

Value function under policy  $\pi$  at state  $s$ .

$$v_\pi(s) = \mathbb{E}_{\pi} \left[ Q_\pi | s_t = s \right] \rightarrow \text{solution to a set of equations.}$$

$$= \sum_a \pi(a|s) \sum_{s',r} P(s',r|s,a) [r + v_\pi(s')] : s \in S$$

$$v^*_\pi(s) = \max_\pi \mathbb{E}_\pi Q_\pi(s) \quad \begin{matrix} \nearrow \text{optimal value function} \\ \Downarrow \end{matrix}$$

under optimal policy.

$$Q^*_\pi(s) = \max_a \sum_{s',r} P(s',r|s,a) [r + \gamma v^*_\pi(s')] \xrightarrow{\text{given } \pi \text{ known optimal policy from } s' \text{ state.}}$$

$$q_\pi(s,a) = \sum_{s',r} P(s',r|s,a) [r + \gamma \max_a q_\pi(s',a)] \quad \begin{matrix} \nearrow \text{not in policy.} \\ \Downarrow \end{matrix}$$

Action-value function: value after taking action  $a$ .

### Policy Improvement:

$$V_\pi(s,a) = \sum_{s',r} P(s',r|s,a) [r + \frac{V_\pi(s')}{\pi}]$$

if  $s \in S$ ,  $q_\pi(s, \pi(s)) \geq V_\pi(s)$ ;  $[\pi(s) \text{ and } \pi'(s) \text{ are two different policy}]$   
 then,  $V_\pi(s) > V_{\pi'}(s)$

Greedy rule:  $\pi'(s) = \arg \max_a q_\pi(s,a)$

$$= \arg \max_a \sum_{s',r} P(s',r|s,a) [r + \gamma V_\pi(s')]$$

$$V_{\pi'}(s) = \max_a \sum_{s',r} P(s',r|s,a) [r + \gamma V_\pi(s')] \quad (\text{Improvement})$$

### Policy Iteration:

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \pi_2 \dots \xrightarrow{E} \pi_k \xrightarrow{I} V_{\pi_k}$$

forall state

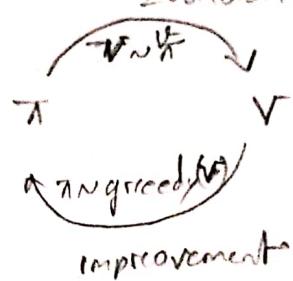
## Value Iteration:

$$\text{Loop} \Rightarrow v_{k+1}(s) = \max_a \sum_{s' \in S} p(s'|s, a) [r + \gamma v_k(s')] \quad \begin{matrix} \text{greedy action} \\ \text{max over } a \end{matrix}$$

$$\text{After loop} \Rightarrow \pi(s) = \arg \max_a \sum_{s' \in S} p(s'|s, a) [r + \gamma v_k(s')]$$

$\{v_k\}$  seq converge  $\rightarrow \lim v_k(s)$  Evaluation

## Generalized Policy Evaluation:



Monte Carlo Prediction  $\rightarrow$  estimate  $V \approx v_\pi$  [episodic]

; DP alternative

### Policy

$$v(s) \in \mathbb{R}$$

return(s)  $\notin$  empty list  $\forall s \in S$ .

loop forever; (for each episode):

generate  $s_0, A_0, R_1, s_1, A_1, R_2, \dots, s_{T-1}, A_{T-1}, R_T$

$$G \leftarrow 0$$

+ in range ( $T-1, T-2, \dots, 0$ )

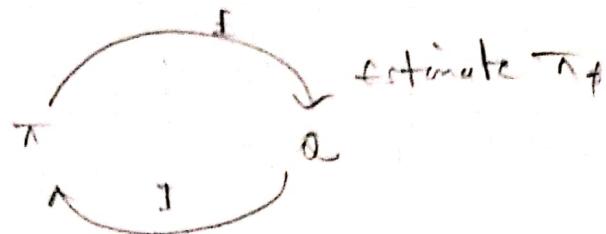
$$G \leftarrow G + \gamma^t R_t + 1$$

unless  $s_t$  appear in  $s_0, \dots, s_{t-1}$ :  
Append a to return  $s_t$

(last  $s_t$  before reach to  $s_0$ )  
(1st  $s_t$  appearance)

$$V(s) \leftarrow \text{avg}(\text{return}(s))$$

## Monte Carlo Control



$$r_0 \xrightarrow{s} q_{\pi_0} \xrightarrow{a} \pi, \xrightarrow{s} q_{\pi_1} \xrightarrow{a} q_{\pi_0} \perp \dots \xrightarrow{\pi_k} \xrightarrow{a} q_{\pi_k}$$

$$\pi(s) = \arg \max_a q_\pi(s, a)$$

$$q_{\pi_k}(s, \pi_{k+1}(s)) = \max_a q_{\pi_k}(s, a)$$

( $\pi_k, \pi_{k+1}$  policies  
may vary in only one state)

MC (Exploring starts) instant control ✓ (also possible  $\epsilon$ -greedy + Alternative for  $\pi(s)$ )

$$\pi(s) \in A(s) \text{ (random) } s \in S$$

$$Q(s, a) \in \mathbb{R} \text{ (random) } a \in A(s)$$

returns  $(s, a) \leftarrow$  empty list  $s \in S, a \in A(s)$

loop forever! (frisidic)

choose  $s_0 \in S, a_0 \in A(s_0)$ ; with prob  $> 0$

generate episode from  $s_0, a_0$  following  $\pi: s_0, a_0, r_1 \dots s_{T-1}, a_{T-1}, r_T$

$$G \leftarrow 0$$

Loop for  $t = 1 \dots T-1, T-2, \dots, 0$

$$G \leftarrow G + r_t P_L$$

unless  $(s_t, a_t)$  pair appears earlier  $(s_t, a_t, \dots, s_{t-1}, a_{t-1})$

