

Reinforcement Learning

Lecture 9

Policy Gradient

Finite Differences

REINFORCE

Actor-Critic Methods

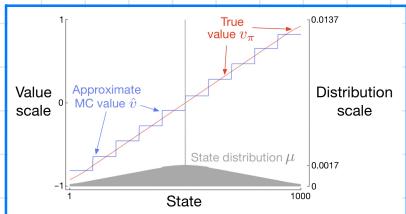
Robert Peharz

Institute of Theoretical Computer Science
Graz University of Technology

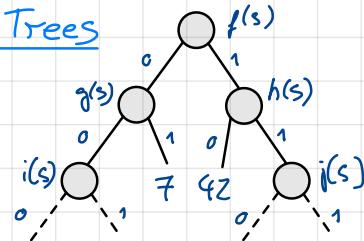
Winter Term 2023/24

Function Approximators

- aggregated states



- Decision Trees

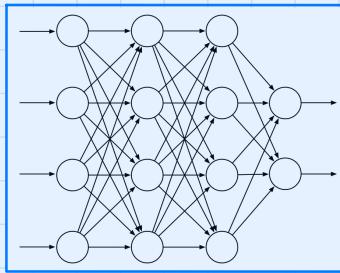


- Linear

$$\hat{v}(s) = \sum_{i=1}^K w_i x_i(s)$$

$$\hat{q}(s, a) = \sum_{i=1}^K w_i x_i(s, a)$$

- Neural Network



Differentiable Approximators

$$\nabla_{\underline{w}} \hat{v}(s, \underline{w}) := \begin{pmatrix} \frac{\partial \hat{v}(s, \underline{w})}{\partial w_1} \\ \frac{\partial \hat{v}(s, \underline{w})}{\partial w_2} \\ \vdots \\ \frac{\partial \hat{v}(s, \underline{w})}{\partial w_d} \end{pmatrix}$$

$$\nabla_{\underline{w}} \hat{q}(s, a, \underline{w}) := \begin{pmatrix} \frac{\partial \hat{q}(s, a, \underline{w})}{\partial w_1} \\ \frac{\partial \hat{q}(s, a, \underline{w})}{\partial w_2} \\ \vdots \\ \frac{\partial \hat{q}(s, a, \underline{w})}{\partial w_d} \end{pmatrix}$$

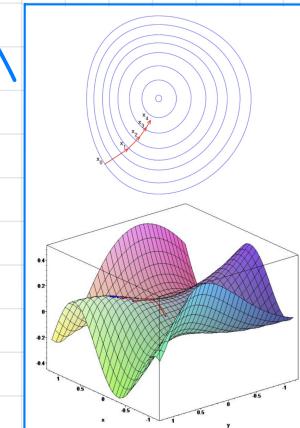
Gradient Descent for Supervised Learning

Stochastic Gradient Descent (SGD)

```

init  $\underline{w}$ 
repeat
   $\underline{w} \leftarrow \underline{w} - \alpha \nabla f_{\underline{w}}(\underline{w})$ 
until  $\nabla f_{\underline{w}}(\underline{w}) \approx 0$ 
  
```

stepsize



Value Function Approximation

$$l(\hat{v}(s, \underline{w}), v_{\pi}(s)) = \frac{1}{2} (v_{\pi}(s) - \hat{v}(s, \underline{w}))^2$$

$$\nabla_{\underline{w}} l = -(v_{\pi}(s) - \hat{v}(s, \underline{w})) \nabla_{\underline{w}} \hat{v}(s, \underline{w})$$

(P)rediction

(C)ontrol

quadratic loss

$$l(\hat{q}(s, a, \underline{w}), q_{\pi}(s, a, \underline{w})) = \frac{1}{2} (q_{\pi}(s, a, \underline{w}) - \hat{q}(s, a, \underline{w}))^2$$

$$\nabla_{\underline{w}} l = -(q_{\pi}(s, a) - \hat{q}(s, a, \underline{w})) \nabla_{\underline{w}} \hat{q}(s, a, \underline{w})$$

$$IE_{\pi}[(v_{\pi}(s) - \hat{v}(s, \underline{w}))^2] / IE_{\pi}[(q_{\pi}(s, a) - \hat{q}(s, a, \underline{w}))^2]$$

Monte Carlo $\underline{w} \leftarrow \underline{w} + \alpha (g - \hat{v}(s_t, \underline{w})) \nabla_{\underline{w}} \hat{v}(s_t, \underline{w})$

