

Mallows meets DPO (Direct Preference Optimization)

In this note, we explain how Mallows ranking models are integrated into DPO [RS23], in the context of RLHF.

① Mallows models.

Given n items indexed by $\{1, 2, \dots, n\}$, $\sigma: \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ specifies the ranking of these items.

E.g. $\sigma(i) = j$ means that item i has rank j . Smaller the rank is, better the item is.

$$\begin{aligned} i > i' &\Leftrightarrow i \text{ wins over } i' \\ &\Leftrightarrow \sigma(i) < \sigma(i') \end{aligned}$$

$$P(\sigma) \propto \phi^{\frac{d(\sigma, \sigma_0)}{\phi}}$$

\uparrow dispersion parameter \leftarrow central ranking.

where $d(\cdot, \cdot)$ is a right invariant metrics. Refer to [T18] for further details.

Examples of "d":

$$\textcircled{a} \quad d(\sigma, \sigma_0) = \sum_i (\sigma(i) - \sigma_0(i))^2$$

Mallows's L^2
or Mallows Θ model

Such d is called Spearman's rank correlation.

$$\textcircled{b} \quad d(\sigma, \sigma_0) = \text{inv}(\sigma_0 \sigma^{-1})$$
 \nearrow
Mallows ϕ model

number of inversions
Such d is called Kendall's tau.

$$\textcircled{c} \quad d(\sigma, \sigma_0) = \text{minimum transportations taking } \sigma \text{ to } \sigma_0$$
Ewens model.

Such d is called Cayley distance.

... Many other choices of "d", see [D88].

Example \textcircled{c} is trivial, as $P(i > i') = \frac{1}{2}$ independent of i, i' .

We focus on Example \textcircled{a} and \textcircled{b} , where we recover the DPO [RS23] from Example \textcircled{a} , and design a novel DPO from Example \textcircled{b} .

Example \textcircled{a} : The Mallows θ model is also Bradley-Terry with

$$\begin{aligned}
 P(i > i') &= P(\sigma(i) < \sigma(i')) \\
 &= \frac{1}{1 + e^{-2(\theta_0(i') - \theta_0(i)) \log \phi}}
 \end{aligned}$$

3 ✓

By ~~RS23~~

$$:= g(-r(x, i') + r(x, i)),$$

where $g(x) = \frac{1}{1+e^{-x}}$ is the link function

$r(x, i) = +2(\sigma_0(i|x))^{\log 2}$ is the reward

By [RS23], RL gives $r(x, i) = \beta \log \frac{\pi(i|x)}{\pi_{\text{ref}}(i|x)} + \beta \log 2$

$$\mathbb{P}(i > i' | x) = g(r(x, i) - r(x, i')) \quad \checkmark$$

This recovers the DPO in [RS23].

Example (b). It is known that [M57]:

$$\mathbb{P}(i > i') = \mathbb{P}(\sigma_0(i) < \sigma_0(i'))$$

$$= \begin{cases} \frac{\sigma_0(i') - \sigma_0(i) + 1}{1 - e^{(\sigma_0(i') - \sigma_0(i) + 1) \log 2}} - \frac{\sigma_0(i') - \sigma_0(i)}{1 - e^{(\sigma_0(i') - \sigma_0(i)) \log 2}} & \text{if } \sigma_0(i') > \sigma_0(i) \\ 1 - \text{"} & \text{if not} \end{cases}$$

$$:= g(-r(x, i') + r(x, i)),$$

where $g(x) = \begin{cases} \frac{x+1}{1-\phi^{x+1}} - \frac{x}{1-\phi^x} & x > 0 \\ 1 - \phi & x < 0 \end{cases}$ is the link function
 $r(x, i) = -\sigma_0(i|x)$ is the reward.

Same as before, replacing $r(x, i)$ with $\beta \log \frac{\pi(i|x)}{\pi_{\text{ref}}(i|x)} + \beta \log Z$

$$\mathbb{P}(i \geq i' | x) = g(r(x, i) - r(x, i')) \leftarrow$$

we derive a novel DPO.

Subtlety: In the link function, there is still the dispersion parameter " ϕ ".

One can estimate $\hat{\phi}$ as in the reward model.

This is not hard, see [MB10, Section 3.2].

Ref:

[RS23]. Rafailov et al, Direct Preference Opt ... Neurips '23.

[T18]. Tang, Mallows ranking model ... ICML '18.

[D88]. Diaconis, Group representation in Prob & Statistics
IMS Lecture note '88.

[M57]. Mallows, Non-null ranking models, Biometrika '57.

[MB10]. Meila & Bao, An exponential model for infinite rankings
JMLR '10.