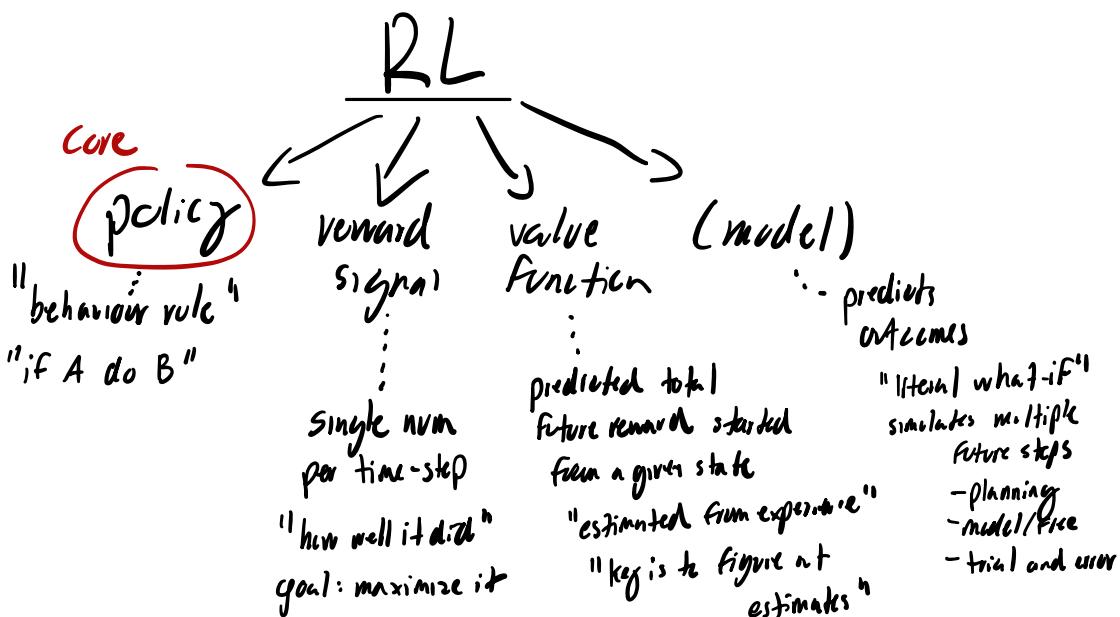
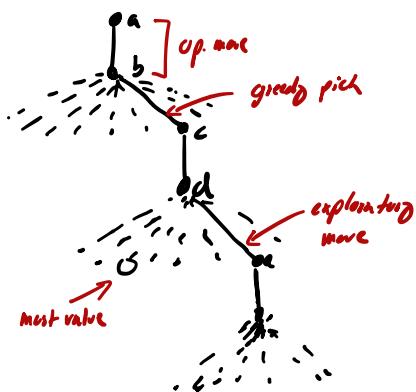


①



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

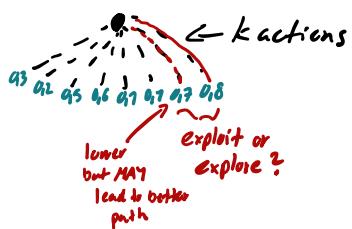
predictor error

+ 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

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

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

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

Incremental Implementation

instead of: $Q_n = \frac{R_1 + R_2 + \dots + R_n}{n} \quad \leftarrow \text{suming } R_i \text{ each time}$

We can: $R_1 + \dots + R_n = n \cdot Q_n \quad \text{so: } Q_{n+1} = \frac{n \cdot Q_n + R_{n+1}}{n+1} = Q_n + \frac{1}{n+1} (R_{n+1} - Q_n)$

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

Non-stationary problem

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

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

ex. $\alpha = 0.1$

\nearrow
 trust 90% of old info and 10% of new $\rightarrow 0.9^k$ rewards from many steps ago const to almost nothing

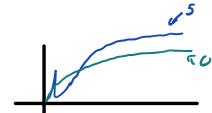
$$(1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

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

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

Cur. avg. rev. estm. \nearrow expl. bonus \searrow
 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 $\sqrt{\frac{1}{N_t(a)} / N_t(a)}$ keeps it play

