# Ch 14 Stochastic Gradient Descent

Friday, March 20, 2020     3:28 PM

SGD = Stochastic Gradient Descent

           (misnomer)

$\min\limits_{w} f(w)$,     $f(w) = L_D(w) \equiv \mathbb{E}\,\ell(w, z)$     SA = Stochastic Approximation

           or $\hat{L}_S(w) = \frac{1}{m}\sum\limits_{i=1}^{m}\ell(w, z_i)$   SAA = Sample Avg. Approximation
                                              = ERM

     Algo, some analysis, covers both cases
         (ie SAA is a special case of SA w/ $D = \text{Uniform}(\{z_i\}_{i=1}^{m})$)

See supplemental notes for basics of <u>gradient descent</u>
     $w^{(t+1)} = w^{(t)} - \eta\, \nabla f(w^{(t)})$
                                   $\eta$ is learning rate / stepsize
take $T$ steps, output   1) $w^{(T)}$
                    2) $\arg\min\limits_{w \in \{w^{(1)}, \dots, w^{(T)}\}} f(w)$   ← can't always evaluate $f$
               or
                    3) $\bar{w} = \frac{1}{T}\sum\limits_{t=1}^{T} w^{(t)}$

**Corollary 14.2** Analysis of GD if $f$ convex, Lipschitz, but <u>not</u> smooth
     Let $f$ be <u>convex</u>, $\rho$-<u>Lipschitz</u>, $w^* \in \arg\min\limits_{w} f(w)$, $\|w^*\| \le B$,
     iterate $w^{(t+1)} = w^{(t)} - \eta\, d_t$, $d_t \in \partial f(w^{(t)})$,
     then choosing $\eta = \sqrt{\dfrac{B^2}{\rho^2 T}}$ gives $f(\bar{w}) - f(w^*) \le \dfrac{B\rho}{\sqrt{T}}$ ← super slow!

## §14.3 SGD

     SGD algo: for $t = 1, 2, \dots, T$
                    Draw r.v. $V_t$ s.t. $\underline{\mathbb{E}(V_t \mid w^{(t)}) \in \partial f(w^{(t)})}$
                    $w^{(t+1)} = w^{(t)} - \eta \cdot V_t$
              output $\bar{w} = \frac{1}{T}\sum\limits_{t=1}^{T} w^{(t)}$   other possibilities too (Polyak-Ruppert averaging, various weights)

What might we want to show?
1. First, pick error metric — or $\bar{w}$
     $e_t = \begin{cases} f(w^{(t)}) - f(w^*) & \text{standard choice if } f \text{ convex} \quad (\text{ok}) \\ \|w^{(t)} - w^*\| & \text{choice if } f \text{ strongly convex} \quad (\text{best}) \\ \|\nabla f(w^{(t)})\| & \text{choice if } f \text{ not convex} \quad (\text{weak}) \end{cases}$

2. $V_t$, hence $w^{(t)}$, hence $e_t$, is a random variable
     A. $e_t \xrightarrow{P} 0$ ← or anything   convergence in probability (measure)

         means $\forall \varepsilon > 0$, $\lim\limits_{t \to \infty} \mathbb{P}(|e_t| > \varepsilon) = 0$      (weak)

     B. $e_t \xrightarrow{L^p} 0$ if $\mathbb{E}|e_t|^p = 0$, $L^p$ convergence    (p=2 aka quadratic mean)
                                                        (ok, p=1)

c. $e_t \xrightarrow{a.s.} 0$ if $\mathbb{P}(\lim e_t = 0) = 1$, almost sure convergence aka w/ probability 1 (best)

other types as well (in distribution ...) see prob. textbook

Ex $e_t = \begin{cases} 1 & \text{w.p. } 1/t \\ 0 & \text{w.p. } 1-1/t \end{cases}$ then $\forall r \geq 1$, $\mathbb{E} e_t = 1/t$

so $e_t \xrightarrow{L^r} 0$ ($\forall r$) and $e_t \xrightarrow{\mathbb{P}} 0$

but (I don't think) $e_t \xrightarrow{a.s.} 0$

$e_t = \begin{cases} t^\alpha & \text{wp } 1/t \\ 0 & \text{wp } 1-1/t \end{cases}$

$\alpha = 1/2$ $\mathbb{E} e_t = 1/\sqrt{t} \to 0$ so $e_t \xrightarrow{L^1} 0$

$\mathbb{E} e_t^2 = 1$ so $e_t \not\xrightarrow{L^2} 0$

(Fact: $e_t \xrightarrow{L^r} 0 \Rightarrow e_t \xrightarrow{\mathbb{P}} 0$)

$\alpha = 1$ $\mathbb{E} e_t = 1$ So doesn't converge in $L^1$ even

but $\forall \varepsilon > 0$, $\mathbb{P}(|e_t| > \varepsilon) \leq 1/t \to 0$

so converges in probability

What type to use?

Most ML results show $L^1$ convergence,

$\mathbb{E}(|e_t|^1) \to 0$

(or, since usually $e_t \geq 0$, $\mathbb{E} e_t \to 0$ )

However, if convergence is fast enough,

or for simplest cases (original Robbins-Monro '50's)

prove $L^2$ or almost sure

Thm 14.8 [SSS] $L^1$ convergence of SGD assuming...

