$$E f(x^{k+1}) \leq E f(x^k) - \tau \left( 1 - \frac{\tau \bar{L}}{2} \right) E |\nabla f(x^k)|^2 + \frac{\tau^2 \bar{L}}{2} E \left( \frac{1}{m} \sum_{i=1}^{m} |\nabla f_i(x^k) - \nabla f(x^k)|^2 \right)$$

Remember

$$f = \frac{1}{m} \sum_{i=1}^{m} f_i \qquad\qquad m \sim 10^5 - 10^{10} \qquad (\text{Imagenet : dataset of images} \sim 10^6)$$

↑ loss function evaluated on a point of the training set $(D = \{(a_1, y_1) \cdots (a_N, y_N)\}$.

All these functions are functions of the "weights" of a NN

We cannot directly use GD (= Gradient Descent) $\qquad x^{k+1} = x^k - \tau \nabla f$
Instead we use SGD

$$x^{k+1} = x^k - \tau \nabla f_{i_k}(x^k) \qquad\qquad i_k \text{ is sampled with prob. } 1/m$$
$$\text{(i.e. uniform probability)}$$
$$\text{from} \quad \{1, \ldots, m\}.$$

"I pick an example at random and on
that example I compute the gradient"

- We already understood that we have to impose a bound on the variance of $\nabla f$

- $$\boxed{\frac{1}{m} \sum_{i=1}^{m} |\nabla f_i(x) - \nabla f(x)|^2 \leq \sigma^2}$$

- $$E f(x^{k+1}) \leq E f(x^k) - \tau \left( 1 - \frac{\tau \bar{L}}{2} \right) E |\nabla f(x^k)|^2 + \frac{\tau^2 \bar{L}}{2} E \left( \frac{1}{m} \sum_{i=1}^{m} |\nabla f_i(x^k) - \nabla f(x^k)|^2 \right)$$

$$\sum_{i=1}^{+\infty} |\nabla f|^2 < +\infty$$

Bounded
$\leq \sigma^2$

1. In order to control the variance term our only hope is $\tau^2$.

2. The variance term is proportional to $\tau^2$ while the gradient term $\left( E |\nabla f(x^k)|^2 \right)$
   is proportional to $\tau$

First of all we understand that we need a variable learning rate $\tau \to \tau_k$

$$E f(x^{k+1}) \leq E f(x^k) - \tau_k \left( 1 - \frac{\tau_k \bar{L}}{2} \right) E |\nabla f(x^k)|^2 + \frac{\tau_k^2 \bar{L}}{2} E \left( \frac{1}{m} \sum_{i=1}^{m} |\nabla f_i(x^k) - \nabla f(x^k)|^2 \right)$$

We recall the bound of the variance and we also choose $\underline{\underline{\tau_k \bar{L} \le 1}}$

$$- \tau_k \left( 1 - \frac{\tau_k \bar{L}}{2} \right) \le - \frac{\tau_k}{2} \quad \bullet$$

$$- \tau_k + \frac{\tau_k^2 \bar{L}}{2} \le - \frac{\tau_k}{2} \quad , \quad - \frac{\tau_k}{2} + \frac{\tau_k^2 \bar{L}}{2} \le 0 \qquad \tau_k \left( -1 + \tau_k \bar{L} \right) \le 0$$

Now since $\tau_k > 0 \quad \Rightarrow \quad -1 + \tau_k \bar{L} \le 0$

$$\mathbb{E} f(x^{k+1}) \le \mathbb{E} f(x^k) - \frac{\tau_k}{2} \mathbb{E} |\nabla f(x^k)|^2 + \frac{\tau_k^2 \bar{L}}{2} \sigma^2$$

