Constrained Policy Optimization for Large Language Model Alignment

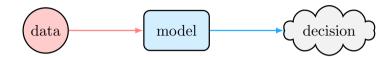
Dongsheng Ding

dongshed@utk.edu

EECS

Statistics and Data Science Seminar, UTK
November 13, 2025

The reality of machine learning



goal - loss / reward / likelihood

■ REQUIREMENTS

harmlessness



safety



robustness



fairness



Image sources: Stanford HAI, Waymo, WIRED, NBC

The risk of machine learning

Al chatbots might be sabotaging women by advising them to ask for lower salaries, study says

New York Post, JUL 29, 2025

NHTSA probes Waymo self-driving cars over school bus safety concerns

Reuters, OCT 21, 2025

Al-Powered Robots Can Be Tricked Into Acts of Violence
WIRED, DEC 4, 2024

Study reveals why Al models that analyze medical images can be biased

MIT News, JUN 28, 2024

Requirement-driven machine learning

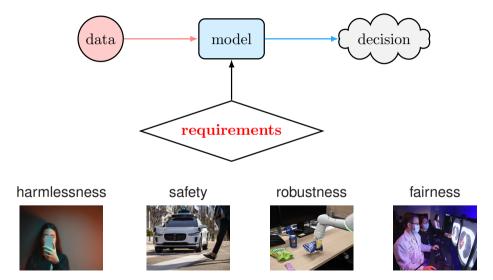


Image sources: Stanford HAI, Waymo, WIRED, NBC

Motivating application: Robotics

■ LLM-CONTROLLED ROBOTS



Figure AI

 $\underset{\mathsf{LLM}}{\operatorname{maximize}} \quad \text{dos} \quad$

subject to don'ts > threshold

Motivating application: Healthcare

■ ALTHERAPY CHATBOTS



Stanford HAI

maximize LLM policy

helpfulness

subject to harmlessness > threshold

REAL-WORLD CHALLENGE

Constraint satisfaction

OBJECTIVE

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

Outline

- CONSTRAINED LLM ALIGNMENT
 - * constrained policy optimization
- ALIGNMENT METHOD & THEORY
 - * non-iterative & iterative methods
 - * duality gap & optimality gap
- EMPIRICAL STUDY
 - * safety-alignment task
- SUMMARY & OUTLOOK

CONSTRAINED LLM ALIGNMENT

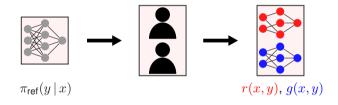
constrained policy optimization

OTTAINED ELM ALIGINMENT

Alignment framework

■ REINFORCEMENT LEARNING FROM HUMAN FEEDBACK

* reward modeling



 $\pi_{\mathsf{ref}} \colon \mathcal{X} \ (\mathsf{prompts}) o \mathcal{Y} \ (\mathsf{responses}) \ - \ \mathsf{reference} \ \mathsf{LLM} \ \mathsf{policy}$

r(x, y), g(x, y) - reward/utility models

* reward/utility models

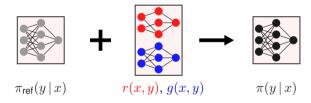
e.g., SafeRLHF: helpfulness and harmlessness

Dai et al., ICLR '24

Alignment framework

■ REINFORCEMENT LEARNING FROM HUMAN FEEDBACK

* policy optimization



 $\pi_{\mathsf{ref}} \colon \mathcal{X} \ (\mathsf{prompts}) o \mathcal{Y} \ (\mathsf{responses}) \ - \ \mathsf{reference} \ \mathsf{LLM} \ \mathsf{policy}$

 π : \mathcal{X} (prompts) $\to \mathcal{Y}$ (responses) – aligned LLM policy

e.g., direct preference optimization (preference-based)

Rafailov et al., NeurIPS '23

Response space

■ RESPONSE SPACE SIZE

e.g., ChatGPT-3.5 / Llama 3:
$$4000^{30}$$
 $\approx 10^{108}$ $\gg 10^{80}$

exponentially large decision space

Constrained alignment problem

$$\begin{array}{ll} \underset{\pi}{\text{maximize}} & \mathbb{E}_{x} \left[\mathbb{E}_{y \sim \pi} \left[r(x,y) \right] - \beta \, D_{\mathsf{KL}}(\pi(\cdot \,|\, x) \, \|\, \pi_{\mathsf{ref}}(\cdot \,|\, x)) \right] \\ \\ \text{subject to} & \mathbb{E}_{x} \left[\mathbb{E}_{y \sim \pi} \left[g(x,y) \right] \right] \, \geq \, 0 \\ \\ \downarrow & \mathsf{KL-regularized objective} \\ \end{aligned}$$

policy constraint

* limit the policy space to an inequality constraint

e.g., harmless policy, safe policy

Constrained policy optimization

policy constraint

- * no transition dynamics
- * concave KL-regularized objective and linear constraint

