598 RL    Fall 2020

Day 2

Bandit problem

k-armed bandit

choose an action $A$ &larr; takes values in $\{1, \ldots, k\}$

get a reward $R$ &larr; takes values in $\mathbb{R}$

want to maximize $E[R]$ $\leftarrow$

in other words, we want:

$$\underset{a \in \{1, \ldots, k\}}{\arg \max} \quad \underbrace{E[R | A = a]}_{Q^*(a)}$$

$$\int_{r \in \mathbb{R}} r \, p(r) \, dr$$

$-$ or $-$

$$\boxed{\sum_r r \, p(r)}$$

$$Q(a) = \frac{r_1 + r_2 + \ldots + r_n}{n}$$

① choose "$a$" uniformly at random from $\{1, \ldots, k\}$

② choose "$a$" as $\underset{a \in \{1, \ldots, k\}}{\arg \max} \quad Q(a)$

EXPLORE $\left\{\rule{0pt}{40pt}\right.$ ① choose "a" uniformly at random from $\{1, \ldots, k\}$

EXPLOIT $\left\{\rule{0pt}{40pt}\right.$ ② choose "a" as $\underset{a \in \{1, \ldots, k\}}{\arg\max} Q(a)$ $\left.\rule{0pt}{40pt}\right\}$ this is "greedy"

"$\epsilon$ - greedy" means do the non-greedy thing with probability $\epsilon$

Example:

if $\epsilon = 0.1$, then

90% of time we exploit

10% of time we explore

$$Q_n(a) = \frac{r_1 + \dots + r_n}{n}$$

$$= \frac{r_n}{n} + \frac{r_1 + \dots + r_{n-1}}{n}$$

$$= \frac{r_n}{n} + \frac{(n-1)}{(n-1)} \frac{r_1 + \dots + r_{n-1}}{n}$$

$$= \frac{r_n}{n} + \left(\frac{n-1}{n}\right) \left(\frac{r_1 + \dots + r_{n-1}}{n-1}\right)$$

$$\underbrace{\qquad\qquad}_{Q_{n-1}(a)}$$

$$= \frac{r_n}{n} + Q_{n-1}(a) - \left(\frac{1}{n}\right) Q_{n-1}(a)$$

$$= \underbrace{Q_{n-1}(a)}_{} + \underbrace{\frac{1}{n}}_{\substack{\text{step} \\ \text{size}}} \left( r_n - Q_{n-1}(a) \right)$$

$$\underset{\text{TARGET}}{\uparrow}$$