Ref}}(x,y_l)}igg)$

Direct Preference Alignment (DPA) approaches

MLE-style closed-form approaches, direct tuning on offline data, alternative to RL.

approaches



 $\log rac{\pi_{ heta}(y_w|x)^{rac{1}{|y_w|}}}{|y_w|}$

 $\pi_{ heta}(y_l|x)^{\,\overline{|y_l|}}$

implementation

details

variations



Preference learning as a discrete reasoning problem

Question: given an existing loss function, can we **derive** a symbolic expression that characterizes its core semantics?



constraints on predictions (variables)

Output distribution with thresholding, delimiting valid vs. invalid output

- Expresses high-level model behavior, structure of output distribution.
- Assumption: every loss has an internal logic, our goal is to uncover that logic, use to formally characterize DPA space.

There are many variations of DPO that modify details of the original loss, account for various shortcomings.

RRHF $\max(0, -\rho_{\theta})$

issue: Not clear how losses relate to one another, the underlying principles.

remo **E.g.**, how many refere losses are there? mode What is the chang structure of the refere space we are policy searching? (optimization add agnostic)

	Loss $\rho_{\theta} := \log \frac{\rho_{\theta}^t}{\rho_{\theta}^b}, \ P_{\{\theta, \text{ref}\}} \approx \pi_{\{\theta, \text{ref}\}}$
	Baselines ρ_{θ}
	$\ell_{CE} \log \frac{P_{\theta}(y_w x)}{1 - P_{\theta}(y_w x)} = $
	$\ell_{\texttt{CEUnl}} \log \frac{P_{\theta}(y_w x)(1 - P_{\theta}(y_l x))}{P_{\theta}(y_l x) + (1 - P_{\theta}(y_w x)))}$
emove	Single model approaches (no reference) P_{θ}
eference nodel	$\ell_{\text{CPO}} \log \frac{P_{\theta}(y_w x)}{P_{\theta}(y_l x)}$
hange	$\ell_{\text{ORPO}} \log \frac{P_{\theta}(y_w x)(1-P_{\theta}(y_l x))}{P_{\theta}(y_l x)(1-P_{\theta}(y_w x))}$
eference olicy	$\ell_{\texttt{SimPO}} \log \frac{P_{\theta}(y_w x) P_{\text{mref}}(y_l x)}{P_{\text{mref}}(y_w x) P_{\theta}(y_l x)}$
, energy	with reference model P_{ref}
add	$\ell_{\text{DPO}} \log \frac{P_{\theta}(y_w x) P_{\text{ref}}(y_l x)}{P_{\text{ref}}(y_w x) P_{\theta}(y_l x)}$
more	$\ell_{\text{DPOP}} \log \frac{P_{\theta}(y_w x)P_{\theta2}(y_w x)P_{\text{ref}}(y_l x)}{P_{\text{ref}}(y_w x)P_{\text{ref2}}(y_w x)P_{\theta}(y_l x)}$

A logic for preference learning and reduction to weighted model counting

Loses expressed **semantically** as symbolic formulas interpreted in probabilistic logic; counting the propositional models of formulas.

Devised mechanical procedure for **translating** DPA losses into logic, novel encoding, preference structure.



Derived representations are **correct**, guaranteed to uniquely compile back into original loss under new logic.

Losses differ in terms of hard constraints employed encoding Symbolic Program \overline{P} Loss CE $M(\mathbf{x},\mathbf{y}_w), \ \mathsf{P}_{\mathbf{H}} := \top$ Unl $\mathsf{P} := \mathsf{And}(\mathsf{M}(\mathbf{x},\mathbf{y}_w), \mathsf{Not}(\mathsf{M}(\mathbf{x},\mathbf{y}_l)))$ $\mathsf{P}_{\mathbf{H}} := \top$ CPO P :=**Implies**(M(x,y_l), M(x,y_w)) one-true constraint $\mathsf{P}_{\mathbf{H}} := \operatorname{Or}(\mathtt{M}(\mathbf{x}, \mathbf{y}_l), \ \mathtt{M}(\mathbf{x}, \mathbf{y}_w))$ ORPO $\mathsf{P} := \mathtt{Implies}(\mathtt{M}(x, y_l), \mathtt{M}(x, y_w))$;; one-hot constraint $P_{\mathbf{H}} := \mathbf{Or}($ $And(M(x,y_l), Not(M(x,y_w))),$





predictions; e.g., ~4.2 billion definable variants of DPO.

Exploring loss space, deriving new losses from first principles



Can use logical semantics to structure DPA space, explore space; **loss** lattice ordered by entailment.

Procedure: Start from successful losses, formalize/modify semantics then experiment; devise novel formulas from scratch.

select citations

Darwiche, Adnan et al. . On probabilistic **inference by weighted model counting** 2008, Journal of Artificial Intelligence

Xu, Jingyi et al. **A semantic loss function for** deep learning with symbolic knowledge ICML

Tang, Yunhao et al. Generalized Preference **Optimization.** ICML 2024

Xu, Haoran et al. Contrastive Preference **Optimization** ICML 2024

Rafailov, Rafael et al. Direct Preference **Optimization.** NeurIPS 2023

Hong, Jiwoo et al. ORPO: Monolithic **Preference Optimization without Reference** Model . EMNLP 2024

Meng, Yu. et al. SimPO: Simple Preference **Optimization with a Reference-free Reward** NeulIPS 2024