

Current Deep Learning Algorithms are based on

1. Back Propagation
2. Stochastic Gradient Descent

ML basics : Empirical Risk minimization

Classical Statistical ML setting :

We have a probability space $(\Omega, \mathcal{F}, \pi)$
Usually, e.g., in Supervised Learning

$\Omega = X \times Y$
data, examples, perceptual info.
;
table

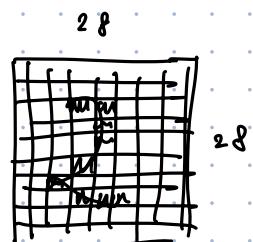
x^t → ext
 π → probability measure.
 \mathcal{F} → measurable subsets of Ω

$$(x, y) \in \Omega$$

x = image

MNIST

y = label

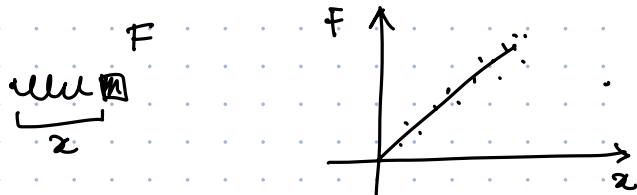


Let's talk about π

$$\begin{aligned} &(x, y) \\ &\left(\begin{array}{|c|c|} \hline \text{#} & \text{#} \\ \hline \text{#} & \text{#} \\ \hline \end{array} \right), 2 \quad \left(\begin{array}{|c|c|} \hline \text{#} & \text{#} \\ \hline \text{#} & \text{#} \\ \hline \end{array} \right), 3 \\ &\left(\begin{array}{|c|c|} \hline \text{#} & \text{#} \\ \hline \text{#} & \text{#} \\ \hline \end{array} \right), 1 \end{aligned}$$

$$X = \mathbb{R}^{784} \quad \frac{11}{2}$$

$$\bar{y} = 2$$



General problem in ML is to minimize the **Functional Risk**

Usually what you want to do is to make predictions on your examples. Meaning that you want to find $f: X \rightarrow Y$ $x \mapsto y$

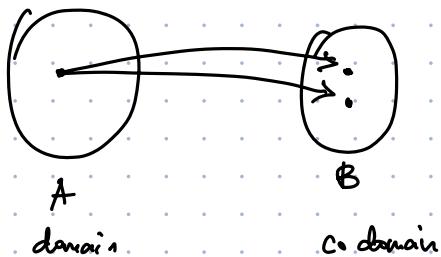
$$f(\boxed{2}) = 2$$

l is called "loss function" and measures the "goodness" of your predictions.

$$l: Y \times Y \rightarrow Y$$

Functions

f is a "rule" that associates elements from a set A to elements of a set B with the following property



f cannot associate one element of A to two distinct elements of B.

$f: A \rightarrow B$

$A \times B$: cartesian product of A and B $A \times B = \{(a, b) : a \in A, b \in B\}$.

$$\mathbb{R}^2 = \mathbb{R} \times \mathbb{R}$$

$$= \{ (x, y) : x \in \mathbb{R}, y \in \mathbb{R} \}$$



$$L(f) = \mathbb{E}_{\substack{\uparrow \\ \Omega}} l \circ (f \times d) = \int_{\Omega} l(f(z, y)) d\pi(z, y) \quad \text{Functional Risk}$$

Example $S\mathbb{R} = \mathbb{R}^2$ A is any subset of \mathbb{R}^2

$$\bullet \quad \pi(A) = C \int_A e^{-(x-y)^2} e^{-y^2} dx dy \quad C \text{ is a normalization constant.}$$

Exercise : Compute G

$$\pi(\mathbb{R}^2) = 1 \quad \text{Normalization of the probability}$$

$$\pi(R^2) = C \int_{R^2} e^{-(x-y)^2} e^{-z^2} dz dy = 1$$

$$\int_{\mathbb{R}} e^{-x^2} \left(\int_{\mathbb{R}} e^{-(z-y)^2} dy \right) dz$$

$$\int_{\mathbb{R}} e^{-(x-y)^2} dy$$

$y-x = t$, $dy = dt$

$$= \int_{\mathbb{R}} e^{-t^2} dt = \sqrt{\pi}$$

$$\int_{\mathbb{R}} e^{-x^2} \sqrt{\pi} dx = \sqrt{\pi} \int_{\mathbb{R}} e^{-x^2} dx = \pi$$

