

Supervised Learning Preference Optimization: Rethinking RLHF and DPO as Supervised Learning

Yaoshiang Ho

YAOSHIANG@GMAIL.COM

Abstract

We propose Supervised Learning Preference Optimization (SLPO), an approach to aligning language models that relies solely on the classic supervised learning loss function, cross-entropy. Unlike Reinforcement Learning from Human Feedback (RLHF), SLPO directly adjusts probability mass for chosen, rejected, and other sequences relative to a reference model, eliminating the need for KL divergence regularization. SLPO also avoids the extensive reparameterization of model outputs into scores and the Bradley-Terry model required by Direct Preference Optimization (DPO). As a result, SLPO is a purely supervised learning approach to alignment, free of reinforcement learning concepts. We further demonstrate how to efficiently implement a targeted probability mass for the intractably large set of sequences that are neither explicitly chosen nor rejected.

Keywords: Reinforcement Learning from Human Feedback (RLHF), Direct Preference Optimization (DPO)

1 Introduction

Alignment is the task of ensuring that the behavior of a Large Language Model (LLM) is consistent with human preferences.

The alignment approach popularized by the commercial success of ChatGPT was Reinforcement Learning from Human Feedback (RLHF), (Ouyang et al., 2022). Despite its effectiveness, RLHF requires training a second model, called a reward model, as well as Proximal Policy Optimization (PPO), resulting in a technique that is more complex than the basic supervised learning. It also requires a Kullback-Leibler (KL) divergence term to regularize the changes to the LLM during alignment training.

Direct Preference Optimization (DPO) is a simpler approach to alignment which does not require a secondary reward model. In the paper introducing DPO, the authors examine the underlying approach of RLHF and propose the DPO objective to align the target LLM directly using maximum likelihood estimation (MLE). The key insight from the DPO paper is that an LLM’s outputs can be reparameterized into a reward model using ratios, logs, and the Bradley-Terry model (Bradley and Terry, 1952).

The specific contribution of this paper is to reframe the alignment phase away from reward modeling entirely and treating it simply as a pure supervised learning problem by training a model to align to a directly modified probability distribution. We call this approach Supervised Learning Preference Optimization (SLPO).

2 Related Work

As discussed in the introduction, the commercial success of GPT3 and ChatGPT popularized RLHF as an alignment technique (Ouyang et al., 2022). To resolve weaknesses of RLHF, DPO was introduced and has been used in recent frontier models like Meta’s Llama 3 family of models (Rafailov et al., 2023; Dubey et al., 2024).

Following DPO, the space of alignment methods has rapidly expanded. Given the pace of innovation in this area, we refer to a recent comprehensive survey on alignment methodologies by Wang et al. (2024).

3 Preliminaries: DPO

We analyze the DPO objective (Rafailov et al., 2023, Eq. 7), defined by its authors as :

$$L_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = \underbrace{-\mathbb{E}_{(x, y_w, y_l) \sim D} \log}_{1} \left[\underbrace{\sigma}_{2} \left(\underbrace{\beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)}}_{3} - \underbrace{\beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)}}_{4} \right) \right] \quad (1)$$

where

π_{θ} is the language model to be aligned,

π_{ref} is the reference language model,

D is the dataset of training examples,

x is the input context,

y_w is the winning sequence,

y_l is the losing sequence,

σ is the sigmoid function,

β is a hyperparameter.

In the underbraced section 1, we see the standard negative log likelihood (NLL) objective. In the underbraced section 2, we see the Bradley-Terry model ¹. In the underbraced sections 3 and 4, we see how the reference and language model’s predictions are reparameterized into a winning and losing score: they are the log of the ratio of the language model to the reference model, for the winning and losing completion, respectively. This score is later described as the reward function.

1. A ranking method that is mathematically equivalent and perhaps more widely understood is the ELO score, used to rank Chess players, and, LLMs in the Chatbot Arena (Elo, 1978; Chiang et al., 2024). Both ELO and Bradley-Terry assign scores to players, and pass the difference through a sigmoid function to assess the probability of the minuend (aka LHS) player of winning and subtrahend (aka RHS) player of losing.

The optimization of y_w , a sequence of tokens, rather than a specific token, is a key difference between the alignment phase and other phases of training an LLM. In particular, the pre-training and supervised fine-tuning (SFT) phases focuses on next-token prediction.

With simple algebraic manipulation, we can rewrite this objective as:

$$L_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = \underbrace{-\mathbb{E}_{(x, y_w, y_l) \sim D} \log}_1 \left[\underbrace{\sigma}_2 \left(\underbrace{\beta \log \frac{\pi_\theta(y_w | x)}{\pi_\theta(y_l | x)}}_3 - \underbrace{\beta \log \frac{\pi_{\text{ref}}(y_w | x)}{\pi_{\text{ref}}(y_l | x)}}_4 \right) \right].$$

We can interpret the undercomponents as follows: The first and second underbraces remain unchanged. The third underbrace is the log of the ratio language models probability of the winner divided by the loser. This ratio is optimized to be bigger, given that the log, sigmoid, and log from underbraces 1, 2, and 3 are all monotonic functions. Since a ratio of probabilities is an odds ratio, we will refer to this as the language model’s *log-odds ratio*.

