

## Uma viagem no tempo: de 1992 até 2017

Participação especial: Reinforce,  
A2C e PPO

Material adicional, 2025

FABRÍCIO J. BARTH

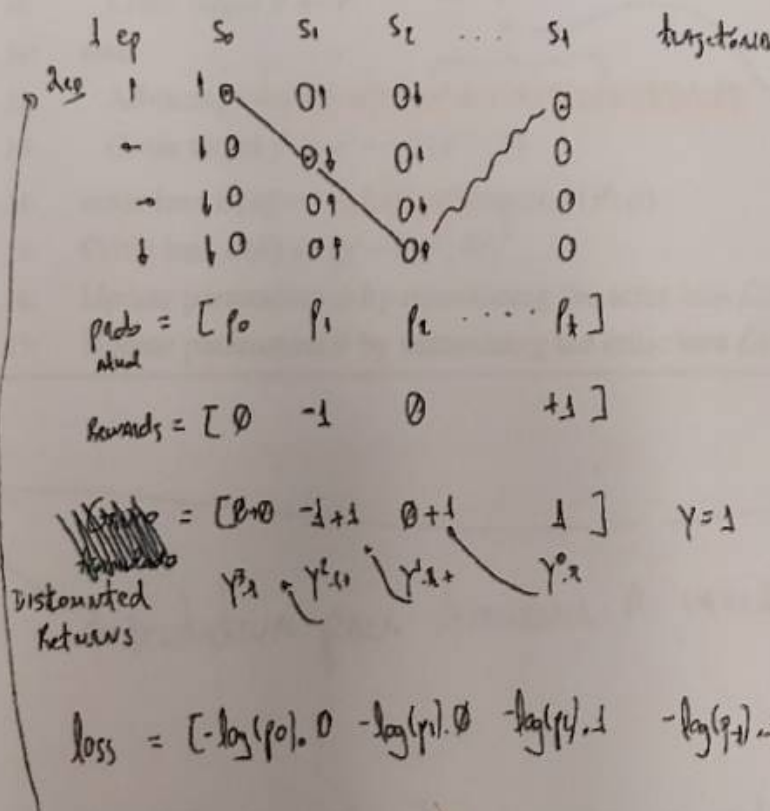
### Algorithm 13 REINFORCE

- 1: Initialize policy network  $\pi$  with random parameters  $\phi$
- 2: Repeat for every episode:
- 3: **for** time step  $t=0, 1, 2, \dots, T-1$  **do**
- 4:   Observe current state  $s^t$
- 5:   Sample action  $a^t \sim \pi(\cdot | s^t; \phi)$
- 6:   Apply action  $a^t$ ; observe reward  $r^t$  and next state  $s^{t+1}$
- 7:   Loss  $\mathcal{L}(\phi) \leftarrow -\frac{1}{T} \sum_{t=0}^{T-1} \left( \sum_{\tau=t}^{T-1} \gamma^{\tau-t} r^{\tau} \right) \log \pi(a^t | s^t; \phi)$
- 8:   Update parameters  $\phi$  by minimizing the loss  $\mathcal{L}(\phi)$

gera uma trajetória  
 $s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T$   
 segundo  $\pi$

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} \times r_k$$

Simulação do aprendizado usando Reinforce:



backpropagation

\* executar o programa  
 exemplo

pode ter alta variância,  
 o que leva a alta variância  
 dos gradientes e instabilidade  
 no treinamento.

#### Algorithm 14 Simplified advantage actor-critic (A2C)

- 1: Initialize actor network  $\pi$  with random parameters  $\phi$
- 2: Initialize critic network  $V$  with random parameters  $\theta$
- 3: Repeat for every episode:
- 4: **for** time step  $t = 0, 1, 2, \dots$  **do**
- 5:   Observe current state  $s^t$
- 6:   Sample action  $a^t \sim \pi(\cdot | s^t; \phi)$
- 7:   Apply action  $a^t$ ; observe reward  $r^t$  and next state  $s^{t+1}$
- 8:   **if**  $s^{t+1}$  is terminal **then**
- 9:     Advantage  $Adv(s^t, a^t) \leftarrow r^t - V(s^t; \theta)$
- 10:    Critic target  $y^t \leftarrow r^t$
- 11:   **else**
- 12:     Advantage  $Adv(s^t, a^t) \leftarrow r^t + \gamma V(s^{t+1}; \theta) - V(s^t; \theta)$
- 13:     Critic target  $y^t \leftarrow r^t + \gamma V(s^{t+1}; \theta)$
- 14:    actor loss  $\mathcal{L}(\phi) \leftarrow -Adv(s^t, a^t) \log \pi(a^t | s^t; \phi)$
- 15:    Critic loss  $\mathcal{L}(\theta) \leftarrow (y^t - V(s^t; \theta))^2$
- 16:    Update parameters  $\phi$  by minimizing the actor loss  $\mathcal{L}(\phi)$
- 17:    Update parameters  $\theta$  by minimizing the critic loss  $\mathcal{L}(\theta)$

Rede usada  
p/ escolher  
A ação dada  
um  
estado

Rede usada p/ o cálculo da Vantagem

cálculo da Vantagem  
menos um  
baseline

Δ tentativa para diminuir a variância.