TD $\underline{w} \leftarrow \underline{w} + \alpha [R + \gamma \hat{v}(s', \underline{w}) - \hat{v}(s, \underline{w})] \nabla_{\underline{w}} \hat{v}(s, \underline{w})$

Linear Function Approximation

$$\hat{v}(s, \underline{w}) = \underline{w}^T \underline{x} = \sum_{i=1}^d w_i x_i(s)$$

$$\nabla_{\underline{w}} \hat{v} = \underline{x}$$

Easy Gradient, Convex Loss

$$\hat{q}(s, a, \underline{w}) = \underline{w}^T \underline{x} = \sum_{i=1}^d w_i x_i(s, a)$$

$$\nabla_{\underline{w}} \hat{q} = \underline{x}$$

Convergence of Algorithms

Prediction

On/Off-Policy	Algorithm	Table Lookup	Linear	Non-Linear
On-Policy	MC	✓	✓	✓
	TD(0)	✓	✓	✗
	TD(λ)	✓	✓	✗
Off-Policy	MC	✓	✓	✓
	TD(0)	✓	✗	✗
	TD(λ)	✓	✗	✗

Control

Algorithm	Table Lookup	Linear	Non-Linear
Monte-Carlo Control	✓	(✓)	✗
Sarsa	✓	(✓)	✗
Q-learning	✓	✗	✗
Gradient Q-learning*	✓	✓	✗

(✓) = chatters around near-optimal value function

Policy Gradient Methods

Learning
values \hat{v}, \hat{q}

Actor-Critic

Learning
policy $\tilde{\pi}$

(last time)

(today)

What are Policy Gradient Methods?

- so far, we focused on value functions v_{π} and q_{π}
- last lecture, we approximated value functions with differentiable approximators $\hat{v}(s, \underline{w})$, $\hat{q}(s, \underline{w})$ and learned them with stochastic gradient descent (SGD)
- Policy Gradient: parametrize policy with a vector $\underline{\theta}$

$$\pi(a|s) := \hat{\pi}(a|s, \underline{\theta})$$

and maximize a suitable performance measure $J(\underline{\theta})$ with gradient ascend:

$$\underline{\theta}_{t+1} \leftarrow \underline{\theta}_t + \alpha \nabla_{\underline{\theta}} J(\underline{\theta})$$

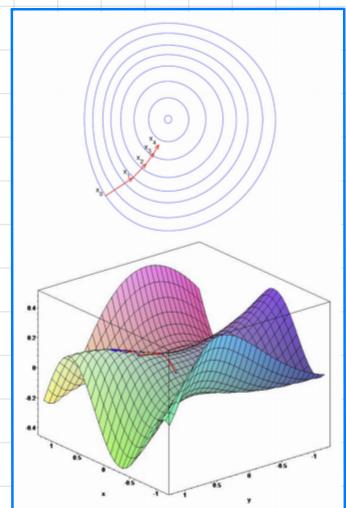


Image: D. Silver

Pros of Policy Gradient Methods

- directly learn $\hat{\pi}$ rather than "indirectly" via q (data efficiency)
- better convergence properties ————— non-linear value function approximators often unstable
- sometimes $\hat{\pi}(a|s, \underline{\theta})$ is easier to specify than $\hat{v}(s, \underline{w}), \hat{q}(s, a, \underline{w})$
- $\hat{\pi}(a|s, \underline{\theta})$ can encode domain knowledge (e.g. physical constraints), or be combined of pre-specified sub-policies (e.g. robot controllers)
- $\hat{\pi}(a|s, \underline{\theta})$ can learn stochastic policies, which is important in partially observed MDPs (POMDPs)

Example: aliased grid world, grey states are indistinguishable

- deterministic policy fails in 50% of the trials
- optimal: go left/right with probability 0.5 on grey states

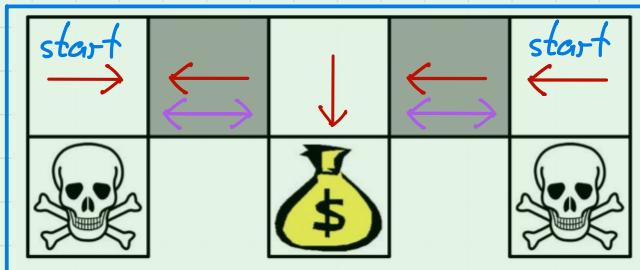


Image: D. Silver

Cons of Policy Gradient Methods

- local maxima — value-based RL (GPI) converges in principle to a globally optimal policy, while policy gradient typically gets stuck in local optima

