

Lecture - 8 (contd)

INTRODUCTION TO STOCHASTIC APPROXIMATION

e.g. ① find
maximizer

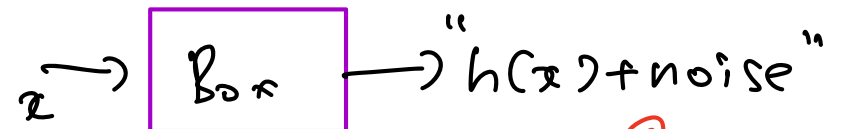
$$\max_x f(x)$$

$$\nabla f(x) = 0 \quad \text{or} \quad h = \nabla f$$

② $f(x) = x \in \text{fixed point}$

$$h(x) = f(x) - x$$

Problem: Solve $h(x) = 0$ given noisy measurements of h , i.e., given access to a black box that, on input $x \in \mathcal{R}^d$, gives as output $h(x) + \text{noise}$.



Robbins-Monro algorithm: Starting with $x_0 \in \mathcal{R}^d$, do:

Stochastic
root
finding

$$x(n+1) = x(n) + a(n)[h(x(n)) + M(n+1)], \quad n \geq 0.$$

stepsize

Here the stepsize sequence (or 'learning parameter') $\{a(n)\}$ satisfies: $a(n) \geq 0$ and

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

(\implies slow decrease to zero, e.g., $\frac{1}{n}$, $\frac{1}{n \log n}$, $\frac{1}{n^{2/3}}$ etc.).

$$(A0) \quad \sum a(n) = \infty, \quad \sum a(n)^2 < \infty$$

(A1) $h : \mathcal{R}^d \mapsto \mathcal{R}^d$ is Lipschitz: $\|h(x) - h(y)\| \leq L\|x - y\|$
 for $x, y \in \mathcal{R}^d$.

\downarrow
Euclidean norm

(A2) $\{M(n)\}$ a square-integrable martingale difference sequence, i.e., for

$$\mathcal{F}_n := \underbrace{\sigma(x_0, M_m, m \leq n)}_{\text{"information up to time 'n'"}} , n \geq 0,$$

we have

$$E[\|M(n)\|^2] < \infty$$

and in addition, it is 'uncorrelated with past',

i.e.,

$$E[M_i(n+1)|\mathcal{F}_n] = 0 \quad \forall i.$$

$M_i(\cdot)$ is a martingale difference

. (Equivalently,

$$E[M_i(n+1)|x_0, M_m, m \leq n] = 0 \quad \forall i.)$$

. Thus

$$\begin{aligned} & E[M_i(n+1)f(x_0, M_1, \dots, M_n)] \\ &= E[E[M_i(n+1)f(x_0, M_1, \dots, M_n)|x_0, M_m, m \leq n]] \\ &= E[E[M_i(n+1)|x_0, M_m, m \leq n]f(x_0, M_1, \dots, M_n)] \\ &= 0. \end{aligned}$$

Hence 'uncorrelated with past'.

Furthermore, we assume that for some $K > 0$, Part of Assumption (A2)

$$E [\|M(n+1)\|^2 | \mathcal{F}_n] \leq K (1 + \|x(n)\|^2) \quad \forall n \geq 0.$$

(equivalently,

$$E [\|M(n+1)\|^2 | x_0, M_m, m \leq n] \leq K (1 + \|x(n)\|^2) \quad \forall n \geq 0.)$$

In particular, if

(A3)

$$\sup_m \|x(n)\| < \infty \text{ a.s.,}$$

"It rate is bounded"

we have

$$\sup_n E [\|M(n+1)\|^2 | \mathcal{F}_n] < \infty \text{ a.s.}$$

Convergence claim: A rough introduction

$$\textcircled{1} \rightarrow x(n+1) = x(n) + a(n) \underbrace{(h(x(n)) + M_{n+1})}_{(*)}$$

$$x(n) \xrightarrow{a.s.} x^* \text{ s.t. } h(x^*) = 0 \quad \left[\begin{array}{l} \text{Assume} \\ x^* \text{ is} \\ \text{unique} \end{array} \right]$$

if (i) h is Lipschitz

$$(ii) \sum a(n) = \infty \quad \sum a(n)^2 < \infty$$

$$(iii) E(M_{n+1} | \mathcal{F}_n) = 0$$

$$E(\|M_{n+1}\|^2 | \mathcal{F}_n) \leq K(1 + \|x(n)\|^2)$$

$$(iv) \sup_n \|x(n)\| < \infty$$

Tip:

Lecture-9

To infer the limit of $\textcircled{1}$, take the conditional expectation of $(*)$ w.r.t \mathcal{F}_n & compare it to zero

$$\begin{aligned} \text{e.g., } x(n+1) &= x(n) - a(n) \left(\hat{\nabla} f(x(n)) \right) \\ &= x(n) - a(n) \left(\nabla f(x(n)) + M_{n+1} \right) \\ M_{n+1} &= \hat{\nabla} (f(x(n)) - \nabla f(x(n))) \end{aligned}$$

Under suitable assumption, $x(n) \xrightarrow{a.s.} x^*$ as $n \rightarrow \infty$, where $\nabla f(x^*) = 0$

This is more general than it appears. Suppose the algorithm is

$$x(n+1) = x(n) + a(n)f(x(n), \xi(n+1)), \quad n \geq 0,$$

where $\{\xi(n)\}$ are IID. This is often how many recursive algorithms are stated.