Append  $a_t \rightarrow$  Returns  $(s_t, a_t)$

$$Q_t(s_t, a_t) = \text{avg return}(s_t, a_t))$$

$$\pi(s_t) \leftarrow \arg \max_a Q_t(s_t, a) \quad \text{deterministic}$$

inc update  $v(s_t) \leftarrow v(s_t) + \alpha (r_{t+1} + \gamma v(s_{t+1}))$

Temporal difference:  $v(s_t) \leftarrow v(s_t) + \alpha [r_{t+1} + \gamma v(s_{t+1}) - v(s_t)]$

TD(0) estimate  $v$  (method)

(value based))

$\rightarrow \alpha \in [0, 1]$

• initialize  $v(s)$ ,  $\forall s \in S^+$ ,  $v(\text{terminal}) = 0$

loop; for each episode

initialize  $S$

Loop

$a \leftarrow \text{action by } \pi$

take action  $a$ , observe  $r, s'$

$v(s) \leftarrow v(s) + \alpha [r + \gamma v(s') - v(s)]$

$s \leftarrow s'$

until  $s$  is terminal

{ TD estimate:  $\delta_t = [r_{t+1} + \gamma v(s_{t+1}) - v(s_t)]$

$E_t - v(s_t) = \sum_{k=1}^{T-1} \gamma^{k-1} \delta_k$

$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)]$

Some: TD on policy

$\alpha \in [0, 1]$

$Q(s, a)$  initialize,  $s \in S^+$ ,  $a \in A(s)$ ;  $Q(\text{terminal}, \cdot) = 0$

loop

init  $S$

choose  $a_t$  from  $S$  using  $\epsilon$ -greedy desire from  $Q$

loop; for each episode

a. Take  $a_t$ , observe  $s', r$

b. choose' from  $S'$ ,  $\epsilon$ -greedy desire from  $Q$

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

$s \leftarrow s'$ ,  $a \leftarrow a'$

until  $s$  is terminal

ε-learning - off-policy TD

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \max_{a \in A} Q(s_{t+1}, a) - Q(s_t, a_t)]$$

$$Q(s, a) \forall s \in S, a \in A(s), S(\text{terminal}), \cdot = 0$$

Loop for each

init  $S$

Loop episode

choose  $a$  from  $s$  using policy derived from  $Q$  ( $ε$ -greedy)

Take action  $a$  & observe reward  $R, s'$

$$Q'(s, a) \leftarrow Q(s, a) + \alpha [R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s, a)]$$

$$s \leftarrow s'$$

until  $s$  is terminal.

Expected SARSA:

$$\begin{aligned} Q(s_t, a_t) &\leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma \mathbb{E}_{\pi}[Q(s_{t+1}, a_{t+1}) | s_{t+1}] - Q(s_t, a_t)] \\ &\leftarrow Q(s_t, a_t) + \alpha [R_{t+1} + \gamma \sum_a \pi(a | s_{t+1}) Q(s_{t+1}, a) - Q(s_t, a_t)] \end{aligned}$$

Double Q learning:

$$Q_1(s_t, a_t) \leftarrow Q_1(s_t, a_t) + \alpha [R_{t+1} + \gamma Q_2(s_{t+1}, \arg \max_a Q_1(s_{t+1}, a)) - Q_1(s_t, a_t)]$$

similar to  $Q$  learning → just change the update equation  
 $(Q_1, Q_2)$  with each prob of 0.5

$$G_t = R_{t,1} + \gamma \mathbb{E}_{s_{t+1}}[V_t(s_{t+1})] \delta R_{t,2} + \gamma^2 \mathbb{E}_{t+2}[V_t(s_{t+2})] + \dots$$

$$G_{t+1} = R_{t+1} + \gamma V_t(s_{t+1})$$

$$G_{t+1:t+2} = R_{t+1} + \gamma \mathbb{E}_{s_{t+2}}[V_{t+1}(s_{t+2})]$$

$\downarrow$   
final state of time

t+1-time policy  
state of H+1

n-step returns:

$$V_{t+n}(s_t) = V_{t+n-1}(s_t) + \alpha [G_{t+1:t+n} - V_{t+n-1}(s_t)]$$

n-step TD:  $v \approx V$

input:  $\pi$

$v \in [0, 1], n$

$V(s)$  arbitrary initialization.

all state and operation (for  $s_t$  and  $\pi_t$ ) and their mod  $(n+1)$

Loop: for each episode

init 1 state  $s_0$  & terminal

$T < \infty$

Loop:  $t = 0, 1, 2, \dots$

if  $t < T$  then

action  $\pi(\cdot | s_t)$

observe  $R_{t+1}, s_{t+1}$

$s_{t+1} = \text{terminal} \Rightarrow t \leftarrow t+1$

$t \leftarrow t-n+1 \quad // \text{so far update.}$

if  $\gamma > 0$

$$\alpha \leftarrow \sum_{i=t-n+1}^{min(n, t)} \gamma^{i-t+1} R_i$$

if  $t+n < T$  then  $\alpha \leftarrow \alpha + \gamma^n V(s_{t+n})$

$$V(s_t) \leftarrow V(s_t) + \alpha [a - V(s_t)]$$

until  $a = n-1$

Error reduction property of n-step returns

$$\max_s | \mathbb{E}_\pi[G_{t+1:t+n} | s_t = s] - V_n(s) | \leq \gamma^n \max_s | V_{t+n-1}(s) - V_n(s) |$$

## n-step off policy learning: (importance sampling)

$$V_{t+n}(s_t) = V_{t+n-1}(s_t) + \alpha \sum_{i=t+1}^n [G_{t+i} - V_{t+n-1}(s_t)]$$

$$\rho_{t+h} = \prod_{i=t}^{min(h, T-1)} \frac{\pi(a_i | s_i)}{b(a_i | s_i)} \quad // \text{importance ratio.}$$

init

loop: Algorithm: n step

init  $Q$ ,  $b(a|s)$

init  $\pi$  with n.t. greedy in  $Q$

$\alpha \in [0, 1]$ ,  $n$  positive int.

