

On the Effect of Negative Gradient in Group Relative Deep **Reinforcement Optimization**





Wenlong Deng^{1,2}, Yi Ren¹, Muchen Li¹, Danica J. Sutherland^{1,3}, Xiaoxiao Li^{1,2*} and Christos Thrampoulidis^{1*}

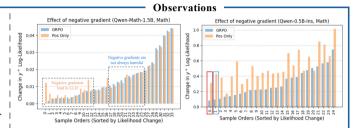




Motivation

Background: We identify Lazy Likelihood Displacement (LLD), wherein the likelihood of correct responses marginally increases or even decreases during training.

Related Work: Several approaches have been proposed to mitigate the reduced probabilities of preferred responses in DPO, which is undesirable as it often degrades model performance by diverting probability mass away from optimal response.



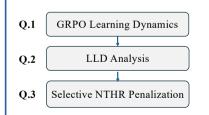
Definition 4.1 Let $\pi_{\theta_{max}}$ and $\pi_{\theta_{max}}$ denote the initial and final language models, before and after optimizing a preference learning objective \mathcal{J} (e.g., Eq. $\overline{(1)}$) over a dataset \mathcal{D} , such that $\mathcal{J}(\theta_{fin})$ < $\mathcal{J}(\theta_{init})$. We say that **LLD** occurs for a tuple $(\mathbf{x}, \mathbf{y}^+) \in \mathcal{D}$ if, for small nonnegative constant $\epsilon \geq 0$, $\ln \pi_{\theta_{e_{-}}}(\mathbf{y}^{+}|\mathbf{x}) < \ln \pi_{\theta_{e_{-}}}(\mathbf{y}^{+}|\mathbf{x}) + \epsilon$.

Focus & Contribution

Questions:

- (1) Why does GRPO training sometimes reduce the probability of correct responses?
- (2) Can we identify the sources of LLD?
- (3) Can we mitigate this misalignment without losing data efficiency?

Our Contribution:



Q.1 Theory Undersntanding

GRPO training's influence of the likelihood change of correct responses:

$$rac{d}{dt} \ln \pi_{ heta(t)}(oldsymbol{y}_i^+ | oldsymbol{x}) = \left\langle
abla \ln \pi_{ heta(t)}(oldsymbol{y}_i^+ | oldsymbol{x}), \; rac{d heta(t)}{dt}
ight
angle$$

Theorem (GRPO learning dynamics, informal): The likelihood change of **correct responses** will be *lazy* if the following quantity *increases*

$$p^{-} \sum_{k=1}^{|\boldsymbol{y}_{i}^{+}|} \sum_{j=1}^{N^{-}} \sum_{k'=1}^{|\boldsymbol{y}_{j}^{-}|} \alpha_{k,k'}^{-} \cdot \langle \mathbf{h}_{\mathbf{x},\boldsymbol{y}_{i,< k}^{+}}, \mathbf{h}_{\mathbf{x},\boldsymbol{y}_{j,< k'}^{-}} \rangle - p^{+} \sum_{k=1}^{|\boldsymbol{y}_{i}^{+}|} \sum_{i'=1}^{N^{+}} \sum_{k''=1}^{|\boldsymbol{y}_{i'}^{+}|} \alpha_{k,k''}^{+} \cdot \langle \mathbf{h}_{\mathbf{x},\boldsymbol{y}_{i,< k}^{+}}, \mathbf{h}_{\mathbf{x},\boldsymbol{y}_{i',< k''}^{+}} \rangle.$$
 (5)

Illustration of sample highlighted by the red block in first figure. Mostly, LLD

samples are:

- (1) Almost Correct
- (2) Format Issue

Math Question: Find the greatest integer value of b for which the expression has a domain of all real numbers.

Ground-truth Answer: 5

Sample: [Math dataset], Model: [Owen-0.5B-Ins]

Answer: To solve this problem, we need to ... Now, we find the approximate value of $\sqrt{32}$: $\sqrt{32} \approx 5.656$ So, we have: b > 5.6565) Identify the greatest integer value of b: The greatest integer greater than 5.656 is 6. Therefore, the greatest integer value of \$b\$ for which the expression

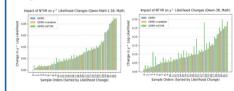
 $9x^3 + 4x^2 + 11x + 7$ has a domain of all real numbers is 6.

O.2 &O.3 Analysis and NTHR

Top-K	Qwen-1.5B-deepseek		Qwen-1.5B-math		
_	GWHES	Random	GWHES	Random	
10	50%	17.5%	60%	21.3%	
15	75%	26.3%	75%	31.9%	

Using GWHES (Eq. (5)) score can effectively identify LLD samples compared with random selection.

NTHR Selective Penalization: We use GWHES score and reduce the penalty for negative tokens whose GWHES scores exceed a given threshold.



