

PROOF: WE USE THE POLICY IMPROVEMENT THEOREM

(CHAPTER 4). LET π' BE THE ϵ -GIBBY, AND π THE ϵ -SOFT POLICIES.

$$q_{\pi'}(s, \pi'(s)) = \sum_{a'} \pi'(a'(s)) q_{\pi}(s, a') =$$

$$\frac{\epsilon}{|A(s)|} \sum_{a'} q_{\pi}(s, a') + (1-\epsilon) \max_{a'} q_{\pi}(s, a')$$

$$\geq \frac{\epsilon}{|A(s)|} \sum_{a'} q_{\pi}(s, a') + (1-\epsilon) \sum_{a'} \left(\frac{\pi(a'(s)) - \frac{\epsilon}{|A(s)|}}{1-\epsilon} \right) q_{\pi}(s, a')$$

(EQUALLY HELDS WHEN $\pi(a')$ IS CHOSEN AS IN (A), pg. 40)

COEFFICIENTS ≥ 0 AND SUM = 1

$$= \frac{\epsilon}{|A(s)|} \sum_{a'} q_{\pi}(s, a') - \frac{\epsilon}{|A(s)|} \sum_{a'} q_{\pi}(s, a') + \sum_{a'} \pi(a'(s)) q_{\pi}(s, a') = v_{\pi}(s)$$

$$\Rightarrow q_{\pi'}(s, \pi'(s)) \stackrel{(A)}{\geq} v_{\pi}(s)$$

THEOREM 2: THE POLICY CALCULATED BY THE ALGORITHM CAN ONLY CONVERGE TO A POLICY THAT IS OPTIMAL AMONG ALL POLICIES THAT ARE ϵ -SOFT

PROOF: CONSIDER A NEW ENVIRONMENT THAT BEHAVES RANDOMLY; WITH PROB ϵ , IT CHOOSES A RANDOM ACTION, UNIFORMLY. OTHERWISE, IT FOLLOWS SPECIFIED ACTION ACCORDING TO POLICY.

THEN OPTIMAL POLICY IN NEW ENVIRONMENT IS SAME AS BEST AMONG ϵ -SOFT POLICIES IN OLD ENVIRONMENT.

CONSIDER BELLMAN OPTIMALITY FOR NEW ENVIRONMENT:

$$\begin{aligned}
V_{\pi}(s) &= \max_{\pi} E [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = \pi] \\
&= \max_{\pi} \left\{ (1-\epsilon) E [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = \pi, \text{NON-RANDOM}] \right. \\
&\quad \left. + \epsilon E [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = \pi, \text{RANDOM}] \right\} \\
&= (1-\epsilon) \max_{\pi} E [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = \pi, \text{NON-RANDOM}]
\end{aligned}$$

$$+ \frac{\epsilon}{|A(s)|} \sum_{\pi} \sum_{s',r} p(s',r | s, \pi) [r + \gamma \tilde{V}_{\pi}(s')] \Rightarrow$$

$$V_{\pi}(s) = (1-\epsilon) \max_{\pi} \sum_{s',r} p(s',r | s, \pi) [r + \gamma \tilde{V}_{\pi}(s')] + \frac{\epsilon}{|A(s)|} \sum_{\pi} \sum_{s',r} p(s',r | s, \pi) [r + \gamma \tilde{V}_{\pi}(s')]$$

} BELLMAN OPTIMALITY FOR NEW ENV.

NOW ASSUME THAT ALGORITHM CONVERGES TO POLICY π . THEN:

$$\begin{aligned}
V_{\pi}(s) &= \sum_{\pi} \pi(\pi(s)) q(s, \pi) \\
&= \sum_{\pi} \frac{\epsilon}{|A(s)|} q(s, \pi) + (1-\epsilon) \max_{\pi} q(s, \pi)
\end{aligned}$$

$$\begin{aligned}
&= (1-\epsilon) \max_{\pi} \sum_{s',r} p(s',r | s, \pi) [r + \gamma V_{\pi}(s')] \\
&\quad + \frac{\epsilon}{|A(s)|} \sum_{\pi} \sum_{s',r} p(s',r | s, \pi) [r + \gamma V_{\pi}(s')]
\end{aligned}$$

SO $V_{\pi}(s)$ ALSO SATISFIES BELLMAN OPTIMALITY FOR NEW ENV. SO IT IS OPTIMAL AMONG ϵ -SOFT POLICIES FOR OLD ENV.

THEOREM 1: ALGORITHM SLOWLY IMPROVES

THEOREM 2: ALGORITHM STOPS WHEN OPTIMALITY IS ACHIEVED. NB: PROOFS NOT STRICT

5.5 OFF-POLICY PREDICTION VIA IMPORTANCE SAMPLING

SAMPLING

UNTIL NOW: SINGLE POLICY π USED BOTH FOR EXPLORING AND APPROXIMATING OPTIMAL \Rightarrow CONFLICTING GOALS.

