# Reinforcement Learning
# Reinforcement Learning: Policy Approach

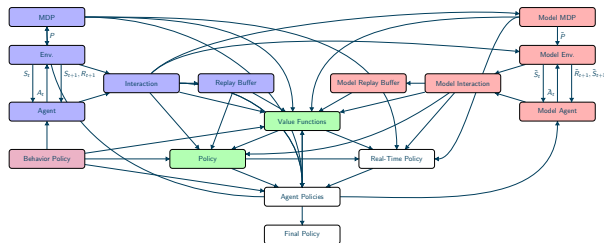Erwan Le Pennec
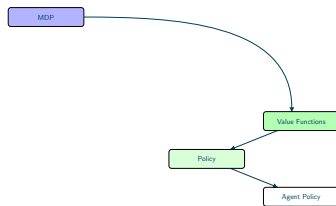Erwan.Le-Pennec@polytechnique.edu

ÉCOLE
**POLYTECHNIQUE**

M2DS - Reinforcement Learning – Fall 2024

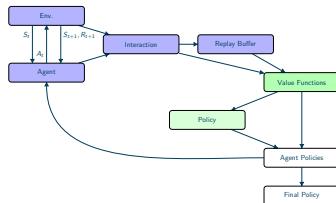# RL: What Are We Going To See?



## Outline

- Operations Research and MDP.
- Reinforcement learning and interactions.
- More tabular reinforcement learning.
- Reinforcement and approximation of value functions.
- Actor/Critic: a Policy Point of View

# Operations Research and MDP



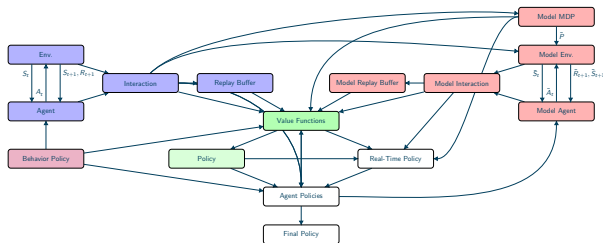## How to find the best policy knowing the MDP?

- Is there an optimal policy?
- How to estimate it numerically?

- Finite states/actions space assumption (tabular setting).
- Focus on interative methods using value functions (dynamic programming).
- Policy deduced by a statewise optimization of the value function over the actions.
- Focus on the discounted setting.

# Reinforcement Learning and Interactions



## How to find the best policy not knowing the MDP?

- How to interact with the environment to learn a good policy?
- Can we use a Monte Carlo strategy outside the episodic setting?
- How to update value functions after each interaction?

- Focus on stochastic methods using tabular value functions ($Q$ learning, SARSA...)
- Policy deduced by a statewise optimization of the value function over the actions.
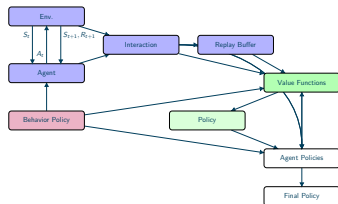
# More Tabular Reinforcement Learning
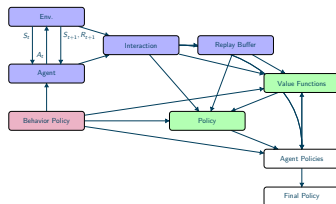


## Can We Do Better?

- Is there a gain to wait more than one step before updating?
- Can we interact with a different policy than the one we are estimating?
- Can we use an estimated model to plan?
- Can we plan in real-time instead of having to do it beforehand?

- Finite states/actions space setting (tabular setting).

# Reinforcement and Approximation of Value Functions



## How to Deal with a Large/Infinite states/action space?

- How to approximate value functions?
- How to estimate good approximation of value functions?

- Finite action space setting.
- Stochastic algorithm (Deep $Q$ Learning...).
- Policy deduced by a statewise optimization of the value function over the actions.

# Actor/Critic: a Policy Point of View



## Could We Directly Parameterized the Policy?

- How to parameterize a policy?
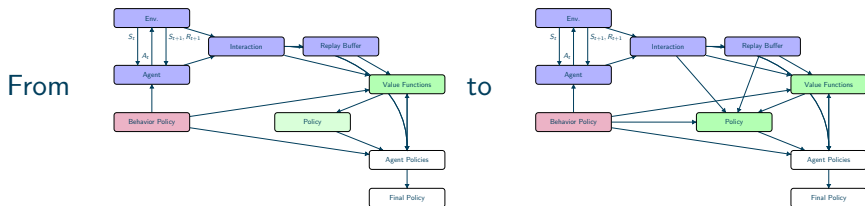- How to optimize this policy?
- Can we combine parametric policy and approximated value function?

- State Of The Art Algorithms (DPG,PPO, SAC...)

# Outline

# Policy Point of View

From  to 

## Policy Point of View

- Optimize policy directely instead of deriving it from a value function.
- Avoid the argmax operator.
- Most natural POV?

- Pontryagin vs Hamilton-Jacobi(-Bellman) in control!

