

8a) ϵ -greedy algorithm: policy = $\begin{cases} \text{explore w.p. } \epsilon \\ \text{exploit otherwise} \end{cases}$

$$p(T) = \sum_{t=0}^{T-1} (r_{\max}(t) - r(t))$$

$$P(r(t) \neq r_{\max}(t)) = \frac{N-1}{N} \text{ with } N = \text{num slots.}$$

average difference between $r(t)$ and $r_{\max}(t)$ is $p_{\max} - \frac{1}{N-1} \sum_{i=1}^{N-1} p_i$
and we explore ϵT times. so.

$$\begin{aligned} \lim_{T \rightarrow \infty} \mathbb{E}[p(t)] &= \lim_{T \rightarrow \infty} \left(\frac{N-1}{N} \cdot \left(p_{\max} - \frac{1}{N-1} \sum_{i=1}^{N-1} p_i \right) \cdot \epsilon \right) \\ &= \frac{N-1}{N} \cdot \epsilon \cdot \left(p_{\max} - \frac{1}{N-1} \sum_{i=1}^{N-1} p_i \right) \end{aligned}$$

\Rightarrow To decrease this, we can increase epsilon which will allow us to explore faster and find the optimal slot quicker.

8b) In i), we are approximating the Q function with θ . The assignment in ii) is meant to achieve convergence, such that comparing Q to Q_0 has a significance. If we did not do it after several iterations, Q would be very similar to Q_0 .

8c) We may want to do this to verify the gradient descent approach. If the gradient descent is well-implemented, it should converge to the same spot always. We may also find more local minimas.