

→ RL-HF and PPO (DPO advised)

- used for alignment
AI alignment.

To force LLM behave in a specific way e.g. assistant to behave in a friendly way.

→ Don't use bad or racist words and in particular style.

RL: RL is concerned with how an intelligent agent should take action in an environment to maximize cumulative reward. (select a policy that maximizes the return when agent acts accordingly to it)

Policy: rules how the agent select the action to perform given the state it is in (probability to move next given current) The goal of RL is to select a policy that maximizes the expected return when agent acts according to it (optimizer)

Connection to LLM.

Agent, LLM State, prompt(Action, next token), my LLM Reward, rewarding for good and not rewarded for bad response How to reward / penalize?

We'll build a reward model

→ Instead of one answer, we generate multiple answers for a prompt (by using high temperature) and then we ask people to choose the answer they prefer. Using the dataset of preference, we create a model that generates numeric reward for each Q and A.

→ we create a NN that gives the highest reward ^{layer}
chosen answer and lower to others ^{in output}

REWARD MODEL ARCHITECTURE.

- we take a pretrained LLM, an input (π/A) after concatenating them, the from last hidden state passes to linear layer of (1-feature) to generate a single value considered as reward for model.
- we also have to tell the model to give high reward to chosen π/A .

Loss - after

$$\text{Loss} = -\log(\sigma(r(\pi, y_w) - r(\pi, y_i)))$$

good ↓ bad ↑
consider we have two answers y_w & y_i .
two cases,

$$r(\pi, y_w) > r(\pi, y_i)$$

→ sigmoid will return value > 0.5 , loss will be small

$$r(\pi, y_w) < r(\pi, y_i)$$

→ sigmoid will return value < 0.5 , loss will be high

Trajectories

+ policy

$$\hat{\pi} = \underset{\pi}{\operatorname{argmax}} J(\pi)$$

The expected return over policy is expected return over all trajectories

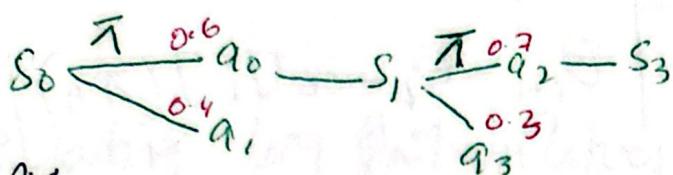
$$J(\pi) = \int P(T|\pi) R(T) = \underset{T \sim \pi}{\mathbb{E}} [R(T)]$$

Trajectory (path) is a series of (action, states), starting from initial state ($\tau = (s_0, a_0, s_1, a_1)$)

We will model the next state as stochastic (can't always be correct always)

$$s_{t+1} \text{ or } P(s_{t+1} | s_t, a_t)$$

→ Policy tells the probability of next state to move in all directions with given probability



Probability of Trajectory

$$P(\tau | \pi) = p_0(s_0) \prod_{t=0}^{T-1} P(s_{t+1} | s_t, a_t) \pi(a_t | s_t)$$

Start of next state
 probability of
 giving current state choosing action given state

| action defined on
the current state
and independent.
we multiply all
to get trajectory

we will always work with discounted Reward (immediate instead of future without discount if it will take long to reach optimum, because $\gamma^8 \rightarrow$ will decay so, intermediate)

$$R(\tau) = \sum_{t=0}^{\infty} \gamma^t r_t$$

$$\leftarrow (0 < \gamma < 1)$$

Trajectories in LLM.

states = prompt, policy to optimize=?
action + next token,

→ Trajectories are series of prompt or token we have chosen.

Policy - Gradient Optimization

we want to maximize this function actually

$J(\pi_\theta) = E_{T \sim \pi} [R(T)] \rightarrow$ to, for max we use gradient ascent

objective function

Expression of gradient.

$$\nabla_\theta J(\theta) = \nabla_\theta E_{T \sim \pi} [R(T)]$$

$$\Rightarrow \nabla_\theta \int_T P(T|\theta) R(T)$$

$$= \int_T \nabla_\theta P(T|\theta) R(T)$$

$$= \int_T P(T|\theta) \nabla_\theta \log P(T|\theta) R(T)$$

$$= E_{T \sim \pi} [\nabla_\theta \log P(T|\theta) R(T)]$$

$$\Rightarrow E_{T \sim \pi} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) R(T) \right]$$