# 1 Policy Gradient Theorems

2 Monte Carlo Based Policy Gradient

3 Actor / Critic Principle

4 3 SOTA Algorithms

5 References

$$J_\mu(\pi) = \sum_s \mu(s) v_\pi(s)$$

Goal: average expected return over the states

- Target used to define the linear programming formulation of an optimal policy in the tabular setting.
- $\mu$ can be the initial distribution of the states (independent of $\pi$)...
- but may also depends on $\pi$ (for instance the associated stationary measure)
- Other choices will appear.

- Goal: optimize $J_\mu(\pi)$ in $\pi$!

# Parametric Policy

$$\pi_\theta(a|s) = \begin{cases} \dfrac{e^{h_\theta(a,s)}}{\sum_{a'} e^{h_\theta(a,s')}} & \text{(softmax)} \\ P_{h_\theta(s)}(a) & \text{(parametric conditional model)} \\ \mathbf{1}_{a=h_\theta(s)} & \text{(deterministic)} \end{cases}$$

## Parametric Policy

- Restriction of the set of policy to a parametrized one.
- Most classical parametrizations:
    - Soft-max with a preference function $h_\theta(a, s)$,
    - Parametric conditional model with parameter $h_\theta(s)$
- To be useful need to be able to sample the distribution.
- $h_\theta$: from linear model to deep learning...
- Most of our result will assume that $\pi_\theta(a|s)$ is differentiable with respect to $\theta$.

- Deterministic policies will be considered with a different analysis.

# Episodic Setting: Gradient of Expected Returns

$$\nabla \mathcal{B} = \mathcal{B} \, \nabla \log \mathcal{B}$$

$$\nabla P = P \times \nabla \log P$$

$$v_{\pi_\theta}(s) = \mathbb{E}_{\pi_\theta}[G_0 | S_0 = s]$$

$$\nabla_\theta v_{\pi_\theta}(s) = \mathbb{E}_{\pi_\theta}\left[\left(\sum_{t=0}^{T_\tau - 1} \nabla \log \pi_\theta(A_t | S_t)\right) G_0 \middle| S_0 = s\right]$$

## Expected Returns

- Rely on $v_{\pi_\theta}(s) = \sum_\tau \mathbb{P}_{\pi_\theta}(\tau | S_0 = s) \, G_0(\tau)$ and

$$\nabla \mathbb{P}_{\pi_\theta}(\tau | S_0 = s) = \mathbb{P}_{\pi_\theta}(\tau | S_0 = s) \, \nabla \log \mathbb{P}_{\pi_\theta}(\tau | S_0 = s)$$
$$= \mathbb{P}_{\pi_\theta}(\tau | S_0 = s) \sum_t \left(\nabla \log \pi_\theta(A_t | S_t) + \nabla p(R_{t+1}, S_{t+1} | S_t, A_t)\right)$$
$$= \mathbb{P}_{\pi_\theta}(\tau | S_0 = s) \sum_t \nabla \log \pi_\theta(A_t | S_t)$$

- In an episodic setting, any trajectory $\tau$ ends at a finite time $T_\tau$.

$$J_{\mu_0}(\pi_\theta) = \sum_s \mathbb{P}(S_0 = s)\, v_{\pi_\theta}(s)$$

$$\nabla J_{\mu_0}(\pi_\theta) = \mathbb{E}_{\pi_\theta}\left[ \left( \sum_{t=0}^{T_\tau - 1} \nabla \log \pi_\theta(A_t | S_t) \right) G_0 \right]$$

*↗ ML*

### Policy Gradient Theorem

- Natural $\mu$: initial state distribution.
- Gradient is an expectation: MC type algorithm...

- Can be interpreted as the gradient of a the maximum likelihood of the actions weighted by the return.
- Favors good actions over bad ones.

$$J_{\mu_0}(\pi_\theta) = \sum_s \mathbb{P}(S_0 = s)\, v_{\pi_\theta}(s)$$

$$\nabla J_{\mu_0}(\pi_\theta) = \mathbb{E}_{\pi_\theta}\left[\left(\sum_{t=0}^{T_\tau - 1} \nabla \log \pi_\theta(A_t|S_t)\right)(G_0 - b)\right]$$

## Variance Reduction and Baseline

- The previous formulae are valid if one replace $G_0$ by any function of $\tau$.
- For any constant $b$, this leads to

$$\nabla \mathbb{E}_{\pi_\theta}[b] = 0 = \mathbb{E}_{\pi_\theta}\left[\left(\sum_{t=0}^{T_\tau - 1} \nabla \log \pi_\theta(A_t|S_t)\right)b\right]$$

- Optimal value for
$$b = \mathbb{E}_{\pi_\theta}\left[\left(\sum_{t=0}^{T_\tau-1}\nabla\log\pi_\theta(A_t|S_t)\right)^2 G_0\right] \Big/ \mathbb{E}_{\pi_\theta}\left[\left(\sum_{t=0}^{T_\tau-1}\nabla\log\pi_\theta(A_t|S_t)\right)^2\right]$$

- Most used value $b = \mathbb{E}_{\pi_\theta}[G_0]$.