The fourth underbrace is reference model’s *log-odds ratio*. As the weights of the reference model are constant during optimization of the language model, so too are these values. However, each row of training data will have a different value for this constant. More specifically, since π_θ is initialized as π_{ref} , the difference between underbraces 3 and 4 starts at zero during training, and the shape of the sigmoid function (underbrace 2) and its gradient are known. During training, as the language model’s *log-odds ratio* increases, the sigmoid function will naturally regularize and decelerate the increase by reducing the magnitude of the gradient (a useful version of the vanishing gradient problem). This is the exact outcome described by the DPO authors in their analysis of the gradient of DPO. *But we have developed a different intuition for the DPO loss: rather than reparameterizing the language model’s output into a reward model, the DPO objective optimizes the language model’s log-odds ratio, with a constant shift so that the shifted value starts at zero during optimization.*

Let’s continue reworking the DPO objective into mathematically equivalent forms to better understand it by looking at tokens and their logits rather than sequences ².

The variable y is a sequence of tokens, defined as

$$\pi_{\text{ref}}(y | x) = \prod_{t=1}^T \pi_{\text{ref}}(y_t | x, y_{<t}), \quad (2)$$

where

$y = (y_1, y_2, \dots, y_T)$ is the output sequence,

x is the input context, and

2. Technically, a logit is the log-odds ratio and it is the output of a feature extractor before a sigmoid activation function in binary classification. The term logit is also used for categorical classification, when a feature extractor’s outputs are activated with the softmax function. However, these softmax logits cannot be exponentiated to calculate an odds ratio. They do have an unfortunate property however: they are shift invariant, leading to numerical instability and the need for log-softmax and its use of the log-sum-exp trick to improve numerical stability. In this paper, since we are dealing with both sigmoid and softmax functions, we will refer to the value passed into a softmax function as a softmax logit.

$y_{<t} = (y_1, y_2, \dots, y_{t-1})$ represents the tokens generated prior to time step t .

$\pi_{\text{ref}}(y_t \mid x, y_{<t})$ denotes the conditional probability of generating token y_t given the input x and the previously generated tokens $y_{<t}$.

Since both π language model are LLMs, they are activated using the softmax function. The softmax operation can be mathematically cumbersome due to its dependence on all other logits for normalization. For simplicity, we assume that the logits are transformed using a log-softmax function, which computes the logarithm of the softmax probabilities in a numerically stable way. This transformation results in log-probs, which are easier to reason about since they can be exponentiated to recover probabilities. Importantly, we do not lose any generality because log-probs can also be passed through a softmax function to recover probabilities, ensuring equivalent behavior.

Let us establish the term f for the layers of the model up to the softmax activation, namely, the feature extractor which outputs logits. The output of g are logprobs, which are a special case of logits: they can be converted to the probability distribution function using either softmax or exponentiaion.

$$g(y \mid x) = \text{logsoftmax}(f(y \mid x)) \quad (3)$$

such that

$$\pi(y \mid x) = \exp(g(y \mid x)) = \text{softmax}(g(y \mid x)) = \text{softmax}(f(y \mid x)).$$

Plugging Equations 2 and 3 into the DPO objective 1,

$$L_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \log \left[\sigma \left(\beta \log \frac{\prod_{t=1}^{T_w} \exp(g_{\theta}(y_{w,t} \mid x, y_{w,<t}))}{\prod_{t=1}^{T_w} \exp(g_{\text{ref}}(y_{w,t} \mid x, y_{w,<t}))} - \beta \log \frac{\prod_{t=1}^{T_l} \exp(g_{\theta}(y_{l,t} \mid x, y_{l,<t}))}{\prod_{t=1}^{T_l} \exp(g_{\text{ref}}(y_{l,t} \mid x, y_{l,<t}))} \right) \right]$$

which simplifies to:

$$L_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \log \left[\sigma \left(\beta \left(\sum_{t=1}^{T_w} [g_{\theta}(y_{w,t} \mid x, y_{w,<t}) - g_{\text{ref}}(y_{w,t} \mid x, y_{w,<t})] - \sum_{t=1}^{T_l} [g_{\theta}(y_{l,t} \mid x, y_{l,<t}) - g_{\text{ref}}(y_{l,t} \mid x, y_{l,<t})] \right) \right) \right] \quad (4)$$

We notice another possibility to rewrite the DPO objective towards something more familiar. Optimizing the sigmoid of a difference $a - b$ through BCE is equivalent to optimizing the softmax of a and b through CCE towards a one-hot encoded target of $y_{\text{true}} = [1, 0]$. This is well known but derived in Appendix A. This allows us to rewrite the DPO objective as follows:

$$L_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \text{CCE} \left[y_{\text{true}} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, y_{\text{pred}} = \text{softmax} \left(\begin{bmatrix} \beta S_w \\ \beta S_l \end{bmatrix} \right) \right] \quad (5)$$

where

$$\begin{aligned}
 S_w &= \underbrace{\sum_{t=1}^{T_w} g_{\theta}(y_{w,t} \mid x, y_{w,<t})}_A - \underbrace{\sum_{t=1}^{T_w} g_{\text{ref}}(y_{w,t} \mid x, y_{w,<t})}_B \\
 S_l &= \underbrace{\sum_{t=1}^{T_l} g_{\theta}(y_{l,t} \mid x, y_{l,<t})}_A - \underbrace{\sum_{t=1}^{T_l} g_{\text{ref}}(y_{l,t} \mid x, y_{l,<t})}_B
 \end{aligned}$$

Key Observation #1: Much of the DPO objective can be simplified to the form of a standard classification task. Equation 5 provides our first key observation. What started off as a complex reparameterization of the language model’s output into a reward model and a Bradley-Terry model, has been simplified into a standard categorical crossentropy loss applied to the softmax, with simple to understand targets. However, the input to the softmax is not a single logit, but rather, the summation of all the winning (or losing) logits in the language model (underbraces A), subtracted by the same summation from the reference model (underbraces B). Our next two observations will interpret those.

