

## Lecture - 10

### "Kushner - Clark lemma & applications"

$$(1) \quad x_{n+1} = x_n + a(n) ( h(x_n) + M_{n+1})$$

↘ step size    ↗ want to find  $h(x) = 0$     ↗ noise

$$(2) \quad x_{n+1} = x_n + a(n) ( h(x_n) + M_{n+1} + \beta_n)$$

↓  
 additional error

Assumptions:-

ODE

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

(A1)  $h: \mathbb{R}^d \rightarrow \mathbb{R}$  Lipschitz

(A2)  $\{\beta_n\}$  bounded (random) seq, s.f.  $\beta_n \rightarrow 0$

(e.g.  $T_{n+1} = T_n + a(n) (-\nabla f(x_n) + M_{n+1} + C, \xi_n^2)$ )

(A3)  $a(n) \rightarrow 0$  and  $\sum_n a(n) = \infty$

(A4)  $\{M_{n+1}\}$  seq s.t.  $\forall \epsilon > 0$

$$\lim_{n \rightarrow \infty} P \left( \sup_{m \geq n} \left| \left| \sum_{i=n}^m a(i) M_{i+1} \right| \right| > \epsilon \right) = 0$$

Note:  $\{M_{n+1}\}$  is not assumed to be a martingale diff.

(A5)

Bounded iterates:  $\sup_n \|x_n\| < \infty$

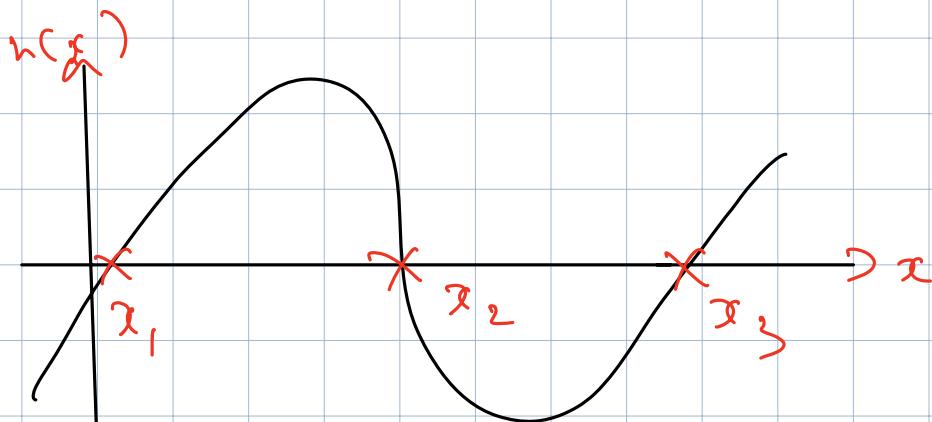
Kushner-Clark lemma

Under

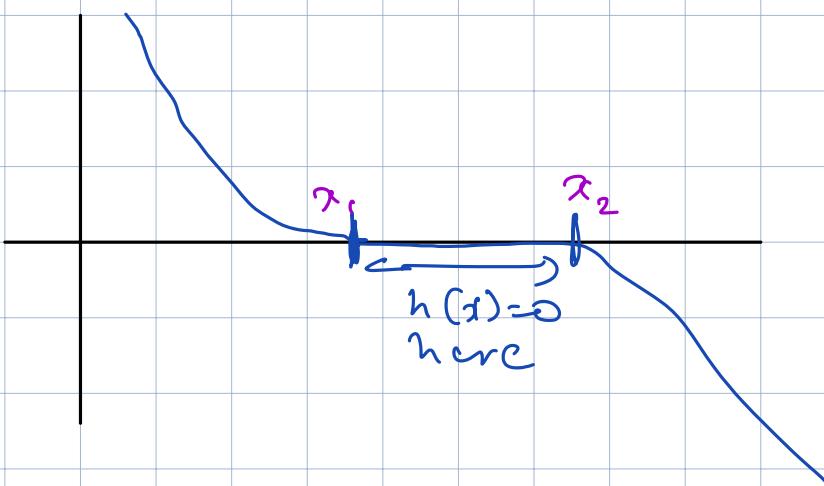
(A1) - (A5),

$x_n \rightarrow \{x^* \mid h(x^*) = 0\}$  a.s. as  $n \rightarrow \infty$

Converge to a set



K-C lemma says conv to one  
of these points  $\{x_1, x_2, x_3\}$



K-C lemma says  
conv to  
 $[x_1, x_2]$

Case of a MDSS  $\{M_{n+1}\}$ :

$$\mathcal{F}_n = \sigma(x_1, \dots, x_n)$$

A6  $E[M_{n+1} | \mathcal{F}_n] = 0$

$$E[|M_{n+1}|^2 | \mathcal{F}_n] \leq C(1 + \|x_n\|^2)$$

(A3')  $\sum a(n) = \infty, \quad \sum a(n)^2 < \infty$

K-C lemma variant:

Under A1, A2, A3', A5, A6,

$x_n \rightarrow \{x^* \mid h(x^*) = 0\}$  a.s. as  $n \rightarrow \infty$

Background "Doob's martingale inequality"

$\{W_m\}$  martingale

(\*)  $P\left(\sup_{m \geq 0} \|W_m\| \geq \epsilon\right) \leq \frac{1}{\epsilon^2} \lim_{n \rightarrow \infty} E\|W_n\|^2$

$$\sum_{n=k}^i a(n) M_{n+1} \quad \leftarrow \text{martingale}$$

Apply

Doob's inequality to

$\downarrow$   
 $a$

$$\lim_{k \rightarrow \infty} P \left( \sup_{n \geq k} \left\| \sum_{n=k}^l a(n) M_{n+1} \right\| \geq \epsilon \right)$$

$$\leq \frac{1}{\epsilon^2} \lim_{k \rightarrow \infty} \sum_{n=k}^{\infty} a(n)^2 E \|M_{n+1}\|^2$$

$$\leq \frac{\text{const}}{\epsilon^2} \lim_{k \rightarrow \infty} \sum_{n=k}^{\infty} a(n)^2$$

Since  
 $\sum a(n)^2 < \infty$

$$= 0$$

So, (A4) of  $k$ -c lemma is satisfied

The claim follows

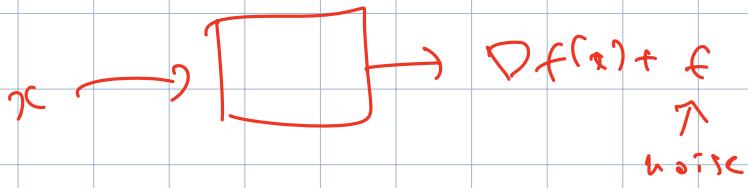
Coming next:

Stochastic gradient algorithm

$$x_{n+1} = x_n - \alpha(n) \hat{\nabla} f(x_n)$$

$$E(\hat{\nabla} f(x_n) | \mathcal{F}_n) = \nabla f(x_n)$$

"unbiased gradient  
Oracle"



$$E(\hat{\nabla} f(x_n) | \mathcal{F}_n) \neq \nabla f(x_n)$$

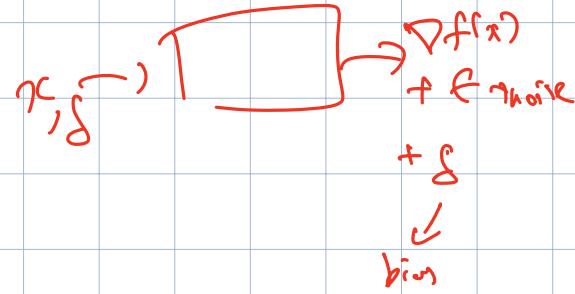
$$E(\hat{\nabla} f(x_n) | \mathcal{F}_n)$$

$$= \nabla f(x_n)$$

$$\leq C \delta_n^2$$

"Zeroth-order  
Optimization"

Biased gradient  
oracle



Recall:

Simultaneous perturbation  
gradient estimation

$$\hat{\nabla}_i f(x_n) = \frac{f(x_n + \delta_n \Delta_n^i) - f(x_n - \delta_n \Delta_n^i)}{2 \delta_n \Delta_n^i} \quad \Delta_n^i \rightarrow i^{\text{th coord}}$$