↓
→ This expectation is over all trajectories, in which can be highly costly. The one way to elide the issue is to approximate it using sample mean by collecting a set of D trajectories.

$$\hat{S} = \sum_{t=0}^T \sum_{\pi \in D} \nabla_\theta \log \pi_\theta(a_t | s_t) R(T)$$

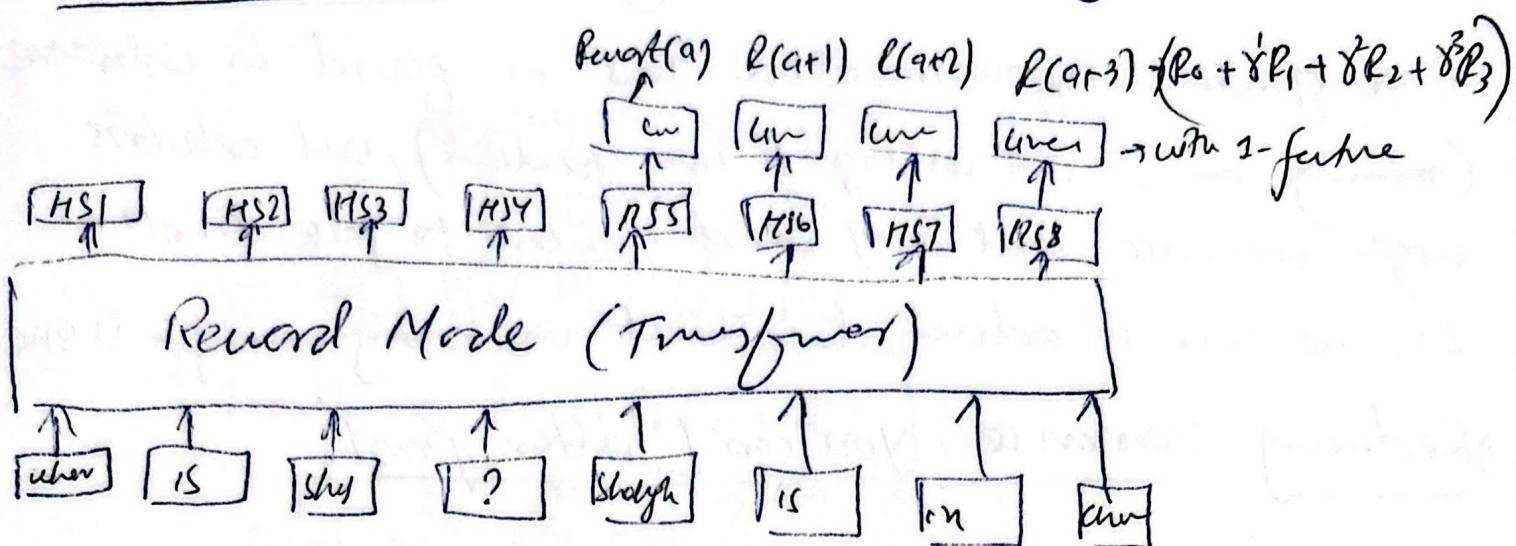
Policy Gradient Algorithm

$$\theta_{t+1} = \theta_t + \alpha \nabla_\theta J(\pi_\theta)|_{\theta_t}$$

$$\therefore \nabla_\theta \log P(T|\theta) = \frac{1}{P(T|\theta)} \cdot \nabla_\theta P(T|\theta)$$

$$\therefore P(T|\theta) = p_0 s_0 \prod_{t=0}^{T-1} P(s_{t+1}|s_t, a_t) \pi_\theta(a_t | s_t)$$
$$\nabla_\theta \log P(T|\theta) = \nabla_\theta \log p_0(s_0) + \sum \nabla_\theta \log P(s_{t+1}) + \nabla_\theta \log \pi_\theta(a_t | s_t)$$

