# Regressing Relative Reward

# Recap: KL-reg RL objective

$$J(\pi) = \mathbb{E}_{x \sim \nu} \left[ \mathbb{E}_{\tau \sim \pi(\cdot|x)} \hat{r}(x, \tau) - \beta \mathsf{KL} \left( \pi(\cdot \mid x) \,\Big|\, \pi_{ref}(\cdot \mid x) \right) \right]$$

$$\hat{\pi}(\tau \mid x) \propto \pi_{ref}(\tau \mid x) \cdot \exp \left( \frac{\hat{r}(x, \tau)}{\beta} \right)$$

Stay close to $\pi_{ref}$          Optimize reward

# Recap: DPO

$$\arg\max_{\theta} \sum_{x,\tau,\tau',z} \ln \frac{1}{1 + \exp\left(-z \cdot \beta \left(\ln \frac{\pi_\theta(\tau \mid x)}{\pi_{ref}(\tau \mid x)} - \ln \frac{\pi_\theta(\tau' \mid x)}{\pi_{ref}(\tau' \mid x)}\right)\right)}$$

# Recap: DPO

$$\arg\max_{\theta} \sum_{x,\tau,\tau',z} \ln \frac{1}{1 + \exp\left(-z \cdot \beta \left(\ln \frac{\pi_\theta(\tau \mid x)}{\pi_{ref}(\tau \mid x)} - \ln \frac{\pi_\theta(\tau' \mid x)}{\pi_{ref}(\tau' \mid x)}\right)\right)}$$

Use policies to model the reward difference (aka your LLM is your secret reward model)

# But DPO's performance isn't as strong as RM+PPO in practice..

Evaluation and Generation gap (aka evaluation is easier than generation..)

When a reward / verifier is easier to learn, RM + PPO can win…

# But DPO's performance isn't as strong as RM+PPO in practice..

Evaluation and Generation gap (aka evaluation is easier than generation..)

When a reward / verifier is easier to learn, RM + PPO can win…

1. DPO uses finite data + gradient descent to learn the generator directly

# But DPO's performance isn't as strong as RM+PPO in practice..

Evaluation and Generation gap (aka evaluation is easier than generation..)

When a reward / verifier is easier to learn, RM + PPO can win…

1. DPO uses finite data + gradient descent to learn the generator directly

2. PPO uses the RM, and can take advantage of unseen prompts and new training data…

# But DPO's performance isn't as strong as RM+PPO in practice..

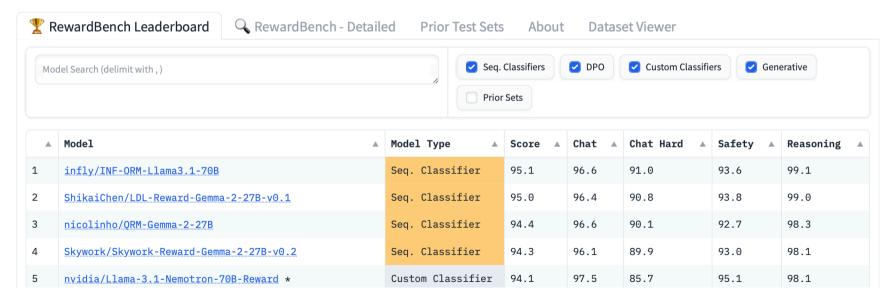PPO can also take advantage of the state-of-art RMs from the community…

## RewardBench: Evaluating Reward Models

**Evaluating the capabilities, safety, and pitfalls of reward models**

Code | Eval. Dataset | Prior Test Sets | Results | Paper | Total models: 165 | * Unverified models | ⚠️ Dataset Contamination |
Last restart (PST): 22:01 PDT, 28 Mar 2025

⚠️ Many of the top models were trained on unintentionally contaminated, AI-generated data, for more information, see this gist.