This can be put in the above form by letting

$$h(x) = E[f(x, \xi(n))] = E[f(x(n), \xi(n+1)) | x(n) = x]$$

$$= E[f(x(n), \xi(n+1)) | \mathcal{F}_n],$$

$$M(n+1) = f(x(n), \xi(n+1)) - h(x(n)).$$

$$E(M(n+1) | \mathcal{F}_n) = 0$$

noise term

$$= x(n) + a(n) \left(E(f(x(n), \xi(n+1)) | \mathcal{F}_n) + M(n+1) \right)$$

Examples: Stochastic Gradient Descent ($h = -\nabla f$),
reinforcement learning algorithms (more later)

Highlights:

1. Typically small amount of computation and memory requirements per iterate (Contrast: batch-mode algorithms)

2. Incremental: makes a small change in the current iterate at each step

$$a(n) \propto \frac{1}{n}$$

3. Slowly decreasing stepsize captures 'exploration' (\approx large steps initially) vs 'exploitation' (\approx small steps later) trade-off

① $x(n) \xrightarrow{a.s.} x^*$ as $n \rightarrow \infty$ "SLLN spirit"

② $\sqrt{n}(x(n) - x^*) \xrightarrow{d} N(0, \Sigma)$ "CLT spirit"

4. Averages out the noise (can be thought of as a generalization of the Strong Law of Large Numbers)

③ $P(\|x(n) - x^*\| > \epsilon) \leq \exp(-cn\epsilon^2)$ "Concentration bound"

1.-3. typical of adaptive behavior \Rightarrow extremely well suited for adaptive algorithms or models of adaptation

One of the two main workhorses of statistical computation, MCMC being the other.

Markov Chain Monte Carlo

portfolio
 $P(\hat{\mu}_n - \mu > \epsilon) \leq \exp(-cn\epsilon^2)$

Applications:

statistics, signal processing, machine learning, adaptive control and communications

Also for models of learning, bounded rationality, herding behavior, etc.

Classical approach for analysis: uses 'almost supermartingales' etc. (Robbins-Siegmund, ...)

Alternative approach: ODE (Ordinary Differential Equations) approach (Meerkov '72, Derevetskii-Fradkov '74, Ljung '77)

ODE approach: Treat the iterates as a noisy discretization of the ODE

$$\dot{x}(t) = h(x(t)).$$

Recall the Euler scheme for this ODE:

$$x(n+1) = x(n) + ah(x(n)), \quad n \geq 0,$$

where $a > 0$ is a small discrete time step.

$$x(n+1) = x(n) + a(n) \left(h(x(n)) + \overset{\text{noise}}{\mu_{n+1}} \right)$$

Then SA can be viewed as an Euler scheme to approximate the ODE with slowly decreasing time steps $\{a(n)\}$ and measurement noise.

Robbins-Monro conditions:

$\sum_n a(n) = \infty \implies$ the entire time axis is covered. This is essential because we want to track the asymptotic (as $t \uparrow \infty$) behavior of the ODE.

$\sum_n a(n)^2 < \infty \implies$ the approximation of the ODE gets better with time:

$a(n) \rightarrow 0$ ensures that errors due to discretization are asymptotically zero

$\sum_n a(n)^2 < \infty$ ensures that errors due to the martingale difference noise are asymptotically zero, a.s. (multiplication by $a(n)$ reduces the (conditional) variance of noise)

Advantages:

1. Once you have mastered the approach, you can often write the limiting ODE by inspection and analyze it.
2. Designing algorithms: any convergent ODE is a template for an algorithm.
3. Finer dynamic phenomena lead to useful results, e.g., avoidance of unstable equilibria a.s. \implies avoidance of 'traps' (undesirable equilibria)

Analogy with SLLN suggests related results for fluctuations, e.g., central limit theorem, law of iterated logarithms, concentration inequalities

Further issues and variations:

stability tests, multiple timescales, distributed and asynchronous implementations, differential inclusion limits, constant stepsizes, other noise models, etc.

"Stochastic fixed point iterations

Covered on blackboard"

An application: Mean estimation

Consider a random variable (r.v.) X with mean μ
& finite variance, say σ^2 .

Suppose we are given iid samples X_1, \dots, X_m

Let r_m be the estimate of μ .