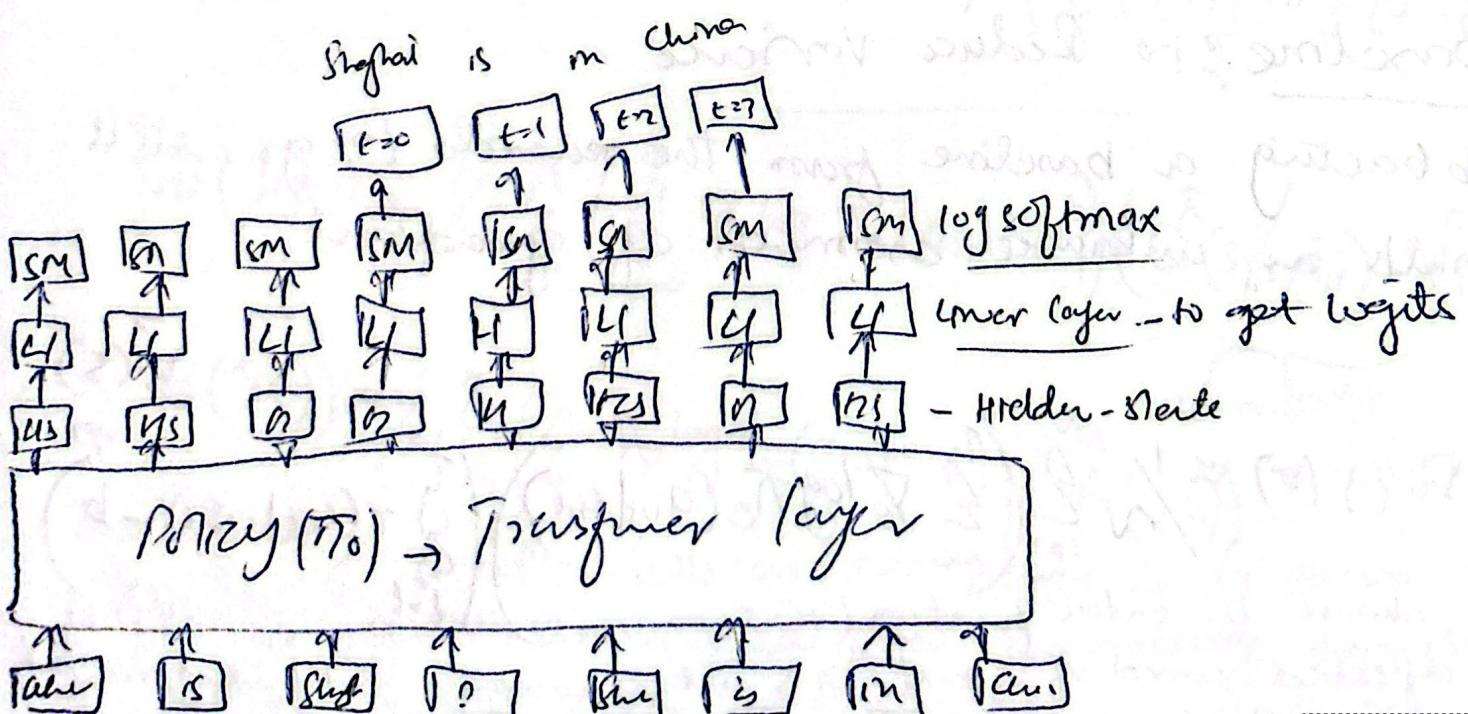
Calculating reward for each trajectory.



Calculating log Probabilities of our Policy

after calculating policy and reward . we run gradient ascent . (to optimize LEM)

$$\hat{g} = \frac{1}{|D|} \sum_{d \in D} \sum_{i \in w} \nabla \log \pi_\theta(a_t(s_t) | R(t))$$



Problems with Gradient Policy Optimization

1) the gradient approximation that we found is unbiased (meaning on it will converge to true gradient), but exhibits high variance that may lead gradient to false direction - so, we have to reduce it without increasing sample D size.

2) Reducing Variance: you can't alter past

→ each action can't alter the reward, we look in front

$$\nabla_{\theta} J(\pi_\theta) = E_{\pi=\pi_\theta} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_\theta(a_t | s_t) R(s) \right]$$

$$\nabla_{\theta} J(\pi_\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^T \nabla_{\theta} \log \pi_\theta(a_{i,t} | s_{i,t}) \right) \underbrace{\left(\sum_{t=0}^T r(s_{i,t}, a_{i,t}) \right)}_{\text{rewards upto}}$$

In each action, we are multiplying rewards that come before the action was taken. It is proved that past terms cancel out in expectation, we remove the

2) Baseline to Reduce Variance.

