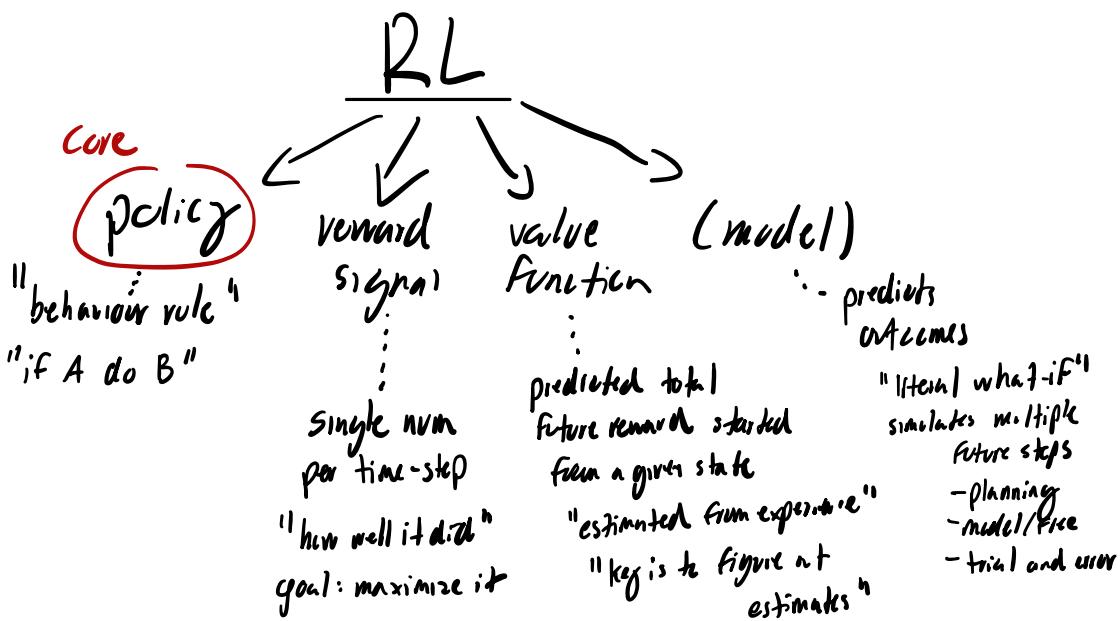
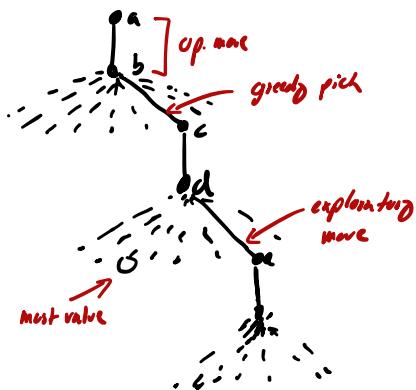


①



Tic Tac Toe

X	O	X
X	O	O
X	U	X



these

TD style

$$V(s_t) \leftarrow V(s_t) + \alpha [V(s_{t+1}) - V(s_t)]$$

learning rate

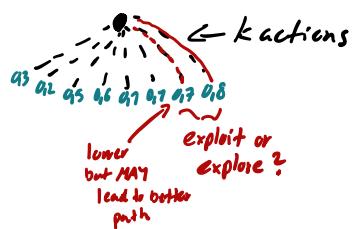
predictor error

+ next state better
- worsePart 1

(2)

K-bandit problem

- simplest RL problem



$q_t(a) = E[R_t | A_t=a]$

for action 'a'
true action value function
"true expected reward"

Expectation (average of all possible outcomes)
actual action at t
"given that" at t, we picked a
random reward at step t

$$Q_t(a) = \frac{\text{sum of rewards when a taken prior to } t}{\text{number of times a taken prior to } t}$$

$$A_t = \underset{\text{exploitation}}{\arg \max} Q_t(a) \quad \text{greedy}$$

ϵ -greedy: ex. $\epsilon=0.2 \rightarrow 80\%$ we exploit
 20% we pick random arm to explore

Incremental Implementation

instead of: $Q_n = \frac{R_1 + R_2 + \dots + R_n}{n}$ ← summing R_i 's each time

We can: $R_1 + \dots + R_n = n \cdot Q_n$ so: $Q_{n+1} = \frac{n \cdot Q_n + R_{n+1}}{n+1} = Q_n + \frac{1}{n+1} (R_{n+1} - Q_n)$ ← we "udge" the old average error

New Estimate = Old estimate + $\alpha [\text{Target} - \text{Old Estimate}]$

Non-stationary problem

Stationary: $\alpha = \frac{1}{n}$

Non-stationary: $\alpha \in (0,1)$

ex. $\alpha = 0.1$

↑
trust 90% of old info and 10% of new → 0.9^k rewards from many steps ago count to almost nothing

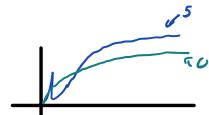
$$(1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha(1-\alpha)^{n-i} R_i$$

Exponentially shrinking remainder of initial guess

R_i decays the faster back in time i.e.

Optimistic initial values

if instead of $Q_t(a)=0$ we start with $Q_t(a)=5 \rightarrow$ more exploration, the final values will always be less than the old.



Upper-Confidence-Bound Action Selection

$$A_t = \underset{a}{\arg \max} [Q_t(a) + C \sqrt{\frac{\ln t}{N_t(a)}}]$$

cur. avg. rev. estm. exp. bonus

how large bonus grows slowly over time
big confidence at start, then gets smaller

We satisfy both goals:

- exploit action with highest estm. average so far
- explore action with uncertain estimates, because bonus $\sqrt{\frac{\ln t}{N_t(a)}}$ keeps it playing

Associative Search

2 slot machine: A: (0.1, 0.1) → usually we don't know B: (0.9, 0.8) with one we are pulling → 0.5

→ but! if we do know which we are pulling learn 2 minimacities → if A pull 2nd → 0.55
if B pull 1st

Gradient bandit algorithms

- rather than estimating each $q_t(a)$, we learn a set of preferences and turn them to probabilities (via soft-max)

- each arm, preference $H_t(a)$

- we apply Boltzmann or Gibbs soft-max:

$$\pi_t(a) = \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \quad \text{IF } H_t \text{ large, } \pi_t \text{ large}$$

Update Rule:

- after chose A_t we: compute $\bar{R}_t = \frac{1}{t} \sum_{i=1}^t R_i$

then:

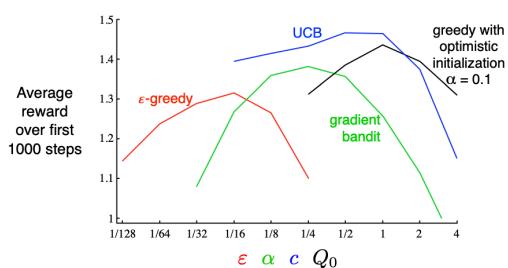
$$H_{t+1}(A_t) = H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \pi_t(A_t)),$$

$$H_{t+1}(a) = H_t(a) - \alpha(R_t - \bar{R}_t)\pi_t(a),$$

reduce

ensures we spread the penalty proportionally

repeat $t+1$



③ Finite Markov Decision Processes

MDP = states + actions + rewards + dynamics

Markov property: Future is independent from the past, given the present state
Ex. in maze, we need past states, to know where we came from \rightarrow not Markov.

Returns and episodes

episodic

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T,$$

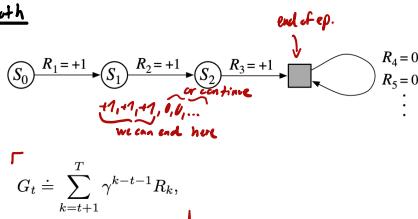
(to converge for infinite tasks)

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1},$$

Continuing \Rightarrow

where γ is a parameter, $0 \leq \gamma \leq 1$, called the discount rate.

both



Bellman Function simply

Value of $s =$ immediate reward + discount \times value of next state

$$V(s) = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t=s]$$

{ it is the expected sum of discounted rewards from that point onward }

Easy ex. $A \xrightarrow{1/2} B \xrightarrow{1/3} C \xrightarrow{1/1} \text{END}$ $\gamma=0.9$

we want: $v(A), v(B), v(C)$

$$v(C) = \frac{1}{2} \cdot 0.9 \cdot v(\text{END}) = 0$$

$$v(B) = 3 + 0.9 \cdot v(C) = 3.9$$

$$v(A) = 2 + 0.9 \cdot v(B) = 3.75$$

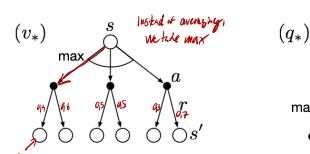
Optimal Policies and optimal value functions

$$V^*(s) = \max_a E[R_{t+1} + \gamma V^*(s_{t+1}) | S_t=s, A_t=a]$$

best achievable value

$$Q^*(s,a) = E[R_{t+1} + \gamma \max_{a'} Q^*(s_{t+1}, a') | S_t=s, A_t=a]$$

not best move



(4) Dynamic Programming

- if we have MDP and enough computation, we can turn

Bellman functions into iterative algorithms that converge to v^* and π^*

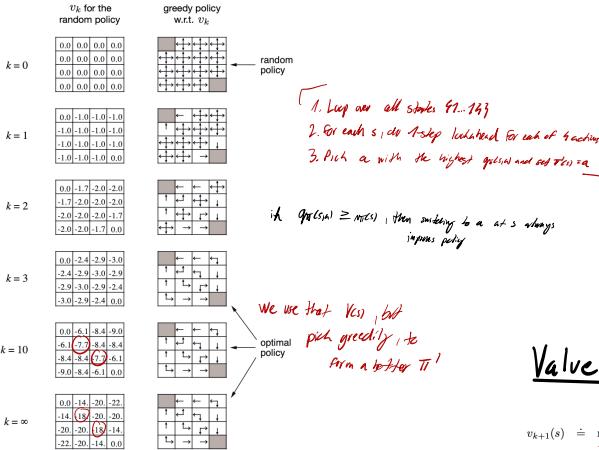
$$v_*(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

$$= \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma v_*(s')], \text{ or}$$

$$q_*(s,a) = \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a]$$

$$= \sum_{s',r} p(s',r|s,a) [r + \gamma \max_{a'} q_*(s',a')],$$

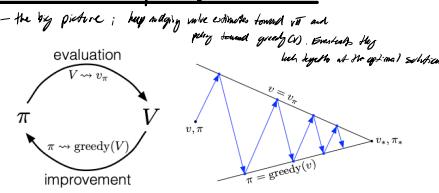
Policy Improvement



Asynchronous DP

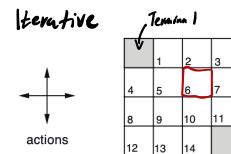
- instead of big sweep over every state, we only pick certain ones, that are 'hot'.

Generalized policy iteration



$\pi_* \longleftrightarrow V_*$

Policy evaluation (Prediction)



$R_t = -1$
on all transitions

exact B. eq.

$$V_E(6) = \frac{1}{\gamma} [(-1 + V_E(1)) + (-1 + V_E(2)) + (-1 + V_E(3)) + (-1 + V_E(5))]$$

but for iterative:

$$V_{k+1}(6) = \frac{1}{\gamma} \sum_{s \in \{1,2,3,5\}} (-1 + V_k(s))$$

Iterations:

$$1. V_{k+1}(6) = 0, \quad 1. V_k(6) = \frac{1}{\gamma} (C_1 + C_2 + \dots + C_4) = -1$$

$$2. V_{k+1}(6) = \frac{1}{\gamma} (-8) = -2, \quad 3. V_k(6) = -3, \dots$$

eventually they converge to a true V(s). (using number of iterations).

We can use 2 ways (old school) or 1 (values propagate faster)

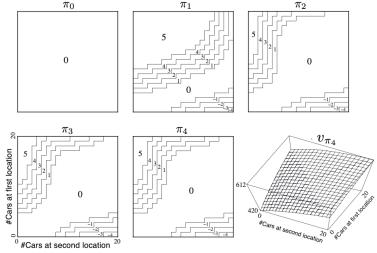
exact number.

