

Text :- ↳ V.S. Borkar, Stochastic Approximation: A dynamical systems viewpoint.
(2008)

2nd Edition (2022)

↳ Bruce Hajek, Random processes for Engineers
Cambridge Univ Press, 2015

↳ V.S. Borkar, Probability theory: Advanced Course, Springer 1995.

Evaluation :- 2 Mid-terms (20% + 20%)

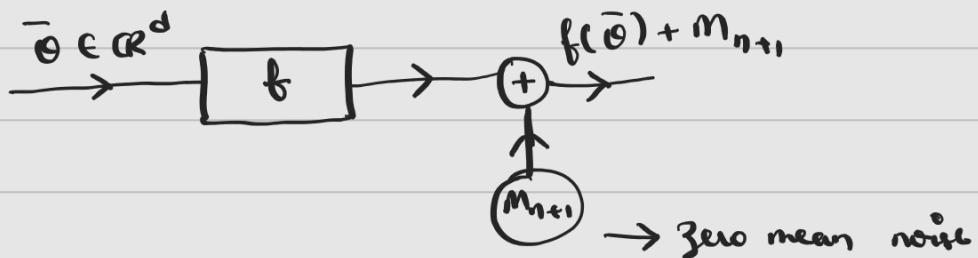
1 Project (20%)

1 Final exam (40%)

Paper :- K.L. Robbins and H. Marro, ,

Annals of Math. Stat., 1951

Suppose $f: \mathbb{R}^d \rightarrow \mathbb{R}^d$ and we don't know the functional form of f .



Q1 Can we find a θ^* for which $f(\theta^*) = 0$

Algorithm :- (Robbins - Monro Algorithm)

$$\theta_{n+1} = \theta_n + \alpha(n)(f(\theta_n) + M_{n+1})$$

Starting from some $\theta_0 \in \mathbb{R}^d$

$\alpha(n)$: learning rate parameter or stepsize.

Typically $\alpha(n) = 1/n$, $n \geq 1$

As $n \rightarrow \infty$, under some conditions, $\theta_n \rightarrow \theta^*$ s.t. $f(\theta^*) = 0$.

Examples where R-M algorithm is useful

↪ Finding fixed points of a function given noisy observations.

Finds θ^* s.t. $\theta^* = f(\theta^*)$

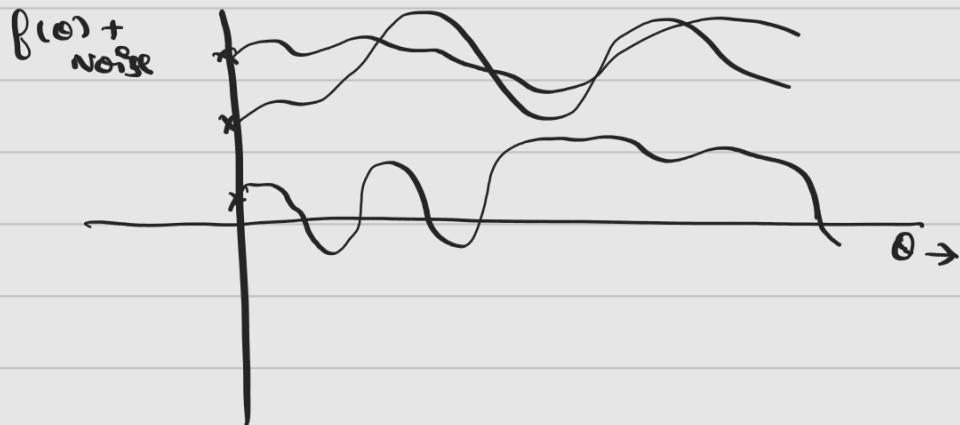
$$\text{let } g(\theta) = f(\theta) - \theta$$

(zeros of g are fixed points of f)

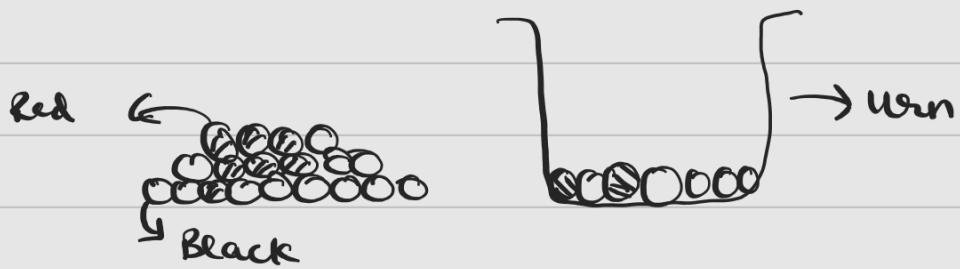
↪ Finding local minima of a function $g: \mathbb{R}^d \rightarrow \mathbb{R}$

$$\text{let } f = \nabla g$$

The zeros of f will be stationary points of g .