🏆 **RewardBench Leaderboard**    🔍 RewardBench - Detailed    Prior Test Sets    About    Dataset Viewer

Model Search (delimit with , )

☑ Seq. Classifiers    ☑ DPO    ☑ Custom Classifiers    ☑ Generative

☐ Prior Sets

| | Model | Model Type | Score | Chat | Chat Hard | Safety | Reasoning |
|---|---|---|---|---|---|---|---|
| 1 | infly/INF-ORM-Llama3.1-70B | Seq. Classifier | 95.1 | 96.6 | 91.0 | 93.6 | 99.1 |
| 2 | ShikaiChen/LDL-Reward-Gemma-2-27B-v0.1 | Seq. Classifier | 95.0 | 96.4 | 90.8 | 93.8 | 99.0 |
| 3 | nicolinho/QRM-Gemma-2-27B | Seq. Classifier | 94.4 | 96.6 | 90.1 | 92.7 | 98.3 |
| 4 | Skywork/Skywork-Reward-Gemma-2-27B-v0.2 | Seq. Classifier | 94.3 | 96.1 | 89.9 | 93.0 | 98.1 |
| 5 | nvidia/Llama-3.1-Nemotron-70B-Reward * | Custom Classifier | 94.1 | 97.5 | 85.7 | 95.1 | 98.1 |

https://huggingface.co/spaces/allenai/reward-bench

# Today's question

PPO can be very expensive, can we develop RL algorithm that is more efficient and may be more effective?

# Outline

1. Mirror descent — reward maximization subject to a KL reg to the old policy

2. Reparametrization trick and REBEL

3. Connections to old algorithms we learned

# Mirror Descent

Let us assume that we are given a reward function $r(x, \tau)$
(e.g., learned by ourselves or an open-source model)

# Mirror Descent

Let us assume that we are given a reward function $r(x, \tau)$
(e.g., learned by ourselves or an open-source model)

$$\text{Want to } \max_{\pi} \mathbb{E}_{x,\tau \sim \pi(\cdot|x)} \left[ r(x, \tau) \right]$$

$\pi_{\theta}$

# Mirror Descent

Mirror descent (MD) incrementally (iteratively) updates the policy:

Given $\pi_t$, we update to $\pi_{t+1}$ as follows:

# Mirror Descent

Mirror descent (MD) incrementally (iteratively) updates the policy:

Given $\pi_t$, we update to $\pi_{t+1}$ as follows:

$$\pi_{t+1} = \arg\max_{\pi} \mathbb{E}_{x,\tau \sim \pi(\cdot|x)} \left[ r(x,\tau) - \beta \mathsf{KL}\left( \pi(\,\cdot\,|\,x)\,|\,\pi_t(\,\cdot\,|\,x) \right) \right]$$

# Mirror Descent

Mirror descent (MD) incrementally (iteratively) updates the policy:

Given $\pi_t$ , we update to $\pi_{t+1}$ as follows:

$$\pi_{t+1} = \arg\max_\pi \mathbb{E}_{x,\tau \sim \pi(\cdot|x)} \left[ r(x,\tau) - \beta \mathsf{KL} \left( \pi(\cdot \mid x) \mid \pi_t(\cdot \mid x) \right) \right]$$

KL to the previous policy (e.g., recall NPG and the logic behind PPO's clipping)

# Mirror Descent

In theory, using the same idea we had from KL-reg RL, $\pi_{t+1}$ has a closed-form:

$$\pi_{t+1} = \arg\max_{\pi} \mathbb{E}_{x,\tau\sim\pi(\cdot|x)} \left[ r(x,\tau) - \beta\text{KL}\left(\pi(\,\cdot\,|x)\,|\,\pi_t(\,\cdot\,|x)\right) \right]$$

$$\Rightarrow \pi_{t+1}(\tau\,|\,x) \propto \pi_t(\tau\,|\,x)\exp\left(r(x,\tau)/\beta\right)$$

