

Recall Traces: Backtracking Models for Efficient Reinforcement Learning

Anirudh Goyal¹, Philemon Brakel², William Fedus¹, Soumye Singhal^{1,4}, Timothy Lillicrap², Sergey Levine⁵, Hugo Larochelle³, Yoshua Bengio¹

ICLR 2019

¹Mila, University of Montreal, ² Google Deepmind, ³ Google Brain, ⁴ IIT Kanpur, ⁵ University of California, Berkeley. Corresponding author :anirudhgoyal9119@gmail.com

Motivation



In many environments only a tiny subset of all states yield high reward. We may want to preferentially train on those high-reward states and the probable trajectories leading to them.

sample efficiency

Producing Intended High Value States

□ Method 1

Rely on picking the most valuable states stored in the replay buffer ${\cal B}$.

(off-policy method) (on-policy method)

Method 2

Goal GAN^[1]: goal states g are produced via a Generative Adversarial Network.

Algorithm 2 Produce High Value States

Require: Critic V(s) **Require:** D; transformation 'decoder' from g to s **Require:** Experience buffer \mathcal{B} with tuples (s_t, a_t, r_t, s_{t+1}) **Require:** gen_state; boolean whether to generate states **Require:** GAN, some generative model trained to model high-value goal states 1: if gen_state then 2: $g \sim GAN$ 3: $D: g \mapsto s$ 4: Return s5: else 6: Return $argmax(V(s)) \forall s \in \mathcal{B}$ 7: end if

[1] Florensa C, Held D, Geng X, et al. Automatic goal generation for reinforcement learning agents[C]//International conference on machine learning. PMLR, 2018: 1515-1528.

D Backtracking Model $p(s_t, a_t | s_{t+1})$

Predict the preceding states that terminate at a given high-reward state.

Start from a high value state (or one that is estimated to have high value), predict and sample which (state, action)-tuples may have led to that high value state.

Recall Traces

use these traces to improve a policy.

Generating Recall Traces



$$q_{\phi}(\Delta s_t, a_t | s_{t+1}) = q(\Delta s_t | a_t, s_{t+1})q(a_t | s_{t+1})$$

Generating Recall Traces

 $a_t \sim q(a_t | s_{t+1})$ $\Delta s_t \sim q(\Delta s_t | a_t, s_{t+1})$ $s_t = \Delta s_t + s_{t+1}$

Improving The Policy From The Recall Traces

Algorithm 1 Improve Policy via Recall Traces and Backtracking Model

Require: RL algorithm with parameterized policy (i.e. TRPO, Actor-Critic) **Require:** Agent policy $\pi_{\theta}(a|s)$ **Require:** Backtracking model $B_{\phi} = q_{\phi}(\Delta s_t, a_t | s_{t+1})$ **Require:** Critic V(s)**Require:** k quantile of best state values used to train backtracking model, k_{trai} number of trajectories filtered by returns. **Require:** N; number of backward trajectories per target state **Require:** α , β ; forward, backward learning rates 1: Randomly initialize agent policy parameters θ 2: Randomly initialize backtracking model parameters ϕ 3: **for** t = 1 to *K* **do** Execute policy to produce trajectory τ 4: Add trajectory $\tau = (s_1, a_1, r_1, \cdots, s_T, a_T, r_T)$ in \mathcal{B} 5: 6: Estimate $\nabla_{\theta} R(\pi_{\theta})$ from RL algorithm $\theta \leftarrow \theta + \alpha \nabla_{\theta} R(\pi_{\theta})$ 7: Compute $\mathcal{L}_{\mathcal{B}}$ via Equation 2, using top k% valuable states from top k_{trai} trajectories in \mathcal{B} 8: $\phi \leftarrow \phi + \beta \nabla_{\phi} \mathcal{L}_{\mathcal{B}}$ 9: Obtain target high value state s (see Algorithm 2 for details) 10: Generate N recall traces $\tilde{\tau}$ for s using $B_{\phi}(s)$ 11: Compute imitation loss $\mathcal{L}_{\mathcal{I}}$ via Equation 3 12: $\theta \leftarrow \theta + \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}}$ 13: 14: end for