equação similar ao Reinforce, com a adição do conceito de Vantagem

---

**Algorithm 15** Simplified proximal policy optimization (PPO)

---

- 1: Initialize actor network  $\pi$  with random parameters  $\phi$
  - 2: Initialize critic network  $V$  with random parameters  $\theta$
  - 3: Repeat for every episode:
  - 4: **for** time step  $t = 0, 1, 2, \dots$  **do**
  - 5:   Observe current state  $s^t$
  - 6:   Sample action  $a^t \sim \pi(\cdot | s^t; \phi)$
  - 7:   Apply action  $a^t$ ; observe reward  $r^t$  and next state  $s^{t+1}$
  - 8:    $\pi_\beta(a^t | s^t) \leftarrow \pi(a^t | s^t; \phi)$
  - 9:   **for** epoch  $e = 1, \dots, N_e$  **do**
  - 10:      $\rho(s^t, a^t) \leftarrow \pi(a^t | s^t; \phi) \div \pi_\beta(a^t | s^t)$
  - 11:     **if**  $s^{t+1}$  is terminal **then**
  - 12:       Advantage  $Adv(s^t, a^t) \leftarrow r^t - V(s^t; \theta)$
  - 13:       Critic target  $y^t \leftarrow r^t$
  - 14:     **else**
  - 15:       Advantage  $Adv(s^t, a^t) \leftarrow r^t + \gamma V(s^{t+1}; \theta) - V(s^t; \theta)$
  - 16:       Critic target  $y^t \leftarrow r^t + \gamma V(s^{t+1}; \theta)$
  - 17:     Actor loss  $\mathcal{L}(\phi) \leftarrow -\min \left( \begin{array}{l} \rho(s^t, a^t) Adv(s^t, a^t), \\ \text{clip}(\rho(s^t, a^t), 1 - \epsilon, 1 + \epsilon) Adv(s^t, a^t) \end{array} \right)$
  - 18:     Critic loss  $\mathcal{L}(\theta) \leftarrow (y^t - V(s^t; \theta))^2$
  - 19:     Update parameters  $\phi$  by minimizing the actor loss  $\mathcal{L}(\phi)$
  - 20:     Update parameters  $\theta$  by minimizing the critic loss  $\mathcal{L}(\theta)$
- 

10)  $\pi(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)}$  como uma forma de estabilizar o treinamento.

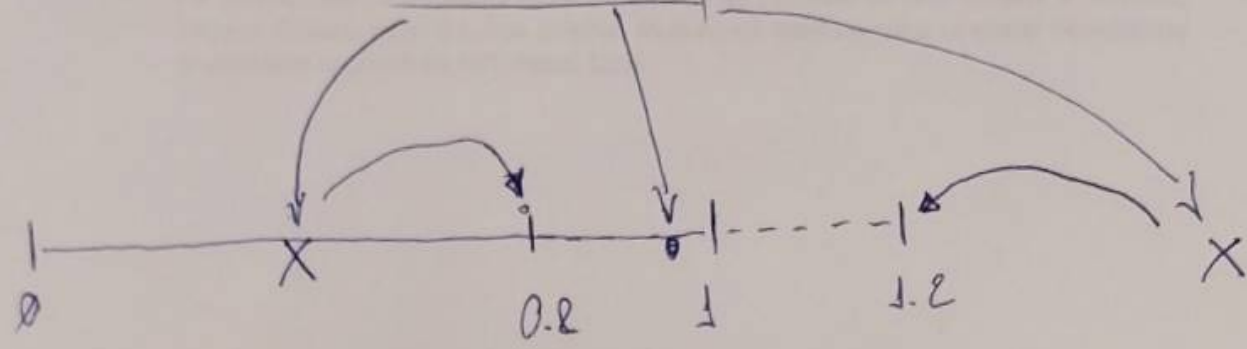
No entanto, mesmo assim podemos ter alguns problemas, por exemplo:

$$\frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)} = \frac{0.4}{0.0001} = 4000 \times \text{GAIN}$$

Vamos ter um número grande para o  $\text{loss}(\theta)$ !



probabilidade atual  $\Rightarrow \pi_{\theta}(a_t | s_t)$   
 probabilidade antiga  $\Rightarrow \pi_{\theta_{old}}(a_t | s_t)$



Probabilidades são valores entre 0 e 1 e neste caso queremos deixar mais próximo 1. pq elas não deveriam mudar muito.

$$\min \left( p(s_t, a_t) \cdot Adv(s_t, a_t), \boxed{\text{clip}(p(s_t, a_t), 1-\epsilon, 1+\epsilon)} \cdot Adv(s_t, a_t) \right)$$

⇓  
linha 17

if  $p(s_t, a_t)$  entre  $1-\epsilon$  and  $1+\epsilon$ , então fica  
 if menor que  $1-\epsilon$ , então  $1-\epsilon$   
 if maior que  $1+\epsilon$ , então  $1+\epsilon$

## Referências

Os pseudo-códigos utilizados neste material foram retirados do livro Stefano V. Albrecht, Filippos Christianos, and Lukas Schäfer. Multi-Agent Reinforcement Learning: Foundations and Modern Approaches. MIT Press, 2024.