

Constrained Policy Optimization for Large Language Model Alignment

Dongsheng Ding

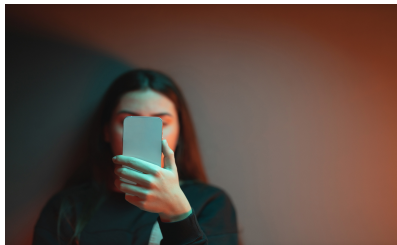


THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

2025 SIAM-NNP, State College, PA; Nov. 1, 2025

Motivating application

■ AI THERAPY CHATBOTS



Stanford HAI

maximize helpfulness
LLM policy

subject to $\text{harmlessness} \geq \text{margin}$

CHALLENGE: Constraint satisfaction

Context

■ POWERFUL CAPABILITIES

question-answering, summarization, translation, etc.

■ LESSONS LEARNED

- ★ importance of **policy optimization** in alignment
reward/preference-based
- ★ **multidimensionality** of human preference; **multi-objective** alignment
- ★ **convex** in distribution space; **nonconvex** in parameter space

■ WHAT NOW ?

- ★ **safety-critical domains**: healthcare, robotics, finance
- ★ **constrained policy optimization**: tremendous advances

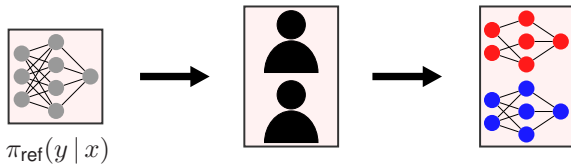
OBJECTIVE

Find a **Large Language Model (LLM)** that
maximizes a performance metric
subject to a constraint on
another performance metric

Alignment framework

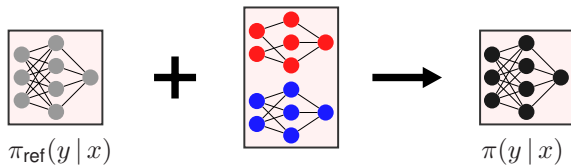
■ REINFORCEMENT LEARNING FROM HUMAN FEEDBACK

★ reward modeling



$r(x, y)$, $g(x, y)$ – reward/utility models

★ policy optimization



$\pi_{\text{ref}}, \pi: \mathcal{X}$ (prompts) $\rightarrow \mathcal{Y}$ (responses) – LLM policies

Constrained alignment problem

$$\underset{\pi}{\text{maximize}} \quad \mathbb{E}_x [\mathbb{E}_{y \sim \pi} [r(x, y)] - \beta D_{\text{KL}}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))]$$

$$\text{subject to} \quad \mathbb{E}_x [\mathbb{E}_{y \sim \pi} [g(x, y)]] \geq 0$$



KL-regularized objective



policy constraint

- ★ limit the policy space to **an inequality constraint**

Convex constrained policy optimization \rightarrow **Strong duality**

Lagrangian approach

■ LAGRANGIAN

$$L(\pi, \lambda) = \mathbb{E}_x [\mathbb{E}_{y \sim \pi} [r(x, y) + \lambda g(x, y)] - \beta D_{\text{KL}}(\pi(\cdot | x) \| \pi_{\text{ref}}(\cdot | x))]$$

- ★ penalize violation via dual variable $\lambda \geq 0$

■ LAGRANGIAN RELAXATION

$$\underset{\pi}{\text{maximize}} \quad L(\pi, \lambda)$$

convex conjugate

- ★ exponentially tilted distribution $\pi^*(\cdot | x; \lambda) \propto \pi_{\text{ref}}(\cdot | x) e^{(r(x, \cdot) + \lambda g(x, \cdot))/\beta}$
- ★ existence of an optimal dual variable λ^*

