

We are interested in proving the convergence of

$$x^{k+1} = x^k - \tau \nabla f(x^k) \quad \text{when } \tau > 0 \quad (\text{Gradient Step})$$

The gradient descent algorithm is the following

1. fix  $x_0 \in X \subseteq \{\mathbb{R}^n\}$  Hilbert space  $\ell^2, L^2, \dots$
2. Compute a gradient step

$$\underline{x^{k+1}} = \underline{x^k - \tau \nabla f(x^k)} \quad \text{for } \tau > 0$$

In this way you compute  $x_1, x_2, \dots$  that defines a sequence of points in  $X$ .  
Theorem Let  $f \in C^1(X; \mathbb{R})$  with  $\nabla f$  L-Lipschitz and  $\inf f > -\infty$ .  
 the function  $f$  is bounded from below.

If you choose  $0 < \tau < 2/L$  then the method converges in the sense that

$$*\nabla f(x^k) \rightarrow 0 \quad \text{up to subsequences.}$$

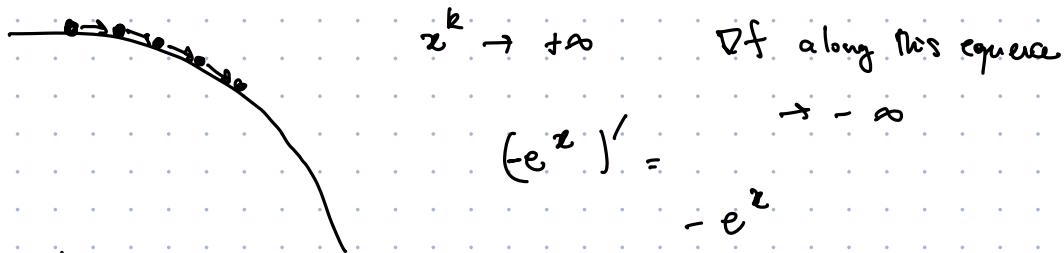
Rem 1. We are working under milder assumptions (no convexity, no coercivity etc.) so we cannot expect to say much about the convergence of the sequence.

Rem 2. In ML we don't need to actually say too much about the convergence of the weights.

Why  $\inf f > -\infty$ ?

$e^x$  is bounded from below  $e^x \geq 0 \quad \inf e^x = 0$

$$\inf -e^x = -\infty$$



But  $\nabla(-e^x)$  is not even Lipschitz

Another example is  $f(z) = -z$

is a  $C^1$  function

$$\text{The } \nabla f = f' = -1$$

$$|-1 - (-1)| \leq L|x-y| \quad \forall x, y \in \mathbb{R}$$

$$0 \leq L|x-y| \quad \forall L$$

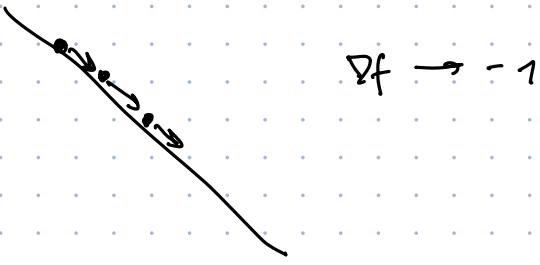
A function  $g: X \rightarrow Y$  is said to be L-Lipschitz if

$$|g(x) - g(y)|_X \leq L|x-y|_X$$

Simple examples of Lipschitz functions are  $C^1$  functions with  $|g'|$  bounded

