

TD - LEARNING

- NO ENV MODEL, LEARN ON ESTIMATES W/O WAITING FOR FINAL STATE. 'MC + DP' → BOOTSTRAP AND SAMPLING | TD AND MC ARE **SAMPLE BACKUPS** (VS FULL BACKUPS) THEY LOOK AT SAMPLE SUCCESSOR.
- $TD(0) = V(s_t) \leftarrow V(s_t) + \alpha [R_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$, UPDATES AT NEXT TIME-STEP
- ERROR PROPORTIONAL TO CHANGE OVER TIME OF PREDICTIONS.
- PROVED TO CONVERGE W/ STATIC SMALL AND STOCHASTICALLY DECREASING α
- **ADVANTAGES**
 - NO MODEL
 - NATURALLY ONLINE, FULLY INCREMENTAL
 - LESS PRONE TO ISSUES DUE TO DISCOUNTING | TRUNCATING
 - EMPIRICALLY FASTER CONVERGENCE THAN MC
 - TD APPROXIMATES MSE WITH N MEMORY, VS N^2 OF CLASSIC FORM!
- **ON BATCH UPDATES**
 - WHEN UPDATES ONLY AFTER EVERY NEW BATCH.
 - $TD(0)$ CONVERGES DETERMINISTICALLY TO MLE. COGNOMY - EQUIVALENCE ESTIMATE
 - Q-MC CONVERGE TO MINIMIZED MSE ON TRAINING DATA
 - **NONBATCH**: NO CONVERGENCE FOR EITHER, BUT STEPS IN THAT DIRECTION.

SARSA

- ON-POLICY TD CONTROL • GENERAL GPI FRAMEWORK • S-A PAIR TO S-A PAIR LEARNING, WE ESTIMATE $Q_\pi(s, a)$ FOR π ON ALL (s, a)
- $Q(s, a) \leftarrow Q(s, a) + \alpha [R + \gamma Q(s', a') - Q(s, a)]$ • $Q_\pi \rightleftharpoons \pi$ • THOUGHT TO CONVERGE IF (s, a) VISITED ∞ TIMES AND IT CONVERGES IN THE LIMIT
- π CAN BE ϵ -GREEDY, ϵ -SOFT, ETC...

Q-LEARNING

- OFF-POLICY TD CONTROL • $Q(s, a): Q(s, a) + \alpha [R_{t+1} + \gamma \max_a Q(s', a) - Q(s, a)]$
- Q DIRECTLY APPROXES q^* • CONVERGENCE: ALL PAIRS CONTINUE TO BE UPDATED • MAXIMIZES OVER ALL POSSIBLE ACTIONS a AT NEXT STATE
- IN ONLINE LEARNING SARSA PREFERABLE DUE TO CONSIDERING ACTION SELECTION • ϵ DECREASING OVER TIME → ON CONVERGENCE

X AFTERSTATE VALUE FUNCTION: MODEL STATE AFTER THE ACTION. USEFUL IN IE GAMES, QUEUES. BETTER BECAUSE REDUCE SMOOTH JERKS. GPI OR FULL THERM.

ACTOR-CRITIC METHODS

- EXPANSION OF REINFORCEMENT COMPARISON TO TD AND FULL RL PROBLEM: ON-POLICY • MINIMAL COMPUTATION TO SELECT ACTIONS; CAN USE EXPLICITLY STOCHASTIC POLICIES
- **ACTOR**: SEPARATE STRUCTURE REPRESENTING POLICY IMPROVEMENT OF VALUE FUNCTION, IMPROVES POLICY • GENERALIZED POLICY ITERATION
- **CRITIC**: ESTIMATED VALUE FUNCTION → CRITIQUE IS TD ERROR, DRIVES LEARNING IN BOTH π AND Q ; EVALUATES CURA POLICY

$$- \overset{C}{s_t} = R_{t+1} + \gamma V(s_{t+1}) - V(s_t) \longrightarrow \overset{A}{P}(s, a) \leftarrow P(s, a) + \beta \delta_t, \quad \overset{A}{P}(s, a) \leftarrow P(s, a) + \beta \delta_t [1 - \pi_t(s, a)] \text{ ALTERNATIVELY}$$

R-LEARNING

- OFF-POLICY CONTROL • NO DISCOUNTED EXPERIENCE, NO EPISODES • MAXIMUM REWARD PER TIMESTEP
- VALUE FUN FOR π DEFINED RELATIVE TO AVG EXPECTED REWARD PER TIMESTEP UNDER π • ASSUMES PROCESS ERGODIC AND π IRREDUCIBLE OF STARTING STATE
- $Q^\pi = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=0}^{N-1} E_\pi \{ R_t \}$ • SHORT-TERM TRANSIENT DETERMINES VALUE STATE $\tilde{V}^\pi(s) = \sum_{t=0}^{\infty} E_\pi \{ R_{t+1} - Q^\pi | s_t = s \}$, $Q^\pi(s, a) = \sum_{t=0}^{\infty} E_\pi \{ R_{t+1} - Q^\pi | s_t = s, a_t = a \}$
- ACTS ON A STOCHASTIC
- **ACTOR**: BEHAVIOR POLICY, MAY BE DETERMINISTIC, COBBOS WAS UNIK, OR GRADIENT DESCENT ON SOME FIXED PARAMETRIC POLICY
- **CRITIC**: TARGET POLICY, STOCHASTIC, IE ϵ -GREEDY, NOT REFINED POLICY • SARSA/Q/... FOR LEARNING
- IMPROVES MAY BE DETERMINISTIC, ϵ -GREEDY WAS CURRENT A-V FUNCTION

DOUBLE-Q LEARNING

- **VANILLA Q** OVERESTIMATES VALUES BECAUSE USE SAME VALUES TO SELECT AND EVALUATE ACTIONS
- **IDEA**: DECOUPLE SELECTION FROM EVALUATION → HAVE 2 VALUE FUNCTIONS Q, Q'
- ASSIGN EXPERIENCE RANDOMLY • FOR EACH UPDATE USE ONE FOR GREEDY POLICY, OTHER FOR VALUE
- $Q = R_{t+1} + \gamma Q(s_{t+1}, \arg \max_a Q(s_{t+1}, a, Q), Q)$ → $Q^{DQ} = R_{t+1} + \gamma Q(s_{t+1}, \arg \max_a Q(s_{t+1}, a, Q), Q')$ • Q IS CURRENT ONLINE
- IN DQN → DOUBLE DQN WE HAVE TARGET NET ON CRITICAL NODE → USE THAT AS Q' ! NO NEED TO INTRODUCE OTHER NETWORKS FOR EVAL OF CURRENT POLICY