We apply this inequality recursively $n$ times

$$\mathbb{E} f(x^1) \le \mathbb{E} f(x^0) - \frac{\tau_0}{2} \mathbb{E} |\nabla f(x^0)|^2 + \frac{\tau_0^2 \bar{L}}{2} \sigma^2$$

$$\mathbb{E} f(x^2) \le \mathbb{E} f(x^1) - \frac{\tau_1}{2} \mathbb{E} |\nabla f(x^1)|^2 + \frac{\tau_1^2 \bar{L}}{2} \sigma^2$$

$$\le \mathbb{E} f(x^0) + \qquad - \frac{1}{2} \left( \tau_0 \mathbb{E} |\nabla f(x^0)|^2 + \tau_1 \mathbb{E} |\nabla f(x^1)|^2 \right)$$

$$+ \frac{\bar{L} \sigma^2}{2} \quad \tau_0^2 + \tau_1^2$$

$$\mathbb{E} f(x^n) \le \mathbb{E} f(x^0) - \frac{1}{2} \sum_{k=0}^{n-1} \tau_k \mathbb{E} |\nabla f(x^k)|^2 + \frac{\bar{L} \sigma^2}{2} \sum_{k=0}^{n-1} \tau_k^2$$

As we did for GD.

$$- \infty < \underline{\underline{\inf f}} \le \mathbb{E} f(x^n) \le \mathbb{E} f(x^0) - \frac{1}{2} \sum_{k=0}^{n-1} \tau_k \mathbb{E} |\nabla f(x^k)|^2 + \frac{\bar{L} \sigma^2}{2} \sum_{k=0}^{n-1} \tau_k^2$$

$$\frac{1}{2} \sum_{k=0}^{n-1} \tau_k \mathbb{E} |\nabla f(x^k)|^2 \le \mathbb{E} f(x^0) - \inf f + \frac{\bar{L} \sigma^2}{2} \sum_{k=0}^{n-1} \tau_k^2$$

The problem now is that $\displaystyle\sum_{k=0}^{n-1} \tau_k^2 \to +\infty$ as $n \to +\infty$. We need to avoid this

and therefore we assume that

$$\boxed{\sum_{k=0}^{+\infty} \tau_k^2 = T < +\infty}$$

$$\frac{1}{2} \sum_{k=0}^{n-1} \tau_k \mathbb{E} |\nabla f(x^k)|^2 \le \underbrace{\mathbb{E} f(x^0) - \inf f + \frac{\bar{L} \sigma^2}{2} \sigma^2 T}_{> 0} \overset{\vee}{}$$

If we define $S_n = \frac{1}{2} \displaystyle\sum_{k=0}^{n-1} \tau_k \mathbb{E} |\nabla f(x^k)|^2$ Then $S_n \nearrow$ and bounded

from above $\Rightarrow S_n$ converges which means that

$$\sum_{k=0}^{+\infty} \tau_k \mathbb{E} |\nabla f(x^k)|^2 < +\infty$$

If $\boxed{\sum_{k=0}^{+\infty} \tau_k = +\infty}$ then we are safe in the sense that

$$\exists \; i_k \quad \text{s.t.} \qquad \mathbb{E}\,|\nabla f(x^{i_k})|^2 \to 0$$

$$\left( \begin{array}{c} \Downarrow \\ \nabla f \to 0 \quad a.s. \end{array} \right)$$

If you take for instance $\tau_k = \frac{1}{k}$

$$\sum_{k=0}^{+\infty} \frac{1}{k^2} = \left( \frac{\pi^2}{6} \; ? \right)$$

$$\sum_{k=0}^{+\infty} \frac{1}{k} = +\infty$$

What are the assumptions we did to prove this convergence result.

