

## tic tac toe

$nr \leftarrow 3$ ,  $nc \leftarrow 4$ ,  $size \leftarrow nr \cdot nc$ .

1. state

$\forall i \in [1, |re|]$

1.2.  $M \leftarrow M(nr, nc) \leftarrow \emptyset$

$x \leftarrow re(i)$

$M$  is data, state,  $M_{ij} \in \{-1, 0, 1\}$

$x_0 \leftarrow nc$  if  $i \leq nr$

$f(M) = h$ ,  $f$  is bijective.  $h$  is hashvalue.

$(nr, \text{else})$

1.2.1  $f: x \leftarrow 0$

if  $x = x_0$

$\forall i \in M, x \leftarrow x + (i+1)$

$w \leftarrow 1, end \leftarrow 1, \text{stop}$

$\therefore h \leftarrow x$

$\text{if } x = -x_0$

2.2.  $\therefore h$  is ternary number:

$w \leftarrow -1, end \leftarrow 1, \text{stop}$

$(M_1+1, M_2+1, \dots, M_{nr, nc}+1)$

$\text{if } ||M||_1 = size$

1.3. is end

$w \leftarrow 0, end \leftarrow 1, \text{stop}$

$\text{if } end \neq -1, \text{return end.}$

$\text{end} \leftarrow 0.$

$\{/ end \in \{-1, 0, 1\}\}$

14. next state ( $i, j$ , sb:symbol)

sequence  $re \leftarrow \emptyset$ .

$s' \leftarrow \text{state}(1)$

$\forall i \in [1, nr],$

$s': M \leftarrow M$

$re \text{ add } \text{sum}(M_{i,:})$

$s'. M_{ij} \leftarrow sb.$

$\forall i \in [1, nc],$

15. print state

$re \text{ add } \text{sum}(M_{:,i})$

init

$\text{if } nr = nc$

$M \leftarrow \emptyset, \in \{-1, 0, 1\}^{nr \times nc}$

$t \leftarrow \text{trace}(M)$

$h \leftarrow -1, \in \{-1\} \cup \{0, 1, \dots, 3^{\frac{size}{2}} - 1\}$

$rt \leftarrow \text{reverse trace}(M)$

$w: \text{winner} \leftarrow 0 \in \{-1, 0, 1\}$

$re \text{ add } t.$

$end \leftarrow -1 \in \{-1, 0, 1\}$

$re \text{ add } rt.$

2. all states S.

2.3 get all states:

csb: current symbol  $\leftarrow 1$ .

// 1 moves first, -1 second.

s  $\leftarrow$  state()

S: all states  $\leftarrow$  dict

S add (s.hash, s)

get all states impl (s, cs, S)

2.4 get all states impl (s, cs, S) 3.4 alternate:

$\forall i \in [1, nr]$ ,

$\forall j \in [1, nc]$ ,

if  $s.M_{ij} = 0$

then  $s' \leftarrow s.\text{next state}(i, j, cs)$

$h' \leftarrow s'.\text{hash}$

if  $h' \notin S$

$\text{end}' \leftarrow s' \text{ is end}$

S add ( $h'$ ,  $s'$ )

if  $\text{end}' = 0$

get all states impl ( $s'$ ,  $-cs$ , S) if is print state

3. Judge (Player p1, p2).

3.2 init:

$p1 \leftarrow p1$ ,  $p2 \leftarrow p2$ .

cp: current player  $\leftarrow 0 \in \{-1, 0, 1\}$

~~psb1~~: pl symbol  $\leftarrow 1$

~~psb2~~:  $\leftarrow -1$

p1 set symbol (~~psb1~~)

p2 set symbol (~~psb2~~)

3.3. reset:

p1 reset

p2 reset

3.4 alternate:

loop:

yield p1

yield p2

3.5 play (is print state  $\leftarrow 0$ )

alt  $\leftarrow$  alternate

reset

s  $\leftarrow$  State

p1 set state s

p2 set state s

s point

loop:

player p  $\leftarrow$  next(alt)

page 2.

```

i,j,sb ← p.act
s' ← s.next.state(i,j,sb)
h' ← s'.hash
s ← S(h') // get s from S.
// not use s', 'cause s has some
// values computed already.
P1 set state s.
P2 set state s.
if s.print
    s.print
if s.isend
    return s.winner.

```

```

4.4 reset
St ← ∅
g ← ∅
4.5 set state(s)
St.add(s,
g.add(1).
4.6 set symbol(sb)
// compute V(s) initially.
sb ← sb.
 $\forall h \in S$ 
s ← S(h)
if s.end
    if s.winner = sb
        V(h) ← 1
    else if s.winner = 0
        V(h) ← 0.5
    else
        V(h) ← 0

```

~~end loop.~~

4. A7 player

4.2 init ( $\alpha$ : step size  $< 0.1$ ,  
 $\varepsilon \leftarrow 0.1$ )

V: estimations ← dict

$\alpha \leftarrow \alpha$ ,  $\varepsilon \leftarrow \varepsilon$

Sequence St ← ∅

Seq greedy g ← ∅

symbol sb ← 0

page 3.

#### 4.7 backup:

states hash  $hs \leftarrow \{ s.\text{hash} \mid s \in S_t \}$

//  $h \leftrightarrow s$

$\forall i \in \{ |hs| - 1, \dots, 2, 1 \}$

$s \leftarrow hs(i), s' \leftarrow hs(i+1)$

td error  $\delta \leftarrow g(i)(V(s') - V(s))$

$V(s) \leftarrow V(s) + \alpha \delta$

$\nabla V_t = \{ V(s) \mid s \in ss \}$

$i \leftarrow \text{argmax } V_t(i) \text{ randomly}$

// if  $\exists$  multiple  $i$ , choose

// randomly.

$a \leftarrow \alpha s(i)$

return  $(a, sb)$

#### 4.8 act

// choose  $a$  based on  $s$ .

$s \leftarrow S_t(-1)$

$as : \text{next action set} \leftarrow \emptyset$

$ss : \text{next state set} \leftarrow \emptyset$

$\forall i \in [1, nr]$

$\forall j \in [1, nc]$

if  $s.M(i, j) = 0$

$as \text{ add } (i, j)$

$s' \leftarrow s, \text{next state}(i, j, sb)$

$ss \text{ add } (s', \text{hash})$

$ran \leftarrow \text{rand()}$

if  $ran < \varepsilon$

choose  $i$  from  $[1, |as|]$  randomly. ~~5.6 set~~ :

at  $as(i)$ , ~~and~~  $g(-1) \leftarrow 0$

return  $(a, sb)$

else

#### 4.9 save policy.

#### 4.10 load policy

#### 5. Human player

##### 5.1 init

$sb \leftarrow 0$

keys  $\leftarrow \{ 'q', 'w', 'e', \dots \}$

state  $s \leftarrow \text{state}$

##### 5.2 reset

pass

##### 5.3 set state ( $s$ )

$s \leftarrow s$

##### 5.4 set symbol ( $sb$ )

$sb \leftarrow sb$

page 4.

5. b act

```
key ← input()
n ← keys.index(key)
i ← (n-1) (n-1)/nc + 1
j ← (n-1) % nc + 1.
return (i, j, sb).
```

6. train (epochs, print every n ← 500)  $\forall i \in [1, \text{turns}]$

```
p1 ← A1 player ( $\epsilon \leftarrow 0.01$ )
p2 ← A2 player ( $\epsilon \leftarrow 0.01$ )
j ← judge(p1, p2)
p1 win ← 0.0
p2 win ← 0.0
```

$\forall i \in [1, \text{epochs}]$

winner w ← j play ( $\beta$  print ← 0)

if  $w = 1$   
p1 win ++

if  $w = -1$   
p2 win ++

if  $i \% \text{print} \neq 0$

print ( $i, p1win/i, p2win/i$ )

end if  
p1 ~~backup~~, p2 ~~backup~~, j reset.

7. compete (turns)

```
p1 ← A1 player ( $\epsilon \leftarrow 0$ )
p2 ← A2 player ( $\epsilon \leftarrow 0$ )
j ← judge(p1, p2)
p1 load policy, p2 load policy.
p1 win ← 0, p2 win ← 0
```

w ← j play
if  $w = 1$ , p1 win ++
if  $w = -1$ , p2 win ++
j reset.

end if.

print (turns, p1win/turns, p2win/turns,

8 play (human ~~symbol~~  $\leftarrow 1$ )  
loop: ~~hsb~~  $\in \{1, -1\}$

p1 if  $hsb = 1$ , p1 ← human player  
p2 ← A1 player  
else p1 ← A1 player, p2 ← human player.

p1.load policy, p2.load policy.  
j ← judge(p1, p2)  
w ← j play.

if  $w = hsb$  print (you win)  
else if  $w = -hsb$  print (you lose)  
else print (tie).

# 1. Bandit

1.3 init(  $k \leftarrow 1$ ,  $\varepsilon \leftarrow 0$ , initial = 0,  
 $\alpha = 0.1$ , is\_using\_sample\_average = 0,  
 $c \leftarrow 0$ , is\_using\_gradient = 0,  
is\_using\_gradient\_baseline = 0,  
true\_reward = 0 )

$A_s \leftarrow A_s, \text{letters} \in N^k$

time  $t \leftarrow 0$

$\bar{R} \leftarrow 0$ .

## 1.4 reset

vec  $q^* \sim N(\text{tr}, 1) \in R^k$

$Q \leftarrow \text{initial}, \in R^k$

$N \leftarrow 0, \in N^k$

$A^* \leftarrow \underset{\alpha}{\operatorname{argmax}} q^*(\alpha)$

$t \leftarrow 0$

## 1.5 act

① if rand() <  $\varepsilon$

return  $A \leftarrow A_s$  randomly.