Let $f$ be convex, $w^*$ a minimizer, $\|w^*\| \leq B$,

$\|v_t\| \leq \rho$ $\forall t \in [T]$ (w.p. 1) (like $\rho$-Lipschitz), then

$$0 \leq \mathbb{E} f(\bar{w}) - f(w^*) \leq \frac{B\rho}{\sqrt{T}},$$

i.e. for $\varepsilon$ error, $T \geq \frac{B^2 \rho^2}{\varepsilon^2}$

proof:

$(v_t)$ is a stochastic process

$F_T := \sigma(v_t : t \leq T)$, then $(F_T)_{T \in \mathbb{N}}$ is a "filtration"

background:

used to help w/ conditional probabilities

write $\mathbb{E}(v_t \mid \{v_{t-1}, v_{t-2}, \ldots, v_0\})$ as $\mathbb{E}(v_t \mid F_t)$

and use "law of total expectation" aka "tower property"

$\mathbb{E}(\mathbb{E}(X \mid F)) = \mathbb{E}(X)$

i.e. simpler notation, $\mathbb{E}_\alpha g(\alpha) = \mathbb{E}_\beta (\mathbb{E}_\alpha [g(\alpha) \mid \beta])$

see, e.g.,

https://ocw.mit.edu/courses/sloan-school-of-management/15-070j-advanced-stochastic-processes-fall-2013/lecture-notes/MIT15_070JF13_Lec9.pdf

By default, write $\mathbb{E}$ to mean $\mathbb{E}[\,\cdot \mid F_T]$

Define $\bar{w} = \frac{1}{T}\sum_t w^{(t)}$, $f(\bar{w}) \leq \frac{1}{T}\sum f(w^{(t)})$ by Jensen,

$$f(\bar{w}) - f^* \leq \frac{1}{T}\sum_{t=1}^{T} f(w^{(t)}) - f^*$$

$$\mathbb{E}\left[\text{``}\underline{\quad}\text{''}\right] \leq \mathbb{E}\left[\text{``}\underline{\quad\quad}\text{''}\right]$$

then via deterministic bounds (lemma 14.1)

if $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$, $\|w^*\| \leq B$, $\|v_t\| \leq \rho$,

$$\mathbb{E}\left[\frac{1}{T}\sum \langle w^{(t)} - w^*, v_t\rangle\right] \leq \frac{B\rho}{\sqrt{T}}$$

Now claim

$$\mathbb{E}\left[\frac{1}{T}\sum f(w^{(t)}) - f^*\right] \leq \underbrace{\qquad}\quad \text{to conclude proof.}$$

$$\mathbb{E}\left[\frac{1}{T}\sum \langle w^{(t)} - w^*, v_t\rangle\right] = \frac{1}{T}\sum \mathbb{E}\langle w^{(t)} - w^*, v_t\rangle \qquad \textcolor{red}{\text{recall } w^{(t)} = w^{(t-1)} - \eta\, v_{t-1}}$$

$$= \frac{1}{T}\sum \mathbb{E}\left(\mathbb{E}\left[\langle w^{(t)} - w^*, v_t\rangle \mid \mathcal{F}_{t-1}\right]\right)$$

$$\textcolor{red}{\underset{\uparrow}{\quad}\text{not random for now}}$$

and $\underbrace{\mathbb{E}[v_t \mid \mathcal{F}_{t-1}]}_{g_t} \in \partial f(w^{(t)})$

so by convexity $f(w^{(t)}) - f^* \leq \langle w^{(t)} - w^*, g_t\rangle$

$$\geq \frac{1}{T}\sum \mathbb{E}\, f(w^{(t)}) - f^*$$

$$= \mathbb{E}\left(\frac{1}{T}\sum f(w^{(t)}) - f^*\right) \quad \textcolor{green}{\square}$$

## §14.5 Learning w/ SGD

ie. let $f(w) = L_D(w) := \underset{z \sim D}{\mathbb{E}}[l(w, z)]$

Can't compute $f(w)$ or $\nabla f(w)$ since we don't know $L_D$!

... but we can draw from $D$ and use SGD.

ie., sample $z_t \sim D$, let $v_t \in \partial l(w^{(t)}, v_t)$

$\underset{\uparrow \text{ w.r.t } w}{}$

So immediate corollary

## Corollary 14.2

$l$ is $\rho$-Lipschitz, $\|w^*\| \leq B$, then

$\forall \epsilon > 0$, running SGD w/ $T \geq \frac{B^2 \rho^2}{\epsilon^2}$ iterations,

w/ $\eta = \sqrt{\frac{B^2}{\rho^2 T}}$ then

$$\mathbb{E}\, L_D(\bar{w}) \leq \min_{w \in H} L_D(w) + \epsilon$$

$\underset{\nearrow}{}$ expected risk, like we discussed in **Stability** chapter

$\textcolor{red}{\uparrow \text{we didn't discuss constraints, but many simple ones (and regularizers) easily fit into SGD/GD}}$

$T$ is like $m$, $= \#$ iid samples

(if someone says "epochs", they are in the SAA/ERM setting, and $T = 4m$ is "4 epochs". In the true SA setting, like above corollary, we never reach a single epoch, ie., $m = \infty$)

Results also hold if $f$ is $\beta$-smooth, and for RLM

For SAA work, a good review is

For SA in optimization Context,
    See Nemirovski "Robust SA"
    or Nesterov "Primal-dual Avg"