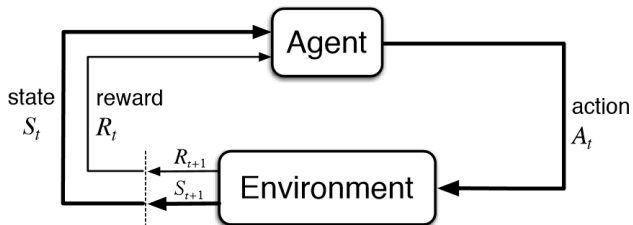


Reinforcement Learning via Policy Optimization

Hanxiao Liu

April 17, 2018

Reinforcement Learning

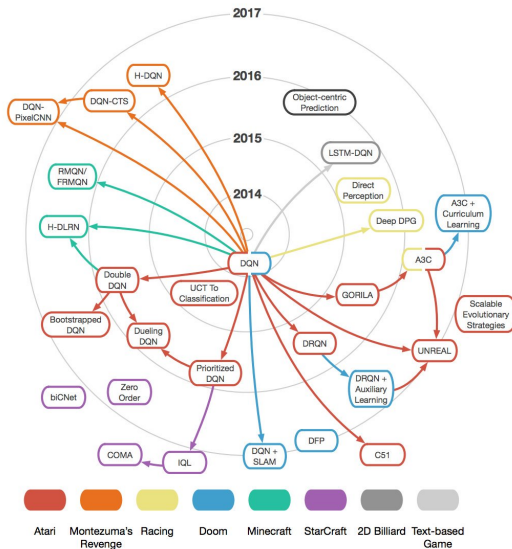


Policy $a \sim \pi(s)$

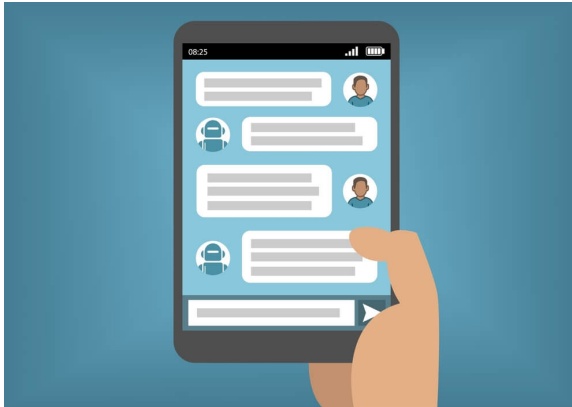
Example - Game



Application - More Games



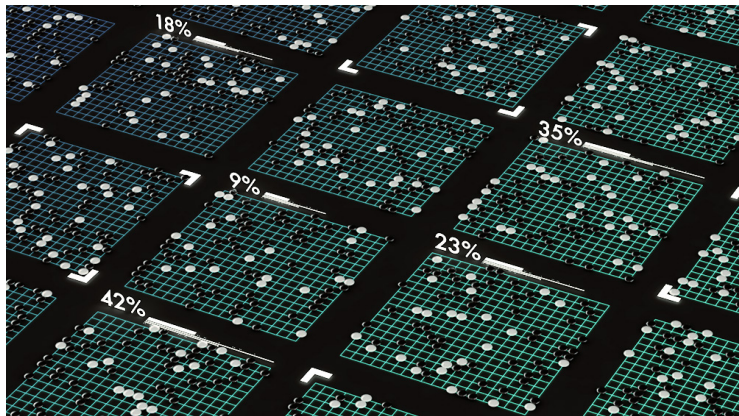
Application - ChatBot



Application - Robotics



Application - Combinatorial Problems



Supervised Learning vs RL

Supervised setup

- ▶ $s_t \sim P(\cdot)$
- ▶ $a_t = \pi(s_t)$
- ▶ immediate reward

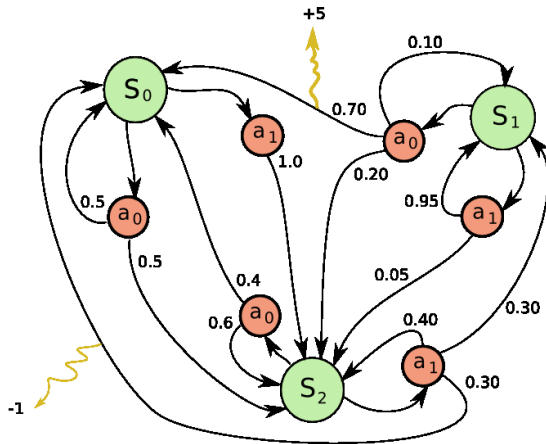
Reinforcement setup

- ▶ $s_t \sim P(\cdot | s_{t-1}, a_{t-1})$
- ▶ $a_t \sim \pi(s_t)$
- ▶ delayed reward

