Convergence of Q-learning Algorithms.

# I. Stochastic Approximation.

Consider algorithm $x_{n+1} = x_n + a_n(h(x_n) + \xi_{n+1})$

where $\mathbb{E}[\xi_{n+1} | \mathcal{F}_n] = 0$ a.s. $\forall n$

$\mathcal{F}_n = \sigma(x_i, i \le n)$.

Assume $h$ is $L$-libschitz.

$$\sum_n a_n = \infty, \quad \sum_n a_n^2 < \infty$$

$$\sup_n \|x_n\| < \infty \quad a.s.$$

$$\mathbb{E}[\|\xi_{n+1}\|^2 | \mathcal{F}_n] \le K(1 + \|x_n\|^2) \quad a.s. \; \forall n$$

Then the algorithm trajectories look like those of

$$\dot{x} = h(x) \qquad (1)$$

If (1) has a unique globally asymptotic stable point $x^*$ then $\lim_{n \to \infty} x_n = x^*$.

# II. Q-learning.

Let $F: \mathbb{R}^{S \times A} \longrightarrow \mathbb{R}^{S \times A}$

$$F(q)_{sa} = r(s,a) + \lambda \mathop{\mathbb{E}}_{s \sim p(\cdot | s,a)} \left[ \max_b q(s,b) \right]$$

$F$ is a contraction $(1\text{-lipschitz})$

Q- learning algorithm :

For any pair $(s,a)$, if we look at
the $n$-th times, state $s$ and action $a$ are
observed :

$$q_{n+1}(s,a) = q_n(s,a) + \alpha_n \Big[ r(s,a)$$
$$+ \partial \max_b q_n(S_{n+1}(s,a), b)$$
$$- q_n(s,a) \Big]$$

⚠ $q_n(s,a)$ is the Q-value after the pair $(s,a)$
has been observed $n$-times.

If $\sum_n \alpha_n = \infty$, $\sum_n \alpha_n^2 < \infty$, we have a stochastic
approximation algorithm, mimicking the dynamical
system :
$$\dot{q} = F(q) - q$$

This system is globally stable because $F$ is
contractive (Exercice), and so
Q- learning converges to the unique fixed
point of $F$

#