→ subtracting a baseline from the rewards to go, still results in unbiased estimator of gradient.

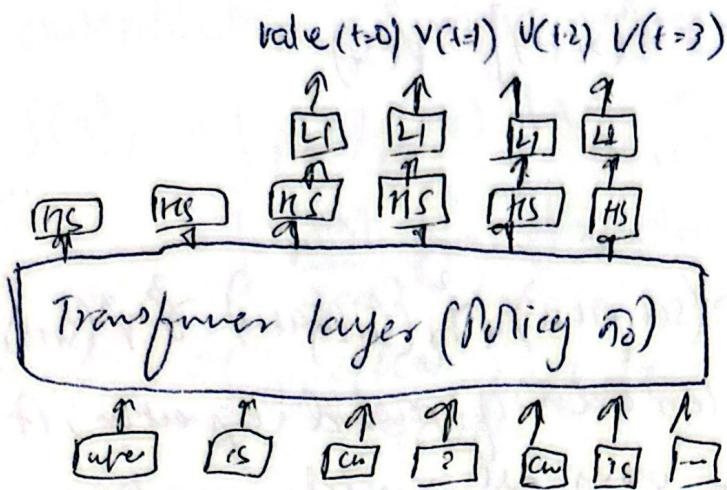
$$\hat{V}(s)$$

$$\nabla_{\theta} J(\pi_\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^T \nabla_{\theta} \log \pi_\theta(a_{i,t} | s_{i,t}) \right) \left(\sum_{t'=t}^T r(s_{i,t'}, a_{i,t'}) - b \right)$$

→ we choose $b = \text{value function}$ (that indicates the expected reward of agent in s act (a) according to policy (π)) $\Rightarrow \hat{V}^\pi(s)$

can also be dependent on state

How do we estimate $\rightarrow V^\pi(s)$:



3) Reducing Variance - introducing φ and V^* (advantage function)

φ = Rewards to go sum

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t}) \right) \underbrace{\left(\sum_{t=t}^T \varphi(s_{i,t} | a_{i,t}) \right)}_{\varphi}$$

$$so, \quad \varphi(s, a) = E_s [r(s, a)] + \delta E_a [\varphi^*(s', a')]$$

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t}) \right) (\hat{\theta}(s_{i,t}, a_{i,t}) - V^*(s))$$

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t}) \right) (\hat{A}(s, a)) \text{ advantage function}$$

- The advantage function tells us how better is to choose a particular action in a state over the average expectation we get by choosing randomly an action in same state.
- how good is particular than the average action in state

Estimating the advantage term.

$$\hat{A}^\pi(s_t, a_t) = \underbrace{[r(s_t, a_t) + \gamma V^\pi(s_{t+1})]}_{\Phi^\pi(s_t, a_t)} - V^\pi(s_t)$$

or

on three-step

$$\hat{A}^\pi(s_t, a_t) = [r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \gamma^2 r(s_{t+2}, a_{t+2}) + \gamma^3 V^\pi(s_{t+3})] - V^\pi(s_t)$$

→ If we stop early, use 1st order of reward equation, it will have high ~~bias-variance~~, while using all reward equations and doing no approximation with value function we get high variance. In order to solve bias-variance, we take weighted sum of terms + get Generalized Advantage Estimator

$$g_t = r_t + \gamma V^\pi(s_{t+1}) - V^\pi(s_t)$$

$$\hat{A}_t = g_t + \gamma \lambda \hat{A}_{t+1}$$

\curvearrowleft decaying parameter

Advantage for LMQ₃

so, we want to maximize the reward for action that are good for a state (s)

This is recursive formulae
that for 1st term equals
the first expansion. For
the second - 1st term
equals to second expansion
weighted by $(1-\lambda)$ etc

Problem #2 : Gradient-Policy Optimization.

$$\theta_{k+1} = \theta_k + \alpha \nabla_\theta J(\pi)|_{\theta_k}$$

→ As we are sampling from model, while applying gradient ascent, we have to do a lot of calculation V, E, Q etc and we update the parameter for all trajectories, hence many times by taking small step (which make is complete ~~all~~ inefficient)

Importance Sampling & off-policy learning

importance sampling allow us to evaluate an expectation over a distribution π using sample taken from distibut γ .

