

# ELIGIBILITY TRACES

- CAN COMBINE WITH MOST TD-METHODS
- FORWARD VIEW:** BRIDGE FROM TD TO MC METHODS, INTERMEDIATE SPECTRUM OF METHODS
- BACKWARD VIEW:** TEMPORAL RECORD OF EVENT OCCURRENCE (STATE VISIT, ACTION TAKING), MARKS THINGS FOR UNDERGOING LEARNING CHANGES  
ONLY MARKERS ARE ASSIGNED REWARD/PUNISHMENT

## M-STEP TD

TD = ONE-STEP AHEAD → M-STEP: N STEPS AHEAD IE 2 SA:  $R_{T+1} + \gamma R_{T+2} + \gamma^2 V_T(S_{T+2})$   $G_T^{T+N}(C) = R_{T+1} + \gamma R_{T+2} + \dots + \gamma^{N-1} R_{T+N} + \gamma^N C$   
MC = ALL STEPS AHEAD TO END, OR

N-STEP BACKUP:  $\Delta(S_T) = \alpha [G_T^{T+N}(V_T(S_{T+N})) - V_T(S_T)]$ , UPDATE:  $V_{T+1}(S) = V_T(S) + \Delta_T(S)$ ; ONLINE; OFFLINE  $V_{T+1}(S) = V_T(S)$   
ERROR-REDUCTION PROPERTY: MUAR N; LOWER ERROR → CONVERGENCE  
 $V_T(S) = V_{T-1}(S) + \sum_0^{T-1} \Delta_T(S)$

- COOL, BUT INCONVENIENT TO IMPLEMENT BECAUSE WAITING.

## FWD VIEW

- IDEA: LET'S AVERAGE DIFFERENT N-STEP RETURNS, ON AS LONG AS WEIGHT SUM TO 1. STILL ERROR REDUCTION ← COMPLEX BACKUPS
- $\lambda$  RETURN:**  $L_T = (1-\lambda) \sum_{n=0}^{\infty} \lambda^n G_T^{T+n}(V_T(S_{T+n}))$ , CONTAINS ALL N-STEP BACKUPS, WITH EXP-DECAYING WEIGHTS. AFTER TERMINAL ALL RETS ARE G-C  
 $\lambda = 1$  SAME AS CONSTANT  $\alpha$  MC RETURN • CAN UPDATE ONLINE OR OFFLINE  
 $\lambda = 0$  SAME AS TD(0) METHOD • FWD VIEW BECAUSE NEVER GO BACK TO PREV STATES  
ALGO:  $\Delta_C(S_T) = \alpha [L_T - V_T(S_T)]$

## BWD VIEW

- ELIGIBILITY TRACE FOR STATES  $\begin{cases} E_T(S) = \gamma \lambda E_{T-1}(S) & S \neq S_T \\ E_T(S) = \gamma \lambda E_{T-1}(S) + 1 & S = S_T \end{cases}$   $\lambda$  = TRACE-DECAY PARAM → ACCUMULATING TRACE
- TD( $\lambda$ ) UPDATES:  $\begin{cases} \delta_T = R_{T+1} + \gamma V_T(S_{T+1}) - V_T(S_T) \\ \Delta V_T(S) = \alpha \delta_T E_T(S) \text{ FOR ALL } S \end{cases}$
- UPDATES PROPORTIONAL TO RECENTLY VISITED STATES, ONLINE OR OFFLINE
- $\begin{cases} \lambda = 0 & \text{TD}(0) \\ \lambda = 1 & \text{CONSTANT } \alpha \text{ MC} \end{cases}$  •  $\lambda = 1, \gamma = 1 \rightarrow$  NO DISCOUNT, NO DELAY ERROR MC • MORE GENERAL FORMULATION OF MC
- BWD BECAUSE PROPAGATES OR UPDATES BACK IN TIME

REPLACEMENT TRACE:  $E_T(S_T) = 1$  DUTCH TRACE:  $E_T(S_T) = (1-\alpha) \gamma \lambda E_{T-1}(S_T) + 1$  → GENERALLY BEST FOR ON-LINE ALSO APPROXIMATING  $\lambda$ -RETURN.

TOTAL EXACT STEP-BY-STEP EQUIVALENCE OF  $\lambda$  BASED FWD AND BWD IMPLEMENTATIONS.

## SARSA( $\lambda$ )

- TRACES FOR S-A PAIRS, ANY VARIANT  $Q_{T+1}(S, a) = Q_T(S, a) + \alpha \delta_T E_T(S, a)$ ,  $\delta_T = R_{T+1} + \gamma Q_T(S_{T+1}, A_{T+1}) - Q_T(S_T, A_T)$  • POLICY IMPROVEMENT AS USUAL, EG Q-GOSSIP

## Q( $\lambda$ )

- LOOKAHEADS ONLY UNTIL FIRST EXPLORATORY ACTION BECAUSE AFTER NO MORE RELATION TO CURRENT POLICY. M-STEP RETURNS ONLY UNTIL THEN.
- TRACES SET TO 0 WHEN EXPLORATORY NON-GREEDY ACTION IS TAKEN. ALWAYS SET CUR (S,A) TRACE TO 1 AFTER ACTION, EVEN WHEN EXPLORATORY.
- $Q_{T+1}(S, a) = Q_T(S, a) + \alpha \delta_T E_T(S, a)$   $\delta_T = R_{T+1} + \gamma \max_{a'} Q_T(S_{T+1}, a') - Q_T(S_T, A_T)$  • MAY BE ONLY LITTLE FASTER THAN 1-STEP Q LEARNING
- BOOSTSTRAPS EVEN WITH  $\lambda = 1$  DUE TO TRACE CUTS AND USING VALUE ESTIMATES VS ACTUAL REWARDS, NOT OPTIMAL → DECOUPLING NEEDED w/ IMPROVABLE SAMPLING?

## TRACES FOR ACTOR-CRITIC

- USE STATE TRACES FOR CRITIC, AND ACTION-STATE TRACES FOR ACTOR. SEPARATE SETS. USE SUITABLE ALGOS FOR THE UPDATES.

## VARIABLE- $\lambda$

IDEA:  $\lambda \rightarrow \lambda_T$ ; VARIES FROM STEP TO STEP IE AS Fcn OF STATE  $\lambda_T = \lambda(S_T)$ . DEPENDING ON CERTAINTY OF ESTIMATE.

## IMPLEMENTATION REMARKS

ON ON SEQUENTIAL PIPELINES BECAUSE TRACES ARE MOSTLY SPREAD ONLY A FEW TIMES MORE TIME THAN NO TRACES.  
IF ANY → ONLY DOUBLE  
MAKES MORE SENSE IN ONLINE THAN IN OFFLINE SETTINGS