

DYNAMIC PROGRAMMING

ALGO FOR OPTIMAL POLICIES IN ENV PERFECT ENVIRONMENT MODEL AS MDP. ASSUME DISCRETE STATE AND ACTION SPACES.

• ONLINE, ITERATIVE APPROXIMATIONS

• $V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) [R + \gamma V_k(s')]$, $V_k = V_{k+1}$ FIXED POINT. FULL BACKUP OPERATION NEW VALUE: OLD SUCCESSOR STATES + EXPECTED IMMEDIATE REWARDS FOR ALL ONE-STEP TRANSITIONS

• SWEEP THROUGH STATE SPACE. • STOP WHEN $|V_{k+1}(s) - V_k(s)| \leq \epsilon$

• POLICY IMPROVEMENT THEOREM: $Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s) \rightarrow V_{\pi'}(s) \geq V_{\pi}(s)$. 'BETTER TO SELECT A IN S, AND THEN π' , THAN π ALL THE TIME!'
 $\forall s \in S$

- STOCHASTIC POLICIES: $Q_{\pi}(s, \pi'(s)) = \sum \pi'(a|s) Q_{\pi}(s, a)$. IF TIES CAN PARTITION PROBABILITIES AMONG MAXIMAL TIES

• POLICY ITERATION: WE CAN ITERATIVELY OBTAIN BETTER POLICIES. EVALUATION-IMPROVEMENT CYCLE. MDP HAS FINITE POLICIES \rightarrow CONVERGENCE IS FAST
 $\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \dots$ • EACH EVALUATION STEP CAN BE DONE ITERATIVE COMPUTATION \rightarrow WE CAN TRUNCATE BEFORE FULL CONVERGENCE

• VALUE ITERATION: COMBINES ONE E AND ONE I STEP IN EACH PASS. $V_{k+1}(s) = \max_a \sum_{s'} P(s', R|s, a) [R + \gamma V_k(s')]$
 • TURNS BELLMAN OPTIMALITY EQUATION INTO UPDATE RULE • IS ONLY MAX ON POLICY EVAL SWEEP. (STATE SPACE)

• ASYNCHRONOUS DP: IN-PACE ITERATIVE DP ALGOs, DO NOT REQUIRE SYSTEMATIC SWEEP OF STATE SPACE. THEY BACKUP AND USE WHATEVER AVAILABLE VALUES IN ANY ORDER \rightarrow FULL BACKUP REQUIRED FOR CONVERGENCE. SUPER-FLEXIBLE. MIXES POLICY EVAL, ITERATION, VALUE ITERATION, TRUNCATION \rightarrow ALLOW TO RUN THE RL AGENT ONLINE, IN REAL-TIME \rightarrow FOCUS ON THE INTERESTING/RELEVANT PORTION OF STATE SPACE

• GENERALIZED POLICY ITERATION: GENERAL IDEA TO MAKE POLICY EVALUATION AND IMPROVEMENT INTERACT.
 { • VF CONSISTENT WITH CURRENT POLICY $\rightarrow E$ • STABILIZED \rightarrow WE HAVE OPTIMALITY. POLICY IS GREEDY WRT ITS EVAL FOR
 • MAKE π GREEDY WRT CURRENT VF $\rightarrow I$ • CONVERGENCE GUARANTEED FOR MDP
 • 'BOOTSTRAPPING'

• EFFICIENCY: DP IS WAY. FINDS OPTIMAL POLICY IN P-TIME EVEN IF $|S|$ IS 10^5 . \rightarrow WAY DIRECT SEARCH. BREAKS DOWN 2 ORDERS OF MAGNITUDE BETTER THAN LINEAR PROGRAMMING. HANDLES MILLIONS OF STATES

MONTE CARLO METHODS

DO NOT ASSUME COMPLETE KNOWLEDGE OF ENVIRONMENT. GENERATE ONLY SAMPLE TRANSITIONS. NOT COMPLETE DISTRIBUTIONS OF ALL TRANSITIONS. AVERAGE SAMPLE RETURNS. ALL TASKS FORMULATED AS EPISODIC. UPDATES ONLY AT EPISODE TERMINATION. SAMPLE AND AVERAGE RETURNS FOREACH STATE-ACTION PAIR. RETURNS OF ACTION DEPEND ON ACTIONS TAKEN LATER IN SAME EPISODE. NONSTATIONARY WRT EARLIER STATE. DO NOT BOOTSTRAP.

• MC PREDICTION: ESTIMATE $V_{\pi}(s)$ GIVEN SET OF EPISODES BY FOLLOWING π AND PASSING THROUGH s . AVERAGE RETURNS OBSERVED AFTER VISITS TO s .
 FIRST-VISIT VS ALL-VISITS: DISCARD RETURN ON 1ST VISIT OR NOT. CONVERGENCE BECAUSE CLT, FINITE VARIANCES

• MC ESTIMATION OF ACTION VALUES: NO MODEL \rightarrow STATE VALUES ALONE NOT SUFFICIENT: LET'S ESTIMATE Q_{π} : BASICALLY SAME FIRST-VISIT. ALL VISITS ALWAYS
 \rightarrow ISSUE: IF π IS DETERMINISTIC WE ONLY GET RETURNS FOR ONE ACTION FROM EACH STATE. OTHER ESTIMATES WILL NOT IMPROVE WITH EXPERIENCE.
 • MAINTAIN THE EXPLORATION \rightarrow IMPROVE EVERY S-A PAIR AS STARTING PAIR. ON EPISODIC BASIS, **EXPLORE NEW STATES**.
 • ONLY COORDINATE STOCHASTIC POLICIES WITH MULTIPLE ACTIONS IN EACH STATE.