Associative Search

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

\rightarrow but! if we do know which we are pulling learn 2 minimizations \rightarrow if A pull 2nd $\rightarrow 0.55$
 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:

$$\pi_t(a) = \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \quad \begin{matrix} \text{IF } H_t \text{ large, } \pi_t \text{ large} \\ \text{prop. of picking } a \text{ in step } t \end{matrix}$$

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

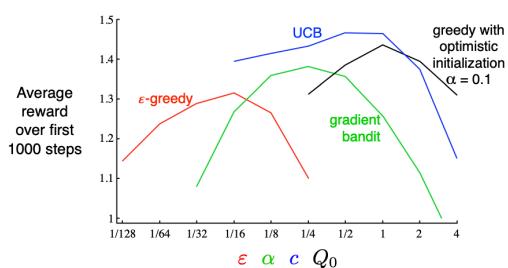
\nearrow normalizing factor

and

for all $a \neq A_t$,

\uparrow
 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

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

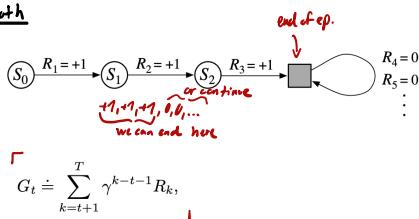
(to converge for infinite tasks)

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

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. $A \xrightarrow{1/2} B \xrightarrow{1/3} C \xrightarrow{1/1} \text{END}$ $\gamma = 0.5$

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

$$v(C) = \frac{1}{2} \cdot 0.5 \cdot v(\text{END}) = 0$$

$$v(B) = 3 + 0.5 \cdot v(C) = 3.5$$

$$v(A) = 2 + 0.5 \cdot v(B) = 3.75$$

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

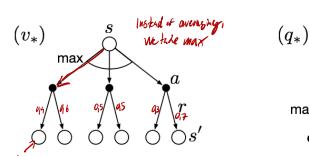


Figure 3.4: Backup diagrams for v_* and q_* .

Central property

the next state (S_{t+1}, R_{t+1}) depends only on (S_t, R_t) not earlier

transition

prop. function $p(s'|r|s, a)$

"prob. of going to s' and seeing reward r after taking a in s "

Decision Making

- consider how action affects not only the next immediate reward, but also the future state you end up in

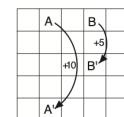
Agent vs Environment

- he may fully know env, but still struggle (Robots)

Rewards

- reward signal is way of communicating to agent what you want it to achieve and how you want it achieved.

Ex. GridWorld



policy $\gamma = 0.9$

$R = 0$

	A	B	A'	B'	
3,3	8,8	4,4	5,3	1,5	
1,5	3,0	2,2	1,9	0,5	
0,1	0,7	0,7	0,-0.4	-0.4	
-1,0	-0,4	-0,4	0,6	1,2	
-1,9	-1,3	-1,2	1,4	2,0	

$$\text{RHS} = 0.25(0.9 \cdot 0.9 + 0.4 \cdot 2.3 + 0.9 \cdot 0.4) - 0.4 \cdot 0.1 = 0.625 - 0.04 = 0.585$$

Adding constants to immediate rewards in Continuing

$$R'_t = R_t + c$$

\vdots

$$R'_T(s) = R_T(s) + \frac{c}{1-\gamma}, \quad \forall s$$

- by adding c , shifts over state value by $\frac{c}{1-\gamma}$

Adding c to All rewards in Episodic

$$G_t = G_t + c \cdot (T-t)$$

No longer in harmless shift, depends on the length of episodes.

Bellman equations in real problems?

We'd need to know: - exact model of the world

- unlimited compute
- perfect Markov state

Instead, we approximate! (to be continued)

- accept that somewhere it might be suboptimal!

