

Analysis of Stochastic Gradient Descent

Review:

approx ERM

$$L_n(f) := \frac{1}{n} \sum_{i=1}^n \ell(f, z_i) \quad \rightarrow \min_{f \in \mathcal{F}}$$

$$(\xi_t) \stackrel{iid}{\sim} \text{Unif}([n])$$

$$f_0 \in \mathcal{F} \quad (\text{init})$$

$$f_t = \Pi \left(f_{t-1} - \alpha_t \underbrace{\nabla \ell(f_{t-1}, \xi_t)}_{\text{Stoch. grad.}} \right) \quad t=1, 2, \dots$$

$$\mathbb{E}_{\xi_t} [\nabla \ell(f_{t-1}, \xi_t)] = \frac{1}{n} \sum_{i=1}^n \nabla \ell(f_{t-1}, z_i)$$

$$t=1, \dots, T: \quad \mathbb{E}[L_n(f_T)] - \min_{f \in \mathcal{F}} L_n(f) \leq ?$$

\uparrow
w.r.t. (ξ_t)

$T \rightarrow \infty$

Stochastic Approximation

— classic paper of Robbins-Monro (1950s)

\mathcal{F} : Hilbert space

$\Gamma: \mathcal{F} \rightarrow \mathbb{R}$: cont. diff., has a finite inf.:

$$\Gamma^* := \inf_{f \in \mathcal{F}} \Gamma(f) < \infty$$

$f_1 \in \mathcal{F}$: init

$(\xi_t)_{t \geq 1}$: iid random elements of some space \mathcal{S}

$g : \mathcal{F} \times \mathcal{S} \rightarrow \mathcal{F}$: stoch. gradients

$$f_{t+1} = f_t - \alpha_t g(f_t, \xi_t) \quad (\text{SA update})$$

$(\alpha_t)_{t \geq 1}$: positive nonincreasing seq. of step sizes

$(f_t)_{t \geq 1}$: random process w/ values in \mathcal{F}

$f_1 \in \mathcal{F}$: deterministic init

$$f_{t+1} = G_t(f_t, \xi_t) \quad \xi_t \text{ iid}$$

$\Rightarrow (f_t)_{t \geq 1}$ is a Markov process

$$\forall V : \mathcal{F} \rightarrow \mathbb{R}$$

$$\mathbb{E}[V(f_t) | f_1, \dots, f_{t-1}] = \mathbb{E}[V(f_t) | f_{t-1}]$$

Goals: i) control expected optimality gap

$$\Delta_t := \mathbb{E}[\Gamma(f_t)] - \Gamma^*$$

ii) control expected sq. norms of gradients:

$$\mathbb{E} \|\nabla \Gamma(f_t)\|^2$$

Examples

$$\Gamma(f) = \frac{1}{n} \sum_{i=1}^n \ell(f, z_i) \quad (\text{emp. loss})$$

$$\text{opt. gap. :} \quad \Gamma(f_t) - \Gamma^* = \left(\underset{f_t}{\text{loss of}} \right) - \left(\underset{\text{loss}}{\text{min. emp}} \right)$$

$$\bullet \sum_t \stackrel{\text{iid}}{\sim} \text{unif}([n]) \quad \mathcal{F} = [n]$$

$$g(f, \mathbb{F}) := \nabla \ell(f, z_{\mathbb{F}})$$

$$\mathbb{E}_{\mathbb{F}} g(f, \mathbb{F}) = \frac{1}{n} \sum_{i=1}^n \nabla \ell(f, z_i) = \nabla \Gamma(f)$$

$$\bullet \sum_t \stackrel{\text{iid}}{\sim} \text{unif} \left(\binom{[n]}{k} \right) \quad 1 \leq k < n$$

$$\mathcal{F} = \binom{[n]}{k} = \text{all subsets of } [n] \text{ of card. } k$$

$$\mathbb{F} = \{i_1, \dots, i_k\} \subset [n]$$

$$g(f, \mathbb{F}) := \frac{1}{k} \sum_{j \in \mathbb{F}} \nabla \ell(f, z_{\mathbb{F}}) \quad [\text{mini-batches}]$$

$$\text{Both cases:} \quad \mathbb{E}_{\mathbb{F}} [g(f, \mathbb{F})] = \nabla \Gamma(f) \\ (\text{unbiased stoch. gradients})$$

$$\text{SGD/SA:} \quad f_{t+1} = f_t - \alpha_t (\nabla \Gamma(f_t) + \text{noise})$$

Assumptions:

1) $\exists 0 < \mu \leq 1$ s.t.

$$\langle \nabla \Gamma(f_t), \mathbb{E}_{\xi_t} [g(f_t, \xi_t)] \rangle \geq \mu \|\nabla \Gamma(f_t)\|^2 \quad \forall t$$

e.g. if $\mathbb{E}_{\xi} [g(f, \xi)] = \nabla \Gamma(f)$, then we have $\mu = 1$ + equality (unbiased stoch. grads)