$$\pi(\mathbb{R}^2) = C\pi = 1 \quad C = \frac{1}{\pi}$$

$$p(x,y) = \frac{1}{\pi} e^{-(x-y)^2} e^{-x^2}$$

$$f : X \rightarrow Y$$

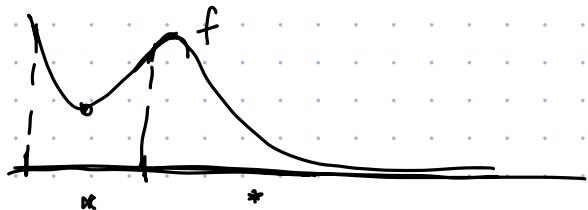
$$\text{Take } l(a,b) = (a-b)^2$$

$$l(f(x), y) = (f(x) - y)^2$$

$$L(f) = \frac{1}{\pi} \int_{\mathbb{R}^2} (f(x) - y)^2 e^{-(x-y)^2} e^{-x^2} dx dy \leftarrow$$

Problem is Minimize L on some space

$$\text{dom}(L) = \{ f : X \rightarrow Y : L(f) < +\infty \}$$



$$f(x) = \sin(x) \quad L(\sin) < +\infty \\ = +\infty$$

Calculus of variations

$$\min_{x \in X} F(x)$$

" X = space of functions." $C^0(\mathbb{R}; \mathbb{R})$ continuous functions from \mathbb{R} to \mathbb{R}

$$f(x) = x$$

$$f \in C^0(\mathbb{R}; \mathbb{R})$$



$$\sin(x) \quad e^x$$

δ is not a function

$$H(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{if } z < 0 \end{cases}$$

And in F which are integrals. In particular we are interested in functionals of the form

$$(*) \quad F(u) = \int_a^b L(u(x), u'(x), x) dx. \quad u \in X$$

$$F(u) = \int_0^1 u(x) dx \quad F(id) = \int_0^1 x dx = \frac{1}{2}$$

Therefore we want to find minima of functionals of the form $(*)$.

Algorithm : 1. Find a candidate solution $u^*(x)$ that satisfies

Euler-Lagrange. $\left(\frac{d}{dx} L_p - L_u = 0 \right) \quad " \nabla f = 0 "$

2. Check if this is indeed a minima, for instance by proving that

$$F(u) \geq F(u^*).$$

L is usually called Lagrangian. $L(u, p, x)$

L_u = partial derivative of L with respect to the I argument

L_p = partial derivative of L with respect to the II argument.

Remark : if L does not depend on p then E-L equations are just

$$\boxed{L_u = 0}$$

$$L(f) = \frac{1}{\pi} \int_{\mathbb{R}^2} (f(x) - g)^2 e^{-\frac{(x-y)^2}{2}} e^{-x^2} dx dy.$$

$$= \int L(f(x), x) dx$$

$$\begin{aligned}
 L(f) &= \frac{1}{\pi} \int_{\mathbb{R}^2} \left(f(x) + y^2 - 2yf(x) \right) e^{-(x-y)^2} e^{-x^2} dx dy \\
 &= \frac{1}{\pi} \int_{\mathbb{R}} \left[f(x) e^{-x^2} \left(\int_{\mathbb{R}} e^{-(x-y)^2} dy \right) + e^{-x^2} \left(\int_{\mathbb{R}} y^2 e^{-(x-y)^2} dy \right) \right. \\
 &\quad \left. - 2f(x) e^{-x^2} \left(\int_{\mathbb{R}} y e^{-(x-y)^2} dy \right) \right] dx
 \end{aligned}$$

$$\int_{\mathbb{R}} e^{-(x-y)^2} dy = \sqrt{\pi}$$

$$\int_{\mathbb{R}} y^2 e^{-(x-y)^2} dy \quad y-x = z \quad y = z+x$$

"

$$\begin{aligned}
 \int_{\mathbb{R}} (z+x)^2 e^{-z^2} dz &= \int_{\mathbb{R}} (z^2 + 2xz + x^2) e^{-z^2} dz \\
 &= \underbrace{\int_{\mathbb{R}} z^2 e^{-z^2} dz}_{\text{0}} + 2x \underbrace{\int_{\mathbb{R}} z e^{-z^2} dz}_{\text{0}} \\
 &\quad + x^2 \underbrace{\int_{\mathbb{R}} e^{-z^2} dz}_{\text{0}}
 \end{aligned}$$

