

What's on your mind?

Choose an option:

- ☒ Lecture Review
- ☐ Ask Questions
- ☐ Test Yourself

TA-Lite: Academic Assistant

Logged in as: [r Rao@iu.edu](#)

Select Lecture

Which lecture would you like help with?

Lecture 08: Reinforcement Learning from Human Feedback

Select Material Type

Lecture Slides

Lecture

Lecture Review

Choose what you'd like to review:

Generate Summary

Generate

Summary from slides

In this article, the authors propose a new approach to training language models (LLMs) using reinforcement learning. The approach aims to improve the reasoning capability in LLMs by practicing solutions that take them from problem statements to final answers. The authors simulate human preferences and train an agent to follow these preferences. The agent receives rewards for good answers and penalties for bad answers, with the goal of maximizing cumulative rewards over time.

The authors describe the problem of training LLMs, which is to take a problem statement (prompt) and generate an answer. However, current methods have limitations, such as relying on hand-crafted rules or relying solely on supervised learning. They propose using reinforcement learning (RL) to train an agent to follow human preferences in generating answers.

The authors introduce the concept of a "preference model," which is a neural network that simulates human preferences. They also define a loss function for the preference model, which is a negative log-likelihood function that measures the difference between the human preference and the LLM's output. The goal is to minimize this loss function to incentivize the LLM to follow human preferences.

The authors then describe their proposed approach, which consists of three steps:

1. Instance the environment: Create a simulated environment that includes the problem statement and the desired output.
2. Learn the agent policy: Use RL to train an agent to take actions in the environment and receive rewards based on the quality of the output.
3. Evaluate the agent: Evaluate the trained agent using various metrics, such as accuracy and fluency, to determine its performance.

The authors also provide examples of how their approach can be applied in practice, such as in natural language processing tasks. They conclude that their proposed approach has the potential to improve the reasoning capability of LLMs and follow user preferences.