Improving The Policy From The Recall Traces



$$\mathcal{L}_{\mathcal{I}} = \sum_{t=0}^{T} \log p(a_t|s_t) = \sum_{t=0}^{T} \log \pi_{\theta}(a_t|s_t), \quad (3)$$

 \Box Variational Inference (RL Model: θ , Backtracking Model: ϕ)

$$p(R > L| heta) = rac{p(R > L, au| heta)}{p(au|R > L, heta)} \ \log p(R > L| heta) = \log rac{p(R > L, au| heta)}{q(au)} - \log rac{p(au|R > L, heta)}{q(au)}$$

左边 =
$$\int q(\tau) \log p(R > L|\theta) d\tau = \log p(R > L|\theta)$$

右边 = $\int q(\tau) \log \frac{p(R > L, \tau|\theta)}{q(\tau)} d\tau - \int q(\tau) \log \frac{p(\tau|R > L, \theta)}{q(\tau)} d\tau$
= ELBO + KL $(q(\tau)||p(\tau|R > L, \theta))$

$$egin{aligned} &\therefore \log p(R > L| heta) = & ext{ELBO} + ext{KL}(q(au)||p(au|R > L, heta)) \ &\geq & ext{ELBO} \end{aligned}$$

 \Box EM Algorithm (RL Model: θ , Backtracking Model: ϕ) (1) E-step: 固定 θ , 优化 ϕ , 通过最小化KL来最大化ELBO $\phi = argmin \operatorname{KL}(q_{\phi}(au) || p(au | R > L, heta))$ $= argmin \operatorname{KL}(p(au | R > L, heta) || q_{\phi}(au))$ $\therefore \operatorname{KL}(p(au|R>L, heta)||q_{\phi}(au)) = \int p(au|R>L, heta)\lograc{p(au|R>L, heta)}{q_{\phi}(au)}\,\mathrm{d} au$ $=\int p(au|R>L, heta)\log p(au|R>L, heta)\,\mathrm{d} au-\int p(au|R>L, heta)\log q_{\phi}(au)\,\mathrm{d} au$ $\therefore \phi = argmax_{\phi} \int p(\tau | R > L, \theta) \log q_{\phi}(\tau) \, \mathrm{d}\tau$ $= \! \mathop{argmax}\limits_{\phi} E_{ au \sim p(au \mid R > L, heta)} [\log q_{\phi}(au)]$

 \blacksquare EM Algorithm (RL Model: θ , Backtracking Model: ϕ)

(2) M-step: 固定 ϕ , 优化 θ , 最大化ELBO

$$\begin{split} \text{ELBO} &= \int q_{\phi}(\tau) \log \frac{p(R > L, \tau | \theta)}{q_{\phi}(\tau)} \, \mathrm{d}\tau \\ &= \int q_{\phi}(\tau) \log p(R > L, \tau | \theta) \, \mathrm{d}\tau - \int q_{\phi}(\tau) \log q_{\phi}(\tau) \, \mathrm{d}\tau \\ &\therefore \theta = \operatorname*{argmax}_{\theta} \text{ELBO} = \operatorname*{argmax}_{\theta} \int q_{\phi}(\tau) \log p(R > L, \tau | \theta) \, \mathrm{d}\tau \\ &\because q_{\phi}(\tau) = p(\tau | R > L, \theta_{t}) \\ &\therefore \theta = \operatorname*{argmax}_{\theta} E_{\tau \sim q_{\phi}(\tau)} [\log p(R > L, \tau | \theta)] \\ &= \operatorname*{argmax}_{\theta} E_{\tau \sim q_{\phi}(\tau)} [\log p(\tau | \theta)] \end{split}$$

Experiments : Access To True Backtracking Model



Figure 2: Training curves from the Four Room Environment for the Actor-Critic baseline (blue) and the backtracking model augmented Actor-Critic (orange). For the size-19 environment, several of the Actor-Critic baselines failed to converge, whereas the augmented recall trace model always succeeded in the number of training steps considered. For additional results see Figure 9 in Appendix.

