

NAME: Andrew Cside  
 CSCI S-89c Deep Reinforcement Learning  
 Part I of Assignment 1

Suppose we run an  $\varepsilon$ -greedy algorithm for the  $k$ -armed Bandit problem, where  $\varepsilon \in (0, 1)$ . Assuming  $q_*(a_1) \neq q_*(a_2)$  for all  $a_1 \neq a_2$ , where  $a_1, a_2 \in \{1, 2, \dots, k\}$ , please express

$$\lim_{t \rightarrow \infty} E[R_t]$$

in terms of  $\varepsilon$  and  $q_*(a)$ ,  $a \in \{1, 2, \dots, k\}$ .

SOLUTION:

$$E[R_t] = \sum \pi_t(a) \cdot a$$

where  $\pi_t = \begin{cases} \varepsilon/k & - \text{non greedy} \\ 1 - \varepsilon + (\varepsilon/k) & - \text{greedy} \end{cases}$

$$E[R_t] = \sum (\varepsilon/k + 1 - \varepsilon + \varepsilon/k) a$$

$$= \sum \left( \frac{2\varepsilon}{k} + 1 - \varepsilon \right) a$$

$$E[R_t | A_t = a] = \sum \left( \frac{2\varepsilon}{k} + 1 - \varepsilon \right) = q_*(a)$$