$$L = \sup_x |g'(x)|$$

But  $f(z) = -z$  has the same problem of  $-e^{-x}$



Proof  $f(z^{k+1}) = f(z^k - \tau \nabla f(z^k))$  by definition of  $z^{k+1}$

$$\phi(t) = f(z + ty) \quad \forall z, y \in X \quad t \in \mathbb{R}$$

Since  $f \in C^1$  then  $\phi$  is differentiable and

$$\phi'(t) = \nabla f(z + ty) \cdot y$$

Now I integrate between 0 and  $\tau$  both sides

$$\int_0^\tau \phi'(s) ds = \int_0^\tau \nabla f(z + sy) \cdot y ds$$

$$\phi(\tau) - \phi(0) = \int_0^\tau \nabla f(z + sy) \cdot y ds -$$

$$f(z + \tau y) = f(z) + \int_0^\tau \nabla f(z + sy) \cdot y ds.$$

Lemma 1 If  $f: X \rightarrow \mathbb{R}$  is  $C^1$  then  $\forall \tau$

$$f(z + \tau y) = f(z) + \int_0^\tau \nabla f(z + sy) \cdot y ds \quad \forall z, y \in X.$$

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

Let us now use Lemma 1 with  $z = z^k$  and  $y = -\nabla f(z^k)$

$$f(z^k - \tau \nabla f(z^k)) = f(z^k) + \int_0^\tau \nabla f(z^k - s \nabla f(z^k)) \cdot (-\nabla f(z^k)) ds.$$

$$= f(z^k) - \int_0^\tau \nabla f(z^k - s \nabla f(z^k)) \cdot \nabla f(z^k) ds.$$

The idea now is to use Lipschitz condition on  $\nabla f$

$$\nabla f(z) - \nabla f(y)$$

$$f(z^k) - \int_0^\tau (\nabla f(z^k - s \nabla f(z^k)) - \nabla f(z^k)) \cdot \nabla f(z^k) ds - \tau |\nabla f(z^k)|^2$$

$$\int_0^\tau (\nabla f(z^k - s \nabla f(z^k)) - \nabla f(z^k) + \underbrace{\nabla f(z^k)}_{|\nabla f|^2} \cdot \nabla f(z^k) ds$$

$$f(x^{k+1}) = f(x^k) - \tau |\nabla f(x^k)|^2 - \int_0^\tau (\underbrace{(\nabla f(x^k) - s \nabla f(x^k)) \cdot \nabla f(x^k)}_{\text{Lipschitz condition}}) ds$$

We use the fact that  $\forall x \in \mathbb{R} \quad x \leq |x|$

$$\leq f(x^k) - \tau |\nabla f(x^k)|^2 + \int_0^\tau |(\nabla f(x^k) - s \nabla f(x^k)) \cdot \nabla f(x^k)| ds.$$

Now by Cauchy-Schwarz

$$(|x \cdot y| \leq \|x\| \|y\|)$$

$$\leq f(x^k) - \tau |\nabla f(x^k)|^2 + \int_0^\tau |\nabla f(x^k) - s \nabla f(x^k)| |\nabla f(x^k)| ds$$

Finally I can use Lipschitz condition on  $\nabla f$

$$\leq f(x^k) - \tau |\nabla f(x^k)|^2 + \int_0^\tau L |x^k - s \nabla f(x^k) \cdot x^k| |\nabla f(x^k)| ds.$$

$$= f(x^k) - \tau |\nabla f(x^k)|^2 + L |\nabla f(x^k)|^2 \int_0^\tau s ds$$

$$= f(x^k) - \tau |\nabla f(x^k)|^2 + L \frac{\tau^2}{2} |\nabla f(x^k)|^2$$

$$= f(x^k) + \tau \left( \frac{L}{2} \tau - 1 \right) |\nabla f(x^k)|^2.$$

$$= f(x^k) - \tau \left( 1 - \frac{L\tau}{2} \right) |\nabla f(x^k)|^2$$

So in the end we had

$$f(x^{k+1}) \leq f(x^k) - \tau \left( 1 - \frac{L\tau}{2} \right) |\nabla f(x^k)|^2. \quad (*)$$

Now I can start from  $x^0$  and use (\*) recursively

$$f(x^1) \leq f(x^0) - \tau \left( 1 - \frac{L\tau}{2} \right) |\nabla f(x^0)|^2$$

$$f(x^2) \leq f(x^1) - \tau \left( 1 - \frac{L\tau}{2} \right) |\nabla f(x^1)|^2$$

$$\leq f(x^0) - \tau \left( 1 - \frac{L\tau}{2} \right) (|\nabla f(x^0)|^2 + |\nabla f(x^1)|^2)$$