$\pi_{\text{ref}}$

$\pi_{t+2}(\tau|x) \propto \pi_{t+1}(\tau|x) \exp(r(x,y)/\beta)$

$\pi_{t+3}$ — —

$\pi_{t+\varphi}$

# Mirror Descent

In theory, using the same idea we had from KL-reg RL, $\pi_{t+1}$ has a closed-form:

$$\pi_{t+1} = \arg\max_{\pi} \mathbb{E}_{x,\tau \sim \pi(\cdot|x)} \left[ r(x,\tau) - \beta\mathsf{KL}\left(\pi(\cdot\,|\,x)\,|\,\pi_t(\cdot\,|\,x)\right) \right]$$

$$\Rightarrow \pi_{t+1}(\tau\,|\,x) \propto \pi_t(\tau\,|\,x)\exp\left(r(x,\tau)/\beta\right)$$

NN

Q: can we easily implement $\pi_{t+1}$?

# Mirror Descent

Ignoring the implementation issue, mirror descent in theory has very good convergence rate

# Mirror Descent

Ignoring the implementation issue, mirror descent in theory has very good convergence rate

$\pi_{x+1} \propto \pi_x \exp(\%_\beta)$

$\frac{1}{T}$

After T iterations, we can find a policy $\hat{\pi}$, s.t.,

$$\left| \mathbb{E}_{x,\tau \sim \hat{\pi}(\cdot|x)} r(x,\tau) - \mathbb{E}_{x,\tau \sim \pi^\star(\cdot|x)} r(x,\tau) \right| \leq O(1/T)$$

(Proof out of the scope, see CS6789)

# Outline

1. Mirror descent — reward maximization subject to a KL reg to the old policy

2. Reparametrization trick and REBEL

3. Connections to old algorithms we learned

# Reparameterization

Mirror descent indicates the following ideal update:

$$\pi_{t+1}(\tau \mid x) = \pi_t(\tau \mid x)\exp\left(r(x, \tau)/\beta\right)/Z(x)$$

$$\sum_\tau \pi_{t+1}(\cdot \mid x) = 1$$

$$Z(x) = \mathbb{E}_{\tau \sim \pi_t(\cdot \mid x)}\, \exp\left(\frac{r(x, \tau)}{\beta}\right)$$

# Reparameterization

Mirror descent indicates the following ideal update:

$$\pi_{t+1}(\tau \mid x) = \pi_t(\tau \mid x)\exp\left(r(x,\tau)/\beta\right)/Z(x)$$

---

1. Take log on both sides and rearrange terms, we get

$$r(x,\tau) = \beta\left(\ln\frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} + \ln Z(x)\right)$$

$$r(x,\tau) - r(x,\tau')$$

# Reparameterization

Mirror descent indicates the following ideal update:

$$\pi_{t+1}(\tau \mid x) = \pi_t(\tau \mid x) \exp\left(r(x,\tau)/\beta\right)/Z(x)$$

1. Take log on both sides and rearrange terms, we get

$$r(x,\tau) = \beta\left(\ln \frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} + \ln Z(x)\right)$$

2. Instead of modeling reward, we model reward difference to cancel $Z(x)$:

$$r(x,\tau) - r(x,\tau') = \beta\left(\ln \frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} - \ln \frac{\pi_{t+1}(\tau' \mid x)}{\pi_t(\tau' \mid x)}\right)$$

*ideal closed-form solution from MD*

# Reparameterization

We obtained the following relationship between $r$ and $\pi_{t+1}$ & $\pi_t$:

# Reparameterization

We obtained the following relationship between $r$ and $\pi_{t+1}$ & $\pi_t$:

$$\forall (x, \tau, \tau') : \; r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} - \ln \frac{\pi_{t+1}(\tau' \mid x)}{\pi_t(\tau' \mid x)} \right)$$

# Reparameterization

We obtained the following relationship between $r$ and $\pi_{t+1}$ & $\pi_t$:

$$\forall (x, \tau, \tau') : \; r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} - \ln \frac{\pi_{t+1}(\tau' \mid x)}{\pi_t(\tau' \mid x)} \right)$$

This indicates $\pi_{t+1}$ is **the minimizer** of the following least square regression problem:

# Reparameterization

We obtained the following relationship between $r$ and $\pi_{t+1}$ & $\pi_t$:

$$\forall (x, \tau, \tau') : \; r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} - \ln \frac{\pi_{t+1}(\tau' \mid x)}{\pi_t(\tau' \mid x)} \right)$$

This indicates $\pi_{t+1}$ is **the minimizer** of the following least square regression problem:

$$\mathbb{E}_{x, \tau, \tau'} \left[ \beta \left( \ln \frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} - \ln \frac{\pi_{t+1}(\tau' \mid x)}{\pi_t(\tau' \mid x)} \right) - \left( r(x, \tau) - r(x, \tau') \right) \right]^2$$

# Reparameterization

We obtained the following relationship between $r$ and $\pi_{t+1}$ & $\pi_t$:

$$\forall (x, \tau, \tau'): \ r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} - \ln \frac{\pi_{t+1}(\tau' \mid x)}{\pi_t(\tau' \mid x)} \right)$$

This indicates $\pi_{t+1}$ is **the minimizer** of the following least square regression problem:

$\pi_{t+1}$ should be the minimizer regardless of what the distribution is;

$$\mathbb{E}_{x, \tau, \tau'} \left[ \beta \left( \ln \frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} - \ln \frac{\pi_{t+1}(\tau' \mid x)}{\pi_t(\tau' \mid x)} \right) - \left( r(x, \tau) - r(x, \tau') \right) \right]^2$$

# Reparameterization

We obtained the following relationship between $r$ and $\pi_{t+1}$ & $\pi_t$:

$$\forall (x, \tau, \tau'): \; r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} - \ln \frac{\pi_{t+1}(\tau' \mid x)}{\pi_t(\tau' \mid x)} \right)$$

This indicates $\pi_{t+1}$ is **the minimizer** of the following least square regression problem:

$\pi_{t+1}$ should be the minimizer regardless of what the distribution is;

$$\mathbb{E}_{x, \tau, \tau'} \left[ \beta \left( \ln \frac{\pi_{t+1}(\tau \mid x)}{\pi_t(\tau \mid x)} - \ln \frac{\pi_{t+1}(\tau' \mid x)}{\pi_t(\tau' \mid x)} \right) - \left( r(x, \tau) - r(x, \tau') \right) \right]^2$$

In pratice, we often use

$$x, \tau \sim \pi_t(\cdot \mid x), \tau' \sim \pi_t(\cdot \mid x)$$

$x \sim \mathcal{V} \in$ Distribution of prompts

# Reparameterization

We obtained the following relationship between $r$ and $\pi_{t+1}$ & $\pi_t$:

$$\forall (x, \tau, \tau') : \; r(x, \tau) - r(x, \tau') = \beta \left( \ln \frac{\pi_{t+1}(\tau \,|\, x)}{\pi_t(\tau \,|\, x)} - \ln \frac{\pi_{t+1}(\tau' \,|\, x)}{\pi_t(\tau' \,|\, x)} \right)$$

This indicates $\pi_{t+1}$ is **the minimizer** of the following least square regression problem:

$\pi_{t+1}$ should be the minimizer regardless of what the distribution is;

$$\mathbb{E}_{x, \tau, \tau'} \left[ \beta \left( \ln \frac{\pi_{t+1}(\tau \,|\, x)}{\pi_t(\tau \,|\, x)} - \ln \frac{\pi_{t+1}(\tau' \,|\, x)}{\pi_t(\tau' \,|\, x)} \right) - \underline{\left( r(x, \tau) - r(x, \tau') \right)} \right]^2$$

In pratice, we often use
$x, \tau \sim \pi_t(\cdot \,|\, x), \tau' \sim \pi_t(\cdot \,|\, x)$

Relative reward

# REBEL algorithm

Put things together, we arrive at the following iterative algorithm:

Given $\pi_t$, we compute $\pi_{t+1}$ via least square regression:

# REBEL algorithm

Put things together, we arrive at the following iterative algorithm:

Given $\pi_t$, we compute $\pi_{t+1}$ via least square regression:

$$\pi_{t+1} = \arg\min_{\pi} \mathbb{E}_{x,(\tau,\tau')\sim\pi_t(\cdot|x)} \left( \beta \left( \ln \frac{\pi(\tau\,|\,x)}{\pi_t(\tau\,|\,x)} - \ln \frac{\pi(\tau'\,|\,x)}{\pi_t(\tau'\,|\,x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

# REBEL algorithm

Put things together, we arrive at the following iterative algorithm:

<span style="color:red">Given $\pi_t$, we compute $\pi_{t+1}$ via least square regression:</span>

$$\pi_{t+1} = \arg\min_{\pi} \mathbb{E}_{x,(\tau,\tau')\sim\pi_t(\cdot|x)} \left( \beta \left( \ln \frac{\pi(\tau\,|\,x)}{\pi_t(\tau\,|\,x)} - \ln \frac{\pi(\tau'\,|\,x)}{\pi_t(\tau'\,|\,x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

sample $\tau, \tau'$ from the latest policy
$\pi_t(\,\cdot\,|\,x)$, independently;

# REBEL algorithm

Put things together, we arrive at the following iterative algorithm:

Given $\pi_t$, we compute $\pi_{t+1}$ via least square regression:

$$\pi_{t+1} = \arg\min_{\pi} \mathbb{E}_{x,(\tau,\tau')\sim\pi_t(\cdot|x)} \left( \beta \left( \ln \frac{\pi(\tau\,|\,x)}{\pi_t(\tau\,|\,x)} - \ln \frac{\pi(\tau'\,|\,x)}{\pi_t(\tau'\,|\,x)} \right) - \underline{\left( r(x,\tau) - r(x,\tau') \right)} \right)^2$$

Relative reward

sample $\tau, \tau'$ from the latest policy $\pi_t(\,\cdot\,|\,x)$, independently;

# REBEL algorithm

Put things together, we arrive at the following iterative algorithm:

Given $\pi_t$, we compute $\pi_{t+1}$ via least square regression:

$$\pi_{t+1} = \arg\min_\pi \mathbb{E}_{x,(\tau,\tau')\sim\pi_t(\cdot|x)} \left( \beta \left( \ln \frac{\pi(\tau\,|\,x)}{\pi_t(\tau\,|\,x)} - \ln \frac{\pi(\tau'\,|\,x)}{\pi_t(\tau'\,|\,x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

$$\underbrace{\phantom{\beta \left( \ln \frac{\pi(\tau\,|\,x)}{\pi_t(\tau\,|\,x)} - \ln \frac{\pi(\tau'\,|\,x)}{\pi_t(\tau'\,|\,x)} \right)}}_{\text{Regressor}} \qquad \underbrace{\phantom{r(x,\tau) - r(x,\tau')}}_{\text{Relative reward}}$$

sample $\tau, \tau'$ from the latest policy $\pi_t(\,\cdot\,|\,x)$, independently;

# Difference between REBEL and DPO

Discussion: what is the difference between DPO, REBEL, and PPO

PPO: $x, \tau, \tau', z$

PPO: RM + online RL (Train Value function)

Rebel: RM + "DPO" like But Reward signal

$\tau_1 > \tau_2 > \tau_3$

$\tau_1 > \tau_2 - - > \tau_k$

$\tau_1 > \tau_2$

$\tau_2 > \tau_3$

$\tau_1 \tau_2$
$\tau_2 \tau_3$
$\tau_1 \tau_3$

$BT: \quad \tau_1 > \tau_3$

**Outline**

1. Mirror descent — reward maximization subject to a KL reg to the old policy

2. Reparametrization trick and REBEL

3. Connections to old algorithms we learned

# REBEL

Consider parameterized policy $\pi_\theta$, recall that rebel solves least square regression every iteration:

$$\theta_{t+1} = \arg\min_\theta \mathbb{E}_{x,(\tau,\tau')\sim\pi_{\theta_t}(\cdot|x)} \left( \beta \left( \ln \frac{\pi_\theta(\tau\,|\,x)}{\pi_{\theta_t}(\tau\,|\,x)} - \ln \frac{\pi_\theta(\tau'\,|\,x)}{\pi_{\theta_t}(\tau'\,|\,x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

