

Stability of Stochastic Gradient Descent

- Learning algo $A: \mathcal{Z}^* \rightarrow \mathcal{F}$
 $\mathcal{Z}^* = \bigcup_{n \geq 1} \mathcal{Z}^n$

Properties:

- consistency:

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left\{ \mathbb{E} [L_P(A(z^n))] - L_P^* \right\} = 0$$

where $L_P(f) = \mathbb{E} [l(f, z)] \quad z \sim P$
 $L_P^* = \inf_{f \in \mathcal{F}} L_P(f)$

- AERM:

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \mathbb{E} [L_n(A(z^n)) - L_n^*] = 0$$

where $L_n(f) := \frac{1}{n} \sum_{i=1}^n l(f, z_i)$
 $L_n^* := \inf_{f \in \mathcal{F}} L_n(f)$

- stability:

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[l(A(z^n), z_i') - l(A(z_{(-i)}^n), z_i') \right] = 0$$

where $z^n = (z_1, \dots, z_n)$, $z^n = (z_1', \dots, z_n')$ iid $\sim P$

$$z_{(-i)}^n = (z_1, \dots, z_{i-1}, z_i', z_{i+1}, \dots, z_n)$$

$$z^n \xrightarrow{A} A(z^n), \quad z_{(-i)}^n \xrightarrow{A} A(z_{(-i)}^n)$$

NB: stability \Leftrightarrow generalization (on avg)

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \left| \mathbb{E} [L_n(A(z^n)) - L(A(z^n))] \right| = 0$$

AERM + stability \Rightarrow consistency

$$\text{Gen.} : \lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} \mathbb{E}_P \left| L_n(A(z^n)) - L(A(z^n)) \right| = 0$$

e.g. $0 \leq l \leq 1$

AERM + gen. on avg. \Rightarrow gen.

Stab. on avg. \Leftarrow uniform stability:

$$\sup_{z \in \mathcal{Z}} \left| \mathbb{E} [l(A(z^n), z)] - \mathbb{E} [l(A(z_{i_1}^n), z)] \right| \rightarrow 0$$

E.g. $f \mapsto l(f, z)$ is Lipschitz, unif. in z

$$\sup_{z \in \mathcal{Z}} |l(f, z) - l(f', z)| \leq L \|f - f'\|$$

so if $A(z^n) \approx A(z_{i_1}^n)$ $\forall i$

$$\text{then } \sup_{z \in \mathcal{Z}} |l(A(z^n), z) - l(A(z_{i_1}^n), z)| \leq L \|A(z^n) - A(z_{i_1}^n)\|.$$

Stochastic Gradient Descent

Goal: prove uniform stability of SGD
(Hardt - Recht - Singer 2016)

ERM:

$$\min_{f \in \mathcal{F}} L_n(f)$$

\mathcal{F} : closed, convex
subset of H.S. \mathcal{H}

$$L_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(f, z_i) \quad z_1, \dots, z_n \stackrel{i.i.d.}{\sim} \mathcal{P}$$

Approximate ERM: gradient descent

- assume $f \mapsto \ell(f, z)$ is diff. ($\forall z$)

$\nabla \ell(f, z)$: gradient

GD: $f_0 \in \mathcal{F}$ (init)

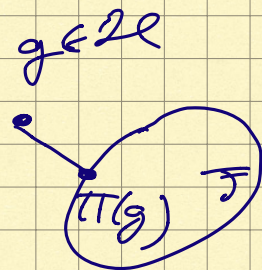
$$f_t = \Pi \left(f_{t-1} - \alpha_t \nabla L_n(f_{t-1}) \right) \quad t=1, 2, \dots$$

where $\nabla L_n(f) := \frac{1}{n} \sum_{i=1}^n \nabla \ell(f, z_i)$

$(\alpha_t)_{t \geq 1}^\infty$: positive step sizes

$\Pi: \mathcal{H} \rightarrow \mathcal{F}$: projection onto \mathcal{F}

$$\Pi(g) := \operatorname{argmin}_{f \in \mathcal{F}} \|g - f\|^2$$



Opt. theory:

$$\lim_{t \rightarrow \infty} L_n(f_t) \rightarrow L_n^*$$

under various regularity assumptions
(including carefully tuned step sizes)