Key Observation #2: The difference terms simply zero-align the input to the softmax function across all examples during training, allowing the training to equally optimize all examples. Since the language model is initialized by the reference model, the difference introduced in underbraces B of Equation 5 simply zero-aligns the input to softmax, allowing the training process to send identical gradients to all examples, regardless of the relative probabilities $\pi(y_w)$ and $\pi(y_l)$ for that example. This difference is the vestige of underbraces 3 and 4 in the original DPO objective in Equation 1 was referred to as a "reward function", but we now consider it simply a zero-aligning of examples. We have successfully eliminated reinforcement learning concepts from the DPO objective.

Key Observation #3: DPO optimizes a joint probability The summation in underbraces A of Equation 5 may not be familiar. In fact, the softmax of a sum optimizes a joint probability. See Appendix B for a derivation. This should not be surprising since the original DPO objective optimizes the joint probability of the winning and losing sequences. Equation 2 expresses this as the joint probability of tokens, represented as a product, and a product of probabilities is equivalent to a sum of log probabilities.

Key Observation #4: Behind the summations and shifts in unbraces A and B of Equation 5, the final objective is still $+\inf$ or $-\inf$ for each token logit. This is obvious from Equation 5. An equivalent statement is that the DPO objective does not consider any sequence other than the winning and losing, and if left unchecked, would ultimately force the model to predict the winning sequence with probability 1 and all other sequences, including the loser, with probability 0. This is obviously destabilizing, e.g. other sequences may have been perfectly reasonable as winning sequences. This is a weakness of the DPO objective.

We have been able to remove much of the machinery of the DPO objective through algebraic manipulation and equivalences to more common deep learning concepts. We have also developed key observations on the strength and weaknesses of the DPO objective. Can we take these learnings and develop a superior objective?

4 Preliminaries: Probabilistic Classification

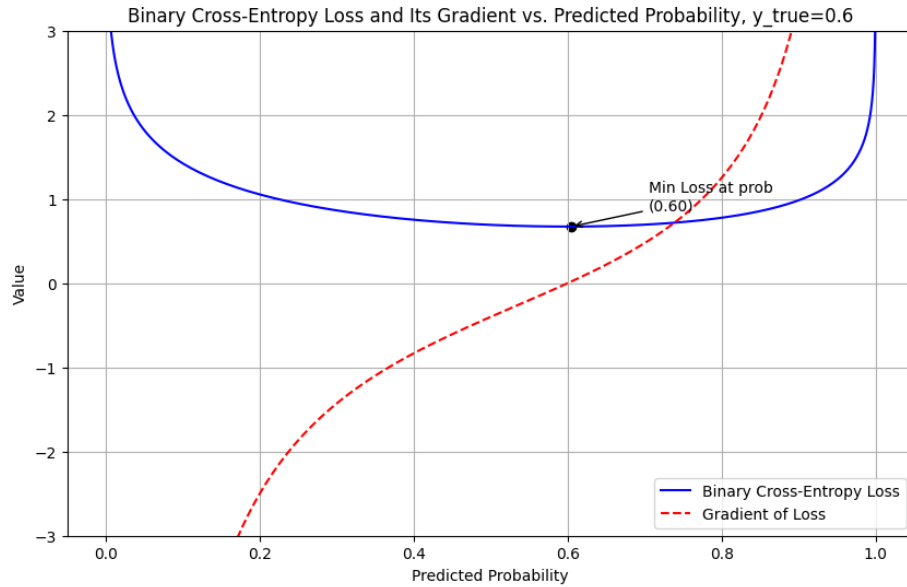


Figure 1: Setting y_{true} to a value other than 0 or 1 is a valid use of the crossentropy loss function and perfectly consistent with both KL divergence and MLE.

A basic task for deep learning has been classification, both binary and categorical.

- Look at this image and decide which of the ten digits it is (LeCun, 1998).
- Classify every pixel in this image as a foreground class or background.
- Decide if this movie review’s sentiment is positive or negative (Pang and Lee, 2004).

Modern LLMs auto-regressively predict a probability distribution for the next token among a finite set of possibilities, conditioned on the previous tokens. Hence, an LLM also solves a classification problem.

Given the prevalence of classification, deep learning oftens assumes target labels can be described as either ordinal numbers or equivalently one-hot encoded vectors. When target labels are intermediate values between one and zero, they were called "soft targets", treated as an regularization technique or even called dark knowledge (Hinton et al., 2015; Szegedy et al., 2016; Hinton et al., 2014).