NOW: TWO POLICIES: TARGET POLICY π } OFF-POLICY
 BEHAVIOR POLICY b } METHODS

ADVANTAGE: MORE GENERAL

DISADVANTAGE: SLOWER

FIRST, WE FOCUS ON PREDICTION OF J USING b

REQUIREMENT: COVERAGE

$$J(\pi(s)) > 0 \Rightarrow b(s)$$

INDEED, WE NEED TO ENCOURAGE WHAT WE MEASURE

WE DEFINE

IMPORTANCE-SAMPLING RATIO $P_{t:T-1}$ AS FOLLOWS

PROBABILITY OF STATE-ACTION TRAJECTORY

$$P_{t:T-1} \{ A_t, S_{t+1}, A_{t+1}, \dots, A_{T-1}, S_T \mid S_t, A_{t:T-1} \sim \pi \}$$

$$= \pi(A_t | S_t) p(S_{t+1} | S_t, A_t) \pi(A_{t+1} | S_{t+1}) \dots p(S_T | S_{T-1}, A_{T-1})$$

$$= \prod_{k=t}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)$$

LIKELIHOOD FOR POLICY b .

HOW MUCH MORE PROBABLE TO SEE THIS TRAJECTORY WHEN USING π INSTEAD OF b ?

$$P_{t:T-1} \triangleq \frac{\prod_{k=t}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)}{\prod_{k=t}^{T-1} b(A_k | S_k) p(S_{k+1} | S_k, A_k)} = \frac{\prod_{k=t}^{T-1} \pi(A_k | S_k)}{\prod_{k=t}^{T-1} b(A_k | S_k)}$$

SO THIS RATIO ONLY DEPENDS ON THE POLICIES, NOT THE ENVIRONMENT

OBSERVE THAT $P_{t:T-1}$ IS ALSO A RANDOM VARIABLE BUT DEPENDS ONLY ON TRAJECTORY ITS USEFULNESS COMES FROM THE FOLLOWING

$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s]$$

(NESTED EXPECTATION) $\sum_{\text{ALL TRAJECTORIES THAT LEAD TO T}} P[\text{TRAJ} | \pi] E[G_t | S_t = s, \text{TRAJ}]$

$= \sum_{\text{ALL TRAJECTORIES THAT LEAD TO T}} P_{t:T-1} P[\text{TRAJ} | b] E[G_t | S_t = s, \text{TRAJ}] =$

$\sum_{\text{ALL TRAJECTORIES THAT LEAD TO T}} P[\text{TRAJ} | b] E[P_{t:T-1} G_t | S_t = s, \text{TRAJ}]$

(NESTED EXPECTATION) $= E_b [P_{t:T-1} G_t | S_t = s] \Rightarrow$

$V_{\pi}(s) = E_b [P_{t:T-1} G_t | S_t = s]$

(NOTE: WE ALSO HAVE $V_{\pi}(s) = E_{\pi} [G_t | S_t = s]$
 $V_b(s) = E_b [G_t | S_t = s]$)

WE INTRODUCE NOTATION:

→ LET t BE THE TIME INDEX OF CONSECUTIVE EPISODES FOLLOWING POLICY b . (SO IF AT $t = \infty$ WE END EPISODE, AT TIME $t+1$ WE START NEW EPISODE)

→ LET $\mathcal{T}(s)$ THE SET OF TIMES WE VISIT STATE s (EITHER FOR FIRST TIME, OR ALL TIMES, ACCORDING TO THE VARIANT USED)

→ LET $T(t)$ THE FIRST TIME OF TERMINATION FOLLOWING t .

THEN, WE CAN ESTIMATE $V_{\pi}(s)$ AS FOLLOWS:

$$V(s) = \frac{\sum_{t \in \mathcal{U}(s)} P_{t:T(t)-1} G_t}{| \mathcal{U}(s) |} \quad \textcircled{A}$$

EXPRESSES HOW MANY MORE OR FEWER TIMES WE WOULD HAVE SEEN THAT STATE UNDER π

THIS IS ORDINARY IMPORTANCE SAMPLING

ANOTHER FORMULA IS FREQUENTLY USED, AND IS ACTUALLY PREFERABLE:

$$V(s) = \frac{\sum_{t \in \mathcal{U}(s)} P_{t:T(t)-1} G_t}{\sum_{t \in \mathcal{U}(s)} P_{t:T(t)-1}} \quad \textcircled{B}$$

WEIGHTED IMPORTANCE SAMPLING

BOTH \textcircled{A} AND \textcircled{B} CONVERGE TO $V_{\pi}(s)$, AS $t \rightarrow \infty$

