

# Policy Gradient Methods

## Reinforcement Learning Seminar

Yingru Li

The Chinese University of Hong Kong, Shenzhen

February 11, 2019

# Table of Contents

- 1 Introduction
  - Policy based Reinforcement Learning
  - Policy Search
  - Finite Difference Policy Gradient
- 2 Monte-Carlo Policy Gradient
  - Likelihood Ratios
  - Policy Gradient Theorem
- 3 Actor-Critic Policy Gradient
  - Compatible Function Approximation
  - Advantage Function Critic
  - Eligibility Traces
  - Natural Policy Gradient

# Table of Contents

- 1 Introduction
  - Policy based Reinforcement Learning
  - Policy Search
  - Finite Difference Policy Gradient
- 2 Monte-Carlo Policy Gradient
  - Likelihood Ratios
  - Policy Gradient Theorem
- 3 Actor-Critic Policy Gradient
  - Compatible Function Approximation
  - Advantage Function Critic
  - Eligibility Traces
  - Natural Policy Gradient

# Policy-Based Reinforcement Learning

- In the last lecture we approximated the Value or Action-value function (Q-factor) using parameters  $\theta$ ,

$$V_{\theta}(s) \approx V^{\pi}(s)$$
$$Q_{\theta}(s, a) \approx Q^{\pi}(s, a)$$

- A policy was generated directly from the value function (e.g. using  $\epsilon$ -greedy)
- In this lecture we will directly parameterize the policy

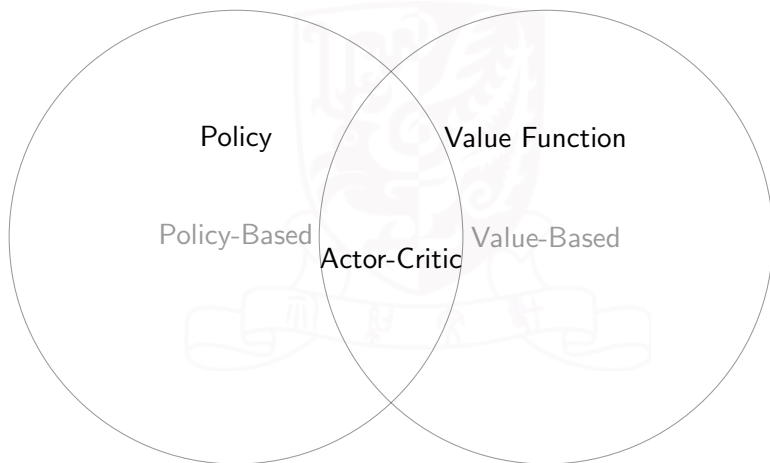
$$\pi_{\theta}(s, a) = \mathbb{P}[a \mid s, \theta]$$

- We will focus again on model-free reinforcement learning

# Value-Based and Policy-Based RL

- Value Based
  - Learnt Value Function
  - Implicit policy ( $\epsilon$ -greedy)
- Policy Based
  - No Value Function
  - Learnt Policy
- Actor-Critic
  - Learnt Value Function
  - Learnt Policy

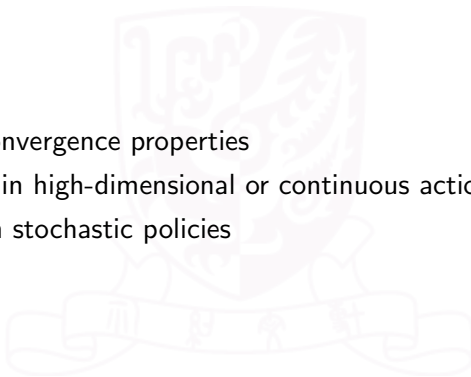
# Value-Based and Policy-Based RL



# Advantages of Policy-Based RL

## Advantages:

- Better convergence properties
- Effective in high-dimensional or continuous action spaces
- Can learn stochastic policies