Episodic

15

$$v_{\pi_\theta}(s) = \mathbb{E}_{\pi_\theta}\left[\sum \gamma^t R_t \Big| S_0 = s\right]$$

$$\nabla v_{\pi_\theta}(s) = \sum_t \gamma^t \mathbb{E}_{\pi_\theta}\left[\left(\sum_{t'=0}^{t-1} \nabla \log \pi_\theta(A_{t'}|S_{t'})\right) R_t \Big| S_0 = s\right]$$

$$= \sum_{t'} \mathbb{E}_{\pi_\theta}\left[\nabla \log \pi_\theta(A_{t'}|S_{t'}) \left(\sum_{t \geq t'} \gamma^t R_t\right) |S_0 = s\right]$$

$$= \sum_{t'} \gamma^{t'} \mathbb{E}_{\pi_\theta}[\nabla \log \pi_\theta(A_{t'}|S_{t'}) q_{\pi_\theta}(S_{t'}, A_{t'})|S_0 = s]$$

$$= \sum_{t'} \gamma^{t'} \mathbb{E}_{\pi_\theta}\left[\nabla \log \pi_\theta(A_{t'}|S_{t'}) \underbrace{(q_{\pi_\theta}(S_{t'}, A_{t'}) - v_{\pi_\theta}(S_{t'}))}_{a_{\pi_\theta}(S_{t'}, A_{t'})} |S_0 = s\right]$$

### From Returns to Value Functions

- Action point of view and use of value functions.

Episodic / Discounted

$$\nabla v_{\pi_\theta}(s) = \sum_{t'} \gamma^{t'} \mathbb{E}_{\pi_\theta}[\nabla \log \pi_\theta(A_{t'}|S_{t'}) q_{\pi_\theta}(S_{t'}, A_{t'})|S_0 = s]$$

$$= \sum_{t'} \gamma^{t'} \mathbb{E}_{\pi_\theta}[\nabla \log \pi_\theta(A_{t'}|S_{t'}) a_{\pi_\theta}(S_{t'}, A_{t'})|S_0 = s]$$

$$= \sum_{s'} \left( \sum_t \gamma^t \mathbb{P}_{\pi_\theta}(S_t = s'|S_0 = s) \right) \left( \sum_a \pi_\theta(a|s') \nabla \log \pi_\theta(a|s') q_{\pi_\theta}(s', a) \right)$$

$$= \sum_{s'} \left( \sum_t \gamma^t \mathbb{P}_{\pi_\theta}(S_t = s'|S_0 = s) \right) \left( \sum_a \pi_\theta(a|s') \nabla \log \pi_\theta(a|s') a_{\pi_\theta}(s', a) \right)$$

### Focus on states

- Even more stochastic gradients!

# Policy Gradient(s)

$$J_{\mu_0}(\pi_\theta) = \sum_s \mu_0(s) v_{\pi_\theta}(s)$$

$$\nabla J_{\mu_0}(\pi_\theta) = \sum_s \left( \sum_t \gamma^t \mathbb{P}_{\pi_\theta}(S_t = s) \right) \left( \sum_a \pi_\theta(a|s) \nabla \log \pi_\theta(a|s) q_{\pi_\theta}(s, a) \right)$$

$$= \sum_s \left( \sum_t \gamma^t \mathbb{P}_{\pi_\theta}(S_t = s) \right) \left( \sum_a \pi_\theta(a|s) \nabla \log \pi_\theta(a|s) (q_{\pi_\theta}(s, a) - v_{\pi_\theta}(s, a)) \right)$$

## Discounted Setting

- Average (discounted) return from the beginning.
- Focus on early steps in discounted setting...

Episodic / Discounted

$$J_{\mu_0}(\pi') - J_{\mu_0}(\pi) = \sum_s \sum_t \gamma^t \mathbb{P}_{\pi'}(S_t = s) \left( \sum_a \left( \pi'(a|s) - \pi(a|s) \right) q_\pi(s, a) \right)$$

$$= \sum_s \sum_t \gamma^t \mathbb{P}_{\pi'}(S_t = s) \left( \sum_a \left( \pi'(a|s) - \pi(a|s) \right) a_\pi(s, a) \right)$$

### Proof

- By construction, if $S_t$ is a trajectory using policy $\pi'$:

$$v_{\pi'}(S_t) - v_\pi(S_t) = \sum_a \left( \pi'(a|S_t) - \pi(a|S_t) \right) q_\pi(S_t, a) + \sum_a \pi'(a|s_t) \left( q_{\pi'}(S_t, a) - q_\pi(S_t, a) \right)$$

$$= \sum_a \left( \pi'(a|s_t) - \pi(a|S_t) \right) v_\pi(S_t, a) + \mathbb{E}_{\pi'}[v_{\pi'}(S_{t+1}) - v_\pi(S_{t+1})|S_t]$$

- Discounted setting shortcut

$$v_{\pi'} - v_\pi = r_{\pi'} + \gamma P^{\pi'} v_{\pi'} - r_\pi - \gamma P^\pi v_\pi = r_{\pi'} - r_\pi + \gamma \left( P^{\pi'} - P^\pi \right) v_\pi + \gamma P^{\pi'} \left( v_{\pi'} - v_\pi \right)$$

$$v_{\pi'} - v_\pi = (I - \gamma P^{\pi'})^{-1} \left( r_{\pi'} - r_\pi + \gamma \left( P^{\pi'} - P^\pi \right) v_\pi \right)$$

