# An Introduction to RL from Human Feedback

① Language as an MDP

② Two Key Challenges of Fine-Tuning

③ Learning from Preferences

④ KL-Regularized RL

⑤ On the Information Geometry of RLHF

Gobul Swamy
gswamy@cmu.edu
www.gobul.dev

---

① $M = \{$

$S$: set of partial generations, e.g. all strings w/ len $\leq H \rightarrow |S| = |A|^H$

$A$: set of tokens, e.g. $\{'a', 'b', \dots\} \rightarrow$ "byte-pair", watch Karpathy

$T$: $P(s'|s,a) = \begin{cases} 1, & s' = s \circ a \\ 0, & o.w. \end{cases} \rightarrow$ "tree-structured", deterministic, known, resets easy

$r$: ? ? ?

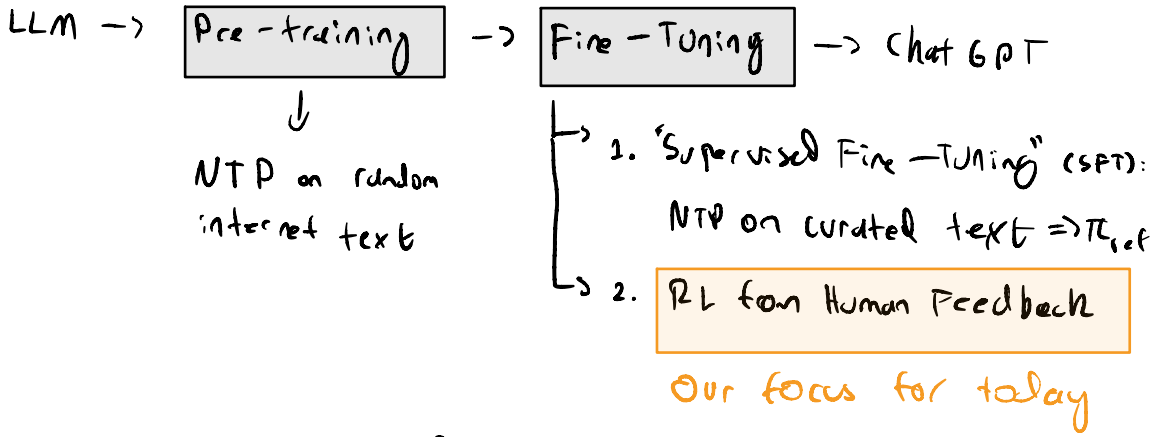$P_0$: prompts, e.g. "Summarize __".

$H$: maximum generation length.

$\}$

$\xi = [s_0, a_0, s_1, a_1, \dots s_H]$

$= [$"We", "The", "We The", "People", "We The People"$]$

"Next-Token Prediction"

②

LLM → **Pre-training** → **Fine-Tuning** → Chat GPT

⟱

NTP on random
internet text

1. "Supervised Fine-Tuning" (SFT):
   NTP on curated text $\Rightarrow \pi_{ref}$

2. RL from Human Feedback

Our focus for today

**Challenge 1:** What is $r$?

- the "reward design" problem
⇒ Leads to ③ : learning from preferences

**Challenge 2:** How to stay "close" to $\pi_{ref}$?
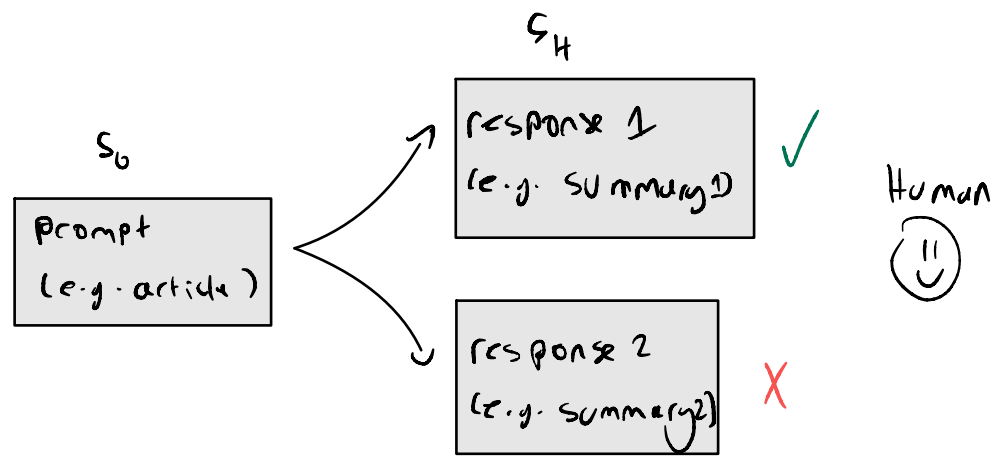
- the "fine-tuning" problem
⇒ Leads to ④ : KL-regularized RL

③ <u>Reward Design</u> is hard for problems where the behavior
we want is for an agent to be "human-like"
(e.g. self-driving, chat-bots).