# Advantages of Policy-Based RL

## Advantages:

- Better convergence properties
- Effective in high-dimensional or continuous action spaces
- Can learn stochastic policies

## Disadvantages:

- Typically converge to a local rather than global optimum
- Evaluating a policy is typically inefficient and high variance



# Example: Rock-Paper-Scissors

- Two-player game of rock-paper-scissors
  - Scissors beats paper
  - Rock beats scissors
  - Paper beats rock
- Consider policies for *iterated* rock-paper-scissors
  - A deterministic policy is easily exploited
  - A uniform random policy is optimal (i.e. Nash equilibrium)

# Policy Objective Functions

- Goal: given policy  $\pi_\theta(s, a)$  with parameters  $\theta$ , find best  $\theta$
- But how do we measure the quality of a policy  $\pi_\theta$ ?
- In episodic environments we can use the start value

$$J_1(\theta) = V^{\pi_\theta}(s_1) = \mathbb{E}_{\pi_\theta}[v_1]$$

- In continuing environments we can use the average value

$$J_{avV}(\theta) = \sum_s d^{\pi_\theta}(s) V^{\pi_\theta}(s)$$

- Or the average reward per time-step

$$J_{avR}(\theta) = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(s, a) \mathcal{R}_s^a$$

- where  $d^{\pi_\theta}(s)$  is stationary distribution of Markov chain for  $\pi_\theta$

# Policy Optimization

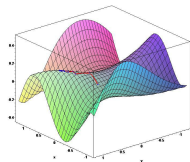
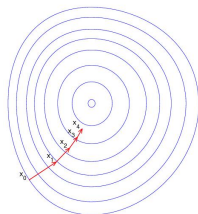
- Policy based reinforcement learning is an optimization problem
- Find  $\theta$  that maximises  $J(\theta)$
- Some approaches do not use gradient (gradient-free)
  - Hill climbing
  - Simplex / amoeba / Nelder Mead
  - Genetic algorithms
- Greater efficiency often possible using gradient
  - Gradient descent
  - Conjugate gradient
  - Quasi-newton
- We focus on gradient descent, many extensions possible
- And on methods that exploit sequential structure

# Policy Gradient

- Let  $J(\theta)$  be any policy objective function
- Policy gradient algorithms search for a local maximum in  $J(\theta)$  by ascending the gradient of the policy, w. r.t. parameters  $\theta$
- Where  $\nabla_{\theta} J(\theta)$  is the policy gradient

$$\nabla_{\theta} J(\theta) = \begin{pmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{pmatrix}$$

- and  $\alpha$  is a step-size parameter



# Computing Gradients By Finite Differences

- To evaluate policy gradient of  $\pi_{\theta}(s, a)$
- For each dimension  $k \in [1, n]$ 
  - Estimate  $k$ -th partial derivative of objective function w.r.t.  $\theta$
  - By perturbing  $\theta$  by small amount  $\epsilon$  in  $k$ -th dimension

$$\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \epsilon u_k) - J(\theta)}{\epsilon}$$

where  $u_k$  is unit vector with 1 in  $k$ th component, 0 elsewhere

- Uses  $n$  evaluations to compute policy gradient in  $n$  dimensions
- Simple, noisy, inefficient - but sometimes effective
- Works for arbitrary policies, even if policy is not differentiable

# Table of Contents

- 1 Introduction
  - Policy based Reinforcement Learning
  - Policy Search
  - Finite Difference Policy Gradient
- 2 Monte-Carlo Policy Gradient
  - Likelihood Ratios
  - Policy Gradient Theorem
- 3 Actor-Critic Policy Gradient
  - Compatible Function Approximation
  - Advantage Function Critic
  - Eligibility Traces
  - Natural Policy Gradient

# Score Function

- We now compute the policy gradient analytically
- Assume policy  $\pi_\theta$  is differentiable whenever it is non-zero
- and we know the gradient  $\nabla_\theta \pi_\theta(s, a)$
- **Likelihood ratios** exploit the following identity

$$\begin{aligned}\nabla_\theta \pi_\theta(s, a) &= \pi_\theta(s, a) \frac{\nabla_\theta \pi_\theta(s, a)}{\pi_\theta(s, a)} \\ &= \pi_\theta(s, a) \nabla_\theta \log \pi_\theta(s, a)\end{aligned}$$