$\vdots$

$$f(x^n) \leq f(x^0) - \underbrace{\tau \left( 1 - \frac{L\tau}{2} \right)}_{>0} \sum_{i=0}^{n-1} |\nabla f(x^i)|^2.$$

Remember that by hypothesis  $\tau < \frac{2}{L}$ , so  $1 - \frac{L\tau}{2} > 0$  since

$$-\frac{L\tau}{2} > -1 \quad \frac{L\tau}{2} < 1 \quad \tau < \frac{2}{L}$$

$$-\infty < \inf f \leq f(x^n) \leq f(x^0) - \tau \left(1 - \frac{L\tau}{2}\right) \sum_{i=0}^{n-1} |\nabla f(x^i)|^2$$

$$\tau \left(1 - \frac{L\tau}{2}\right) \sum_{i=0}^{n-1} |\nabla f(x^i)|^2 \geq f(x^0) - \inf f$$

Notice that  $f(x^0) - \inf f$  is a fixed positive quantity that does not depend on  $N$ .

This means that if we define

$$S_n = \sum_{i=0}^{n-1} |\nabla f(x^i)|^2$$

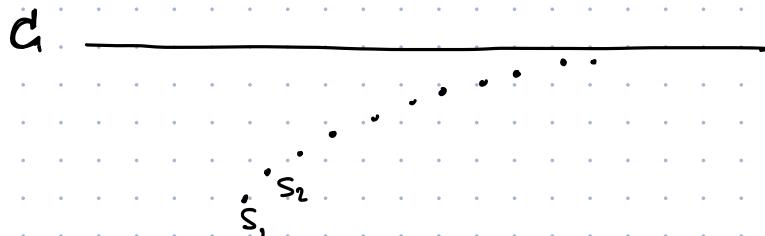
This quantity as  $n \rightarrow +\infty$  is an increasing sequence meaning that

$$S_{n+1} \geq S_n.$$

Moreover  $(S_n)$  is bounded from above by

$$C = \frac{f(x_0) - \inf f}{\tau \left(1 - \frac{L\tau}{2}\right)} \geq 0$$

The situation is this



By Bolzano-Weierstrass Theorem, an increasing sequence bounded from above converges.

What is  $\lim_{n \rightarrow +\infty} S_n$  ?  $S_n = \sum_{i=0}^{n-1} |\nabla f(x^i)|^2$

By definition

$$\lim_{n \rightarrow +\infty} S_n = \sum_{i=0}^{+\infty} |\nabla f(x^i)|^2 < +\infty$$

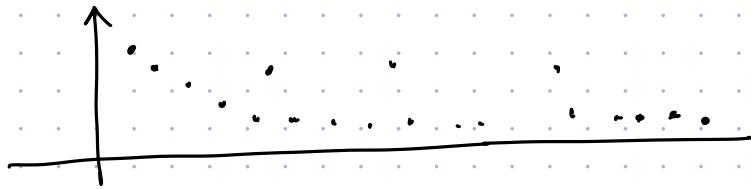
And the series converges.

$$\sum_{i=0}^{+\infty} a_k < +\infty \text{ Then it is not properly true to say } a_k \rightarrow 0.$$

This means that  $\exists$  a subsequence  $a_{k_k}$   $k=1, 2, \dots$  such that

$$|\nabla f(x^{k_k})|^2 \rightarrow 0 \text{ as } k \rightarrow +\infty$$

Run: From our perspective convergence up to subsequences is perfectly fine.



