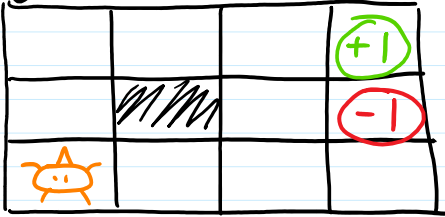


Temporal Difference (TD) Learning

GridWorld



at each step go $\uparrow, \downarrow, \leftarrow, \rightarrow$.8 go intended dir
 .2 go to side

$S =$ set of states
 state (row, col)
 next state

model $\left\{ \begin{aligned} P(s, s', a) &= \text{probability of action } a \text{ in state } s \text{ resulting in } s' \\ A &= \{\uparrow, \downarrow, \leftarrow, \rightarrow\} \\ R(s, s', a) &= \text{reward from } s \xrightarrow{a} s' \end{aligned} \right.$

policy: function $\pi: S \rightarrow A$

$V_{\pi}(s)$ = expected solve discounted reward given that in state s using π

π_{opt} : policy that maximizes $V_{\pi_{\text{opt}}}(s_0)$ initial state

$$V(s) = V_{\pi_{\text{opt}}}(s)$$

Value iteration (\equiv introduce turn limit)

$$\text{generally } V(s) = \max_a \sum_{s'} P(s, s', a) \cdot \left(\underbrace{R(s, s', a)}_{\text{immediate reward}} + \gamma \underbrace{V(s')}_{\text{discount future reward}} \right)$$

$$\pi_{\text{opt}}(s) = \underset{a}{\text{argmax}} Q(s, a)$$

$V(r, c, n)$ = value of state (r, c) w/n steps left

$$= \begin{cases} -1 & \text{if } r, c = (1, 3) \\ +1 & \text{if } r, c = (0, 3) \\ 0 & \text{if } n = 0 \\ \max_a \begin{aligned} &.8 \cdot V(r, c) + a, n-1 \\ &+.1 V(r, c) + \uparrow, n-1 \\ &+.1 V(r, c) + \downarrow, n-1 \end{aligned} & \end{cases}$$

$V(r, c, n)$ converges to $V(r, c)$ as $n \uparrow$

TD Value Learning

observe $s \xrightarrow{\pi(s)} s'$

giving ^{immediate} reward $R(s, s', a)$

get discounted future reward $\gamma \cdot V^\pi(s')$

$$\text{sample of } V^\pi(s) = R(s, s', a) + \gamma V^\pi(s')$$

$$\text{update estimate } V^\pi(s) \leftarrow V^\pi(s) + \underbrace{\alpha}_{\text{learning rate}} \underbrace{(R(s, s', a) + \gamma V^\pi(s') - V^\pi(s))}_{\text{error}}$$

$$= (1 - \alpha) V^\pi(s) + \alpha R(s, s', a) + \gamma V^\pi(s')$$

Q Learning

$Q(s, a)$ = estimate of expected discounted future reward

$$V(s) \approx \max_a Q(s, a)$$

initialize $Q(s, a) = \begin{cases} R(s) & \text{for terminal } s \\ 0 & \text{otherwise} \end{cases}$ *if you know this*

while not done

$s \leftarrow s_0$ *initial state*

while s not terminal

choose action a

ϵ -greedy is good
prob ϵ choose a randomly uniformly
 $1-\epsilon$ choose $\arg\max_a Q(s, a)$

observe transition (s, a, r, s')

immediate reward $(R(s, s', a))$

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

learning rate
(can \downarrow as episodes \uparrow)

surprise
(error)

episode

Q-Yahtzee

$Q(s, a) =$ estimate of expected value of action a in state s

unused categories *rerolls* $\approx 2^{12} \cdot 64 \cdot 3 \cdot 3 \cdot 2^8 \approx 2^{29} \approx 512 \text{ million}$
upper table *current roll*
Yahtzee bonus

6... 32 reroll *1... 13 choose cat* $avg \approx 2^4$
size of Q table

$\approx 2^{33}$ 8 billion
 will take a long time

Generalization :

Ones	1
Twos	4
Threes	
Fours	
Fives	20
Sixes	18
3K	23
4K	0
FH	
SS	30
LS	
C	
Y	50

Ones	2
Twos	4
Threes	
Fours	16
Fives	
Sixes	18
3K	17
4K	0
FH	
SS	0
LS	
C	
Y	50



after roll

$s_1 =$



after roll 1

$s_2 =$

$Q(s_1, a) \approx Q(s_2, a)$

Function Approximators

- learn a fn that approximates $Q(s, a)$

Linear Approximator

Define features of states and possibly actions

can ignore action for features of state only
on pace to earn upper bonus
is chance unused
both LS, SS unused
upper category a is unused

$$\begin{matrix} f_1(s, a) \\ f_2(s, a) \\ \vdots \end{matrix}$$

$$Q(s, a) = w_1 \cdot f_1(s, a) + \dots + w_n f_n(s, a)$$

In state s

Choose action a using exploit/explore policy

Observe transition (s, a, r, s')

Update

$$\begin{aligned} & \cancel{Q(s, a) \leftarrow Q(s, a) + \alpha (\max_{a'} Q(s', a') - Q(s, a))} \\ & w_i \leftarrow w_i + \alpha (r + \gamma \max_{a'} Q(s', a') - Q(s, a)) \cdot f_i(s, a) \end{aligned}$$

Given feature of states $f: S \rightarrow \mathbb{R}$

$$f_a: S \times A \rightarrow \mathbb{R} \quad f_a(s, a') = \begin{cases} f(s) & \text{if } a' = a \\ 0.0 & \text{otherwise} \end{cases}$$

Beware: α -learning may diverge
 decrease α as time \uparrow
 keep features on same scale as reward

Football: which position is better

$$(30, 4, 10, 6)$$

or

$$(50, 4, 10, 20)$$

1-10 from 30 w/ 30 sec remaining
 1-10 from 50 w/ 100 sec remaining

$$(35, 3, 20, 12)$$

or

$$(35, 2, 5, 12)$$

2-20 from 35 w/ 60 sec
 3-5 from 35 w/ 60 sec