Policy Iteration

- evaluate and improve until you can't do better

1. π_0 - random actions
2. repeat for $k = 0, \dots$
 1. Evaluate $\pi_k \rightarrow V^{\pi_k}$ by iterative sweeps
 2. Improve $\pi_k \rightarrow \pi_{k+1}$ by greedy one-step lookahead
 3. If $\pi_{k+1} = \pi_k$, converged, optimal
3. Output π^* and V^*

Ex. Rental problem



Value Iteration

- simply plugs the Bellman optimality equation directly into update rules, with convergence

- faster than full policy iteration

$$\begin{aligned} v_{k+1}(s) &\stackrel{\text{bellman}}{=} \max_a \mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = a] \\ &= \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')], \end{aligned}$$

Value ID. update

On gambler problem done

Goal: reach \$5 (win=1)
States = 1..5
Actions = 0,1..min(5-s,s)
Car: heads=0.4 (on heads we go up 1 if heads, or just +1 if tails)
 $s-a$
 $\gamma=1$

$$\begin{aligned} 1. V_0(0) &= 0 \\ 2. \forall j \in \{1, \dots, 4\} \quad V_1(j) &= V_0(j) \\ Q(j,a) &= Q_1(j)(V_1(j)) = 0.4(V_1(j+1)) + 0.6(V_1(j)) = 0 \Rightarrow V_1(j) = 0 \\ Q(5,a) &= 0 \quad Q(5,0) = Q_1(5) = V_1(6) = 0.4 \\ V_1(5) &= \max(Q_1(5)) = 0.4 \\ V_1(1) &= 0.4 \\ V_1(0) &= 0 \end{aligned}$$

$$2. \forall j \in \{1, \dots, 4\} \quad V_2(j) = Q_1(j)(V_1(j)) + Q_2(j)(V_1(j)) = 0.76$$

$$V_2 = [0.76, 0.76, 0.76, 0.76]$$

3...

slowly converges to true V^* $\rightarrow V_{\infty} = \arg \max Q_{\infty}$

$$\begin{aligned} S_1 &= 5-1, \text{ start} \\ S_2 &= 1 \\ S_3 &= 2 \\ S_4 &= 3 \\ S_5 &= 4 \end{aligned}$$

Efficiency of DP

- at first glance hopelessly expensive
- but when compared to computing any feasible policy π^*
- DP's cost is polynomial
- linear programming solutions DP is linear, but large scale LP hard
- real world is how we model our environment (MDP)

Sys: $S=1$, $A=1$

$S=2$: 7

$S=3$: 21

$S=4$: 63

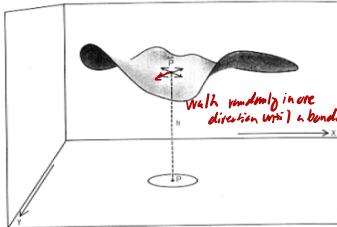
5 Monte Carlo Methods

Value Functions and Policies from Experience

- episodic tasks
- a way to learn exactly Bellman fixed-point solutions π^* and q^* that DP would compute but using only samples of experience

Monte Carlo Prediction (MCs)

Idea: run entire episodes under your π . Whenever you see 's' for the first time in an episode, record the full return from that point. Then $V(s)$ is the average of all these returns.



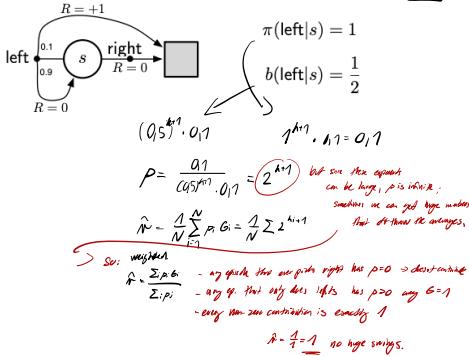
A bubble on a wire loop.

MC Control without exploring starts

- dropping unrealistic "exploring starts" assumption
- in real terms, we can't start agents in every state, so how do we ensure we try all actions?

E-greedy exploration

- same as MC control but with $\frac{E}{|A(s)|}$ prob. to explore.



Weighted IS is always along the gradient's choice in MC

Incremental Implementation

- For weighted off-policy MC maintains. For each state (s, a) constant-action pair:

 1. a running total of weights C_s , and
 2. the current weighted average Q_s , and an early return G with weight W :

$$C_s \leftarrow C_s + W, Q_s \leftarrow Q_s + \frac{W}{C_s} (G - Q_s)$$

'batch' weighted average: each return is with a weight W

$$V_n \leftarrow \frac{\sum_{k=1}^{n-1} W_k G_k}{\sum_{k=1}^{n-1} W_k}, n \geq 2,$$

Estimation of true value V_n by averaging all returns

Instead of re-running all steps, we keep track of:

$$C_n \leftarrow \sum_{k=1}^n W_k, V_n \leftarrow \frac{W}{C_n} (C_n - V_n), n \geq 1,$$

When we get next return, we update:

$$V_{n+1} \leftarrow V_n + \frac{W}{C_n} (C_n - V_n), n \geq 1,$$

and

$$C_{n+1} \leftarrow C_n + W_{n+1},$$

Off-policy MC Control

- learn optimal policy π^* (Chernoffistic, greedy wrt Q)
- act under separate soft policy π
- use weighted importance sampling - we've got into about T , and do little policy improvement toward goodness wrt Q update.

- Alg: 1. gather epis. value avg b
2. might not return if learned action/best b, would have taken other actions.
3. incrementally update π and make greedy wrt new Q
4. but stop early if each of when you've an action it wouldn't have taken.

*Discounting-aware IS

$$V(s) = \frac{\sum_{t \in T(s)} (1-\gamma)^{T(t)-1} \sum_{a \in A(s)} \pi_{t,a}(s, a) G_{t,a}}{|T(s)|}, \quad (5.9)$$

and a weighted importance-sampling estimator, analogous to (5.6), wrt

$$V(s) = \frac{\sum_{t \in T(s)} ((1-\gamma)^{T(t)-1} \sum_{a \in A(s)} \pi_{t,a}(s, a) G_{t,a} + \gamma^{T(t)-1} \mu_{t,T(t)-1}(G_{t,T(t)}))}{\sum_{t \in T(s)} ((1-\gamma)^{T(t)-1} \sum_{a \in A(s)} \pi_{t,a}(s, a) G_{t,a} + \gamma^{T(t)-1} \mu_{t,T(t)-1})}, \quad (5.10)$$

*Per-decision IS

- it avoids changing the entire IS ratio, massively cuts variance.

$$\text{estimator: } \hat{G}_t = \frac{1}{|T(s)|} \sum \tilde{G}_t$$

\downarrow independent action is

$$\tilde{G}_t = \rho_{t,1} R_{t+1} + \gamma \rho_{t,2} R_{t+2} + \gamma^2 \rho_{t,3} R_{t+3} + \dots + \gamma^{T-t-1} \rho_{t,T-1} R_T.$$

\downarrow standard ratios

$$\rho_{t,i} = \frac{1}{|T(s)|} \sum_{a \in A(s)} \pi_{t,a}(s, a)$$

Summary

- learns from full episodes
- estimates $V(s)$ by averaging returns G , stored after visiting a state
- each update depends on empirical returns, not on other value estimates.