// (  $A_s$  is vector, but here is seen as  
a set )

② if  $c \neq 0$ , ( i.e.,  $c > 0$  )

if  $t < k$ , if  $N(a) = 0$ . return  $A \leftarrow a, \forall a$ .

else  $Q_{UCB}(a) \leftarrow Q(a) + c \sqrt{\frac{\log t}{N(a)}}, \forall a$

return  $A \leftarrow \underset{a}{\operatorname{argmax}} Q_{UCB}(a)$  randomly.

③ if is\_using\_gradient

$H \leftarrow Q$

$e^H \leftarrow e^H$

$\pi \leftarrow \frac{e^H}{\sum e^H} \in R^k$

return  $A \leftarrow \text{draw from } A_s$   
by  $\pi$ .

④  $A \leftarrow \operatorname{argmax} Q(a)$  randomly.  
return  $A^*$ .

## 1.6 step( action: A )

$R \leftarrow N(q^*(A), 1)$ ,

$t \leftarrow t+1, N(A) \leftarrow N(A)+1$ .

$\bar{R} \leftarrow \bar{R} + \frac{1}{t}(R - \bar{R})$

If is sample average

$Q(A) \leftarrow Q(A) + \frac{1}{N(A)}(R - Q(A))$

else if is\_using\_gradient

onehot  $\leftarrow 0, \in N^k$

onehot(A)  $\leftarrow 1$ .

if is\_u-g-b., baseline  $\leftarrow \bar{R}$   
else baseline  $\leftarrow 0$ .

$H \leftarrow H + \alpha(R - \text{baseline})(\text{onehot} - \pi)$

else  $Q(A) \leftarrow Q(A) + \alpha(R - Q(A))$

return  $R$ .

2. simulate(runs ← 2000, time ← 1000, bandits)

rewards rs ← 0 ∈  $R^{|\text{bandits}| \times \text{runs} \times \text{time}}$

best action counts bac ← 0, same shape with rs.

~~for~~  $\forall i, \text{bandit} \in \text{bandits}$

$\forall \text{run} \in [1, \text{runs}]$ ,

bandit.reset

$\forall t \in [1, \text{time}]$

a ← bandit.act

r ← bandit.step(a)

rs(i, run, t) ← r

if a = bandit.best\_action A\*

bac(i, run, t) ← 1

end  $\forall t$ .

end  $\forall i, \text{run}$ .

mean best action counts mbac ← bac.mean(axis ← 2)

mean rewards mrs ← rs.mean(axis ← 2)

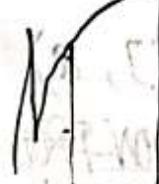
return mbac, mrs, ∈  $R^{|\text{bandits}| \times \text{time}}$ .

3. bandit ← Bandit( $\epsilon=0$ , UCB para c ← 2,

is\_using\_sample\_average ← 1)

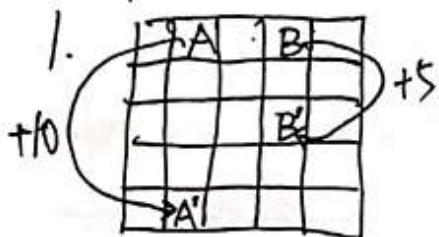
average rewards rs ← simulate(2000, 1000, bandit)

plot(rs).;

UCB: 

( $\epsilon$ )

## Gridworld



1.3 reward  $\in \{-1, 0, 5, 10\}$

action  $\in \{\text{left, up, right, down}\}$

1.4 suppose the agent selects all four actions with equal probability in all states.

$v_h$  is the value function for this random policy, for the discounted reward case with  $\gamma = 0.9$ .

$v_h$  can be computed by solving the system of linear equations,

i.e., Bellman equation for  $v_h$ .

$$v_h(s) \doteq \underset{\pi}{E}(G_t | S_t = s) \\ = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) (r + \gamma v_h(s')), \quad \forall s \in S$$

2.  $X_m \leftarrow 5, Y_m \leftarrow 5$

~~Pos\_A~~  $\leftarrow (2, 5)$

$Pos\_A' \leftarrow (2, 1)$

$Pos\_B \leftarrow (4, 5)$

$Pos\_B' = (4, 3)$

$R_A = 10, R_B = 5$

$\gamma = 0.9$

$A = \{(1, 0), (0, 1), (1, 0), (0, -1)\}$

Action\_Prob  $\leftarrow 0.25$

2.3 step(states, action a)

if  $s = Pos\_A$

return  $Pos\_A', R_A$ .

if  $s = Pos\_B$

return  $Pos\_B', R_B$

$s' \leftarrow s + a$

$x, y \leftarrow s'$

if  $x < 1$  or  $x > X_m$  or  $y < 1$  or  $y > Y_m$

1.5 let  $v_*$  be optimal value function.

reward  $r \leftarrow -1$

We can solve the Bellman optimality equation for  $v_*$ .

$$v_*(s) = \max_{a \in A(s)} \sum_{s', r} p(s', r | s, a) (r + \gamma v_*(s'))$$

$s' \leftarrow s$

else  $r \leftarrow 0$   
return  $s', r$ .

3. figure 3.2(c)

Mat  $v \leftarrow 0, \in R^{X_m \times Y_m}$

Mat  $v' \leftarrow v$

loop:

$v' \leftarrow 0$

$\forall (i, j) \in [1, X_m] \times [1, Y_m],$

$s \leftarrow (i, j)$

$\forall a \in A$

$s', r \leftarrow \text{step}(s, a)$

$\pi_{a|s} \leftarrow \text{Action-Pr}$

$v'(s) += \pi_{a|s} \cdot (r + \gamma v(s'))$

// Bellman equation for  $v'$

end  $\forall a$

end  $\forall i, j$ .

if  $\|v - v'\|_1 < 10^{-4}$

draw  $v$ .

stop

$v \leftarrow v'$

page 2.

4. figure 3.5(c)

Mat  $v \leftarrow 0(X_m, Y_m)$

Mat  $v' \leftarrow v$

$vs \leftarrow 0 \in R^{|A|}$

loop:

$v' \leftarrow 0$

$\forall s=(i, j) \in [1, X_m] \times [1, Y_m]$

$vs \leftarrow 0$

$\forall a \in A$

$s', r \leftarrow \text{step}(s, a)$

$vs(a) \leftarrow r + \gamma v(s')$

end  $\forall a$

$\hat{v}'(s) \leftarrow \max_a (vs(a))$

// value iteration.

end  $\forall s$ .

if  $\|v - v'\|_1 < 10^{-4}$

draw  $v$

stop

$v \leftarrow v'$

5. draw  $v$ :

$v\_draw \leftarrow v$ .

$v\_draw$  rotate 90° clockwise

$v\_draw$  axis show.

## Gambler's problem

1.

The problem can be formulated as all corresponding to ties for an undiscounted, episodic, finite MDP. the argmax action selection  $s \in S = \{1, 2, \dots, 99\} \cup \{0, 100\}$  with respect to the optimal value function.  $s=0$  and  $s=100$  are two dummy states corresponding to termination, giving them values of 0 and 1 respectively.

2.  $\text{Goal} = 100$ ,

$$a \in A(s) = \{0, 1, \dots, \min(s, 100-s)\} \quad S = \{0, 1, \dots, \text{Goal}\},$$

$$r \in \{0, 1\}.$$

$$S^+ = \{1, 2, \dots, \text{Goal}-1\}$$

$$S^+ = \{1, 2, \dots, 99\},$$

$$p_h = 0.4$$

$$S^+ \cup \{0, 100\} = S.$$

3. figure 4.3()

1.3 the state-value function  $V(s), s \in S$  Map  $V: S \rightarrow \mathbb{R}$ , gives the probability of winning from each state.

$$V(S) \leftarrow 0, V(\text{Goal}) = 1$$

$$\text{sweeps\_history} \leftarrow \emptyset$$

The optimal policy maximizes the probability loop: // value iteration of reaching the goal.

$$\text{sweeps\_history}.add(V)$$

$$\Delta \leftarrow 0$$

1.4. let  $p_h$  denote the probability of the coin coming up heads. If  $p_h$  is known, the

$$v \leftarrow V(s)$$

the entire problem is known and it can be solved, for instance, by value iteration.

$$\text{Vector } A \leftarrow A(S) = \{0, \dots, \min(s, \text{Goal}-s)\}$$

1.5 for the case of  $p_h = 0.4$ , there is

$$\text{Vector } v \in \mathbb{R}^{|\mathcal{A}|}$$

a whole family of optimal policies,

$$\forall i_a \in |\mathcal{A}|, (\text{here we mean } i_a \in \{1, \dots, |\mathcal{A}|\})$$

$a \leftarrow A(i_a)$

$$vs(i_a) \leftarrow p_h \cdot V(s+a) + (1-p_h) \cdot V(s-a)$$

end  $\forall i_a$

$$V(s) \leftarrow \max_{i_a} (vs)$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

end  $\forall s$ .

if  $\Delta < 10^{-9}$ : sweephistory add V. break loop.

end loop.

figure:

$\forall V \in \text{sweephistory}$

plot V.

xlabel (Capital)

ylabel (value estimations)

plot ( $\pi$ )

xlabel (Capital)

ylabel (final policy)

// compute optimal policy.

Map  $\pi^* : S \rightarrow N$ ,

$$\pi(s) \leftarrow 0.$$

$$\forall s \in S^+$$

$$\text{Vec } A \leftarrow A(s) = \{0, 1, \dots, \min(s, \text{Goal}-s)\}$$

$$\text{Vec } vs \leftarrow 0, \in R^{|A|}$$

$$\forall i_a \in |A|$$

$$a \leftarrow A(i_a)$$

$$vs(i_a) \leftarrow p_h \cdot V(s+a) + (1-p_h) \cdot V(s-a)$$

end  $\forall i_a$

$$TC(s) \leftarrow A(\operatorname{argmax}_{i_a} (vs)) + 1$$

// avoid choosing  $a=0$  which doesn't change state.

end  $\forall s$ .

