

# Note

Wednesday, August 12, 2020

8:48 PM

Define:

$r(s, a)[i]$  as the reward obtained by performing action  $a$  at state  $s$  for the  $i$ th time

$n(s, a)$  is the number of times performing action  $a$  at state  $s$

$R^\wedge(s, a) = \frac{1}{n(s, a)} \sum_{t=0} r(s, a)[i]$ , is the average reward obtained by performing action  $a$  at state  $s$

In many reinforcement learning problems, taking the same action  $a$  at state  $s$  can result in different reward  $r(s, a)$ .

Objective of a reinforcement learning problem is to maximize the discounted returns of average rewards

$$E_{s_0, a_0, \dots \sim \pi} \left( \sum_{t=0} \gamma^t R^\wedge(s_t, a_t) \right)$$

$R^\wedge(s_t, a_t)$  here has a confidence interval

$CI(R^\wedge(s_t, a_t)) := [R^\wedge(s_t, a_t) - \epsilon_{n(s, a)}^R, R^\wedge(s_t, a_t) + \epsilon_{n(s, a)}^R]$   
that we are  $1 - \delta_{n(s, a)}^R$  confident with.

By Hoeffding's Inequality:

$$\epsilon_{n(s, a)}^R = \sqrt{\frac{\ln \frac{2}{\delta_{n(s, a)}^R}}{n(s, a)}}$$

The count-based exploration bonus method proposed to maximize the discounted returns of upper tail of the confidence interval average rewards

$$E_{s_0, a_0, \dots \sim \pi} \left( \sum_{t=0} \gamma^t (R^\wedge(s_t, a_t) + \epsilon_{n(s, a)}^R) \right)$$

which transforms the estimate of  $Q^\sim(s, a)$  from

$$Q^\sim(s, a) = R^\wedge(s, a) + \gamma \sum_{s'} T(s'|s, a) \max_{a'} Q^\sim(s', a')$$

to be

$$Q^\sim(s, a) = (R^\wedge(s, a) + \epsilon_{n(s, a)}^R) + \gamma \sum_{s'} T(s'|s, a) \max_{a'} Q^\sim(s', a')$$

And we use  $\frac{\beta}{\sqrt{n(s, a)}}$  to represent  $\epsilon_{n(s, a)}^R := \sqrt{\frac{\ln \frac{2}{\delta_{n(s, a)}^R}}{n(s, a)}}$  in computation