=> key idea: <u>learn</u> the reward function from data via MLE

=> for RLHF, we will learn from preference feedback

$S_H$



$S_0$

Prompt
(e.g. article)

response 1
(e.g. summary1) ✓

response 2
(e.g. summary2) ✗

Human

$S_H'$     => $D = \{(S_0, S_H^+, S_H^-)\}$

=> "Bradley-Terry" Model: $P_{r^*}(S_H^+ > S_H^-) = \dfrac{1}{1 + \exp(r^*(S_H^-) - r^*(S_H^+))}$

★ "polite fiction", ask me later

$= \sigma(r^*(S_H^+) - r^*(S_H^-))$

$\hat{r}_{BT} = \underset{r \in R}{\arg\min} \ D_{KL}(P_{r^*} \| P_r)$  (forward KL proj)

$= \underset{r \in R}{\arg\min} \ \mathbb{E}_{S_H^+, S_H^- \sim P_{r^*}}[\underbrace{\log P_{r^*}(S_H^+ > S_H^-)}_{constant} - \log P_r(S_H^+ > S_H^-)]$

$\approx \underset{r \in R}{\arg\min} \ \mathbb{E}_D[-\log P_r(S_H^+ > S_H^-)]$  (FKL = MLE)

$= \underset{r \in R}{\arg\min} \ \mathbb{E}_D[-\log \sigma(r(S_H^+) - r(S_H^-))]$

④ Let us use $\pi_{ref}$ to refer to the output of SFT.
We want to stay "close" in policy space to $\pi_{ref}$
during RLHF to not "exploit" the reward model.
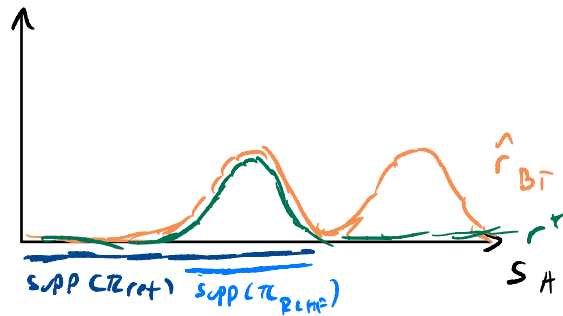This leads to the KL-Regularized RL Problem:

"reverse KL"

$$\pi_{RLHF} = \underset{\pi \in \Pi}{\arg\max} \; \underset{s \sim \pi}{\mathbb{E}} \left[ \hat{r}_{BT}(s_H) \right] - D_{KL}(\pi \| \pi_{ref})$$

$$= \underset{\pi \in \Pi}{\arg\max} \; \underset{s \sim \pi}{\mathbb{E}} \left[ \hat{r}_{BT}(s_H) - \sum_{h}^{H} \log \frac{\pi(a_h | s_h)}{\pi_{ref}(a_h | s_h)} \right]$$

↳ never pick a
token ref. policy
wouldn't

Solve using,
e.g., PPO.



supp($\pi_{ref}$)  supp($\pi_{RLHF}$)  $s_H$

$\hat{r}_{BT}$

=> Intuitively, RLHF is doing "mode-selection" on top of SFT

Surprisingly enough, there is a closed form answer to the
above KL-regularized RL problem at the trajectory level.

Let's work w/ trajectory-level distributions w/o prompts:

$$= \arg\max_{P} \, \mathbb{E}_{\xi \sim P}\left[\hat{r}_{\theta T}(\xi) - \log P(\xi) + \log P_{ref}(\xi)\right]$$

s.t. $\sum_{\xi} P(\xi) = 1$ $\Rightarrow$ P is a valid prob. distribution

$$= \arg\max_{P} \min_{\lambda} \sum_{\xi} P(\xi)\left[\hat{r}_{\theta T}(\xi) - \log P(\xi) + \log P_{ref}(\xi)\right] + \lambda\left(\sum_{\xi} P(\xi) - 1\right)$$

$$= \arg\max_{P} \min_{\lambda} \sum_{\xi} P(\xi)\left[\hat{r}_{\theta T}(\xi) - \log P(\xi) + \log P_{ref}(\xi) + \lambda\right] - \lambda$$

$$= \arg\max_{P} \min_{\lambda} L(P, \lambda)$$

Let's apply the stationarity condition from KKT: $\forall \xi$,

$$\nabla_{P(\xi)} L(P^*, \lambda^*) = 0$$

$$\Rightarrow \left[\hat{r}_{\theta T}(\xi) - \log \overset{*}{P}(\xi) + \log P_{ref}(\xi) + \lambda^*\right] + P^*(\xi)\left[\frac{-1}{P^*(\xi)}\right] = 0$$

