## Lecture 9: RLHF and Guest Lecture on DPO

Emma Brunskill

CS234 Reinforcement Learning.

Spring 2024

## Midterm

- In class on Wednesday
- You are allowed 1 side of 1 8.5" x 11" sheet of notes
- All material through today's lecture (Monday) is eligible for the exam
- See Ed post for additional details and past related midterm/quizzes
- Good luck!

Select all that are true

T

- The Bradley Terry model expresses the probability that someone will select option $b_i$ over $b_j$
- Using preference tuples and the Bradley Terry model, one can learn a model of the reward function

  T
- The resulting reward function can be shifted by any constant and will not change the resulting preferences

  T
- The resulting reward function can be multiplied by any constant and will not change the resulting preferences

  F
- In RLHF we update the reward model after each PPO roll out   F
- Not sure

Select all that are true

- The Bradley Terry model expresses the probability that someone will select option $b_i$ over $b_j$
- Using preference tuples and the Bradley Terry model, one can learn a model of the reward function
- The resulting reward function can be shifted by any constant and will not change the resulting preferences
- The resulting reward function can be multiplied by any constant and will not change the resulting preferences
- In RLHF we update the reward model after each PPO roll out
- Not sure

1,2,3 are true. 4 is false: for example, we cannot multiply the rewards by -1 and preserve the ordering.

## Class Structure

- Last time: Imitation Learning (Max Entropy IRL) and RLHF
- This time: RLHF and Direct Preference Optimization (best paper runner up at top ML conference) guest lecture
- Next time: Midterm

- RLHF for LLM
- Direct Preference Optimization

- Often easier for people to make than hand writing a reward function
- Often easier than providing scalar reward (how much do you like this ad?)

- Consider $k$-armed bandits[1]: $K$ actions $b_1, b_2, \ldots b_k$. No state/context.
- Assume a human makes noisy pairwise comparisons, where the probability she prefers $b_i \succ b_j$ is

$$P(b_i \succ b_j) = \frac{\exp\left(r(b_i)\right)}{\exp\left(r(b_i)\right) + \exp\left(r(b_j)\right)} = p_{ij} \qquad (1)$$

- Assume have $N$ tuples of form $(b_i, b_j, \mu)$ where $\mu(1) = 1$ if the human marked $b_i \succ b_j$, $\mu(1) = 0.5$ if the human marked $b_i = b_j$, else 0 if $b_j \succ b_i$
- Maximize likelihood with cross entropy

$$loss = - \sum_{(b_i, b_j, \mu) \in \mathcal{D}} \mu(1) \log P(b_i \succ b_j) + \mu(2) \log P(b_j \succ b_i) \qquad (2)$$

- Use learned reward model, and do PPO with this model
- See prior lecture for notes on doing this over trajectories

---

[1]We will see more on bandits later in the course

- How is this used in ChatGPT?
- Next set of slides are from part of Tatsu Hashimoto's Lecture 11 in CS224N

# High-level instantiation: 'RLHF' pipeline



**Step 1**
**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3 with supervised learning.

**Step 2**
**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.

**Step 3**
**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

- First step: instruction tuning!
- Second + third steps: maximize reward (but how??)

# 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 [Phelps et al., 2015; Clark et al., 2018]

An earthquake hit San Francisco. There was minor property damage, but no injuries.

$>$

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

$>$

The Bay Area has good weather but is prone to earthquakes and wildfires.

$s_1$     1.2        $s_3$          $s_2$



Reward Model ($RM_\phi$)

The   Bay   Area   ...   ...   wildfires

Bradley-Terry [1952] paired comparison model

$$J_{RM}(\phi) = -\mathbb{E}_{(s^w, s^l) \sim D}\left[\log \sigma(RM_\phi(s^w) - RM_\phi(s^l))\right]$$
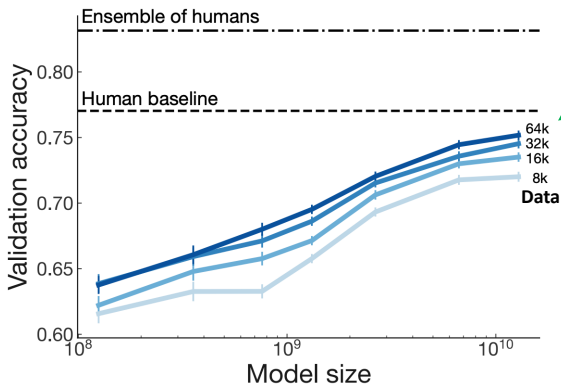
"winning" sample    "losing" sample    $s^w$ should score higher than $s^l$