$\Delta$

Regressor

Target:

By Running Adam on $\theta$

# REBEL

Consider parameterized policy $\pi_\theta$, recall that rebel solves least square regression every iteration:

$$\theta_{t+1} = \arg \min_\theta \mathbb{E}_{x,(\tau,\tau') \sim \pi_{\theta_t}(\cdot|x)} \left( \beta \left( \ln \frac{\pi_\theta(\tau \mid x)}{\pi_{\theta_t}(\tau \mid x)} - \ln \frac{\pi_\theta(\tau' \mid x)}{\pi_{\theta_t}(\tau' \mid x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

In practice, hard to solve it exactly

# REBEL

Consider parameterized policy $\pi_\theta$, recall that rebel solves least square regression every iteration:

$$\theta_{t+1} = \arg\min_{\theta} \mathbb{E}_{x,(\tau,\tau')\sim\pi_{\theta_t}(\cdot|x)} \left( \beta \left( \ln \frac{\pi_\theta(\tau\,|\,x)}{\pi_{\theta_t}(\tau\,|\,x)} - \ln \frac{\pi_\theta(\tau'\,|\,x)}{\pi_{\theta_t}(\tau'\,|\,x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

In practice, hard to solve it exactly

What happens if we solve it approximately? What happens if we perform

# REBEL

Consider parameterized policy $\pi_\theta$, recall that rebel solves least square regression every iteration:

$$\theta_{t+1} = \arg\min_\theta \mathbb{E}_{x,(\tau,\tau')\sim\pi_{\theta_t}(\cdot|x)} \left( \beta \left( \ln\frac{\pi_\theta(\tau\,|\,x)}{\pi_{\theta_t}(\tau\,|\,x)} - \ln\frac{\pi_\theta(\tau'\,|\,x)}{\pi_{\theta_t}(\tau'\,|\,x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

In practice, hard to solve it exactly

What happens if we solve it approximately? What happens if we perform

1. one step of gradient descent

2. one step of Gauss-newton method

# REBEL recovers variance-reduced policy gradient

Approximately solve the least square regression problem via just one step of gradient descent

$$\theta_{t+1} = \arg\min_{\theta} \mathbb{E}_{x,(\tau,\tau')\sim\pi_{\theta_t}(\cdot|x)} \underbrace{\left( \beta \left( \ln \frac{\pi_\theta(\tau\,|\,x)}{\pi_{\theta_t}(\tau\,|\,x)} - \ln \frac{\pi_\theta(\tau'\,|\,x)}{\pi_{\theta_t}(\tau'\,|\,x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2}_{:\ell(\theta)}$$

$$\ell(\theta)$$

# REBEL recovers variance-reduced policy gradient

Approximately solve the least square regression problem via just one step of gradient descent

$$\theta_{t+1} = \arg\min_{\theta} \mathbb{E}_{x,(\tau,\tau')\sim\pi_{\theta_t}(\cdot|x)} \left( \beta \left( \ln\frac{\pi_\theta(\tau\,|\,x)}{\pi_{\theta_t}(\tau\,|\,x)} - \ln\frac{\pi_\theta(\tau'\,|\,x)}{\pi_{\theta_t}(\tau'\,|\,x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad}_{:\ell(\theta)}$$

$\theta_{t+1} \Leftarrow \theta_t - \eta\,\nabla_\theta \ell(\theta_t)$   Let's try this out!

