$$J(\theta) = E[r(x,y)]$$

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

$$= \begin{cases} r(x,y) P_{\theta}(y|x) \\ d \end{cases}$$

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

L) Q: can we just sample a few y and compute the average of  $r(x,y) \frac{d}{d\theta} P_{\theta}(y|x)$ ?

L) A: no, be cause  $\begin{bmatrix}
r(x,y) & d & P_{\theta}(y|x) \\
d\theta & & \text{from sampling} \\
replied$   $= \begin{cases}
\sqrt{P_{\theta}(y|x)} & r(x,y) & d & P_{\theta}(y|x) \\
\sqrt{P_{\theta}(y|x)} & r(x,y) & d\theta
\end{cases}$ 

dividing by this now yields log trick!

So from last class.

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

baseline: 1. man butch Stats

2. group skits

3. trained value for or critic by RLHF

Value for :

b) takes prompt x and partial response  $Y_{ci}$ , predicts Y(x,y)

b) just like reward midd, scalar head trained on top of LM, 12 loss

Louslike RM, Value for trained alongside policy model!

be're not done yet! Lo want to make small, stable thanges when updating Po La policy ratio p= Po(4/x) Pous (4/x) P:=1

D=000 for first step on butch, different after subsequent steps,

Is how much morelless likely is this token under the new policy vs. the policy that generaled it?

Lour con meight our advantage by Pi, and control update by clipping P;

Let's denote AKIN=A(X-Y)-BDER(PO 11PSFT)  $J(e) = E \left( \sum_{i=1}^{11} \min \left( S_i A_i^{kl}(x_i y) \right) \right)$ clip(Pi, 1-E, 1+E)Ak(KIY)) Telipping
hyperperam
= 0.1, 0.2

intuition. L) if  $A_t^{KL} > 0$ , we like this token and want to increase its prob.

L) max increase of  $(1+E)A_t^{KL}$ b) if AKCO, we don't like the token L> max de crease (1-ε)A<sub>ξ</sub><sup>kl</sup> Lono gradient than clipped term! example: E=0.2, so we clip (P;, 0.8, 1.2) A=2.0, P=1.3 Loundaged = 2.6 dipped = 2.4 => choose min, no gradients

A:= 2.0, P:= 1.1 token already more likely under New & likely under New &

Gnormal gradient, in weuse prob of token  $A_i = -2.0$ ,  $\rho_i = 0.6$ Ly unclipped = -1.2 Clipped = -1.6 > Choose min, no gradient, token already less likely under new &  $A_i = -2.0$ ,  $Q_i = 0.4$ normal grad, 4 unclipped = -1,8 = de cruse pob of them clipped = -1.6 7 tost A:70 A; CD PPO algo: Schulmone#2017

GRPO: group-relative PPO La baselne comes from group Statistics Ly for each prompt X, sample a group of responses Y, 1/2, ..., YG 1) compte mean reward 1/x and Add dev of reward ox over group b) for any specific y; in groups A(x,4;)= ((x,4;)-Nx

12-510re

Deepseek R1: 5 GRPO by remove reward models instead use simple verifiable tasts Leg. Math problems b prompt x: 37+53=3 by response 7: 32 => reward 0 response 42: 47 Drewed D response 13: 90 > reward 1 1) incentivize thinking / reasoning

before answering