2) $\exists B \geq 0, B_G \geq 0$ s.t.

$$\mathbb{E}_{\xi_t} \|g(f_t, \xi_t)\|^2 \leq B + B_G \|\nabla \Gamma(f_t)\|^2 \quad \forall t$$

(can show: $B_G \geq \mu^2$)

- can verify directly for above example

$$\xi \sim \text{unif}([0, 1]), \quad g(f, \xi) = \nabla \ell(f, z_{\xi})$$

$$\mathbb{E}_{\xi} \|g(f, \xi)\|^2 = \frac{1}{n} \sum_{i=1}^n \|\nabla \ell(f, z_{\xi_i})\|^2$$

$$\|\nabla \Gamma(f)\|^2 = \left\| \frac{1}{n} \sum_{i=1}^n \nabla \ell(f, z_i) \right\|^2$$

B, B_G would depend on $\max_i \|\nabla \ell(f, z_i)\|^2$

Convex Functions

- M -smooth: $\|\nabla\Gamma(f) - \nabla\Gamma(f')\| \leq M\|f - f'\|$
- m -strongly convex:

$$\Gamma(f') - \left(\Gamma(f) + \langle \nabla\Gamma(f), f' - f \rangle \right) \geq \frac{m}{2} \|f' - f\|^2$$

Lemma Γ M -smooth \Rightarrow

$$\Gamma(f') - \left(\Gamma(f) + \langle \nabla\Gamma(f), f' - f \rangle \right) \leq \frac{M}{2} \|f' - f\|^2$$

$$\Rightarrow m \leq M$$

Fixed step sizes: $\alpha_t = \alpha \quad \forall t$

$$\begin{aligned} \Gamma(f_{t+1}) - \Gamma(f_t) & \quad f_{t+1} - f_t = -\alpha g_t \\ & = \Gamma(f_t - \alpha g_t) - \Gamma(f_t) \quad g_t := g(f_t, \xi_t) \end{aligned}$$

$$\begin{aligned} & \leq \langle \nabla\Gamma(f_t), f_{t+1} - f_t \rangle + \frac{M}{2} \|f_{t+1} - f_t\|^2 \\ & = -\alpha \langle \nabla\Gamma(f_t), g_t \rangle + \frac{M}{2} \alpha^2 \|g_t\|^2 \end{aligned}$$

$$\mathbb{E}_{\xi_t} \left\{ \Gamma(f_{t+1}) - \Gamma(f_t) \right\}$$

$$= -\alpha \langle \nabla\Gamma(f_t), \mathbb{E}_{\xi_t} g_t \rangle + \frac{M}{2} \alpha^2 \mathbb{E}_{\xi_t} \|g_t\|^2$$

$$\leq -\alpha \mu \|\nabla\Gamma(f_t)\|^2 + \frac{M}{2} \alpha^2 (\beta + B_G \|\nabla\Gamma(f_t)\|^2)$$

$$= \frac{\alpha^2 MB}{2} - \alpha \left(\mu - \frac{\alpha MB_G}{2} \right) \|\nabla\Gamma(f_t)\|^2$$

want to be > 0

relate to $\Gamma(f_t) - \Gamma^*$

$$- \text{choose } \alpha M B_G \leq \mu \quad (\Rightarrow) \quad \alpha \leq \frac{\mu}{M B_G}$$

$$\mu - \frac{\alpha M B_G}{2} \geq \frac{\mu}{2}$$

$$\mathbb{E}_{\xi_t} \{ \Gamma(f_{t+1}) - \Gamma(f_t) \} \leq \frac{\alpha^2 M B}{2} - \frac{\alpha \mu}{2} \|\nabla \Gamma(f_t)\|^2$$

By strong convexity, $\forall f \in \mathcal{F}$

$$\Gamma(f) \geq \Gamma(f_t) + \langle \nabla \Gamma(f_t), f - f_t \rangle + \frac{m}{2} \|f - f_t\|^2$$

$$\geq \Gamma(f_t) + \min_{g \in \mathcal{F}} \{ \langle \nabla \Gamma(f_t), g \rangle + \frac{m}{2} \|g\|^2 \}$$

$$= \Gamma(f_t) - \frac{1}{2m} \|\nabla \Gamma(f_t)\|^2$$

$$\Rightarrow -\|\nabla \Gamma(f_t)\|^2 \leq 2m (\Gamma(f) - \Gamma(f_t)), \forall f$$

$$-\|\nabla \Gamma(f_t)\|^2 \leq 2m (\Gamma^* - \Gamma(f_t))$$

$$\mathbb{E}_{\xi_t} \{ \Gamma(f_{t+1}) - \Gamma(f_t) \}$$

$$\leq \frac{\alpha^2 M B}{2} - \alpha \mu m (\Gamma(f_t) - \Gamma^*)$$

$$\Delta_t := \mathbb{E} \{ \Gamma(f_t) \} - \Gamma^*$$

$$\Delta_{t+1} - \Delta_t \leq \frac{\alpha^2 M B}{2} - \alpha \mu m \Delta_t$$

$$\Delta_{t+1} \leq \frac{\alpha^2 MB}{2} + (1 - \alpha \mu m) \Delta_t \quad \Delta_1 = \Gamma(\xi_1) - \gamma^*$$