## Chapter 4

### 1. Value iteration

$$V = 0 \in R^{|S|}$$

loop:

$$\Delta \leftarrow 0$$

$$\forall s \in S$$

$$v \leftarrow V(s).$$

$$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)(r + \gamma V(s')) \text{ end } \forall s$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

end  $\forall s$ .

if  $\Delta < \theta$  ( $\theta$  a small positive number)

stop loop

end loop

$$\Delta \leftarrow 0$$

$$\forall s \in S$$

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s))(r + \gamma V(s'))$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

$$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)(r + \gamma V(s')) \text{ end } \forall s$$

if  $\Delta < \theta$ , stop loop

3) policy improvement

policy stable  $\leftarrow 1$

$$\forall s \in S$$

$$\text{old } a \leftarrow \pi(s)$$

$\therefore \pi(s) = \arg\max_a \sum_{s',r} p(s',r|s,a)(r + \gamma V(s'))$

$$\pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r|s,a)(r + \gamma V(s'))$$

### 2. correction: Gambler's problem

$$S = \{1, 2, \dots, 99\}$$

$$S^+ = S \cup \{0, 100\}$$

i.e.  $S \leftrightarrow S^+$

if  $\text{old } a \neq \pi(s)$ : policy stable  $\leftarrow 0$

end  $\forall s$ .

if policy stable

stop and return  $V \approx V_*$ ,

$$\pi \approx \pi_*$$

### 3. Policy iteration (using iterative policy evaluation).

1)  $V(s) \leftarrow R$ ,  $\pi(s) \in A(s)$ ,  $\forall s \in S$ . go to 2).

initialize  $V$ ,  $\pi$  arbitrarily.

2) policy evaluation

loop:

## Car rental

- max cars at 1:  $M_1 = 20$
- max cars at 2:  $M_2 = 20$ .
- max move:  $M_m = 5$
- expected rental request at  $i$ :  $\lambda_{i1}$
- expected ~~rental~~ returns at  $i$ :  $\lambda_{i2}$   
 $i=1, 2.$
- $\therefore \begin{pmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \end{pmatrix} = \begin{pmatrix} 3 & 3 \\ 4 & 2 \end{pmatrix}$

discount  $\gamma = 0.9$

rent reward  $R_r = R_{rent} = 10$ ,

move cost  $c_m = 2$ .

$$A = \{-M_m, \dots, M_m\}$$

Poisson upper bound  $M_p = 11$

poisson cache : dict

- poisson pr( $n, \lambda$ ) //  $\lambda < 10$   
key =  $n/10 + \lambda + 1$

if key  $\notin$  poisson cache pc

$$pc(key) \leftarrow \text{poisson. pmf}(n, \lambda) \quad ; \quad n_{iv} \leftarrow \min(n'_1, req_1)$$

return pc(key)

- Bellman equation for  $q_1$

Bellman for  $q$ -SA :

input:  $s, a, V$ , is constant returned cars

output:

$$q(s, a) \leftarrow E_a \leftarrow \sum_{s', r} p(s', r | s, a) (r + \gamma V(s'))$$

process:

expected return :  $E_a \leftarrow 0$

$$E_a \leftarrow E_a - |a| \cdot c_m$$

number of cars at  $i$ :  $n_i$

$$n_1 \leftarrow \min(s_1 - a_1, M_1)$$

$$n_2 \leftarrow \min(s_2 + a_2, M_2)$$

$\forall$  rental request  $req_1 \in \{0, \dots, M_p\}$

$\forall req_2 \in \{0, \dots, M_p\}$

$p \leftarrow \text{poisson-pr}(req_1, \lambda_{11})$

•  $\text{poisson-pr}(req_2, \lambda_{21})$

valid rental at  $i$  :  $n_{iv}$

$$n'_1 \leftarrow n_1, \quad n'_2 \leftarrow n_2$$

$$n_{iv} \leftarrow \min(n'_1, req_1)$$

$$n_{iv} \leftarrow \min(n'_2, req_2)$$

$$r \leftarrow (n_{iv} + n_{bv}) R_r$$

$$n'_1 \leftarrow n'_1 - n_{iv}, \quad n'_2 \leftarrow n'_2 - n_{bv}$$

2.3 poisson. pmf :

$$f(n, \lambda) = e^{-\lambda} \cdot \lambda^n / n!$$

$$\text{if } e^\lambda = \sum_{n=0}^{\infty} \frac{\lambda^n}{n!}$$

if is constant returned cars

returned cars  $ret_1 \leftarrow \lambda_{12}$

$ret_2 \leftarrow \lambda_{22}$

$n'_1 \leftarrow \min(n'_1 + ret_1, M_1)$

$n'_2 \leftarrow \min(n'_2 + ret_2, M_2)$

$s' \leftarrow (n'_1, n'_2)$

$E_a \leftarrow E_a + p(r + \gamma V(s'))$

else

//  $E_a \leftarrow E_a + p(r + \gamma \sum_s p'(V(s))$

$E'_a \leftarrow 0$

$\forall ret_1, ret_2 \in \{0, \dots, M_p\}$

$p' \leftarrow \text{poisson\_pr}(ret_1, \lambda_{12}) \cdot \text{poisson\_pr}(ret_2, \lambda_{22})$   $\Delta \leftarrow \Delta V | v - V(s)|$

$n'_1 \leftarrow \min(n'_1 + ret_1, M_1)$

$\forall s \in S, v \leftarrow V(s),$

$V(s) \leftarrow \text{Bellman\_forq-sa}$

$n'_2 \leftarrow \min(n'_2 + ret_2, M_2)$

$(s, \pi(s), V, \text{is constant}_m)$

$s' \leftarrow (n'_1, n'_2)$

//  $V(s) \leftarrow q(s, \pi(s))$

$\Delta \leftarrow \Delta V | v - V(s)|$

$E'_a \leftarrow E'_a + p'V(s')$

end if  $\Delta < 10^{-4}$ : stop loop.

end if  $ret_1, ret_2$ .

end loop. // innerloop.

$\therefore E_a \leftarrow E_a + p(r + \gamma E'_a)$

$qsa \in R^{|\mathcal{A}|}$

end if freq<sub>1</sub>, freq<sub>2</sub>.

$\forall i \in \{1, \dots, |\mathcal{A}|\}$

return  $E_a$ .

$a \in A_{ia}$

4. figure 4-2 (is constant returned cars  $\leftarrow 1$ )

if  $a \in [-s_2, 0] \cup [0, s_1]$

$q_{sa}(i_a) \leftarrow \text{Bellman for } q_{sa}(s, a, V, \text{is constant...})$

else

$q_{sa}(i_a) \leftarrow -\infty$

end  $\forall i_a$

1.  $\pi(s) \leftarrow A(\text{argmax}(q_{sa}))$

if is policy stable and  $a_{old} \neq \pi(s)$

    is policy stable  $\leftarrow 0$

end  $\forall s$

if is policy stable

    draw  $V$ , stop

$k \leftarrow k + 1$

5. We formulate this as a  
continuing finite MDP.

## Car rental (synchronous)

1.  $M_1=20, M_2=20, c_m=2.$

3.3 solve:

additional park cost  $c_{park}=4$

iterations  $k \leq 0$

$R_r=10$ , rent reward.

loop:

$$\begin{pmatrix} \lambda_{11} & \lambda_{12} \\ \lambda_{21} & \lambda_{22} \end{pmatrix} = \begin{pmatrix} 3, 3 \\ 4, 2 \end{pmatrix}$$

$V \leftarrow$  policy evaluation ( $V, \pi$ )

is policy change,  $\pi \leftarrow$

$$S \in \{0, \dots, M_1\} \times \{0, \dots, M_2\}$$

policy improvement ( $A, V, \pi$ )

poisson cache pc: dict

if is policy change = 0

2. poisson ( $n, \lambda$ )

stop

$$\text{key} \leftarrow n \cdot 10 + \lambda$$

$$k \leftarrow k+1$$

if key  $\notin$  pc

3.4 (out-place)

$$pc(\text{key}) \leftarrow e^{-\lambda} \cdot \lambda^n / n!$$

policy evaluation ( $V, \pi$ )

return pc(key)

process:

3. class: policy iteration

$V_{\text{new}}: S \rightarrow R$

3.2 init { is solve 4-5  $\leftarrow 0$ }

loop:

$$M_p \leftarrow 8$$

$$V_{\text{new}}(s) \leftarrow V_{\pi \text{-} \text{pe}}(\pi, V, s), \forall s \in S$$

$$\delta_M \leftarrow 10^{-1}$$

$$\text{if } \|V_{\text{new}} - V\|_1 < \delta_M$$

$$\gamma \leftarrow 0.9$$

$V \leftarrow V_{\text{new}}$ , return  $V$ .

$$A \leftarrow \{-M_2, \dots, 0\} \cup \{0, \dots, M_1\}$$

$V \leftarrow V_{\text{new}}$

Map  $V: S \rightarrow R$ ,  $V(S) \leftarrow 0$

3.5. policy improvement ( $A, V, \pi$ )

Map  $\pi: S \rightarrow A \subset Z$

$\pi_{\text{new}}: S \rightarrow A$

is solve 4-5  $\leftarrow$  is solve 4-5

expected action returns qsa:  $S \times A \rightarrow R$

$q_{\pi}(s, a) \leftarrow V \leftarrow q_{\pi}(s, a) - \pi(V, s, a)$        $p \leftarrow \text{poisson}(\text{req}_1, n_1) \cdot \text{poisson}(\text{req}_2, n_2)$   
 $\forall a \in A, s \in S.$        $n'_1 \leq n_1, n'_2 \leq n_2$

$\pi_{\text{new}}(s) \leftarrow \text{argmax}(q_{\pi}(s, :)), \forall s \in S.$        $n_{uv} \leftarrow n'_1 \wedge \text{req}_1, n_{bu} \leftarrow n'_2 \wedge \text{req}_2$

if  $\pi_{\text{new}} \neq \pi$        $r \leftarrow (n_{uv} + n_{bu})R_r$

$\pi \leftarrow \pi_{\text{new}}, \text{return } 1$

if is solve 4-5

return 0 // is stable, is changed = 0.      if  $n'_i \geq 10$

3.6. //  $O(n^4)$  for all req, ret.