But the two roots of the crossentropy loss functions, KL divergence and Maximum Likelihood Estimation (MLE) applied to the probability distribution function of a Multi-Bournoulli distribution, are not limited to binary and categorical outcomes. For example, the crossentropy loss function can be used directly to condition a model to predict to be true 60% of the time and false 40% of the time, as demonstrated in Figure 4.

The term we will use for the idea of predicting a probability distribution is *probabilistic classification*.

5 Supervised Learning Preference Optimization

Armed with intuition of the DPO loss and probabilistic classification, we can now turn to a simpler, supervised learning approach to alignment. Our goals are:

- Apply only MLE under a probabilistic classification approach. That is, we only want to use a CCE loss function.
- Avoid any RL concepts like rewards.
- Avoid any extra ratios, logs, or the Bradley-Terry model.
- Continue optimizing sequences rather than individual tokens.
- Continue giving comparable weight to all examples during optimization.
- Continue stabilizing the optimization to prevent divergence.
- Instead of having an unconstrained objective, specify an objective with a natural "bowl shape" at a target, to assist in optimization.

The SLPO objective achieves these goals by directly optimizing probabilities for the winning, losing, and all other sequences. It is defined as follows:

$$L_{\text{SLPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \left(\underbrace{w_w \cdot \log \pi_{\theta}(y_w \mid x)}_1 + \underbrace{0 \cdot \log \pi_{\theta}(y_l \mid x)}_2 + \underbrace{(1 - w_w) \cdot \log \pi_{\theta}(y_{\overline{w \cup l}} \mid x)}_3 \right) \quad (6)$$

where

x is the input context,

y_w is the winning token sequence,

y_l is the losing token sequence,

π_{θ} is the language model to be aligned,

π_{ref} is the reference language model,

$w_w = \pi_{\text{ref}}(y_w \mid x) + \pi_{\text{ref}}(y_l \mid x)$, and

$y_{\overline{w \cup l}}$ represents all other possible token sequences.

Intuitively, the goal is to take the probabilities of the winning and losing sequence and reallocate them to the winning sequence (underbrace 1). The losing sequence is optimized to zero probability (underbrace 2). Importantly, the other tokens are directly optimized to maintain the same probability as the reference model (underbrace 3). This is essentially the SLPO’s version of the KL divergence regularization term in traditional RLHF.

The SLPO objective is explicitly constructed as a CCE loss function and can be easily rewritten, similarly to Equation 5:

$$L_{\text{SLPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \text{CCE} \left[y_{\text{true}} = \begin{bmatrix} w_w \\ 0 \\ 1 - w_w \end{bmatrix}, y_{\text{pred}} = \begin{bmatrix} \pi_\theta(y_w | x) \\ \pi_\theta(y_l | x) \\ \pi_\theta(y_{w \cup l} | x) \end{bmatrix} \right] \quad (7)$$

Alternatively, a softer version of the SLPO objective would shift only some of the probability mass of the rejected sequence to the winning sequence.

$$\begin{aligned} L_{\text{SLPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} & \left(w_w \cdot \log \pi_\theta(y_w | x) \right. \\ & + w_l \cdot \log \pi_\theta(y_l | x) \\ & \left. + (1 - w_w - w_l) \cdot \log \pi_\theta(y_{w \cup l} | x) \right) \end{aligned} \quad (8)$$

where

$$w_w = \pi_{\text{ref}}(y_w | x) + \alpha \cdot \pi_{\text{ref}}(y_l | x),$$

$$w_l = (1 - \alpha) \cdot \pi_{\text{ref}}(y_l | x), \text{ and}$$

$$0 \leq \alpha \leq 1.0.$$

As before, we can rewrite this in the form of a CCE loss:

$$L_{\text{SLPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \text{CCE} \left[y_{\text{true}} = \begin{bmatrix} w_w \\ w_l \\ 1 - w_w - w_l \end{bmatrix}, y_{\text{pred}} = \begin{bmatrix} \pi_\theta(y_w | x) \\ \pi_\theta(y_l | x) \\ \pi_\theta(y_{w \cup l} | x) \end{bmatrix} \right] \quad (9)$$

6 Implementation

For a standard language model, the vocabulary size v is typically in the range of 10,000 to 100,000 tokens. The sequence length T is typically in the range of 128 to 1024 tokens. So the $y_{w \cup l}$ set of sequences is intractably large to enumerate.

Let us consider y_w independently of y_l . Then we can rewrite the SLPO Equation 6 in two halves and consider them separately:

$$L_{\text{SLPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim D} \left(\underbrace{w_w \cdot \log \pi_\theta(y_w | x)}_1 + \underbrace{(1 - w_w) \cdot \log \pi_\theta(y_{\bar{w}} | x)}_2 + \underbrace{0 \cdot \log \pi_\theta(y_l | x)}_3 + \underbrace{1 \cdot \log \pi_\theta(y_{\bar{l}} | x)}_4 \right) \quad (10)$$

The term focused on the winning sequence (underbrace 1) is simple: it is just the sum of the logits of the winning token sequences, passed through softmax and CCE with the target weight.