Episodic / Discounted

19

$$\left| J_{\mu_0}(\pi') - J_{\mu_0}(\pi) - \sum_s \sum_t \gamma^t \mathbb{P}_\pi(S_t = s) \left( \sum_a \left( \pi'(a|s) - \pi(a|s) \right) a_\pi(s, a) \right) \right|$$

$$= \left| \sum_s \sum_t \gamma^t \left( \mathbb{P}_{\pi'}(S_t = s) - \mathbb{P}_\pi(S_t = s) \right) \left( \sum_a \left( \pi'(a|s) - \pi(a|s) \right) a_\pi(s, a) \right) \right|$$

$$\leq \frac{2\gamma}{(1-\gamma)^2} \max_s \| \pi'(\cdot|s) - \pi(\cdot|s) \|_1^2 \max_{s,a} |a_\pi(s, a)|$$

---

### Approximate Policy Improvement Lemma

- If $\max_s \| \pi'(\cdot|s) - \pi(\cdot|s) \|_1 \leq \epsilon$

$$\mathbb{P}_{\pi'}(S_t = s) = (1 - \epsilon)^t \mathbb{P}_\pi(S_t = s) + (1 - (1 - \epsilon)^t) \mathbb{P}_{\text{mistake}}(S_t = s)$$

$$\rightarrow |\mathbb{P}_{\pi'}(S_t = s) - \mathbb{P}_\pi(S_t = s)| \leq 2(1 - (1 - \epsilon)^t) \leq 2\epsilon t$$

- $\sum_t 2\gamma^t t = \frac{2\gamma}{(1-\gamma)^2}$

Discounted

$$\left| J_{\mu_0}(\pi') - J_{\mu_0}(\pi) - \sum_s \sum_t \gamma^t \mathbb{P}_\pi(S_t = s) \left( \sum_a \left( \pi'(a|s) - \pi(a|s) \right) a_\pi(s,a) \right) \right|$$

$$\leq \frac{2\gamma}{(1-\gamma)^2} \max_s \|\pi'(\cdot|s) - \pi(\cdot|s)\|_1^2 \max_{s,a} |a_\pi(s,a)|$$

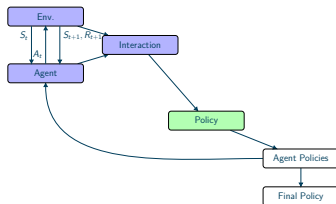## Approximate Policy Improvement Lemma and Policy Gradient Theorem

- Let $\pi' = \pi_{\theta+h}$ and $\pi_\theta$
  - $\pi_{\theta+h}(a|s) - \pi_\theta(a|s) = \pi_\theta(a|s)\langle \nabla \log \pi_\theta(a|s), h \rangle + O(\|h\|^2)$
  - $\|\pi_{\theta+h}(\cdot|s) - \pi_\theta(\cdot|s)\|_1 \leq \|h\| \max_a \|\nabla \log \pi_\theta(a|s)\| + O(\|h\|^2)$
- Implies Policy Gradient Theorem:

$$J_{\mu_0}(\pi_{\theta+h})$$

$$= J_{\mu_0}(\pi_\theta) + \sum_s \sum_t \gamma^t \mathbb{P}_{\pi_\theta}(S_t = s) \left( \sum_a \pi_\theta(a|s)\langle \nabla \log \pi_\theta(s,a), h \rangle a_\pi(s,a) \right) + O(\|h\|^2)$$

Discounted

# Outline

# Monte Carlo Approach

$$G_t = \sum_{t' \geq t} R_{t+1}$$

$$Q_{t,\pi_\theta}(s, a) = \mathbb{E}[G_t | S_t = s, A_t = a]$$

## Monte Carlo

- Replace every return by an empirical estimate along episodes.
- Need to wait until the end of the episods.

# REINFORCE: Monte Carlo Based Policy Gradient

$$J_{\mu_0}(\pi_\theta) = \sum_s \mathbb{P}(S_0 = s)\, v_{\pi_\theta}(s)$$

$$\nabla J_{\mu_0}(\pi_\theta) = \mathbb{E}_{\pi_\theta}\left[\left(\sum_{t=0}^{T_\tau - 1} \nabla \log \pi_\theta(A_t|S_t)\right) G_0\right]$$

$$= \sum_s \left(\sum_t \mathbb{P}_{\pi_\theta}(S_t = s)\right)\left(\sum_a \pi_\theta(a|s)\nabla \log \pi_\theta(a|s) q_{\pi_\theta}(s,a)\right)$$

$$\widehat{\nabla} J_{\mu_0}(\pi_\theta) = \left(\sum_{t=0}^{T_\tau - 1} \nabla \log \pi_\theta(A_t|S_t)\right) G_0 \quad \text{or} \quad \widehat{\nabla} J_{\mu_0}(\pi_\theta) = \sum_t \nabla \log \pi_\theta(A_t|S_t) G_t$$

### REINFORCE

- Plain MC (SGD) algorithm.
- Need to wait until the end of the episods.
- Convergence guarantees (even in off-line setting with importance sampling).