store and access  $(s_t, a_t, r_t)$

Loop:

init  $s_0 \neq \text{terminal}$

action  $a_t \in b(\cdot | s_t)$

$T \leq \infty$

Loop:  $t = 0, 1, 2, \dots$

| if  $t < T$  then

$a_t$

observe  $R_{t+1}, s_{t+1}$

| if  $s_{t+1}$  terminal

$T \leftarrow t+1$

| else:  $a_{t+1} \sim b(\cdot | s_{t+1})$

$\tau \leftarrow t-n+1$

| if  $\tau \geq 0$ :  $\min(\tau, T-1) \frac{\pi(a_i | s_i)}{b(a_i | s_i)}$

$\rho \leftarrow \prod_{i=\tau+1}^{\min(\tau, T)} \frac{\pi(a_i | s_i)}{b(a_i | s_i)}$

$G_t \leftarrow \prod_{i=\tau+1}^{\min(\tau, T)} \gamma^{i-\tau-1} R_i$

| if  $\tau+n < T$  then  $G_t \leftarrow G_t + \gamma^{\tau-n} Q(s_{\tau+n}, a_{\tau+n})$

$Q(s_\tau, a_\tau) \leftarrow Q(s_\tau, a_\tau) + \alpha \rho [G_t - Q(s_\tau, a_\tau)]$

1-step SARSA

$$G_{t+n} = R_{t+1} + \gamma P_{t+1} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q_{t+n-1}(s_{t+n}, a_{t+n})$$

$\downarrow$  ;  $t+n \geq T$  ;  $n \geq 1; 0 \leq t \leq T-n$

$$Q_{t+n}(s_t, a_t) := Q_{t+n-1}(s_t, a_t) + \alpha \left[ G_{t+n} - Q_{t+n-1}(s_t, a_t) \right]$$

1-step SARSA:  $Q = q_\pi$  or  $Q_\pi$ : "ig" on policy

init  $Q(s, a)$ ;  $s \in S, a \in A$ ; arbitrary

$\pi$  init =  $\epsilon$ -greedy.

Param:  $\alpha \in [0, 1]$ ; small  $\epsilon > 0$ ;  $n$

store and access (for  $s_t, a_t, R_t$ ) to take  $(n+1) \bmod n$

Loop:

ini  $s_0 \neq$  terminal

select  $\pi$  store  $\pi_{0:n}( \cdot | s_0 )$

$T \in \mathbb{N}$

loop for  $t = 0, 1, 2, \dots$

if  $t < T$  then

Action  $a_t$

observe  $R_{t+1}, s_{t+1}$

$s_{t+1} = : \text{terminal} \Rightarrow T = t+1$  end state.

else:  $\pi_{t+1} \sim \pi( \cdot | s_{t+1} )$

$t \leftarrow t+n-1$  ; (worst updated time)

if  $t \geq 0$ :  $\min(t+n, T), i = t \rightarrow T$

$a_i \leftarrow \sum_{i' \in C(i)} R_{i'}$

if  $t+n < T$  then  $a \leftarrow a + \gamma^n Q_\pi(s_{t+n}, a_{t+n})$

$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [a - Q(s_t, a_t)]$

if  $\pi$  is being learned,  $\pi(\cdot | s_t)$  is  $\epsilon$ -greed w.r.t  $Q$ .

until  $t = T-1$

Q-Learning & Expected SARSA  $\rightarrow$  1 step algo. } off policy.

Tree backup  $\rightarrow$  multistep alg.

$$1\text{-step expected SARSA: } G_{t+1|t+1} = R_{t+1} + \gamma \sum_a \pi(a|s_{t+1}) Q_{t+1}(s_{t+1}, a)$$

$$2\text{-step: } G_{t+2|t+1} = R_{t+1} + \gamma \sum_a \pi(a|s_{t+1}) Q_{t+1}(s_{t+1}, a) \\ + \gamma \sum_a \pi(a|s_{t+2}) \underbrace{\left[ R_{t+2} + \gamma \sum_b \pi(b|s_{t+2}) Q_{t+2}(s_{t+2}, b) \right]}_{G_{t+2|t+2}}$$

$$+ \gamma \sum_a \pi(a|s_{t+2}) \left[ R_{t+2} + \gamma \sum_b \pi(b|s_{t+2}) Q_{t+2}(s_{t+2}, b) \right]$$

$$\vdots$$

$$G_{t+n|t+1}$$

$$n\text{-step: } G_{t+n|t+1} = R_{t+1} + \gamma \sum_{a \in A_{t+1}} \pi(a|s_{t+1}) Q_{t+n|t+1}(s_{t+1}, a) + \gamma^{n-1} Q_{t+n|t+1}$$

$$\text{update value: } Q_{t+n}(s_t, a_t) = Q_{t+n|t+1}(s_t, a_t) + \alpha \left[ G_{t+n|t+1} - Q_{t+n|t+1}(s_t, a_t) \right]$$

Tree back alg:  $Q \approx q_{\pi^*}$ , or  $q_{\pi}$

init  $Q(s, a) \in \mathbb{R}$ ,  $a \in A$ ,  $\pi$ ,  $\alpha \in (0, 1]$ ,  $n \mod (n+1)$