Bellman( $V, s, a$ ) for  $q(s, a)$

process:

$E_a \leftarrow 0$

$n'_1 \leftarrow n'_1 - n_{uv}, n'_2 \leftarrow n'_2 - n_{bu}$

if is solve 4-5

$E'_a \leftarrow 0$

if  $a > 0$

$\forall \text{ret}_1, \text{ret}_2 \in \{0, m, M_p\}$

// free shuttle to loc2.

$p' \leftarrow \text{poisson}(\text{ret}_1, \lambda_{12}) \cdot \text{poisson}(\text{ret}_2, \lambda_{22})$

$E_a \leftarrow E_a - (a-1)c_m$

$n'_1 \leftarrow (n'_1 + \text{ret}_1) \wedge M_1,$

else

$n'_2 \leftarrow (n'_2 + \text{ret}_2) \wedge M_2$

$E_a \leftarrow E_a - |a|c_m$

// classical Bellman equation

else

// for  $V(s)$

$E'_a \leftarrow E'_a + p'V(s')$ ,

here  $s' = (n'_1, n'_2)$

$n_1 \leftarrow (s_1 - a) \wedge M_1$

end  $\forall \text{ret}_1, \text{ret}_2$ .

$n_2 \leftarrow (s_2 + a) \wedge M_2$

$\therefore E_a \leftarrow E_a + p(r + \gamma E'_a)$

$\forall \text{req}_1, \text{req}_2 \in \{0, \dots, M_p\}$

end  $\forall \text{req}_1, \text{req}_2$ .

return  $E_a$ .

3.7  $v_{\pi\text{-spe}}(\pi, V, s) // \pi(s)$

$a \in \pi(s)$ ,

$E_a \leftarrow \text{Bellman}(V, s, a)$

return  $E_a$ .

3.8  $q_{\pi\text{-api}}(V, s, a) // q_{\pi}(s, a)$

if  $a \in [-s_2, 0] \cup [0, s_1]$

$E_a \leftarrow \text{Bellman}(V, s, a)$

else

$E_a \leftarrow -\infty$

return  $E_a$ .

4. test:

solver  $\leftarrow$  Policy Iteration (is\_sole 45 < 1)

solver.solve

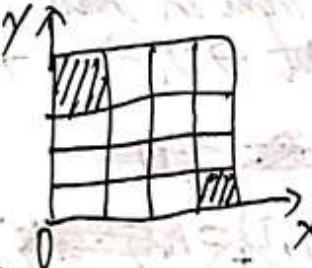
plot  $\pi$ .

## grid world

1.  $X_m \leq 4, Y_m \leq 4$   
 $A = \{(0, -1), (-1, 0), (0, 1), (1, 0)\}$

2.  $\backslash$  is terminal ( $s$ )  
if  $s = (1, Y_m)$  or  $s = (X_m, 1)$   
return 1  
else

return 0.



3. step( $s, a$ )

if is terminal ( $s$ )

return  $s, 0$ .

$s' \leftarrow s + a$

if  $s'_1 < 1$  or  $s'_1 > X_m$  or  $s'_2 < 1$  or  $s'_2 > Y_m$

$s' \leftarrow s$

$r \leftarrow -1$

return  $s', r$

4. compute  $V$  (in place  $\leq 1, \gamma \leq 1$ )

Map  $V: S \rightarrow \mathbb{R}$

$R \leftarrow 0$ .

if in place

$V_{old} \leftarrow V$

denote  $V$  as  $V_{new}$

else

$V_{new} \leftarrow V$

denote  $V$  as  $V_{old}$

loop:

if in place

$V_{old} \leftarrow V$

$\forall s \in S,$

$v \leftarrow 0$

$\forall a \in A,$

$s', r \leftarrow \text{step}(s, a)$

$v \leftarrow v + \pi(s) \left( r + \gamma \cdot V(s') \right)$

end  $\forall a$ .

$V_{new}(s) \leftarrow v$

end  $\forall s$ .

$\Delta \leftarrow \max(|V_{old} - V_{new}|)$

if not in place

$V \leftarrow V_{new}$

if  $\Delta < 10^{-4}$ ; stop loop

$k \leftarrow k+1$ , end loop.

return  $V, k$ .

# 11. Blackjack

## Blackjack

5.  $\pi_d: N \rightarrow A$

$A = \{H, S\}^2 \cup \{A_h, A_s\}$ ,  
(hit, stick)

$$a_h \leftarrow 1, a_s \leftarrow 2$$

$\text{Vec}(2) \pi_d: \{1, \dots, 11\} \leftarrow a_h$

$\pi_d: \{12, \dots, 16\} \rightarrow a_h$

$\{17, \dots, 21\} \rightarrow a_s$ .

$\pi_p$

~~2.  $\pi_p: N \rightarrow A$~~

6. get\_cards()

$\text{Vec}(2), \pi_p: \{1, \dots, 11\} \rightarrow a_h$  card ~ uniform(1, 13)

$\pi_p: \{12, \dots, 19\} \rightarrow a_h$

$$\text{card} = \text{card} \times 10$$

$\{20, 21\} \rightarrow a_s$

return card.

3.  $\pi_p: S \rightarrow A$

7. card\_value(card)

asable-ace-player  $a_{sp}$ ,

if card = 1 return 11

player sum  $s_p$

else return card

dealer card  $c_d$

8. play( $\pi$ , initial  $s_0$ , initial  $a_0$ )

$\leftarrow s$

$s_0 \leftarrow \text{none}$   $a_0 \leftarrow \text{none}$

return  $\pi_p(s_p)$

//  $\pi$  is policy for player.

4.  $\mu_p: S \rightarrow A$

sum of player  $\sum_{i=1}^{sum_p} = 0$

$a_{sp}, s_p, c_d \leftarrow s$

trajectory of player  $traj_p \leftarrow \text{none}$

if binomial(1, 0.5) = 1

$a_{sp} \leftarrow \text{hit} \quad 0.5 \cdot 0 \text{ or } 1$

return  $a_s$

$\{c_{1d} = 0, c_{2d} = 0\}$  (dealer states)

else return  $a_h$ .

$sum_p = 0$

4.3. binomial( $n, p$ )

if  $s_0$  is none:

$$= \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k}$$

while  $sum_p < 12$ :

$sum_p$

card ← get card

sum<sub>p</sub> ← sum<sub>p</sub> + card\_value(card)

if card = 1 // if sum > 21

    sum<sub>p</sub> = sum<sub>p</sub> - 10

else

    up<sub>p</sub> = 1

    if card = 1 and up<sub>p</sub> = 1  
        up<sub>p</sub> = 0  
    else  
        up<sub>p</sub> = 1

    if up<sub>p</sub> = 1 (card = 1)

end while s<sub>p</sub> < 12

G<sub>d</sub> ← get card.

G<sub>2d</sub> ← get card

else // s<sub>0</sub> is not none.

    up<sub>p</sub>, sum<sub>p</sub>, G<sub>d</sub> ← s<sub>0</sub>

    G<sub>2d</sub> ← get card

    assert sum<sub>p</sub> ≤ 21

    s<sub>0</sub> ← (up<sub>p</sub>, sum<sub>p</sub>, s<sub>p</sub>, G<sub>d</sub>)

    sum<sub>d</sub> ← card\_value(G<sub>d</sub>) + card\_value(G<sub>2d</sub>)

    black = { i | i in {G<sub>d</sub>, G<sub>2d</sub>} }

assert sum<sub>d</sub> = 21

// now we have sum<sub>d</sub> ∈ {1, ..., 21}

// game starts!

// player's turn.

loop: s ← (up<sub>p</sub>, sum<sub>p</sub>, q<sub>d</sub>)

if initial action a<sub>0</sub> ≠ none

    s ← (up<sub>p</sub>, sum<sub>p</sub>, a<sub>0</sub>)

initial action a<sub>0</sub> ← none

else

    s ← (up<sub>p</sub>, sum<sub>p</sub>, a<sub>0</sub>)

    a ← π<sub>s</sub>(up<sub>p</sub>, sum<sub>p</sub>, a<sub>0</sub>)

    traj<sub>p</sub>.add(s, a)

    if a = a<sub>0</sub>: break loop

    card ← get card

    out<sub>dice</sub> ← int(a<sub>0</sub>)

    if card = 1

        count\_ace ← count\_ace + 1

    sum<sub>p</sub> ← sum<sub>p</sub> + card\_value(card)

If sum<sub>d</sub> > 21

assert sum<sub>d</sub> = 22. (l, ∵ dealer only has 2 cards)

    sum<sub>d</sub> = sum<sub>d</sub> - 10.