- The **score function** is  $\nabla_\theta \log \pi_\theta(s, a)$

# Softmax Policy

- We will use a softmax policy as a running example
- Weight actions using linear combination of features  $\phi(s, a)^T \theta$
- Probability of action is proportional to exponential weight

$$\pi_{\theta}(s, a) \propto e^{\phi(s, a)^T \theta}$$

- The score function is

$$\nabla_{\theta} \log \pi_{\theta}(s, a) = \phi(s, a) - \mathbb{E}_{\pi_{\theta}}[\phi(s, \cdot)]$$



# Gaussian Policy

- In continuous action spaces, a Gaussian policy is natural
- Mean is a linear combination of state features  $\mu(s) = \phi(s)^T \theta$
- Variance may be fixed  $\sigma^2$ , or can also be parameterized
- Policy is Gaussian,  $a \sim \mathcal{N}(\mu(s), \sigma^2)$
- The score function is

$$\nabla_{\theta} \log \pi_{\theta}(s, a) = \frac{(a - \mu(s))\phi(s)}{\sigma^2}$$

# One-Step MDPs

- Consider a simple class of one-step MDPs
- Starting in state  $s \sim d(s)$
- Terminating after one time-step with reward  $r = \mathcal{R}_{s,a}$
- Use likelihood ratios to compute the policy gradient

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\pi_\theta}[r] \\ &= \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} \pi_\theta(s, a) \mathcal{R}_{s,a} \\ \nabla_\theta J(\theta) &= \sum_{s \in \mathcal{S}} d(s) \sum_{a \in \mathcal{A}} \pi_\theta(s, a) \nabla_\theta \log \pi_\theta(s, a) \mathcal{R}_{s,a} \\ &= \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) r] \end{aligned}$$

# Policy Gradient Theorem

- The policy gradient theorem generalises the likelihood ratio approach to multi-step MDPs
- Replaces instantaneous reward  $r$  with long-term value  $Q^\pi(s, a)$

## Policy Gradient Theorem

For any differentiable policy  $\pi_\theta(s, a)$ , for any of the policy objective functions  $J = J_1, J_{avR}$ , or  $\frac{1}{1-\gamma} J_{av}$ , the policy gradient is

$$\begin{aligned}\nabla_\theta J(\theta) &= \sum_{s \in \mathcal{S}} d^{\pi_\theta}(s) \sum_a \nabla_\theta \pi_\theta(s, a) Q^{\pi_\theta}(s, a) \\ &= \mathbb{E}_{\pi_\theta} \nabla_\theta \log \pi_\theta(s, a) Q^{\pi_\theta}(s, a)\end{aligned}$$

# Monte-Carlo Policy Gradient (REINFORCE)

- Update parameters by stochastic gradient ascent
- Using return  $v_t$  as an unbiased sample of  $Q^{\pi_\theta}(s_t, a_t)$

$$\Delta\theta_t = \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) v_t$$

---

## Algorithm 1 REINFORCE

---

- 1: Init  $\theta$  arbitrarily
  - 2: **for** each episode  $\{s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_\theta$  **do**
  - 3:   **for**  $t = 1$  to  $T - 1$  **do**
  - 4:      $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) v_t$
  - 5:   **end for**
  - 6: **end for**
  - 7: **return**  $\theta$
-

# Table of Contents

- 1 Introduction
  - Policy based Reinforcement Learning
  - Policy Search
  - Finite Difference Policy Gradient
- 2 Monte-Carlo Policy Gradient
  - Likelihood Ratios
  - Policy Gradient Theorem
- 3 Actor-Critic Policy Gradient
  - Compatible Function Approximation
  - Advantage Function Critic
  - Eligibility Traces
  - Natural Policy Gradient

