

What we usually do in ML is to minimize the empirical risk

$$L_N(f) = \frac{1}{N} \sum_{i=1}^N l(f(x_i), y_i).$$

on what we called \mathcal{P} . From now on we will always assume that

\mathcal{P} = "set of NN with a fixed architecture"

$f \in \mathcal{P}$ can always be represented as

$$f(\cdot; w)$$

↗ weights of the NN

Therefore if we use this representation in the empirical risk we obtain the following

$$L_N(f(\cdot; w)) = \frac{1}{N} \sum_{i=1}^N l(f(x_i; w), y_i)$$

$l(a, b) = \frac{1}{2} |a - b|^2$

Rem. Once l is fixed, for instance it's a quadratic loss or Cross Entropy loss or ...

since (x_i, y_i) are also given because they are the points of a given training set then

\bullet $=$ is a function solely of the weights.

$$L_N(f(\cdot; w)) = \hat{L}_N(w) \quad \text{if } w \in \mathbb{R}^n$$

n = number of weights of the NN

Side note

$$\begin{aligned} z_i &= \sigma \left(\sum_j w_{ij} x_j \right) \\ &= \sigma \left(\sum_j w_{ij} x_j + b_i \right) \end{aligned}$$

↑ bias.

Therefore the minimization of L_N is equivalent to the minimization of \hat{L}_N .

The problem that we need to solve is

$$\min_{w \in \mathbb{R}^n} \hat{L}_N(w)$$

Rem We are going to solve tackle this problem using "gradient-based" methods, which of course involve the computation of $\nabla \hat{L}_N$

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

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

↓ yesterday we called this $C^1(w)$

1. So ℓ is usually given. For example : $f \quad \ell(a, b) = \frac{1}{2} (a - b)^2$
 $\ell(a, b) = (a - b)$

2. $\nabla_w f(x_i, w)$ is the gradient of a NN with respect to its weights.
And we saw that this can be computed efficiently $O(n)$ using backpropagation

Continuous Optimization

New Notation



Brief introduction on Hilbert spaces (and informal)

We say that X is an Hilbert space if it is a vector space with an inner product
(
you know how to
+ and
multiply by scalars)

$$\cdot : X \times X \rightarrow \mathbb{R}$$

$$(x, y) \mapsto x \cdot y$$

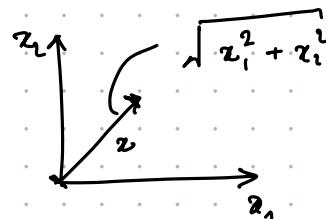
if $X = \mathbb{R}^n$ a natural choice for \cdot is the standard scalar product in \mathbb{R}^n

$$x \cdot y = \sum_{i=1}^n x_i \cdot y_i$$

Once you have an inner product then you can define a norm $\| \cdot \|_X$

$$\| x \|_X = \sqrt{x \cdot x}$$

$$\text{In } \mathbb{R}^n \quad \| x \| = \sqrt{\sum_{i=1}^n (x_i)^2}$$



X is Hilbert if it is complete with respect to this norm

complete = Cauchy sequences converge.

A sequence (x_n) is a Cauchy sequence if $\forall \varepsilon > 0 \exists M > 0$ s.t. $\forall m, n > M$

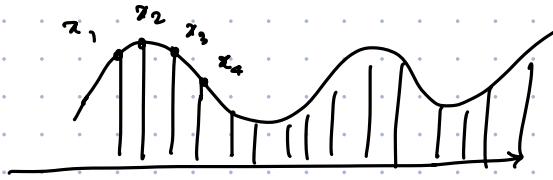
$$\| x_n - x_m \|_X < \varepsilon.$$

Example 1. $\ell^2 = \{ (x_n) : \sum_{n=1}^{+\infty} x_n^2 < +\infty \}.$

$x_1, x_2, x_3, x_4, \dots$

and $x \cdot y = \sum_{i=1}^{+\infty} x_i y_i$

Ex. Prove that ℓ^2 is complete and therefore it is Hilbert.



$$L^2(a, b) \subset \{ f : (a, b) \rightarrow \mathbb{R} : \int_a^b f^2 dx < +\infty \}.$$

$$f \cdot g = \int_a^b f \cdot g \, dx$$

Now, the problem that we want to solve is

$$\min_{x \in X} f(x) \quad \text{where } X \text{ is Hilbert and } f : X \rightarrow \mathbb{R}.$$

