RLHF and Instruction-tuning



Few-shot Learning Thus Far

Thus far, we have talked about using LMs "out-of-the-box" for few-shot
o surprising emergent property

Questions:

- Can we directly train models to do few-shot learning?
- Can we directly train models to follow arbitrary user instructions?
- · Can we directly train models to obey toxicity & safety constraints?

Lecture Overview

- Instruction Finetuning
- Reinforcement Learning from Human Feedback (RLHF)
- Open challenges with RLHF

Lecture Overview

Instruction Finetuning

Reinforcement Learning from Human Feedback (RLHF)

Open challenges with RLHF

Language modeling ≠ assisting users

PROMPT Explain the moon landing to a 6 year old in a few sentences.

Language modeling ≠ assisting users

- PROMPT Explain the moon landing to a 6 year old in a few sentences.
- COMPLETION GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences. Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

Language modeling ≠ assisting users

PROMPT Explain the moon landing to a 6 year old in a few sentences.

COMPLETION GPF-3 Explain the theory of gravity to a 6 year old. Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

Language models are not aligned with user intent [Ouvang et al., 2022].

Language modeling ≠ assisting users

- PROMPT Explain the moon landing to a 6 year old in a few sentences.
- COMPLETION GPT-3
 - Explain the theory of gravity to a 6 year old.
 - Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old. Explain evolution to a 6 year old.

Language models are not *aligned* with user intent [<u>Ouvang et al., 2022</u>]. Finetuning to the rescue!

Language modeling ≠ assisting users

PROMPT Explain the moon landing to a 6 year old in a few sentences.

COMPLETION Human

A giant rocket ship blasted off from Earth carrying astronauts to the moon. The astronauts landed their spaceship on the moon and walked around exploring the lunar surface. Then they returned safely back to Earth, bringing home moon rocks to show everyone.

Language models are not *aligned* with user intent [<u>Ouvang et al., 2022</u>]. Finetuning to the rescue!

Instruction Finetuning • Collect examples of (instruction, output) pairs across many tasks and finetune an LM Please answer the following specific Medicine fo

[FLAN-T5; Chung et al., 202









Qualitative Results

Model input (Disambiguation QA)	Before instruction finetuning	
Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous.	The reporter and the chef will discuss their favorite dishes. The reporter and the chef will discuss the reporter's favorite dishes. The reporter and the chef will discuss the chef's favorite dishes. The reporter and the chef will discuss the reporter's and the chef's favorite dishes. # (doesn't answer question)	
Sentence: The reporter and the chef will discuss their favorite dishes.		
Options: (A) They will discuss the reporter's favorite dishes (B) They will discuss the chef's favorite dishes (C) Ambiguous		
Highly recommend trying FLAN-T	75 out to get a sense of its capabilities:	
15 https://huggingtac	ce.co/google/flan-t5-xxl [Chung et al., 2022	



Lecture Plan: From Language Models to Assistants

1. Instruction finetuning + Simple and straightforward, generalize to unseen tasks - ?

- ?

2. Reinforcement Learning from Human Feedback (RLHF)

3. What's next?

Limitations of instruction finetuning?

- Problem 1: it's expensive to collect ground-truth data for tasks
- Provide me five active research areas in April 2023 for LLMs
- Problem 2: tasks like open-ended creative generation have no right answer.
 Write me a story about a dog and her pet grasshopper.
- Problem 3: Even with instruction tuning, you are not directly "maximizing human preferences"
- · Can we explicitly attempt to satisfy human preferences?

Lecture Overview

- Instruction Finetuning
- Reinforcement Learning from Human Feedback (RLHF)
- Open challenges with RLHF

Optimizing for human preferences

• Let's say we were training a language model on some task (e.g. summarization).

Optimizing for human preferences

- Let's say we were training a language model on some task (e.g. summarization). For each LM sample s, imagine we had a way to obtain a human reward of that summary: R(s) ∈ ℝ, higher is better.



•	Let's say we were training a language model on some task (e.g. summarization). For each LM sample s , imagine we had a way to obtain a <i>human reward</i> of that summary: $R(s) \in \mathbb{R}$, higher is better.			
	SAN FRANCISCO, California (CNN) A magnitude 4.2 earthquake shook the San Francisco overturn unstable objects.	An earthquake hit San Francisco. There was minor property damage, but no injuries. S_1 $R(s_1) = 8.0$	The Bay Area has good weather but i prone to earthquakes and wildfires. S_2 $P(c_1) = 1.2$	





How do we model human preferences?

• Awesome: now for any arbitrary, non-differentiable reward function R(s), we can train our language model to maximize expected reward.