(4) Dynamic Programming

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

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

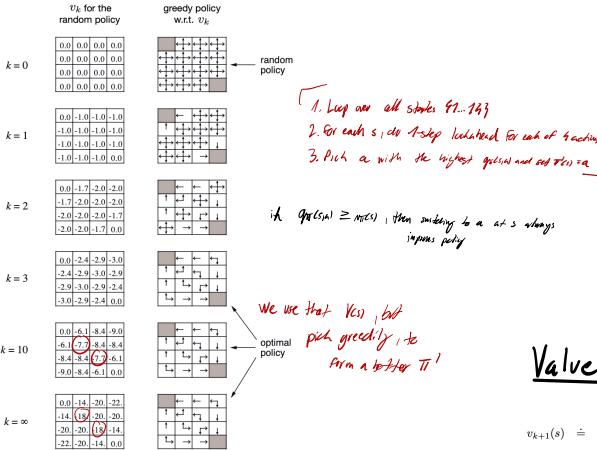
$$v_*(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

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

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

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

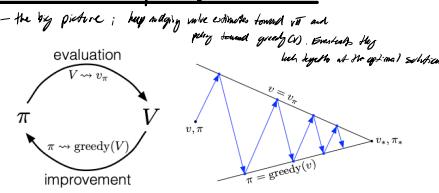
Policy Improvement



Asynchronous DP

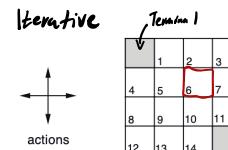
- instead of big sweep over every state, we only pick certain ones, that are 'hot'.

Generalized policy iteration



$\pi_* \longleftrightarrow V_*$

Policy evaluation (Prediction)



$R_t = -1$
on all transitions

exact B. eq.

$$V_E(6) = \frac{1}{\gamma} [(-1 + V_E(1)) + (-1 + V_E(2)) + (-1 + V_E(3)) + (-1 + V_E(7))]$$

but for iterative:

$$V_{k+1}(6) = \frac{1}{\gamma} \sum_{s' \in S, a \in A} (-1 + V_k(s'))$$

Iterations:

$$1. V_{k+1}(6) = 0, \quad 1. V_k(6) = \frac{1}{\gamma} (C_1 + C_2 + \dots + C_6) = -1$$

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

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

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

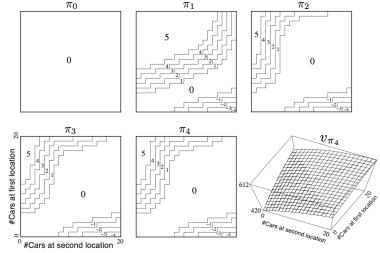
exact number.

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 updates
 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, ..., max(5-s, s)
Carried funds = 0.4 (in heads we go up 0.4, in tails -0.4)
Tails = 3 - a

$\gamma = 1$

$$\begin{aligned} 1. V_0(0) &= 0 \\ 2. \forall j &= 1 \dots 4: V_1(j) \\ Q(j,0) &= 0.4(V_1(j+1) + 0.6(V_1(j))) = 0 \Rightarrow V_1(0) = 0 \\ Q(j,1) &= 0.4(V_1(j) + 0.6(V_1(j+1))) = 0.4 \\ V_1(1) &= \max(0, 0.4) = 0.4 \\ V_1(2) &= 0.4 \\ V_1(3) &= 0.4 \end{aligned}$$

$$2. \forall j < 5: V_2(j) = \dots = 0.76$$

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

3...

slowly converges to true V^* \rightarrow $\pi_{k+1} = \arg \max \pi_{k+1}$

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

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 linear model and apparently enormous MDP

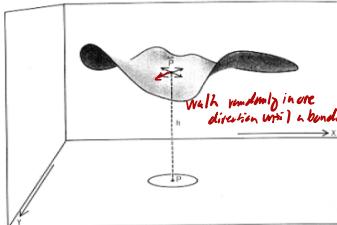
5 Monte Carlo Methods

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.



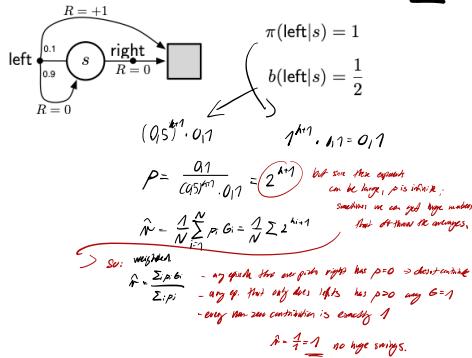
A bubble on a wire loop.

