

Mariia Aleksandrovych, Yaffa Atkins, Atoosa Lotfi, Olga Kliuchnik

Exercises 3.2, 4.2, 4.3, 4.4 and 4.7

EXERCISE 3.2. This is a simple version of a *multi-armed bandit problem*. A slot machine has two arms and the winning probability of the two arms, denoted by p_A and p_B , respectively, are not known. It is assumed that the event of winning on each arm is independent of past history and also independent of the other arm's results. Let θ_n denote the belief that we have that arm A has higher winning probability. Consequently, we choose arm A at step n with probability θ_n (called "exploitation" in machine learning). The goal is to learn the correct value $\theta^* = 1$ if $p_A > p_B$, and $\theta^* = 0$ otherwise (assuming that $p_A \neq p_B$).

A straightforward learning approach is

$$\theta_{n+1} = \begin{cases} \theta_n & \text{if } n\text{-th game was a loss,} \\ \theta_n + \epsilon_n (1 - \theta_n) & \text{if win on arm A,} \\ \theta_n - \epsilon_n \theta_n & \text{if win on arm B.} \end{cases}$$

- (a) Write the above recursion as stochastic approximation, specifying the feedback Y_n .
- (b) Provide the target vector-field, assuming no bias term is involved.
- (c) Argue that the target vector field is coercive for the learning problem.
- (d) *Optional:* Program the procedure to test various schemes for the learning rates $\{\epsilon_n\}$, including the case of constant ϵ . Discuss your results.

Felisa Vázquez-Abad and Bernd Heidergott.

107

- a. Stochastic approximation of the problem is $\theta_{n+1} = \theta_n + \epsilon_n Y_n$ where $Y_n = \begin{cases} 0 \\ 1 - \theta_n \\ -\theta_n \end{cases}$. $Y_n = 0$ when

n th game was a loss. $Y_n = 1 - \theta_n$ if there is a win on arm A. This causes θ_n to increase which makes sense since A winning encourages our belief that $p_A > p_B$. $Y_n = -\theta_n$ if there is a win on arm B. This causes θ_n to decrease which makes sense, since B winning encourages our belief that p_A is not greater than p_B .

- b. By definition, a vector field is the expected value of feedback function on sigma algebra

$$\sigma(\mathfrak{F}_{n-1}) \quad G(\theta_n) = E[Y_n | \mathfrak{F}_{n-1}] = \begin{cases} E(0) \\ E(1 - \theta_n) \\ E(-\theta_n) \end{cases} \quad G(\theta_n) = \begin{cases} 0 \\ 1 - \theta_n \\ -\theta_n \end{cases}$$

Since there is no bias, and constant condition on the feedback, the expected value is simply the value.

- c. θ^* is bound between 0 and 1 so it cannot blow up, the vector field is locally Lipschitz continuous because cannot vary more than 0 to 1, no matter how many times A wins in a row, or how many times B wins in a row. θ moves a little bit with each iteration. The possible solutions are $\theta^* = 1$ and $\theta^* = 0$. You probably will never reach a solution but whichever solution you are moving towards is an asymptotically stable point of the ODE because if you see a lot more wins for A, P_A is probably bigger and you will keep moving towards 1, and vice versa. Whichever point you are moving towards is a KKT of the NLP because it is stationary since it is once you reach 1 or zero you stop, and they satisfy the equality and inequality constraint of being between 0 and 1.

EXERCISE 4.2. Let $\{X_k\}$ be an iid sequence with mean μ and define

$$S_n = \sum_{k=1}^n (X_k - \mu), n \geq 1,$$

Show that S_n is a Martingale with respect to the natural filtration of the process, where $\mathcal{F}_n = \sigma(X_1, \dots, X_n)$.

4.2

$\{X_k\}$ iid, $S_n = \sum_{k=1}^n (X_k - \mu), n \geq 1$, show S_n is martingale.

S_n is martingale if $E[S_{n+1} | X_1, \dots, X_n] = S_n$.

$$E[S_{n+1} | X_1, \dots, X_n] = E[S_n + (X_{n+1} - \mu) | X_1, \dots, X_n]$$

$$S_n \text{ is constant on } X_1, \dots, X_n = S_n + E[X_{n+1} - \mu | X_1, \dots, X_n]$$

$$X_{n+1} \text{ is independent of } X_1, \dots, X_n = S_n + E[X_{n+1} - \mu]$$

$$= S_n \Rightarrow \text{it's Martingale}$$

EXERCISE 4.3. Let $\delta M_i = Y_i - \mathbb{E}[Y_i | \mathcal{F}_{i-1}]$ be defined as in the martingale difference noise model. Show that the process $M_n \stackrel{\text{def}}{=} \sum_{i=0}^n \epsilon_i \delta M_i$ is a martingale process on $(\Omega, \mathbb{P}, \{\mathcal{F}_n\})$. Show that $\mathbb{E}[\delta M_n \delta M_m] = 0$. for the basic definition and properties of martingale processes we refer to the Appendix.