$$\Rightarrow \log P^*(\xi) = \hat{r}_{\theta T}(\xi) + \log P_{ref}(\xi) + \lambda^* - 1$$

$$\Rightarrow P^*(\xi) = \frac{P_{ref}(\xi) \cdot \exp\left(\hat{r}_{\theta T}(\xi)\right)}{\exp(1 - \lambda)}$$

$$\nabla_{\lambda^*} L(P^*, \lambda^*) = \sum_{\xi} P^*(\xi) - 1 = 0$$

$$\Rightarrow \frac{1}{\cancel{}} \sum_{\xi} P_{ref}(\xi) \cdot \exp\left(\hat{r}_{\theta T}(\xi)\right) = \exp(1 - \lambda^*)$$

$$\Rightarrow \lambda^* = 1 - \log\left(\sum_{\xi} P_{ref}(\xi) \exp\left(\hat{r}_{\theta T}(\xi)\right)\right)$$

"partition function"

"exponential family" / "maximum entropy"

$$\boxed{\lambda^* = 1 - \log(2), \quad P^*(\xi) = \frac{P_{ref}(\xi) \cdot \exp\left(\hat{r}_{\theta T}(\xi)\right)}{Z}}$$
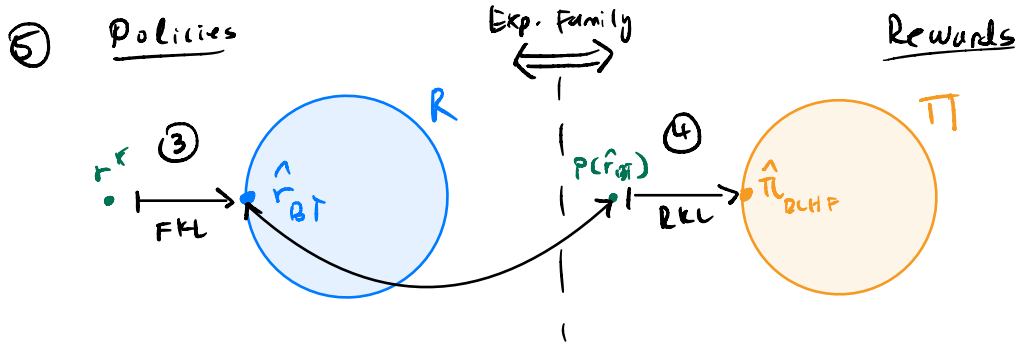
Now, consider the RKL projection of $p^*$ onto some class $\mathbb{P}$:

$$\arg\min_{p \in \mathbb{P}} D_{KL}(p \| p^*)$$

$$= \arg\min_{p \in \mathbb{P}} \sum_{\zeta} p(\zeta)\left[\log p(\zeta) - \log p^*(\zeta)\right]$$

$$= \arg\min_{p \in \mathbb{P}} \sum_{\zeta} p(\zeta)\left[\log p(\zeta) - \log p_{ref}(\zeta) - \hat{r}_{BT}(s_H) + \underset{\text{constant}}{\underbrace{\log Z}}\right]$$

$$= \arg\min_{p \in \mathbb{P}} \sum_{\zeta} p(\zeta)\left[\log p(\zeta) - \log p_{ref}(\zeta) - \hat{r}_{BT}(s_H)\right]$$

$$= \arg\max_{p \in \mathbb{P}} \sum_{\zeta} p(\zeta)\left[\hat{r}_{BT}(s_H) - \log \frac{p(\zeta)}{p_{ref}(\zeta)}\right]$$

$$= \arg\max_{p \in \mathbb{P}} \underset{\zeta \sim p}{\mathbb{E}}\left[\hat{r}_{BT}(s_H)\right] + D_{KL}(p \| p_{ref})$$

Now, let's set $\mathbb{P} = \left\{p(\zeta) = \prod_h \pi(a_h | s_h) \mid \pi \in \Pi\right\}$ and swap:

$$= \arg\max_{\pi \in \Pi} \underset{\zeta \sim \pi}{\mathbb{E}}\left[\hat{r}_{BT}(s_H)\right] + D_{KL}(\pi \| \pi_{ref})$$

=> Thus, Max Ent / soft RL is an RKL projection

⑤ <u>Policies</u>                    Exp. Family          <u>Rewards</u>



<u>Aside</u>: We can easily incorporate contexts/prompts/int. states:

$$P^*(s_H | s_0) = \frac{\prod_h^H \pi_{ref}(a_h | s_h) \cdot \exp(\hat{r}_{BT}(s_H))}{Z(s_0)}$$

If you'd like to learn:

- how you actually solve the RL problem in practice
- why we can't just optimize the policy on pref. data
- how to do the above w/o Bradley Terry assumption

Take  17-740  w/ Drew, Steven, and I next semester!

Website: www.interactive-learning-algos.github.io