$sum_p \leftarrow sum_p + \text{card value}(\text{card})$

if  $sum_d > 21$

if  $\text{card} = 1$

if  $uap = 1$

$if uap = 1$

$sum_d \leftarrow sum_d - 10$

$sum_p = sum_p - 10$  (we can't have

~~more than one~~ ~~useable ace(s)~~)

more than one if  $sum_d > 21$

$uad = 0$

else

$uap = 1$

return  $(s_0, 1, \text{traj}_p)$

if  $sum_p > 21$

end loop (dealer)

if  $uap = 1$

$sum_p = sum_p - 10$

$uap = 0$

// compare  $sum_p$  and  $sum_d$ .  
(here we have  $sum_p \leq 21$ ,  $sum_d \leq 21$ )

if  $sum_p > 21$

if  $sum_p > sum_d$

return  $(s_0, 1, \text{traj}_p)$

~~R~~  $\leftarrow 1$   
return ~~s\_0, 1, traj\_p~~

end loop (player)

else if  $sum_p = sum_d$

~~R~~  $\leftarrow 0$   
return ~~s\_0, 0, traj\_p~~

else ~~R~~  $\leftarrow -1$

return ~~s\_0, R, traj\_p~~

// dealer's turn.

loop:

$a \leftarrow \pi_d(sum_d)$

if  $a = a_s$ : break loop

card  $\leftarrow \text{get card}()$ ,  $sum_d \leftarrow sum_d + \text{card value(card)}$

if  $card = 1$

if  $uad = 1$ :  $sum_d \leftarrow sum_d - 10$

else:  $uad = 1$

9. MC complete with on-policy

random (see book)  
 $\pi(s) \leftarrow \operatorname{arg\max}_a Q(s,a)$

MC policy evaluation ~~(episodes)~~

end  $V(s,a)$ .

end  $V$ .

return  $Q, \pi$ .

$\forall i \in \{1, \dots, \text{episodes}\}$

~~A, traj<sub>p</sub> ← play( $\pi_p$ , none, none)~~ control

~~A(s, a) ∈ traj<sub>p</sub> // MC off-policy (episodes)~~

~~traj<sub>p</sub>, C<sub>p</sub>~~

$n(s) \leftarrow n(s) + 1$

MC-policy estimate state-value

$V(s) \leftarrow V(s) + \frac{1}{n(s)}(G - V(s))$  off-policy (episodes)

return  $V$

//  $V_{n+1}(s) = V_n(s) + \frac{1}{C_n}(G_n - V_n(s))$

10. // MC with exploring starts.

//  $V_{n+1}(s) = V_n(s) + \frac{1}{n}(W_n G_n - V_n(s))$

MC-es (episodes)

$s_0 = (\alpha_0 = 1, s_0 = 1, C_0 = 2)$

$Q(s,a) \leftarrow 0$

~~N ← 0, V ← 0, W ← 0, C ← 1, OS, WS.~~

$\pi \leftarrow \pi_p$

~~A, traj<sub>p</sub> ← play( $\mu_p, s_0, \text{none}$ )~~

$n(s, a) \leftarrow 0$

$\forall i \in \{1, \dots, \text{episodes}\}$

~~if  $a \in \pi_p(s)$~~

$s_0 \leftarrow \text{random.choice}(S)$

~~if  $a = \pi_p(s)$~~

$a_0 \leftarrow \text{random.choice}(A)$

$N \leftarrow N + \frac{1}{2} \cdot \mu_p(a|s), \text{ here } = 0.5$

~~A, traj<sub>p</sub> ← play( $\pi, s_0, a_0$ ) else~~

$A(s, a) \in \text{traj}_p$

$W \leftarrow 0, \text{break } A(s, a)$

$n(s, a) \leftarrow n(s, a) + 1$

~~end  $A(s, a)$~~

$Q(s, a) \leftarrow Q(s, a) + \frac{1}{n(s, a)}(G - Q(s, a))$  ④

$V_0 \leftarrow V_0 + \frac{1}{n}(W G - V_0)$  ⑤

return OS, WS

$V_W \leftarrow V_W + \frac{W}{C}(G - V_W)$  ⑥

$Q(s(i)) \leftarrow V_0, W(s(i)) \leftarrow N, \text{ end } t_i$

## 12. figures\_1()

$V_1 \leftarrow \text{MC\_policy\_evaluation\_on\_policy}(10^4)$

$V_2 \leftarrow \text{MC\_policy\_evaluation\_on\_policy}(5 \times 10^5)$

$V_1(:, :, u_{ap}=1). \log()$

14. figure 5\_3()

$V_1(:, :, u_{ap}=0). \log()$

true\_value = -0.27726

$V_2(:, :, u_{ap}=1). \log()$

episodes =  $10^4$

$V_2(:, :, u_{ap}=0). \log()$

rans = 100

## 13. figure 5\_2()

$e_{ew} \leftarrow 0$  (episodes),  $e_{ew} \leftarrow 0$  (episodes)

$Q \leftarrow \text{MC\_es}(5 \times 10^5)$

$\forall i \in \{1, \dots, rans\}$

$V(s) \leftarrow \max_a Q(s, a), \forall s$

OS, WS  $\leftarrow \text{MC\_estimate\_state\_value}$

$f(x_1, x_2, x_3)$

$f(x_1, x_2, x_3)$

- off-policy (episodes)

$g(x_1, x_2) = \max_{x_3} f(x_1, x_2, x_3)$

$g_1 = \max_{x_3} f(x_1, 1, x_3)$

$e_0 \leftarrow e_0 + (OS - \text{true value})^2$

$g(x_1, 1)$

$g_2(x_1, 2)$

$g_2 = \max_{x_3} f(x_1, 2, x_3)$

// vector

$e_w \leftarrow e_w + (WS - \text{true value})^2$

$V(:, :, u_{ap}=1). \log(\text{"optimal value"})$  end  $\forall i$

$V(:, :, u_{ap}=0). \log(\text{"optimal value"})$   $e_0 \leftarrow e_0 / rans$ ,  $e_w \leftarrow e_w / rans$

$\pi(:, :, u_{ap}=1). \log(\text{"optimal policy"})$

plot(e\_0)

$\pi(:, :, u_{ap}=0). \log(\text{"optimal policy"})$

title(ordinary importance sampling)

title(weight importance sampling)

$\pi$  is deterministic policy

$$v_{\pi}(s)$$

$$v_{\pi}, q_{\pi}, \pi, q^*, q_{\pi(s,a)}$$

3.2 the backup diagram

for MC estimation of  $v_{\pi}$ .

1. First-visit MC policy evaluation

$\pi, V, n(s), v(s) \forall s \in S$

Repeat:

generate an episode using  $\pi$

$s$  is in the episode

$G \leftarrow \text{return following the first occurrence of } s$

1

2

3

4

5

$$n(s) \leftarrow n(s) + 1$$

$$V(s) \leftarrow \frac{\sum G(s)}{n(s)}$$

$$V(s) \leftarrow V(s) + \frac{1}{n(s)}(G - V(s))$$

$$V(s) \leftarrow \frac{\sum G(s)}{n(s)}$$

$$V(s) \leftarrow \frac{\sum G(s)}{n(s)}$$

should use  
incremental  
implementation

4. Without a model,  
state values alone are  
not sufficient to  
determine a policy.

2. Blackjack: 2.2 playing blackjack is naturally  
2.3 consider the policy that sticks if the sum is  
player's sum is 20 or 21, and  
otherwise hits.

$$5. \pi_0 \xrightarrow{E} q_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} q_{\pi_1}$$

2.4 to find  $v$  for this policy by a  
MC approach.

$$\xrightarrow{..} \pi_* \xrightarrow{E} q_*$$

2.5 there is no difference between  
first-visit and every-visit MC methods.

2.6 the terminal reward is also the return.

3. entire trajectory of transitions ~~along~~

a ~~single~~ particular single episodes

6. ~~MC control~~:

MC ES ( MC control)

initialize  $Q(s, a) \leftarrow \text{arbitrary}$

$\pi(s) \leftarrow \text{arbitrary}$

~~$s \in A(s, a), n(s, a) \in \mathbb{N}$~~ ,  $n(s, a) \in \mathbb{N}$

~~returns $(s, a) \leftarrow \text{empty list}$~~

$\forall s \in S, a \in A(s)$

7.3: "the episodes are all simulated games, ~~7.7~~  
it is easy to arrange for exploring starts  
that include all possible"

7.4 pick  $s_0$  at random  
with equal probability.

repeat:

① choose  $s_0 \in S, a_0 \in A(s_0)$ , s.t. 7.5 initial policy is the  
all pairs have probability  $> 0$  one in 2.

② generate an episode starting from  $s_0, a_0$ , \* the initial  $Q(s, a) = 0$   
following  $\pi$ . 7.6  $V$  is computed from  $Q$  found by

③  ~~$V$  pairs, a in the episode~~ 7.7:  $\pi^*$ : one  ~~$s, a$~~

$G \leftarrow \text{return following the first}$

occurrence of  $s, a$ .

usable ace:

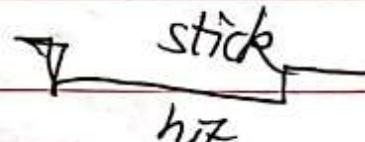
$n(s, a) \leftarrow n(s, a) + 1$

~~$Q(s, a) \leftarrow Q(s, a) + \sum_{t=0}^{T-1} G_t / n(s, a)$~~

~~$\pi(s, a) \leftarrow \frac{\sum_{t=0}^{T-1} G_t / n(s, a)}{n(s, a)}$~~

~~$\pi(s, a) \leftarrow \frac{\sum_{t=0}^{T-1} G_t / n(s, a)}{n(s, a)}$~~

~~$Q(s, a) \leftarrow \frac{\sum_{t=0}^{T-1} G_t / n(s, a)}{n(s, a)}$~~



no usable ace:

~~$Q(s, a) \leftarrow \frac{\sum_{t=0}^{T-1} G_t / n(s, a)}{n(s, a)}$~~



④  $\pi(s) \leftarrow \underset{a}{\operatorname{argmax}} Q(s, a)$ ,  $\forall s$  in the episode

7.8. The policy is the same as  
the basic strategy of Thorp

7. solving Blackjack

7.2 apply MC ES to blackjack. with the ~~the~~ exception of  
the leftmost notch

$$\begin{aligned}
 n(s,a) &\leftarrow n(s,a) + 1 \\
 Q(s,a) &\leftarrow Q(s,a) + \frac{n(s,a)}{\sum a' n(s,a')} (G - Q(s,a)) \\
 \text{sum } G &\leftarrow \text{sum } G + G \\
 n(s,a) &\leftarrow n(s,a) + 1 \\
 Q(s,a) &\leftarrow Q(s,a) + \frac{\text{sum } G(s,a)}{\text{sum } n(s,a')} \\
 \text{sum } G &\leftarrow \text{sum } G + G \\
 \text{sum } n(s,a') &\leftarrow \text{sum } n(s,a') + 1
 \end{aligned}$$

(ESR incremental implementation)

in the policy for a usable ace, which is not present in Thorp's strategy. We are uncertain of the reason for this discrepancy.