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

Lagrangian relaxation

■ LAGRANGIAN

$$L(\pi, \lambda) = \mathbb{E}_x \left[\mathbb{E}_{y \sim \pi} \left[r(x, y) + \frac{\lambda}{\lambda} g(x, y) \right] - \beta D_{\mathsf{KL}} (\pi(\cdot \mid x) \parallel \pi_{\mathsf{ref}}(\cdot \mid x)) \right]$$

* penalize violation via dual variable $\lambda \geq 0$

■ LAGRANGIAN MAXIMIZATION

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

convex conjugate

* exponentially tilted distribution $\pi^*(\cdot \mid x; \lambda)$

$$\pi^{\star}(\cdot \mid x; \frac{\lambda}{\lambda}) \propto \pi_{\mathsf{ref}}(\cdot \mid x) \, \mathrm{e}^{(r(x,\cdot) + \frac{\lambda}{\lambda}g(x,\cdot))/\beta}$$

Existence of an optimal dual variable λ^*

Lagrangian dual function

■ UPPER ENVELOPE FUNCTION

$$\begin{split} D(\lambda) &:= & \underset{\pi}{\text{maximize}} \ L(\pi, \lambda) \\ &= & \beta \, \mathbb{E}_x \left[\log \mathbb{E}_{y \, \sim \, \pi_{\mathsf{ref}}} \left[\, \mathrm{e}^{(r(x,y) + \lambda g(x,y))/\beta} \, \right] \right] \end{split}$$

cumulant-generating function

- * convex, and smooth function
- * strictly convex, and locally strongly convex function

$$\nabla^2 D(\lambda) \;\; \simeq \;\; \mathbb{E}_x \left[\, \mathsf{Var}_{y \, \sim \, \pi^\star(\cdot \, | \, x; \lambda)} \left[\, g(x,y) \, \right] \, \right]$$

Lagrangian dual problem

■ LAGRANGAIN DUAL MINIMIZATION

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

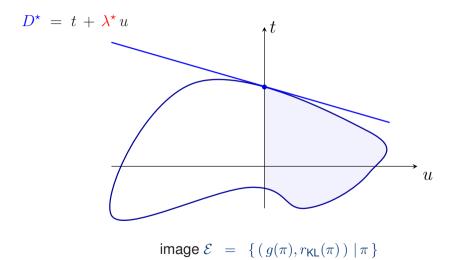
convex and smooth optimization

- * gradient descent finds an optimal dual variable λ^*
- * uniqueness of an optimal dual variable λ^*

Recovery of an optimal constrained policy π^*

$$\pi^{\star}(\cdot \mid x) \; = \; \pi^{\star}(\cdot \mid x; \textcolor{red}{\lambda^{\star}}) \; \propto \; \pi_{\mathsf{ref}}(\cdot \mid x) \, \mathrm{e}^{(r(x,\cdot) \, + \, \textcolor{red}{\lambda^{\star}} g(x,\cdot))/\beta}$$

■ GEOMETRIC INTERPRETATION OF STRONG DUALITY



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

Overview of our results

ZLHBD†R, NeurIPS '25 H†L†DBLHD†, NeurIPS '24

- DUALIZATION-BASED ALIGNMENT METHODS
 - * non-iterative & iterative methods
 - ⋆ duality gap
 - * optimality gap

objective and constraint

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

ALIGNMENT METHOD & THEORY

non-iterative method

Dualization-based alignment

H[†]L[†]DBLHD[†], NeurIPS '24

■ STAGE #1: FIND AN OPTIMAL DUAL VARIABLE

$$\lambda^* = \underset{\lambda \ge 0}{\operatorname{argmin}} D(\lambda)$$

convex and smooth optimization

■ STAGE #2: SEARCH FOR AN LLM POLICY

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

unconstrained alignment

Computational efficiency

Search for an optimal dual variable

Lagrangian maximizer:
$$\pi^*(\lambda) = \underset{\pi}{\operatorname{argmax}} L(\pi, \lambda)$$

Gradient:
$$\nabla D(\lambda) = \nabla_{\lambda} L(\pi, \lambda) \mid_{\pi = \pi^{*}(\lambda)}$$