Gradients

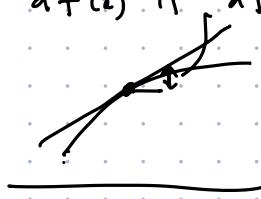
X^* is called the dual space of X and it is the space of all continuous linear functionals on X .

$$X^* = \{ L : X \rightarrow \mathbb{R} : L \text{ is linear and continuous} \}.$$

$$L(\alpha x + \beta y) = \alpha L(x) + \beta L(y)$$

Dif. The differential of f at a point x that we denote $df(x)$ is df if it exists an element of X^* such that

$$f(y) = f(x) + \langle df(x), y-x \rangle + o(|x-y|)$$



$$\langle L, x \rangle = L(x) \quad \text{if } L \in X^* \text{ and } x \in X$$

If the differential $df(x)$ exists for every $x \in X$ then we call f differentiable and we write

$$f \in C^1(X; \mathbb{R})$$

Many times we say that f is Fréchet-differentiable.

Theorem (Rietz) If X is Hilbert then every element $L \in X^*$ can be represented as follows

$$\forall z \in X \quad \langle L, z \rangle = \pi \cdot z \quad (*)$$

$$L(z)$$

in the sense that it exists a vector $\pi \in X$ such that $(*)$ holds.

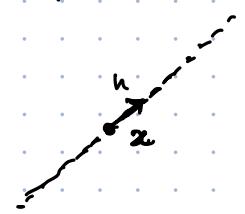
Since X is Hilbert $df(z) \in X^*$ by Rietz \exists a vector in X that we denote $\nabla f(z)$ (and we call it gradient of f in the point z) such that

$$\langle df(z), z \rangle = \nabla f(z) \cdot z \quad \forall z \in X$$

In infinite dimensions (in Hilbert) the generalization of directional derivative is called Gâteau derivative.

Fix $z \in X$ and define $\phi(t) = f(z + th)$ where $h \in X$ and $h \neq 0$.

\uparrow
I'm looking at how f behaves around z in the direction suggested by h



$\underline{df(z; h)} = \underline{\phi'(0)}$ is defined to be the "Gâteau derivative/differential" along the direction h .

In the Calculus of Variations this is also what is usually called first variation.

Observation: if f is differentiable in z , then

$$f(z + th) = f(z) + \underline{\langle df(z), th \rangle} + o(|th|)$$

$(o(1) = o(|h|)$ for functions of the variable t)

$$\frac{f(z + th) - f(z)}{t} = \underline{\langle df(z), h \rangle} + o(1)$$

If you take the limit $t \rightarrow 0$, $o(1) \rightarrow 0$

$$\boxed{f(z + th) = \underline{\phi(t)}} \quad \therefore \quad \frac{f(z + th) - f(z)}{t} = \frac{\underline{\phi(t) - \phi(0)}}{t}$$

$$\text{as } t \rightarrow 0 \quad \frac{\underline{\phi(t) - \phi(0)}}{t} \rightarrow \underline{\phi'(0)} = d$$

So we obtain the usual relation between directional derivative and partial derivatives

$$\boxed{df(z; h) = \langle df(z), h \rangle}$$

Proposition 1. Let X be Hilbert and $f \in C^1(X; \mathbb{R})$, $z \in X$ such that $d f(z) \neq 0$.

Let $\nabla f(z)$ be the gradient of f at z (Rietz). Then

$$\frac{\nabla f(z)}{\|\nabla f(z)\|} \in \operatorname{argmax}_{\|h\|=1} \{ |d f(z; h)| : h \in X \}$$

The proof of this is basically an application of Cauchy-Schwarz inequality.

Cauchy-Schwarz: If $x, y \in X$

$$|x \cdot y| \leq \|x\| \cdot \|y\|$$

and $=$ holds only if $x = \alpha y$ where $\alpha \in \mathbb{R}$.

Let us now try to prove Proposition 1.

Choose $h = \frac{\nabla f(z)}{\|\nabla f(z)\|}$

$$d f(z; h) = \nabla f(z) \cdot \frac{\nabla f(z)}{\|\nabla f(z)\|} = \frac{\|\nabla f(z)\|^2}{\|\nabla f(z)\|} = \|\nabla f(z)\|.$$

Now for any other direction

$$|d f(z; h)| = |\nabla f(z) \cdot h| \leq \|\nabla f(z)\| \cdot \|h\| = \underline{\|\nabla f(z)\|}.$$

This is basically the observation that motivates gradient based optimization algorithms

whose most famous representative is GrD. (gradient descent).

$$\boxed{z^{k+1} = z^k - \tau \nabla f(z^k)} \quad \text{Gradient Descent.}$$

Curves of maximal slope. (Ennio De Giorgi)

Are a way to extend gradient flows in metric spaces.

