

# Class Tutorial 10

---

## 1. Stochastic Approximation ODE Approach

Consider an initially empty urn to which balls, either red or black, are added one at a time.

Let  $y_n$  denote the *number* of red balls at time  $n$ , and  $\theta_n \triangleq y_n / n$  the *fraction* of red balls at time  $n$ .

The probability of adding a red ball at time  $n+1$  is a function of  $\theta_n$  alone, and we denote it by  $p(\theta_n)$ . We are interested in the limit  $\lim_{n \rightarrow \infty} \theta_n$ . Let  $\xi_{n+1}$  denote the following random variable

$$\xi_{n+1} = \begin{cases} 1 & , \quad (n+1)\text{st ball is red} \\ 0 & , \quad (n+1)\text{st ball is black} \end{cases}$$

a. Write an update equation for  $\theta_n$  in the form of a stochastic approximation

$$\theta_{n+1} = \theta_n + a_n [h(\theta_n) + \omega_n], \text{ where } \omega_n \text{ is a martingale difference noise.}$$

b. Write down the corresponding ODE. What are the asymptotically stable equilibria?

c. Assume  $\theta^*$  is a unique stable equilibrium of the ODE, and that  $p(\theta)$  is Lipschitz. Show that  $\lim_{n \rightarrow \infty} \theta_n = \theta^*$ .

## Solution

a. We have  $y_{n+1} = y_n + \xi_{n+1}$ , therefore

$$\begin{aligned} \theta_{n+1} &= y_{n+1} / (n+1) = \\ &= \frac{y_n + \xi_{n+1}}{n+1} = \frac{ny_n}{n(n+1)} + \frac{\xi_{n+1}}{n+1} = \\ &= \frac{n}{n+1} \theta_n + \frac{\xi_{n+1}}{n+1} = \theta_n + \frac{1}{n+1} (\xi_{n+1} - \theta_n) \end{aligned}$$

We now transform it to a stochastic approximation form:

$$\begin{aligned} \theta_{n+1} &= \theta_n + \frac{1}{n+1} (\xi_{n+1} + \mathbb{E}(\xi_{n+1} | \theta_n) - \mathbb{E}(\xi_{n+1} | \theta_n) - \theta_n) \\ &= \theta_n + \frac{1}{n+1} (p(\theta_n) - \theta_n + \xi_{n+1} - p(\theta_n)) \end{aligned}$$

Here,  $\omega_n \equiv \xi_{n+1} - p(\theta_n)$ , and obviously  $\mathbb{E}(\omega_n | \theta_n) = 0$ . Also,  $\alpha_n = \frac{1}{n+1}$  and

$$h(\theta_n) = p(\theta_n) - \theta_n.$$

b.  $\dot{\theta}(t) = h(\theta(t)) = p(\theta(t)) - \theta(t)$ . The equilibria satisfy  $p(\theta^*) = \theta^*$ . Stability criterion: by linearization around  $\theta^*$ :

$$\begin{aligned} p'(\theta^*) - 1 &< 0 \\ p'(\theta^*) &< 1 \end{aligned}$$

c. We will use the stochastic approximation convergence theorem (Theorem 1 in the lecture notes). We need to show that:

Assumption G1:  $\sum_n a_n = \infty$ ,  $\sum_n a_n^2 < \infty$ , clearly holds for  $\alpha_n = \frac{1}{n+1}$ .

Assumption N1:  $\omega_n$  is a martingale difference, and  $\mathbb{E}[\|\omega_n\|^2 | \mathcal{F}_n] \leq A + B\|\theta_n\|^2$ . Clearly holds since the noise is bounded.

The sequence  $\theta_n$  is bounded w.p. 1. Here this holds by definition of  $\theta_n$ .

## 2. Stochastic Approximation Contraction Mapping Approach

Consider the Q-learning algorithm where  $\pi$  is some 'explorative' policy.

- Show that the updating equations of the Asynchronous Q-learning algorithm can be written as a sum of two terms, a  $TQ_n - Q_n$  and a martingale noise.
- Assume that the step-size requirements satisfy the assumption

$$\forall s, a \text{ (w.p. 1)} \sum \alpha_n(s, a) = \infty, \sum \alpha_n^2(s, a) < \infty,$$

Prove the convergence of Q-learning to the optimal Q-function.

Q-learning
<p><b>Initialize:</b> <math>Q_0(s, a)</math>  <b>For</b> <math>k=0,1,\dots</math> <b>do</b>            Sample <math>a_n \sim \pi(s_n)</math>            Sample <math>s_{n+1} \sim P(\cdot   s_n, a_n), r(s_n, a_n, s_{n+1})</math>            <math>\delta_n = r(s_n, a_n, s_{n+1}) + \gamma \max_{a'} Q_n(s_{n+1}, a') - Q_n(s_n, a_n)</math>            <math>Q_{n+1}(s_n, a_n) = Q_n(s_n, a_n) + \alpha_n(s_n, a_n) \delta_n</math></p>

## Solution

- a) We consider updating the full vector of the Q-function. For all the  $s, a$  entries which are not the current state-action pair we set

$$\alpha_n(s_n, a_n) = 0.$$

Furthermore, we have that

$$\delta_n = (TQ_n)(s_n, a_n) - Q_n(s_n, a_n) + \omega_n$$

Where

$$(TQ_n)(s_n, a_n) = E_{s' \sim P(\cdot | s_n, a_n)} \left[ r(s_n, a_n, s') + \gamma \max_{a'} Q_n(s', a') \right]$$

$$\omega_n = r(s_n, a_n, s_{n+1}) + \gamma \max_{a'} Q_n(s_{n+1}, a') - (TQ_n)(s_n, a_n)$$

It is easy to see that  $\omega_n$  is a martingale noise by verifying that

$$E[\omega_n | \mathcal{F}_n] = 0$$

Where  $\mathcal{F}_n$  is the entire history, meaning the value of  $Q_n$ , the current state and action,  $s_n, a_n$ , and  $\alpha_n(s_n, a_n)$ .

- b) We use Theorem 9.6 from the lectures to show the convergence. We start by observing that  $T$  is a contraction operator. Furthermore,
- The step-size requirements holds by assumption.
  - As we saw,  $\omega_n$  is a martingale noise. It is also satisfied that

$$E[\omega_n^2 | \mathcal{F}_n] \leq |Q|_\infty^2$$

Thus, using Theorem 9.6 we get that Q-learning converges w.p 1 to the optimal Q function.