For the term focused on the losing sequence (underbrace 3), the loss at every token position is also simple: like underbrace 1, we sum the logits before we pass them into softmax and assign a target value of zero in the CCE loss. The terms in underbrace 2 and underbrace 4 are the tricky ones, since they deal with complementary sets. However, we can see a recursive pattern by visualizing a tree. For the first position in the sequence, there are essentially two options: the winning token (or losing token) and everything else. The target probability of everything else is simply $1 - w_w$, and the predicted probability is

$$\pi_\theta(y_{\bar{w},1} | x) = \sum_{v \neq y_{w,1}} \exp(\text{logsoftmax}(z_{v,1})).$$

Obviously, we can calculate this with the logsumexp trick for both numerical stability and to stay in log space. Importantly, *we do not need to recursively explore sequence position 2 and beyond for sequences that start with a non-winning token because the probabilities of all possible continuations sum to 1.*

Similarly, for position two, we can directly calculate the predicted probability as

$$\pi_\theta(y_{w,1}) \cdot \pi_\theta(y_{\bar{w},2}).$$

And so on for positions 3, 4, and 5:

$$\begin{aligned} & \pi_\theta(y_{w,1}) \cdot \pi_\theta(y_{w,2}) \cdot \pi_\theta(y_{\bar{w},3}) \\ & \pi_\theta(y_{w,1}) \cdot \pi_\theta(y_{w,2}) \cdot \pi_\theta(y_{w,3}) \cdot \pi_\theta(y_{\bar{w},4}) \\ & \pi_\theta(y_{w,1}) \cdot \pi_\theta(y_{w,2}) \cdot \pi_\theta(y_{w,3}) \cdot \pi_\theta(y_{w,4}) \cdot \pi_\theta(y_{\bar{w},5}) \end{aligned}$$

Summing across all possible positions results in

$$\log \pi_\theta(y_{\bar{w}} | x),$$

which we compute efficiently in a vectorizable way.

Formally, the term $\log \pi_\theta(y_{\bar{w}} | x)$ is calculated as:

$$\log \sum_{t=1}^T \underbrace{\left(\prod_{k=1}^{t-1} \pi_\theta(y_{w,k} | x, y_{w,<k}) \right)}_1 \cdot \underbrace{\pi_\theta(y_{\bar{w},t} | x, y_{w,<t})}_2.$$

Underbrace 1 represents the path down the tree of all possible sequences that follows the winning sequence, including a zero-length path. Underbrace 2 represents the probability of the next token not being a winner. As stated before, underbrace 2 is inclusive of all possible continuations after position t , since the probabilities of all continuations sum to 1.

6.1 Code

An implementation of the SLPO objective on a single sequence is provided below. The full implementation with type hints, comments, and vectorized batched operations is available at <https://github.com/yaoshiang/slpo>.

```
import torch
from torch.nn import functional as F

def _logdiffexp(t1, t2):
    retval = t1 + torch.log1p(-torch.exp(t2 - t1))
    retval[(t1 == float("-inf")) & (t2 == float("-inf"))] = float("-inf")
    return retval

def _logsumexp(t1, t2):
    return torch.logsumexp(torch.stack([t1, t2], dim=0), dim=0)

def y_ybar_single(logps, y_tokens):
    T, _ = logps.shape
    device = logps.device
    dtype = logps.dtype

    logp_y_t = logps[torch.arange(T), y_tokens]
    logp_y = logp_y_t.sum()

    zeros = torch.zeros(1, device=device, dtype=dtype)
    logp_y_until_t = torch.cat([zeros, logp_y_t[:-1]], dim=-1)
    logp_y_joint_until_t = logp_y_until_t.cumsum(dim=-1)

    logp_t = torch.logsumexp(logps, dim=-1)
    logp_y_bar_t = _logdiffexp(logp_t, logp_y_t)

    logp_not_y_through_t = logp_y_joint_until_t + logp_y_bar_t
    logp_not_y = torch.logsumexp(logp_not_y_through_t, dim=-1)

    return logp_y, logp_not_y

def slpo_loss_single(input, target):
    N = target["y"].size(0)
    input = input.to(torch.float64)

    logps = F.log_softmax(input, dim=-1)
    logp_y, logp_bar_y = y_ybar_single(logps, target["y"])
    logp_p = torch.log_softmax(torch.stack([logp_y, logp_bar_y]), dim=-1)

    if target["winner"].item():
```

```

logp_w_w = _logsumexp(
    target["logp_ref_w"].double(),
    target["logp_ref_l"].double()
)
logp_1_minus_w_w = torch.log(-torch.expm1(logp_w_w))
pre_rooted_logp_q = torch.stack([logp_w_w, logp_1_minus_w_w])
logp_q = torch.log_softmax(pre_rooted_logp_q / N, -1)

loss = F.kl_div(logp_p, logp_q, log_target=True, reduction="batchmean")
else:
    logp_q = torch.tensor([-torch.inf, 0.0]).double()
    loss = F.kl_div(logp_p, logp_q, log_target=True, reduction="batchmean")

return loss

```

7 Experiments

To establish a strong baseline for comparison with SLPO, we replicate the DPO experimental setup using the Anthropic HH dataset and Pythia 2.8B model. Following standard practice, we first perform supervised fine-tuning (SFT) on the chosen responses before alignment training.

We used the original DPO source code from <https://github.com/eric-mitchell/direct-preference-optimization> for all DPO experiments to ensure faithful reproduction of the baseline results.

