

# REVISION, RL

MARCOV = STATE TRANSITIONS DEPEND ONLY ON CURRENT STATE, NOT PAST

STATIONARY = TRANSITION PROBABILITIES DON'T CHANGE OVER TIME

MARCOV REWARD PROCESS = MRP =  $(S, P(\text{TRANSITION PROBABILITIES}), R, \gamma)$

RETURN = TOTAL DISCOUNTED REWARD FROM TIME  $t = \sum_{k=0}^{\infty} \gamma^k V_{t+k+1}$

STATE VALUE FUNCTION,  $V^*(s) = \text{EXPECTED RETURN FROM } s, V = R_s + \gamma \sum P_{ss'} V^*(s')$

DETERMINISTIC POLICY = ONLY ONE ACTION POSSIBLE

MARCOV DECISION PROCESS = MDP =  $(S, A, P_{ss'}, \gamma, R_{ss'})$  <sup>POLICY</sup>  $\pi$

STATE-ACTION VALUE FUNCTION =  $Q^*(s, a) = E[R_t | s, a] \quad V^*(s) = \sum \pi(s, a) Q^*(s, a)$

BELLMAN OPTIMALITY EQUATION  $\Rightarrow V^*(s) = \max_a \sum_{s'} P_{ss'} (R_{ss'} + \gamma V^*(s'))$

DYNAMIC PROGRAMMING = COMPUTE OPTIMAL POLICY GIVEN A MODEL  
- BREAK UP THE PROBLEM INTO SMALLER IDENTICAL PROBLEMS

POLICY ITERATION = EVALUATE POLICY UNTIL CONVERGENCE  $\rightarrow$  UPDATE POLICY  
- MULTIPLE EVALUATIONS PER POLICY

VALUE ITERATION = EVALUATE POLICY ONCE  $\rightarrow$  UPDATE POLICY

- CAN BE VIEWED AS EVALUATING 4/0 POLICY, UPDATE WITH OPTIMAL GREEDY ACTION

SYNCHRONOUS BACKUPS = ALL STATES DO AT ONCE, 2 COPIES OF VAL FUNC. = EXPENSIVE

ASYNCHRONOUS BACKUPS = DO SELECTED STATES, 1 COPY OF VAL FUNC. = CHEAPER

BOOTSTRAPPING = UPDATE ESTIMATES BASED ON OTHER ESTIMATE - LESS VARIANCE, MORE BIAS

MODEL-FREE LEARNING

MONTÉ CARLO - GET FULL TRACE OF EXPERIENCES, THEN UPDATE

- NO BOOTSTRAPPING

- MUST BE EPISODIC WITH TERMINAL STATES

- CAN AVERAGE EVERY VISIT OR FIRST VISIT (PER TRACE)

- DOES NOT EXPLOIT MARKOV PROPERTY  $\Rightarrow$  MORE EFFECTIVE IN NON-MARKOV ENV.

- MCCAN UPDATE AFTER EACH TRACE OR MC CAN BATCH THE TRACES

(CAN AVERAGE ALL VISITS EQUALLY  $V(s) \leftarrow V(s) + \frac{1}{N} (R - V(s))$ )

OR GRADUALLY FORGET OLD EPISODES  $V(s) \leftarrow V(s) + \alpha (R - V(s))$  <sup>GOOD FOR NON-STATIONARY ENV</sup>

TEMPORAL DIFFERENCE - UPDATE AT EACH STEP, BOOTSTRAPPING <sup>LESS VARIANCE, MORE BIAS</sup>

$$V(s_t) \leftarrow V(s_t) + \alpha (\underbrace{r_{t+1} + \gamma V(s_{t+1})}_{\text{TD TARGET}} - \underbrace{V(s_t)}_{\text{TD ERROR}})$$