THE CONVERGENCE OF \textcircled{A} FOLLOWS BY THE SLLN FOR THE CONVERGENCE OF \textcircled{B} WE NOTE:

$$\frac{\sum_{t \in \mathcal{U}(s)} P_{t:T(t)-1} G_t}{| \mathcal{U}(s) |} \rightarrow \frac{V_{\pi}(s)}{1} = \frac{\sum_{t \in \mathcal{U}(s)} P_{t:T(t)-1} G_t}{\sum_{t \in \mathcal{U}(s)} P_{t:T(t)-1}}$$

(CONSIDER A MODIFIED ENVIRONMENT WHERE ALL REWARDS ARE ZERO, EXCEPT THE REWARD MOVING TO A TERMINAL STATE, WHICH IS 1)

PROS AND CONS:

(A) IS UNBIASED, I.E. ITS MEAN IS $V_{\pi}(s)$

(A) HAS LARGE VARIANCE, BECAUSE THE RATIOS MIGHT BE LARGE

(B) IS BIASED, I.E. ITS MEAN $\neq V_{\pi}(s)$

(B) HAS MUCH SMALLER VARIANCE, BUT...

OVERALL, (B) IS PREFERABLE

FOR EXAMPLE, COMPARE (A), (B) WHEN $|r(s)| = 1$
(FOR FIRST-VISIT CASE) AND THERE IS A SINGLE RETURN G_t

(A) $V(s) = \sum_{t: T(t)=1}^{\text{SAT } 10} P_t G_t$ WHICH IS UNBIASED, BUT HAS LARGE VARIANCE

(B) $V(s) = G_t$ WHICH IS BIASED, BUT HAS SMALLER VARIANCE

(NOTE: EVERY-VISIT ESTIMATES ARE BOTH BIASED)

5.6 INCREMENTAL IMPLEMENTATION

WE DEVELOP ALGORITHM THAT SIMULATES ACCORDING TO POLICY π AND COMPUTES

$$V(s) = \frac{\sum_{t \in T(s)} P_{t: T(t)=1} G_t}{\sum_{t \in T(s)} P_{t: T(t)=1}}$$

WE HAVE SEQUENCE OF

$$G_1, G_2, \dots, G_{n-1}, \dots$$

$$W_1, W_2, \dots, W_{n-1}, \dots$$

$$(W_i = P_{t_i} = T/t_i) - 1$$

WE NEED TO FORM ESTIMATE

$$V_n \stackrel{\Delta}{=} \frac{\sum_{k=1}^{n-1} W_k G_k}{\sum_{k=1}^{n-1} W_k}, \quad n \geq 2, \quad V_1 = \text{ARBITRARY}$$

WE ALSO DEFINE

$$C_n = \sum_{k=1}^n W_k, \quad C_0 = 0$$

WE NOTE THAT

$$V_{n+1} = \frac{\sum_{k=1}^n W_k G_k}{\sum_{k=1}^n W_k} = \frac{\sum_{k=1}^{n-1} W_k G_k}{C_{n-1}} + \frac{W_n G_n}{C_n}$$

$$= \frac{\sum_{k=1}^{n-1} W_k G_k}{C_{n-1}} + \frac{W_n V_n}{C_n} + \frac{W_n}{C_n} [G_n - V_n]$$

$$= \frac{\sum_{k=1}^{n-1} W_k G_k}{C_{n-1}} \cdot \frac{C_{n-1}}{C_n} + \frac{W_n V_n}{C_n} + \frac{W_n}{C_n} [G_n - V_n]$$

$$= V_n \left[\frac{C_{n-1} + W_n}{C_n} \right] + \frac{W_n}{C_n} [G_n - V_n] \Rightarrow$$

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

THIS LEADS TO FOLLOWING ALGORITHM (Pg. 110)

OFF-POLICY MC PREDICTION FOR ESTIMATION $Q \approx q_{\pi}$

INPUT: TARGET POLICY π

INITIALIZE, $\forall s \in S, w \in A(s)$

$Q(s, a) \in \mathbb{R}$ (ARBITRARY)

$C(s, a) \leftarrow 0$

LOOP FOREVER (FOR EACH EPISODE)

$b \leftarrow$ ANY POLICY WITH COVERAGE OF π

GENERATE AN EPISODE FOLLOWING $b: S_0, A_0, R_0, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

$W \leftarrow 1$

LOOP FOR EACH STEP OF EPISODE, $t = T-1, T-2, \dots, 0$
WHILE $W \neq 0$

$G \leftarrow \gamma G + R_{t+1}$

$C(s_t, A_t) \leftarrow C(s_t, A_t) + W$

$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \frac{W}{C(s_t, A_t)} [G - Q(s_t, A_t)]$

$W \leftarrow W \frac{\pi(A_t | S_t)}{b(A_t | S_t)}$

⊗

THIS CAN BE EASILY MODIFIED TO BECOME A CONTROL ALGORITHM, AS WE SHOW NEXT

⊗ SHOULD PROBABLY BE UP!