• MC CONTROL: TO APPROXIMATE POLICIES. $\pi \xrightarrow{Q \rightarrow Q_{\pi}} Q$ • E LINE BEFORE
 $\pi \xrightarrow{\pi \rightarrow \text{GREEDY}(Q)} \pi$ • I FOR ANY Q, OPTIMAL π IS ONE CHOOSING ACTION WITH MAXIMAL ACTION-VALUE $\pi(s) = \text{ARGMAX}_a Q(s, a)$

\rightarrow P.I. THEOREM $\rightarrow Q_{\pi}(s, \pi_{k+1}(s)) = Q_{\pi}(s, \text{ARGMAX}_a Q_{\pi}(s, a)) = V_{\pi}(s)$

MC-ES ALGO: RETURNS FOREACH S-A PAIR ARE ACCUMULATED AND AVERAGED REGARDLESS OF ENFORCED π . NOT SURE IF CONVERGES. BUT LIKELY

HOW TO GUARANTEE ACTIONS SELECTED INFINITELY OFTEN? (WITHOUT EXPLORING STATES)
 - ON-POLICY: E/I DECISION MAKING POLICY \rightarrow MAKE POLICIES ϵ -SOFT $\pi(a|s) \geq \epsilon/|A(s)| \forall s, a \rightarrow \epsilon$ -GREEDY $\in \epsilon$ -SOFT. P.I THEOREM GUARANTEES CONVERGENCE POLICY-ITERATION WORKS
 \rightarrow BEST ϵ -SOFT π
 - OFF-POLICY E/I ANOTHER POLICY

• OFF-POLICY, IMPORTANCE SAMPLING

WE HAVE EPISODES GENERATED BY DIFFERENT POLICY μ , $\pi \neq \mu$ IS TARGET POLICY. μ IS BEHAVIOR POLICY

COVERAGE ASSUMPTION: $\pi(a|s) > 0 \rightarrow \mu(a|s) > 0$. ~~usually~~ μ MUST BE STOCHASTIC. IT CAN BE DETERMINISTIC. IT USUALLY GREEDY, μ EXPLORATORY

IMPORTANCE SAMPLING: ESTIMATE SAMPLES FOR A DISTRIBUTION GIVEN SAMPLES FROM ANOTHER. **IS-RATIO:** $\rho_t^T = \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)}$ ONLY DEPENDS ON THE 2 POLICIES! IE G-GOOD

• $V_\pi(s) = V(s) = \frac{\sum_{t \in \tau(s)} \rho_t^T G_t}{|\tau(s)|}$, $\tau(s)$ IS SET OF TIME STEPS WE VISIT s . $T(t)$ IS 1ST TERMINATION TIME AFTER t . G_t IS RETURN UP TO $T(t)$.

→ $\frac{\sum_{t \in \tau(s)} \rho_t^{T(t)} G_t}{\sum_{t \in \tau(s)} \rho_t^{T(t)}}$ **ORDINARY IS** UNBIASED BUT VARIANCE IS EXTREME, UNBOUNDED, MAY NOT CONVERGE IF SCALES RETURNS HAVE DIFFERENT VARIANCES.
WEIGHTED IS: BIASED, VARIANCE LIMITED, FAST TO CONVERGE **OBS!** THIS IS INCREMENTAL, NO EPISODIC 'RESET'

INCREMENTAL

EPISODIC METHODS: SAME AS NON MC ALGOS, BUT AVG RETURNS INSTEAD THAN REWARDS → IF ~~ON-POLICY~~ **ON-POLICY**

IMPLEMENTATION • **OFF-POLICY:** ORDINARY IS → USE STD INCREMENTAL METHODS BUT USE ~~SCALD~~ **SCALD** RETURNS

WEIGHTED IS → $V_{n+1} = V_n + \frac{W_n}{C_n} [G_n - V_n]$, $C_{n+1} = C_n + W_{n+1}$. • $W=1$ → CAN USE FOR ON-POLICY

• OFF-POLICY METHODS MIGHT ONLY LEARN FROM THE TAIL, AFTER 1ST NON-GREEDY ACTION → LEARNING COULD GET SLOW.

• **IS, TRUNCATED RETURNS:** IF HIGH DISCOUNTING, FACTORS OF RETURN SERIES γ^n CO MAINTEN LITTLE BUT INCREASE VARIANCE.

→ THINK OF γ AS DEGREE OF PARTIAL TERMINATION. • $G_t = \gamma^{T-t-1} \bar{G}_t + (1+\gamma) \sum_{h=t+1}^{T-1} \gamma^{h-t-1} \bar{G}_h$ \bar{G}_t IS HORIZON FOR RETURNS

FLAT PARTIAL RETURNS

• **FPR** HAS TO BE SCALED BY SIMILARLY TRUNCATED IS RATIO (SECTION-BASED P. 134 2 ED) → WEIGHTED, ORDINARY ESTIMATOR FORMULAS

$$\bar{G}_t^h = \sum_{i=t}^h R_i$$