Example [Jan Peters, scholarpedia.org]

Policy initialized with "straddle jump" won't find the "Fosbury flop"

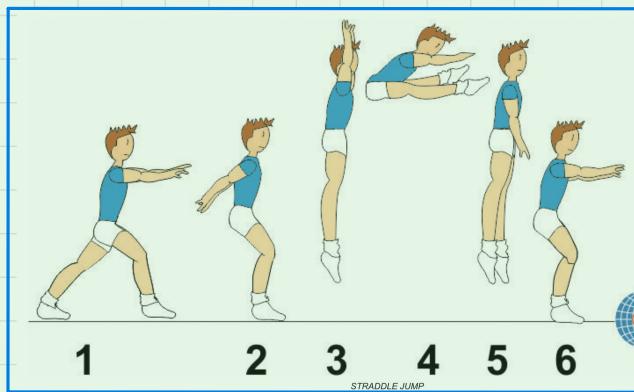


Image: fig-aerobic.com

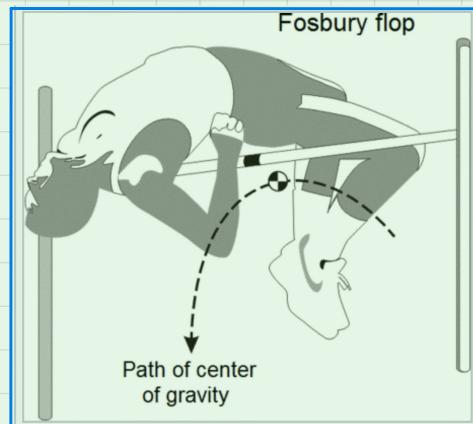


Image: [wikipedia.org](https://en.wikipedia.org)

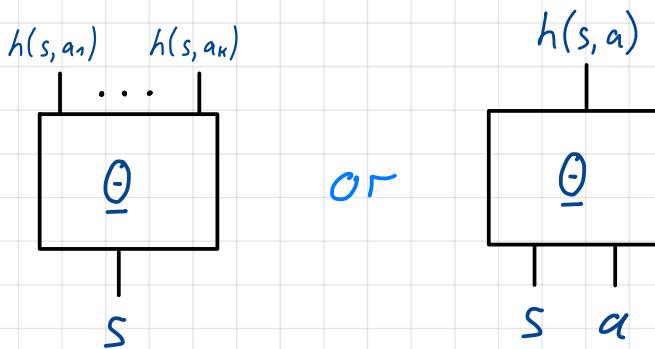
- ignores information about the environment (robustness, transferability)
- thus, might generalize poorly to (slightly) changed environments

Softmax Policy

generic parametrization for finite A

- specify "numeric preference score" $h(s, a, \underline{\theta})$, which is large if a should be preferred in state s , and low if it should be avoided.
- linear I: specify non-linear features $x_i(s, a)$ and let
$$h(s, a, \underline{\theta}) = \sum_i \theta_i x_i(s, a)$$
- linear II: specify non-linear features $x_i(s)$ and let
$$h(s, a, \underline{\theta}) = \sum_i \theta_{ia} x_i(s)$$

- neural network:



Softmax policy

$$\hat{\pi}(a|s, \underline{\theta}) := \frac{e^{h(s, a, \underline{\theta})}}{\sum_{a'} e^{h(s, a', \underline{\theta})}}$$

(= softmax layer in NNs)

(note that $\hat{\pi} > 0$ and $\sum_a \hat{\pi}(a|s) = 1$)

Behavioural Cloning - Supervised Pre-training

- we have now a parametrized policy $\tilde{\pi}(a|s, \theta)$
- say we have some data $D = \{(s_1, a_1), (s_2, a_2), \dots, (s_N, a_N)\}$ collected from (human) experts with policy $\tilde{\pi}_{\text{expert}}$
- learning from this data is just vanilla classification*!
- maximize conditional log-likelihood (* regression for continuous a)