$$\nabla_\theta \ell(\theta) = \left( \beta\left[ \ln\frac{\pi_\theta(\tau|x)}{\pi_{\theta_t}(\tau|x)} - \ln\frac{\pi_\theta(\tau'|x)}{\pi_{\theta_t}(\tau'|x)} \right] - (r(x,\tau) - r(x,\tau')) \right)\left( \beta\cdot\nabla\ln\pi_{\theta'}(\tau|x) - \beta\nabla\ln\pi_\theta(\tau'|x) \right)$$

set $\theta = \theta_t$

$$= \beta\left( -(r(x,\tau) - r(x,\tau')) \right)\left( \nabla\ln\pi_{\theta_t}(\tau|x) - \nabla\ln\pi_{\theta_t}(\tau'|x) \right)$$

$$= \beta\left( \nabla\ln\pi_{\theta_t}(\tau|x)(r(x,\tau) - r(x,\tau')) \right) \ominus \beta\,\nabla\ln\pi_{\theta_t}(\tau'|x)(r(x,\tau') - r(x,\tau))$$

# REBEL recovers variance-reduced NPG

Approximately solve the least square regression problem via just one step of Gauss-newton

$$\theta_{t+1} = \arg\min_{\theta} \mathbb{E}_{x,(\tau,\tau') \sim \pi_{\theta_t}(\cdot|x)} \left( \beta \left( \ln \frac{\pi_\theta(\tau\,|\,x)}{\pi_{\theta_t}(\tau\,|\,x)} - \ln \frac{\pi_\theta(\tau'\,|\,x)}{\pi_{\theta_t}(\tau'\,|\,x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

# REBEL recovers variance-reduced NPG

Approximately solve the least square regression problem via just one step of Gauss-newton

$$\theta_{t+1} = \arg\min_\theta \mathbb{E}_{x,(\tau,\tau')\sim\pi_{\theta_t}(\cdot|x)} \left( \beta\left( \ln\frac{\pi_\theta(\tau|x)}{\pi_{\theta_t}(\tau|x)} - \ln\frac{\pi_\theta(\tau'|x)}{\pi_{\theta_t}(\tau'|x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

GN approximate the non-linear part inside the square ***via first-order Taylor expansion at*** $\theta_t$

$$\ln\frac{\pi_\theta(\tau|x)}{\pi_{\theta_t}(\tau|x)} - \ln\frac{\pi_\theta(\tau'|x)}{\pi_{\theta_t}(\tau'|x)} \approx \left( \nabla\ln\pi_{\theta_t}(\tau|x) - \nabla\ln\pi_{\theta_t}(\tau'|x) \right)^\top (\theta - \theta_t)$$

first-ord

Term

# REBEL recovers variance-reduced NPG

Approximately solve the least square regression problem via just one step of Gauss-newton

$$\theta_{t+1} = \arg\min_{\theta} \mathbb{E}_{x,(\tau,\tau')\sim\pi_{\theta_t}(\cdot|x)} \left( \beta \left( \ln\frac{\pi_\theta(\tau|x)}{\pi_{\theta_t}(\tau|x)} - \ln\frac{\pi_\theta(\tau'|x)}{\pi_{\theta_t}(\tau'|x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

GN approximate the non-linear part inside the square *via first-order Taylor expansion at* $\theta_t$

$$\ln\frac{\pi_\theta(\tau|x)}{\pi_{\theta_t}(\tau|x)} - \ln\frac{\pi_\theta(\tau'|x)}{\pi_{\theta_t}(\tau'|x)} \approx \left( \nabla\ln\pi_{\theta_t}(\tau|x) - \nabla\ln\pi_{\theta_t}(\tau'|x) \right)^\top (\theta - \theta_t)$$

Plug the linear approximation back and solve for $\theta$:

$$\theta_{t+1} = \arg\min_{\theta} \mathbb{E}_{x,(\tau,\tau')\sim\pi_{\theta_t}(\cdot|x)} \left( \beta \left( \left( \nabla\ln\pi_{\theta_t}(\tau|x) - \nabla\ln\pi_{\theta_t}(\tau'|x) \right)^\top (\theta - \theta_t) \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

feature      linear in $\theta$

# REBEL recovers variance-reduced NPG

Approximately solve the least square regression problem via just one step of Gauss-newton

$$\theta_{t+1} = \arg\min_{\theta} \mathbb{E}_{x,(\tau,\tau')\sim\pi_{\theta_t}(\cdot|x)} \left( \beta \left( \ln \frac{\pi_{\theta}(\tau|x)}{\pi_{\theta_t}(\tau|x)} - \ln \frac{\pi_{\theta}(\tau'|x)}{\pi_{\theta_t}(\tau'|x)} \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