# Make sure your reward model works first!

Evaluate RM on predicting outcome of held-out human judgments



**Large enough RM trained on enough data approaching single human perf**

[Stiennon et al., 2020]

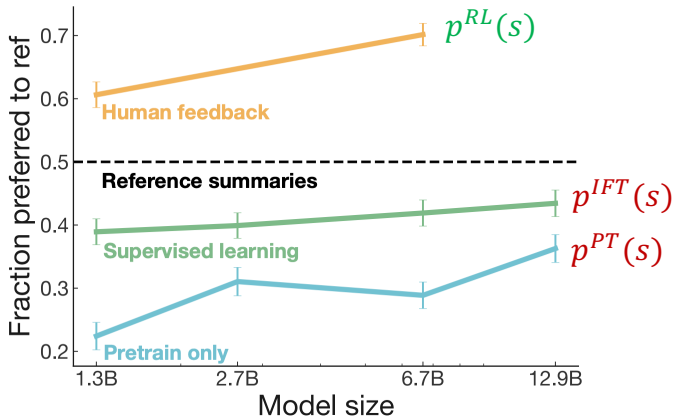# RLHF: Putting it all together [Christiano et al., 2017; Stiennon et al., 2020]

- 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.
- Now to do RLHF:
  - Initialize a copy of the model $p_\theta^{RL}(s)$, with parameters $\theta$ we would like to optimize
  - Optimize the following reward with RL:

$$R(s) = RM_\phi(s) - \beta \log \left( \frac{p_\theta^{RL}(s)}{p^{PT}(s)} \right)$$

Pay a price when
$p_\theta^{RL}(s) > p^{PT}(s)$

This is a penalty which prevents us from diverging too far from the pretrained model. In expectation, it is known as the **Kullback-Leibler** (**KL**) divergence between $p_\theta^{RL}(s)$ and $p^{PT}(s)$.

# RLHF provides gains over pretraining + finetuning



[Stiennon et al., 2020]

43

# InstructGPT: scaling up RLHF to tens of thousands of tasks



**Step 1**
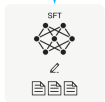**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

This data is used to fine-tune GPT-3 with supervised learning.

**Step 2**
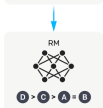**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

A labeler ranks the outputs from best to worst.

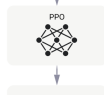This data is used to train our reward model.

**Step 3**
**Optimize a policy against the reward model using reinforcement learning.**

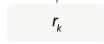A new prompt is sampled from the dataset.

The policy generates an output.

The reward model calculates a reward for the output.

The reward is used to update the policy using PPO.

**30k tasks!**

[Ouyang et al., 2022]

44

# Controlled comparisons of "RLHF" style algorithms

| Method | Simulated win-rate (%) | Human win-rate (%) |
|---|---|---|
| GPT-4 | $79.0 \pm 1.4$ | $69.8 \pm 1.6$ |
| ChatGPT | $61.4 \pm 1.7$ | $52.9 \pm 1.7$ |
| PPO | $46.8 \pm 1.8$ | $55.1 \pm 1.7$ |
| Best-of-$n$ | $45.0 \pm 1.7$ | $50.7 \pm 1.8$ |
| Expert Iteration | $41.9 \pm 1.7$ | $45.7 \pm 1.7$ |
| SFT 52k (Alpaca 7B) | $39.2 \pm 1.7$ | $40.7 \pm 1.7$ |
| SFT 10k | $36.7 \pm 1.7$ | $44.3 \pm 1.7$ |
| Binary FeedME | $36.6 \pm 1.7$ | $37.9 \pm 1.7$ |
| Quark | $35.6 \pm 1.7$ | - |
| Binary Reward Conditioning | $32.4 \pm 1.6$ | - |
| Davinci001 | $24.4 \pm 1.5$ | $32.5 \pm 1.6$ |
| LLaMA 7B | $11.3 \pm 1.1$ | $6.5 \pm 0.9$ |

- Many works study RLHF behaviors using GPT-4 feedback (**Simulated**) as a surrogate for **Human** feedback.
- PPO (method in InstructGPT) does work
- Simple baselines (Best-of-n, Training on 'good' outputs) works well too

[Dubois et al 2023]

- RLHF for LLM
- **Direct Preference Optimization**

## Learning More

- Learning and making decisions from human preferences is a rich area intersecting social choice, computational economics and AI
- New course at Stanford on this topic: Koyejo's CS329H: Machine Learning from Human Preferences

# Class Structure

- Last time: Imitation Learning (Max Entropy IRL) and RLHF
- This time: RLHF and Direct Preference Optimization (best paper runner up at top ML conference) guest lecture
- Next time: Midterm