3. Core Algorithms

Setting	Key Idea	Policy Improvement
On-policy MC Control	e-greedy episodes, first-visit MC on Q	Make it e-greedy wrt current Q
Off-policy Prediction	Importance sampling of returns under π to learn π^* (No improvement— π is fixed)	
Off-policy Control	Weighted IS + GPI: evaluate in π then greedify	Stop each episode at first non- π

4. Importance Sampling Variants

- Exploring Starts (MC-ES): force random (s_0, a_0) at episode start.
- ϵ -soft / e-greedy: keep $b(a|s) > 0 \forall a$ as automatic coverage.
- Off-policy: behavior b explores, target π is refined toward optimal.

Estimator	Bias	Variance	Use When
Ordinary IS (Eq 5.5)	unbiased	potentially w/ variance	theoretic analysis, small M
Weighted IS (Eq 5.6)	small bias \rightarrow as $M \rightarrow \infty$	boundless	practical off-policy MC
Per-decision IS (Eq 5.10)	unbiased	much lower	long episodes / $M=1$
Discount-aware IS (Eq 5.9/5.10)	unbiased or small bias	lower for $M=1$	deep discounts, reduce variance

⑥ Temporal-Difference Learning

Monte Carlo + DP
directly from our experience
After each step we get 'partial return' $R_{t+1} + \gamma V(S_{t+1})$
one step lookahead of the return

TD Prediction

Monte Carlo baseline full return
 $V(S_t) \leftarrow V(S_t) + \alpha[R_t + \gamma V(S_{t+1})]$

TD update: we stay bootstraped target, no sum as we step
 $R_{t+1} + \gamma V(S_{t+1})$

and update immediately:

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

TD error: signed difference between our step target and old estimate

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

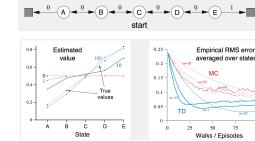
Alg.
Input: the policy π to be evaluated
Algorithm parameter: step-size $\alpha \in [0, 1]$
Initialize $V(s)$, for all $s \in S^*$, arbitrarily except that $V(\text{terminal}) = 0$
Loop for each episode:
 Loop for each step of episode:
 $A \sim \pi(s)$
 Take action A , observe R, S'
 $V(S') \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$
 $S \leftarrow S'$
 until S is terminal

TD(0)
updates after each step.
cautious!

Advantages of TD prediction methods

- updates immediately rather than waiting for opt. to end
- works on continuing problems (no episode)
- bootstraps, propagates info rapidly through state-action network
- model-free
- proven to converge, learns faster than MC

Ex Random walk



Optimality of TD(0)

Under batch updating: Comparing with convergence

- with MC batch: TD(0) converges in the correct values faster
than MC

MC batch - finds the sample-mean solution for stuck

TD(0) batch - finds solution consistent with maximum likelihood model of TD(0),
generates noise to force check

Expected Sarsa

- on-policy, but explores nextaction, takes with expected value under policy

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \mathbb{E}_\pi(Q(S_{t+1}, A_{t+1}) \mid S_{t+1}) - Q(S_t, A_t)]$$

$$\leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \sum_a \pi(a|S_{t+1})Q(S_{t+1}, a) - Q(S_t, A_t)]$$

- averages over all nextactions under Expected Sarsa



Figure 6.4: The backup diagrams for Q-learning and Expected Sarsa.

+ lower variance than Sarsa, simpler, more stable and safer than other?

- extra computation

Summary

- one-line update, fully online, model-free
- converge under broad conditions
- works for prediction, control, planning, controller domains, beyond games

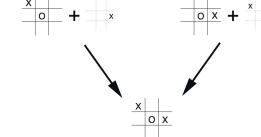
Games, Action States, special cases.

- in many domains you know exactly what the immediate effect of your action will be, but you don't know what happens after that (opponent move..)

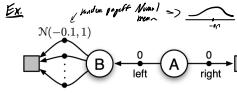
Workarounds - many (s,a) pairs may be seen in S, we can average the

- Policy learning
- clearer backups

$$\text{to just: } V_{\pi}(s) \leftarrow V_{\pi}(s) + \alpha[r + \gamma V_{\pi}(s')] - V_{\pi}(s)]$$



Maximization Bias and Double Learning



- Q-learning picks left far more than expected Expected Sarsa, it thinks action B is better than A. \rightarrow This is a **bias**.

Solver: Double Learning - decouple action selection logic from value estimation by maintaining 2 independent estimates.

On each step: - ADP pair

$$Q_1(s, a) \leftarrow Q_1(s, a) + \alpha[r + \gamma Q_2(s', \arg \max_b Q_2(s', b)) - Q_1(s, a)]$$

$$- if fails: Q_2(s, a) \leftarrow Q_2(s, a) + \alpha[r + \gamma Q_1(s', \arg \max_b Q_1(s', b)) - Q_2(s, a)]$$

- Over time: Q_1, Q_2 are uncorrelated \rightarrow bias cancels out

7 n-step Bootstrapping

- lets you see n steps future rewards before bootstrapping

n-step TD Prediction

n-step Return

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

$$\text{if } n=1 \Rightarrow G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1})$$

$n \rightarrow \infty \Rightarrow$ becomes full G_t return

update rule

- more collected rewards and reward act S_{t+n} , update the old G_t :

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t^{(n)} - V(S_t)] = \sum_{k=0}^{n-1} \gamma^k G_k$$

- at the end of n , we "flash" at perform the rest of n updates

- small $n \rightarrow$ low variance, higher bias

- large $n \rightarrow$ low bias, higher variance

- often looking for n it's unknown, which leaves higher than both.