Need $0 < 1 - \alpha \mu m < 1$

$$\alpha \leq \frac{\mu}{MB_G}$$

$$2\mu m \leq \frac{m}{M} \left(\frac{\mu^2}{BG_G} \right) \leq \frac{m}{M} < 1$$

≤ 1

$$\Delta_{t+1} \leq (1 - \alpha \mu m) \Delta_t + \frac{\alpha^2 MB}{2}$$

$$\Rightarrow \Delta_t \leq \underbrace{(1 - \alpha \mu m)^{t-1}}_{\rightarrow 0 \text{ as } t \rightarrow \infty} \Delta_1 + \frac{\alpha^2 MB}{2}$$

$$\frac{\alpha^2 MB}{2} = \frac{\alpha MB}{2} \cdot \alpha \leq \frac{\alpha MB}{2} \cdot \frac{1}{\mu m}$$

$$= \frac{\alpha B}{2\mu} \left(\frac{M}{m} \right)$$

if $\Gamma(\cdot)$ is C^2 , $\frac{M}{m} \sim \text{cond. \# of } \nabla^2 \Gamma$

Goal: $\mathbb{E} \Delta_t \leq \varepsilon$

1) choose α s.t. $\frac{\alpha B}{2\mu} \cdot \frac{M}{m} \leq \frac{\varepsilon}{2}$

$$\alpha \leq \frac{\mu m \varepsilon}{MB}$$

2) choose t s.t. $(1 - \alpha \mu m)^{t-1} \leq \frac{\varepsilon}{2}$

- Diminishing stepsize

$\alpha_t \rightarrow 0$ as $t \rightarrow \infty$

Robbins Monro:
$$\left. \begin{aligned} \sum_{t=1}^{\infty} \alpha_t &= \infty \\ \sum_{t=1}^{\infty} \alpha_t^2 &< \infty \end{aligned} \right\} \begin{aligned} \text{e.g. } \alpha_t &= \frac{c}{\gamma+t} \\ c, \gamma &> 0 \end{aligned}$$

Take $\alpha_t = \frac{c}{\gamma+t}$

$$\alpha_1 > \alpha_2 > \alpha_3 > \dots$$

Choose c, γ carefully:

1) $\alpha_1 \leq \frac{\mu}{MB_G} \Rightarrow \alpha_t \leq \frac{\mu}{MB_G}$

$$\begin{aligned} \Delta_{t+1} &\leq (1 - \alpha_t m \mu) \Delta_t + \frac{\alpha_t^2 MB}{2} \\ &= \left(1 - m \mu \frac{c}{\gamma+t}\right) \Delta_t + \frac{c^2 MB}{2(\gamma+t)^2} \end{aligned}$$

want: $\Delta_t \leq \frac{\nu}{\gamma+t}$ (ν to be tuned)

$$\begin{aligned} \Delta_{t+1} &\leq \left(1 - \frac{m \mu c}{\gamma+t}\right) \frac{\nu}{\gamma+t} + \frac{c^2 MB}{2(\gamma+t)^2} \\ &= \left(\frac{\gamma+t - m \mu c}{(\gamma+t)^2}\right) \nu + \frac{c^2 MB}{2(\gamma+t)^2} \end{aligned}$$

$$= \left(\frac{\bar{t} - m \mu c}{\bar{t}^2}\right) \nu + \frac{c^2 MB}{2\bar{t}^2} \quad \bar{t} := \gamma+t$$

$$= \frac{\nu(\bar{t}-1)}{\bar{t}^2} + \frac{1}{\bar{t}^2} \left[\nu(1 - m \mu c) + \frac{c^2 MB}{2} \right]$$

Choose $\frac{c}{\gamma}$ s.t. $\rightarrow \leq 0$

$$\left. \begin{aligned} \frac{\bar{t}-1}{\bar{t}^2} &\leq \frac{1}{\bar{t}+1} \\ \frac{(\bar{t}-1)(\bar{t}+1)}{\bar{t}^2} &\leq 1 \end{aligned} \right\} c, \gamma, \nu \text{ tuned}$$

Const. step.: $\Delta_t \leq (1 - \alpha m \mu)^{t-1} \Delta_1 + K \alpha$
 $\log(\frac{1}{\epsilon})$ comp. to give ϵ -opt.

Diminishing
step:

$$\Delta_t \leq \frac{\nu}{\gamma+t}$$

$\frac{1}{\epsilon}$ - comp. to give ϵ -opt.

$$\frac{\nu}{\gamma+t} \leq \epsilon$$

$$\gamma+t \geq \frac{\nu}{\epsilon}$$

$$t \geq \frac{\nu}{\epsilon} - \gamma$$