From a numerical prospective you basically fix a threshold  $\varepsilon = 10^{-7}$  (this is an example) and you stop your algorithm when

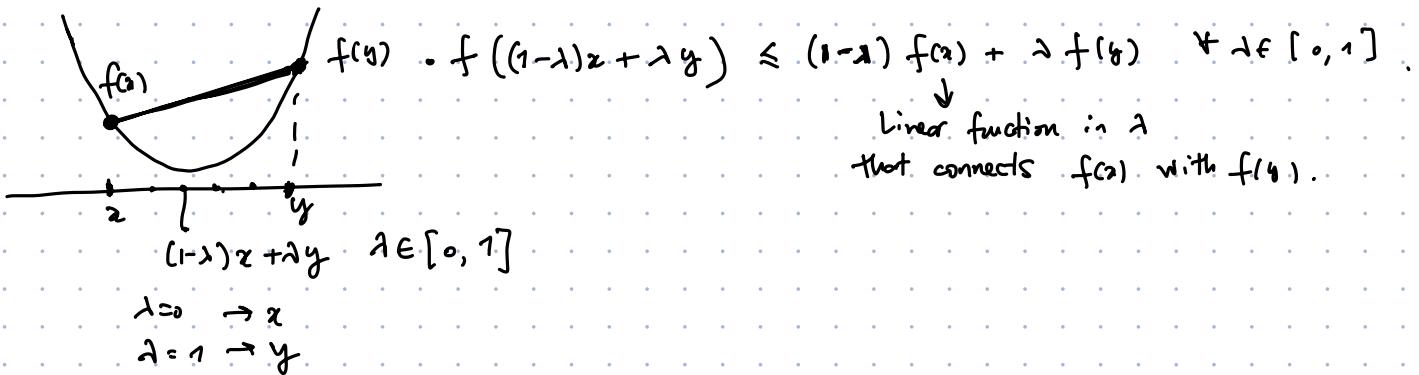
$$|\nabla f|^2 < \varepsilon.$$

This theorem says that you will reach that point.

But in general we don't (and we can't) have an estimate on the time (= number of steps) we will need to reach that condition.

### Convex Case

Disclaimer: now we work in  $\mathbb{R}^n$ , so  $X \subset \mathbb{R}^n$  (i.e. in a finite dimensional Hilbert space)



$f$  is convex if and only if  $\forall x, y \in X$

$$\bullet \quad f((1-\lambda)x + \lambda y) \leq (1-\lambda)f(x) + \lambda f(y) \quad \forall \lambda \in [0, 1], \forall x, y \in X$$

You should remember that

$g \in C^2(a, b; \mathbb{R})$  then  $g$  is convex iff  $g'' \geq 0$  on  $(a, b)$   
 $\leftarrow$  now we are in finite dim

Proposition: If  $f \in C^2(\mathbb{R}^n, \mathbb{R})$  then  $f$  is convex if and only if  $D^2f \geq 0$

Hessian:  $(D^2f)_{ij} = \frac{\partial^2}{\partial x_i \partial x_j} f \quad \forall i, j = 1, \dots, n$ .

Proof: Fix  $x, y \in \mathbb{R}^n$  and define  $\psi(s) = f(x + sy)$ .

Since  $f$  is convex then  $\psi$  is convex too.  $\quad (1-\lambda)(x+ty)$

$$\begin{aligned} \psi((1-\lambda)t + \lambda s) &= f(x + [(1-\lambda)t + \lambda s]y) = f(x - \underline{\lambda}x + \lambda x + [(1-\lambda)t + \lambda s]\underline{y}) \\ &= f((1-\lambda)x + (1-\lambda)ty + \lambda(x + sy)) = f((1-\lambda)(x+ty) + \lambda(x+sy)) \leq (1-\lambda)\psi(t) + \lambda\psi(s) \end{aligned}$$

$\leftarrow$  by convexity of  $f$

$f(x+ty) \quad f(x+sy)$

this proves that if  $f$  is convex then  $\psi$  is convex. But if  $\forall z, y \in \mathbb{R}^n$   $\psi$  is convex, by definition of convexity also  $f$  is convex

But  $\psi$  is convex

if and only if  $\psi''(s) \geq 0$

$$\psi'(s) = \nabla f(z + sy) \cdot y, \quad \psi''(s) = \boxed{y \cdot D^2 f(y) \geq 0} \quad \forall y \in \mathbb{R}^n$$