Metric spaces = (\mathcal{G}, d)
 \uparrow \wedge
 space distance

Gradient Flows

Are solutions of

$$\begin{cases} u' = -\nabla f(u) \\ u(0) = u_0 \end{cases} \quad \text{It is an ODE on } (0, \infty)$$

Gradient flows are "the continuous version" of Gradient descent.

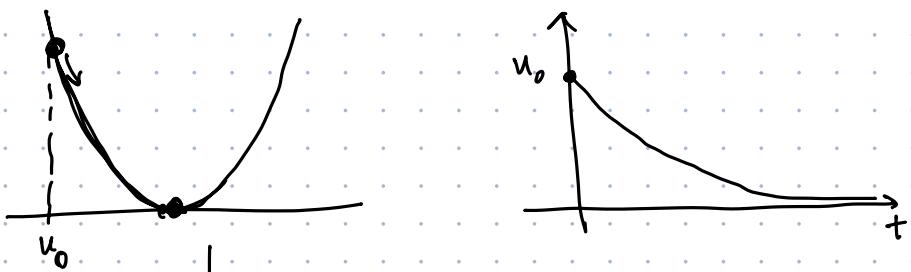
As often happens it is easier to analyze dynamics in continuous time.
because often you have simple closed form solutions..

Example of GF

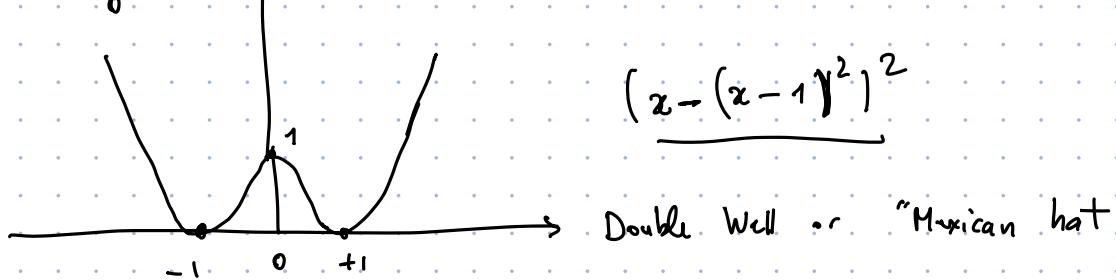
$$\left\| \begin{array}{l} f(u) = \frac{u^2}{2} \\ \left(\begin{array}{l} u'(t) = -u(t) \\ u(0) = u_0 \end{array} \right) \quad t > 0 \end{array} \right.$$

The general solution is $u(t) = A e^{-t}$
 $u(0) = A = u_0$

The solution of (*) is $u(t) = \underline{\underline{u_0 e^{-t}}}$



Exercise



$f(u(t))$ this decreases and indeed

$$\boxed{\frac{d}{dt} f(u(t)) = \nabla f(u(t)) \cdot u'(t) = -(\nabla f(u(t)))^2 \leq 0}$$

Curves of maximal slope ("Gradient Flows in Metric Spaces" Ambrosio, Gigli, Savaré)

↑
Very hard / difficult book

$$u'(t) = -\nabla \phi(u(t)) \quad \text{Gradient Flow}$$

In order to write this equation you need two things:

1. u' derivative of $t \mapsto u(t)$
 2. $D^k u(t)$ you need the gradient. (for instance if $t \in \mathbb{C}^n$ in hilbert or banach ok!)

If you don't have this two objects you can't even write the equation.

In metric spaces you don't have this.

IDEA Even you can't define n' and $\nabla\phi$ you can define something that generalizes
 $|n'|$ and $|\nabla\phi|$.

For instance if $t \rightarrow \sigma(t)$ is a curve in \mathcal{P} and it is absolutely continuous, then you can define

$$|v'|_r(t) = \lim_{h \rightarrow 0} \frac{d(v(t+h), v(t))}{h} \leftarrow \text{metric derivative.}$$

Obs. You do not define r' and then you take the norm. (You can't do this)
you directly define $|r'|$

In the same way you can define what is called an upper gradient which is a renegate of the $|Df|$.

Usually the interesting object is the "Local slope"

$$|\partial \phi|(r) = \limsup_{w \rightarrow r} \frac{(\phi(w) - \phi(r))^+}{d(w, r)}$$



Ok I have metric analogues of $|u'|$ and $|\nabla \phi(u)|$ but the gradient flow equation is not just a statement in terms of u' and $\nabla \phi$, direction are important.