How do we model human preferences?

- Awesome: now for any **arbitrary, non-differentiable reward function** R(s), we can train our language model to maximize expected reward.
- Not so fast! (Why not?)
 Problem 1: human-in-the-loop is expensive!

4.

How do we model human preferences?

- Awesome: now for any **arbitrary, non-differentiable reward function** R(s), we can train our language model to maximize expected reward.
- Not so fast! (Why not?)
- Problem 1: human-in-the-loop is expensive!
- Solution: instead of directly asking humans for preferences, model their preferences as a separate (NLP) problem! [Knox and Stone, 2009]

How do we model human preferences?

- Awesome: now for any **arbitrary**, **non-differentiable reward function** R(s), we can train our language model to maximize expected reward.
- Not so fast! (Why not?)

đ.

- Problem 1: human-in-the-loop is expensive!
 Solution: instead of directly asking humans for preferences, model their
- preferences as a separate (NLP) problem! [Knox and Stone, 2009] An earthquake hit The Bay Area has San Francisco. good weather but is Train an LM $RM_{\phi}(s)$ to There was minor prone to predict human property damage, earthquakes and preferences from an



How do we model human preferences?

- Problem 2: human judgments are noisy and miscalibrated!
- Solution: instead of asking for direct ratings, ask for pairwise comparisons, which can be more reliable [<u>Phelos et al., 2015; Clark et al., 2018</u>]

A 4.2 magnitude earthquake hit San Francisco, resulting in massive damage.

 s_3

 $R(s_3) = 4.1? 6.6? 3.2?$

ł.



Make sure your reward model works first!



RLHF: Putting it all together [Christiano et al., 2017; Stiennon et al., 2020]

- Finally, we have everything we need:
- A pretrained (possibly instruction-finetuned) LM p^{PT}(s)
- + A reward model $RM_{\phi}(s)$ that produces scalar rewards for LM outputs, trained on a dataset of human comparisons
- A method for optimizing LM parameters towards an arbitrary reward function.







2

InstructGPT

	PROMPT	Explain the moon landing to a 6 year old in a few sentences.
	COMPLETION	OPF.3 Explain the theory of gravity to a 6 year old. Explain the theory of relativity to a 6 year old in a few sentences. Explain the big bang theory to a 6 year old. Explain evolution to a 6 year old.
		InstructOPT People went to the moon, and they took pictures of what they saw, and sent them back to the earth so we could all see them.
37		

InstructGPT				
	PROMPT	Write a short poem about a wise frog.		
	COMPLETION	$_{\rm OPP-3}$ Write a short story in which a character has two different names.		
		Write a short story in which you try to get something back that you have lost.		
		Write a short story in which a character has a bad dream.		
		InstructGPT		
		The frog is a symbol of wisdom He knows all the secrets of the world		
		He is a master of disguise		
		And a great teacher of life He is a symbol of transformation		
		And the bringer of change		
		He is the frog who has seen it all And knows the meaning of it all		
18				

ChatGPT: Instruction Finetuning + RLHF for dialog agents



Feedback (RLHF), using the same methods as InstructGPT, but with slight differences in the data collection setup. We trained an initial model using supervised fine-tuning: human AI trainers provided conversations in which they played both sides-the user and an AI assistant. We gave the trainers access to model-written suggestions to help them compose their responses. We mixed this new dialogue dataset with the InstructGPT dataset, which we transformed into a

(Instruction finetuning!)

https:

ChatGPT: Instruction Finetuning + RLHF for dialog agents



Note: OpenAI (and similar companies) are keeping more details secret about ChatGPT training (including data, training parameters, model size)perhaps to keep a competitive edge...

To create a reward model for reinforcement learning, we needed to collect comparison data, which consisted of two or more model responses ranked by quality. To collect this data, we took conversations that AI trainers had with the chatbot. We randomly selected a model-written message, sampled several alternative completions, and had AI trainers rank them. Using these reward models, we can fine-tune the model using Proximal Policy Optimization. We performed several iterations of this process.

(RLHF!)

https://openai.com/blog/chatept

Lecture Overview

- Instruction Finetuning
- Reinforcement Learning from Human Feedback (RLHF)
- Open challenges with RLHF

Limitations of RL + Reward Modeling

 Human preferences are unreliable! "Reward hacking" is a common problem in RL



Limitations of RL + Reward Modeling

Human preferences are unreliable! "Reward hacking" is a common problem in RL

- Chatbots are rewarded to produce responses that seem authoritative and helpful, regardless of truth
- This can result in making up facts
 + hallucinations