On the other hand saying that

$$y \cdot D^2 f y \geq 0 \quad \forall y \in \mathbb{R}^n \text{ means } \underline{\underline{D^2 f \geq 0}} \quad \square$$

Theorem (Baillon - Haddad) If  $f$  is convex and  $C^2(\mathbb{R}^n, \mathbb{R})$  and  $\nabla f$  is  $L$ -Lipschitz then

$$(\nabla f(z) - \nabla f(y)) \cdot (z - y) \geq \frac{1}{L} \| \nabla f(z) - \nabla f(y) \|^2.$$

Lemma 1 If  $f: X \rightarrow \mathbb{R}$  is  $C^1$  then  $\forall \tau$

$$f(z + \tau y) = f(z) + \int_0^\tau \nabla f(z + sy) \cdot y \, ds \quad \forall z, y \in X.$$

$$\nabla f(z) - \nabla f(y) = \int_0^1 D^2 f(y + s(z-y)) \cdot (z-y) \, ds.$$

$$X = z + y \quad Y = X - Y \Rightarrow z + sy = Y + s(X - Y)$$

$$Y = z$$

Then

$$\nabla f(z) - \nabla f(y) = \int_0^1 D^2 f(y + s(z-y)) \cdot (z-y) \, ds.$$

Then we know that  $0 \leq D^2 f \leq L \text{Id}$

Define  $A = \int_0^1 D^2 f(y + s(z-y)) \, ds$  which is symmetric since

$$A = \int D^2 f \quad \text{and} \quad 0 \leq A \leq L \text{Id}.$$

$$(*) \quad \nabla f(z) - \nabla f(y) = A(z-y) \quad \text{by def of } \|z\|^2.$$

$$|\nabla f(z) - \nabla f(y)|^2 = |A(z-y)|^2 = A(z-y) \cdot A(z-y)$$

Then  $A^{1/2}$  exists ( $A^{1/2} \cdot A^{1/2} = A$ ) and commutes with  $A$

$$[A, A^{1/2}] = 0$$

$$AA^{1/2} = A^{1/2}A$$

$$|\nabla f(x) - \nabla f(y)|^2 = A^{1/2}A(x-y) \cdot A^{1/2}(x-y) = A A^{1/2}(x-y) \cdot A^{1/2}(x-y)$$

Since  $A \in L\mathbb{I}d$

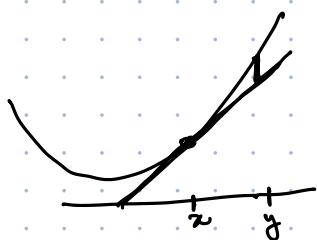
$$\leq L [A^{1/2}(x-y)] \cdot [A^{1/2}(x-y)] = L[A(x-y)] \cdot (x-y)$$

$$\stackrel{(*)}{=} L(\nabla f(x) - \nabla f(y)) \cdot (x-y)$$

$$|\nabla f(x) - \nabla f(y)|^2 \leq L (\nabla f(x) - \nabla f(y)) \cdot (x-y)$$

$$f \in C^2 \text{ is convex} \iff D^2 f \geq 0$$

$$f \in C^1 \text{ is convex} \iff \boxed{f(y) \geq f(x) + \nabla f(x) \cdot (y-x)}$$



$$\nabla f = [0, 1]$$

$$\Psi(t) = f(x + t(y-x))$$

$$\Psi(0) = f(x)$$

$$\Psi(1) = f(y)$$

$$\Psi'(t) = \nabla f(x + t(y-x)) \cdot (y-x)$$

$$\Psi'(0) = \nabla f(x) \cdot (y-x)$$

$$\Psi(1) - \Psi(0) \geq \Psi'(0) \quad \text{since } \Psi \text{ is convex.}$$

$$f(y) - f(x) \geq \nabla f(x) \cdot (y-x).$$

We need the following fact:

$$\text{Define } T_\tau = id - \tau \nabla f$$

$$T_\tau(x) = x - \tau \nabla f(x)$$

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

$$x^{k+1} = T_\tau(x^k)$$

$T_\tau$  is the operator that generates the sequence of points of GD.