# REINFORCE with Baseline

$$\nabla J_{\mu_0}(\pi_\theta) = \mathbb{E}_{\pi_\theta}\left[\left(\sum_{t=0}^{T_\tau-1} \nabla \log \pi_\theta(A_t|S_t)\right)(G_0 - b)\right]$$

$$= \sum_s \left(\sum_t \mathbb{P}_{\pi_\theta}(S_t = s)\right)\left(\sum_a \pi_\theta(a|s)\nabla \log \pi_\theta(a|s)\left(q_{\pi_\theta}(s, a) - b(s)\right)\right)$$

$$\widehat{\nabla} J_{\mu_0}(\pi_\theta) = \left(\sum_{t=0}^{T_\tau-1} \nabla \log \pi_\theta(A_t|S_t)\right)(G_0 - b)$$

or $\quad \widehat{\nabla} J_{\mu_0}(\pi_\theta) = \sum_t \nabla \log \pi_\theta(A_t|S_t)\left(G_t - b(S_t)\right)$

## REINFORCE with baseline

- Several choices for $b$...
- and for $b(s)$ which can be any function (a crude estimate of $V_{t,\pi}(s)$ for instance)!
- Convergence guarantees (even in off-line setting with importance sampling).

Episodic

# Discounted REINFORCE?

$$\nabla J_{\mu_0}(\pi_\theta) = \mathbb{E}_{\pi_\theta}\left[\left(\sum_{t=0}^{T_\tau-1}\nabla\log\pi_\theta(A_t|S_t)\right)(G_0-b)\right]$$

$$= \sum_s\left(\sum_t\gamma^t\mathbb{P}_{\pi_\theta}(S_t=s)\right)\left(\sum_a\pi_\theta(a|s)\nabla\log\pi_\theta(a|s)\left(q_{\pi_\theta}(s,a)-b(s)\right)\right)$$

$$\widehat{\nabla}J_{\mu_0}(\pi_\theta) = \left(\sum_{t=0}^{T_\tau-1}\nabla\log\pi_\theta(A_t|S_t)\right)(G_0-b)$$

$$\text{or}\quad \widehat{\nabla}J_{\mu_0}(\pi_\theta) = \sum_t\gamma^t\nabla\log\pi_\theta(A_t|S_t)\left(G_t-b(S_t)\right)$$
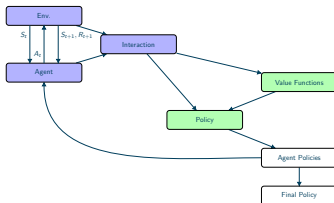
---

### Discounted REINFORCE

- Can be defined. . .
- but still requires an episodic setting for the discounted return $G_t$ to be computed.

# Discounted Measure?

$$\widehat{\nabla} J_{\mu_0}(\pi_\theta) = \sum_t \gamma^t \nabla \log \pi_\theta(A_t|S_t) \left( G_t - b(S_t) \right)$$

$$\longrightarrow \widehat{\nabla} J_{\mu_{\pi_\theta}}(\pi_\theta) = \frac{1}{1 - \gamma} \nabla \log \pi_\theta(A_t|S_t) \left( G_t - b(S_t) \right)?$$

### Discounted Measure?

- Much less weights for later states if $\mu$ corresponds to the initial state distribution!
- Equal weights corresponds to an averaged probability independent $t$, which is well defined if the initial distribution is the stationary distribution $\mu_{\pi_\theta}$ corresponding to $\pi_\theta$ (it it exists).
- Approximately true after a burning stage if we reach stationarity. . .
- Better handled by the average return!

- More on this later. . .

# Actor/Critic

## Actor/Critic

- Actor: Parametric policy $\pi_\theta$ used.
- Critic: $Q$-value function $Q_{\boldsymbol{w}}(\cdot, \cdot)$ approximating $Q_{\pi_\theta}$.
- Critic follows the Actor, which is optimized using the Critic.

- In Value Approximation, the Actor follows the Critic (through the argmax operator).
- In on-line methods, the Actor is used to interact with the environment.

# Actor/Critic

$$J_{(\underline{\mu_0})}(\pi_\theta) = \sum_s \mu_0(s) v_{\pi_\theta}(s)$$

$$\nabla J_{\mu_0}(\pi_\theta) = \sum_s \left( \sum_t \gamma^t \mathbb{P}_{\pi_\theta}(S_t = s) \right) \left( \sum_a \pi_\theta(a|s) \nabla \log \pi_\theta(a|s) (q_{\pi_\theta}(s, a) - v_{\pi_\theta}(s, a)) \right)$$

$$\widehat{\nabla} J_{\mu_0}(\pi_\theta) = \sum_t \gamma^t \pi_\theta(A_t|S_t) \nabla \log \pi_\theta(A_t|S_t) \left( q_{\pi_\theta}(S_t, A_t) - \sum_a \pi_\theta(a|S_t) q_{\pi_\theta}(S_t, A_t) \right)$$

$$\simeq \sum_t \gamma^t \pi_\theta(A_t|S_t) \nabla \log \pi_\theta(A_t|S_t) \left( Q_w(S_t, A_t) - \sum_a \pi_\theta(a|S_t) Q_w(S_t, A_t) \right)$$

## Actor/Critic

- Critic update: Stochastic Policy Gradient with plugin.
- Actor update: any $Q$-value methods estimating $q_{\pi_\theta}$.
- Requires a two-scales algorithm so that $Q_w$ is always a good estimate of $q_{\pi_\theta}$.