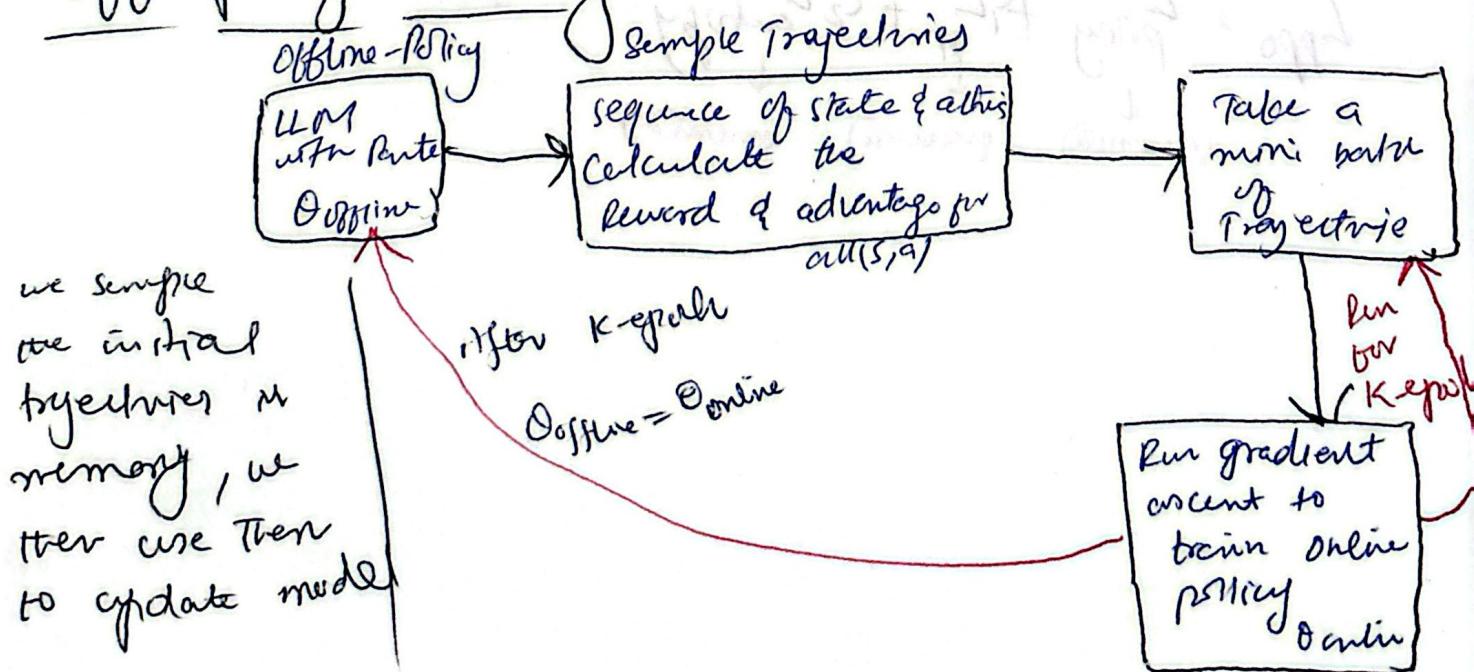
$$\begin{aligned} \underset{\pi \sim p(\pi)}{E}[f(a)] &= \int p(\pi) f(a) d\pi \\ &= \int \frac{q(a|\pi)}{q(a|\pi)} \pi(a) \cdot q(a) d\pi \\ &= \int q(a) \frac{p(a)}{q(a)} f(a) da \\ &= E_{a \sim q(a)} \left[\frac{p(a)}{q(a)} f(a) \right] \end{aligned}$$

So,

$$V_\theta(\pi(0)) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^T \gamma^t \log \pi_\theta(a_{i,t} | s_{i,t}) \right) A^\pi(s_i, a_i) \text{ becomes.}$$

$$V_{\theta_{\text{online}}}(\pi_{\text{online}}) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=0}^T V_{\theta_{\text{online}}} \left(\frac{\log \pi_{\text{online}}(a_{i,t} | s_{i,t})}{\log \pi_{\text{offline}}(a_{i,t} | s_{i,t})} \right) A^\pi(s_i, a_i) \right)$$

→ Off-policy learning.