Dual problem

■ DUAL FUNCTION

$$D(\lambda) := \underset{\pi}{\text{maximize}} \quad L(\pi, \lambda)$$

$$= \beta \mathbb{E}_x \left[\log \left(\sum_y \pi_{\text{ref}}(y | x) e^{(r(x,y) + \lambda g(x,y)) / \beta} \right) \right]$$

cumulant-generating function

★ **convex**, **smooth**, and **local strongly convex** function

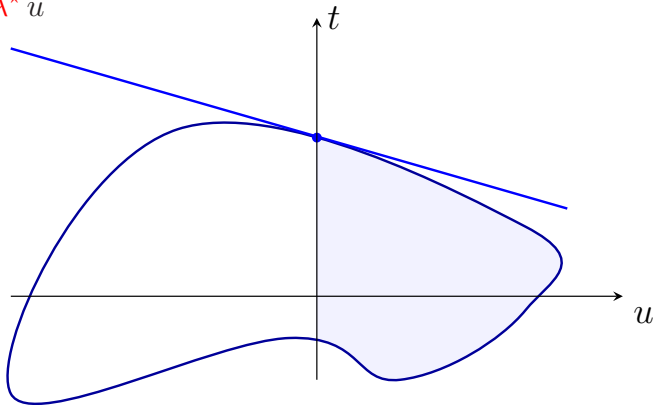
■ DUAL PROBLEM

$$\underset{\lambda \geq 0}{\text{minimize}} \quad D(\lambda)$$

Gradient descent finds **an optimal dual variable** λ^*

■ GEOMETRIC INTERPRETATION OF STRONG DUALITY

$$D^* = t + \lambda^* u$$



$$\text{image } \mathcal{E} = \{ (g(\pi), r_{\text{KL}}(\pi)) \mid \pi \}$$

Optimal hyperplane touches \mathcal{E} at an optimal policy: $D^* = r_{\text{KL}}(\pi^*) := P^*$

Dualization-based alignment

$H^\dagger L^\dagger D B L H D^\dagger$, NeurIPS '24

■ STAGE #1: FINDING AN OPTIMAL DUAL VARIABLE

$$\lambda^* = \operatorname{argmin}_{\lambda \geq 0} D(\lambda)$$

smooth convex optimization

■ STAGE #2: SEARCHING FOR AN LLM POLICY

$$\pi^* = \operatorname{argmax}_{\pi} L(\pi, \lambda^*)$$

unconstrained alignment

Computational efficiency

Constrained parameter optimization

$$\underset{\theta}{\text{maximize}} \quad \mathbb{E}_x \left[\mathbb{E}_{y \sim \pi_\theta} [r(x, y)] - \beta D_{\text{KL}}(\pi_\theta(\cdot | x) \| \pi_{\text{ref}}(\cdot | x)) \right]$$

$$\text{subject to} \quad \mathbb{E}_x \left[\mathbb{E}_{y \sim \pi_\theta} [g(x, y)] \right] \geq 0$$



KL-regularized objective



policy constraint

★ search model parameter θ

CHALLENGE

Nonconvex constrained optimization → **Lack of strong duality**

Overview of our results

ZLHBD[†]R, NeurIPS '25

■ ITERATIVE DUALIZATION-BASED ALIGNMENT

- ★ duality gap

- ★ optimality gap

objective and constraint

Dual methods find **an optimal constrained LLM policy**,
up to **a parametrization gap**

Iterative dualization-based alignment

Parametrized dual problem

■ DUAL FUNCTION

$$D_p(\lambda) := \max_{\theta} L(\pi_{\theta}, \lambda)$$

★ convex, and nondifferentiable function

■ DUAL PROBLEM

$$\min_{\lambda \geq 0} D_p(\lambda)$$

Subgradient descent finds **an optimal parametrized dual variable** λ_p^*

Iterative dualization-based alignment

Lagrangian maximizer: $\theta^*(\lambda) \in \operatorname{argmax}_{\theta} L(\pi_{\theta}, \lambda)$

Subgradient: $u(\lambda) = \nabla_{\lambda} L(\pi_{\theta}, \lambda) |_{\theta = \theta^*(\lambda)}$

■ SUBGRADIENT DESCENT

$$\lambda(t+1) = [\lambda(t) - \eta u(\lambda(t))]_{+}$$

★ explicit subgradient $u(\lambda(t)) = \mathbb{E}_x[\mathbb{E}_{y \sim \pi_{\theta^*(\lambda(t))}}[g(x, y)]]$

QUESTION: Optimality of λ_p^* -recovered model $\pi_{\theta^*(\lambda_p^*)} := \pi_p^*(\lambda_p^*)$?