III/III

$$\int_{\mathbb{R}} z^2 e^{-z^2} dz = -2 \underbrace{\int_{\mathbb{R}} z \cdot \frac{d}{dz} e^{-z^2} dz}_{\text{0}} = 2 \underbrace{\int_{\mathbb{R}} e^{-z^2} dz}_{\text{0}} = \sqrt{\pi}$$

$$\boxed{\int_{\mathbb{R}} y^2 e^{-(x-y)^2} dy = 2\sqrt{\pi} + x^2 \sqrt{\pi}}$$

$$\int_{\mathbb{R}} y e^{-(x-y)^2} dy = \int_{\mathbb{R}} (x+z) e^{-z^2} dz = x \sqrt{\pi}$$

$$\begin{aligned}
 L(f) &= \frac{1}{\pi} \int_{\mathbb{R}} \left[f(x) \sqrt{\pi} e^{-x^2} + e^{-x^2} (2\sqrt{\pi} + x^2 \sqrt{\pi}) - 2f(x) \sqrt{\pi} x e^{-x^2} \right] dx \\
 &\quad \mathcal{L}(f)
 \end{aligned}$$

So the stationarity conditions for L , which are the F-L equations

$$\partial_f = 0$$

$$2f(x)\sqrt{\pi}e^{-x^2} - 2\sqrt{\pi}x e^{-x^2} = 0$$

$$e^{-x^2} > 0$$

$$\boxed{f(x) = x}$$

We still need to check that this is a minimum.

In this example I gave you the **True** distribution of data and therefore I could find explicitly a minimum of the functional risk.

However, in general, we **NEVER** know the true distribution π but we only can aspire to know independent sample from π . Meaning that we usually know a dataset

$$\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N) : (x_i, y_i) \in \Omega\}.$$

The main idea then is to approximate π with

$$\pi_N = \left(\frac{1}{N} \sum_{i=1}^N \delta_{(x_i, y_i)} \right)$$

empirical distribution

What happens when you use π_N instead of π in the definition of the risk

$$\int_{\Omega} l(f(x), y) d\pi_N(x, y) = \frac{1}{N} \sum_{i=1}^N \int_{\Omega} l(f(x), y) d\delta_{(x_i, y_i)}$$

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

Empirical Risk

Functional Risk \rightarrow Empirical Risk

Functional Space $\rightarrow \mathbb{R}^n$ (finite dimensional space)

Basically we restrict the domain of optimization only to **PARAMETRIC FUNCTIONS**

P

$$\mathcal{P} = \{ g \in \text{dom}(L) : g(z) = f(z; w) \quad w \in \mathbb{R}^n \}.$$

We will choose \mathcal{P} to be the set of let us say all neural networks with some architecture.

$$\min_{w \in \mathbb{R}^n} L_N(f(\cdot; w))$$

We have shown that $f^*(z) = z$ is a stationary point (a candidate to be a minimum) of

$$L(f) = \frac{1}{\pi} \int_{\mathbb{R}^2} (f(z) - y)^2 e^{-(z-y)^2} e^{-z^2} dz dy$$



$$L(f^*) ?$$

$$L(f^*) = \frac{1}{\pi} \int_{\mathbb{R}^2} (z-y)^2 e^{-(z-y)^2} e^{-z^2} dz dy$$

$$= \frac{1}{\pi} \int_{\mathbb{R}} e^{-z^2} \left(\int_{\mathbb{R}} (z-y)^2 e^{-(z-y)^2} dy \right) dz$$

$$2\sqrt{\pi}$$

$$z-y = t$$

We have already computed

$$\int_{\mathbb{R}} z^2 e^{-z^2} dz = \frac{\sqrt{\pi}}{2}$$

$$L(f^*) = 2 \frac{\sqrt{\pi}}{\pi} \int_{\mathbb{R}} e^{-z^2} dz = \frac{1}{2}$$

$$L(f^*) = \frac{1}{2}$$

$$L(f) \geq \frac{1}{2} \quad \forall f \in \text{dom}(L)$$