INITIALIZE $\forall s \in S$, OR $A(s)$:

$Q(s, a) \in \mathbb{R}$ (ARBITRARY)

$C(s, a) \leftarrow 0$ (SUM OF WEIGHTS)

$\pi(s) \leftarrow \text{ARGMAX}_a Q(s, a)$ (TIES BROKEN CONSISTENTLY)

LOOP FOREVER (FOR EACH EPISODE)

$b \leftarrow$ ANY SOFT POLICY

GENERATE AN EPISODE USING $b: S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}$

$G \leftarrow 0$

$w \leftarrow 1$

LOOP FOR EACH STEP OF EPISODE, $t = T-1, T-2, \dots, 0$:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + w$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{w}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \text{ARGMAX}_a Q(S_t, a) \text{ (TIES BROKEN CONSISTENTLY)}$$

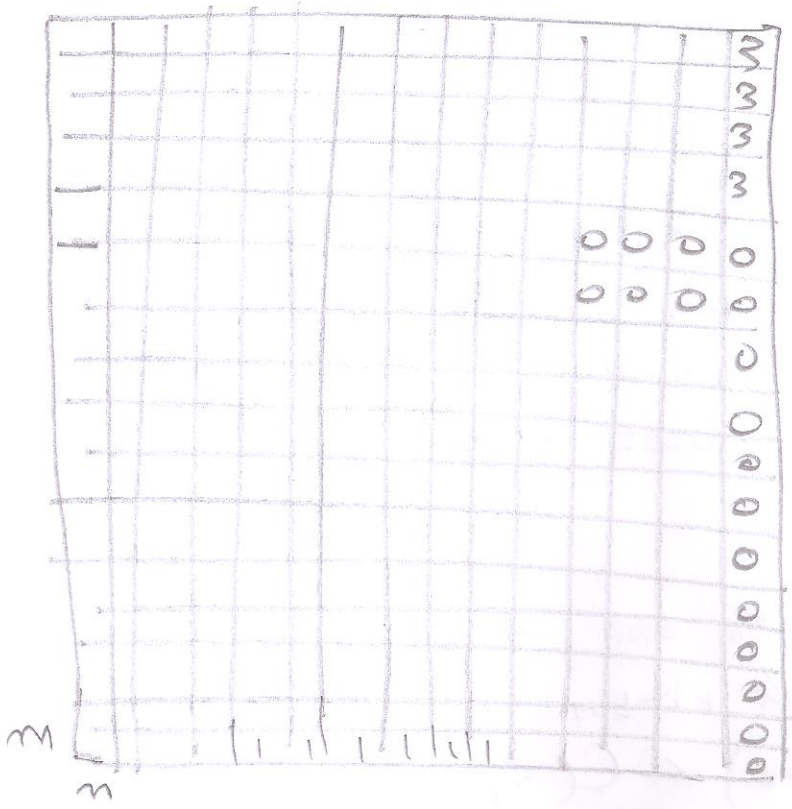
IF $A_t \neq \pi(S_t)$ THEN EXIT INNER LOOP (AND GO TO NEXT EPISODE)

$$w \leftarrow w \cdot \frac{1}{b(A_t | S_t)}$$

⊕ PROBABLY PLACED WRONG IN THE BOOK

ONE PROBLEM: WE NEED TO NOT START THE EPISODE WHEN WE CHANGE THE POLICY

HW #5: SOLVE THE RACE TRACK EXERCISE, S.12



0: NO GO
 1: START
 2: PAD
 3: END

- STATE: LOCATION AND VELOCITY
 - VELOCITY: ENDS LOCATIONS UP/DOWN AND LEFT/RIGHT (MAX +5) ⇒ 11x11 CASES
 - ACTION: CHANGE VELOCITY COMPONENTS +1, -1, 0
 - REWARD: -1 PER STEP NOT FINISHING
 - IF YOU HIT BOUNDARY, YOU GO BACK AT START WITH ZERO VELOCITY
 - WITH PROB 0.1, VELOCITY INCREMENTS ARE ZERO (OIL ON THE TRACK!)
 - USE OFF-POLICY METHOD OF PREVIOUS PAGE
- COMMENT: MOST EFFICIENT WAY TO SOLVE THIS IS WRITE IT AS A ROUTING PROBLEM, ON A GRAPH, WHERE THERE IS A GRAPH NODE FOR EACH STATE BUT OUR METHOD GENERALIZES BETTER