$$J(\theta) = \sum_{i=1}^N \log \tilde{\pi}(a_i | s_i, \theta)$$

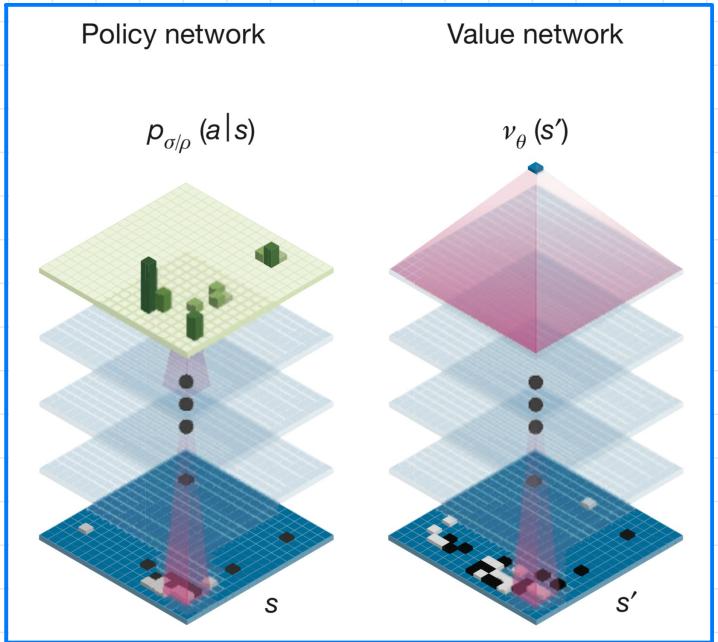
- behavioral cloning / supervised pre-training is often as an effective initialization

Pre-Training in Alpha Go

Mastering the game of Go with deep neural networks and tree search

David Silver^{1*}, Aja Huang^{1*}, Chris J. Maddison¹, Arthur Guez¹, Laurent Sifre¹, George van den Driessche¹, Julian Schrittwieser¹, Ioannis Antonoglou¹, Veda Panneershelvam¹, Marc Lanctot¹, Sander Dieleman¹, Dominik Grewe¹, John Nham², Nal Kalchbrenner¹, Ilya Sutskever², Timothy Lillicrap¹, Madeleine Leach¹, Koray Kavukcuoglu¹, Thore Graepel¹ & Demis Hassabis¹

NATURE | VOL 529 | 28 JANUARY 2016



- Alpha Go, which defeated the European Go Champion Fan Hui and later Grand Master Lee Sedol, used a "policy network", i.e. a parametrized policy based on a convolutional neural network
- the policy network was pre-trained on 30M positions from the KGS Go Server and was then refined using policy gradient
- Later versions (Alpha Go Zero) didn't use pre-training any more

Performance Metric for Interactive Setting

- for RL (interactive setting) we need a performance metric $J(\theta)$
- for episodic environments, we can use the value of a dedicated start state s_1

$$J_1(\theta) = v_{\pi}(s_1) = E_{\pi}[G_1 \mid s_1 = s_1]$$

- for continuing environments, we might use
 - expected value per time step

$$J_{EV}(\theta) = E_{\pi}[v_{\pi}(s)] = \sum_s \mu_{\pi}(s) v_{\pi}(s)$$

- or expected reward per time step

$$J_{ER}(\theta) = \sum_s \mu_{\pi}(s) r(s)$$

where μ_{π} is the stationary state distribution under π

Methods to Estimate the Policy Gradient

Finite Differences

REINFORCE (Williams 1992)

Actor-Critic Methods

Finite Differences

$$J(\underline{\theta}) = J_1(\underline{\theta}) = v_{\pi}(s_1)$$

- definition of partial derivative

$$\frac{\partial J}{\partial \theta_i} := \lim_{h \rightarrow 0} \frac{J(\dots, \theta_i + h, \dots) - J(\dots, \theta_i, \dots)}{h}$$

- finite difference: for small $h > 0$

$$\frac{\partial J}{\partial \theta_i} \approx \frac{J(\dots, \theta_i + h, \dots) - J(\dots, \theta_i, \dots)}{h}$$

- run M roll-outs for current $\underline{\theta}$ and for each $\underline{\theta} + h \in \overset{\leftarrow}{e}_i^{th}$ unit vector

Monte Carlo

$$\text{compute } \frac{\partial J}{\partial \theta_i} \approx \frac{\frac{1}{n} \sum g_m - \frac{1}{n} \sum g_{m,i}}{h}$$

g_m : m^{th} return with $\underline{\theta}$

$g_{m,i}$: m^{th} return with $\underline{\theta} + h \in \overset{\leftarrow}{e}_i$