7.1 Supervised Fine-Tuning

Since the original DPO paper does not specify the number of SFT training epochs, we trained Pythia 2.8B for 10 epochs on the Anthropic HH chosen responses and monitored the evaluation loss to identify the optimal checkpoint. As shown in Figure 2, the evaluation loss reached its minimum at step 497,232 (approximately 3 epochs) with a loss of 121.05. We use this checkpoint as the initialization for both DPO and SLPO training to ensure a fair comparison.

7.2 Alignment Training

For DPO training, we sweep beta values of 0.1, 0.3, and 0.5 as suggested by the original implementation, training for 10 epochs to allow thorough exploration of the hyperparameter space. SLPO training follows the same protocol for fair comparison.

8 Future Work

Most importantly, for future work we will complete the experimental comparison of SLPO against DPO following the setup outlined by Rafailov et al. (2023).



Figure 2: Evaluation loss during SFT training on Anthropic HH dataset. The checkpoint at step 497,232 (marked) was selected for subsequent DPO and SLPO training.

We will also explore custom loss functions that can bypass the autograd of PyTorch, since much of the reasoning behind the SLPO objective was in the gradient space, and only brought back into loss space for implementation.

9 Conclusion

In this paper, we have analyzed the DPO objective, found its strengths and weaknesses, and proposed a new objective that is a pure supervised learning approach to alignment.

Given the fact that DPO appeared to be superior to RLHF, if either a formal proof or empirical experiments can prove that the SLPO is superior to DPO, we would consider whether Reinforcement Learning concepts were ever necessary in alignment.

Appendix A. Equivalence of Cross-Entropy on Sigmoid of Difference and Categorical Cross-Entropy on Softmax for Two Classes

Logistic regression uses a single logit α , transforming it through the *sigmoid* function to obtain a probability $\sigma(\alpha) = \frac{1}{1+e^{-\alpha}}$. The *binary cross-entropy* (BCE) loss for a label y and predicted probability $\hat{y} = \sigma(\alpha)$ is:

$$\text{BCE}(y, \hat{y}) = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})].$$

Alternatively, single-label multiclass classification of two outcomes uses two softmax logits z_1 and z_2 (one per class) and applies the softmax function to obtain probabilities:

$$p(y = 1 \mid z_1, z_2) = \frac{e^{z_1}}{e^{z_1} + e^{z_2}}, \quad p(y = 2 \mid z_1, z_2) = \frac{e^{z_2}}{e^{z_1} + e^{z_2}}.$$

The associated categorical cross-entropy (CCE) loss is:

$$\text{CCE}(y, \hat{y}) = - \sum_{k=1}^2 y_k \log(\hat{y}_k),$$

where $\hat{y}_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2}}$, $\hat{y}_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2}}$, $y_1 = 1$, and $y_2 = 0$.

If we set

$$\alpha = z_1 - z_2,$$

then

$$\frac{e^{z_1}}{e^{z_1} + e^{z_2}} = \frac{1}{1 + e^{z_2 - z_1}} = \sigma(\alpha).$$

Therefore, the predicted probability \hat{y}_1 from softmax matches the $\sigma(\alpha)$ in logistic regression. This means that the BCE loss on the sigmoid of a difference of two numbers ($z_1 - z_2$) is equivalent to the CCE loss on the softmax of the first class (z_1) (Goodfellow et al., 2016).

Appendix B. Softmax as a sum as a joint probability

We prove this for the special case of two examples, each with N possible classes, where softmax probabilities are given by:

$$\hat{y}_j^1 = \frac{e^{z_j^1}}{\sum_m e^{z_m^1}}, \quad \hat{y}_k^2 = \frac{e^{z_k^2}}{\sum_n e^{z_n^2}}$$

The joint probability of selecting class j for the first example and class k for the second is:

$$P(j, k) = \hat{y}_j^1 \hat{y}_k^2$$

We define the target joint probability as:

$$y_{\text{true},j,k} = C, \quad \sum_{(m,n) \neq (j,k)} y_{\text{true},m,n} = 1 - C.$$

We aim to show that the cross-entropy loss:

$$L = -C \log(\hat{y}_j^1 \cdot \hat{y}_k^2) - (1 - C) \log(\Omega)$$

where Ω represents all other possible outcomes, a set with cardinality $C^2 - 1$. can be rewritten as the categorical cross-entropy of a softmax over two terms. Define the logits for two terms:

$$\begin{aligned} s_1 &= \log \hat{y}_j^1 + \log \hat{y}_k^2 = z_j^1 - \log \sum_m e^{z_m^1} + z_k^2 - \log \sum_n e^{z_n^2}, \\ s_2 &= \log(\Omega). \end{aligned}$$

Applying the softmax function to these logits:

$$p_1 = \frac{e^{s_1}}{e^{s_1} + e^{s_2}}, \quad p_2 = \frac{e^{s_2}}{e^{s_1} + e^{s_2}}.$$

Since $e^{s_1} = \hat{y}_j^1 \hat{y}_k^2$ and $e^{s_2} = \Omega$, we obtain:

$$p_1 = \frac{\hat{y}_j^1 \hat{y}_k^2}{\hat{y}_j^1 \hat{y}_k^2 + \Omega}, \quad p_2 = \frac{\Omega}{\hat{y}_j^1 \hat{y}_k^2 + \Omega}.$$

