▷ <u>*Assumptions:*</u>

(1) $|u^{T} Hf(x) u| \leq L \|u\|^2 \quad \forall x$ (bound on the second deriv.)

(2) $\underbrace{\mathbb{E}[\|\nabla f_{i(k)}(x) - \mathbb{E}[\nabla f_{i(k)}(x)]\|^2]}_{=} \leq \sigma^2 \quad \forall x$

$$Var[\nabla f_{i(k)}(x)]$$

$\iff$ is about $Var[\nabla f_{i(k)}(x)] \leq \sigma^2$

$\iff \mathbb{E}[\|\nabla f_{i(k)}(x) - \nabla f(x)\|^2] \leq \sigma^2$

$$\mathbb{E}[\|\nabla f_{i(k)}(x) - \nabla f(x)\|^2]$$

Expectation
is linear $\rightsquigarrow = \mathbb{E}[\|\nabla f_{i(k)}(x)\|^2] - 2\mathbb{E}[\langle \nabla f_{i(k)}(x), \nabla f(x)\rangle] + \mathbb{E}[\|\nabla f(x)\|^2]$

$$= \mathbb{E}[\|\nabla f_{i(k)}(x)\|^2] - 2\langle \underbrace{\mathbb{E}[\nabla f_{i(k)}(x)]}_{\nabla f(x^{(k)})}, \nabla f(x)\rangle$$

$$+ \underbrace{\mathbb{E}[\|\nabla f(x)\|^2]}_{\|\nabla f(x)\|^2} \qquad \underbrace{}_{-2\|\nabla f(x)\|^2}$$

$$= \mathbb{E}[\|\nabla f_{i(k)}(x)\|^2] - \|\nabla f(x)\|^2 \leq \sigma^2$$

In conclusion, the bound on variance gives

$$\mathbb{E}[\|\nabla f_{i(k)}(x)\|^2] \leq \sigma^2 + \|\nabla f(x)\|^2$$

By Taylor, $\exists \xi \in \mathbb{R}^n$ s.t.

$$f(x^{(k+1)}) \overset{\uparrow}{=} f(x^{(k)} - a^{(k)} \nabla f_{i(k)}(x^{(k)}))$$

$$= f(x^{(k)}) - a^{(k)} \langle \nabla f_{i(k)}(x^{(k)}), \nabla f(x^{(k)})\rangle$$

$$+ \frac{1}{2}(a^{(k)} \nabla f_{i(k)}(x^{(k)}))^T Hf(\xi)(a^{(k)} \nabla f_{i(k)}(x^{(k)}))$$

Ass.(1)
$\Rightarrow f(x^{(k+1)}) \leq f(x^{(k)}) - a^{(k)}\langle \nabla f_{i(k)}(x^{(k)}), \nabla f(x^{(k)})\rangle$

$$+ \frac{L}{2}(a^{(k)})^2 \|\nabla f_{i(k)}(x^{(k)})\|^2$$

$$\mathbb{E}[f(x^{(k+1)})] \leq \mathbb{E}[f(x^{(k)})] - a^{(k)}\|\nabla f(x^{(k)})\|^2$$

$$+ \frac{(a^{(k)})^2 L}{2}(\sigma^2 + \|\nabla f(x^{(k)})\|^2)$$

$$= \mathbb{E}[f(x^{(k)})] - \frac{a^{(k)}}{2}\underbrace{(2 - a^{(k)} L)}_{\geq 1}\|\nabla f(x^{(k)})\|^2$$

$$+ \frac{(a^{(k)})^2 L \cdot \sigma^2}{2} \qquad \text{if } a^{(k)} \text{ is sufficiently}$$
$$\text{small } (a^{(k)} \leq \frac{1}{L}) \longrightarrow$$

$$\leq \mathcal{E}\left[f(x^{(k)})\right] - \frac{a^{(k)}}{2}\left\|\nabla f(x^{(k)})\right\|^2 + \frac{(a^{(k)})^2 L \cdot \sigma^2}{2}$$

$$\Rightarrow a^{(k)}\left\|\nabla f(x^{(k)})\right\|^2 \leq \left(\mathcal{E}\left[f(x^{(k)})\right] - \mathcal{E}\left[f(x^{(k+1)})\right]\right) + (a^{(k)})^2 L \cdot \sigma^2$$

$$\Rightarrow \sum_{k=0}^{T-1} a^{(k)}\left\|\nabla f(x^{(k)})\right\|^2 = 2\left(\underbrace{\mathcal{E}\left[f(x^{(0)})\right]}_{f(x^{(0)})} - \underbrace{\mathcal{E}\left[f(x^{(T)})\right]}_{}\right) + L\sigma^2 \sum (a^{(k)})^2$$