+ simple to implement

+ can be applied to non-differentiable π as well

- expensive, only feasible for low-dimensional $\underline{\theta}$

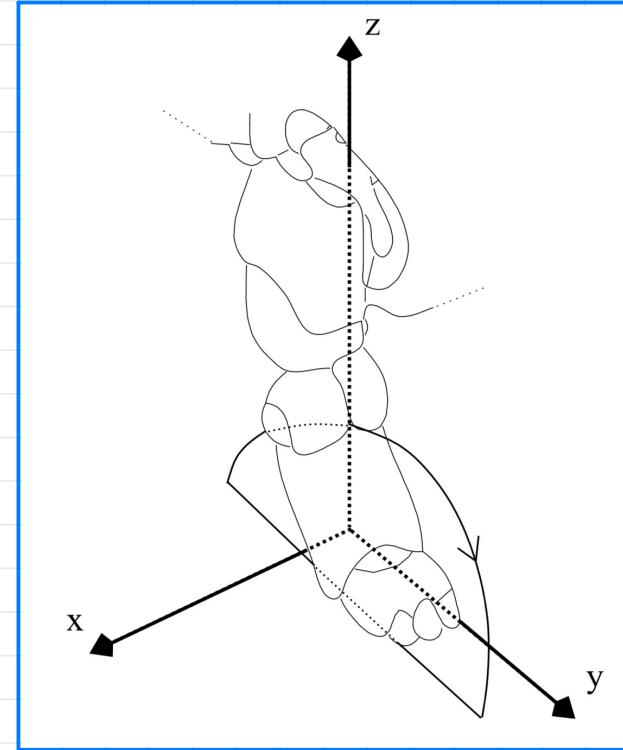
in total, we need $M(d+1)$ rollouts per gradient estimate

Teaching Aibo Fast Walk

Policy Gradient Reinforcement Learning for Fast Quadrupedal Locomotion

Nate Kohl and Peter Stone

Department of Computer Sciences, The University of Texas at Austin



- 12 parameters
- J : walking speed
- policy gradient with finite differences
- parallelized over several robots

All told, the following set of 12 parameters define the Aibo's gait [10]:

- The front locus (3 parameters: height, x -pos., y -pos.)
- The rear locus (3 parameters)
- Locus length
- Locus skew multiplier in the x - y plane (for turning)
- The height of the front of the body
- The height of the rear of the body
- The time each foot takes to move through its locus
- The fraction of time each foot spends on the ground

REINFORCE

Williams, 1992

- consider episodic tasks and $J_1(\underline{\theta}) = v_{\pi}(s_1) = \mathbb{E}_{\pi}[g_1 | s_1 = s_1]$
- how to compute $\nabla_{\underline{\theta}} J_1(\underline{\theta})$?
- let $s_1, a_1, r_2, s_2, \dots, s_{T-1}, a_{T-1}, r_T, s_T$ be an episode
- split episode into $\mathcal{T} = s_1, a_1, s_2, a_2, \dots, s_{T-1}, a_{T-1}, s_T$ and
 $\bar{R} = r_2, r_3, \dots, r_T$
- now, we can write

$$J_1(\underline{\theta}) = \mathbb{E}_{\mathcal{T}} \left[\mathbb{E}_{\bar{R}} [g_1 | \mathcal{T}] \right] = \mathbb{E}_{\mathcal{T}} [\bar{g}(\mathcal{T})]$$

→ just use observed return

$$\bar{g}(\mathcal{T}) \approx \sum_{k=1}^{T-1} \gamma^k r_{t+1} \quad \checkmark$$

- the gradient is tricky: does not depend on $\underline{\theta}$

$$\nabla_{\underline{\theta}} J_1(\underline{\theta}) = \nabla_{\underline{\theta}} \mathbb{E}_{\mathcal{T}} [\bar{g}(\mathcal{T})] = \nabla_{\underline{\theta}} \sum_{\mathcal{T}} p(\mathcal{T}; \underline{\theta}) \bar{g}(\mathcal{T})$$



1. \mathcal{T} is a (long) sequence

2. $p(\mathcal{T}; \underline{\theta})$ is determined by both the unknown MDP and π

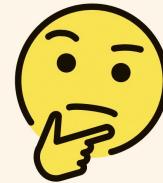
The Log-Derivative Trick

- let $p(x|\theta)$ be a distribution over some random variable X
- note that $\nabla_{\theta} \log p(x|\theta) = \frac{1}{p(x|\theta)} \nabla_{\theta} p(x|\theta)$
- thus, $\nabla_{\theta} p(x|\theta) = p(x|\theta) \nabla_{\theta} \log p(x|\theta)$
- now, note that

$$\begin{aligned}\underline{\nabla_{\theta} [E_x[f(x)]]} &= \nabla_{\theta} \sum_x p(x,\theta) f(x) \\ &= \sum_x \nabla_{\theta} p(x,\theta) f(x) \\ &= \sum_x p(x|\theta) \nabla_{\theta} \log p(x|\theta) f(x) \\ &= \underline{[E_x[\nabla_{\theta} \log p(x|\theta) f(x)]]}\end{aligned}$$