Chapter - 1



At any instant, we randomly pick a ball and drop it in urn.
By time n , we drop n balls in the urn.

Objective :- Find the long-run fraction of red balls in the urn.

Let $y_n = \# \text{ of red balls in urn by time } n$

$x_n = \text{fraction of red balls in urn at time } n$.

$$\therefore x_n = \frac{y_n}{n}$$

Let $\xi_{n+1} = I\{\text{1st ball is red in colour}\} = \begin{cases} 1 & \text{if } n+1^{\text{th}} \text{ ball is red} \\ 0 & \text{o.w.} \end{cases}$

$$\text{Then } y_{n+1} = y_n + \xi_{n+1}$$

Assumption :- Probability that $(n+1)^{\text{st}}$ ball dropped is red in colour given the entire history of red and black balls dropped until time n depends only on $p(x_n)$, where

$$p: [0, 1] \rightarrow [0, 1]$$

$$P(\xi_{n+1} = 1 | \xi_1, \xi_2, \dots, \xi_n) = p(x_n) \text{ a.s.}$$

$$\text{Note :- } x_n = \frac{y_n}{n}$$

$$\begin{aligned}
 \hat{x}_{n+1} &= \frac{y_{n+1}}{(n+1)} = \frac{(y_n + \hat{\xi}_{n+1})}{(n+1)} \\
 &= \frac{n}{n+1} \cdot \frac{y_n}{n} + \frac{1}{n+1} \hat{\xi}_{n+1} \\
 &= \frac{n}{n+1} x_n + \frac{\hat{\xi}_{n+1}}{n+1} \\
 x_{n+1} &= x_n + \frac{1}{n+1} (\hat{\xi}_{n+1} - x_n) \\
 &= x_n + \frac{1}{n+1} (p(x_n) - x_n) + \frac{1}{n+1} (\hat{\xi}_{n+1} - p(x_n))
 \end{aligned}$$

$$\text{Let } M_{n+1} = \hat{\xi}_{n+1} - p(x_n), n \geq 0$$

$$E[M_{n+1} | \hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_n]$$

$$\begin{aligned}
 &= E[\hat{\xi}_{n+1} | \hat{\xi}_1, \hat{\xi}_2, \dots, \hat{\xi}_n] - E[p(x_n) | \hat{\xi}_1, \dots, \hat{\xi}_n] \\
 &= I(\hat{\xi}_{n+1} = 1 | \hat{\xi}_1, \dots, \hat{\xi}_n) - E[p(x_n) | \hat{\xi}_1, \dots, \hat{\xi}_n] \\
 &= p(x_n) - p(x_n) \\
 &= 0
 \end{aligned}$$

$$E[M_{n+1}] = E[E[M_{n+1} | \hat{\xi}_1, \dots, \hat{\xi}_n]] = 0$$

Algorithm :- $x_{n+1} = x_n + \frac{1}{n+1} (p(x_n) - x_n) + \frac{1}{n+1} M_{n+1}, n \geq 0$

The ODE approach

$\left[\begin{array}{l} 1977 - \text{d'yury} \\ 1977-78 - \text{Drozdovsky} \\ 1978 - \text{Kushner and Clark} \end{array} \right]$

Under some conditions we can approximate the algorithm via an ODE

$$\dot{x}(t) = p(x(t)) - x(t) \quad \xrightarrow{(*)} \text{Numerical scheme for solving ODE}$$

Euler discretization of the ODE

From (*), we have

$$x(t) \doteq x_0 + \int_0^t p(x(\tau)) - x(\tau) d\tau$$

Suppose $h > 0$ is a small time element.

$$\text{Then } \int_0^t p(x(\tau)) - x(\tau) d\tau = \sum_{m=0}^n (p(x(mh)) - x(mh)) h$$

$t \approx nh$

$$x(nh) \doteq x_0 + \sum_{m=0}^{n-1} p(x(mh)) - x(mh) h$$

$$x((n+1)h) = x_0 + \sum_{m=0}^n p(x(mh)) - x(mh) h$$

Then $x((n+1)h) = x(nh) + h(p(x(nh)) - x(nh))$ (Euler discretization)

□

↪ (2)