The categorical cross-entropy loss with target distribution $[C, 1 - C]$ is:

$$\begin{aligned} L &= -C \log p_1 - (1 - C) \log p_2 \\ &= -C \log \left(\frac{\hat{y}_j^1 \hat{y}_k^2}{\hat{y}_j^1 \hat{y}_k^2 + \Omega} \right) - (1 - C) \log \left(\frac{\Omega}{\hat{y}_j^1 \hat{y}_k^2 + \Omega} \right). \end{aligned}$$

Expanding the logarithms:

$$\begin{aligned} L &= -C (\log(\hat{y}_j^1 \hat{y}_k^2) - \log(\hat{y}_j^1 \hat{y}_k^2 + \Omega)) \\ &\quad - (1 - C) (\log(\Omega) - \log(\hat{y}_j^1 \hat{y}_k^2 + \Omega)) \\ &= -C \log(\hat{y}_j^1 \hat{y}_k^2) - (1 - C) \log(\Omega). \end{aligned}$$

Thus, we have rewritten the original loss function as the categorical cross-entropy of a softmax over two terms:

$$L = \text{CCE}(\text{softmax}([s_1, s_2]), y_{\text{true}} = [C, 1 - C]).$$

Appendix C. Approximate Equivalence of Softmax of Sum and Temperature Scaling

Here we show that the two-class softmax of a sum of inputs $\alpha_0 \dots \alpha_T$, holding all other inputs constants, is approximated by temperature scaling, under simplifying assumptions:

- We consider the case where the first class is true and the second class is false.

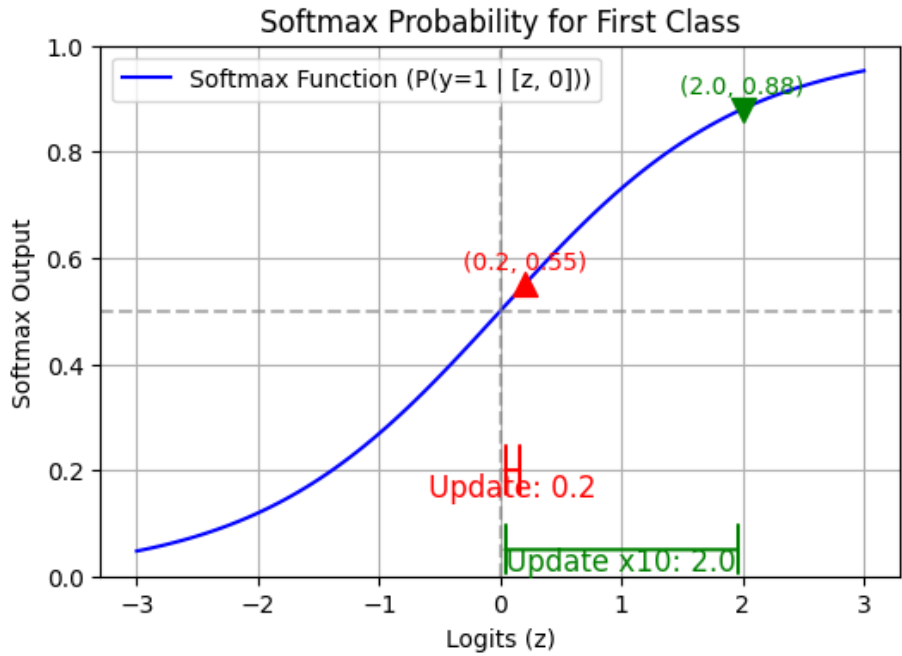


Figure 3: The red update represent the actual update to a logit. But on the next iteration of backprop, the gradient is not calculated based on the sigmoid of logit itself, but the sigmoid of the sum of the logit and nine others, which we assume to be approximately equal. This pushes the gradient calculation out to a less steep part of the curve. This "step function" like behavior is similar to a low temperature scaling of a logit during training.

- we assume the output of softmax is passed to Categorical Crossentropy (CCE) loss.
- We assume the second logit to the softmax is zero.
- We assume the inputs are approximately equal to each other.
- We assume the many inputs, such as 100.

The softmax of the sum of T inputs is:

$$\text{softmax}([\alpha_1 + \alpha_2 + \dots + \alpha_T, 0]) = \frac{e^{\alpha_1 + \alpha_2 + \dots + \alpha_T}}{e^{\alpha_1 + \alpha_2 + \dots + \alpha_T} + e^0}$$

The gradient of the softmax function with a CCE loss is well known to be:

$$\frac{\partial \text{CCE}}{\partial z_i} = \hat{y}_i - y_i.$$

Equivalently, since \hat{y}_i is the softmax of the sum of the inputs,

$$\frac{\partial \text{CCE}}{\partial z_i} = \text{Softmax}(z_i) - y_i,$$

where z_i is the input to softmax. Applying this to our unusual softmax of a sum formulation by setting $z_i = \sum_{i=1}^T \alpha_i$:

$$\frac{\partial \text{CCE}}{\partial z_i} = \text{Softmax}([\sum_{i=1}^T \alpha_i, 0]) - y_i.$$