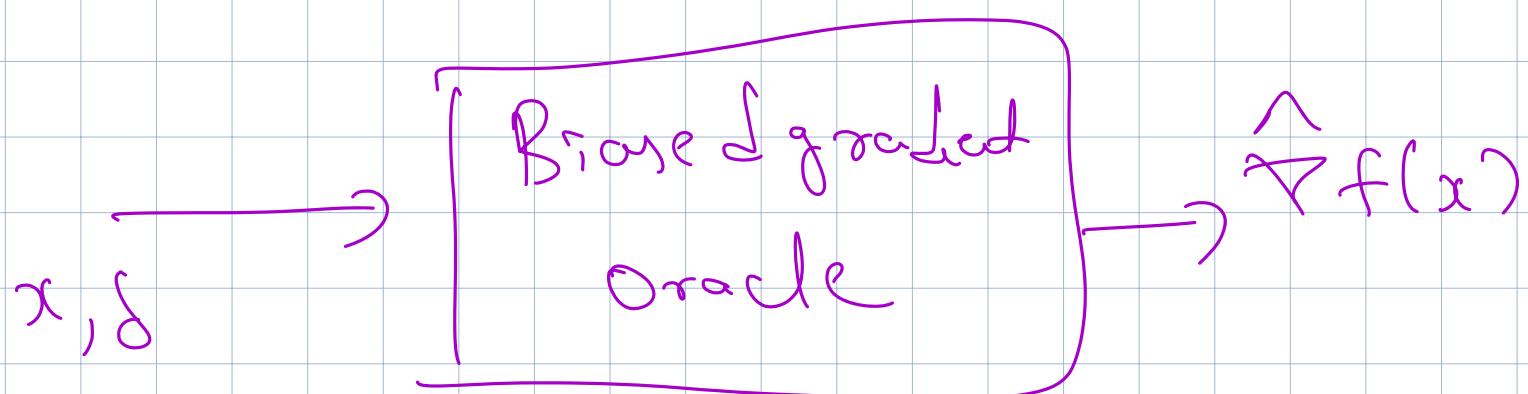
$$\Delta_n^i = (\Delta_{n,1}^i, \dots, \Delta_{n,d}^i) \quad \Delta_n^i \rightarrow \text{hadamard}$$

$$E(\hat{\nabla} f(x_n) | \mathcal{F}_n) - \nabla f(x_n) \leq C \delta_n^2$$

(assuming  $f$  is smooth, etc)

$$\hat{\nabla}_i f(x_n) = \underbrace{f(x_n + \delta_n \Delta_n^i) + \xi_n^+ - (f(x_n - \delta_n \Delta_n^i) + \xi_n^-)}_{2 \delta_n \Delta_n^i \rightarrow 0 \text{ word}}$$

Assume  $E(\xi_n^+ - \xi_n^- | \mathcal{F}_n) = 0$



$$E[\hat{\nabla} f(x) | x] - \nabla f(x) \leq C \delta^2$$

$$E[||\hat{\nabla} f(x) - E(\hat{\nabla} f(x) | x)||^2] \leq \frac{C}{\delta^2}$$

K-C lemma variant (MDS)

$$x_{n+1} = x_n + \alpha(n)(h(x_n) + M_{n+1} + \beta_n)$$

Cone

(B1)  $h: \mathbb{R}^d \rightarrow \mathbb{R}$  Lipschitz

(B2)  $\{\beta_n\}$  bounded (random) seq s.f.  $\beta_n \rightarrow 0$

(B3)  $E[M_{n+1} | \mathcal{F}_n] = 0$

$$E[\|M_{n+1}\|^2 | \mathcal{F}_n] \leq C(1 + \|x_n\|^2)$$

(B4)  $\sum \alpha(n) = \infty, \quad \sum \alpha(n)^2 < \infty$

(B5) Bounded iterates:  $\sup_n \|x_n\| < \infty$

Kushner-Clark lemma

Under (B1) - (B5),

$x_n \rightarrow \{x^* \mid h(x^*) = 0\}$  a.s. as  $n \rightarrow \infty$

# Sto-grad algo with noise

gradient information!

$$\begin{aligned}x_{n+1} &= x_n + \alpha(n) (-\hat{\nabla} f(x_n)) \\&= x_n + \alpha(n) (-\nabla f(x_n) + M_{n+1})\end{aligned}$$

(C1)  $f$  is  $L$ -smooth ( $\Leftrightarrow \nabla f$  is  $L$ -Lipschitz)

(C2)  $E(M_{n+1} | \mathcal{F}_n) = 0$  (or  $E(\hat{\nabla} f(x_n) | \mathcal{F}_n) = \nabla f(x_n)$ )

$$E(\|\hat{\nabla} f(x_n) - E(\hat{\nabla} f(x_n) | \mathcal{F}_n)\|^2) \leq \sigma^2$$

(C3)  $\sum \alpha(n) = \infty$ ,  $\sum \alpha(n)^2 < \infty$

(C4)  $\sup_n \|x_n\| < \infty$  a.s.