SOFT POLICY = FINITE PROBABILITY TO EXPLORE ALL STATES

EG.  $\epsilon$ -GREEDY = RANDOM  $\epsilon$  % OF THE TIME, OTHERWISE OPTIMAL

$\epsilon$ -GREEDY IS ONE IF  $\epsilon$  REDUCES TO ZERO  $\epsilon_t = \frac{1}{t}$

GLIE =  $\epsilon$ -GREEDY IN THE LIMIT w/ INFINITE EXPLORATION

- ALL STATE/ACTION PAIRS EXPLORED INFINITELY MANY TIMES

- CONVERGES TO GREEDY POLICY

EXPLORING STARTS = START EPISODES WITH PARDON STATE/ACTION, EXPLORE ALL

ON-POLICY = OPTIMIZE THE POLICY YOU ARE FOLLOWING

OFF-POLICY = OPTIMIZE A POLICY WHILE FOLLOWING A DIFFERENT ONE

$\pi$  = TARGET POLICY

$\pi'$  = BEHAVIOUR POLICY



SARSA = ON-POLICY TD  $Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma Q(s',a) - Q(s,a))$

SARSA CONVERGES TO OPTIMAL IF 1) POLICY,  $\pi$ , IS GLIE

2) ROBBINS-MUNRO STEP SIZES FOR  $\alpha_t$   $\begin{cases} \sum \alpha_t = \infty \\ \sum \alpha_t^2 < \infty, \text{ e.g. } \frac{1}{2k} \end{cases}$

IF SPARSE REWARDS USE SARSA-LAMBDA

= AVERAGE REWARDS OF ALL TRACES, FROM LENGTH 1, 2, 3... TERMINATION

Q-LEARNING = OFF-POLICY TD, EG  $\pi$  IS GREEDY,  $\pi'$  IS  $\epsilon$ -GREEDY

UPDATE EACH STEP  $Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma \max_{a'} Q(s',a') - Q(s,a))$

		FULL BACKUP (DP)	SAMPLE BACKUP (TD)
PROBLEM	BELLMAN EXACT Eg. $V_\pi(s)$	ITERATIVE POLICY EVALUATION $V(s) \leftarrow E[R + \gamma V(s')   s]$	TD LEARNING $V(s) \leftarrow V(s) + \alpha [R + \gamma V(s') - V(s)]$
CONTROL	BELLMAN EXACT Eg. $Q_\pi(s,a)$	CONVERGE VALUE, THEN ADJUST POLICY Q-POLICY ITERATION $Q(s,a) \leftarrow E[R + \gamma Q(s',a')   s,a]$	SARSA $Q(s,a) \leftarrow Q(s,a) + \alpha [R + \gamma Q(s',a) - Q(s,a)]$
CONTROL	BELLMAN OPTIMAL Eg. $Q^*(s,a)$	ALTERNATE - EVAL W/ POLICY IMPROVEMENT Q-VALUE ITERATION $Q(s,a) \leftarrow E[R + \gamma \max_{a'} Q(s',a')   s,a]$	Q-LEARNING $Q(s,a) \leftarrow Q(s,a) + \alpha [R + \gamma \max_{a'} Q(s',a') - Q(s,a)]$

FUNCTION APPROXIMATION - GOOD FOR LARGE OR CONTINUOUS STATE SPACE

$V^*(s), Q^*(s,a) \approx \hat{V}(s, w), \hat{Q}(s,a, w)$  W APPROXIMATES STATE

USE DIFFERENTIABLE APPROX. UPDATE W WITH GRADIENT DESCENT

LINEAR APPROX.  $\hat{V} = X^T w = \sum x_j(s) w_j \Rightarrow \Delta w = \alpha (V^*(s) - \hat{V}(s, w)) X(s) = LR + \epsilon$  (PROBLY FEEDBACK VAL)

COARSE CODING = REPRESENT A STATE W/ OVERLAPPING FEATURES,  $\begin{cases} 1 \text{ IF FEATURE IN STATE} \\ 0 \text{ IF NOT} \end{cases}$