In other words

$$u' = -\nabla \phi(u) \quad \xrightarrow{\text{take } |\cdot|} \quad |u'| = |\nabla \phi|$$

(•)

(• •)

(∞) is surely not equivalent to (\circ)

The informative content of (∞) is less than that of (\circ)

De Giorgi idea: hilbert \uparrow you are recovering the sign of (\circ)

$$\begin{aligned}\frac{d}{dt} \phi(u(t)) &= \nabla \phi(u(t)) \cdot u'(t) = -|\nabla \phi(u(t))| |u'(t)| \\ &= -\frac{1}{2} |\nabla \phi|^2 - \frac{1}{2} |u'|^2 \\ &\uparrow \text{by using } (\infty)\end{aligned}$$

Now (\circ) is equivalent to

$$\begin{aligned}\circ \quad \frac{d}{dt} \phi(u(t)) &= -\frac{1}{2} |u'|^2 - \frac{1}{2} |\nabla \phi(u)|^2 \\ \int_s^t \left(\frac{1}{2} |u'(r)|^2 + \frac{1}{2} |\nabla \phi(u(r))|^2 \right) dr &= \phi(u(s)) - \phi(u(t))\end{aligned}$$

You say that u is a curve of maximal slope if it exists g such that

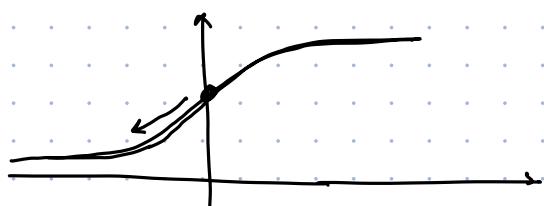
$$\int_s^t \left(\frac{1}{2} |u'(r)|^2 + \frac{1}{2} |g(u(r))|^2 \right) dr \leq \phi(u(s)) - \phi(u(t))$$

$\forall L^1$ -a.e. t, s with $s < t$.

Proof of convergence of GD under assumptions that are ML-friendly.

Basically we want to give a proof of convergence of GD without using the usual convexity assumption.

Next time we'll try to consider also the convex case to see what we can say more than the non-convex



✓ Hilbert.

What is GD? It is a procedure that computes a sequence of points $(x^n) \subset X$ defined as follows

$$(GD) \quad \begin{cases} 1. \quad x^0 = x_0 \in X \quad x_0 \text{ is a fixed point.} \\ 2. \quad x^{k+1} = x^k - \tau \nabla f(x^k) \quad k = 0, 1, 2, 3, \dots \quad \tau > 0 \end{cases}$$

\downarrow in ML is called "learning rate"

Obs1. $x^{k+1} = x^k - \tau \nabla f(x^k)$ is an explicit Euler scheme with step size τ for the gradient flow $\dot{x} = -\nabla f(x)$

Explicit Euler = Forward Euler
Implicit Euler = Backward Euler

$$\begin{aligned} \dot{x} = F(x) \rightarrow x^{k+1} &= x^k + \tau F(x^k) \\ &\downarrow \quad \text{explicit} \\ x^{k+1} &= x^k - \tau F(x^k) \quad \text{implicit.} \end{aligned}$$

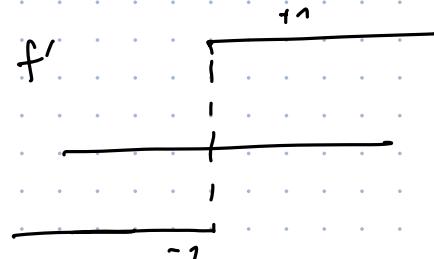
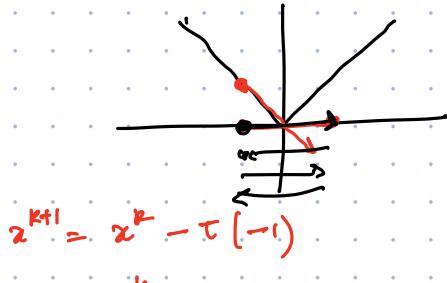
Theorem (Convergence of GD): Let $f \in C^1(X; \mathbb{R})$ and suppose that

$\boxed{\nabla f \text{ is L-Lipschitz}}$, then if $0 < \tau < \frac{2}{L}$ then (GD)

converges in the sense that

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

$$Ex \quad f(x) = |x|$$



$$x^{k+2} = x^{k+1} - \tau(1) = x^{k+1} - \tau = x^k$$