error reduction property - in worst case the error shrinks by at least γ times γ^n

$$\max_s |\mathbb{E}_\pi[G_{t+n}|S_t=s] - v_\pi(s)| \leq \gamma^n \max_s |V_{t+n-1}(s) - v_\pi(s)|,$$



$$G_{t+3} = R_{t+4} + \gamma R_{t+5} + \gamma V(S_{t+6}) = 7+2+3+\gamma=6$$

$$V(A) \leftarrow V(A) + \alpha [G_{t+3} - V(A)]$$

$$\text{then: } (3+3+\gamma+1) \rightarrow 7+2+3 \rightarrow y_0 \rightarrow \text{update } V(S_t=6)$$

$$G_{t+4} = R_{t+5} + \gamma R_{t+6} + V(S_{t+7}) = 2+3+\gamma+0=7$$

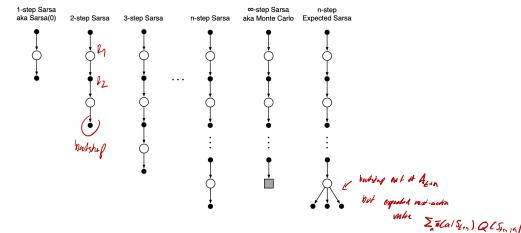
$$V(B) \leftarrow V(B) + \alpha [G_{t+4} - V(B)]$$

$$\text{then: } (4+3+\gamma+1) \rightarrow 9 \rightarrow \text{update } V(C)$$

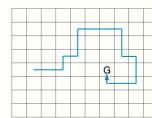
$$G_{t+5} = R_{t+6} + \gamma R_{t+7} = 3+1=4$$

↑ we run out of rewards, we stop the recycling

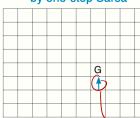
n-step Sarsa



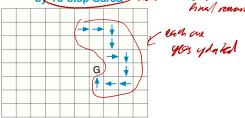
Path taken



Action values increased by one-step Sarsa



Action values increased by n-step Sarsa how far we search when final reward arrives



per-decision Methods with Control Variates

- we can do the update where $\rho=0$, don't think about return to 0, just saying that

$$G_{t+h} = R_{t+h} + \gamma G_{t+h+1}$$

$$G_{t+h} \doteq \rho_t(R_{t+h}) + \gamma G_{t+h+1} + (1-\rho_t)V_{t+h-1}(S_h), \quad t < h < T,$$

$$\rho \rightarrow 0 \rightarrow G_t = V(S_t)$$

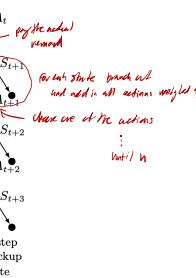
off-policy Learning without IS: n-step tree backup alg.

return

$$G_{t+h} = R_{t+h} + \gamma \sum_{a \in A_h} \pi(a|S_h) Q_{t+h-1}(S_h, a) + \gamma \pi(A_h|S_h) G_{t+h+1}$$

update

$$Q_{t+h}(S_h, A_h) \leftarrow Q_{t+h-1}(S_h, A_h) + \alpha [G_{t+h} - Q_{t+h-1}(S_h, A_h)]$$



Unifying Algorithm: n-step Q(σ)

- all put in unified family of n-step alg., called Q(σ), which interpolates between:

$\sigma=1 \rightarrow$ pure sampling and every step \rightarrow becomes n-step Sarsa

$\sigma=0 \rightarrow$ pure expectation w/ every step \rightarrow n-step tree-backup

$\sigma \in (0,1)$ \rightarrow sample a random σ in the tree and take the policy according to σ for the tree

\rightarrow hybrid, bias/variance

Let the horizon be $h = t + n$. Define the n-step $Q(\sigma)$ return recursively as

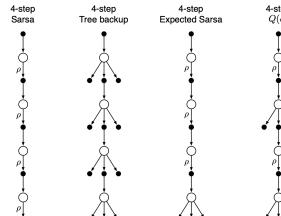
$$G_{t+h} = R_{t+1} + \gamma \left(\sigma_{t+1} \rho_{t+1} + (1 - \sigma_{t+1}) \pi(A_{t+1} | S_{t+1}) \right) (G_{t+1,h} - Q_{t-1}(S_{t+1}, A_{t+1})) + \gamma \bar{V}_{h-1}(S_{t+1})$$

where

$$V_{h-1}(S_{t+1}) = \sum_a \pi(a | S_{t+1}) Q_{h-1}(S_{t+1}, a),$$

and at the "leaf" (when $t+1 \geq T$ or $t+n \geq T$) you bootstrap to R_T or to 0 in the usual way.

$$Q_{t+h}(S_t, A_t) \leftarrow Q_{t+h-1}(S_t, A_t) + \alpha [G_{t+h} - Q_{t+h-1}(S_t, A_t)].$$



(8) Planning and Learning with Tabular Methods

Models and Planning

What is a model? - any mechanism by which an agent can predict the consequences of its actions.

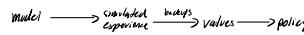
- outputs prediction + R and s'
- if it returns full probability dist over state $Q(s, a)$ \rightarrow distribution model
- if single sample $C(s, a)$, according to state probabilities \rightarrow sample model

Simulating experience via models - by feeding states and actions into model separately

What's planning? model $\xrightarrow{\text{planning}} \text{policy}$

- state-space planning: - search directly over states and not trees.
- plan-space planning: - search over entire action sequences or partial plans

Unified structure



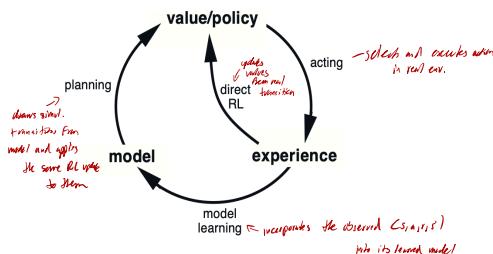
Random-sample one-step tabular Q-learning

```

Loop forever:
1. Select state  $S$  at random
2. Send  $S, A$  to sample model, and obtain a sample next reward,  $R$ , and state,  $S'$ 
3. Apply one-step tabular Q-learning to  $S, A, R, S'$ :
 $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$ 