8. among  $\epsilon$ -soft policies,  $\epsilon$ -greedy policies are in some sense those that are closest to greedy.  $A^* \leftarrow \arg\max Q(s,a)$
- ~~8.3 next page~~ I does not require that the  $\pi(a|s) \leftarrow \begin{cases} 1-\epsilon & \text{if } a = A^*(s) \\ \frac{\epsilon}{|A(s)|-1} & \text{otherwise} \end{cases}$
9. policy be taken all the way to a greedy policy, only that it be moved toward a greedy ~~10.3~~ ~~next page~~ policy.

10. on-policy first-visit MC control (for  $\epsilon$ -soft policies)

$r(s,a) \leftarrow$  arbitrary.

~~Return(s,a)~~ ~~empty list~~  $n(s,a) \leftarrow 0$

$\pi(a|s) \leftarrow$  an arbitrary ~~soft~~  $\epsilon$ -soft policy

Repeat:

- ① generate an episode using  $\pi$

- ② If pair  $s,a$  in the episode

-  $G \leftarrow$  return following the first ~~occurrence~~ <sup>occurrence</sup> of  $s,a$

8.3  $\epsilon$ -soft policies

$$\Leftrightarrow \pi(a|s) \geq \frac{\epsilon}{|A(s)|}, \forall s, a.$$

for some  $\epsilon > 0$

8.4

any  $\epsilon$ -soft policy  $\pi$ ,  
any  $\epsilon$ -greedy policy with respect  
to  $q_\pi$  is guaranteed to be  
better than or equal to  $\pi$ .

3. On-policy first-visit MC

control (for  $\epsilon$ -soft policies)

init:  $\pi(a|s)$ , an arbitrary  $\epsilon$ -soft policy

10.3

In the on-policy method we  
move it only to an  $\epsilon$ -greedy policy.

11. off-policy prediction via  
importance sampling

11.3 In practice, the weighted importance-sampling estimator usually has dramatically lower variance and is strongly preferred.

Repeat:

$$\pi'(a|s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A(s)|} \\ \frac{\epsilon}{|A(s)|} \end{cases}$$

10.3 here  $\pi(a|s)$  is  $\epsilon$ -soft policy,

(not greedy policy)

If  $\pi'$  is partial  $\epsilon$ -greedy policy,

$$q_{\pi'}(s, \pi'(s)) = \begin{cases} v_n(s), \text{ if } \pi'(a|s) = \pi(a|s) \\ v_n(s), \text{ if } \pi'(a|s) \neq \pi(a|s) \end{cases}$$

$$\Rightarrow q_{\pi'}(s, \pi'(s)) = \begin{cases} v_n(s) \text{ if } \pi'(a|s) = \pi(a|s) \\ \geq v_n(s) \text{ if } \pi'(a|s) \neq \pi(a|s) \end{cases}$$

11.4

$$V_t(s) = \frac{\sum_{t \in J(s)} p_t^{T(t)} G_t}{|J(s)|}$$

ordinary importance sampling.

11.5 weighted importance sampling:

$$V_t(s) = \frac{\sum_{t \in J(s)} p_t^{T(t)} G_t}{\sum_{t \in J(s)} p_t^{T(t)}}, \quad V_t(s) = 0 \text{ if } \sum_{t \in J(s)} p_t^{T(t)} = 0$$

$$11.6. p_t^T = \dots = \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)}$$

11.7 12. off-policy estimation of a blackjack state value.

12.2.

4

$$Q_t(s, a) = \frac{\sum_{t \in J(s, a)} p_t^{T(t)} G_t}{\sum_{t \in J(s, a)} p_t^{T(t)}}$$

12.2. We evaluate the state in which the dealer is showing a deuce, the sum of the player's cards is 13, and the player has a usable ace.

12.3. the data ~~is~~ generated by starting in this state then choosing to hit or stick at random with equal probability (the behavior policy). The target policy is to stick only on a sum of 20 or 21.

$$V_3 = \frac{G_1 + G_2}{2}$$

$$\therefore V_{n+1} = \frac{G_{n+1} + G_n}{n}$$

$$\therefore V_{n+1} = G_1 + \dots + G_n$$

$$\therefore V_{n+1} = \frac{(n-1)V_n + G_n}{n}$$

$$= \frac{nV_n + G_n - V_n}{n}$$

$$= V_n + \frac{1}{n}(G_n - V_n)$$

~~$$B.3. \quad \therefore V_{n+1} = \frac{\sum W_k G_k}{\sum W_k}, n \geq 1 \quad (L)$$~~

$$\therefore C_n V_{n+1} = \sum W_k G_k$$

$$= C_n V_n + W_n G_n, C_0 = 0$$

~~$$\therefore V_{n+1} = \frac{C_n V_n + W_n G_n - W_n V_n}{C_n}$$~~

$$\therefore V_{n+1} = V_n + \frac{W_n}{C_n}(G_n - V_n) \quad (n \geq 1)$$

here  $V_1$  is arbitrary.

~~$$B.4. \quad V_2 = \frac{W_1 G_1}{W_1} = G_1, \quad 13.4. \quad \text{let } W_1 = 1, \therefore C_1 = 1$$~~

$$V_3 = \frac{W_1 G_1 + W_2 G_2}{W_1 + W_2}$$

$$\therefore V_{n+1} = V_n + \frac{1}{n}(G_n - V_n), n \geq 1$$

$$13.5 \quad V_{n+1} = \frac{\sum W_k G_k}{n} \Rightarrow V_{n+1} = V_n + \frac{1}{n}(W_n G_n - V_n)$$

~~$Q \rightarrow Q_{\pi}$~~  / If incremental off-policy (1) : If nongreedy actions  
 every-visit MC policy evaluation. are common, then learning  
 process:  $Q(s, a) \leftarrow \text{arbitrary}, C(s, a) \leftarrow 0$  will be slow, particularly for  
 $\mu(a|s) \leftarrow \text{an arbitrary soft behavior policy}$  states appearing in the  
 $\pi(a|s) \leftarrow \text{an arbitrary target pol.} \forall s \in S, a \in A(s)$  early portions of long episodes.

Repeat:

1) generate an episode using  $\mu$ : 15. off-policy every-visit MC control  
 $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ . (returns  $\pi \approx \pi^*$ ) ~~(15)~~

2)  $A \leftarrow 0$

3)  $W \leftarrow 1$

4)  $\forall t = T-1, T-2, \dots, 0$

$$A \leftarrow \gamma A + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W \quad \text{loop:}$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}))$$

$$W \leftarrow W \cdot \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)}$$

process:

$$Q(s, a) \leftarrow \text{arbitrary}, C(s, a) \leftarrow 0.$$

$\pi(s) \leftarrow \text{a deterministic policy}$   
~~(15)~~ that is greedy w.r.t.  $Q$ .

If  $W=0$  : exit for loop.

① generate an episode using any soft policy  $\mu$ :

$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$ .

15. off-policy MC control

15.3  $\pi$  converges to optimal at all encountered states.

②  $A \leftarrow 0$

③  $W \leftarrow 1$

④  $\forall t = T-1, T-2, \dots, 0$

15.4 : off-policy MC control method

learns only from the tails of episodes.

$$A \leftarrow \gamma A + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{N}{C(S_t, A_t)}(G - Q(S_t, A_t))$$

$\pi(S_t) \leftarrow \operatorname{argmax} Q(S_t, a)$  (with ties broken consistently)

if  $A_t \neq \pi(S_t)$  then exit for loop.

$$N \leftarrow N \frac{1}{P(A_t | S_t)}$$

( $\varepsilon$ )  
why put  
here, not  
last  
speed up convergence

# random walk

true value  $\hat{V}_\pi$ .

$$1. \quad V_\pi(C) = 0.5$$

$$V_\pi(i) = \frac{i}{6}, \quad i=1, \dots, 5.$$



$$0 \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad s \in S$$

1.5 All episodes start in C, and proceed either left or right by one state on each step, with equal probability.

~~$$1.3 \quad \hat{V}(S) = 0.5$$~~

~~$$1.4 \quad \hat{V}(S) = 0, \quad V(0)=0$$~~

~~$$2. \quad V_\pi : S \rightarrow R, \quad 1.5 \quad \text{start in}$$~~

true value

$$V_\pi(s) = \frac{s}{6}, \quad s=1, \dots, 5, \quad V_\pi(0)=0$$

~~$$2.3 \quad \text{action left } a_l = 0$$~~

~~$$\text{action right } a_r = 1$$~~

~~$$3. \quad \text{temporal difference TD}(0)$$~~

~~$$(V, \alpha=0.1, \text{is batch}=0)$$~~

state  $s=3$ .

trajectory traj add s.

