

## 1 UCB iteration

(Note: I reverse the notations  $Q(s, a)$  and  $\hat{Q}(s, a)$  in Page 4 of “Dong K, Wang Y, Chen X, et al. Q-learning with ucb exploration is sample efficient for infinite-horizon mdp[J]. arXiv preprint arXiv:1901.09311, 2019.”)

for  $t=1,2,\dots$  do:  
 Take action  $a_t \leftarrow \arg \max_a Q(s_t, a)$ ;  
 Receive reward  $r_t$  and transit to  $s_{t+1}$ ;  
 $N(s_t, a_t) \leftarrow N(s_t, a_t) + 1$ ; (note:  $N(s_t, a_t)$  is the state\_action\_counter in codes);  
 $k \leftarrow N(s_t, a_t)$ ,  $b_k \leftarrow f(k)$ ; (note:  $f(k) = \sqrt{\frac{C_1 \log(k+1) + C_0}{k}}$  is the bonus,  $C_1, C_0$  are the tuning parameters bonus\_coef\_1, bonus\_coef\_0 in codes)  
 $\hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) + \alpha_k [(r_t + b_k + \max_a Q(s_t, a)) - \hat{Q}(s_t, a_t)]$ ; (note:  $\alpha_k = \frac{H+1}{H+k}$  is the learning\_rate in codes)  
 $Q(s_t, a_t) \leftarrow \min(Q(s_t, a_t), \hat{Q}(s_t, a_t))$ ; (note:  $Q(s_t, a_t), \hat{Q}(s_t, a_t)$  are Q\_table, Q\_hat\_table in codes)

### 1.1 Summary of tuning parameter in above iteration

$C_1, C_0$  in bonus;  $H$  in learning rate

## 2 index of state space for python codes

### 2.1 transition prob matrix

Note that

$$P_X(\Delta) = I + \Delta Q$$

where

$$Q_{\_i1\_i} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 1(\lambda_{2,1}) & -1 & 0 & 0 & 0 \\ 0 & 1 & \dots & \dots & 0 \\ 0 & 0 & 1 & -1 & 0 \\ 0 & 0 & 0 & 1 & -1 \end{bmatrix} \text{ and } Q_{\_i\_i1} = \begin{bmatrix} -1 & 1(\lambda_{1,2}) & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & \dots & \dots & 0 \\ 0 & 0 & 0 & -1 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$P_X$  is equal to

$$\begin{bmatrix} 1 - \lambda_{1,2}\Delta & \lambda_{1,2}\Delta & 0 & 0 & 0 \\ \lambda_{2,1}\Delta & 1 - (\lambda_{2,1} + \lambda_{2,3})\Delta & \lambda_{2,3}\Delta & 0 & 0 \\ 0 & \dots & \dots & \dots & 0 \\ 0 & 0 & \dots & \dots & \dots \\ 0 & 0 & 0 & \lambda_{|S_X|-1, |S_X|-2}\Delta & 1 - (\lambda_{|S_X|-1, |S_X|-2} + \lambda_{|S_X|-1, |S_X|})\Delta \\ 0 & 0 & 0 & 0 & \lambda_{|S_X|, |S_X|-1}\Delta \end{bmatrix} \begin{matrix} \lambda_{|S_X|} \\ 1 - \lambda_{|S_X|} \end{matrix}$$

$\{\lambda_{i, i+1}\}_{i=1,2,\dots, |S_X|-1}$  and  $\{\lambda_{i+1, i}\}_{i=2,3,\dots, |S_X|-1, |S_X|}$

## 2.2 index of state space

state variables:

x is an integer in  $\mathcal{S}_X/(\frac{\delta}{2}) := \{1, 2, \dots, (2N_P - 1)\} = \{1, 2, \dots, \dim\_X\} = \{0, 1, \dots, \dim\_X - 1\} + 1$ , and  $\dim\_X = |\mathcal{S}_X|$ .

so “python index”  $\text{idx\_x}$  for Q\_table equals  $x-1$

y is an integer in  $\mathcal{S}_Y := \{-N_Y, -(N_Y - 1), \dots, -2, -1, 0, 1, 2, \dots, (N_Y - 1), N_Y\} = \{0, 1, \dots, \dim\_Y - 1\} - N_Y$

so “python index”  $\text{idx\_y}$  for Q\_table equals  $y + N_Y$

Here, the set starting from 0 can always be used as python index set.

For quoted price:  $p\_a, p\_b$  are integers in  $\mathcal{S}_X/\delta := \{0, 1, 2, \dots, N_P\}$  no need to modify because it has already been index set starting from 0. However, since we add an action “doing nothing” in the state space, we make the action state be  $(p\_a, p\_b) \in \mathcal{S}_X^+/\delta \times \mathcal{S}_X^+/\delta$ , where

$$\mathcal{S}_X^+/\delta = \{0, 1, 2, \dots, N_P, N_P + 1\} \text{ and } N_P + 1 \text{ means “doing nothing”}.$$

So, when we specify the available action space based on the state variable  $(x, y)$ , we need to do:

$p\_a\_action\_list$  is in  $\mathcal{S}_X/\delta > x/2$ , for example:

if  $x = (2k)$  which means  $X_i = x = 2k\frac{\delta}{2} = k\delta$ , then  $\{\mathcal{S}_X/\delta > x/2\} = \{\mathcal{S}_X/\delta > k\delta\} = \{(k+1)\delta, (k+2)\delta, \dots, (2N_P - 1)\delta\}$ .

if  $x = (2k + 1)$  which means  $X_i = x = (2k + 1)\frac{\delta}{2} = (k + \frac{1}{2})\delta$ , then  $\{\mathcal{S}_X/\delta > x/2\} = \{\mathcal{S}_X/\delta > (k + \frac{1}{2})\delta\} = \{(k+1)\delta, (k+2)\delta, \dots, (2N_P - 1)\delta\}$ .

## 2.3 iteration of the true value function satisfying the Bellman equation

(1) the expectation of the reward at step  $i$  is  $p\_ask\_fill * (-x*tick/2 + p\_a*tick) + p\_buy\_fill * (x*tick/2 - p\_b*tick)$  is

$$E_{s' \sim P(s, a)}[r(s', s, a)]$$

(2) We have that,  $\text{np.dot}(\text{vec\_prob\_X.T}, \text{np.dot}(V\_star, \text{vec\_prob\_Y}))$  is

$$\sum_{i,j} p_i^X p_j^Y V_{i,j}^* = E_{s' \sim P(s, a)}[V^*(s')]$$

where  $\text{vec\_prob\_X} = (p_i^X)_i$ ,  $\text{vec\_prob\_Y} = (p_j^Y)_j$  captures  $s' \sim P(s, a)$  because the state variable  $s = (x, y)$ .