Loop: episode

init stores  $s_0$  determined

action  $a_0$ , arbitrary  $s_0$ , store  $a_0$

$T \leftarrow \infty$

Loop  $t = 0, 1, 2, \dots$

if  $t < T$ :

action  $a_t$  obj  $R_{t+1}, s_{t+1}$

if  $s_{t+1}$  == terminal  $T \leftarrow t+1$

else: arbitrary  $a_{t+1}$  as  $s_{t+1}$ ; store  $a_{t+1}$

$t \leftarrow t+1-n$

if  $n \geq 0$ :

if  $t+1 \geq T$ :  $a \in A_t$

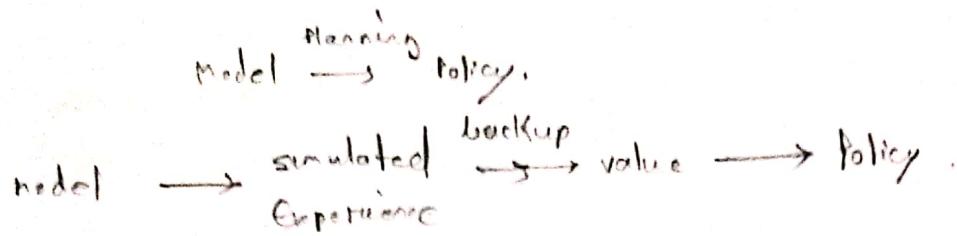
else:  $a \in A_{t+1} + \gamma \sum_a \pi(a|s_{t+1}) Q(s_{t+1}, a)$

Loop:  $k = \min(t, T-1)$  down through  $t+1$

$a \in A_k + \gamma \sum_{a' \in A_{k+1}} \pi(a'|s_k) Q(s_{k+1}) + \gamma \pi(a_k|s_k) Q(s_k, a_k)$

$Q(s_k, a_k) \leftarrow Q(s_k, a_k) + \alpha [a - Q(s_k, a_k)] \dots$

Planning & Learning Method: model based  $\rightarrow$  DP, heuristic search (Planning)  
 model free  $\rightarrow$  MC, TD, (Learning)



### random sample one-step tabular Q-planning:

Loop forever:

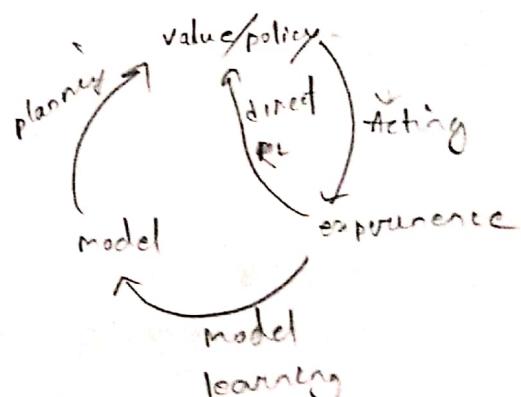
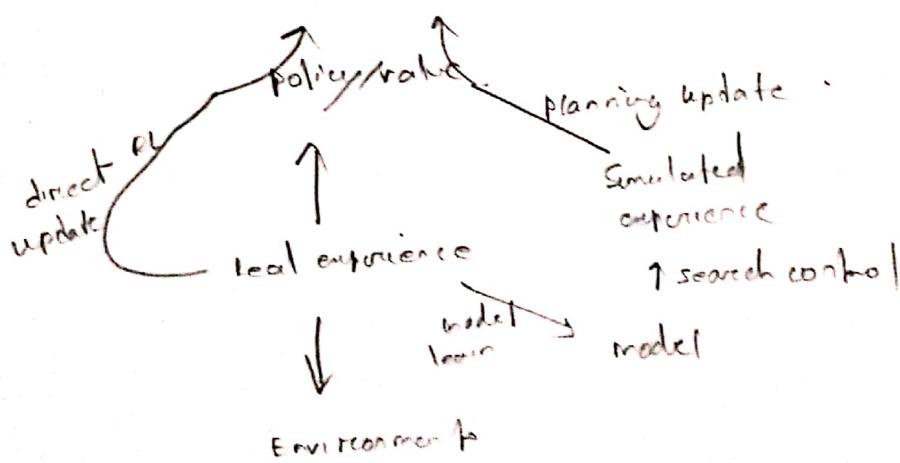
1. select state  $s \in S$  and  $a \in A(s)$  at random

2. Send:  $s, a$  to sample model & obtain  
 sample  $R$ , sample  $s'$

3. apply one step tabular Q learning.

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

### Dyna: Planning, Acting, Learning:



### Tabular Dyna Q

INIT:  $Q(s, a)$ , model  $(s, a)$ ,  $s \in S$ ,  $a \in A(s)$

Loop forever:

a)  $s \leftarrow$  current(non-terminal) state  
 b)  $a \leftarrow$  greedy( $s, Q$ )

c) Take action  $a \rightarrow$  observe  $R, s'$   
 $d) Q(s, a) \leftarrow Q(s, a) + \alpha [R + \gamma \max_a Q(s', a) - Q(s, a)]$

e) model  $(s, a) \leftarrow R, s'$  (Deterministic env)

f) loop  $n$  times  $\rightarrow s \leftarrow$  random state, and  $a \leftarrow$  action taken  $R, s' \leftarrow$  model  $(s, a)$   
 $Q(s, a) \leftarrow R + \gamma \max_a Q(s', a) - Q(s, a)$

## Prioritized Sweeping: Dyna Q

Init  $Q(s, a)$ , model  $(s, a) \rightarrow s'$ , Pqueue to empty.

Loop forever:

a)  $s \in$  nonterminal state

i)  $a \in \text{policy}(s, Q)$

ii) Take  $A$  & observe  $R, s'$

iii)  $\text{model}(s, a) \leftarrow R, s'$

iv)  $p \leftarrow |R + \gamma \max_a Q(s', a) - Q(s, a)|$

v) if  $p > \theta$  then insert  $s, a$  into Pqueue with priority  $p$ .

vi) loop repeat  $n$  times while Pqueue not empty:

$s, a \leftarrow \text{first(Pqueue)}$

$R, s' \leftarrow \text{model}(s, a)$

$Q(s, a) \leftarrow Q(s, a) + \alpha [\bar{R} + \gamma \max_a Q(s', a) - Q(s, a)]$

loop for  $\forall \bar{s}, \bar{a}$  predicted to lead to  $s$ :

$\bar{R} \leftarrow \text{predicted reward } \bar{s}, \bar{a}, s$

$p \leftarrow |\bar{R} + \gamma \max_a Q(s, a) - Q(\bar{s}, \bar{a})|$

if  $p > \theta$  then insert  $\bar{s}, \bar{a}$  into Pqueue with priority  $p$