$\underset{\uparrow}{\text{summing up}}$
order $T-1$
stops

$f(x^{(T)}) \geq f^*$ is the local minimum
$\Rightarrow \mathcal{E}\left[f(x^{(T)})\right] \geq f^*$

Hence:

$$\sum_{k=0}^{T-1} a^{(k)}\left\|\nabla f(x^{(k)})\right\|^2 \leq 2\left(f(x^{(0)}) - f^*\right) + L\sigma \sum_{k=0}^{T-1} (a^{(k)})^2$$

Now assume that we do random number $t$ of steps between
$0 \dots T-1$ with probability

$$P(c=t) = \frac{a^{(t)}}{\sum_{i=0}^{T-1} a^{(k)}}$$

| 0 | 1 | 2 | 3 | | T-1 |
|---|---|---|---|---|---|
| $\frac{a^{(0)}}{\sum a^{(i)}}$ | $\frac{a^{(1)}}{\sum a^{(i)}}$ | $\frac{a^{(i)}}{\sum a^{(i)}}$ | | | $\frac{a^{(T-1)}}{\sum a^{(i)}}$ |

$$\mathcal{E}\left[\left\|\nabla f(x^{(T)})\right\|^2\right] \overset{\underset{def}{}}{=} \sum_{i=0}^{T-1}\left\|\nabla f(x^{(i)})\right\|^2 \cdot P(c=i)$$

$\underset{\uparrow}{\text{new randomness}}$

$$= \sum_{i=0}^{T-1} \frac{a^{(i)}\left\|\nabla f(x^{(i)})\right\|^2}{\sum_{i=0}^{T-1} a^{(i)}} = \left(\sum_{i=0}^{T-1} a^{(i)}\right)^{-1} \sum_{i=0}^{T-1} a^{(i)}\left\|\nabla f(x^{(i)})\right\|^2$$

$$\leq \frac{2(f(x^{(0)}) - f^*)}{\sum_{i=0}^{T-1} a^{(i)}} + L\sigma^2 \frac{\sum_{i=0}^{T-1} (a^{(i)})^2}{\sum_{i=0}^{T-1} a^{(i)}}$$

if $a^{(i)} = a \; \forall i$, then $\mathcal{E}\left[\left\|\nabla f(x^{(T)})\right\|^2\right] \leq \underset{\underset{T\to\infty}{\downarrow}}{\frac{2(f(x^{(0)}) - f^*)}{aT}} \underset{0}{} + \underset{\underset{0}{\downarrow}}{\underbrace{L\sigma^2 a}_{\text{Noise ball}}}$

How to assure convergence (in expectation)

- $\sum_{i=0}^{T-1} a^{(i)} \xrightarrow[T\to\infty]{} \infty$

- $\sum_{i=0}^{T-1} a^{(i)}$ goes to $\infty$ faster than $\sum_{i=0}^{T-1} (a^{(i)})^2$

$\rightarrow$

for example: $a^{(k)} = \frac{1}{\sqrt{k+1}}$ does the job!

$$\sum_{k=0}^{T-1} \frac{1}{\sqrt{k+1}} \sim \int_{0}^{T} \frac{1}{\sqrt{x}} dx = 2\sqrt{T}$$

$$\sum_{k=0}^{T-1} \frac{1}{k+1} \sim \int_{0}^{T} \frac{1}{x} dx = \log T$$

$\sqrt{T}$ grows much faster than $\log T$.

Remark: If $f$ is strictly convex, then SGD performs much better: it converges (in expectation) even with constraint learning rate.