The nuance here is that due to the chain rule, the gradient on the sum of the inputs is passed to each input independently. This is the heart of the "multiplication" effect which we will soon see is equivalent to temperature scaling. For every input α_k , the gradient is:

$$\frac{\partial \text{CCE}}{\partial \alpha_k} = \text{Softmax}([\sum_{i=1}^T \alpha_i, 0]) - y_i.$$

But since we assumed all the inputs are approximately equal to each other, this is equivalent to:

$$\frac{\partial \text{CCE}}{\partial \alpha_k} = \text{Softmax}([T \cdot \alpha_k, 0]) - y_i.$$

for all k . Rewriting this with the definition of softmax

$$\frac{\partial \text{CCE}}{\partial \alpha_k} = \frac{e^{T \cdot \alpha_k}}{e^{T \cdot \alpha_k} + e^0} - y_i. \quad (11)$$

shows a multiplier on the input α_k of T , representing the cardinality of the elements of the sum (e.g. tokens in a sequence).

Now let's turn to the temperature scaling. The canonical definition of temperature scaling per (Hinton et al., 2015, Equation 2) is the following, where we use the term temp^{-1} to distinguish it from our use of T as the number of inputs:

$$\frac{\partial \mathcal{C}}{\partial z_i} = \text{temp}^{-1} \left(\frac{e^{z_i \cdot \text{temp}^{-1}}}{\sum_j e^{z_j \cdot \text{temp}^{-1}}} - \frac{e^{v_i \cdot \text{temp}^{-1}}}{\sum_j e^{v_j \cdot \text{temp}^{-1}}} \right)$$

where v_i are logits of the teacher and z_i are the logits of the student. Equivalently:

$$\frac{\partial \mathcal{C}}{\partial z_i} = \left(\frac{e^{z_i \cdot \text{temp}^{-1} - \log(\text{temp})}}{\sum_j e^{z_j \cdot \text{temp}^{-1} - \log(\text{temp})}} - \frac{e^{v_i \cdot \text{temp}^{-1} - \log(\text{temp})}}{\sum_j e^{v_j \cdot \text{temp}^{-1} - \log(\text{temp})}} \right).$$

Considering we only have two classes, and the second class has a fixed logit of zero, we can simplify this to:

$$\frac{\partial \mathcal{C}}{\partial z_1} = \left(\frac{e^{z_1 \cdot \text{temp}^{-1} - \log(\text{temp})}}{e^{z_1 \cdot \text{temp}^{-1} - \log(\text{temp})} + e^{-\log(\text{temp})}} - \frac{e^{v_1 \cdot \text{temp}^{-1} - \log(\text{temp})}}{e^{v_1 \cdot \text{temp}^{-1} - \log(\text{temp})} + e^{-\log(\text{temp})}} \right),$$

As the logits are optimized up to reasonable non-zero values such as 1.0, and with the assumption that the term temp^{-1} is large (e.g. the number of tokens in a sequence such as 100), the terms $\log(\text{temp})$ is dominated by the term temp^{-1} and the term $.$ We can replace the terms that use v with y_i since we are using hard targets, not soft labels from a teacher. Hence, we can approximate the gradient under KD as

$$\frac{\partial \mathcal{C}}{\partial z_i} = \left(\frac{e^{z_i \cdot \text{temp}^{-1}}}{e^{z_i \cdot \text{temp}^{-1} + C}} - y_i \right) \quad (12)$$

where C is a small constant that is dominated by the term $e^{z_1 \cdot \text{temp}^{-1}}$. For example, at reasonable values of $z_1 = 1.0$ and $\text{temp}^{-1} = 100$, the term $e^{z_1 \cdot \text{temp}^{-1}} = e^{100}$, while the term $e^{-\log(\text{temp})} = 1/\text{temp} = 0.01$.

This exactly of the form of Equation 11 with the temperature scaling term temp^{-1} replacing the number of inputs (such as the number of tokens in a sequence) T . Hence, as optimization proceeds, a softmax of a sum of many inputs becomes similar to optimizing the softmax of each input with a low temperature.

See Figure C for an illustration of this concept.

References

- R. A. Bradley and M. E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. doi: 10.2307/2334029. URL <https://doi.org/10.2307/2334029>.
- Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E Gonzalez, et al. Chatbot arena: An open platform for evaluating llms by human preference. *arXiv preprint arXiv:2403.04132*, 2024.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

- Arpad E. Elo. *The Rating of Chessplayers, Past and Present*. Arco Publishing, New York, 1978. ISBN 9780668047210.
- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. URL <http://www.deeplearningbook.org>.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Dark knowledge. Presented at the University of Chicago, 2014. URL <https://www.ttic.edu/dl/dark14.pdf>.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. URL <https://arxiv.org/abs/1503.02531>.
- Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Bo Pang and Lillian Lee. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the ACL*, 2004.
- 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 A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 53728–53741. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/a85b405ed65c6477a4fe8302b5e06ce7-Paper-Conference.pdf.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826, 2016.
- Zhichao Wang, Bin Bi, Shiva Kumar Pentiyala, Kiran Ramnath, Sougata Chaudhuri, Shubham Mehrotra, Zixu, Zhu, Xiang-Bo Mao, Sitaram Asur, Na, and Cheng. A comprehensive survey of llm alignment techniques: Rlhf, rlaif, ppo, dpo and more, 2024. URL <https://arxiv.org/abs/2407.16216>.