# Reinforcement Learning with Human Feedback for Aligning LLMs

Suhas Palwai and Stuart Powers

November 8, 2024

# Outline

# What is Reinforcement Learning (RL)?

- RL is a framework where an **agent** learns to make decisions by interacting with an **environment**.
- The goal is to maximize cumulative **rewards**.
- Applications include robotics, game playing, and natural language processing.

# Key Components of RL

- **Agent**: Learns and makes decisions.
- **Environment**: The system the agent interacts with.
- **State Space** ($\mathcal{S}$): All possible states $s$.
- **Action Space** ($\mathcal{A}$): All possible actions $a$.
- **Transition Dynamics** ($P(s'|s, a)$)
  - Inputs: Current state $s$, action $a$, next state $s'$.
  - Output: Probability of transitioning to $s'$.
- **Reward Function** ($R(s, a)$)
  - Inputs: State $s$, action $a$.
  - Output: Reward $r$.

## Policy and Objective

- **Policy** $(\pi(a|s))$
    - Maps states to action probabilities.
    - Defines the agent's behavior.
- **Trajectory** $\tau = (s_0, a_0, s_1, a_1, \ldots, s_T)$
- **Probability of Trajectory**

$$P(\tau|\pi) = \rho_0(s_0) \prod_{t=0}^{T-1} \pi(a_t|s_t) P(s_{t+1}|s_t, a_t)$$

- **Return along Trajectory**

$$R(\tau) = \sum_{t=0}^{T-1} \gamma^t R(s_t, a_t)$$

- **Objective Function**

$$J(\pi) = \mathbb{E}_{\tau \sim P(\cdot|\pi)}[R(\tau)]$$

# Central Optimization Problem

$$\pi^* = \arg\max_{\pi} J(\pi)$$

- Find the optimal policy $\pi^*$ that maximizes expected return.

# State-Value Function

- **Definition**:

$$V^{\pi}(s) = \mathbb{E}_{\tau \sim P(\cdot|\pi)} \left[ R(\tau) \middle| s_0 = s \right]$$

- Expected return starting from state $s$ and following policy $\pi$.

# Action-Value Function

- **Definition**:

$$Q^\pi(s, a) = \mathbb{E}_{\tau \sim P(\cdot|\pi)} \left[ R(\tau) \middle| s_0 = s, a_0 = a \right]$$

- Expected return starting from state $s$, taking action $a$, then following $\pi$.

# Relationship between $V^\pi(s)$ and $Q^\pi(s, a)$

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) Q^\pi(s, a)$$

- The value of a state is the expected value of action-values under policy $\pi$.

# What is RLHF?

- Incorporates human feedback into the RL framework.
- Useful for tasks where rewards are hard to define.
- Aligns AI behavior with human values and preferences.

# Policy Gradient Methods

- Directly adjust policy parameters to maximize expected return.
- Policy gradient theorem:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(a_t|s_t) A^{\pi_\theta}(s_t, a_t) \right]$$

- High variance and instability can be issues.

# Proximal Policy Optimization (PPO)

- Addresses instability in policy gradients.
- Introduces a clipped surrogate objective.
- Limits the magnitude of policy updates.

# PPO Objective Function

$$L^{\text{PPO}}(\theta) = \mathbb{E}\left[\min\left(r_t(\theta)\hat{A}_t, \text{clip}\left(r_t(\theta), 1 - \epsilon, 1 + \epsilon\right)\hat{A}_t\right)\right]$$

- $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$
- $\hat{A}_t$ is the advantage estimate.
- $\epsilon$ controls the clipping range.

# Intuition Behind PPO

- Prevents large policy updates that could destabilize training.
- Maintains policy within a "trust region."
- Balances exploration and exploitation.

# What is DPO?

- Directly optimizes policy based on human preferences.
- Uses pairwise comparisons instead of scalar rewards.
- Effective for subjective tasks.

# DPO Objective Function

$$L^{\text{DPO}}(\theta) = \sum_i \log \sigma \left( f_\theta(y_i^+ | x_i) - f_\theta(y_i^- | x_i) \right)$$

- $f_\theta(y|x) = \log \pi_\theta(y|x)$
- $\sigma(z) = \frac{1}{1+e^{-z}}$ is the sigmoid function.
- Maximizes likelihood of preferred responses.

# Intuition Behind DPO

- Directly aligns model outputs with human preferences.
- Simplifies training without explicit reward functions.
- Captures nuanced human judgments.

# Experiment

- Use a model like GPT-2.
- Collect human-annotated preferences.
- Fine-tune the model using PPO and DPO.

# Key Takeaways

- RL is about maximizing expected cumulative rewards.
- RLHF incorporates human feedback for alignment.
- PPO provides stable policy updates.
- DPO directly optimizes for human preferences.
- Both methods enhance LLM alignment with human values.

# Questions

Thank you!
Questions?