```

Dyna: Integrated planning, Acting, and Learning



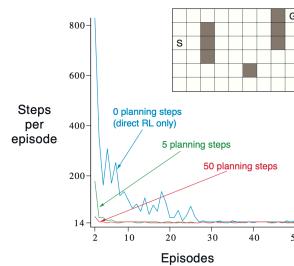
Tabular Dyna-Q

```

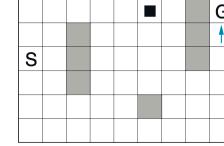
Initialize  $Q(s, a)$  and  $Model(s, a)$  for all  $s \in S$  and  $a \in A(s)$ 
Loop forever:
(a)  $S \leftarrow$  current (nonterminal) state
(b)  $A \leftarrow \epsilon\text{-greedy}(S, Q)$ 
(c) Take action  $A$ ; observe resultant reward,  $R$ , and state,  $S'$ 
(d)  $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$  - direct RL
(e)  $Model(S, A) \leftarrow R, S'$  (assuming deterministic environment) - learning
(f) Loop repeat  $n$  times:
  S  $\leftarrow$  random previously observed state
  A  $\leftarrow$  random action previously taken in  $S$ 
   $R, S' \leftarrow Model(S, A)$ 
   $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$  - planning

```

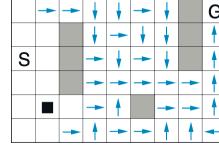
Ex. maze



WITHOUT PLANNING ($n=0$)



WITH PLANNING ($n=50$)

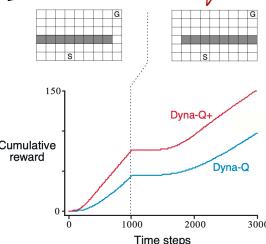


it's a lucky tuple of past outcomes - memory!

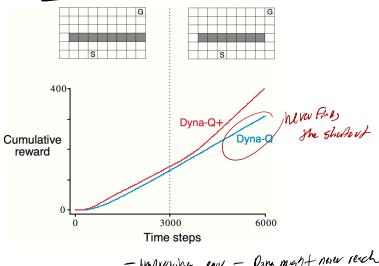
When the Model is Wrong

- When model is incorrect, can produce policies that exploit "phantom" opportunities.

Ex1



Ex2



Dyna-Q+

- exploration bonus heuristic

if trans. hasn't been tried in T steps, then add it to reward during planning

$$r + kT$$

small constant

Expected vs Sample Updates

\hookrightarrow complete exact Bellman backup by summing over all possible next events

Distribution model vs sampled update:

$$Q_{\text{exact}} \leftarrow \sum_s \beta(s|s, a) [r + \gamma \max_a Q_{\text{exact}}(s, a)]$$

- instead of zero variance updates, but costs b of possible successor states

+ much less expensive (O(b))

7 (classic TD(0) backup) down to single (O(1)) from sample model or real env and do:

$$Q_{\text{exact}} \leftarrow Q_{\text{exact}} + \alpha [R + \gamma \max_a Q_{\text{exact}}(s', a) - Q_{\text{exact}}]$$

- each update costs constant time, but reduces sampling noise

+ brings back cheap O(1)

Quantifying the tradeoff

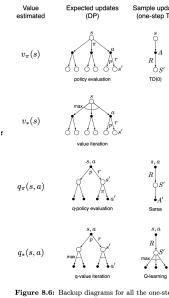
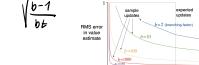


Figure 8.6: Backup diagrams for all the one-step updates considered in this book.

Trajectory Sampling

2 ways of distributing planning updates

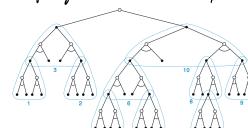
1. Exhaustive sweeps - split length long share, compute exact Bellman backup

2. Trajectory sampling - instead of blind sweep, sample certain episodes from model following policy

Heuristic Search

decision time planning that:

- builds a tree of possible continuations from current state
- uses one-step Bellman backups to propagate leaf-node continuations back to root
- discards world of action after selecting the best one
- focuses all planning effort on current decision, prioritizes when longer lookahead for single choice.



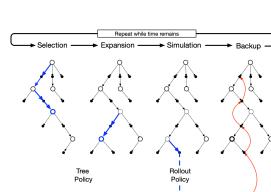
Rollout Algorithms

decision time planning method that use Monte Carlo control on model to choose action

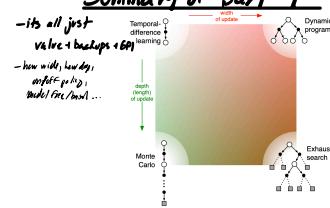
roll right:
1. right \rightarrow cost 1, $r = -1$
2. under 5 pcts of $P(Cost, r = -1)$
 $\Rightarrow 6^{10} = 2$

roll 2:
1. $\rightarrow C(1, 1)$
2. $\rightarrow C(1, 2)$
 $\Rightarrow 6^{10} = 2$

avg = -2
we going choose the best one.



Summary of Part 1



Part 2: Approximate Solution Methods

- from small, tabular problems to real world, enormous and continuous
- for robotics, vision, games - number of possible states explosive
- parameterized function instead of exact values - supervised learning

⑨ On-policy Prediction with approximation

Value function approximation

- returns each Bellman-style update as a supervised learning example and plugs in a general function approximator in place of lookup table.

Every RL update is: $s \rightarrow a$

- we treat (s, a) as a training pair for function approximator
with aim to reward state s from action a to target r

Weight vector $w \in \mathbb{R}^d$ and mapping $\hat{V}(s, a) \approx V(s, a)$

Here choices are: linear functions

- neural networks

- decision trees, linear numbers.

Generalization vs. bias:

↳ update only changes $V(s)$

update at s changes $\hat{V}(s, a)$ for all a , which is good

depends on the same weight - one sample "represents" whole region of state space

(while standard supervised learning - chain carries responsibility, see the Fig.)