Duality gap

$$\text{Duality gap: } |P^* - D_p^*|$$

Theorem (informal)

★ **Duality gap** is dominated by

ν

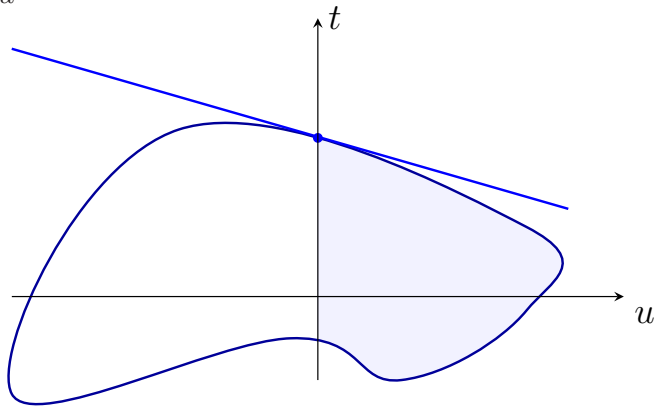
parametrization gap $\nu := \max_{\pi} \min_{\theta} \text{dist}_1(\pi, \pi_{\theta})$

★ ν -parametrization gap yields ν -duality gap

linear independence

■ GEOMETRIC INTERPRETATION OF DUALITY GAP

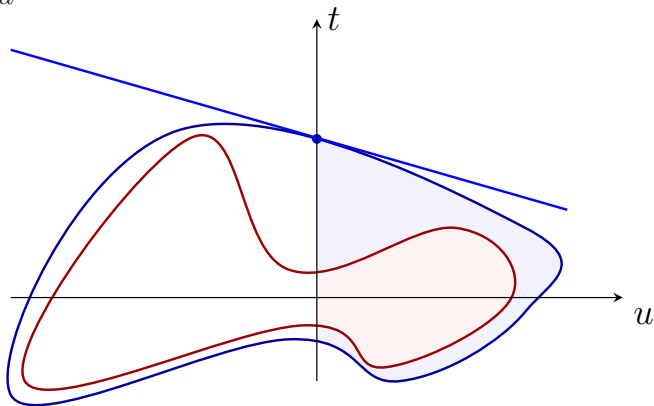
$$D^* = t + \lambda^* u$$



$$\text{image } \mathcal{E} = \{ (g(\pi), r_{\text{KL}}(\pi)) \mid \pi \}$$

■ GEOMETRIC INTERPRETATION OF DUALITY GAP

$$D^* = t + \lambda^* u$$



$$\text{image } \mathcal{E} = \{ (g(\pi), r_{\text{KL}}(\pi)) \mid \pi \}$$

$$\text{image } \mathcal{E}_p = \{ (g(\pi_\theta), r_{\text{KL}}(\pi_\theta)) \mid \theta \}$$