(Comparing (1) and (2))

$$\hookrightarrow h = 1/(n+1)$$

↪ In (1), additional term = $\frac{M_{n+1}}{n+1}$

↪

$$x_{n+1} = x_0 + \sum_{m=0}^n \underbrace{\frac{1}{m+1} (p(x_m) - x_m)}_{n \downarrow} + \sum_{m=0}^n \underbrace{\frac{1}{m+1} M_{m+1}}_{\text{As } n \rightarrow \infty \text{ } \sum_{m=0}^n \frac{M_{m+1}}{m+1} < \infty \text{ w.p. 1}}$$

Let's consider the ODE

$$\dot{x}(t) = p(x(t)) - x(t)$$

Assume $p: [0, t] \rightarrow [0, 1]$ is Lipschitz continuous

$$\exists L > 0, \forall x, y \in [0, 1] \quad |p(x) - p(y)| \leq L |x - y|$$

We can show that the ODE is well-posed.

Suppose, we start ODE (*) s.t.

$$x_0 \in [0, 1]$$

Note :- When $\pi(t) = 0$, $\dot{\pi}(t) \geq 0$

When $\pi(t) = 1$, $\dot{\pi}(t) \leq 0$

Suppose $\dot{\pi}(t) > 0$ when $\pi(t) = 0$

⇒ The ODE will move in the right direction.

However $\dot{\pi}(t) \leq 0$ when $\pi(t) = 1$

↪ Note that at some point $x(t) \in [\pi(0), 1]$, $\dot{\pi}(t) = 0$ by continuity of trajectory

↪ Let $H = \{ \pi | p(\pi) = \pi \} \equiv$ set of equilibria of the ODE (*).

↪ This also means that H is nonempty.

↪ Note :- $x(t) \rightarrow H$ as $t \rightarrow \infty$.

We then argue under some conditions that algorithm (1) satisfies

$$x_n \rightarrow H \text{ as } n \rightarrow \infty \text{ w.p. 1.}$$

Assumption

1) $p : [0, 1] \rightarrow [0, 1]$ is Lipschitz.

2) $a_n > 0 \quad \forall n \quad [a(n) = \frac{1}{n+1} \text{ here}]$

$$\downarrow \quad \overbrace{\quad}^{\quad} \quad \downarrow$$

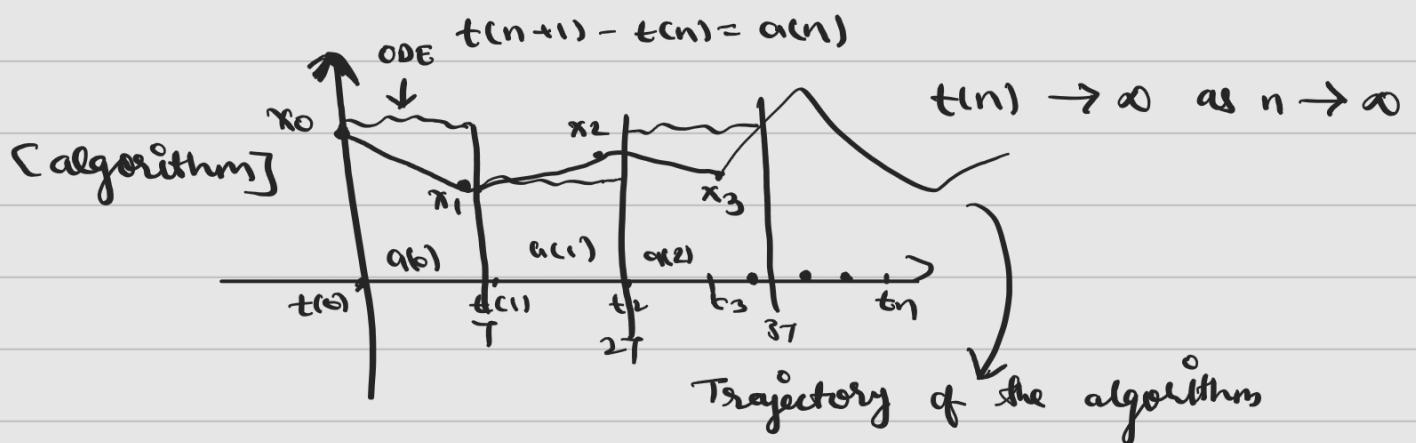
$$\sum_n a(n) = \infty, \quad \sum_n (a(n))^2 < \infty$$

Example :- $a(n) = \frac{1}{n+1}, \quad \frac{1}{\ln(n+1)}, \quad \alpha \in (0.5, 1]$

$$\frac{\log(n+1)}{(n+1)}, \quad \frac{1}{(n+1)\log(n+1)}, \quad n \geq 1$$

Suppose $t(n) = \sum_{i=0}^{n-1} a(i), \quad n \geq 0$

with $x(0) = 0$



Let $\gamma > 0$ be fixed.