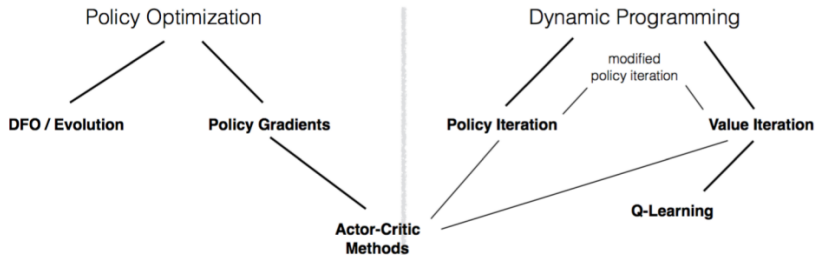
Both aims to maximize the total reward.

Markov Decision Process

What is the underlying assumption?



Landscape



Monte Carlo

- ▶ Return $R(\tau) = \sum_{t=1}^T R(s_t, a_t)$.
- ▶ Trajectory $\tau = s_0, a_0, \dots, s_T, a_T$

Goal: finding a good π such that R is maximized.

Policy evaluation via MC

1. Sample a trajectory τ given π .
2. Record $R(\tau)$

Repeat the above for many times and take the average.

Policy Gradient

Let's parametrize π using a linear function/neural net.

► $a \sim \pi_{\theta}(s)$

Expected return

$$\max_{\theta} U(\theta) = \max_{\theta} \sum_{\tau} P(\tau; \pi_{\theta}) R(\tau) \quad (1)$$

Heuristic: Raise the probability of good trajectories.

1. Sample a trajectory τ under π_{θ}
2. Update θ using gradient $\nabla_{\theta} \log P(\tau; \pi_{\theta}) R(\tau)$.

Policy Gradient

The log-derivative trick

$$\nabla_{\theta} U(\theta) = \sum_{\tau} \nabla_{\theta} P(\tau; \pi_{\theta}) R(\tau) \quad (2)$$

$$= \sum_{\tau} \frac{P(\tau; \pi_{\theta})}{P(\tau; \pi_{\theta})} \nabla_{\theta} P(\tau; \pi_{\theta}) R(\tau) \quad (3)$$

$$= \sum_{\tau} P(\tau; \pi_{\theta}) \nabla_{\theta} \log P(\tau; \pi_{\theta}) R(\tau) \quad (4)$$

$$= \mathbb{E}_{\tau \sim P(\cdot; \pi_{\theta})} \nabla_{\theta} \log P(\tau; \pi_{\theta}) R(\tau) \quad (5)$$

$$\approx \nabla_{\theta} \log P(\tau; \pi_{\theta}) R(\tau) \quad \tau \sim P(\cdot; \pi_{\theta}) \quad (6)$$

$$\nabla_{\theta} U(\theta) \approx \nabla_{\theta} \log P(\tau; \pi_{\theta}) R(\tau) \quad \tau \sim P(\cdot; \pi_{\theta}) \quad (7)$$

- ▶ Analogous to SGD (so variance reduction is important), but data distribution here is a moving target (so we may want a trust region).

Policy Gradient

We can subtract any constant from the reward

$$\nabla_{\theta} \sum_{\tau} P(\tau; \pi_{\theta})(R(\tau) - \textcolor{red}{b}) \quad (8)$$

$$= \nabla_{\theta} \sum_{\tau} P(\tau; \pi_{\theta})R(\tau) - \nabla_{\theta} \sum_{\tau} P(\tau; \pi_{\theta})b \quad (9)$$

$$= \nabla_{\theta} \sum_{\tau} P(\tau; \pi_{\theta})R(\tau) - \nabla_{\theta} b \quad (10)$$

$$= \nabla_{\theta} \sum_{\tau} P(\tau; \pi_{\theta})R(\tau) \quad (11)$$

Policy Gradient

The variance is magnified by $R(\tau)$

$$\nabla_{\theta} U(\theta) \approx \nabla_{\theta} \log P(\tau; \pi_{\theta}) R(\tau) \quad (12)$$

More fine-grained baseline?

$$\sum_{t=0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{t'=0} R(s_{t'}, a_{t'}) \quad (13)$$

$$\rightarrow \sum_{t=0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \underbrace{\sum_{t'=t} R(s_{t'}, a_{t'})}_{R_t} \quad (14)$$

$$= \sum_{t=0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R_t \quad (15)$$

$$\rightarrow \sum_{t=0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (R_t - b_t) \quad (16)$$

Policy Gradient

$$\sum_{t=0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (R_t - b_t) \quad (17)$$

Good choice of b_t ?

$$b_t^* = \operatorname{argmin}_b \mathbb{E} (R_t - b)^2 \quad (18)$$

$$= \mathbb{E} [R_t] \quad (19)$$

$$\approx V_{\phi}(s_t) \quad (20)$$

$$\sum_{t=0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \underbrace{(R_t - V_{\phi}(s_t))}_{A_t} \quad (21)$$

- Promote actions that lead to positive advantage.