## Approximate Solution Method:

Parameter estimation: Approximate  $v_{\pi}$  given  $\pi$

$$\hat{v}(s, w) \approx v_{\pi}(s) \quad // \text{as memory}$$

under state  $s \rightarrow u$  target update.

$$m \rightarrow s_t \rightarrow a_t ; TD(0) \Rightarrow s_t \rightarrow r_{t+1} + \gamma \hat{v}(s_{t+1}, w)$$

$$; TD(n) \Rightarrow s_t \rightarrow a_{t+n}$$

$$\text{DP policy evaluation } s \rightarrow E_{\pi} [r_{t+1} + \gamma \hat{v}(s_{t+1}, w) | s_t = s]$$

Prediction Objective: mean square value error  $\bar{VE}(w) = \sum_{s \in S} \mu(s) [\hat{v}_n(s) - \hat{v}(s, w)]^2$   
 ↓  
 (-time spent in  $s$ )  
 (if no policy distribution under  $\pi$ )

$$\begin{cases} \text{policy estimation: } \eta(s) = b(s) + \sum_{\tilde{s}} \pi(\tilde{s}) \sum_a \pi(a|\tilde{s}) p(s|\tilde{s}, a) ; \tilde{s} \in S \\ \mu(s) = \frac{\eta(s)}{\sum_{\tilde{s} \in S} \eta(\tilde{s})} ; \end{cases}$$

$$\text{Goal } \bar{VE}(w^*) \leq \bar{VE}(w) ; \forall w$$

stochastic gradient & semi-gradient:

$$w_{t+1} = w_t + \alpha \nabla \left[ v_{\pi}(s_t) - \hat{v}(s_t, w) \right]^2$$

$$= w_t + \alpha \nabla \left[ v_{\pi}(s_t) - \hat{v}(s_t, w) \right] \nabla_w \hat{v}(s_t, w)$$

$$\Rightarrow w_{t+1} = w_t + \alpha \left[ v_{\pi}(s_t) - \hat{v}(s_t, w) \right] \nabla_w \hat{v}(s_t, w)$$

$$\hat{v}_{\pi}(s_t) = E \left[ v_{\pi} | s_t = s \right] // \text{unbiased case}$$

## Gradient Monte Carlo algo. $\hat{v} \in \mathbb{R}^d$

Input:  $\pi$  policy

input : differentiable  $\hat{v}: \mathcal{S} \times \mathbb{R}^d \rightarrow \mathbb{R}$

param  $\alpha > 0$

init  $w \in \mathbb{R}^d$

Loop forever:

Generate episode  $s_0, a_0, r_1, s_1, A_1, \dots, R_T, g_T$  using  $\pi$

loop for each step of episode,  $t = 0, 1, \dots, T-1$

$$w \leftarrow w + \alpha [r_t + \gamma \hat{v}(s_{t+1}, w) - \hat{v}(s_t, w)] \nabla \hat{v}(s_t, w)$$

Semi-gradient: target:  $v_t = r_{t+1} + \gamma \hat{v}(s_{t+1}, w)$

D(0) semi-gradient algo:

Loop for each episode:

init  $s$

loop for each step of episode:

choose  $a \sim \pi(a|s)$

Take action  $a$ , observe  $R, s'$

$$w \leftarrow w + \alpha [R + \gamma \hat{v}(s', w) - \hat{v}(s, w)] \nabla \hat{v}(s, w)$$

$s \leftarrow s'$

until  $s$  is terminal.

## Linear methods

$$\hat{v}(s, w) = w^\top \mathbf{x}(s) = \sum_i w_i x_i(s) \quad ; \quad \mathbf{x}(s) = (x_1(s), x_2(s), \dots, x_d(s))$$

$$\nabla \hat{v}(s, w) = \mathbf{x}(s)$$

$$\therefore w_{t+1} = w_t + \alpha [r_t + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t)] \mathbf{x}(s_t)$$

## Semi-gradient TD(0) method: Linear.

$$w_{t+1} = w_t + \alpha [r_{t+1} + \gamma w_t^\top x_{t+1} - w_t^\top x_t] x_t \\ = w_t + \alpha [r_{t+1} x_t - x_t (x_t - \gamma x_{t+1}) w_t]$$

$$E[w_{t+1} | w_t] = w_t + \alpha [b - A w_t];$$

$$b = E[r_{t+1} x_t] \in \mathbb{R}^d$$

$$A = E[x_t (x_t - \gamma x_{t+1})^\top] \in \mathbb{R}^d \times \mathbb{R}^d$$

## n-step semi-gradient TD

input:  $T$

$\hat{V}$ : state  $\mathbb{R}^d \rightarrow \mathbb{R}$ :  $\hat{V}(\text{terminal}) = 0$

Algo. param:  $\alpha > 0$

$w \leftarrow \text{init}$

+11 store and access ( $s_t$  &  $r_t$ ) can

Loop for each episode:

init  $s_0$ ,  $\text{terminal}$

$T \leftarrow \infty$

Loop for  $t = 0, 1, 2, \dots$ :

if  $t < T$ , then:

Take an action  $\pi(\cdot | s_t)$

observe  $r_{t+1}$ , b next state  $s_{t+1}$

if  $s_{t+1}$  'is terminal'  $T \leftarrow t+1$

$\tau \leftarrow t-n+1$

if  $\tau \geq 0$ :  $\min(\tau, T)$

$b_t \leftarrow \sum_{i=\tau+1}^{T+1} \gamma^{i-\tau-1} r_i$

$\nabla \hat{V}(s_t, w)$

until  $\tau \geq 1$ , if  $\tau < T$  then:  $b_t \leftarrow b_t + \delta^n \hat{V}(s_{\tau+n}, \bar{w}) + \bar{w} \leftarrow \bar{w} + \alpha [b_t - \hat{V}(s_\tau, w)] e$

$$V(s, a, w) \in \mathbb{R}^T \times (s, a) = \sum_{i=1}^d w_i v_i(s, a)$$

Semi-gradient n-step SARSA

$$G_{t+1:n} := R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n \hat{q}(s_{t+1:n}, a_{t+1:n}, w_{t+1:n})$$

$$w_{t+1:n} := w_{t+1:n-1} + \alpha [G_{t+1:n} - \hat{q}(s_t, a_t, w_{t+1:n-1})] \nabla \hat{q}(s_t, a_t, w_{t+1:n})$$