- Is this a real algorithm in a non-episodic setting?

# Actor/Critic

$$J_{\mu_{\pi_\theta}}(\pi_\theta) = \sum_s \mu_{\pi_\theta}(s) v_{\pi_\theta}(s)$$

$$\nabla J_{\mu_{\pi_\theta}}(\pi_\theta) = \sum_s \frac{1}{1-\gamma} \mathbb{P}_{\pi_\theta}(S_t = s) \left( \sum_a \pi_\theta(a|s) \nabla \log \pi_\theta(a|s)(q_{\pi_\theta}(s,a) - v_{\pi_\theta}(s,a)) \right)$$

$$\widehat{\nabla} J_{\mu_{\pi_\theta}}(\pi_\theta) \simeq \frac{1}{1-\gamma} \pi_\theta(A_t|S_t) \nabla \log \pi_\theta(A_t|S_t) \left( Q_{\mathbf{w}}(S_t, A_t) - \sum_a \pi(a|S_t) Q_{\mathbf{w}}(S_t, A_t) \right)$$

## Actor/Critic

- Critic update: Stochastic Policy Gradient with plugin.
- Actor update: any $Q$-value methods estimating $q_{\pi_\theta}$.
- Requires a two-scales algorithm so that $Q_{\mathbf{w}}$ is always a good estimate of $q_{\pi_\theta}$.

- Require the existence of a stationary measure... and that this stationary measure is reached *quickly*.
- Much harder to do off-policy algorithm as the stationary measure is not known!

Discounted

$$Q_{\boldsymbol{w}} \simeq q_{\pi_\theta}$$

### Critic

- On-line TD learning with interaction following $\pi_\theta$.
- Off-Policy TD learning is possible if the policy used for any action is stored.
- Approximate off-policy TD learning is possible using a replay buffer providing $\pi_\theta$ is changing slowly.

- May lead to 3 scales algorithm (Actor/Critic Target/Critic)
- As mentionned in the previous slide, much harder to do off-line update for the actor.

Discounted

$$J'_\mu(\pi) = \sum_s \mu(s) v_\pi(s)$$

### Off-Line Actor

- Idea proposed in 2012.
- Key lemma in the paper
  $$\nabla J'_\mu(\pi_\theta) \simeq \sum_s \mu(s) \sum_a \pi_\theta(a|s) \nabla \pi_\theta(a|s) q_{\pi_\theta}(s, a)$$
  turns out to be wrong!
- Still used as a heuristic justification of many algorithms!
- Explicit formula for $\nabla J'_\mu(\pi_\theta)$ can be obtained but much harder to use. . .

$$J_{\mu_0}(\pi') \geq J_{\mu_0}(\pi) + \sum_t \gamma^t \mathbb{P}_\pi(S_t = s) \left( \sum_a \left( \pi'(s|a) - \pi(s|a) \right) a_\pi(s,a) \right)$$
$$- \frac{2\gamma}{(1-\gamma)^2} \max_s \|\pi'(\cdot|s) - \pi(\cdot|s)\|_1^2 \max_{s,a} |a_\pi(s,a)|$$

### Ideal Minorize-Majorization Algorithm

- At step $k$, find $\theta_{k+1}$ maximizing

$$J_{\mu_0}(\pi_\theta | \pi_{\theta_k}) = \sum_s \sum_t \gamma^t \mathbb{P}_{\pi_{\theta_k}}(S_t = s) \left( \sum_a \left( \pi_\theta(s|a) - \pi_{\theta_k}(s|a) \right) a_{\pi_{\theta_k}}(s,a) \right)$$
$$- \frac{2\gamma}{(1-\gamma)^2} \max_s \|\pi_\theta(\cdot|s) - \pi_{\theta_k}(\cdot|s)\|_1^2 \max_{s,a} |a_{\pi_{\theta_k}}(s,a)|$$

- By construction, $J_{\mu_0}(\pi_{\theta_{k+1}}) \geq J_{\mu_0}(\pi_{\theta_k})$

- Sample efficient algorithm as the same trajectory can be (re)used in the optimization.

Discounted

$$J_{\mu_0}(\pi_\theta) \geq J_{\mu_0}(\pi_{\theta_k}) + \sum_s \sum_t \gamma^t \mathbb{P}_{\pi_{\theta_k}}(S_t = s) \left( \sum_a \left( \pi_\theta(s|a) - \pi_{\theta_k}(s|a) \right) a_{\pi_{\theta_k}}(s, a) \right)$$

$$- \frac{2\gamma}{(1-\gamma)^2} \max_s \|\pi_\theta(\cdot|s) - \pi_{\theta_k}(\cdot|s)\|_1^2 \max_{s,a} |a_{\pi_{\theta_k}}(s, a)|$$

### Optimization

- Gradient descent is possible.
- Gradient of the first term can be approximated using a critic by

$$\sum_s \sum_t \gamma^t \mathbb{P}_\pi(S_t = s) \left( \sum_a \pi_\theta \nabla \pi_\theta(s|a) A_{\pi_{\theta_k}}(s, a) \right)$$

- Gradient of the second term more involved.