## Reducing Variance Using a Critic

- Monte-Carlo policy gradient still has high variance
- We use a critic to estimate the action-value function,

$$Q_w(s, a) \approx Q^{\pi_\theta}(s, a)$$

- Actor-critic algorithms maintain *two* sets of parameters
  - Critic** Updates action-value function parameters  $w$
  - Actor** Updates policy parameters  $\theta$ , in direction suggested by critic
- Actor-critic algorithms follow an *approximate* policy gradient

$$\begin{aligned}\nabla_\theta J(\theta) &\approx \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)] \\ \Delta\theta &= \alpha \nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)\end{aligned}$$

# Estimating the Action-Value Function

- The critic is solving a familiar problem: policy evaluation
- How good is policy  $\pi_\theta$  for current parameters  $\theta$ ?
- This problem was explored in previous two lectures, e.g.
  - Monte-Carlo policy evaluation
  - Temporal-Difference learning
  - $TD(\lambda)$
- Could also use e.g. least-squares policy evaluation

# Action-Value Actor-Critic

Using linear value fn approx.  $Q_w(s, a) = \phi(s, a)^T w$

**Critic** Updates  $w$  by linear  $TD(0)$

**Actor** Updates  $\theta$  by policy gradient

---

## Algorithm 2 QAC

---

- 1: Init  $s, \theta$ , Sample  $a \sim \pi_\theta$
  - 2: **for** each step **do**
  - 3:   Sample reward  $r = \mathcal{R}_s^a$ ,  $s' \sim \mathcal{P}_s^a$ , and  $a' \sim \pi_\theta(s', a')$
  - 4:    $\theta = \theta + \alpha \nabla_\theta \log \pi(s, a) Q_w(s, a)$
  - 5:    $w \leftarrow w + \beta (r + \gamma Q_w(s', a') - Q_w(s, a)) \phi(s, a)$
  - 6:    $a \leftarrow a', s \leftarrow s'$
  - 7: **end for**
-



# Bias in Actor-Critic Algorithms

- Approximating the policy gradient introduces bias
- A biased policy gradient may not find the right solution
- Luckily, if we choose value function approximation carefully
- Then we can avoid introducing any bias
- i.e. We can still follow the exact policy gradient

# Compatible Function Approximation

## Compatible Function Approximation Theorem

If the following two conditions are satisfied:

- Value function approximator is **compatible** to the policy

$$\nabla_w Q_w(s, a) = \nabla_\theta \log \pi_\theta(s, a)$$

- Value function parameters  $w$  minimize the mean-squared error

$$\epsilon = \mathbb{E}_{\pi_\theta} [(Q^{\pi_\theta}(s, a) - Q_w(s, a))^2]$$

Then the policy gradient is exact,

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$$

# Proof of Compatible Function Approximation Theorem

If  $w$  is chosen to minimise mean-squared error, gradient of  $\epsilon$  w r.t.  $w$  must be zero,

$$\nabla_w \epsilon = 0$$

$$\mathbb{E}_{\pi_\theta} [(Q^\theta(s, a) - Q_w(s, a)) \nabla_w Q_w(s, a)] = 0$$

$$\mathbb{E}_{\pi_\theta} [(Q^\theta(s, a) - Q_w(s, a)) \nabla_\theta \log \pi_\theta(s, a)] = 0$$

$$\mathbb{E}_{\pi_\theta} [Q^\theta(s, a) \nabla_\theta \log \pi_\theta(s, a)] = \mathbb{E}_{\pi_\theta} [Q_w(s, a) \nabla_\theta \log \pi_\theta(s, a)]$$

So  $Q_w(s, a)$  can be substituted directly into the policy gradient,

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)]$$

# Reducing Variance Using a Baseline

- We subtract a baseline function  $B(s)$  from the policy gradient
- This can reduce variance, without changing expectation

$$\begin{aligned}\mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) B(s)] &= \sum_{s \in \mathcal{S}} d^{\pi_{\theta}}(s) \sum_a \nabla_{\theta} \pi_{\theta}(s, a) B(s) \\ &= \sum_{s \in \mathcal{S}} d^{\pi_{\theta}} B(s) \nabla_{\theta} \sum_{a \in \mathcal{A}} \pi_{\theta}(s, a) = 0\end{aligned}$$