■ PROJECTED GRADIENT DESCENT

$$\lambda^+ \leftarrow [\lambda - \eta \nabla D(\lambda)]$$

* π -independent gradient $\nabla D(\lambda)$

$$\nabla D(\lambda) = \frac{\mathbb{E}_{y \sim \pi_{\mathsf{ref}}} \left[e^{(r(x,y) + \lambda g(x,y))/\beta} g(x,y) \right]}{\mathbb{E}_{y' \sim \pi_{\mathsf{ref}}} \left[e^{(r(x,y') + \lambda g(x,y'))/\beta} \right]}$$

Offline biased estimate

Search for an LLM policy

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

■ POLICY PARAMETRIZATION

$$\pi(y \mid x) \leftarrow \pi_{\theta}(y \mid x)$$

model parameter θ

■ PARAMETRIZED LAGRANGIAN MAXIMIZER

$$\pi_{\theta^{\star}(\lambda^{\star})} \in \underset{\theta}{\operatorname{argmax}} L(\pi_{\theta}, \lambda^{\star})$$

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

Constrained parameter optimization

maximize
$$\mathbb{E}_{x} \left[\mathbb{E}_{y \sim \pi_{\theta}} \left[r(x, y) \right] - \beta D_{\mathsf{KL}} (\pi_{\theta} (\cdot \mid x) \parallel \pi_{\mathsf{ref}} (\cdot \mid x)) \right]$$
subject to
$$\mathbb{E}_{x} \left[\mathbb{E}_{y \sim \pi_{\theta}} \left[g(x, y) \right] \right] \geq 0$$

KL-regularized objective

policy constraint

 \star decision of model parameter θ

CHALLENGE

Nonconvex constrained optimization \rightarrow Lack of strong duality

ALIGNMENT METHOD & THEORY

iterative method

Parametrized Lagrangian dual problem

LAGRANGIAN DUAL FUNCTION

$$D_{\mathsf{p}}(\lambda) := \underset{\theta}{\operatorname{maximize}} L(\pi_{\theta}, \lambda)$$

* convex, and nondifferentiable function

LAGRANGIAN DUAL MINIMIZATION

$$\underset{\lambda \geq 0}{\text{minimize}} \ D_{\mathsf{p}}(\lambda)$$

Existence of an optimal parametrized dual variable $\lambda_{\rm p}^{\star}$

Search for an optimal parametrized dual variable

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

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

■ PROJECTED SUBGRADIENT DESCENT

$$\lambda^+ \leftarrow [\lambda - \eta u(\lambda)]_+$$

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

Online unbiased estimate

Iterative dualization-based alignment

ZLHBD†R, NeurIPS '25

■ ITERATION #1: COMPUTE A LAGRANGIAN MAXIMIZER

$$\theta^{\star}(\lambda) \in \underset{\theta}{\operatorname{argmax}} L(\pi_{\theta}, \lambda)$$

■ ITERATION #2: PERFORM A SUBGRADIENT DESCENT STEP

$$\lambda^{+} \leftarrow \left[\lambda - \eta \mathbb{E}_{y \sim \pi_{\theta^{\star}(\lambda)}} [g(x, y)] \right]_{+}$$

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

ALIGNMENT METHOD & THEORY

duality gap & optimality gap

WEITH METHOD & THEOTH

Duality gap

Duality gap: $|P^{\star} - D_{\mathsf{p}}^{\star}|$

Theorem (informal)

★ Duality gap is dominated by

 ν

parametrization gap
$$\nu \coloneqq \max_{\pi} \min_{\theta} \operatorname{dist}_1(\pi, \pi_{\theta})$$