TILE CODING = COARSE CODING WITH STATE SPACE PARTITIONED INTO OVERLAPPING

RADIAL BASIS FUNCTIONS = RBFs = COARSE CODING FOR CONTINUOUS SPACES  
RECTANGULAR TILES  
- EACH STATE IS A LINEAR WEIGHTED COMBO OF THE RBFs

MC FUNCTION APPROX - CONVERGES TO LOCAL OPTIMUM, LINEAR OR NON-LINEAR APPROX

TD FUNCTION APPROX - LINEAR APPROX CONVERGES CLOSE (BIASED) BUT NOT NON-LINEAR APPROX

DEEP RL TRIAD = OFF-POLICY, FUNC. APPROX., BOOTSTRAPPING TD

## DQN FEATURES

EXPERIENCE REPLAY = SAVE EXPERIENCES IN BUFFER FOR BATCHES OF RANDOM SAMPLES

- MORE EFFICIENT USE OF DATA, DATA IS REUSED
- CAN MAKE USE OF OLDER EXPERIENCES, AVOID CATASTROPHIC FORGETTING
- REDUCES CORRELATION BETWEEN SAMPLES IN TRAINING BATCH
- REQUIRES THAT IMMEDIATE OUTCOMES FOR S/A ARE SOMEWHAT STABLE  
LOW VARIANCE, LOW ENTROPY

PROBLEM: UPDATES ARE TOO FAST, TOO NOISY, AND CAN LEAD TO RUNWAY BIAS

SOLN = TARGET NETWORK = WHEN UPDATING THE MAIN NETWORK,  $Q$ ,

CALCULATE THE ERRORS BASED ON TARGET NETWORK,  $Q'$ ,

$Q'$  IS UPDATED INFREQUENTLY TO MATCH  $Q$   $\Rightarrow$  STABILITY

CLIPPING REWARDS - DIFFERENT SCENARIOS HAVE DIFFERENT

SCALE OF REWARDS, SO CLIP THEM ALL  $\{-1, 0, 1\}$

SKIPPING FRAMES - COMPUTER CAN REACT FASTER THAN PEOPLE

- SKIPPING TO EVERY 4<sup>TH</sup> FRAME  $\Rightarrow$  LOWER COMPUTATION COST, FASTER TRAINING
- CAN STACK THE FOUR FRAMES AS A SINGLE INPUT  $\Rightarrow$  MORE INFO ON CHANGES IN CURRENT STATE

PROBLEM DQN NETWORKS OVERESTIMATE BECAUSE THEY TAKE MAX  $Q$

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_a Q'(s, a') - Q(s, a)]$$

$Q' = \text{TARGET NETWORK}$

DDQN - LET THE MAIN NETWORK,  $Q$ , CHOOSE THE ACTION TO QUERY, BUT GET THE VALUE FROM THE TARGET NETWORK,  $Q'$

- DE-CORRELATES THE ERRORS BETWEEN CHOOSING ACTIONS & EVALUATING
- LOWERS THE BIAS, MORE STABLE

## ACTOR / CRITIC METHODS

CRITIC = VALUE BASED = APPROXIMATE VALUE FUNC.  $\Rightarrow$  SAMPLE EFFICIENT, STEADY  
EXAMPLES = Q-LEARNING, DQN, DDQN

ACTOR = POLICY BASED = FIND POLICY WHO COMPUTING  $Q$  OR  $V$   $\Rightarrow$  FASTER CONVERGENCE  
 $\Rightarrow$  BETTER FOR CONTINUOUS AND STOCHASTIC ENVIRONMENTS  
EXAMPLES = POLICY GRADIENTS, REINFORCE



## POLICY GRADIENTS

= LEARN A PARAMETERIZED POLICY BASED ON POLICY WEIGHTS,  $\theta$  OPTIMIZE  $\theta$

OPTIMAL POLICY  $\pi^*(a|s, \theta)$  = PROBABILITY OF ACTION  $a$  IN STATE  $s$