- Simpler (TRPO like) strategy: optimize

$$\sum_s \sum_t \gamma^t \mathbb{P}_{\pi_{\theta_k}}(S_t = s) \left( \sum_a \left( \pi_\theta(s|a) - \pi_{\theta_k}(s|a) \right) a_{\pi_{\theta_k}}(s, a) \right)$$

under $\max_s \|\pi_\theta(\cdot|s) - \pi_{\theta_k}(\cdot|s)\|_1^2 \leq \epsilon$ and reduce $\epsilon$ there is no gain.

Discounted

$$J_{\mu_0}(\pi_\theta) \geq J_{\mu_0}(\pi_{\theta_k}) + \sum_s \sum_t \gamma^t \mathbb{P}_{\pi_{\theta_k}}(S_t = s) \left( \sum_a \left( \pi_\theta(s|a) - \pi_{\theta_k}(s|a) \right) a_{\pi_{\theta_k}}(s, a) \right)$$

$$- \frac{2\gamma R_{\max}}{(1 - \gamma)^2} \max_s \mathsf{KL}(\pi_{\theta_k}(\cdot|s), \pi_\theta(\cdot|s))$$

### TRPO/PPO Optimization

- Replace the $\ell_1$ norm by a KL divergence.
- In practice, replace the max by an average and replace $\frac{2\gamma R_{\max}}{(1-\gamma)^3}$ by parameter $\beta$ and replace the $a_{\pi_k}$ by an estimate $A_{\pi_k}$.
- PPO: Gradient descent of the relaxed goal.
- TRPO: Constrained optimization.

- Adaptive scheme to set $\beta$.
- Can be used with continuous action.

Discounted

$$\sum_s \sum_t \gamma^t \mathbb{P}_{\pi_{\theta_k}}(S_t = s) \left( \sum_a \pi_{\theta_k}(s|a) \min \left( \frac{\pi_\theta(s|a)}{\pi_{\theta_k}(s,a)} a_{\pi_{\theta_k}}(s,a), \text{clip}(1-\epsilon, \frac{\pi_\theta(s|a)}{\pi_{\theta_k}(s,a)}, 1+\epsilon) a_{\pi_{\theta_k}}(s,a) \right) \right)$$

## Clipped Objective

- Insight by (re)substracting $\sum_a \pi_{\theta_k}(s|a) a_{\theta_k}(s,a) = 0$:

$$\sum_a \min \left( (\pi_\theta(s|a) - \pi_{\theta_k}(s,a)) \, a_{\pi_{\theta_k}}(s,a), \text{clip}(-\epsilon, \pi_\theta(s|a) - \pi_{\theta_k}(s,a), \epsilon) a_{\pi_{\theta_k}}(s,a) \right)$$

$$= \sum_a \text{clip}(-\epsilon \pi_{\theta_k}(s,a), \pi_\theta(s|a) - \pi_{\theta_k}(s,a), \epsilon \pi_{\theta_k}(s,a)) a_{\pi_{\theta_k}}(s,a)$$

$$- \max \left( 0, -(\pi_\theta(s|a) - \pi_{\theta_k}(s,a)) a_{\pi_{\theta_k}}(s,a) - \epsilon \pi_{\theta_k}(s,a)|a_{\pi_{\theta_k}}(s,a)| \right)$$

- First term amount to replace $\pi_\theta$ by a policy

$$\tilde{\pi}_\theta(a|s) = \text{clip}(\pi_{\theta_k}(a|s)(1-\epsilon), \pi_\theta(a|s), \pi_{\theta_k}(a|s)(1+\epsilon)) + \eta_s \pi_{\theta_k}(a|s)$$

where $\eta$ is so that $\tilde{\pi}$ is a probability for all $s$ and $\|\tilde{\pi}_\theta(\cdot, s) - \pi_{\theta_k}(\cdot, s)\|_1 \le \epsilon$

- Second term: hinge loss type penalization of policy $\pi_\theta$ penalizing *bad* actions.

- Very efficient for discrete actions.                                      38

Discounted

$$\sum_{s,t} \mathbb{P}_{\pi_{\theta_k}}(S_t = s) \left( \sum_a \left( \pi_\theta(s|a) - \pi_{\theta_k}(s|a) \right) a_{\pi_{\theta_k}}(s, a) \right) - \beta \max_s \mathsf{KL}(\pi_{\theta_k}(\cdot|s), \pi_\theta(\cdot|s))$$

$$\sum_{s,t} \mathbb{P}_{\pi_{\theta_k}}(S_t = s) \left( \sum_a \pi_{\theta_k}(s|a) \min \left( \frac{\pi_\theta(s|a)}{\pi_{\theta_k}(s, a)} a_{\pi_{\theta_k}}(s, a), \mathsf{clip}(1 - \epsilon, \frac{\pi_\theta(s|a)}{\pi_{\theta_k}(s, a)}, 1 + \epsilon) a_{\pi_{\theta_k}}(s, a) \right) \right)$$

### Stationary Objective

- Amount to replace $J_{\mu_0}(\pi)$ by $J_{\mu_\pi}(\pi)$
- Most common implementation of PPO. . .
- But no way to justify it mathematically!
- May explain the (possible) instabilities of PPO.