- A good baseline is the state value function  $B(s) = V^{\pi_{\theta}}(s)$
- So we can rewrite the policy gradient using the **advantage function**  $A^{\pi_{\theta}}(s, a)$

$$\begin{aligned}A^{\pi_{\theta}}(s, a) &= Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s) \\ \nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) A^{\pi_{\theta}}(s, a)]\end{aligned}$$

# Estimating the Advantage Function (1)

- The advantage function can significantly reduce variance of policy gradient
- So the critic should really estimate the advantage function
- For example, by estimating *both*  $V^{\pi_{\theta}}(s)$  and  $Q^{\pi_{\theta}}(s, a)$
- Using two function approximators and two parameter vectors,

$$\begin{aligned}V_v(s) &\approx V^{\pi_{\theta}}(s) \\ Q_w(s, a) &\approx Q^{\pi_{\theta}}(s, a) \\ A(s, a) &= Q_w(s, a) - V_v(s)\end{aligned}$$

- And updating *both* value functions by e.g. TD learning

## Estimating the Advantage Function (2)

- For the true value function  $V^{\pi_{\theta}}(s)$ , the TD error  $\delta^{\pi_{\theta}}$

$$\delta^{\pi_{\theta}} = r + \gamma V^{\pi_{\theta}}(s') - V^{\pi_{\theta}}(s)$$

- is an unbiased estimate of the advantage function

$$\begin{aligned}\mathbb{E}_{\pi_{\theta}}[\delta^{\pi_{\theta}} | s, a] &= \mathbb{E}_{\pi_{\theta}}[r + \gamma V^{\pi_{\theta}}(s') | s, a] - V^{\pi_{\theta}}(s) \\ &= Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s) = A^{\pi_{\theta}}(s, a)\end{aligned}$$

- So we can use the TD error to compute the policy gradient

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}}[\nabla_{\theta} \log \pi_{\theta}(s, a) \delta^{\pi_{\theta}}]$$

- In practice we can use an approximate TD error

$$\delta_v = r + \gamma V_v(s') - V_v(s)$$

- This approach only requires one set of critic parameters  $v$

# Critics at Different Time-Scales

- Critic can estimate value function  $V_\theta(s)$  from many targets at different time-scales

- For MC, the target is the return  $v_t$

$$\Delta\theta = \alpha(v_t - V_\theta(s))\phi(s)$$

- For TD(0), the target is the TD target  $r + \gamma V(s')$

$$\Delta\theta = \alpha(r + \gamma V(s') - V_\theta(s))\phi(s)$$

- For forward-view  $TD(\lambda)$ , the target is the  $\lambda$ -return  $v_t^\lambda$

$$\Delta\theta = \alpha(v_t^\lambda - V_\theta(s))\phi(s)$$

- For backward-view  $TD(\lambda)$ , we use eligibility traces

$$\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

$$e_t = \gamma\lambda e_{t-1} + \phi(s_t)$$

$$\Delta\theta = \alpha\delta_t e_t$$

# Actors at Different Time-Scales

- The policy gradient can also be estimated at many time-scales

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A^{\pi_{\theta}}(s, a)]$$

- Monte-Carlo policy gradient uses error from complete return

$$\Delta\theta = \alpha (V_t - V_v(s_t)) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$

- Actor-critic policy gradient uses the one-step TD error

$$\Delta\theta = \alpha (r + \gamma V_v(s_{t+1}) - V_v(s_t)) \nabla_{\theta} \log \pi_{\theta}(s_t, a_t)$$