GN approximate the non-linear part inside the square **via first-order Taylor expansion at** $\theta_t$

$$\ln \frac{\pi_{\theta}(\tau|x)}{\pi_{\theta_t}(\tau|x)} - \ln \frac{\pi_{\theta}(\tau'|x)}{\pi_{\theta_t}(\tau'|x)} \approx \left( \nabla \ln \pi_{\theta_t}(\tau|x) - \nabla \ln \pi_{\theta_t}(\tau'|x) \right)^\top (\theta - \theta_t)$$

Plug the linear approximation back and solve for $\theta$:

$$\theta_{t+1} = \arg\min_{\theta} \mathbb{E}_{x,(\tau,\tau')\sim\pi_{\theta_t}(\cdot|x)} \left( \beta \left( \left( \nabla \ln \pi_{\theta_t}(\tau|x) - \nabla \ln \pi_{\theta_t}(\tau'|x) \right)^\top (\theta - \theta_t) \right) - \left( r(x,\tau) - r(x,\tau') \right) \right)^2$$

Claim: this recovers the NPG update procedure (try this out after class!)

$$F_{\theta_t}^{-1} \left( \nabla \ln \pi_{\theta}(\tau|x) \cdot r \right)$$

# Using RL to optimize 7B size model on TL;DR

| Model size | Algorithm | Winrate ($\uparrow$) |
|---|---|---|
| | SFT | 45.2 ($\pm$2.49) |
| | DPO | 68.4 ($\pm$2.01) |
| | REINFORCE | 70.7* |
| 6.9B | PPO | 77.6‡ |
| | RLOO ($k = 2$) | 74.2* |
| | RLOO ($k = 4$) | <u>77.9*</u> |
| | REBEL | **78.1** ($\pm$1.74) |

\* directly obtained from Ahmadian et al. (2024)

‡ directly obtained from Huang et al. (2024)

[REBEL: Reinforcement Learning via Regressing Relative Rewards, Neurips 2024]

# Using RL to optimize 7B size model on TL;DR

| Model size | Algorithm | Winrate ($\uparrow$) |
|---|---|---|
| | SFT | 45.2 ($\pm$2.49) |
| | DPO | 68.4 ($\pm$2.01) |
| | REINFORCE | 70.7* |
| 6.9B | PPO | 77.6‡ |
| | RLOO ($k = 2$) | 74.2* |
| | RLOO ($k = 4$) | 77.9* |
| | REBEL | **78.1** ($\pm$1.74) |

1. Online RL + RM is often much better than DPO

\* directly obtained from Ahmadian et al. (2024)

‡ directly obtained from Huang et al. (2024)

[REBEL: Reinforcement Learning via Regressing Relative Rewards, Neurips 2024]

# Using RL to optimize 7B size model on TL;DR

| Model size | Algorithm | Winrate ($\uparrow$) |
|---|---|---|
| | SFT | 45.2 ($\pm$2.49) |
| | DPO | 68.4 ($\pm$2.01) |
| | REINFORCE | 70.7* |
| 6.9B | PPO | 77.6$\ddagger$ |
| | RLOO ($k=2$) | 74.2* |
| | RLOO ($k=4$) | 77.9* |
| | REBEL | **78.1** ($\pm$1.74) |

\* directly obtained from Ahmadian et al. (2024)

$\ddagger$ directly obtained from Huang et al. (2024)

[REBEL: Reinforcement Learning via Regressing Relative Rewards, Neurips 2024]

RLOO (k=0)

x $\tau_1, \tau_2, \tau_3, \tau_4$

1. Online RL + RM is often much better than DPO

2. REBEL is in par w/ PPO and RLOO (k=4), but much more computation and memory efficient

# Swimmer experiments in openAI Gym

Given the same amount of **labeled** preference data, rebel can continue learning using fresh online data