- happens as often change over time

- applies rewards and updates value function maintaining targets

The prediction objective (VE)

- change to parameter vector w needs many states; we must choose the most important ones,

Choosing state-weighting distribution $MCS \geq 0$ and $\sum_i MCS = 1$ to express how much we care about a in each state

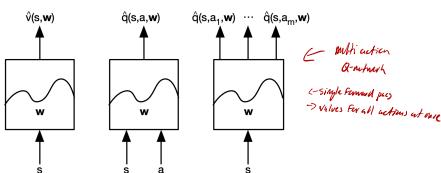
(Continuing with David Silver's RL course slides)

The problem:

We represent V by lookup tables \rightarrow for real problems, impossible \rightarrow too many states

Solution: estimate value function with approximation:

$$\begin{aligned}\hat{V}(s, w) &\approx V(s) \\ \hat{q}(s, a, w) &\approx q(s, a)\end{aligned}$$



Gradient descent

Goal: find parameter vector w minimising mean-squared error between approximate value for (s, a, w) and true value $r_{\pi}(s)$

Loss function: $J(w) = E_{\pi}[(\hat{V}(s, w) - r_{\pi}(s))^2]$

Gradient of $J(w)$: diff. function of w

$$\nabla_w J(w) = \begin{pmatrix} \frac{\partial J}{\partial w_1} \\ \vdots \\ \frac{\partial J}{\partial w_d} \end{pmatrix}$$

$$\Delta w = -\frac{1}{2} \lambda \nabla_w J(w)$$

"more the weight, larger the gradients"

$$E_{\pi} J(w) = \lambda w^2 + 5$$

$$\lambda' = 6w$$

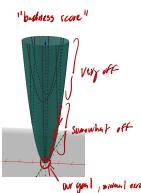
$$E_{\pi} J(w_1, w_2) = 3w_1^2 + 5w_2^2 \quad \text{inital training weights}$$

$$J(w_1) = 6w_1^2$$

$$J(w_2) = 5w_2^2$$

$$(w_1) \leftarrow (w_1) - \lambda \frac{\partial J}{\partial w_1} = (1 - 6w_1)w_1$$

$$(w_2) \leftarrow (w_2) - \lambda \frac{\partial J}{\partial w_2} = (1 - 5w_2)w_2$$



Feature Vectors

$$x(s) = \begin{pmatrix} x_1(s) \\ x_2(s) \\ \vdots \\ x_n(s) \end{pmatrix}$$

real-valued summary or representation of one agent's state

Linear VFA

represent $V(s, w)$ as linear combination of basis functions:

$$V(s, w) = x(s)^T w = \sum_{j=1}^n x_j(s) w_j$$

$$J(w) = E_{\pi} [(v_{\pi}(s) - x(s)^T w)^2]$$

because $\hat{r}_{\pi}(s, a)$ is linear, finally: $\nabla_w J(w) = x(s)$

Stochastic gradient update:

$$\Delta w = d(\hat{r}_{\pi}(s) - \hat{r}_{\pi}(s, a)) \nabla_w J(w) = d(\hat{r}_{\pi}(s) - x(s)^T w) x(s) \quad \text{replace by sample target}$$

One-step TD update: $\Delta w = d(\hat{r}_{\pi}(s) + \gamma \hat{r}_{\pi}(s_{t+1}, a) - \hat{r}_{\pi}(s, a)) \nabla_w J(w)$

Full gradient descent step: $\Delta w = d(\hat{r}_{\pi}(s) - \hat{r}_{\pi}(s, a)) \nabla_w J(w)$

Stochastic GD step - sampling over all the states

$$\Delta w = d(\hat{r}_{\pi}(s) - \hat{r}_{\pi}(s, a)) \nabla_w J(w)$$

$$\text{we don't have this! Instead, we play sampled return: } [MC: E_{\pi} \quad TDCI: E_{\pi} + \gamma \hat{r}_{\pi}(s_{t+1}, a)]$$

$$\text{One-step TD update: } \Delta w = d(\hat{r}_{\pi}(s) + \gamma \hat{r}_{\pi}(s_{t+1}, a) - \hat{r}_{\pi}(s, a)) \nabla_w J(w)$$

In practice, substitute MCS for sample target

$$\Delta w = d(\hat{r}_{\pi}(target) - \hat{r}_{\pi}(sample)) \nabla_w J(w)$$

prediction error

MC $G_t = R_t + \gamma G_{t+1} \dots$

We call it supervised learning: treat each (s_t, G_t) as training example, goal is to $\hat{r}_{\pi}(s_t, a) \approx G_t$

Linear MC policy cost: $\Delta w = d(G_t - \hat{r}_{\pi}(s_{t+1}, a)) x(s_t)$

\rightarrow converges to local optimum, even with non-linear FA

TD same thing: $\langle S_1, R_1, \gamma \hat{r}_{\pi}(S_2, a) \rangle, \dots$

$$\Delta w = d(R_1 + \gamma \hat{r}_{\pi}(S_2, a) - \hat{r}_{\pi}(S_1, a)) \nabla_w J(w)$$

TDCI $\langle S_1, G_1 \rangle, \dots \langle S_{T-1}, G_{T-1} \rangle$

Forward pass in: $\Delta w = d(G_1 + \gamma \hat{r}_{\pi}(S_2, a) - \hat{r}_{\pi}(S_1, a)) \nabla_w J(w)$

Backward pass in: $\begin{cases} G_2 = R_2 + \gamma \hat{r}_{\pi}(S_3, a) \dots \\ E_{\pi} = \gamma A E_{\pi-1} + \hat{r}_{\pi}(S_2, a) \end{cases}$

$\Delta w = d(E_{\pi} - \hat{r}_{\pi}(S_1, a)) \nabla_w J(w)$

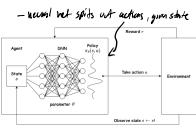
{ equivalent}

Policy Gradient

Policy based RL

- instead of going through value function, we directly represent the policy itself with parameters θ

$$\pi_\theta(s,a) = P(a|s, \theta)$$



- we stay model-free

Advantages of policy RL

- + better convergence behavior - no more gradient vanishing
- + effective in continuous or high-dimensional action spaces
- + learns stochastic policies naturally
- typically converge to local rather than global optimum
- inefficient evaluating policy in high variance

Policy Optimisation

Goal: Find θ^* that maximises $J(\theta)$

Gradient-Free: Hill Climbing, Simplex, Genetic Alg.

Gradient-based: Stochastic gradient descent/ascent

- Conjugate Gradient

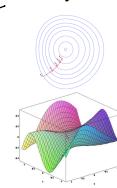
- Quasi-newton

Simple recipe

1. Pick policy objective $J(\theta)$ - start std. random - and view it as function of policy parameters θ

2. Compute policy gradient

$$\nabla_\theta J(\theta) = \left(\frac{\partial J}{\partial \theta_1}, \frac{\partial J}{\partial \theta_2}, \dots \right)$$



3. Take small ascent step: $\Delta\theta = \eta \nabla_\theta J(\theta)$

update: $\theta \leftarrow \theta + \Delta\theta$
repeat.