~~rewards add()~~

loop:

~~$$x \sim \text{binomial}(1, 0.5)$$~~

if  $x=a_l$

$$s' = s - 1$$

~~$$\text{else } s' = s + 1$$~~

~~$$\text{if } s'=6, r=1$$~~

~~$$\text{else } r=0$$~~

~~assume all rewards are 0~~

~~if is batch = 0~~

$$V(s) \leftarrow V(s) + \alpha \cdot (r + V(s') - V(s))$$

~~if traj.add(s), rs.add(r)~~

~~if  $s=6$  or  $s=0$~~

break loop

~~rewards add~~

end loop

return traj, rewards.

~~$$4. \quad \text{constant-}\alpha \text{ MC}$$~~

~~$$(V, \alpha=0.1, \text{is batch} \neq 0)$$~~

$$s=3$$

traj add s.

loop:

$$x \sim \text{binomial}(1, 0.5)$$

if  $x = a_1$

$s' = s - 1$

else

$s' = s + 1$

too add  $s'$

if  $s' \neq 0$

return  $a_i = 1$

else if  $s = 0$

$a_i = 0$

break loop.

$s \leftarrow s'$

end loop.

if is batch = 0

$V(s) \leftarrow V(s) + \alpha \cdot (G - V(s))$ ,  $\forall s \in \text{traj}$ ,  $\alpha \leftarrow \alpha_s$ , method  $\leftarrow MC$

return traj,  $G \leftarrow \max_{s \in G} r_s$ , rs.setSize(traj),  $r \in [r_{\min}, r_{\max}]$

5. example 6.2 left.

compute state value:

episodes = {0, 1, 10, 100}

$V \leftarrow V_0$

$\forall i \in \{0, m, 100\}$

if  $i \in \text{episodes}$

plot  $V$ , with title:  $i + \text{episodes}$ .

$TD(0)(V)$

end if

plot  $V_\pi$

6 example 6.2 right

rms\_error():

$TD(0)\text{-ds} = \{0.05, 0.1, 0.05\}$

$MC\text{-ds} = \{0.01, 0.02, 0.03, 0.04\}$

episodes = 10

$n \times S = 100$ , total errors  $E = 0$  (episodes)

$\forall i \in \{0, \dots, 9\} | TD(0)\text{-ds} + MC\text{-ds} |$

$\alpha \leftarrow \alpha_s(i)$

if  $i \leq | TD(0)\text{-ds} |$

method  $\leftarrow TD$

else

$\alpha \leftarrow \alpha_s(i)$

method  $\leftarrow MC$

errors  $\in \mathbb{R}^{\text{episodes}}$

$V \leftarrow V_0$

$\forall j \in \{0, m, \text{episodes}\}$

$\text{error} = \sqrt{\frac{(V_j - V_\pi)^2}{S(\text{size})}}$ , here = 5

errors.add(error)

if method = "TD"

$TD(0)(V, \alpha)$

else

$MC(V, \alpha)$

end if

totalErrors  $\leftarrow$  totalErrors + errors.

end for

totalErrors  $\leftarrow$  totalErrors / runs.

plot( totalErrors , name: method + 'x' )

7. // figure 6.3

batch\_updating( method, episodes,  $\alpha=10^3$  )

runs = 100

totalErrors E  $\leftarrow$  0(episodes)

$\forall r \in \{1, \dots, \text{runs}\}$

$V \leftarrow V_0$

errors  $\in \mathbb{R}^{\text{episodes}}$

trajs  $\leftarrow \emptyset$ , rewards rss  $\leftarrow 0$

$\forall ep \in \{1, \dots, \text{episodes}\}$

if method = 'TD'

traj, rs  $\leftarrow$  TD(0)(V, isBatch < 1)

else

traj, rs  $\leftarrow$  constantAlphaMC(V, isBatch < 1)

trajs  $\leftarrow$  add(traj)

rss  $\leftarrow$  add(rs)

loop:

sums  $\leftarrow$  0(\$size)

$\forall (\text{traj}, rs) \in (\text{trajs}, rss)$

$\forall i \in \{1, \dots, |\text{traj}|\}$

if method = 'TD'

~~sums~~  $\leftarrow V^n + (rs(i) + V(\text{traj}(i)))$

see  
10

else

$V^n \leftarrow V^n + (rs(i) - V(\text{traj}(i)))$

end if, end if(traj, rs)

~~sums~~  $\leftarrow$  ~~sums~~

if  $|\text{traj}|, |\text{rss}| < 10^3$

break loop.

$V \leftarrow V + \frac{1}{\alpha} \cdot \text{sums}$

end loop.

err  $\leftarrow \sqrt{\frac{\sum (V - V_{\text{true}})^2}{S.size}}$ , here 5.

errors.add(err)

end if ep.

E  $\leftarrow E + \text{errors}$

end for.

E  $\leftarrow \frac{E}{\text{runs}}$

return E.

8 example\_6-2()

~~return~~)

compute\_state\_value()  
rms\_error().

9. figure 6-3()

episodes = 10

TD\_errors  $\leftarrow$  batch\_updating('TD', episodes)      if  $s=1$  or  $s=6$

MC\_errors  $\leftarrow$  batch\_updating('MC', episodes)      break loop

plot(TD\_errors)

plot(MC\_errors)

end loop

~~for batch~~

10.  $s \leftarrow \text{traj}(i), s' \leftarrow \text{traj}(i+1), r \leftarrow \text{rewards}$

~~if~~  $\delta \leftarrow r + \gamma V(s') - V(s)$

~~sum~~  $\delta(s) \leftarrow \text{sum}(\delta(s)) + \delta$

if method = "TD"

~~r~~  $\leftarrow r_s(i)$

$\delta \leftarrow r + \gamma V(s') - V(s)$

else

~~else~~

$a \leftarrow \text{arg}(s)$

$\delta \leftarrow a - V(s)$

$\text{sum}(\delta(s)) \leftarrow \text{sum}(\delta(s)) + \delta$

11. if  $s' = 6$

$a \leftarrow 1$

else if  $s' = 0$

$a \leftarrow 0$

~~else~~  $a \leftarrow 0$

$s \leftarrow s$

traj.add(s), ~~return~~

~~if s=1 or s=6~~

~~break loop~~

~~end loop~~

~~for batch~~

## TD(0)

$$1. V(S_t) \leftarrow V(S_t) + \alpha (A_t - V(S_t)) \quad // i: Q(s,a): s \in S^+, a \in A(s)$$

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t)) \quad \text{Repeat for each episode:}$$

2. Tabular TD(0) for estimating  $V_\pi$

Input:  $\pi$

initialize  $V(s)$  arbitrarily,  $\forall s \in S^+$  derived from  $Q$  (e.g.,  $\epsilon$ -greedy)  
 e.g.  $V(s)=0, \forall s \in S^+$

Process:

loop for each episode

initialize  $S \in S$

repeat for each step of episode

$A \sim \pi(A|S)$

take  $A$ , observe  $R, S'$

$$V(S) \leftarrow V(S) + \alpha (R + \gamma V(S') - V(S)) \quad S \leftarrow S', A \leftarrow A'$$

$S \leftarrow S'$

④ until  $S$  is terminal.

until  $S$  is terminal.

⑤ end repeat for each episode.

end loop.

$$3. ① S_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

$$6. Q(S_t, A_t) \leftarrow Q(S_t, A_t)$$

$$② A_t - V(S_t) = \sum_{k=0}^{T-t-1} \gamma^k S_{t+k}, \text{ here } + \alpha (R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t))$$

$$A_T = 0, V(S_T) = 0.$$

$$4. Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

5. Sarsa : an on-policy TD control algorithm.

initialize  $Q(s,a)$  arbitrarily,  $\forall s \in S, a \in A(s)$ .  $Q(s, \cdot) = 0, \forall s$  is terminal state.

7. Q-learning: an off-policy TD control algorithm.
- Initialize  $S$   
Repeat for each step of episode  
initialize  $Q(s, a)$  arbitrarily,  $\forall s \in S, a \in A(s)$ , ① choose  $A$  from  $S$  using policy derived from  $Q_1, Q_2$ .  
 $Q(s, \cdot) = 0, \forall s \in S^+ \setminus S$   
 $(e.g., \varepsilon\text{-greedy in } Q_1 + Q_2)$
- Repeat for each episode:  
 initialize  $S \in S$   
 repeat for each step of episode :  
 choose  $A$  from  $S$  using policy derived from  $Q$ ,  $Q_1(S, A) < Q_2(S, A)$   
 $(e.g., \varepsilon\text{-greedy})$   
 take  $A$ , observe  $R, S'$   
 $Q(S, A) \leftarrow Q(S, A) + \alpha(R + \gamma \max_a Q(S', a))$  else  
 $-Q(S, A))$   
 $+ \alpha(R + \gamma Q_2(S', \arg \max_a Q_2(S', a)) - Q_1(S, A))$   
 $S \leftarrow S'$   
 $\vdots$   
 until  $S$  is terminal.  
 end Repeat for each episode.
- ④  $S \leftarrow S'$   
 $+ \alpha(R + \gamma Q_1(S', \arg \max_a Q_1(S', a)) - Q_2(S, A))$
8.  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma E(Q(S_{t+1}, A_{t+1}) | S_t) - Q(S_t, A_t))$   
 $\leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \sum_a \pi(a | S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t))$
9.  $Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha(R_{t+1} + \gamma Q_2(S_{t+1}, \arg \max_a Q_1(S_{t+1}, a)) - Q_1(S_t, A_t))$
10. double Q-learning  
 Initialize  $Q_1(s, a), Q_2(s, a)$  arbitrarily,  $\forall s \in S, a \in A(s)$ .  
 $Q_1(s, \cdot) = Q_2(s, \cdot) = 0, \forall s \in S^+ \setminus S$  ⑤ until  $S$  is terminal  
 Repeat for each episode :  
 end repeat for each episode.