Second condition for

$$\sum_n a(n)^2 < \infty$$



Errors due to noise or discretization asymptotically vanish.

generally

$$Z_n = \sum_{m=0}^n \left(\frac{a(m)}{m+1} M_{m+1} \right), n \geq 0$$



Martingale sequence.

↪ Due to the fact that $\sum a(n)^2 < \infty$.

$$Z_n \xrightarrow{\text{a.s.}} Z_\infty \text{ as } n \rightarrow \infty \Rightarrow \sum_{m=n}^\infty a(m) M_{m+1} \rightarrow 0 \text{ as } n \rightarrow \infty \text{ a.s.}$$

3) $E[M_{n+1} | x_m, M_m, m \leq n] = 0 \quad \forall n$

$$E[\|M_{n+1}\|^2 | x_m, M_m, m \leq n] \leq k(1 + \|x_n\|^2)$$

4) (Crucial requirement [nontrivial])

$$\sup_n \|x_n\| < \infty \text{ w.p. 1}$$
$$\sup_n C(\omega) < \infty$$

Does not mean $\sup_n C(\omega) < \infty$.

2nd setting :- Finding a local minimum.

Consider a repeated experiment with $i/p - O/p$ pairs (x_n, y_n) with $x_n \in \mathbb{R}^m$, $y_n \in \mathbb{R}^k$.

\downarrow
 i/p O/p

Assume $(x_n, y_n), n \geq 0$ are i.i.d.

Goal :- Find a best fit

$$y_n = f_w(x_n) + \epsilon_n$$

Here we are given a parameterized family of functions

$$f_w : \mathbb{R}^m \rightarrow \mathbb{R}^k$$

with ϵ_n as the error.

Let $g(w) = \frac{1}{2} E[\|\epsilon_n\|^2]$

$$g(w) = \frac{1}{2} E[\|y_n - f_w(x_n)\|^2]$$

Goal :- Find a w^* s.t. $g(w)$ is minimized at w^* .

Let $\|\cdot\|$ be the euclidean norm.

$$\text{Thus, } g(w) = \frac{1}{2} E[(y_n - b_w x_n)^T (y_n - f_w(x_n))]$$

Assuming that the gradient can be taken inside the expectation,
we have

$$\nabla_w g(w) = \frac{1}{2} \nabla_w E[(y_n - f_w(x_n))^T (y_n - f_w(x_n))]$$

$$= \frac{1}{2} E[\nabla_w ((y_n - f_w(x_n))^T (y_n - f_w(x_n)))]$$

$$= \frac{1}{2} E[-2 \nabla f_w(x_n)^T (y_n - f_w(x_n))]$$

$$= -E[\nabla f_w(x_n)^T (y_n - f_w(x_n))]$$

A gradient descent scheme could be :-

$$w_{n+1} = w_n + \alpha E[(y_n - f_w(x_n))^T \nabla f_w(x_n)]$$

Problem :- We don't know the distribution of (x_n, y_n) pairs.

Alternative :- $w_{n+1} = w_n + \alpha \epsilon_n ((y_n - f_w(x_n))^T \nabla f_w(x_n))$

$$w_{n+1} = w_n + \alpha \epsilon_n [E[(y_n - f_w(x_n))^T \nabla f_w(x_n)] + M_{n+1}]$$

$$\begin{aligned} M_{n+1} &= (y_n - f_w(x_n))^T \nabla f_w(x_n) \\ &\quad - E[(y_n - f_w(x_n))^T \nabla f_w(x_n)] \end{aligned}$$

$$= w_n + \alpha \epsilon_n [-\nabla_w g(w_n) + M_{n+1}]$$

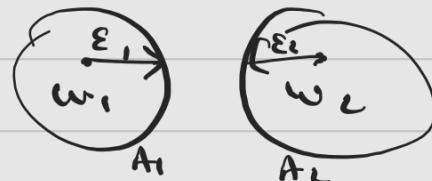
SGD (Stochastic gradient descent scheme)

$$\dot{w}(t) = -\nabla_w g(w(t)) \rightarrow (+)$$

let $H = \{w \mid \nabla_w g(w) = 0\}$



Stationary points of the ODE (+)



$$A_1 \cap A_L = \emptyset$$

Note :- Stationary points are those for which gradient is zero.

