

18/03/2025

Definition of the functional Risk for a function $f: X \rightarrow Y$

it can also happen that $Y \subseteq X$
in some problems.

X : space of labels
 \downarrow
"space of examples
perceptual info"

X = "space of handwritten characters"

$$Y = \{0, 1, 2, \dots, 9\}$$

$$f \times id(x, y) = (f(x), y)$$

Functional Risk

$$L(f) = \underset{\text{X} \times \text{Y}}{\mathbb{E}} l(f \times id) = \int_{X \times Y} l(f(x), y) d\pi(x, y)$$

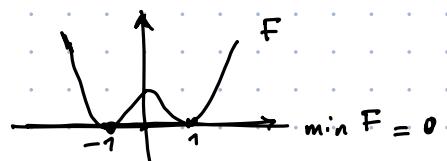
If you know the probability distributions (or measure) of your data, then what L measures is the mean error that the model $f: X \rightarrow Y$ does.
with respect to that distribution (π)

(*)

- Unless π is given (which in practice never happens) you cannot compute L and therefore you can't even try to minimize it.

(*) What is the goal of SML is to find the best model $f: X \rightarrow Y$, where best means a function that "minimizes" L

$$\bullet f^* \underset{\leftarrow}{=} \operatorname{argmin} L(f).$$



$$\operatorname{argmin} F = \{-1, 1\}$$

$$F: H \rightarrow \mathbb{R} \quad \min_{h \in H} F(h) = m$$

$$m: F(h) \geq m \forall h \in H$$

$$F(-1) = F(1) = 0$$

$$\operatorname{argmin} F = \{h \in H : F(h) = \min F\}.$$

Instead of L we can write and optimize the Empirical Risk

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

And the assumption is that (x_i, y_i) $i=1, \dots, N$ are independent samples from the same distribution π .

The other assumption that we typically do is that instead of minimizing L_N on a general functional space, we assume that f belongs to a parametric family of functions: \mathcal{P}

If you want you can think of \mathcal{P} as a family of NN that share some fixed architecture.

So at this point you can define

$$f^{(N)} \in \arg\min \{L_N(f) : f \in \mathcal{P}\}.$$

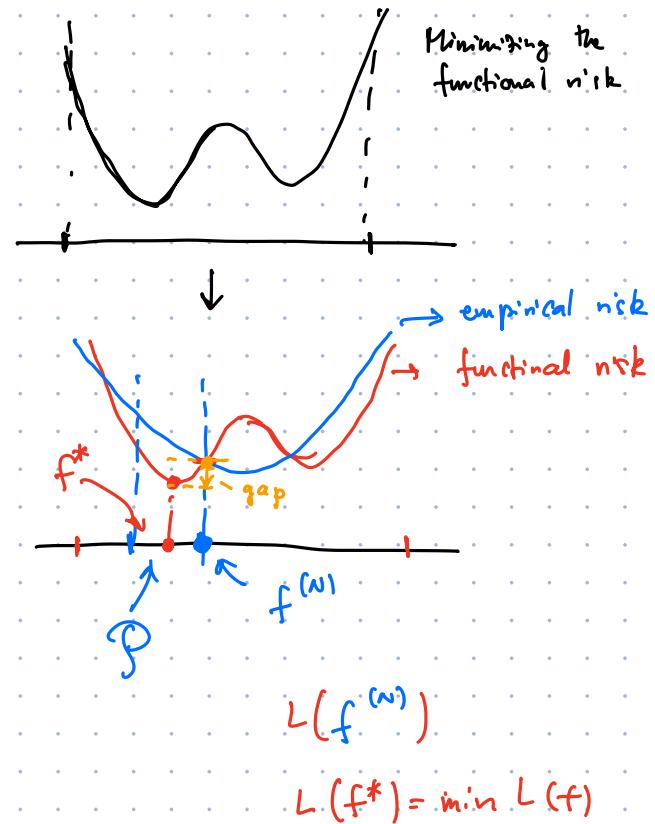
This is what is actually computed by optimization algorithms.

But what we wanted to do is instead compute $\inf \{L(f) : f \in \text{dom}(L)\}$.

$$\text{gap} = L(f^{(N)}) - \inf \{L(f) : f \in \text{dom}(L)\}$$

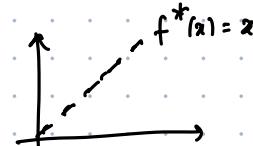
$$= (L(f^N) - \inf \{L(f) : f \in \mathcal{P}\})$$

$$+ (\inf \{L(f) : f \in \mathcal{P}\} - \inf \{L(f) : f \in \text{dom}(L)\})$$



Ex. (Continuing the exercise of Lecture 1).

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



We found a candidate to be a minimum $f^*(x) = x$.