Score Function

$$V_\theta(s,a) = \pi_\theta(s,a) \frac{\nabla_\theta J(s,a)}{\nabla_\theta J(s,a)}$$

score function: how does the probability of taking an action change if we change θ ?

One-step MDPs

1. start w/ dist. over all actions
2. pick an action
3. receive R, update dist.

$$J(\theta) = \mathbb{E}_\theta [R] = \sum_a (\sum_s \pi_\theta(s|a) R_{sa})$$

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

Policy objective Functions

Goal: given parametric policy $\pi_\theta(a|s)$ find best parameters θ^*

- we need w/ $J(\theta)$ tells us "how good" a candidate policy is
- then we do gradient ascent on $J(\theta)$ to refine policy π_θ

Episodic: Start state value $V_\theta(s)$ compute G from s , following policy π_θ

$$J(\theta) = \mathbb{E}_\theta [V_\theta(s)] = \mathbb{E}_\theta [G]$$

↳ actions begin in start state s_0

↳ $\pi_\theta \rightarrow$ better Policy

1. Continuous: average-value (discounted future value)

- we care about how poly does in average over time
- value \mathbb{E}_θ : the Markov chain induced over time until settle into stationary distribution $\pi_\theta(s)$

(It is a long-run fraction of the agent's visits in each state s)
expected discounted return

$$\text{Objective: } J_\theta(s) = \mathbb{E}_\theta [G] V_\theta(s) \leftarrow \mathbb{E}_\theta [G] \pi_\theta(s)$$

for each state s :
roughly how much we will get if we explore in s
but only if we're optimistic are one state

2. Continuous: long-run average reward per time-step (Conditional immediate reward rate)

$$J_\theta(s) = \sum_t \mathbb{E}_\theta [r_t | s_t] \sum_t \pi_\theta(s_t) \text{ expected immediate reward given state } s_t$$

- why? we want a steady performance rather than discounted sum
- we use no discounting

Software policy

- we use probability of action regardless of θ

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

Score function is

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

$$\sum_a \pi_\theta(a|s) \phi(s,a)$$

Gaussian Policy

$$\mu_\theta(s) = \phi(s)^T \theta$$

- Sampling rule: $a \sim \mathcal{N}(\mu_\theta(s), \sigma^2)$

- continuous spaces

- easy to sample and differentiate

$$\text{Score function: } \nabla_\theta \log \pi_\theta(a|s) = \frac{(a - \mu_\theta(s)) \phi(s)}{\sigma^2}$$

Computing gradients by Finite differences

- we have policy π_θ and objective $J(\theta)$ that we can easily calculate (by running gradient) w/ differentials only once.

We still want Taylor

Finite difference approximation: 1) For each k :
 $\theta^{(k+1)} = \theta + \epsilon_k$
diff w.r.t. dimension k

2) New things include: $\theta^{(k)}$ to estimate $J(\theta^{(k)})$
similarly estimate $J(\theta^{(k+1)})$

3) Derivatives approx: $\frac{\partial J}{\partial \theta_k} \approx \frac{J(\theta^{(k+1)}) - J(\theta^{(k)})}{\epsilon_k}$

- simple, wrong, inefficient, sometimes different

$\nabla_\theta J(\theta)$

One-step MDPs

1. start w/ dist. over all actions
2. pick an action
3. receive R, update dist.

$$J(\theta) = \mathbb{E}_\theta [R] = \sum_a (\sum_s \pi_\theta(s|a) R_{sa})$$

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

Policy Gradient Theorem

- undiscounted

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

(higher λ limit log mean value)

Actor-Critic methods (lower variance)

- an actor - policy (MC, TD)

- a critic - learned action value estimate (QLearn), for lower variance estimate

$Q(s,a) \approx Q^\pi(s,a)$

↳ (MC) just policy eval (QLearn, TD), last square..

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

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

- Simple actor-critic algorithm based on action-value critic
- Using linear value fn approx. $Q_w(s,a) = \phi(s,a)^T w$
- Critic: Updates w by linear TD(0)
- Actor: Updates θ by policy gradient

function QAC

Initialise s , θ

Sample $a \sim \pi_\theta$

for each step do

Sample reward $r = R_s^a$; sample transition $s' \sim \mathcal{P}_s^a$.

Sample action $a' \sim \pi_\theta(s', a')$

$\delta = r + \gamma Q_w(s', a') - Q_w(s, a)$

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

$w \leftarrow w + \beta \delta \phi(s, a)$

$a \leftarrow a', s \leftarrow s'$

end for

end function

Monte Carlo Policy Gradient (REINFORCE)

- simplest actor-only method; full episode returns multiplied by score-function gradients

- in discrete/continuous; but high variance!

function REINFORCE

Initialise θ arbitrarily

for each episode $\{s_1, a_1, r_1, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_\theta$ do

for $t=1$ to $T-1$ do

$$\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t)$$

end for

return θ

end function

Reducing Variance with Baseline

- δ 's can be very noisy from sample to sample; gradients become jiggly and

- we can subtract any random or static $B(s)$ from our return besides the expectation, then expected gradient remains unchanged:

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

So we can safely do:

$$\Delta\theta \leftarrow \mathbb{E}_\theta [\nabla_\theta \log \pi_\theta(s,a) (\delta - B(s))]$$

The best baseline - is state-value function $B(s) = V^\pi(s)$

The advantage Function: $A^\pi(s,a) = Q^\pi(s,a) - V^\pi(s)$

and rewrite policy gradient theorem:

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

Estimating Advantage Function

1) Learn Q and V , then $A^\pi(s,a) = Q^\pi(s,a) - V^\pi(s)$

2) ~~compute~~ TD errors as advantage:

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

has expectation:

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

So we can just:

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

Critics at different time scales

■ For MC, the target is the return v_t

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

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

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

■ For forward-view TD(λ), the target is the λ-return v_t^λ

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

■ For backward-view TD(λ), we use eligibility traces

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

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

$$\Delta\theta = \alpha e_t \phi(s_t)$$

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

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

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

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

■ TD Actor-Critic

■ Natural Actor-Critic



Summary of Policy Gradient Algorithms

■ The **policy gradient** has many equivalent forms

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

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 $V^\pi(s)$