⇒ local maxima, local minima, saddle points

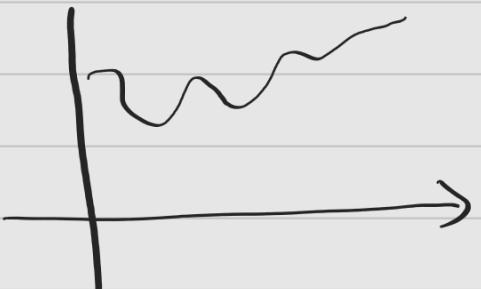
Algorithm :-

$$x_{n+1} = x_n + \alpha(n)(h(x_n) + M_{n+1}), n \geq 0 \quad -(4)$$

↳ Analyze convergence of (4) assuming stability (Chapter 2)
 $(\sup_k \|x_k\| < \infty \text{ w.p. } 1)$

↳ Provide sufficient conditions for stability of (4) (Chapter 3)

↳ loc. -in- probability of SA. (Chapter 4)



↳ Consider processes such as

$$x_{n+1} = x_n + \alpha(n)(f(x_n, \xi_n) + M_{n+1}), n \geq 0$$

Here $\{\xi_n\}$ is parameter dependent Markov process.

$\Rightarrow P(\xi_{n+1} = y' | \xi_n = y, x_n = x) = p_x(y, y')$ Under ergodicity of
 (Chapter 7) Markov process, $h(x) = \int f(x, y) v_x(dy)$
 $\Rightarrow x(t) = h(x)$ ← stationary distribution of $\{\xi_n\}$ for a fixed x .

↳ Algorithm with set-valued maps. (Stochastic recursive inclusions)

Chapter 5-6

$$h(x_n) \subset \mathbb{R}^d$$

↳ Differential inclusion

$$\dot{x}(t) \in h(x(t))$$

↳ If process is not ergodic \Rightarrow stationary distribution is not unique.

$$H(x) = \left\{ \int f(x, y) v_x(dy) \mid y \right\}$$

$$DI : \dot{x}(t) \in H(x(t)), \quad t \geq 0$$

Some basics on Martingales

[Chapter 10 of Bruce Hajek's book]

Suppose X and Y be two discrete r.v.s.

Then

$$g(i) \triangleq E[X|Y=i] = \sum_j m_j p(X=m_j | Y=i) \quad \text{--- (1)}$$

Here $X \in \{m_1, m_2, \dots\}$

Note :- (1) is well defined if $p(Y=i) > 0$

Also, either the sum restricted $\{j | u_j \geq 0\}$ or $\{j | u_j < 0\}$ is convergent.

↳ Define r.v. $E[X|Y] = g(Y)$ where $g(i) = E[X|Y=i]$

Similarly, if X and Y are continuous r.v.s having a joint PDF.

$$E[X|Y=y] = \int x f_{X|Y}(x|y) dx = g(y)$$

and $E[X|Y] = g(Y)$.

↳ Given a set \mathcal{F} , a set of subsets \mathcal{D} of Ω is called a σ -algebra if

(i) $\emptyset \in \mathcal{D}$

(ii) If $A \in \mathcal{D} \Rightarrow A^c \in \mathcal{D}$

(iii) If $A_1, A_2, \dots, A_n \in \mathcal{D} \Rightarrow \bigcup A_i \in \mathcal{D}$

↳ In particular, we require that the set of events \mathcal{F} in a probability space (Ω, \mathcal{F}, P) to be a σ -algebra.

↳ We say that \mathcal{D} is a sub σ -algebra of \mathcal{F} if \mathcal{D} is a σ -algebra and $\mathcal{D} \subset \mathcal{F}$.

↳ We say that a r.v. Z is \mathcal{D} measurable if

$f(z) \in \mathcal{G} \subset D \quad \forall z \in \mathbb{R}.$

Definition :- σ -algebra generated by a set of r.v's $(Y_i, i \in I)$ denoted by $\sigma(Y_i : i \in I)$ is the smallest σ -algebra containing all sets of the form $\{Y_i \leq c\}, i \in I, c \in \mathbb{R}$.

Fact :- The smallest σ -algebra exists and equals the intersection of σ -algebras $\sigma(Y_i), i \in I$.

↪ σ -algebra generated by a random variable Y will be denoted by $\sigma(Y)$ or \mathcal{F}^Y .

↪ Borel - measurable functions

Let $g : \mathbb{R} \rightarrow \mathbb{R}$. We say that g is Borel measurable if $\{x | g(x) \leq c\} \in \mathcal{B}$ $\forall c \in \mathbb{R}$.

Here \mathcal{B} is the σ -algebra of open subsets of \mathbb{R} .

$$g((-\infty, c]) \in \mathcal{B} \quad \forall c \in \mathbb{R}$$