Policy Gradient

$$\sum_{t=0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_t \quad (22)$$

- ▶ Sample a trajectory τ
- ▶ For each step t in τ
 - ▶ $R_t = \sum_{t'=t} r_{t'}$
 - ▶ $A_t = R_t - V_{\phi}(s_t)$
- ▶ take a gradient step $\nabla_{\phi} \sum_{t=0} \|V_{\phi}(s_t) - R_t\|^2$
- ▶ take a gradient step $\nabla_{\theta} \sum_{t=0} \log \pi_{\theta}(a_t | s_t) A_t$

In PG, our opt objective a moving target defined based on trajectory samples given the current policy

- ▶ each sample is only used once

An opt objective that reuses all historical data?

Policy gradient

$$\sum_{t=0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) A_t \quad (23)$$

Objective (on-policy)

$$\mathbb{E}_{\pi} [A(s, a)] \quad (24)$$

Objective (off-policy), via importance sampling

$$\mathbb{E}_{\pi_{old}} \left[\frac{\pi(a|s)}{\pi_{old}(a|s)} A_{old}(s, a) \right] \quad (25)$$

$$\max_{\pi} \mathbb{E}_{\pi_{old}} \left[\frac{\pi(a|s)}{\pi_{old}(a|s)} A_{old}(s, a) \right] \approx \sum_{n=1}^N \frac{\pi(a_n|s_n)}{\pi_{old}(a_n|s_n)} A_n \quad (26)$$

$$s.t. \text{ KL}(\pi_{old}, \pi) \leq \delta \quad (27)$$

- ▶ Quadratic approximation is used for the KL.
- ▶ Use conjugate gradient for the natural gradient direction $F^{-1}g$.

$$\max_{\pi} \sum_{n=1}^N \frac{\pi(a_n|s_n)}{\pi_{old}(a_n|s_n)} A_n - \beta \text{KL}(\pi_{old}, \pi) \quad (28)$$

Fixed β does not work well

- ▶ shrink β when $\text{KL}(\pi_{old}, \pi)$ is small
- ▶ increase β otherwise

Alternative ways to penalize large change in π ? Modify

$$\frac{\pi(a_n|s_n)}{\pi_{old}(a_n|s_n)} A_n = r_n(\theta) A_n \quad (29)$$

As

$$\min(r_n(\theta) A_n, \text{clip}(1 - \epsilon, 1 + \epsilon, r_n(\theta)) A_n) \quad (30)$$

- ▶ being pessimistic whenever “a large change in π leads to a better obj.”
- ▶ same as the original obj otherwise.

Advantage Estimation

Currently

$$\hat{A}_t = r_t + r_{t+1} + r_{t+2} + \dots - V(s_t) \quad (31)$$

More bias, less variance (ignoring long-term effect)

$$\hat{A}_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots - V(s_t) \quad (32)$$

Bootstrapping

$$\hat{A}_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots - V(s_t) \quad (33)$$

$$= r_t + \gamma(r_{t+1} + \gamma r_{t+2} + \dots) - V(s_t) \quad (34)$$

$$= r_t + \gamma V(s_{t+1}) - V(s_t) \quad (35)$$

We may unroll for more than one steps.

Advantage Actor-Critic

- ▶ Sample a trajectory τ

- ▶ For each step in τ

- ▶ $\hat{R}_t = r_t + \gamma V(s_{t+1})$

- ▶ $\hat{A}_t = \hat{R}_t - V(s_t)$

- ▶ take a gradient step using

$$\nabla_{\theta, \phi} \sum_{t=0} \left[-\log \pi_{\theta}(a_t | s_t) \hat{A}_t + \|V_{\phi}(s_t) - \hat{R}_t\|^2 \right]$$

May use TRPO, PPO for the policy part.

Variance reduction via multiple actors

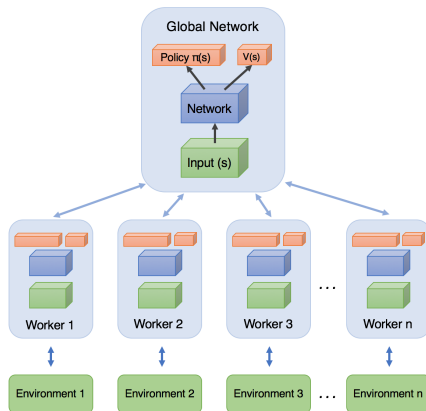


Figure : Asynchronous Advantage Actor-Critic