Solution:

$M_n \triangleq \sum_{i=0}^n \epsilon_i \delta M_i$ is a martingale process ^{on filtered probability space of $(\Omega, \mathcal{F}, \mathbb{P})$} when

$$E[\delta M_n \delta M_m] = 0 \quad \text{for all } m, n \geq 0 \text{ with } m \neq n$$

Let's assume that $0 \leq m < n$ then:

$$\begin{aligned} \sum \delta M_n \delta M_m &= \sum (M_n - M_{n-1})(M_m - M_{m-1}) \\ &= \sum M_n M_m - \sum M_n M_{m-1} - \sum M_{n-1} M_m + \sum M_{n-1} M_{m-1} \\ &= \sum E(M_n M_m | \mathcal{F}_m) - \sum E(M_n M_{m-1} | \mathcal{F}_{m-1}) \\ &\quad - \sum E(M_{n-1} M_m | \mathcal{F}_m) + \sum E(M_{n-1} M_{m-1} | \mathcal{F}_{m-1}) \\ &= \sum M_m E(M_n | \mathcal{F}_m) - \sum M_{m-1} E(M_n | \mathcal{F}_{m-1}) \\ &\quad - \sum M_m E(M_{n-1} | \mathcal{F}_m) + \sum M_{m-1} E(M_{n-1} | \mathcal{F}_{m-1}) \\ &= \sum M_m \cdot M_m - \sum M_{m-1} M_{m-1} - \sum M_m M_m + \sum M_{m-1} M_{m-1} \\ &= \sum M_m^2 - \sum M_{m-1}^2 - \sum M_m^2 + \sum M_{m-1}^2 \\ &= 0 \end{aligned}$$

EXERCISE 4.4. Refer to the model in Example 3.4. Show that for a random variable X with finite variance,

$$\nabla J(\theta) = \begin{pmatrix} -\mathbb{E}[Z(X) - \theta_1 - \theta_2 X] \\ -\mathbb{E}[X Z(X) - \theta_1 X - \theta_2 X^2] \end{pmatrix}. \quad (4.19)$$

For each experimental point x_n we obtain a random observation $\xi_n = Z(x_n)$ with $\mathbb{E}(Z(x_n)) = h(x_n)$. The feedback function is $Y_n = (\xi_n - \theta_n(1) - \theta_n(2)x_n)(1, x_n)^T$. Use (4.19) to show the claim that $\mathbb{E}[Y_n | \mathcal{F}_{n-1}] = -\nabla J(\theta_n)$.

The model from the Exercise 3.4 says that we have cost function as

$$J(\theta) = \frac{1}{2} \mathbb{E} \left[(Z(X) - (\theta_1 + \theta_2 X))^2 \right]$$

Then we take partial derivatives over θ_1 and θ_2 to get the gradient

Then

$$\begin{aligned}\frac{\partial}{\partial \theta_1} J(\theta) &= \frac{\partial}{\partial \theta_1} \left\{ \frac{1}{2} E \left[\left(Z(X) - (\theta_1 + \theta_2 X)^2 \right) \right] \right\} = \frac{\partial}{\partial \theta_1} \left\{ \frac{1}{2} \int_{\Omega} \left[\left(Z(X) - (\theta_1 + \theta_2 X)^2 \right) \right] P(dX) \right\} \\ \frac{\partial}{\partial \theta_2} J(\theta) &= \frac{\partial}{\partial \theta_2} \left\{ \frac{1}{2} E \left[\left(Z(X) - (\theta_1 + \theta_2 X)^2 \right) \right] \right\} = \frac{\partial}{\partial \theta_2} \left\{ \frac{1}{2} \int_{\Omega} \left[\left(Z(X) - (\theta_1 + \theta_2 X)^2 \right) \right] P(dX) \right\}\end{aligned}$$