Experiments : Comparison With Prioritized Experience Replay (PER)



Figure 3: Visitation count visualization of trained policies for PER (left) and Recall Traces (right) for two 4-room grid sizes.



Figure 4: Plots for reward vs. time steps, comparing the performance of recall traces (labeled Back-trackingModel), PER and baseline Actor Critic (AC).

Experiments : Learning Backtracking Model



Figure 6: Our model as compared to TRPO. For TRPO baselines, except walker, we ran with 5 different random seeds. For our model, we ran with 5 different random seeds.



Figure 7: Our model as compared to SAC. We ran SAC baselines with 2 different random seeds. For our model, we ran with 5 different random seeds.

Discussion
$$q = \{a'_1, a'_2, a'_3, \dots\}, \text{ where } a'_t = a_{E_t} + \nu$$

and $\tau_E = \{(s_{E_1}, a_{E_1}), (s_{E_2}, a_{E_2}), \dots\}.$

(3)
$$L_{\mathbf{u}} = -\mathbb{E}_{\pi_{E}} \left[\log D_{\mathbf{u}}(s, a) \right] - \mathbb{E}_{\pi_{\phi}} \left[\log(1 - D_{\mathbf{u}}(s, a)) \right], \quad (1)$$
$$L_{\phi} = \mathbb{E}_{\pi_{\phi}} \left[\log(1 - D_{u}(s, a)) \right] + \lambda ||a - a'||_{2}^{2}. \quad (5)$$

Corrected Augmentation for Trajectories (CAT)^[2]



Fig. 3. Detailed overview of stage 1, which performs Corrected Augmentation For Trajectories. The architecture is semi-supervised since it is guided by unlabelled distorted actions.



9

10 end

Discussion

 $egin{aligned} L_u &= -E_{\pi_E}[log D_u(s,a)] - E_{\pi_\phi}[log(1-D_u(s,a))] \ L_\phi &= E_{\pi_\phi}[log(1-D_u(s,a))] \end{aligned}$



HalfCheetah-v2



	方法	增广轨迹 条数	增广avgret	BC训练是否使用 专家轨迹	BC训练后验证100条轨 迹的avgret
	expert	0	2439.8168651239403	是	991.2917141601323
	cat	497	2221.0755368804244	否	2906.3392535252096
	noise network($\lambda=0.9$)	497	1172.8915687165966	否	1341.52380444484
•	mixup($lpha=0.1$)	42	1124.827231400378	否	1240.4956561966935

 $\lambda \sim Beta(lpha, lpha); \lambda' = max(\lambda, 1 - \lambda) \ a' = \lambda' a_1 + (1 - \lambda') a_2;$

Discussion

$$egin{aligned} L_u &= -E_{\pi_E}[log D_u(s,a)] - E_{\pi_\phi}[log(1-D_u(s,a))] \ L_\phi &= E_{\pi_\phi}[log(1-D_u(s,a))] \end{aligned}$$

16/17



	方法	增广轨迹 条数	增广avgret	BC训练是否使用 专家轨迹	BC训练后验证100条轨 迹的avgret
	expert	0	2439.8168651239403	是	991.2917141601323
tail	cat	497	2221.0755368804244	否	2906.3392535252096
head terminal state	noise network($\lambda=0.9)$	497	1172.8915687165966	否	1341.52380444484
expert trajectory	mixup($lpha=0.1$)	42	1124.827231400378	否	1240.4956561966935
initial state middle	BDM	497		否	319.0855923529666
	BDM_CL	497		否	1095.62363152283

Thanks