Episodic semi-gradient n-step SARSA:  $\hat{q} \approx q_\pi$

Input:  $\hat{q}: \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Output: a policy  $\pi$

$\alpha > 0, \epsilon > 0$

init value-function  $w \in \mathbb{R}^d$  arbitrarily (e.g.,  $w = 0$ )

store and access  $(s_{t+1:n}, R_t)$  take mode ( $n+1$ )

Loop: for each episode:

init & store  $s_0$  & terminal

select & store  $a_{0:n} \sim \pi(\cdot | s_0)$  or  $\epsilon$ -greedy w.r.t.  $\hat{q}(s_0, w)$

$T \leq \infty$

Loop for  $t = 0, 1, 2, \dots$

if  $t < T$  then:

take action  $a_t$

observe & store next reward  $R_{t+1}$ , state  $s_{t+1}$

if  $s_{t+1}$  is terminal

$T \leftarrow t+1$

else

select & store  $a_{t+1:n} \sim \pi(\cdot | s_{t+1})$  or  $\epsilon$ -greedy

$n \leftarrow t+1$

if  $\gamma \geq 0$ ;  $\min_{i=t+1}^n \gamma^{i-t-1} R_i$

$a \leftarrow \sum_{i=t+1}^n \gamma^{i-t-1} R_i$

if  $a < T$  then  $a \leftarrow a + \gamma^n \hat{q}(s_{t+1:n}, a_{t+1:n}, w)$

$w \leftarrow w + \alpha [a - \hat{q}(s_t, a_t, w)] \nabla \hat{q}(s_t, a_t, w)$

until  $a = T-1$

$$w_{t+n} \doteq w_{t+n-1} + \alpha [c_{t+n} \cdot \nabla \hat{q}(s_t, w_{t+n-1})] \nabla \hat{q}(s_t, w_{t+n-1})$$

$$\Rightarrow \hat{q}_{\text{critic}} = r_{t+1} + \gamma q_{t+1} + \dots + \gamma^{n-1} r_{t+n} + \gamma^n \hat{q}(s_{t+n}, w_{t+n-1}), \quad 0 \leq n \leq T$$

On-policy Control with Approximation :-

$$q_{\pi_k}(s, a, w) = q_{\pi^*}(s, a); \quad w \in \mathbb{R}^d$$

Periodic Semi-gradient Control

$$w_{t+1} \doteq w_t + \alpha [u_t - \hat{q}(s_t, a_t, w_t)] \nabla \hat{q}(s_t, a_t, w)$$

One step SARSA :

$$w_{t+1} \doteq w_t + \alpha [r_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}, w_t) - \hat{q}(s_t, a_t, w)] \nabla \hat{q}(s_t, a_t, w)$$

Semi-gradient SARSA:  $\hat{q} \approx q^*$

Input: action value  $\hat{q}: S \times A \times \mathbb{R}^d \rightarrow \mathbb{R}$

$\alpha > 0, \epsilon > 0$

$w \in \mathbb{R}^d$  arbitrary.

Loop for each episode :

$s, a \leftarrow$  init state & action of episode ( $\epsilon$ -greedy)

Loop for each step:

Take action  $a'$  observe  $R, s'$

if  $s'$  is terminal

$$w \leftarrow w + \alpha [r - \hat{q}(s, a, w)] \nabla \hat{q}(s, a, w)$$

go next episode.

choose  $a'$  as function  $\hat{q}(s', \cdot, w)$ ;  $\epsilon$ -greedy

$$w \leftarrow w + \alpha [r + \gamma \hat{q}(s', a', w) - \hat{q}(s, a, w)] \nabla \hat{q}(s, a, w)$$

$s \leftarrow s'$

$a \leftarrow a'$

$$\begin{aligned}
 r(\pi) &= \lim_{h \rightarrow \infty} \frac{1}{h} \sum_{t=1}^h E \left[ R_t \mid s_0, a_0, a_1, \dots, a_{t-1}, s_t \right] \\
 &= \lim_{h \rightarrow \infty} E \left[ R_t \mid s_0, a_0, \dots, a_{t-1}, s_t \right] \\
 &= \sum_s \mu_\pi(s) \sum_a \pi(a|s) \sum_{s', r} p(s'|s, a) r
 \end{aligned}$$

$$\sum_s \mu_\pi(s) \sum_a \pi(a|s) p(s' | s, a) = \mu_\pi(s')$$

$\nearrow$  differential return value.

$$g_t = R_{t+1} - r(t) + R_{t+2} - r(t) + \dots + R_{t+n} - r(t)$$

$$q_\pi(s) = E_\pi [g_t \mid s_t = s]$$

$$q_{\pi_t}(s, a) = E_\pi [g_t \mid s_t = s, a_t = a]$$

$\Rightarrow$  Some definitions:

$$\psi_\pi(s) = \sum_a \pi(a|s) \sum_{r, s'} p(s', r | s, a) [r - \bar{r} + \psi_\pi(s')]$$

$$q_{\pi_t}(s, a) = \sum_{r, s'} p(s', r | s, a) [r - \bar{r} + \sum_{a'} \pi(a'|s') q_\pi(s', a')]$$

$$\bar{\psi}_t(s) = \max_a \sum_{r, s'} p(s', r | s, a) \left[ r - \max_\pi \bar{r} + \psi_\pi(s') \right]$$

$$q_{\pi_t}(s, a) = \sum_{r, s'} p(s', r | s, a) \left[ r - \max_\pi \bar{r} + \max_{a'} q_\pi(s', a') \right]$$

$$\delta_{-1} = R_{t+1} - \bar{r}_t + \hat{\psi}(s_{t+1}, w_t) - \hat{\psi}(s_t, w_t)$$

$$\delta_t = R_{t+1} - \bar{r}_t + \hat{q}_\pi(s_{t+1}, a_{t+1}, w_t) - \hat{q}_\pi(s_t, a_t, w_t)$$

$$\text{update: } w_{t+1} = w_t + \alpha \delta_t \nabla \hat{q}_\pi(s_t, a_t, w)$$