Drawback: computation of f_t takes $\mathcal{O}(n)$ steps
(compute $\nabla L_n(f_{t-1})$)

Compromise: stochastic approx. (SA)

$$\nabla L_n(f) = \frac{1}{n} \sum_{i=1}^n \nabla l(f, z_i)$$

$$I \sim \text{Unif}([n])$$

$$\nabla L_n(f) = \mathbb{E}_I [\nabla l(f, z_I)]$$

$\Rightarrow \nabla l(f, z_I)$ is an unbiased est.
of $\nabla L_n(f)$
(only I is random!)

GD:

$$f_{t-1} - \alpha_t \nabla L_n(f_{t-1})$$

SGD:

$$f_{t-1} - \alpha_t \nabla l(f_{t-1}, z_I)$$

$$f_{t-1} - \alpha_t \nabla l(f_{t-1}, z_I)$$

$$= f_{t-1} - \alpha_t [\nabla L_n(f_{t-1}) + \xi_t]$$

where $\xi_t := \nabla l(f_{t-1}, z_I) - \nabla L_n(f_{t-1})$

$$\mathbb{E}_I [\xi_t] = 0$$

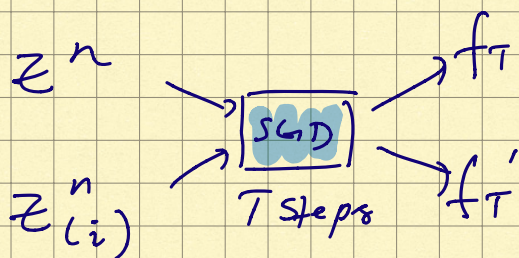
Details matter; this is just heuristic.

Gain: each iteration now takes $\mathcal{O}(1)$ steps.

Let $(I_t)_{t \geq 1}$ be iid from $\text{Unif}([n])$
 (other choices are possible)

SGD: $f_0 \in \mathcal{F}$ (init, indep of z^n)
 $t=1, 2, \dots, T$:
 $f_t = \Pi (f_{t-1} - \alpha_t \nabla l(f_{t-1}, z_{I_t}))$

$A(z^n) = f_T$ - stability?



Notation: $G_{\varphi, \alpha}(f) := \Pi(f - \alpha \nabla \varphi(f))$
 $\varphi: \mathcal{F} \rightarrow \mathbb{R}$ differentiable
 $\alpha > 0$
 Π : proj. onto \mathcal{F}

$(f_t)_{t=0}^T$
 (on z^n)

$(f_t^i)_{t=0}^T$
 (on $z_{(i)}^n$)

$f_0 = f_0^i$ (same init)

$f_t = G_t(f_{t-1})$

$f_t^i = G_t^i(f_{t-1}^i)$

where $G_t = G_{l(\cdot, z_{I_t}), \alpha_t}$, G_t^i analogous

$$z^n = (z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n)$$

$$z_{(i)}^n = (z_1, \dots, z_{i-1}, z_i', z_{i+1}, \dots, z_n)$$

$$\text{Let } \delta_k := \|f_k - f_k'\| \quad \delta_0 = 0$$

$$\delta_k = \|f_k - f_k'\|$$

$$= \|G_k(f_{k-1}) - G_k'(f_{k-1}')\|$$

$$\text{two cases: } \begin{cases} I_k \neq i & (G_k = G_k') \quad \textcircled{1} \\ I_k = i & (G_k \neq G_k') \quad \textcircled{2} \end{cases}$$

$$\textcircled{1} \quad \|G_k(f_{k-1}) - G_k'(f_{k-1}')\|$$

$$= \|G_k(f_{k-1}) - G_k(f_{k-1}')\|$$

$$\leq \eta_k \|f_{k-1} - f_{k-1}'\|$$

$$= \eta_k \delta_{k-1}$$

$$\text{where } \eta_k := \sup_{f, f' \in \mathcal{F}} \frac{\|G_k(f) - G_k(f')\|}{\|f - f'\|}$$