Discounted

$$J_{\mu_0}(\pi_\theta) = \sum_s \mu_0(s) v_{\pi_\theta}(s) \quad \text{with deterministic policy } \pi_\theta(a|s) = \mathbf{1}_{a = h_\theta(s)}$$

$$\nabla J_{\mu_0}(\pi_\theta) = \sum_s \sum_t \gamma^t \mathbb{P}_{\pi_\theta}(S_t = s) \, \nabla_a q(S_t, h_\theta(S_t)) \nabla h_\theta(S_t)$$

### Deterministic Policy Gradient

- Deterministic policy replaced by a randomized one centered on $h_{\theta(s)}$ in the interactions!
- Critic trained with a TD variant of DQN.
- Same formula by using a policy $\pi_\theta = N(h_\theta(s), \sigma^2 \mathrm{Id})$ and letting $\sigma$ goes to 0.
- Off-Policy as claimed?
- Yes for the actor but no theoretical justification for the critic!
- In practice, the buffer contains only samples using a policy close to the current one. . .

Discounted

40

$$R_t \to R_t + \lambda \mathcal{H}(\pi(S_t))$$

## A Modified Reward

- Modification of the reward to favor high entropy policy:
$$R_t \to R_t + \lambda \mathcal{H}(\pi(S_t))$$

- Goal:
$$J(\pi) = \mathbb{E}_\pi \left[ \sum_t \gamma^t \left( R_t + \lambda \mathcal{H}(\pi(S_t)) \right) \right]$$

- Soft value function implicitly defined as the fixed point of
$$\mathcal{T}^\pi q_\pi(s, a) = r_\pi(s, a) + \gamma \sum_{s'} p(s'|s, a) v_\pi(s')$$

$$\text{where} \quad v_\pi(s, a) = \sum_a \pi(a|s) \left( q_\pi(s, a) - \log \pi(a|s) \right)$$

Discounted

# SAC: Policy Improvement and Optimal Policy

$$R_t \to R_t + \lambda \mathcal{H}(\pi(S_t))$$

## A Modified Policy Improvement Lemma

- Policy improvement rule:
$$\pi^+(\cdot|s) = \underset{\pi(\cdot|s)}{\mathrm{argmax}} \sum_a \pi(a|s)\left(q(s,a) - \lambda \log(\pi(a|s))\right)$$

$$\pi^+(a|s) \propto \exp(-\frac{1}{\lambda} q(s,a))$$

implies $G_{\pi^+}(s,a) \geq G_\pi(s,a)$.

- At convergence, $J(\pi^\star)$ is optimal!
- Convergence in the finite setting.

Discounted

# SAC: Parametrization

$$\pi \sim \pi_\theta \quad \text{and} \quad q(s,a) \sim Q_w$$

## SAC Choices

- Fitted TD learning for $Q$:
$$w \simeq \text{argmin} \sum_{(S,A,R,S') \in \mathcal{B}} \left( R + \mathbb{E}_{\pi_\theta}\left[ \gamma Q_{\overline{w}}(S',a) - \lambda \log \pi_\theta(a|S') \right] - Q_w(S,A) \right)^2$$

  where the trajectory pieces are samples from a replay buffer and $\overline{w}$ is a slowdown version of $w$ (two-scales algorithm).

- Online version rather than batch...

- Fitted KL for $\pi$:
$$\theta \simeq \text{argmin} \sum_{(S,A,R,S') \in \mathcal{B}} \text{KL}(\pi_\theta(\cdot|S)|\exp - \lambda Q_{[}\overline{w}](S,\dot) / Z_{\overline{w}}(S))$$
$$\simeq \sum_{(S,A,R,S') \in \mathcal{B}} \mathbb{E}_{\pi_\theta}\left[ \frac{1}{\lambda} \log \pi_\theta(a|S) - Q_\theta(a|s) \right]$$

Discounted

1. Policy Gradient Theorems

2. Monte Carlo Based Policy Gradient

3. Actor / Critic Principle

4. 3 SOTA Algorithms

5. References

# References

R. Sutton and A. Barto.
*Reinforcement Learning, an Introduction (2nd ed.)*
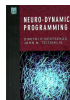MIT Press, 2018

O. Sigaud and O. Buffet.
*Markov Decision Processes in Artificial Intelligence.*
Wiley, 2010

M. Puterman.
*Markov Decision Processes. Discrete Stochastic Dynamic Programming.*
Wiley, 2005

D. Bertsekas and J. Tsitsiklis.
*Neuro-Dynamic Programming.*
Athena Scientific, 1996

W. Powell.
*Reinforcement Learning and Stochastic Optimization: A Unified Framework for Sequential Decisions.*
Wiley, 2022

S. Meyn.
*Control Systems and Reinforcement Learning.*
Cambridge University Press, 2022

V. Borkar.
*Stochastic Approximation: A Dynamical Systems Viewpoint.*
Springer, 2008

T. Lattimore and Cs. Szepesvári.
*Bandit Algorithms.*
Cambridge University Press, 2020

# Licence and Contributors

## Contributors

- Main contributor: E. Le Pennec
- Contributors: S. Boucheron, A. Dieuleveut, A.K. Fermin, S. Gadat, S. Gaiffas, A. Guilloux, Ch. Keribin, E. Matzner, M. Sangnier, E. Scornet.