REINFORCE cont'd

$$\nabla_{\underline{\theta}} J_1(\underline{\theta}) = \nabla_{\underline{\theta}} \mathbb{E}_{\tau} \left[\bar{g}(\tau) \right] = \nabla_{\underline{\theta}} \sum_{\tau} p(\tau; \underline{\theta}) \bar{g}(\tau)$$



- $p(\tau) = \prod_{t=1}^{T-1} \pi(a_t | s_t) p(s_t | s_{t-1}, a_{t-1}) \prod_{t=1}^{T-1} \pi(a_t | s_t) \dots = \prod_{t=1}^{T-1} \pi(a_t | s_t, \underline{\theta}) p(s_{t+1} | s_t, a_t)$
- $\nabla_{\underline{\theta}} J_1(\underline{\theta}) = \nabla_{\underline{\theta}} \mathbb{E}_{\tau} \left[\bar{g}(\tau) \right]$

$$= \mathbb{E}_{\tau} \left[\nabla_{\underline{\theta}} \log p(\tau) \bar{g}(\tau) \right]$$

log-derivative trick

$$= \mathbb{E}_{\tau} \left[\nabla_{\underline{\theta}} \left(\sum_{t=1}^{T-1} \log \pi(a_t | s_t, \underline{\theta}) + \log p(s_{t+1} | s_t, a_t) \right) \bar{g}(\tau) \right]$$

does not depend on $\underline{\theta}$, $\nabla_{\underline{\theta}} \log p(\dots) = 0$

$$= \mathbb{E}_{\tau} \left[\left(\sum_{t=1}^{T-1} \overbrace{\nabla_{\underline{\theta}} \log \pi(a_t | s_t, \underline{\theta})}^{\text{automatic differentiation}} \right) \bar{g}(\tau) \right]$$

estimate with MC

REINFORCE cont'd

initialize $\underline{\theta}$

repeat

generate episode $S_1, A_1, R_2, S_2, A_2, \dots, S_{T-1}, A_{T-1}, R_T$ with $\pi(a|s, \underline{\theta})$

$$\tilde{\nabla}_{\underline{\theta}} J_1(\underline{\theta}) \leftarrow \left(\sum_{k=1}^{T-1} \gamma^k R_{k+1} \right) \left(\sum_{k=1}^{T-1} \nabla_{\underline{\theta}} \log \pi(A_t | S_t, \underline{\theta}) \right)$$

$$\underline{\theta} \leftarrow \underline{\theta} + \alpha \tilde{\nabla}_{\underline{\theta}} J_1(\underline{\theta})$$

- $\tilde{\nabla}_{\underline{\theta}} J_1(\underline{\theta})$ an unbiased estimate of $\nabla_{\underline{\theta}} J_1(\underline{\theta})$
- $E[\tilde{\nabla}_{\underline{\theta}} J_1(\underline{\theta})] = \nabla_{\underline{\theta}} J_1(\underline{\theta})$
- thus, under usual SGD conditions $\left(\sum_{i=1}^{\infty} \alpha_i = \infty, \sum_{i=1}^{\infty} \alpha_i^2 < \infty \right)$

REINFORCE converges to a local maximum of $J_1(\underline{\theta})$

REINFORCE with Baselines

- the REINFORCE estimator

$$\tilde{\nabla}_{\theta} J_1(\underline{\theta}) \leftarrow \left(\sum_{k=1}^{T-1} \gamma^k R_{k+1} \right) \left(\sum_{k=1}^{T-1} \nabla_{\theta} \log \tilde{\pi}(A_t | S_t, \underline{\theta}) \right)$$