$$(\text{want } \eta_k \leq 1)$$

$$\textcircled{2} \quad \|G_k(f_{k-1}) - G_k'(f_{k-1}')\| \quad (G_k \neq G_k')$$

$$\leq \underbrace{\|G_k(f_{k-1}) - G_k(f_{k-1}')\|}_{\leq \eta_k \delta_{k-1}} + \underbrace{\|G_k(f_{k-1}') - G_k'(f_{k-1}')\|}_{\textcircled{1}}$$

$$\begin{aligned}
 (D) &\leq \|G_t(f_{t-1}') - f_{t-1}'\| + \|G_t'(f_{t-1}') - f_{t-1}'\| \\
 &\leq 2c_t
 \end{aligned}$$

where $c_t := \sup_{f \in \mathcal{F}} \max \{ \|G_t(f) - f\|, \|G_t'(f) - f\| \}$

$$\therefore \delta_t \leq \eta_t \delta_{t-1} + 2c_t \mathbb{1}_{\{I_t = \mathcal{Z}\}}$$

Assume $\eta_t \leq 1 \quad \forall t$ (guaranteed by choice of α_t)

Assumptions:

on $f \mapsto \ell(f, z)$

1) convex

2) L -Lipschitz: $\sup_{z \in \mathcal{Z}} |\ell(f, z) - \ell(f', z)| \leq L \|f - f'\|$

3) M -smooth: $\sup_{z \in \mathcal{Z}} \|\nabla \ell(f, z) - \nabla \ell(f', z)\| \leq M \|f - f'\|$

4) m -strongly convex: ($m \geq 0$)

$$f \mapsto \ell(f, z) - \frac{m}{2} \|f\|^2 \quad \text{convex}$$

Thm Assume 1) - 3); then, for $\alpha_t \leq 2/M$,

$$\mathbb{E}[\delta_T] \leq \frac{2L^2}{\eta} \sum_{t=1}^T \alpha_t \quad (\text{expect. w.r.t. } (I_t))$$

Proof (sketch)

$$1) - 3) : \quad \eta_t \leq 1$$

$$c_t \leq \alpha_t L$$

$$\begin{aligned} \|G_{\varphi, \alpha}(f) - f\| &= \|\Pi(f - \alpha \nabla \varphi(f)) - \Pi(f)\| \\ &\leq \alpha \|\nabla \varphi(f)\| \end{aligned}$$

$$\delta_t \leq \delta_{t-1} + 2\alpha_t L \mathbb{1}_{\{I_t = \bar{i}\}} \quad \delta_0 = 0$$

$$\delta_T \leq 2L \sum_{t=1}^T \alpha_t \mathbb{1}_{\{I_t = \bar{i}\}}$$

$$\mathbb{E}[\delta_T] \leq 2L \sum_{t=1}^T \alpha_t \mathbb{P}[I_t = \bar{i}]$$

$$= \frac{2L}{n} \sum_{t=1}^T \alpha_t. \quad \square$$

Corollary $\sup_{z \in \mathcal{Z}} \mathbb{E}[\ell(f_T, z) - \ell(f_T', z)] \leq \frac{2L^2}{n} \sum_{t=1}^T \alpha_t.$

E.g. $T = n$

$$\alpha_t = \frac{2}{M\sqrt{n}}$$

$$t = 1, \dots, T$$

$$\sum_{t=1}^T \alpha_t = \frac{2\sqrt{n}}{M}$$

$$\Rightarrow \mathbb{E}[\ell(f_T, z) - \ell(f_T', z)] \leq \frac{4L^2}{M\sqrt{n}}.$$

Thm Assume 1) - 4). If $\alpha_t \leq \frac{1}{M}$, then

$$E[S_T] \leq \frac{4L}{m\eta} \quad \checkmark T$$

Key idea: $\alpha_t = 1/M \Rightarrow \eta_t \leq 1 - \frac{\alpha m}{2}$
 $= 1 - \frac{m}{2M}$

$$s_t \leq \underbrace{\left(1 - \frac{m}{2M}\right)}_{< 1} s_{t-1} + \frac{2}{M} \mathbb{1}_{\{I_t = i\}}$$