The $J(\theta)$, $\frac{\partial}{\partial \theta_2} J(\theta)$, $\frac{\partial}{\partial \theta_1} J(\theta)$ are integrable on (Ω, \mathfrak{F}_n) Then using Leibniz internal rule, we can

interchange derivative and integration operations. We have:

$$\begin{aligned}\frac{1}{2} \int_{\Omega} \left[\frac{\partial}{\partial \theta_1} \left(Z(X) - (\theta_1 + \theta_2 X)^2 \right) \right] P(dX) &= \frac{1}{2} \int_{\Omega} \left[-2(Z(X) - (\theta_1 + \theta_2 X)) \right] P(dX) = -E \left[(Z(X) - (\theta_1 + \theta_2 X)) \right] \\ \frac{1}{2} \int_{\Omega} \left[\frac{\partial}{\partial \theta_2} \left(Z(X) - (\theta_1 + \theta_2 X)^2 \right) \right] P(dX) &= \frac{1}{2} \int_{\Omega} \left[-2X(Z(X) - (\theta_1 + \theta_2 X)) \right] P(dX) = -E \left[(Z(X)X - (\theta_1 X + \theta_2 X^2)) \right]\end{aligned}$$

Where:

$$\nabla J(\theta) = \begin{pmatrix} -E \left[(Z(X) - \theta_1 - \theta_2 X) \right] \\ -E \left[(Z(X)X - \theta_1 X - \theta_2 X^2) \right] \end{pmatrix}$$

Now let take expectation value of the feedback function $Y_n = \begin{pmatrix} \xi_n - \theta_{n,1} - \theta_{n,2}x_n \\ \xi_n x_n - \theta_{n,1}x_n - \theta_{n,2}x_n^2 \end{pmatrix}$

$$E \left[Y_n | \theta_n \right] = \begin{pmatrix} E \left[\xi_n - \theta_{n,1} - \theta_{n,2}x_n \right] \\ E \left[\xi_n x_n - \theta_{n,1}x_n - \theta_{n,2}x_n^2 \right] \end{pmatrix} = \begin{pmatrix} E \left[Z(x_n) - \theta_{n,1} - \theta_{n,2}x_n \right] \\ E \left[Z(x_n)x_n - \theta_{n,1}x_n - \theta_{n,2}x_n^2 \right] \end{pmatrix} = -\nabla_{\theta} J(\theta_n)$$

EXERCISE 4.7. Consider the following well known Nash equilibrium problem in transportation. The transportation time along an arc i is denoted by t_i and it depends on the vector of normalized traffic flow x , where $x_{(k,l)}$ is the total traffic on the arc (k, l) . In Figure 4.1 there are three possible routes from origin to destination, $\{(OAD), (OABD), (OBD)\}$ and the corresponding travel times per arc are shown. The fraction of traffic on route r is denoted by θ_i , so that $\sum_{i=1}^3 \theta_i = 1.0$. Therefore the fraction of traffic on the arc (O, A) is given by $x_{(O,A)} = \theta_1 + \theta_2$ and similarly for other arcs. The total time on route i is the sum of the travel times along the arcs that form the route and is denoted by $T_i, i = 1, 2, 3$.

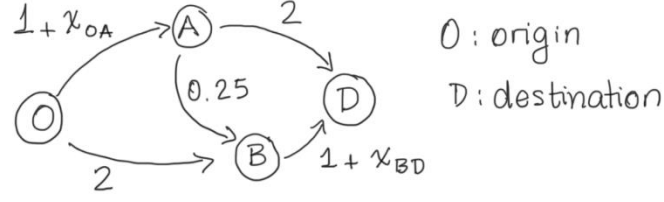


Figure 4.1: Traffic flow equilibrium

A Nash equilibrium will be an allocation such that if a driver on route i decides to change his/her path, then he/she will experience a larger delay than staying at equilibrium. The vector of travel times per route is given by:

$$T(\theta) = \begin{pmatrix} 3 + \theta_1 + \theta_2 \\ 2.25 + \theta_1 + 2\theta_2 + \theta_3 \\ 3 + \theta_2 + \theta_3 \end{pmatrix}$$

- Show that there is a unique Nash equilibrium, by showing that there is a unique value θ^* such that $T_i(\theta^*) = \text{constant}$, is independent of i .
- In order to attain equilibrium, the flow on path i should aim to "equalise" the delay. Suppose now that we do not know the various constants in the delay function, but can only just estimate the delays by running a simulation, which yields an unbiased estimator $\widehat{T}(\theta)$.

$$\theta_{n+1,i} = \theta_{n,i} - \epsilon_n \left(\widehat{T}_i(\theta_n) - \frac{1}{3} \sum_k \widehat{T}_k(\theta_n) \right). \quad (4.20)$$

Show that the algorithm is mass preserving, that is, $\sum_i \theta_{0,i} = 1$ then $\sum_i \theta_{n,i} = 1$.

