Maximize expected reword

$$J(0,x) = \begin{cases} P_{\theta}(x|x) \cdot (x,y) \\ \end{cases}$$

$$\frac{dJ}{d\theta} = \begin{cases} \langle (x,y) \frac{d}{d\theta} P_{\theta} (y|x) \\ \langle (x,y) \frac{d}{d\theta$$

LI REINFORCE (Williams, 1992)

1) "log trick", multiply / divide by
Po (41x)

$$= \begin{cases} \langle \langle (x,y) \rangle \rangle \langle (y|x) \rangle & \frac{d}{d\theta} \langle (y|x) \rangle \\ & \frac{d}{d\theta} \langle (y$$

by de log Po (Y/X)

Li remember
$$\frac{d}{d\theta} \log \upsilon = \frac{1}{\upsilon} \cdot \frac{d\upsilon}{d\theta}$$

$$= dJ = \begin{cases} (x,y) \log(y|x) & d \log(\log(y|x)) \\ d\theta & y \end{cases}$$

Ly this itself is an expectation!

- 1. Given a prompt x, we sample a response y from our current model Po
- 2. Compute reward ((x14) from our frozen reward model
- 3, Adjust & in direction do log & (YIX)
 Scaled by our reward

Ly usually, we sample 1 to 16/32 y's per prompt x to compute the update

b) don't generally use the raw reward r(x,y), we usually subtract a baseline Grst, which reduces vorience

r(x,y)-b } advantage

baseline

braw reword

b) if positive, increase the prob

b) if negative, decrease prob. of y

advantage A(x,y) = r(x,y) - b(x)

Ly if A(x,y) > 0, better than usual Ly if A(x,y) < 0, worse than usual

batch:

X, Y, Y ((x, Y,))

(compute b

(so mean

as mean

butch

foutth

reword

Ly this strategy doesn't consider prompt &

5) to get a better estimate of b(x)L) Sample multiple y; for that x

(5) use b(x) = mean reward over Y;L) RLOO, GRPO

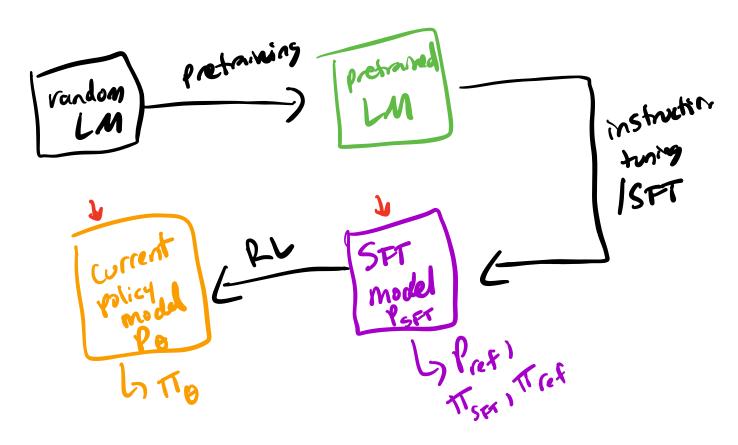
Reword hacking:

1) model finds unintended ways to maximize the reward

La solutions

1. train a better reward model
1. pat more | higher quality
pref. data
1. Start ul a bigger LM

2. penalize large deviations from SFT modul



Louse KL divergence to measure difference between Por and PSPT

DKL (P||Q) = E[log P(4|X)-log Q(4|X)] Y~P(.14)

by if P and Q are equal, Dkc (P||Q)=0 Lithe more P places probability on responses 4 where Q doesn't place prob, the higher the KL penalty

A (X,Y)- & D (PB || PSFF)

Ly hyperparameter that controls strength of penalty Ly RL's updates are scaled by the prob. of the model Po producing Y, it penalizes low-reward responses that the model actually produces

DIST has no such scaling, pushing model towards y it wouldn't currently produce