Semi gradient SARSA:  $q \approx \hat{q}_{\pi}$

Input  $\hat{q}: S \times A \times R^d \rightarrow \mathbb{R}$

$\alpha, \gamma > 0$

init  $w \in \mathbb{R}^d$ ,  $[w \neq 0]$

init  $R \in \mathbb{R}$

init state  $s$ , action  $a$

Loop for each step:

Take action  $a$ , observe  $r', s'$

choose  $a'$  as function of  $\hat{q}(s, w)$  //  $\epsilon$ -greedy

$$\delta \leftarrow r - R + \hat{q}(s', a', w) - q(s, a, w)$$

$$\bar{R} \leftarrow R + \beta \delta$$

$$w \leftarrow w + \alpha \delta \nabla \hat{q}(s, a, w)$$

$$s \leftarrow s'$$

$$a \leftarrow a'$$

Differential semi-gradient:

$$g_{t+1:n} = r_{t+1} - \bar{R}_{t+1:n-1} + \dots + r_{t+n} - \bar{R}_{t+n-1} + \hat{q}(s_{t+n}, a_{t+n}, w_{t+n-1})$$

$$t+n \geq T \text{ then } g_{t+1:n} = g_T$$

$$\delta_t = g_{t+1:n} - \hat{q}(s_t, a_t, w)$$

Differential semi-gradient algorithm:

store & access  $(s_t, a_t, R_t)$  in model  $(t+1)$

init  $s_0, a_0$

Loop:  $t=0, 1, \dots$

take action  $a_t$

obtains store next reward  $R_{t+1}$  & next stage  $s_{t+1}$

select  $a_{t+1} \sim \pi_C(s_{t+1})$  //  $\epsilon$ -greedy

$$\bar{r} \leftarrow t-n+1$$

$$\text{if } \bar{r} \geq 0 : \delta \leftarrow \sum_{i=\bar{r}+1}^{t+n} (r_i - \bar{R}) + \hat{q}(s_{t+n}, a_{t+n}, w) - \hat{q}(s_n, a_n, w)$$

$$\bar{R} \leftarrow \bar{R} + \beta \delta$$

$$w \leftarrow w + \alpha \delta \nabla \hat{q}(s_n, a_n, w)$$

Function Approximation: Linear / non

→ Bootstrapping: TD, IP, MC (update target & estimates)

→ Off-policy training: others than target policy.

combination of these cause instability.

Eligibility trace:  $TD(\lambda) \rightarrow Q$  & SARSA:

$\lambda \in [0, 1]$ ;  $\lambda = 0$ : TD(0) method;  $\lambda = 1$ : MC method; balance between TD & MC.

Eligibility trace; start term vector parallel; long-term weight  $w_t$

Trace decay  $\propto [0, 1]$ .

$$\lambda \text{ returns}, \quad G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t+n}$$

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t+n} + \lambda^{T-t-1} G_T$$

$$w_{t+1} = w_t + \alpha \left[ G_t^\lambda - \hat{v}(s, w) \right] \nabla_w \hat{v}(s, w)$$

Semi gradient TD ( $\lambda$ ):

Input:  $\hat{v}$

Input: differentiable  $\hat{v}: s^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$   $\hat{v}(\text{terminal}; \cdot) = 0$

$\alpha > 0, \lambda \in [0, 1]$

$w \in \text{Arbitrary}$ ,

loop:  
init  $s$

$\hat{z} \leftarrow 0$

loop:

choose  $A \sim \pi(\cdot | s)$

Take  $A$ , observe  $r, s'$

$$\begin{cases} \hat{z} \leftarrow \gamma \lambda \hat{z} + \hat{v}(s, w) \\ \delta \leftarrow r + \hat{v}(s', w) - \hat{v}(s, w) \\ w \leftarrow w + \alpha \delta \hat{z} \\ s \leftarrow s' \end{cases}$$

until  $s'$  is terminal.

## n-step truncated $\lambda$ -return methods: TD( $\lambda$ )

$$G_{t:h}^\lambda = (1-\lambda) \sum_{n=1}^{h-t+1} \lambda^{n-1} G_{t+1:t+n} + \lambda^{h-t+1} G_{t:h}$$

$$\underline{w}_{t+n} = \underline{w}_{t+n-1} + \alpha \left[ G_{t+1:t+n} - \hat{v}(s_t, \underline{w}_{t+n-1}) \right] \nabla \hat{v}(s_t, \underline{w}_{t+n-1})$$

$$G_{t:t+k}^\lambda = \hat{v}(s_t, \underline{w}_{t+1}) + \sum_{i=t}^{t+k-1} (\gamma \lambda)^{i-t} \delta_i'$$

$$\text{where, } \delta_t' = r_{t+1} + \gamma \hat{v}(s_{t+1}, \underline{w}_t) - \hat{v}(s_t, \underline{w}_{t-1})$$

## online $\lambda$ -return algorithm:

$$\underline{w}_{t+1}^h = \underline{w}_t^h + \alpha \left[ G_{t:h}^\lambda - \hat{v}(s_t, \underline{w}_t) \right] \nabla \hat{v}(s_t, \underline{w}_t); 0 \leq t < h \leq T$$