MAXIMIZE OUR PERFORMANCE MEASURE  $J(\theta)$  W/ GRADIENT ASCENT

REINFORCE = 3-STEP POLICY GRADIENT IMPLEMENTATION

- 1) RUN THE POLICY & SAMPLE  $\{s_i, a_i\}$  FROM  $\pi_\theta(a_i|s_i)$
- 2)  $\nabla_\theta J(\theta) = \frac{1}{N} \sum_i \left( \sum_t \nabla_\theta \log \pi_\theta(a_t|s_t) \right) \underbrace{\left( \sum_t r(s_t, a_t) \right)}_{\text{RETURN}}$
- 3)  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

$N$  = NUMBER OF TRACES

- REINFORCE HAS NO BIAS BUT HIGH VARIANCE

GAUSSIAN POLICY DISTRIBUTIONS FOR  $\pi$  = GOOD FOR ROBOTICS, REAL-WORLD APPLICATIONS

PROBLEM: POLICY GRADIENT TAKES MANY ACTIONS DURING AN EPISODE

- HARD TO ASSIGN CREDIT TO SPECIFIC ACTIONS THAT LED TO REWARDS
- SO MAY TAKE MANY UPDATES TO CONVERGE
- POLICY GRAD ONLY VALID IN EPISODIC ENVIRONMENT

SOLVE A2C = ADVANTAGE ACTOR-CRITIC  
= ACTOR IS POLICY BASED, CRITIC IS VALUE BASED

- 1) RUN AND SAMPLE  $\{s_i, a_i\}$  FROM  $\pi_\theta(s_i, a_i)$
- 2) CRITIC FITS  $\hat{V}_\theta(s)$  TO SAMPLED REWARDS
- 3) EVALUATE  $\hat{A}^\pi(s_i, a_i) = r(s_i, a_i) + \hat{V}_\theta(s_{i+1}) - \hat{V}_\theta(s_i)$
- 4)  $\nabla_\theta J(\theta) \approx \sum_i \nabla_\theta \log \pi_\theta(a_i|s_i) \hat{A}^\pi(s_i, a_i)$
- 5) UPDATE POLICY  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

AGENT LEARNS ADVANTAGE,  $A(s, a) = Q(s, a) - V(s) = r + \gamma V(s') - V(s)$

- BOOTSTRAPPING  $\Rightarrow$  LOWER VARIANCE W/O ADDING BIAS
- FASTER POLICY CONVERGENCE
- CAN BE USED IN NON-EPISODIC DOMAINS

A3C = ASYNCHRONOUS ADVANTAGE ACTOR-CRITIC

= MULTIPLE INDEPENDENT AGENTS INTERACTING IN PARALLEL W/ DIFFERENT COPIES OF THE ENVIRONMENT  
 $\Rightarrow$  EXPLORE MORE OF THE STATE-SPACE QUICKER

- EACH AGENT PERIODICALLY UPDATES A GLOBAL NETWORK (NOT SYNCHRONIZED W/ OTHER AGENTS)
- AT EACH UPDATE, ALL AGENTS RESET THEIR PARAMETERS TO THE GLOBAL ONE

## TRICKS OF THE TRADE

FRAME STACKING = STACK MULTIPLE FRAMES AS ONE INPUT TO SEE

POSITION  
VELOCITY  
ACCELERATION

INPUT NORMALIZATION = REDUCE COMPLEXITY OF INPUT (B/W VS COLOR)

FRAME OF REFERENCE = CHANGE COORDINATE SYSTEM TO GLOBALIZE VECTOR

GRADIENT CLIPPING = SCALE THE EUCLIDEAN NORM OF HODDGRAD1 DOW UNDER THRESHOLD

LARGE MINI-BATCH SIZE = BIAS-FREE METHOD OF REDUCING VARIANCE