<span style="color:green">**Assumptions on $f$**</span>

1. $f \in \mathcal{C}^1(\mathbb{R}^n; \mathbb{R})$  <span style="color:orange">( $f$ is continuously differentiable )</span>

2. $\nabla f_i$ is $L_i$-lipshitz $\quad \forall \; i = 1, \dots, m$

3. $\inf f > -\infty$  <span style="color:orange">( $f$ is bounded from below )</span>

4. $\frac{1}{m} \sum_{i=1}^{m} |\nabla f_i(x) - \nabla f(x)|^2 < \sigma^2$  <span style="color:orange">( Variance of the gradients bounded )</span>

<span style="color:green">**Assumption on $\tau_k$**</span>

1. $\sum_{k=0}^{+\infty} \tau_k^2 < +\infty$

2. $\sum_{k=0}^{+\infty} \tau_k = +\infty$ .

Remember that the whole algorithm works because the indices $i_k$

$$x^{k+1} = x^k - \tau_k \nabla f_{i_k}(x^k)$$

are choosen uniformly at random from $\{1, \dots, m\}$

SGD is <u>batch mode</u> optimization method It is NOT an online method

<span style="color:red">Minimizing Movements. ( Actually This is the discrete version of MM )</span>

$$x^{k+1} = x^k - \tau \nabla f(x^k)$$

$$x^{k+1} \in \underset{z \in X}{\arg\min} \quad f(z) + \frac{1}{2\tau} |z - z^k|^2 \qquad (*)$$

" I have to choose the next point ($x^{k+1}$) in such a way that I minimize $f$ but <u>also</u> I don't go too far away from $z^k$ "

1. Notice that there is no gradient here

2. Instead of $|x - x^k|^2$ I can put a generic distance or score
$$d(x, z^k)$$

Now if $f \in e^1(\mathbb{R}^n, \mathbb{R})$ then

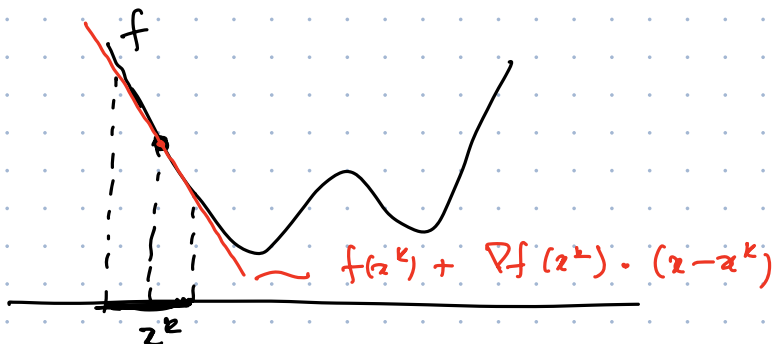$$\nabla \left( f + \frac{1}{2t} |\cdot - z^k|^2 \right) (x^{k+1}) = 0$$

$$\nabla f(x^{k+1}) + \frac{1}{t}(x^{k+1} - x^k) = 0 \quad \Rightarrow \quad x^{k+1} = x^k - \tau \nabla f(x^{k+1})$$

<div align="right">implicit GD.</div>

How do we recover explicit G D.

Since we are looking for solutions close to $x^k$ and since $f$ is diff. we can approximate $f$ around $z^k$ with

$$f(x) \approx f(x^k) + \nabla f(x^k) \cdot (z - x^k)$$



$$\widetilde{\phantom{xx}} f(z^k) + \nabla f(z^k) \cdot (z - z^k)$$

$z^k$

So I can say that instead of $(*)$ I use

$$x^{k+1} \in \underset{z \in X}{\arg\min} \quad \underbrace{f(z^k) + \nabla f(z^k) \cdot (x - z^k)}_{f(x)} + \frac{1}{2\tau} |x - x^k|^2$$

Now the minimality cond. of $\xrightarrow{\quad}$ is

$$\nabla \left( f(z^k) + \nabla f(z^k) \cdot (z - z^k) + \frac{1}{2\tau} |z - z^k|^2 \right) (z^{k+1}) = 0$$

$$\nabla f(z^k) + \frac{1}{t}(z^{k+1} - z^k) \quad \Rightarrow \quad z^{k+1} = x^k - \tau \nabla f(z^k)$$