Effectively mitigates the LLD issue: green lines show a larger increased likelihood

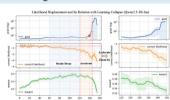
Base model + Method	AIME24	AMC	MATH500	Minerva	Olympiad	Avg.
Qwen2.5-Math-1.5B						
Base	3.3	20.0	39.6	7.7	24.9	19.10
GRPO	13.3	57.5	71.8	29.0	34.1	41.14
Pos Only	10.0	57.5	70.6	30.1	31.0	39.84
NTHR	16.7	57.5	70.8	30.5	34.2	41.94
Qwen2.5-0.5B-Ins						
Base	0.0	2.5	33.4	4.4	7.0	9.46
GRPO	0.0	7.5	33.8	9.2	8.1	11.72
NTHR	0.0	10.0	36.6	8.1	8.6	12.66
Qwen2.5-1.5B-Ins						
Base	0.0	22.5	53.0	19.1	20.7	23.06
GRPO	3.3	32.5	57.2	18.8	23.0	26.96
NTHR	6.7	35.0	58.8	21.0	20.9	28.48
Qwen2.5-Math-1.5B (deepscaler)						
Base	3.3	20.0	39.6	7.7	24.9	19.10
GRPO	10.0	42.5	72.4	32.4	31.9	37.80
NTHR	16.7	47.5	73.2	29.4	31.4	39.60
Qwen2.5-3B						
Base	10.0	37.5	58.6	26.1	24.6	31.36
GRPO	6.7	35.0	66.6	31.2	29.9	33.88
NTHR	10.0	47.5	65.6	31.6	26.8	36.30
				_		

Improved greedy decoding performance.

Borader Impact -

We take this opportunity to highlight two of our recent related works.

[1] LLD is the source of GRPO collapse in Tool-integrated mult-turn RL.



Case study using Search-R1, demonstrating strong LLD.



[2] We unify direct optimization and advantage shaping of recent GRPOvariants.

Algorithm	Advantage Scores (A^+, A^-)	Weighted Empirical Gradient	Population Surrogate Reward
Targets 0/1 coalastic	on (K=1)		
RLOO [KHW19]	$\left(\frac{N(1-\beta)}{N-1}, -\frac{N\beta}{N-1}\right)$	$\hat{\rho}(1 - \hat{\rho})[\hat{\nabla}_{+} - \hat{\nabla}_{-}]$	ρ
GRPO [Sha+24]	(√1 2 2√ 1 2)	$\sqrt{\hat{\rho}(1-\hat{\rho})} [\tilde{\nabla}_{+} - \tilde{\nabla}_{-}]$	$2 \arcsin(\sqrt{\rho})$
Skew-R [This work] ¹	$((1-\hat{\rho})\sqrt{\frac{1-\hat{\rho}}{\hat{\rho}}}, -(1-\hat{\rho})\sqrt{\frac{\hat{\rho}}{1-\hat{\rho}}})$	$(1 - \bar{\rho}) \sqrt{\bar{\rho}(1 - \bar{\rho})} [\hat{\nabla}_{+} - \hat{\nabla}_{-}]$	$\operatorname{arcsin}(\sqrt{\rho}) + \sqrt{\rho(1-\rho)}$
Targets Pass@K con	bation (A≥2)		
$RLOO_N$ [This work]	$(\hat{f}_{K-1}^+, \frac{N(1-\hat{p})}{N-1}, -\hat{f}_{K-1}^-, \frac{N\hat{p}}{N-1})^2$	$\hat{f}_{K-1}^{+} \hat{\rho}(1-\hat{\rho}) \hat{\nabla}_{+} - \hat{f}_{K-1}^{-} \hat{\rho}(1-\hat{\rho}) \hat{\nabla}_{-}$	PK
$\widetilde{\mathrm{GRPO}_K}$ [Che+25a]	$\left(\widetilde{\omega}_{K}\sqrt{\frac{k-\delta}{\delta}}, -\widetilde{\omega}_{K}\sqrt{\frac{\delta}{1-\delta}}\right)^{3}$	$\sqrt{\frac{1-\hat{\rho}_N}{\hat{\rho}_N}} \hat{\rho} [\hat{\nabla}_+ - \hat{\nabla}]$	$\frac{2}{K} \arcsin(\sqrt{\rho_K})$
$GRPO_K$ [This work]	$(\hat{f}_{K-1}^+, \sqrt{\frac{1-\mu}{4}}, -\hat{f}_{K-1}^-, \sqrt{\frac{\mu}{4-3}})$	$\hat{f}_{K-1}^{+}\sqrt{\hat{\rho}(1-\hat{\rho})}\hat{\nabla}_{+} - \hat{f}_{K-1}^{-}\sqrt{\hat{\rho}(1-\hat{\rho})}\hat{\nabla}_{-}$	$B(1 - (1 - \rho_K)^{1/K}; \frac{1}{2}, K - \frac{1}{2})^{-4}$

[1] On GRPO Collapse in Search-R1: The Lazy Likelihood-Displacement Death Spiral.

[2] Advantage Shaping as Surrogate Reward Maximization: Unifying Pass@K Policy Gradients