is unbiased but usually has high variance

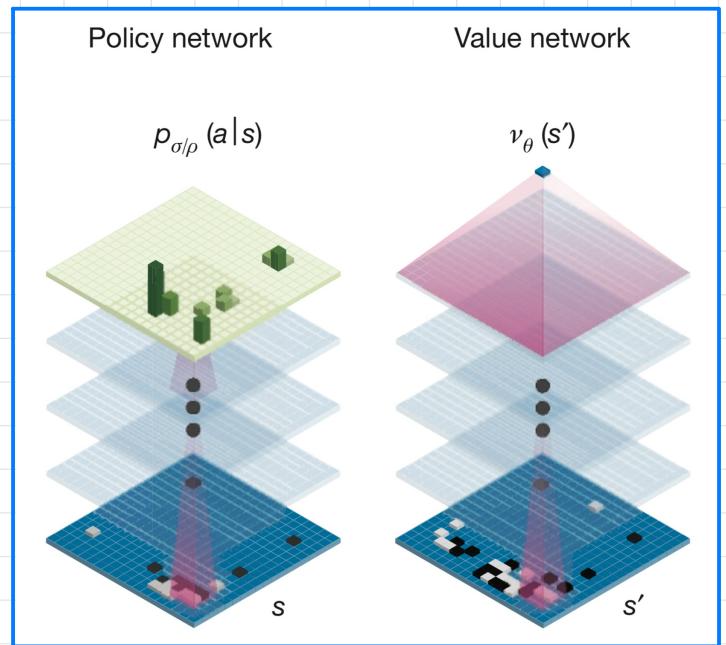
- we can subtract a constant "baseline" b from $\bar{g} = \sum_k \gamma^k R_{k+1}$
 $= \underline{b}$
- $\nabla_{\theta} J_1(\underline{\theta}) = \nabla_{\theta} E_{\tau} [\bar{g}(\tau)] - \nabla_{\theta} \underline{b} = \nabla_{\theta} E_{\tau} [\bar{g}(\tau) - b]$
- for example, we can estimate $b = \hat{v} \approx v_{\pi}(S_1)$, yielding

$$\tilde{\nabla}_{\theta} J_1(\underline{\theta}) \leftarrow \left(\left(\sum_{k=1}^{T-1} \gamma^k R_{k+1} \right) - \hat{v} \right) \left(\sum_{k=1}^{T-1} \nabla_{\theta} \log \tilde{\pi}(A_t | S_t, \underline{\theta}) \right)$$

- still unbiased, but lower variance
- hence, faster learning

REINFORCE in Alpha Go

- the pre-trained policy network (30M positions on KGS server) was then trained with REINFORCE
- self-play: one network is trained, while opponent is a previous version of the policy network
- the resulting policy network won 85% of the games against "Pachi," the strongest open source Go program back then
- (the final Alpha Go system used Monte Carlo Tree search combining the policy network with value function approximation)

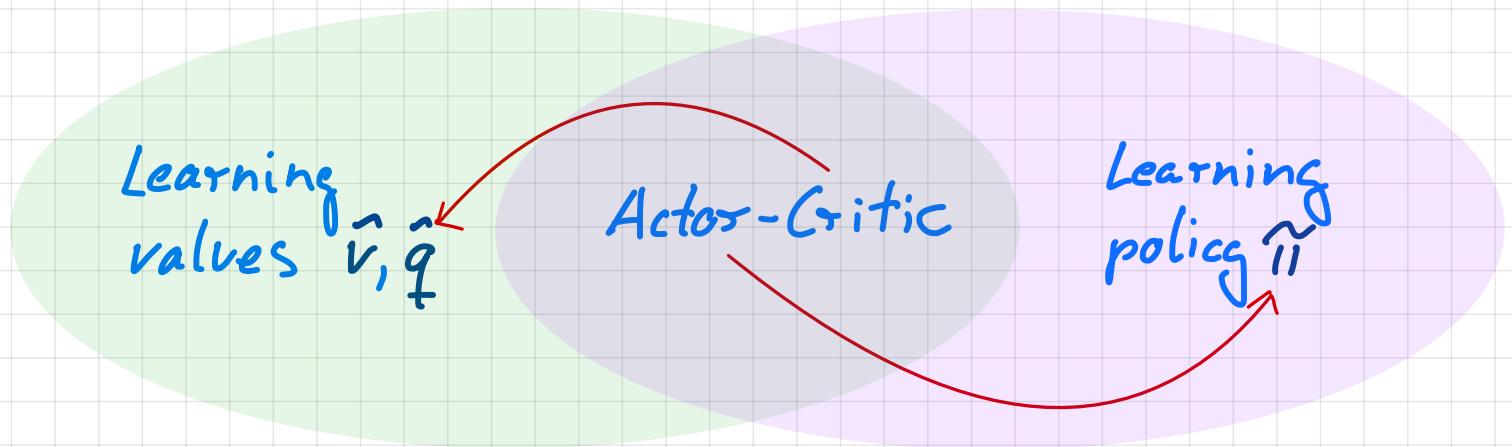


Towards Online Policy Gradients

$$\tilde{\nabla}_{\theta} J_1(\theta) \leftarrow G_1 \left(\sum_{t=1}^{T-1} \nabla_{\theta} \log \tilde{\pi}(A_t | S_t, \theta) \right)$$