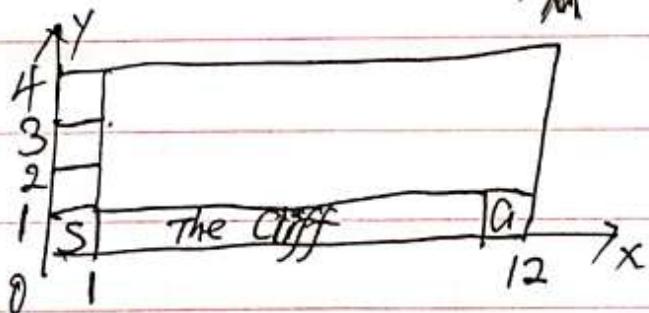
# ~~cliff walking~~

1.  ~~$X_M = 12$~~ ,  $X_M = 12$ ,  $\epsilon = 0.1$ ,  $\alpha = 0.5$ ,  $\gamma = 1$ .

$a_L$ ; ~~left set~~.  $a_r$ ,  $a_u$ ,  $a_d \leftarrow 1, 2, 3, 4$ .

$A = \{a_L, a_r, a_u, a_d\}$

START  $s = (1, 1)$ , GOAL =  $(12, 1)$



2. step(state s, action a):

assert  $s \neq \text{GOAL}$ .

$x, y \leftarrow s$

if  $a = a_L$ :  $x \leftarrow x - 1$

else if  $a = a_r$ :  $x \leftarrow x + 1$

else if  $a = a_u$ :  $y \leftarrow y + 1$

else  $y \leftarrow y - 1$

if  $x < 1$ :  $x \leftarrow 1$

else if  $x > X_M$ :  $x \leftarrow X_M$

if  $y < 1$ :  $y \leftarrow 1$

else if  $y > Y_M$ :  $y \leftarrow Y_M$

$s' \leftarrow (x, y)$

~~$r \leftarrow -1$~~

~~$\text{if } 1 < x < X_M \text{ and}$~~

~~$(x, y) \in [1, X_M] \times [1, Y_M]$~~

~~$s' \leftarrow \text{START}, r \leftarrow -100$~~

~~return  $r, s'$ .~~

$Q(s, a) \in \mathbb{R}$   
 $Q(s, a) \leftarrow -1, \forall s, a$

$Q(s, a) \leftarrow -100, s \in [1, X_M] \times [1, Y_M]$   
 $\times \text{far}$

~~$Q(s, a) \leftarrow -100, s = \text{START}, a = a_L$~~

4 choose action ( $s, Q$ )

~~$\text{sample} \leftarrow \text{binomial}(1, \epsilon)$~~

if  $\text{sample} = 1$



$i_a \sim \text{uniform}(1, 4)$

$a \leftarrow A(i_a)$

else

~~$Q(s \leftarrow Q(s, a))$~~

$a \leftarrow \underset{a}{\operatorname{argmax}} Q(s, a)$  randomly

return  $a$ .

5.  ~~$s \leftarrow \text{start}$~~   $s \leftarrow \text{start}$  ( $Q$ ,  $\epsilon$  expected  $\leftarrow 0$ ,

$\alpha \leftarrow \text{ALPHA}$ )

$s \leftarrow \text{START}$

$a \leftarrow \text{choose action}(s, Q)$

$\text{sum\_r} \leftarrow 0$

loop:

$r, s' \leftarrow \text{step}(s, a)$

$a' \leftarrow \text{choose\_action}(s', Q)$

$\text{sum\_r} \leftarrow \text{sum\_r} + r$

if  $\text{is\_expected} = 0$

~~$t_{\text{tmp}}^a \leftarrow Q(s', a')$~~

else

~~$t_{\text{tmp}}^a \leftarrow 0$~~

$\text{set\_a} \leftarrow \text{argwhere}(Q(s', :) = \max_a Q(s', :))$

~~$A$~~   $\in \text{set\_a}$

if  $a \in \text{set\_a}$

$$p \leftarrow \frac{1-\varepsilon}{|\text{set\_a}|} + \frac{\varepsilon}{|A|}$$

$s \leftarrow s'$

~~$t_{\text{tmp}}^a \leftarrow t_{\text{tmp}}^a + p \cdot Q(s', a)$~~

if  $s = \text{GOAL}$

stop loop

else

$$p \leftarrow \frac{\varepsilon}{|A|}$$

end loop

return  $\text{sum\_r}$ .

$$t_{\text{tmp}}^a \leftarrow t_{\text{tmp}}^a + p \cdot Q(s', a)$$

end if

$$t_{\text{tmp}}^a \leftarrow t_{\text{tmp}}^a + r + \gamma Q$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha(A - Q(s, a))$$

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

if  $s = \text{GOAL}$  : stop loop.

end loop.

return  $\text{sum\_r}$ .

## 6. Q-learning ( $Q$ , $\alpha = \text{ALPHA}$ )

$s \leftarrow \text{START}$

$\text{sum\_r} \leftarrow 0$

loop:

$a \leftarrow \text{choose\_action}(s, Q)$

$r, s' \leftarrow \text{step}(s, a)$

$\text{sum\_r} \leftarrow \text{sum\_r} + r$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \cdot (r + \gamma \max Q(s', :))$$

$$- Q(s, a))$$

$s \leftarrow s'$

if  $s = \text{GOAL}$

stop loop

end loop

return  $\text{sum\_r}$ .

7. print\_optimal\_policy( $Q$ )

7. print optimal policy ( $Q$ )       $\text{sum\_rs\_sarsa} / \approx \text{runs}$   
 $\pi(s) = \underset{a}{\operatorname{argmax}} Q(s, a), \forall s.$        $\text{sum\_rs\_Q\_learning} / \approx \text{runs}.$

print  $\pi$       ,       $\text{plot}_{\text{sum}}(\text{rs\_sarsa}) \text{ name} = "Sarsa"$   
 $\text{plot}_{\text{sum}}(\text{rs\_Q\_learning}), \text{ name} = "Q\_learning"$

i.e. ~~if  $\pi(s)$~~   $a \leftarrow \pi(s)$ ,  
~~if  $s = \text{GOAL}$ , print "G", continue.~~  
~~if  $a = a_i$ : print: " $\leftarrow$ "~~      savefigure("sum\_r")  
~~else if  $a = a_r$ : print: " $\rightarrow$ "~~  
~~else if  $a = a_u$ : print: " $\uparrow$ "~~      print optimal policy ( $Q$ -sarsa)  
~~else : print " $\downarrow$ "~~      print optimal policy ( $Q$ -Q-learning)

## 8. figure 6.4()

episodes = 500

runs = 50

$\text{sum\_rs\_sarsa} \leftarrow \emptyset$  episodes  
 $\text{sum\_rs\_Q\_learning} \leftarrow \emptyset$  episodes  
 $\forall r_{run} \in \{1, \dots, \text{runs}\}$   
 $Q_{\text{sarsa}} \leftarrow R^{X_m \times Y_m \times 4}$   
 $Q_{\text{-Q\_learning}} \leftarrow R^{X_m \times Y_m \times 4}$

$Q_{\text{sarsa}} \leftarrow 0, Q_{\text{-Q\_learning}} \leftarrow 0$

$\forall i \in \{1, \dots, \text{episodes}\}$

$\text{sum\_rs\_sarsa}(i) \leftarrow \text{rs\_sarsa}(i) + \text{sarsa}(Q_{\text{sarsa}})$

$\text{sum\_rs\_Q\_learning}(i) \leftarrow \text{rs\_Q\_learning} + Q_{\text{-Q\_learning}}(Q_{\text{-Q\_learning}})$

end  $\forall i, \text{rs}$

end  $\forall run$

## 9. figure 6.6()

episodes = 1000, runs = 10

$\alpha_S \leftarrow \{0.1, 0.2, \dots, 1\}$

asymptotic-sarsa,  $\alpha_{\text{sarsa}} \leftarrow 1$

$\alpha_{\text{expected sarsa}} \leftarrow$

$\alpha_{\text{-Q\_learning}} \leftarrow$

interim ~~intermediate~~ -sarsa  $\leftarrow 4$

~~initial~~  $\alpha_{\text{expected sarsa}} \leftarrow$

initial  $\alpha_{\text{-Q\_learning}} \leftarrow 6$

performance  $\leftarrow \emptyset$

trun  $\in \{1, \dots, \text{runs}\}$

i  $\in \{1, \dots, |\alpha S|\}$

$\alpha \leftarrow \alpha S(i)$

$Q_{\text{sarsa}} \leftarrow \emptyset^{X_m \times Y_m \times 4}$

$Q_{\text{expected-sarsa}} \leftarrow \emptyset^{X_m \times Y_m \times 4}$

$Q_{\text{-Q-learning}} \leftarrow \emptyset^{X_m \times Y_m \times 4}$

ep  $\in \{1, \dots, \text{episodes}\}$

sum\_rs\_sarsa  $\leftarrow \text{sarsa}(Q_{\text{sarsa}}, \text{is expected} \leftarrow 0, \alpha)$

sum\_rs\_expected\_sarsa  $\leftarrow \text{sarsa}(Q_{\text{expected-sarsa}}, 1, \alpha)$

sum\_rs\_Q\_learning  $\leftarrow \text{Q-learning}(Q_{\text{-Q-learning}}, \alpha)$

performance(asy-sarsa, i)  $\leftarrow$  sum\_rs\_sarsa

performance(asy-expected-sarsa, i)  $\leftarrow$  sum\_rs\_expected\_sarsa

performance(asy-Q-learning, i)  $\leftarrow$  sum\_rs\_Q-learning

if ep  $\leq 100$

performance(~~asy-interim sarsa~~, i)  $\leftarrow$  sum\_rs\_sarsa

...  
...

End & ep, i, run.

performance(1:3,:) / = episodes . runs

(4:6,:) / = 100 . runs

plot(performance, name='sum\_reward per episode')