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.



So: weighted $\hat{P} = \frac{\sum p_i}{\sum p_i} G_i$

- avg prob that our first right has p=0 > don't have to sample
- avg p. that only does left has p=0 any G=1
- avg non-zero contribution is exactly 1

$$\hat{P} = \frac{1}{2} = 1/2$$

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 empty new 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:

$$G_n \leftarrow \frac{W_n}{C_n} V_n, V_n$$

When we get next return, we update:

$$V_{n+1} \leftarrow V_n + \frac{W_n}{C_n} [G_n - V_n], n \geq 1,$$

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 by chance, but loss likely. Should have taken those actions.
 3. incrementally update R and make greedy wrt new Q
 4. but stop early if each of when you took an action it wouldn't have taken.

Exploring Starts (MC-ES): force random (s_0, a_0) at episode start.

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.

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

Exploration Guarantees

- 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.

Monte Carlo estimation of Action Values (Q_{MCs})

- same as MCs but with (s, a) pairs

↳ for each episode, first time you saw (s_t, a_t) , return is: $G_t = R_{t+1} + \gamma R_{t+2} + \dots + R_T$

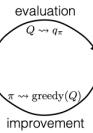
↳ append b to list of returns $\hat{G}_{(s,a)} = \frac{1}{N_{(s,a)}} \sum_{i=1}^{N_{(s,a)}} G^{(i)}$

↳ average over them.

Number of first visits so far

Monte Carlo Control

1. init $G(s,a), N(s,a)=0$, policy π
2. begin each episode in mode chosen (s,a) with prop.
3. $(s,a) \rightarrow r_1 \rightarrow (s_1, a_1) \rightarrow \dots \rightarrow$ terminal
4. for each First-visit (s,a) compute return G_t and update $N_{(s,a)} = 1, Q_{(s,a)} = \frac{1}{N_{(s,a)}} [G_t - Q_{(s,a)}]$
5. for every visited state s_t :
 $T(s_t) \leftarrow \text{argmax}_a Q_{(s,a)}$
6. repeat; $Q \rightarrow q, \pi \rightarrow \pi$



Off-policy prediction via Importance Sampling

Goal: learn about π we really care about, while under more exploratory b

Importance Sampling

coverage condition: $T(s,a) > 0 \Rightarrow b(s,a) > 0 \quad b$ must sometimes pick any action π might pick

importance sampling ratio: $P_{b,T-1} = \prod_{k=1}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)}$ - if π picks actions twice as much as b , the ratio = 2 ...

Ordinary IS estimator

$$V(s) \doteq \frac{\sum_{t \in T(s)} \rho_{t:T(t)-1} G_t}{|\mathcal{T}(s)|}.$$

full return from first visit
all time steps in visited state s

(for zero bias, high variance)

Weighted IS Estimates

$$V(s) \doteq \frac{\sum_{t \in T(s)} \rho_{t:T(t)-1} G_t}{\sum_{t \in T(s)} \rho_{t:T(t)-1}},$$

low variance, small bias that vanishes with more data

- so we re-weight each return G_t by ρ

(bias vanishes with more data)

*Discounting-aware IS

$$V(s) \doteq \frac{\sum_{t \in T(s)} (1-\gamma)^{T(t)-1} \rho_{t:T(t)-1} b_{t:T(t)-1} G_{t:T(t)}}{|\mathcal{T}(s)|}, \quad (5.9)$$

and a weighted importance-sampling estimator, analogous to (5.8), w/

$$V(s) \doteq \frac{\sum_{t \in T(s)} ((1-\gamma)^{T(t)-1} \rho_{t:T(t)-1} b_{t:T(t)-1} G_{t:T(t)})}{\sum_{t \in T(s)} ((1-\gamma)^{T(t)-1} \rho_{t:T(t)-1} b_{t:T(t)-1})}, \quad (5.10)$$

