

REINFORCEMENT LEARNING

POLICY: STATE \rightarrow ACTION MAPPING

REWARD FUNCTION: AGENT STATE (ENV. STATE, ACTION) \rightarrow REWARD MAPPING, SHORT TERM

VALUE FUNCTION: FOR A STATE, AMOUNT OF REWARD EXPECTED TO ACCUMULATE OVER FUTURE, LONG TERM PROSPECT

MODEL: MIMICS ENVIRONMENT BEHAVIOR (STATE, ACTION) \rightarrow NEXT STATE, REWARD PRODUCTION

SIMPLE TO RULE: $V(s) \leftarrow V(s) + \alpha [V(s') - V(s)]$, α STEPSIZE IS PROGRESSIVELY ADJUSTED TO AUGMENT $V(s)$ AFTER GREEDY MOVE.

- EVALUATIVE VS INSTRUCTIVE FEEDBACK \rightarrow NO WAY OF EXPLICITLY KNOWING 'OPTIMAL ACTION' AND DO INFERENCE ABOUT CORRECT ACTIONS
- EXPLOITATION VS EXPLORATION, HOW TO BALANCE.

ACTION-VALUE METHODS:

SAMPLE-AVERAGE: $\hat{Q}_t(a) = \frac{1}{K_a} \sum_{k=1}^K R_{k,a}$ SAMPLE AVG OF REWARDS WHEN ACTION SELECTED $\rightarrow Q_T(a^*) = \max_a \hat{Q}_t(a)$

ϵ -GREEDY: $Q_T(a^*) = (1-\epsilon) \max_a \hat{Q}_t(a) + \epsilon [\text{ANY } a]$, ϵ -GREEDY EXPLORES.

- MORE VARIANCE IN REWARDS (THEY ARE R.V.) \rightarrow MORE EXPLORATION IS GOOD
- HETEROSCEDASTIC REWARDS \rightarrow EXPLOITATION PAYS OFF.

SOFTMAX SELECTION: UNIFORM SELECTION FOR NON-GREEDY MOVES IS BAD, $P(a_i) = \frac{e^{Q_T(a_i)/\tau}}{\sum_n e^{Q_T(a_n)/\tau}}$ τ TEMPERATURE $\tau \rightarrow 0$ FULLY GREEDY $\tau \rightarrow \infty$ FULLY EQUIPROBABLE

BINARY BANDIT TASK: 2 POSSIBLE ACTIONS, EACH HAS EQUAL SUCCESS/FAILURE PROBABILITY. NOT SUM TO 1, INFER IF FAILURE \rightarrow OTHER WAS CORRECT

- LRI ALGO: BANDIT INCREASES PROB OF ACTION BEING BEST, ALWAYS UPDATES
- LRI ALGO: LRP BUT ONLY UPDATES WHEN SUCCESS

FOR NONSTATIONARY PROBLEMS: DO NOT WEIGHT SAME. DO WEIGHTED/EXPONENTIAL MOVING AVG $Q_n = (1-\alpha)Q_0 + \sum_1^n \alpha(1-\alpha)^{n-1} R_i$. SUM OF WEIGHTS IS 1

- α MAY BE VARIED FROM STEP TO STEP

HOW TO SET INITIAL ESTIMATES $Q_0(a)$?

THE MORE OPTIMISTIC \rightarrow THE MORE WE ARE FAVORING EXPLORATION IN THE BEGINNING

REINFORCEMENT COMPARISON:

- ASSESS REWARD BY COMPARISON W/REFERENCE REWARD, IE AVG OF PREVIOUSLY RECEIVED REWARDS, MORE EFFECTIVE THAN ACTION-VALUE. PREFERRED TO MOTOR-CORRECTION
- HAVE PREFERENCE VALUES FOR ACTIONS $\pi_t(a) = P(a_t = a) = \frac{e^{P_t(a)}}{\sum_n e^{P_t(n)}}$ PREFERENCE UPDATES $P_{t+1}(a) = P_t(a_t) + \beta [R_t - \bar{R}_t]$
- REFERENCE UPDATE: $\bar{R}_{t+1} = \bar{R}_t + \alpha [R_t - \bar{R}_t]$

PURSUIT:

HAVE BOTH A-V ESTIMATES AND ACTION PREFERENCES • PREFERENCE 'PURSUES' GREEDY A-V ESTIMATED ACTION

a_{t+1}^* GREEDY ACTION, $\begin{cases} \pi_{t+1}(a_{t+1}^*) = \pi_t(a_{t+1}^*) + \beta [1 - \pi_t(a_{t+1}^*)] \\ \pi_{t+1}(a) = \pi_t(a) + \beta [0 - \pi_t(a)] \end{cases}$ GREEDY ACTION UPDATED INCREASINGLY OTHER ACTIONS DECREASED.

- ACTION-VALUES $Q_{t+1}(a)$ UPDATED IE VIA SAMPLE AVG

ASSOCIATIVE TASK: IE MULTIPLE N-ARMED BANDIT TASK. TASK CHANGES RANDOMLY AT EACH PLAY. WE COULD TREAT IT AS A SINGLE NONSTATIONARY TASK BUT NOT VERY WELL. IF WE ARE GIVEN INFO ON WHAT TASK AT EACH PLAY, ASSOCIATIVE TASK. SEARCH FOR BEST ACTIONS AS WELL AS ASSOCIATING THEM WITH SITUATION/TASK