The PPO loss: Proximal Policy Optimization

$$L_{\text{proxy}} = \min_{\tilde{\pi}} \left(\frac{\pi_{\theta}(a_t | s_t)}{\tilde{\pi}_{\theta}(a_t | s_t)} \hat{A}_t, \text{clip} \left(\frac{\pi_{\theta}(a_t | s_t)}{\tilde{\pi}_{\theta}(a_t | s_t)}, 1-\epsilon, 1+\epsilon \right) \hat{A}_t \right)$$

if High
→ not
→ too
Clip
to this

→ Because we want the model to not be very sure or unsure to select this word or not at all we want a little certainty (too sure will make it select often)

$$L_{\text{VF}} = \frac{1}{2} \| V_{\theta(s)} - \left(\sum_{t=0}^T r_{t+1} | s_0 = s \right) \|_2^2$$

output
value
scale actual value
based on trajectory
we sampled

$$L_{\text{entropy}} = - \sum_a p(a) \log p(a) \rightarrow$$

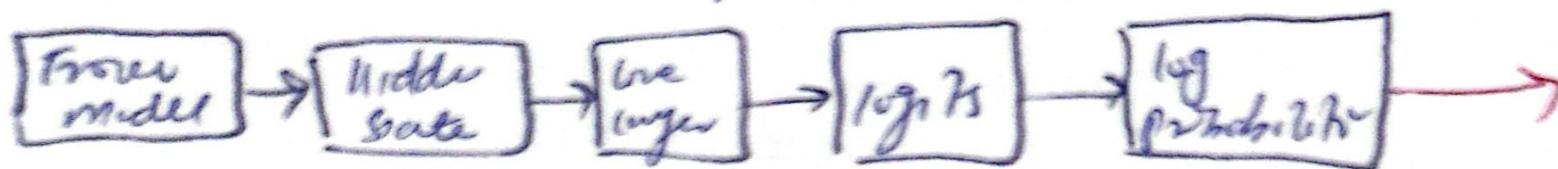
to force our model to explore more options, makes model rigid (select often to explore)
Other options

$$L_{\text{PPO}} = L_{\text{proxy}} + \gamma L_{\text{VF}} + \alpha L_{\text{entropy}}$$

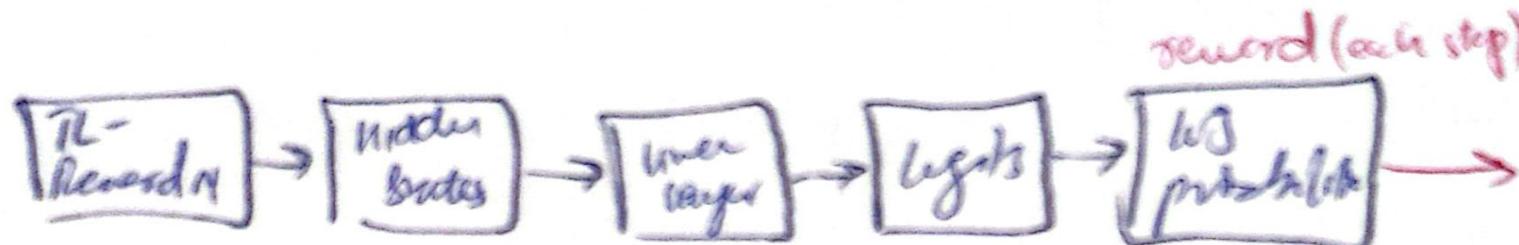
(maximize) (minimize) (maximize)

Reward hacking

We also want model to given answer that only gives a good reward but also makes sense and one similar to unaligned model.



we calculate
KL divergence
to penalize
the reward
by



$$\gamma_0 - \text{KL}(LP_F || LP_A)$$

Direct-Preference Optimization (DPO) - (Bradley-Terry)

Bradley-Terry Model

our goal is to learn a model that gives a numeric score (score) to pair of Q/A that resembles to the preference chosen by people.

$$P(y_w > y_i) = \frac{e^{\gamma^*(x, y_w)}}{e^{\gamma^*(x, y_w)} + e^{\gamma^*(x, y_i)}}$$