*Per-decision IS

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

$$\text{estimator: } \hat{G}_t = \frac{1}{|\mathcal{T}(s)|} \sum_{i=1}^{|\mathcal{T}(s)|} \tilde{G}_i$$

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

sampled ratios

$$\rho_{t:T-1} = \frac{1}{|\mathcal{T}(s)|} \sum_{i=1}^{|\mathcal{T}(s)|} \frac{\pi(A_i | S_i)}{b(A_i | S_i)}$$

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

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

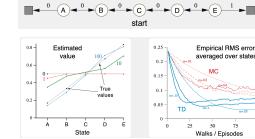
Alg.
Input: the policy π to be evaluated
Algorithm parameter: step-size $\alpha \in [0, 1]$
Initialize $V(s)$, for all $s \in S^*$, arbitrarily except that $V(\text{terminal}) = 0$
Loop for each episode:
 Loop for each step of episode:
 $A \sim \pi(s)$
 Take action A , observe R, S'
 $V(S') \leftarrow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

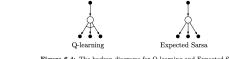


Figure 6.4: The backup diagram for Q-learning and Expected Sarsa.

+ 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

Sarsa: ON-policy TD Control

- learn action value estimates and update policy at the same time using:

$$S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1} \quad \text{it sees the next action of } \text{Exploring} \quad \text{it picks}$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]. \quad \text{TD error}$$

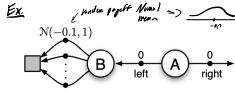
Horizon: - evaluate the current policy via the Sarsa update

- repeat by moving to next greedy, and do update of

- repeat until it's good

- it's like Exploring and gradually reduces to 0, Sarsa will converge to optimal Q^*

Maximization Bias and Double Learning



- Q-learning picks left far more than expected Exploring Plan, it thinks since action B is better than O. \rightarrow This is another bias

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

On each step: - ADP pair

$$\rightarrow 1 \text{ pair: } Q_1(s, a) \leftarrow Q_1(s, a) + \alpha[r + \gamma Q_2(s', \arg \max_b Q_2(s', b)) - Q_1(s, a)],$$

$$\rightarrow 2 \text{ pairs: } Q_2(s, a) \leftarrow Q_2(s, a) + \alpha[r + \gamma Q_1(s', \arg \max_b Q_1(s', b)) - Q_2(s, a)],$$

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

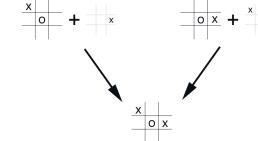
Games, Action States, special cases.

- 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
- Cheaper backups

$$\rightarrow \text{first: } V_{\text{old}}(s) \leftarrow V_{\text{new}}(s) + \alpha[r + \gamma V_{\text{new}}(s') - V_{\text{old}}(s)]$$



7 n-step Bootstrapping

- lets you see n steps future rewards before bootstrapping

n-step TD Prediction

n-step Return

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

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

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

update rule

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

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

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

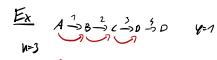
- small $n \rightarrow$ low variance, higher bias

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

- often looking for n in between, which learns faster than both.

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

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



$$G_{t,3} = R_{t+1} + \gamma R_{t+2} + \gamma V(S_{t+3}) = 7+2+3+\gamma=6$$

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

$$\text{Then: } (3+3+\gamma+1) \rightarrow 7+2+0 \rightarrow y_0 \rightarrow \text{update } V(S_t=6)$$

$$G_{t,4} = R_{t+1} + \gamma R_{t+2} + \gamma V(S_{t+4}) = 2+3+5+\gamma=9$$

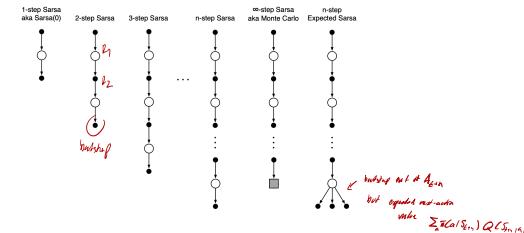
$$V(B) \leftarrow V(B) + \alpha [G_{t,4} - V(B)]$$

$$\text{Then: } (4+3+\gamma+1) \rightarrow 9+0 \rightarrow \text{update } V(C)$$

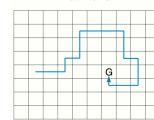
$$G_{t,5} = R_{t+1} + \gamma R_{t+2} + \gamma V(S_{t+5}) = 3+4+\gamma$$

