

Assumptions:

- If the agent executes ϵ -greedy with constant ϵ , the environment does not need to be continuous
- If the agent executes a policy that converges to ϵ -greedy with constant ϵ in the limit $t \rightarrow \infty$ then the environment needs to be continuous around $p_i = \epsilon$ and $p_i = 1 - \epsilon$

If the agent behaves like ϵ -greedy, for at fixed $\epsilon > 0$, in the limit $t \rightarrow \infty$ then the set of explored action probabilities will eventually be dominated by values p_i such that $p_i \in \{\epsilon, 1 - \epsilon\}$.

For some q , we have that

$$\langle p \rangle = q(1 - \epsilon) + (1 - q)\epsilon = q - 2q\epsilon + \epsilon \quad (1)$$

$$\langle p^2 \rangle = q(1 - \epsilon)^2 + (1 - q)\epsilon^2 = q - 2q\epsilon + \epsilon^2 \quad (2)$$

$$= \langle p \rangle - \epsilon(1 - \epsilon) \quad (3)$$

This allows us to simplify the expected learned utility difference, Δ (see table) as a function of only $\langle p \rangle$ and the constants of the problem. Note that because $\langle p \rangle$ is dominated by $p_i \in \{\epsilon, 1 - \epsilon\}$, we must get $\epsilon \leq \langle p \rangle \leq (1 - \epsilon)$

Because of law of large numbers (?), Δ will eventually describe the agents belief arbitrary well. Let $\Delta = Q(a_1) - Q(a_2)$ and $p_i = P(a_1|t = i)$. If currently $\Delta > 0$, the agent will execute $p_i = 1 - \epsilon$ and drive $\langle p \rangle$ towards this value. If currently $\Delta < 0$, the agent will execute $p_i = \epsilon$ and drive $\langle p \rangle$ towards that value. Therefore the only possible stable points are

$$\langle p \rangle = \epsilon \quad \& \quad \Delta \leq 0 \quad (4)$$

$$\langle p \rangle = 1 - \epsilon \quad \& \quad \Delta \geq 0 \quad (5)$$

$$\Delta = 0 \quad \& \quad \frac{d\Delta}{d\langle p \rangle} \leq 0 \quad (6)$$

(4) and (5) represent constant p_i at either extreme, ϵ or $1 - \epsilon$. (6) represent fluctuating p_i but with a stable $\langle p \rangle$ in between the two extremes.

Consider an environment with consists of the agent playing a game against it self with an arbitrary payoff matrix. The copy will always use the same action probabilities as the agent, for any given round.

		copy	
		a_1	a_2
self	a_1	M_{11}	M_{12}
	a_2	M_{21}	M_{22}

This generalisation covers, for example, *Prisoners' dilemma against copy* and *Death in Damascus* but not *Absent minded driver* and *Evidential blackmail*.