Claim: Under (C1)-(C4),

$$x_n \rightarrow \{x^* \mid \nabla f(x^*) = 0\} \text{ a.s as } n \rightarrow \infty$$

Proof:-

(1)  $C \Rightarrow B$  "Def Lipschitz"

(2)  $\beta_n = 0$  ( $B_2$ ) trivial

(3)  $(C_2) \Rightarrow (B_3) \{M_{n+1}\}$

(4)  $(C_3) \Rightarrow (B_4)$  step sizes

(5)  $(C_4) \Rightarrow (B_5)$  bounded  
iterates

Hence proved.



H.W. We K-C lemma to

establish asymp. convergence

of a stochastic fixed point iteration

$$x_{n+1} = x_n + \alpha(n) (f(x_n) - x_n) + M_{n+1}$$

"Make reasonable assumptions" for K-C lemma invocation.

c-s.  $\|f(x) - f(y)\| \leq \lambda \|x - y\|$   
 $0 < \lambda < 1$

---

Sto-gradient algo with biased gradient info:

$$x_{n+1} = x_n + \alpha(n) (-\hat{\nabla} f(x_n)) \quad (\text{xx})$$

$$x_n, \delta_n \rightarrow \boxed{\text{BGO}} \rightarrow \hat{\nabla} f(x_n)$$

$$(D1) \quad \|E(\hat{\nabla} f(x_n) | \mathcal{F}_n) - \nabla f(x_n)\| \leq C_1 \delta_n^2$$

$$\left\{ E \|\hat{\nabla} f(x_n) - E(\hat{\nabla} f(x_n) | \mathcal{F}_n)\|^2 \leq \frac{C_2}{\delta_n^2} \right.$$

bias variance  
tradeoff

for some constants  $C_1, C_2, \forall n \geq 1$

Rewriting  $(\star \star)$

$$x_{n+1} = x_n + \alpha(n) \left( -\hat{\nabla} f(x_n) \right)$$

$$x_{n+1} = x_n + \alpha(n) \left( -\nabla f(x_n) + M_{n+1} + \beta_n \right)$$

So that

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

$$\beta_n \subseteq C_1 S_n^2$$

Side-note: Unbiased Case

$$x_{n+1} = x_n + \alpha(n) \left( -\hat{\nabla} f(x_n) \right)$$

$$= x_n + \alpha(n) \left( -\nabla f(x_n) \right)$$

$$- (\hat{\nabla} f(x_n) - \nabla f(x_n))$$

$M_{n+1}$

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

B60 Core :

$$x_{n+1} = x_n - \alpha(n) (\hat{\nabla} f(x_n))$$

$$\begin{aligned} x_{n+1} &= x_n - \alpha(n) \\ &\quad \left. \begin{aligned} & \hat{\nabla} f(x_n) - E(\hat{\nabla} f(x_n) | \mathcal{F}_n) \\ & + E(\hat{\nabla} f(x_n) | \mathcal{F}_n) - \nabla f(x_n) \\ & + \nabla f(x_n) \end{aligned} \right\} \\ &= x_n - \alpha(n) (\nabla f(x_n) + M_{n+1} + B_n) \end{aligned}$$

### Lecture - 11

Today → ① Asymptotic analysis of 200-SGD

② Variants of zeroth order gradient estimation

$$x_{n+1} = x_n - \alpha(n) \left[ \hat{f}(x_n) - E(\hat{f}(x_n) | \mathcal{F}_n) + E(\hat{f}(x_n) | \mathcal{F}_n) - \nabla f(x_n) + \nabla f(x_n) \right]$$

$$= x_n - \alpha(n) (\nabla f(x_n) + M_{n+1} + \beta_n)$$

(D1) + the following assumptions:

(D2)  $\alpha(n), \delta_n \rightarrow 0$  as  $n \rightarrow \infty$

$$\sum_n \alpha(n) = \infty$$

$$\sum_n \left( \frac{\alpha(n)}{\delta_n} \right)^2 < \infty$$

H.W.