↑ We run out of rewards, we stop the recycling

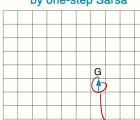
n-step Sarsa



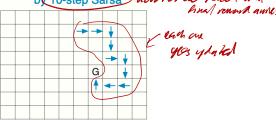
Path taken



Action values increased by one-step Sarsa



Action values increased by n-step Sarsa



per-decision Methods with Control Variates

- we can use IS if only where $\rho=0$, don't visit final return to 0, just skip that update

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

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

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

off-policy Learning without IS:

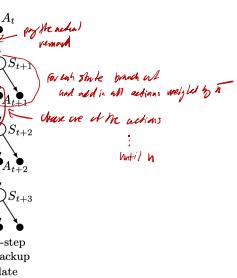
n-step tree backup alg.

return

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

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 $Q(\sigma)$, 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 $1-\sigma$ vs the tree
→ hybrid, bias/variance

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

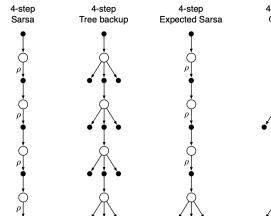
$$G_{t,h} = R_{t+1} + \gamma \left(\sigma_{t+1} \rho_{t+1} + (1 - \sigma_{t+1}) \pi(A_{t+1} | S_{t+1}) \right) (G_{t+1,h} - Q_{t-1}(S_{t+1}, A_{t+1})) + \gamma \bar{V}_{h-1}(S_{t+1})$$

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.

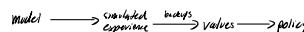
- outputs prediction + R and S'
- if it returns full probability dist over state $Q(S, A)$ \rightarrow distribution model
- if single sample $Q(S, A)$, according to state probabilities \rightarrow sample model

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



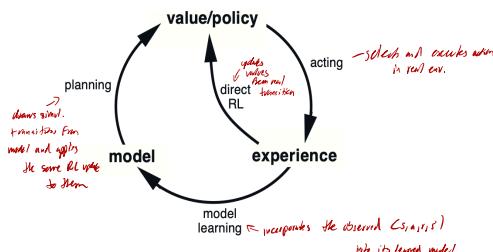
Random-sample one-step tabular Q-learning

```

Loop forever:
1. Select state  $S$  at random
2. Send  $S, A$  to sample model, and obtain  $R, S'$ 
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



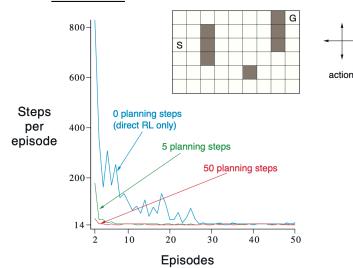
Tabular Dyna-Q