Derivation of loss that maximize P_r .

$$\frac{e^A}{e^A + e^B} = \frac{e^A / e^A}{e^A + e^B / e^A} = \frac{1}{e^A / e^B + 1 - 1}$$

$$\Rightarrow \frac{1}{1 + \left(\frac{e^A}{e^A + e^B} - 1\right)} \Rightarrow \frac{1}{1 + (e^B / e^A)}$$

$$\Rightarrow \frac{1}{1 + e^{B-A}} \Rightarrow \frac{1}{1 + e^{-(A-B)}} \Rightarrow \boxed{\sigma(A-B)}$$

Hence,

$$L_{\text{obj}} = -E_{(x, y_w, y_i) \sim D} [\log \sigma(\gamma_0(x, y_w) - \gamma_0(x, y_i))]$$

RLHF-objective score constraint the model to be not so different from original one

$$\int_{\text{RLHF}} \max_{\pi_0} \mathbb{E}_{x \sim D, y \sim \pi_0(y|x)} [\gamma_0(x, y) - \beta \mathbb{D}_{KL}[\pi_0(y|x) || \pi_{\text{org}}(y|x)]]$$

a temperature is a parameter during test generation that how deterministic model output is
 (0.0 = deterministic)
 (0.3-0.7 = Balanced)
 (1 = very random)

one way to maximize over preference correctly is use Maximum Likelihood Estimation.

MLE can estimate the parameters of score model such that the probability of ranking preference is correctly maximized.

What does it mean to maximize an objective function?

The Optimization Problem

- means to find values of one or more variables such that the fun is maximized. (for unconstrained problem)
(we find derivative of fun, set to zero, get values.)

constraint-optimization. (\mathcal{I}_{PL+g})

These are not differentiable so, we can not run gradient descent and forced to used algos like PPO.

Optimal Policy.

It is possible mathematically but very computationally heavy.

$$\pi_G(y|x) = \frac{1}{Z(x)} \pi_{ref}(y|x) \exp(\gamma_B \pi(x,y))$$

$$\log Z(x) = \sum_y \pi_{ref}(y|x) \exp(\gamma_B \pi(x,y))$$

But suppose, magically we have optimal policy
would be reward fun for optimal policy?

By leveraging:

$$\pi(x,y) = \beta \log \frac{\pi^*(y|x)}{\pi_{ref}(y|x)} + \beta \log Z(x)$$

$$\therefore \text{Bi-model} \rightarrow P(y_2 > p_1) = e^{\pi(x,y_2)}$$

we plug the reward
and get.

$$\frac{e^{\pi(x,y_2)}}{e^{\pi(x,y_2)} + e^{\pi(x,y_1)}} = \sigma(A - B)$$

$\pi^*(x)$, what
is in actual
we don't
have optimal
policy & don't
we pretend
we have
policy.

$$P(y_2 > y_1) = \sigma \left(\beta \log \frac{\pi^*(y_2|x)}{\pi_{ref}(y_2|x)} + \beta \log Z(x) - \beta \log \frac{\pi^*(y_1|x)}{\pi_{ref}(y_1|x)} - \beta \log Z(x) \right)$$

(maximize this)

$$L(\pi_\theta; \pi_{ref}) = -\mathbb{E}_{(y_1, y_2) \sim D} \left[\log \sigma \left(\beta \log \frac{\pi_\theta(y_1 | a)}{\pi_{ref}(y_1 | a)} - \beta \log \frac{\pi_\theta(y_2 | a)}{\pi_{ref}(y_2 | a)} \right) \right]$$

so, instead of optimizing reward function, we are optimizing policy (which defines an optimal reward fun).

Computing log-Probabilities:-

$$L_{DDO}(\pi_\theta; \pi_{ref}) = -\mathbb{E}_{\substack{\text{prompt} \\ \text{closed} \\ \text{true} \\ \text{false} \\ \text{answ} \\ \text{answ}}} \left[\log \left(\beta \log \frac{\pi_\theta(y_1 | a)}{\pi_{ref}(y_1 | a)} + \beta \log \frac{\pi_\theta(y_2 | a)}{\pi_{ref}(y_2 | a)} \right) \right]$$