Suppose

$$\alpha(n) = \frac{1}{n^b}$$

$$\delta_n = \frac{1}{n^c}$$

Under what  
choice for  $b, c$

is (D2)  
Satisfied?

(P3)  $f$  is  $L$ -smooth

(P4)  $E \|M_{n+1}\|^2 \leq \frac{C}{\delta_n^2}, \forall n$

Claim: Under some assumptions,

$$x_n \rightarrow \{x^* \mid \nabla f(x^*) = 0\} \text{ a.s as } n \rightarrow \infty$$

Pf:

- ①  $f$  is  $L$ -smooth  $\Rightarrow$  (B1) of K-C lemma satisfied
- ②  $\beta_n = O(\zeta_n^2)$  & by assumption (D2),  
 $\zeta_n \rightarrow 0$ . So,  $\beta_n \rightarrow 0 \Rightarrow$  (B2) of  
K-C lemma satisfied.
- ③ K-C lemma required

$$\lim_{n \rightarrow \infty} P \left( \sup_{m \geq n} \left\| \sum_{i=n}^m a_i c_i M_{i+1} \right\| \geq \epsilon \right) = 0$$

↓  
Denote this as  $W_m$   $\rightarrow (\star\star)$

$\{W_m\}$  is a martingale

$\{W_m\}$  martingale

$$\textcircled{X} P \left( \sup_{m \geq 0} \|W_m\| \geq \epsilon \right) \leq \frac{1}{\epsilon^2} \lim_{n \rightarrow \infty} E \|W_n\|^2$$

Apply this to  $W_m = \sum_{i=n}^m a(i) M_{i+1}$

$$P \left( \sup_{m \geq n} \left\| \sum_{i=n}^m a(i) M_{i+1} \right\| \geq \epsilon \right)$$

$$\leq \frac{1}{\epsilon^2} E \left\| \sum_{i=n}^{\infty} a(i) M_{i+1} \right\|^2$$

$$\leq \frac{1}{\epsilon^2} \sum_{i=n}^{\infty} a(i)^2 \underbrace{E \|M_{i+1}\|^2}_{\text{Cov}}$$

(Used)

$$E [M_m M_n] \quad \text{for } m < n$$

$$= E [M_m E(M_n | \mathcal{F}_m)]$$

$\Rightarrow$

Var  
(D.L.)

$$\leq \frac{1}{\epsilon^2} \sum_{i=n}^{\infty} \frac{a(i)^2}{\sigma_i^2}$$

Since

$$\sum_{i=r}^{\infty} \frac{a(c_i)^2}{\delta_i^2} < \infty$$

$$\lim_{n \rightarrow \infty} \sum_{i=n}^{\infty} \frac{a(c_i)^2}{\delta_i^2} = 0$$

So, K-C lemma condition (\*\*) holds.

(4)

Since  $\sup_n \|x_n\| < \infty$ ,

Last condition of K-C lemma holds.

hence, the claim follows.



Why is assumption (P4) reasonable:

For  $i=1, \dots, d$ ,

$$\hat{\nabla}_i f(x_n) = \frac{(f(x_n + \delta_n D_n) + \xi^+) - (f(x_n - \delta_n D_n) + \xi^-)}{2 \delta_n D_n}$$

SPSA gradient estimate

$\Delta_n$  = vector of backtracking  $\tau, v, s$ .

$$E \parallel \hat{\nabla} f(x_n) - E(\hat{\nabla} f(x_n) | \mathcal{F}_n) \parallel^2$$

$$\leq E \parallel \hat{\nabla} f(x_n) \parallel^2 \quad (\text{---})$$

$$\hat{\nabla}_i f(x_n)^2$$

$$= \left( \frac{(f(x^+) + \xi^+) - (f(\bar{x}) + \xi^-)}{2\delta \Delta_i} \right)^2$$

If  $f(x^+)^2, f(\bar{x})^2 \leq C_2$

$$E(\xi^+)^2 \leq C_3, E(\xi^-)^2 \leq C_4$$

$$E \hat{\nabla}_i f(x)^2 \leq \frac{\text{const}}{\delta_n^2}$$

So, my (\*\*)  $\oplus$

$$E \left( \left\| \nabla f(x_s) \right\|^2 \right)$$

Depends on  $\sigma^2$   
const  
 $\frac{1}{n}$