■ GEOMETRIC INTERPRETATION OF DUALITY GAP

$$D^* = t + \lambda^* u$$

$$D_p^* = t + \lambda_p^* u$$

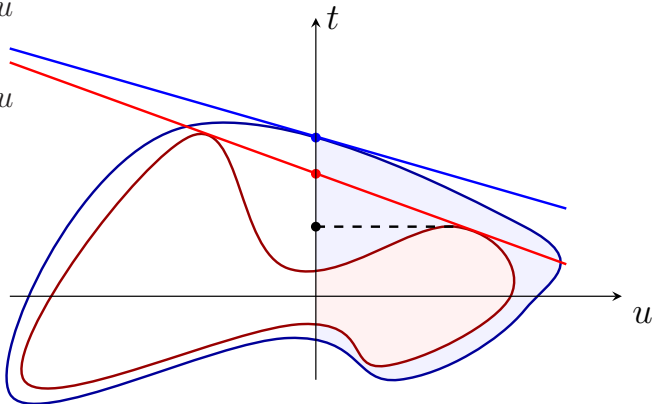


image $\mathcal{E} = \{ (g(\pi), r_{\text{KL}}(\pi)) \mid \pi \}$

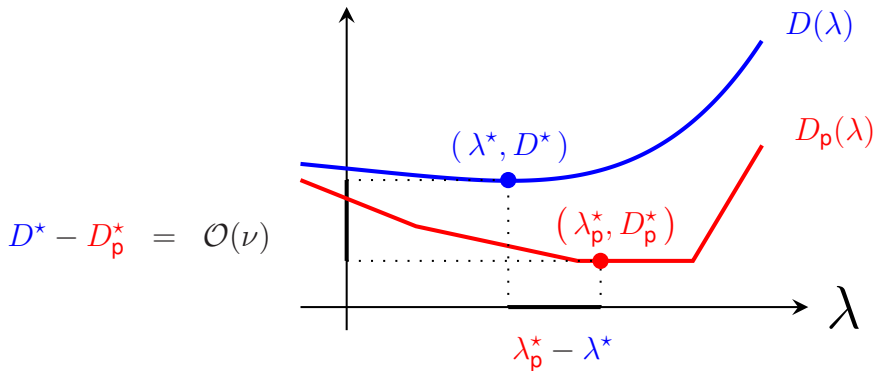
image $\mathcal{E}_p = \{ (g(\pi_\theta), r_{\text{KL}}(\pi_\theta)) \mid \theta \}$

Optimal hyperplane touches \mathcal{E}_p w/ t -intercept D_p^* :

$$D^* - D_p^* = \mathcal{O}(\nu)$$

Gap between optimal dual variables

■ GAP BETWEEN (UN)PARAMETRIZED DUAL FUNCTIONS



Optimal dual variables: λ^* , λ_p^* are close:

$$\|\lambda^* - \lambda_p^*\| = \mathcal{O}(\sqrt{\nu})$$

Optimality gap

Objective optimality: $\left| r_{\text{KL}}(\pi_{\text{p}}^*(\lambda_{\text{p}}^*)) - r_{\text{KL}}(\pi^*) \right|$

Constraint feasibility: $\left| g(\pi_{\text{p}}^*(\lambda_{\text{p}}^*)) - g(\pi^*) \right|$

Implication (informal)

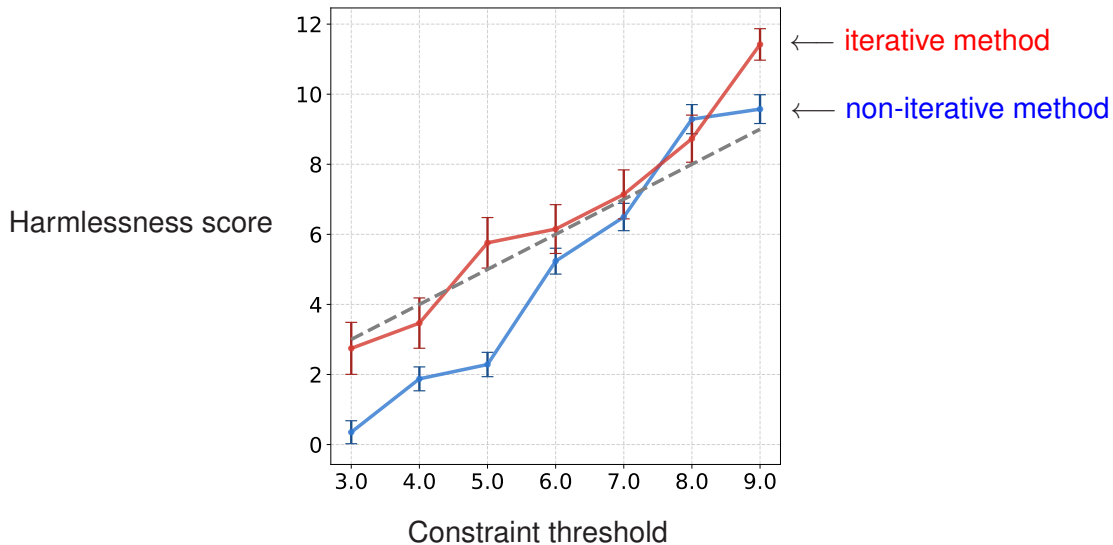
★ **Objective optimality & Constraint feasibility** are dominated by

$$\sqrt{\nu}$$

parametrization gap $\nu := \max_{\pi} \min_{\theta} \text{dist}_1(\pi, \pi_{\theta})$