 \star ν -parametrization gap yields ν -duality gap

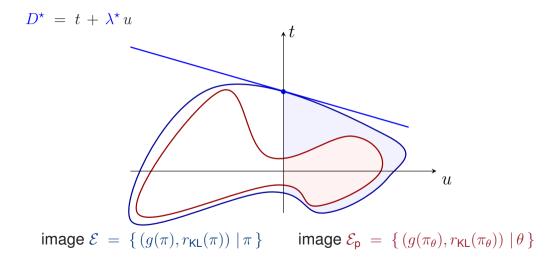
linear independence

■ GEOMETRIC INTERPRETATION OF DUALITY GAP

$$D^{\star} = t + \lambda^{\star} u$$

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

■ GEOMETRIC INTERPRETATION OF DUALITY GAP



■ GEOMETRIC INTERPRETATION OF DUALITY GAP

$$D_{\mathbf{p}}^{\star} = t + \lambda_{\mathbf{p}}^{\star} u$$

$$D_{\mathbf{p}}^{\star} = t + \lambda_{\mathbf{p}}^{\star} u$$

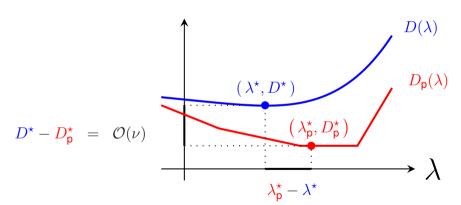
$$\lim \operatorname{age} \mathcal{E} = \{ (g(\pi), r_{\mathsf{KL}}(\pi)) \mid \pi \} \quad \operatorname{image} \mathcal{E}_{\mathbf{p}} = \{ (g(\pi_{\theta}), r_{\mathsf{KL}}(\pi_{\theta})) \mid \theta \}$$

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

$$D^{\star} - D_{\mathsf{p}}^{\star} = \mathcal{O}(\nu)$$

Gap between optimal dual variables

■ GAP BETWEEN (UN)PARAMETRIZED DUAL FUNCTIONS



Optimal dual variables: λ^* , λ_0^* are close:

$$\|\lambda^{\star} - \lambda_{\mathsf{n}}^{\star}\| = \mathcal{O}(\sqrt{\nu})$$

Optimality gap for iterative method

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

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

Implication (informal)

★ Objective optimality & Constraint feasibility are dominated by

$$\sqrt{\nu}$$

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

Root-scaling of parametrization gap

Optimality gap for non-iterative method

Objective optimality: $\left| r_{\mathsf{KL}} \left(\pi_{\mathsf{p}}^{\star}(\lambda^{\star}) \right) - r_{\mathsf{KL}}(\pi^{\star}) \right|$

Constraint feasibility: $\left| g \left(\pi_{\mathsf{p}}^{\star}(\lambda^{\star}) \right) - g(\pi^{\star}) \right|$

Implication (informal)

★ Objective optimality & Constraint feasibility are dominated by

$$\sqrt{\nu}$$

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

Root-scaling of parametrization gap

EMPIRICAL STUDY

safety-alignment task

Practical implementation

NON-ITERATIVE DUALIZATION ALIGNMENT

model-based setting

offline model-based dual pseudo-preference optimization

offline preference-based dual pseudo-preference optimization

■ ITERATIVE DUALIZATION ALIGNMENT

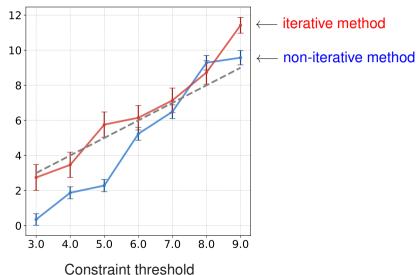
model-based setting

online model-based dual pseudo-preference optimization

preference-based setting online preference-based dual

pseudo-preference optimization

Constraint satisfaction

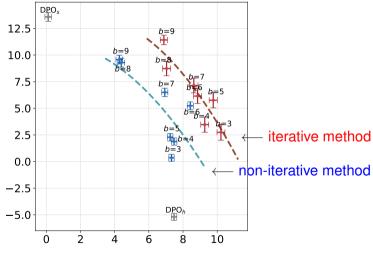


Harmlessness score

Better constraint satisfaction

Helpfulness and Harmlessness tradeoff

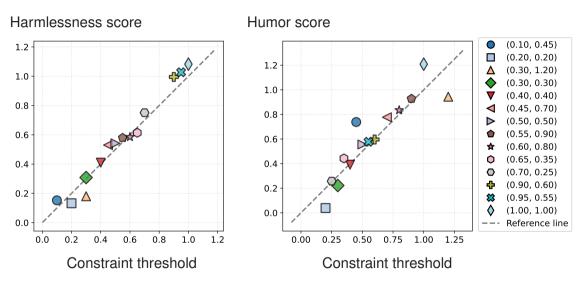
Harmlessness score



Helpfulness score

Higher Pareto frontier

Harmlessness and Humor constraints



Multi-constraint satisfaction

Summary & outlook

ZLHBD†R, NeurIPS '25 H†L†DBLHD†, NeurIPS '24

■ DUALIZATION-BASED ALIGNMENT METHODS

- * non-iterative & iterative methods
- * duality gap & optimality gap

■ OPEN CHALLENGES

- ⋆ optimal sample complexity
- ⋆ multi-turn alignment

Thank you for your attention.