Lemma  $T_\tau$  is 1-Lipschitz.

$$|T_\tau(x) - T_\tau(y)| \leq 1 \cdot |x-y|$$

$$\begin{aligned}
\text{Proof.} \quad & |T_\tau(x) - T_\tau(y)|^2 = |x - \tau \nabla f(x) - y + \tau \nabla f(y)|^2 \\
&= |x - y - \tau (\nabla f(x) - \nabla f(y))|^2 = |x - y|^2 - 2\tau(x - y) \cdot (\nabla f(x) - \nabla f(y)) \\
&\quad + \tau^2 |\nabla f(x) - \nabla f(y)|^2 \\
&\leq |x - y|^2 + \tau^2 |\nabla f(x) - \nabla f(y)|^2 - \underbrace{\frac{2\tau}{L}}_{\text{Since } \tau < \frac{2}{L}} |\nabla f(x) - \nabla f(y)|^2 \\
&= |x - y|^2 - \underbrace{\frac{2\tau}{L} \left(1 - \frac{\tau L}{2}\right)}_{\text{Since } \tau < \frac{2}{L}} |\nabla f(x) - \nabla f(y)|^2.
\end{aligned}$$

Now we have everything to give a bound on the convergence of GD

Assume the  $f$  has a minimum  $x^*$ .

$$f(x^*) \geq f(x^k) + \nabla f(x^k) \cdot (x^k - x^*) \quad \text{by the usual characterization of convex functions}$$

$$f(x^k) - f(x^*) \geq \nabla f(x^k) \cdot (x^k - x^*)$$

$$f(x^k) - f(x^*) \leq \|\nabla f(x^k)\| \|x^k - x^*\| \quad (*)$$

Now, since  $x^*$  is a minimum  $T_\tau(x^*) = x^*$

$$|x^{k+1} - x^0| = |T_\tau(x^k) - T_\tau(x^*)| \leq \|x^k - x^*\|$$

$$|x^k - x^*| \leq |x^0 - x^*|$$

$$\text{From } (*) \quad \Delta_k = f(x^k) - f(x^*)$$

$$\Delta_{k+1} \leq \Delta_k - C \frac{\Delta_k^2}{\|x_k - x^*\|^2} \quad (**)$$

If  $\Delta_k$  satisfies  $(**)$  then

$$\boxed{\Delta_k \leq \frac{|x_0 - x^*|^2}{C(k+1)}}$$

Chambolle Antoine, Leon Bottou

Robbins-Monroe '60 (very very hard to read)

Remember that  $f$  in ML is the empirical risk so

$$f(x) = \frac{1}{m} \sum_{i=1}^m f_i(x)$$

$$\hat{L}_N(w) = \frac{1}{m} \sum_{i=1}^m l(f(x_i, w), y_i)$$

$f$  here is The model.

The functions  $f_i$  are basically the loss function computed on the  $i$ -th example of the training set.

$$l(f(x_i, \cdot), y_i) \rightarrow f_i$$

The dimension of the training set  $m$  is very very large

For modern ML problems it is not feasible to compute  $\nabla f$

The reason why we can still optimize  $f$  with a gradient based method is SGD

SGD:

1. Choose  $x_0 \in \mathbb{R}^n$  randomly and you set  $x^0 = x_0$
2. For every  $k \geq 0$  select  $i_k \in \{1, \dots, m\}$  sampled with probability  $1/m$  and you compute

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

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

$$f(x^k) = \frac{1}{m} \sum_{i=1}^m f_i(x^k)$$

- $\inf f > -\infty$
- $f_i \in C^1$
- $\nabla f_i$  should be  $L_i$ -Lipschitz

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

$$\nabla f(x^k) = \frac{1}{m} \sum_{i=1}^m \nabla f_i(x^k)$$

Now choose  $j \in \{1, \dots, m\}$

$$f_j(x^{k+1}) = f_j(x^k - \tau \nabla f_{I_k}(x^k))$$

Lemma 1 If  $f: X \rightarrow \mathbb{R}$  is  $C^1$  then  $\forall \tau$

$$f(x + \tau y) = f(x) + \int_0^\tau \nabla f(x + sy) \cdot y \, ds \quad \forall x, y \in X.$$