- REINFORCE is a typical Monte Carlo algorithm
 - need to wait until end of episode
- can we get an online algorithm, similar to TD?
- under certain assumptions (ergodicity), also the following is an unbiased estimator of the policy gradient:
$$\tilde{\nabla}_{\theta} J_1(\theta) \leftarrow \sum_{t=1}^{T-1} G_t \nabla_{\theta} \log \tilde{\pi}(A_t | S_t, \theta)$$
- replace G_t with TD-style backups \rightarrow actor-critic methods

Actor-Critic Methods



- learn both $\hat{q}(\underline{w})$ (last lecture) and $\tilde{\pi}(\underline{\theta})$ simultaneously
- $\tilde{\pi}(\underline{\theta})$ follows approximate policy gradient given by $\hat{q}(\underline{w})$
- inherits instability of value function approximation
- however, actor-critic methods tend to be relatively stable in practice, since each update changes the policy only slightly
- in generalized policy iteration, on the other hand, small changes in \hat{q} can lead to dramatic changes in the greedy policy

TD Actor-Critic

input: stepsizes α^θ, α^w

initialize policy parameters $\underline{\theta}$ and \hat{q} parameters \underline{w}

repeat // for each episode (one for continuing tasks)

 initialize S

$A \sim \tilde{\pi}(a|S, \underline{\theta})$

 repeat until S is terminal

 take action A , observe S', R

$A' \sim \tilde{\pi}(a|S', \underline{\theta})$

$\delta \leftarrow R + \gamma \hat{q}(S', A', \underline{w}) - \hat{q}(S, A, \underline{w})$

// TD delta

$\underline{w} \leftarrow \underline{w} - \alpha^w \delta \nabla_{\underline{w}} \hat{q}(S, A, \underline{w})$

// update critic

$\underline{\theta} \leftarrow \underline{\theta} + \alpha^\theta \hat{q}(S, A, \underline{w}) \nabla_{\underline{\theta}} \log \tilde{\pi}(A|S, \underline{\theta})$

// update policy

$S, A \leftarrow S', A'$

Advantage Function - Baseline for Actor-Critic

- TD actor-critic is a one-step algorithm and has less variance than REINFORCE ✓
- still, \hat{q} might have significant variance
- we can use a baseline $b(s)$ and subtract it from \hat{q}
- a natural baseline is v_{π} , leading to the advantage function

$$a_{\pi}(s, a) := q_{\pi}(s, a) - v_{\pi}(s)$$

- the TD error $\delta = R + \gamma v_{\pi}(s') - v_{\pi}(s)$ estimates $a_{\pi}(s, a)$:
- $$\mathbb{E}[\delta | s, A] = \mathbb{E}[R + \gamma v_{\pi}(s') - v_{\pi}(s) | s, A] = q_{\pi}(s, A) - v_{\pi}(s)$$
- thus, instead of \hat{q} , learn \hat{v} as critic and use update

$$\theta \leftarrow \theta + \alpha^{\theta} \delta \nabla_{\theta} \log \pi(A | s, \theta)$$

$$\delta \approx R + \gamma \hat{v}(s') - \hat{v}(s)$$

Super Mario



Asynchronous Methods for Deep Reinforcement Learning

Volodymyr Mnih¹

Adrià Puigdomènech Badia¹

Mehdi Mirza^{1,2}

Alex Graves¹

Tim Harley¹

Timothy P. Lillicrap¹

David Silver¹

Koray Kavukcuoglu¹

¹ Google DeepMind

² Montreal Institute for Learning Algorithms (MILA), University of Montreal

VMNIH@GOOGLE.COM

ADRIAP@GOOGLE.COM

MIRZAMOM@IRO.UMONTREAL.CA

GRAVESEA@GOOGLE.COM

THARLEY@GOOGLE.COM

COUNTZERO@GOOGLE.COM

DAVIDSILVER@GOOGLE.COM

KORAYK@GOOGLE.COM

- actor-critic method using advantage function leading to super-human play of the classic game "Super Mario Bros"
- asynchronous: multiple agents in parallel environments
- Compared best against asynchronous Sarsa and Q-learning