```

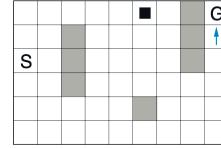
Initialize  $Q(s, a)$  and  $Model(s, a)$  for all  $s \in S$  and  $a \in A(s)$ 
Loop forever:
(a)  $S \leftarrow$  current (nonterminal) state
(b)  $A \leftarrow \epsilon\text{-greedy}(S, Q)$ 
(c) Take action  $A$ ; observe resultant reward,  $R$ , and state,  $S'$ 
(d)  $Q(S, A) \leftarrow Q(S, A)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

```

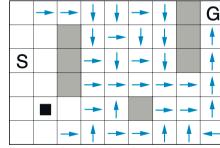
Ex. maze



WITHOUT PLANNING ($n=0$)



WITH PLANNING ($n=50$)

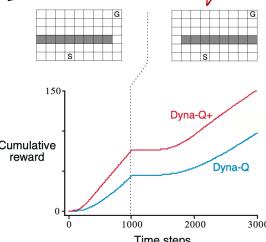


It's a lucky tuple of past outcomes - 'memory'

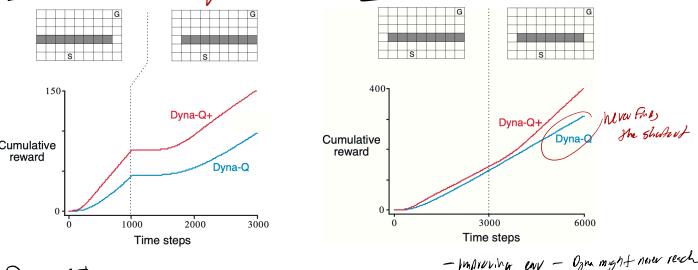
When the Model is Wrong

- 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

Complete exact Bellman backup by summing over all possible next events

Distribution model (either Q or value func):

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

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

+ could be expensive (O(b))

(classic TD(0) backup) does single $Q(s, a)$ from sample model or real env and do:

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

- each update costs constant time, but reduces sampling noise

+ brings back cheap O(1)

Quantifying the tradeoff

$\sqrt{\frac{b-1}{b}}$

Probability of error in value function

Number of type $(s, a) \rightarrow s'$ computations

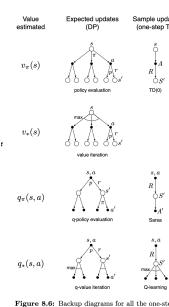


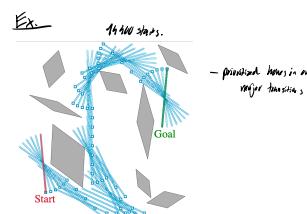
Figure 8.8: Backup diagrams for all the one-step updates considered in this book.

Trajectory Sampling

2 ways of distributing planning updates

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

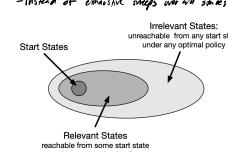
2. Trajectory sampling - instead of blind sweep; instead, draw episodes from model following policy



- prioritized hours in major basins

Real-time Dynamic Programming

- instead of exhaustive sweeps over all states, RTDP updates only the relevant states



Rollout Algorithms

Decision time planning method that uses Monte Carlo control on model to choose action



1. right \rightarrow cost 1, $r = 1$
2. middle child of 1 \rightarrow cost 1, $r = 1$
 $\Rightarrow 6^{11} = 2$

roll 2
1. \rightarrow cost 1
2. \rightarrow cost 2
 $\Rightarrow 6^{11} = 2$

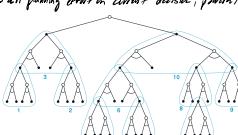
avg = 2

$\frac{-2 - 2}{2} = -2$
we gonna chose the best one.

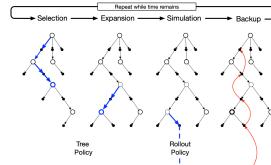
Heuristic Search

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
- focus all planning effort on current decision; prioritized when longer lookahead for single choice.



Selection \rightarrow Expansion \rightarrow Simulation \rightarrow Backup



it's all just value + backups + exp

- how big, how deep, what's policy, model lookahead...

Monte Carlo tree search

Dynamic programming

Exhaustive search

Monte Carlo control

Temporal difference learning

Monte Carlo control

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

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

There choices are:

- linear functions

- neural networks

- decision trees, kernel methods.

Generalization vs. bias

↳ update not only change $V(s)$

↳ update all s' changes $\delta V(s')$ for all s' , which requires

Explicit in the same weight - one example "represents" whole region

of state space

(while standard supervised learning - don't care representability. See the Fig.)

- happens in open loop change over time

- applies rewards and γ inputs online fitting underlying targets

The prediction objective (VE)

- change to parameter vector w across many states, we must choose the most important? ones.

Choosing state-weighting distribution $\mu(s) \geq 0$ and $\sum_s \mu(s) = 1$ to express how much we care about a in each state