## True online TD( $\lambda$ ) [ $w^T x = v_\pi$ ]

input  $\pi$

input:  $x: s^T \rightarrow \mathbb{R}^d$

$\alpha > 0, \lambda \in [0, 1]$

Loop:  
init state and get feat vector  $x$

$z \leftarrow 0$

$v_{\text{old}} \leftarrow 0$

Loop:

choose  $a \sim \pi$

take action,  $a$ , observe  $R, x'$

$v \leftarrow w^T x$

$v' \leftarrow w^T x'$

$\delta \leftarrow R + \gamma v' - v$

$z \leftarrow \gamma \lambda z + (1 - \alpha \gamma \lambda) z^T x$

$w \leftarrow w + \alpha (\delta + v - v_{\text{old}}) z - \alpha (v - v_{\text{old}}) x$

$v_{\text{old}} \leftarrow v'$

$x \leftarrow x'$

until  $x' = 0$  (terminal)

SARSA(λ)  $w^T x = q_n$  or  $q_{\pi}$

input:  $f(s, a) \approx$  return featr. for  $s \& a$ .

input: policy  $\pi$

Algo param  $\alpha > 0$ ,  $\lambda \in [0, 1]$

init  $w = (w_1 \dots w_d)^T \in \mathbb{R}^d$ ,  $\varepsilon \in \mathbb{R}^d$

Loop: for episode.

init  $s$

choose  $A \sim \pi(\cdot | s)$  ε-greedy to  $\hat{q}(s, \cdot, w)$

$\varepsilon \leftarrow 0$

Loop for each step:

Take action  $A \rightarrow r, s'$

$\delta \leftarrow r$

Loop for  $i$  in  $f(s, A)$ :

$\delta \leftarrow \delta - w_i$

$z_i \leftarrow z_i + 1$  // accumulating trace

or  $z_i \leftarrow 1$  // replacing trace

if  $s'$  is terminal

$w \leftarrow w + \alpha \delta z$

Go to next episode

choose  $A \sim \pi(\cdot | s')$  or near greedily  $\hat{q}(s', \cdot, w)$

Loop for  $i$  in  $f(s', A)$ :  $\delta \leftarrow \delta + \gamma w_i$

$w \leftarrow w + \alpha \delta z$

$\varepsilon \leftarrow \gamma \varepsilon$

$s \leftarrow s'; A \leftarrow A'$

## True Online SARSA ( $\lambda$ ) w/ $x$ off

Input:  $x : s \times A \rightarrow \mathbb{R}^d$      $x(\text{terminal}, \cdot) = 0$

Input:  $\pi$

$\alpha > 0$ ;  $\gamma \in [0, 1]$

init  $w \in \mathbb{R}^d$

Loop for each episode  $e$ :

init  $s$

choose  $a \sim \pi(\cdot | s)$

$x \leftarrow x(s, a)$

$z \leftarrow 0$

$Q_{\text{old}} \leftarrow 0$

Loop for each step:

Take  $a' \rightarrow R, s'$

choose  $a'' \sim \pi(\cdot | s')$  near greedily  $s'$  using  $w$

$x' \leftarrow x(s', a')$

$Q \leftarrow w^T x$

$Q' \leftarrow w^T x'$

$\delta \leftarrow R + \gamma Q' - Q$

$z \leftarrow \gamma z + (1 - \alpha \gamma) z^T x$

$w \leftarrow w + \alpha (\delta + Q - Q_{\text{old}}) z - \alpha (Q - Q_{\text{old}}) x$

$Q_{\text{old}} \leftarrow Q'$

$x \leftarrow x'$

$A \leftarrow a'$

until  $s$  is terminal.

semi gradient method:

$$\rho_t = \rho_{t:t} = \frac{\pi(a_t | s_t)}{b(\pi_t | s_t)}$$

$$w_{t+1} = w_t + \alpha \rho_t \delta_t \nabla \hat{v}(s_t, w_t)$$

$$\delta_t = r_{t+1} + \gamma \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \quad // \text{continuous}$$

$$\delta_t = r_{t+1} - \bar{r}_t + \hat{v}(s_{t+1}, w_t) - \hat{v}(s_t, w_t) \quad // \text{episodic}$$

expected SARSA:

$$w_{t+1} = w_t + \alpha \delta_t \nabla \hat{q}(s_t, a_t, w_t)$$

$$\delta_t = r_{t+1} + \gamma \sum_a \pi(a | s_{t+1}) \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \quad // \text{episodic}$$

$$\delta_t = r_{t+1} - \bar{r}_t + \sum_a \pi(a | s_{t+1}) \hat{q}(s_{t+1}, a, w_t) - \hat{q}(s_t, a_t, w_t) \quad // \text{continuous}$$

n-step version of semi gradient expected SARSA:

$$w_{t+n} = w_{t+n-1} + \alpha \rho_{t+n} - \rho_{t+n-1} [c_{t+n-1} - \hat{q}(s_t, a_t, w_{t+n-1}) \nabla \hat{q}(s_t, a_t, w_{t+n-1})]$$

$$c_{t+n-1} = r_{t+1} + \dots + \gamma^{n-1} r_{t+n} + \gamma^n \hat{q}(s_{t+n-1}, a_{t+n-1}, w_{t+n-1}) \text{ on / episode}$$

$$\delta_{t+n-1} = r_{t+1} - \bar{r}_t + \dots + r_{t+n} - \bar{r}_{t+n-1} + \hat{q}(s_{t+n-1}, a_{t+n-1}, w_{t+n-1}) \text{ / cont.}$$

n-step-free - backup style:

$$w_{t+n} = w_{t+n-1} + \alpha [c_{t+n-1} - \hat{q}(s_t, a_t, w_{t+n-1})] \nabla \hat{q}(s_t, a_t, w_{t+n-1})$$

$$c_{t+n-1} = \hat{q}(s_t, a_t, w_{t-1}) + \sum_{k=t}^{t+n-1} \delta_k \prod_{i=t+1}^k \gamma \pi(a_i | s_i)$$