Sample mean $\rightarrow r_m = \frac{1}{m} \sum_{k=1}^m X_k$

$$r_{m+1} = \frac{1}{m+1} \sum_{k=1}^{m+1} X_k$$

$$= \frac{m}{m+1} \left(\frac{1}{m} \sum_{k=1}^m X_k \right) + \frac{1}{m+1} X_{m+1}$$

$$r_{m+1} = \frac{m}{m+1} r_m + \frac{1}{m+1} X_{m+1} \quad \leftarrow \text{iterative scheme for updating sample mean}$$

$$r_{m+1} = r_m + \frac{1}{m+1} (X_{m+1} - r_m) \quad \begin{matrix} \text{new sample} \\ \text{prev. average} \end{matrix}$$

$$r_{m+1} = r_m + \beta_m (X_{m+1} - r_m) \quad \begin{matrix} \text{step size} \\ \text{new sample} \end{matrix} \quad \text{--- } (**)$$

$\beta_m = \frac{1}{m+1}$

where $\beta_m = \frac{1}{m+1}$

(**) resembles the sto. iter. algo (*)

What does Strong Law of large numbers say about r_m ?

$$r_m \rightarrow \mu \text{ a.s. as } m \rightarrow \infty$$

with step size $\beta_m = \frac{1}{m+1}$

$$r_{m+1} = r_m + \beta_m (x_{m+1} - r_m)$$

$$= r_m + \beta_m ((\mu - r_m) + (x_{m+1} - \mu))$$

$$r_{m+1} = (1 - \beta_m) r_m + \beta_m (\underbrace{\mu}_{Hr_m} + \underbrace{(x_{m+1} - \mu)}_{w_m})$$

$$E w_m = 0, \quad E w_m^2 < \infty, \quad \cancel{r^* = r^*} \Leftrightarrow \cancel{r^* = \mu}.$$

So, $(**)$ is really a sto. iterative algorithm.

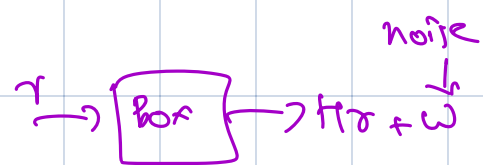
And, we get $r_m \rightarrow r^* (= \mu)$ a.s. as $m \rightarrow \infty$.

The theory of sto. itr. algo that we shall develop ensures $r_m \rightarrow \mu$ a.s. as $m \rightarrow \infty$ for more general step sizes that satisfy

$$\sum_m \beta_m = \infty, \quad \sum_m \beta_m^2 < \infty \quad \text{e.g. } \beta_m = \frac{1}{m}$$

Question: Why do we need these conditions?

On step-size requirements:



$$r_{m+1} = (1 - \beta_m) r_m + \beta_m (Hr_m + w_m)$$

$$r_{m+1} = r_m + \beta_m (\underbrace{Hr_m + w_m - r_m}_{\text{increment}})$$

Suppose $\{w_i\}$ is independent of $r_m, \forall m$
& has variance σ^2 , e.g. $w \sim N(0, \sigma^2 I)$

$$\begin{aligned} \text{Variance of } r_{m+1} \\ &= \text{Var}(r_{m+1}) \end{aligned}$$

$$= \text{Var}[(1 - \beta_m) r_m + \beta_m Hr_m] + \beta_m^2 \text{Var}(w_m)$$

$$= \text{Var}[(1 - \beta_m) r_m + \beta_m Hr_m] + \beta_m^2 \sigma^2$$

$$\geq \beta_m^2 \sigma^2$$

Now, if $\beta_m = \beta \forall m$, then

$$\text{Var}(r_{m+1}) \geq \beta^2 \sigma^2$$

So, fixing $\beta > 0$, $r_m \not\rightarrow r^*$ for any r^*

$\text{Var}(X+Y) = \text{Var}X + \text{Var}Y$
if X, Y independent.
 $\text{Var}(cX) = c^2 \text{Var}X$

$\beta_m = \beta$ "constant step size"
 $\beta_m = \frac{1}{m}$ "diminishing step size"

Let $r_m \rightarrow r^*$ s.t.
 $Ur^* = r^*$

Aside: constant step size stor-approx algos \rightarrow convergence guaranteed are in "distribution" & not "a.s."

So, the step-size has to vanish asymptotically. i.e., $\beta_m \rightarrow 0$ as $m \rightarrow \infty$.

But, β_m cannot go down too fast.

$$r_{m+1} = r_m + \beta_m (Hr_m + w_m - r_m)$$

increment

$$|r_m - r_0| \leq \sum_{\tau=0}^{m-1} \beta_\tau |Hr_\tau - r_\tau + w_\tau|$$

triangle inequality

sum of increments

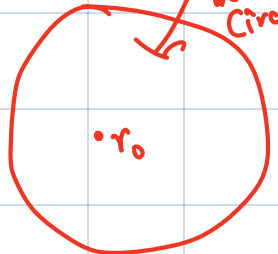
So, if $|Hr_\tau - r_\tau + w_\tau| \leq C_1$ and

if $\sum_{\tau=0}^{\infty} \beta_\tau \leq C_2 < \infty$, then

$|r_m - r_0|$ is bounded above $\forall m$

r_m is restricted to be inside Circle

outside circle



(\Rightarrow) r_m is within a certain radius of r_0

problematic if r^* lies outside the radius.

So, we need $\sum_{\tau} \beta_\tau = \infty$

" $\sum \beta_t = \infty$ & $\beta_t \rightarrow 0$ as $t \rightarrow \infty$ "

$$\beta_t = \frac{1}{t^{1/4}}$$