We use this lemma with  $x = x^k$   $y = -\nabla f_{I_k}(x^k)$   $f = f_j$

$$f_j(x^k - \tau \nabla f_{I_k}(x^k)) = f_j(x^k) + \int_0^\tau \nabla f_j(x^k - s \nabla f_{I_k}(x^k)) \cdot (-\nabla f_{I_k}(x^k)) \, ds$$

$$= f_j(x^k) - \int_0^\tau \nabla f_j(x^k - s \nabla f_{I_k}(x^k)) \cdot \nabla f_{I_k}(x^k) ds.$$

$$= f_j(x^k) - \int_0^\tau (\nabla f_j(x^k - s \nabla f_{I_k}(x^k)) - \nabla f_j(x^k) + \nabla f_j(x^k)) \cdot \nabla f_{I_k}(x^k) ds.$$

$$= f_j(x^k) - \tau \nabla f_j(x^k) \cdot \nabla f_{I_k}(x^k) - \int_0^\tau (\nabla f_j(x^k - s \nabla f_{I_k}(x^k)) - \nabla f_j(x^k)) \cdot \nabla f_{I_k}(x^k) ds.$$

$\checkmark$  Cauchy-Schwarz + Lipschitz of  $\nabla f_j$

$$\leq f_j(x^k) - \tau \nabla f_j(x^k) \cdot \nabla f_{I_k}(x^k) + \int_0^\tau \|L_j\| \|s \nabla f_{I_k}(x^k)\| \| \nabla f_{I_k}(x^k)\| ds.$$

$$= f_j(x^k) - \tau \nabla f_j(x^k) \cdot \nabla f_{I_k}(x^k) + \frac{\tau^2 L_j}{2} \| \nabla f_{I_k}(x^k) \|^2$$

$$f_j(x^{k+1}) \leq f_j(x^k) - \tau \nabla f_j(x^k) \cdot \nabla f_{I_k}(x^k) + \frac{\tau^2 L_j}{2} \| \nabla f_{I_k}(x^k) \|^2.$$

Now we sum both sides for  $j=1 \dots m$  and divide by  $m$

$$\frac{1}{m} \sum_{j=1}^m f_j(x^{k+1}) = f(x^{k+1})$$

$$f(x^{k+1}) \leq f(x^k) - \tau \nabla f(x^k) \cdot \nabla f_{I_k}(x^k) + \frac{\tau^2 \bar{L}}{2} \| \nabla f_{I_k}(x^k) \|^2$$

$$\text{where } \bar{L} = \frac{1}{m} \sum_{i=1}^m L_i$$

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

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

$$E(x - Ex)^2 = E x^2 - (Ex)^2$$

$$\frac{1}{m} \sum_{i=1}^m \| \nabla f_i(x^k) - \nabla f(x^k) \|^2 = \frac{1}{m} \sum_{i=1}^m \| \nabla f_i(x^k) \|^2 + \| \nabla f(x^k) \|^2 - 2 \nabla f_i(x^k) \cdot \nabla f(x^k)$$

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

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

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

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

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

$$E(E(f(x)|Y)) = EX$$

$$E(E(f(x^{k+1})|x^k)) = E f(x^{k+1})$$

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

In order to prove the convergence of SGD we need to add an hypothesis on the variance of the gradients.

$$\frac{1}{m} \sum_{i=1}^m |\nabla f_i(x) - \nabla f(x)|^2 \leq \sigma^2 \quad \forall x \in \mathbb{R}^n \quad \text{This is a reasonable assumption}$$

$$E f(x^{k+1}) \leq E f(x^k) - \tau \left(1 - \frac{\tau L}{2}\right) E|\nabla f(x^k)|^2 + \frac{\tau^2 L}{2} \sigma^2$$

Idea is that  $\tau \rightarrow \tau_k$

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

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