# Policy Gradient with Eligibility Traces

- Just like forward-view  $TD(\lambda)$ , we can mix over time-scales

$$\Delta\theta = \alpha(v_t^\lambda - V_v(s_t))\nabla_\theta \log \pi_\theta(s_t, a_t)$$

- where  $v_t^\lambda - V_v(s_t)$  is a biased estimate of advantage fn
- Like backward-view  $TD(\lambda)$ , we can also use eligibility traces
- By equivalence with  $TD(\lambda)$ , substituting  
 $\phi(s) = \nabla_\theta \log \pi_\theta(s, a)$

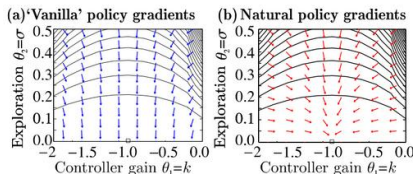
$$\begin{aligned}\delta &= r_{t+1} + \gamma V_v(s_{t+1}) - V_v(s_t) \\ e_{t+1} &= \gamma\lambda e_t + \nabla_\theta \log \pi_\theta(s, a) \\ \Delta\theta &= \alpha\delta e_t\end{aligned}$$

- This update can be applied online, to incomplete sequences

# Alternative Policy Gradient Directions

- Gradient ascent algorithms can follow any ascent direction
- A good ascent direction can significantly speed convergence
- Also, a policy can often be reparametrized without changing action probabilities
- For example, increasing score of all actions in a softmax policy
- The vanilla gradient is sensitive to these reparametrizations

# Natural Policy Gradient



- The **natural policy gradient** is parametrization independent
- It finds ascent direction that is closest to vanilla gradient, when changing policy by a small, fixed amount

$$\nabla_{\theta}^{nat} \pi_{\theta}(s, a) = G_{\theta}^{-1} \nabla_{\theta} \pi_{\theta}(s, a)$$

- where  $G_{\theta}$  is the Fisher information matrix

$$G_{\theta} = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a)^T]$$

# Natural Policy Gradient

- Using compatible function approximation,

$$\nabla_w A_w(s, a) = \nabla_\theta \log \pi_\theta(s, a)$$

- So the natural policy gradient simplifies,

$$\begin{aligned}\nabla_\theta J(\theta) &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) A^{\pi_\theta}(s, a)] \\ &= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \nabla_\theta \log \pi_\theta(s, a)^T w] \\ &= G_\theta w\end{aligned}$$

$$\nabla_\theta^{\text{nat}} J(\theta) = w$$

- i.e. update actor parameters in direction of critic parameters

# Summary of Policy Gradient Algorithms

- The policy gradient has many equivalent forms

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) v_t] \quad \text{REINFORCE}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a)] \quad \text{Actor-Critic}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) A_w(s, a)] \quad \text{Advantage actor-Critic}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta] \quad \text{TD Actor-Critic}$$

$$= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \delta e] \quad \text{TD}(\lambda) \text{ Actor-Critic}$$

$$G_{\theta}^{-1} \nabla_{\theta} J(\theta) = w \quad \text{Natural Actor-Critic}$$

- Each leads a stochastic gradient ascent algorithm
- Critic uses policy evaluation (e.g. MC or TD learning) to estimate  $Q^{\pi}(s, a)$ ,  $A^{\pi}(s, a)$  or  $V^{\pi}(s)$

# State of the art not yet covered

- Deep Deterministic Policy Gradient (DDPG)
- Asynchronous Advantage Actor-Critic Algorithm (A3C), Importance Weighted Actor-Learner Architectures (IMPALA)
- Trust Region Policy Optimization (TRPO), Proximal Policy Optimization (PPO)
- Soft Actor-Critic

# References

- Sutton, Richard S., and Andrew G. Barto. Reinforcement learning: An introduction. MIT press, 2018.
- Sutton, Richard S., et al. "Policy gradient methods for reinforcement learning with function approximation." Advances in neural information processing systems. 2000.
- <https://github.com/dalmia/David-Silver-Reinforcement-learning>, Reinforcement Learning Courses at UCL, David Silver