- Characterise the behaviour of the stochastic approximation (4.20) and specify your assumptions. In particular, show that the target vector field is coercive for the equilibrium problem.

By the way, this example is classic to show that the Nash equilibrium does not minimise overall travel time. For the specific model, $\theta^* = (0.25, 0.50, 0.25)^T$, and $T_i(\theta^*) = 3.75$ for all i . However, for $\theta = (0.50, 0.0, 0.50)$, then $T_1(\theta) = T_3(\theta) = 3.5$.

4.7

t_i = transportation time along arc i

fraction of traffic on route = θ_i where $\sum_{i=1}^3 \theta_i = 1$

a) vector of travel times is given in the problem statement

$$T(\theta) = \begin{pmatrix} 3 + \theta_1 + \theta_2 \\ 2.25 + \theta_1 + 2\theta_2 + \theta_3 \\ 3 + \theta_2 + \theta_3 \end{pmatrix}$$

regardless of i , we want constant θ^*

show that $T(\theta_1) = T(\theta_2) = T(\theta_3)$:

$$\bullet T(\theta_1) = T(\theta_3)$$

$$= 3 + \theta_1 + \theta_2 = 3 + \theta_2 + \theta_3$$

$$= \theta_1 = \theta_3$$

$$\bullet T(\theta_1) = T(\theta_2)$$

$$= 3 + \theta_1 + \theta_2 = 2.25 + \theta_1 + 2\theta_2 + \theta_3$$

$$= .75 = \theta_2 + \theta_3$$

$$\theta_1 + \theta_2 + \theta_3 = 1$$

$$\theta_2 + \theta_3 = .75$$

$$\text{So } \theta_1 = .25$$

$$\theta_1 = \theta_3 = .25$$

$$\text{So } \theta_2 = .5$$

b) we want to equalise the delay, we only have an estimator.

$$\theta_{n+1,i} = \theta_{n,i} - \epsilon_n \left(\underbrace{\widehat{T_i(\theta_n)}}_{\substack{\text{estimate of current} \\ \text{route time}}} - \underbrace{\frac{1}{3} \sum_k \widehat{T_k(\theta_n)}}_{\substack{\text{average of estimate} \\ \text{of total route time}}} \right)$$

Show that this algorithm is mass preserving,

meaning at any iteration, the routes still sum to 1 :

Although individual $T_i(\theta)$ might change, $\sum_i \theta_i$ is always 1. so sum of average $\theta_i = 1$

$1 - 1 = 0$ so the average change is always 0.

$$\text{So } \sum_{i=1}^3 \left(\widehat{T_i(\theta_n)} - \frac{1}{3} \sum_k \widehat{T_k(\theta_n)} \right) = 0$$

$$\text{and } \sum_{i=1}^3 \widehat{T_i(\theta_n)} - 3 * \frac{1}{3} \sum_k \widehat{T_k(\theta_n)} = 0$$

$\epsilon_n * 0 = 0$ so θ_n isn't changing

c) Behaviour of the above stochastic optimization:

If $T_i(\theta) > \text{average } T_i(\theta)$

then $\theta_{n+1} < \theta_n$: moving towards less traffic.

and vice versa, equalizing the delays:

Our assumptions : $\epsilon \rightarrow 0$ and $\sum \epsilon_i = \infty$

This is because as we get closer to equalized delay

we need more accuracy, but we may never actually reach it. we'll just keep getting closer with infinite iterations.

we also assume the solution is a stationary point

because once we reach zero, θ stays the same.

Coercive: $\sum_{i=1}^3 \theta_i = 1$ so θ is bound,

for each θ : $0 < \theta < 1$

Since θ is bound, the target vector field doesn't blow up in finite time, rather it keeps getting closer to the right answer.

The target vector field is locally Lipschitz since θ is bound and cannot blow up, it moves in very small steps which tend to zero

we reach an asymptotically point when all delays are equalized. this is the solution.

$$T_i(\theta_i) - \frac{1}{3} \sum_k T_k(\theta_n) = 0$$

$$\text{so } \theta_{n+1} = \theta_n = \text{stationary}$$

solution vector θ is within the bounds of $0 < \theta < 1$

so it's a KKT point.