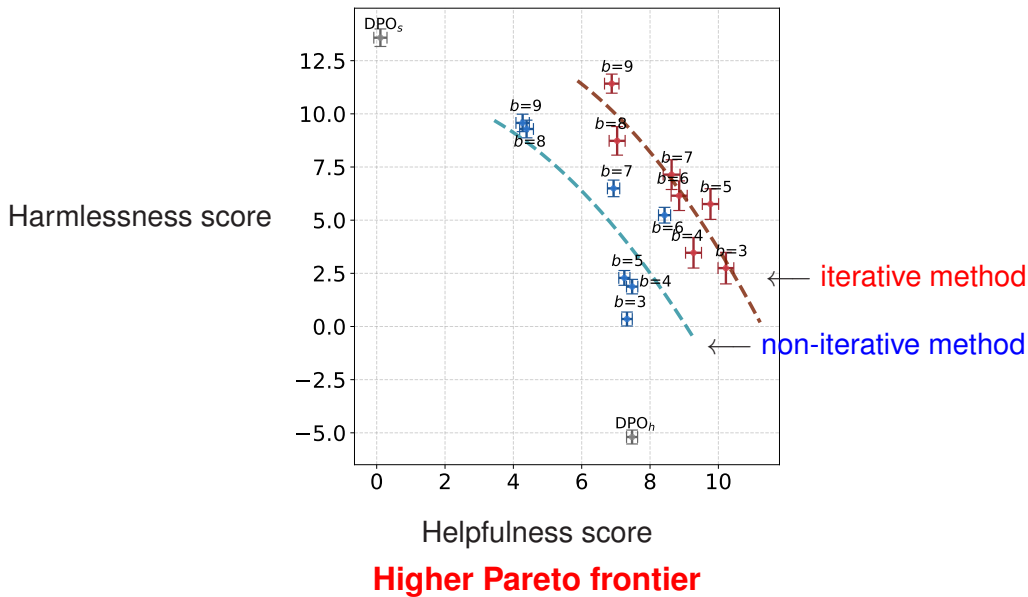
Root-scaling of parametrization gap

Constraint satisfaction



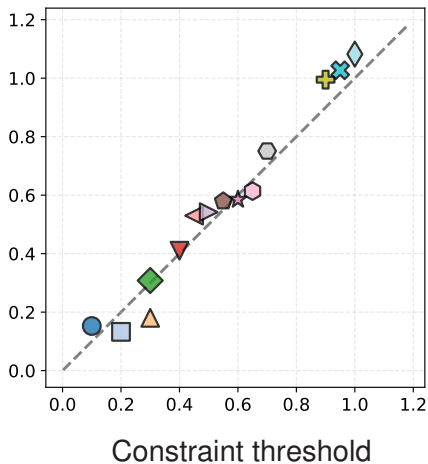
Better constraint satisfaction

Helpfulness and Harmlessness tradeoff

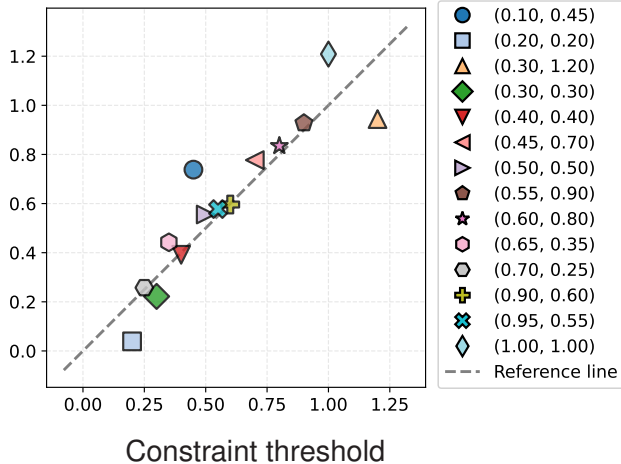


Harmlessness and Humor constraints

Harmlessness score



Humor score



Multi-constraint satisfaction

Summary

ZLHB \mathbf{D}^\dagger R, NeurIPS '25

H † L † DBLH \mathbf{D}^\dagger , NeurIPS '24

■ DUALIZATION-BASED ALIGNMENT METHODS

- ★ iterative dualization-based alignment
- ★ duality gap and optimality gap

■ ON-GOING EFFORTS

- ★ sample complexity and convergence
- ★ other complex constraints

Thank you for your attention.