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

We have to prove that $L(f) \geq \frac{1}{2} \quad \forall f \in \text{dom}(L)$

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

$$f(x) = e^x$$

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

$$+ \frac{1}{\pi} \int_{\mathbb{R}} y^2 e^{-(x-y)^2} e^{-x^2}$$

③

① $\geq 0 \leftarrow$

$$\begin{aligned} ② &= -\frac{2}{\pi} \int_{\mathbb{R}} f(x) e^{-x^2} \left(\int_{\mathbb{R}} y e^{-(x-y)^2} dy \right) dx \\ &= -\frac{2}{\pi} \int_{\mathbb{R}} f(x) e^{-x^2} x \sqrt{\pi} dx = -\frac{2}{\sqrt{\pi}} \int_{\mathbb{R}} f(x) x e^{-x^2} dx. \\ &\quad y-x=z \quad y=x+z \\ \int_{\mathbb{R}} y e^{-(x-y)^2} dy &= \int_{\mathbb{R}} (x+z) e^{-z^2} dz = x \sqrt{\pi} \end{aligned}$$

$$③ = \frac{1}{\pi} \int_{\mathbb{R}} y^2 e^{-(x-y)^2} e^{-x^2} = \frac{1}{2}$$

$$L(f) \geq -\frac{2}{\sqrt{\pi}} \underbrace{\int_{\mathbb{R}} f(x) x e^{-x^2} dx}_{?} + \frac{1}{2} \geq \frac{1}{2}$$

$$\begin{aligned} &\frac{1}{\sqrt{\pi}} \int_{\mathbb{R}} f^2(x) e^{-x^2} dx - \frac{2}{\sqrt{\pi}} \int_{\mathbb{R}} f(x) x e^{-x^2} dx \\ &= +\frac{1}{\sqrt{\pi}} \int_{\mathbb{R}} (f^2(x) - 2f(x)x) e^{-x^2} dx + \frac{1}{2} \\ &\quad \text{VI} \\ &\quad 0 \end{aligned}$$

————— * —————

ATTEMPT 2.

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

$$f(x) = x$$

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

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

$$\begin{aligned}
 &= \frac{1}{\pi} \int_{\mathbb{R}^2} (f(x) - z)^2 e^{-(x-y)^2} e^{-z^2} dx dy + \frac{1}{\pi} \int_{\mathbb{R}^2} (z-y)^2 e^{-(z-y)^2} e^{-z^2} dx dy \\
 &\quad (1) \qquad \qquad \qquad (2) \\
 &+ \frac{2}{\pi} \int_{\mathbb{R}^2} (x-y)(f(x) - z) e^{-(x-y)^2} e^{-z^2} dx dy \\
 &\quad (3)
 \end{aligned}$$

(1) ≥ 0 Because the integrand is positive.

$$y - z = z$$

$$\begin{aligned}
 (2) \quad &\frac{1}{\pi} \int_{\mathbb{R}^2} (z-y)^2 e^{-(z-y)^2} e^{-z^2} dx dy = \frac{1}{\pi} \int_{\mathbb{R}} e^{-z^2} \left(\int_{\mathbb{R}} (z-y)^2 e^{-(z-y)^2} dy \right) dz \\
 &= \frac{1}{\pi} \int_{\mathbb{R}} e^{-z^2} \left(\int_{\mathbb{R}} z^2 e^{-z^2} dz \right) dz = \frac{1}{\pi} \int_{\mathbb{R}} e^{-z^2} \frac{\sqrt{\pi}}{2} dz = \frac{\pi}{2\pi} = \frac{1}{2}.
 \end{aligned}$$

$$\begin{aligned}
 (3) \quad &\frac{2}{\pi} \int_{\mathbb{R}} (f(x) - z) e^{-z^2} \left(\int_{\mathbb{R}} (z-y) e^{-(z-y)^2} dy \right) dz \\
 &\qquad \qquad \qquad \text{||} \qquad \qquad \qquad \text{even function} \\
 &\qquad \int_{\mathbb{R}} \overset{\rightarrow}{e^{-z^2}} dz = 0 \\
 &\qquad \qquad \qquad \text{odd function}
 \end{aligned}$$

$$(1) + (2) + (3) \geq 0 + \frac{1}{2} + 0 = \frac{1}{2} \Rightarrow L(f) \geq \frac{1}{2} \text{ for } f \in \text{dom}(L)$$

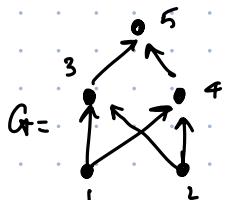
Since we know that $L(f^*) = \frac{1}{2}$ with $f^*(x) = z$, then f^* is a minimizer of L and $\frac{1}{2}$ is the minimum.

FEED FORWARD NEURAL NETWORKS.

A feed forward neural network is a parametric function [its parameters are called "weights"] whose computational structure is based on a dag.

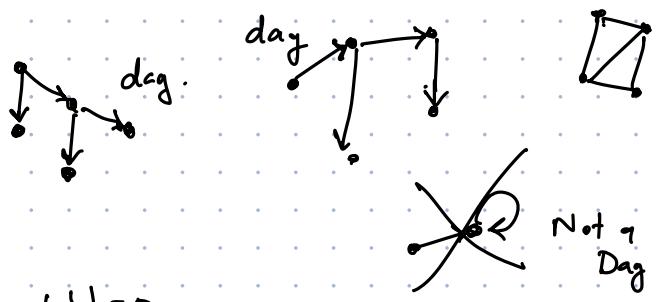
$$c_p : \mathbb{R}^p \rightarrow \mathbb{R}^d$$

1. You have a day $G = (V, A)$



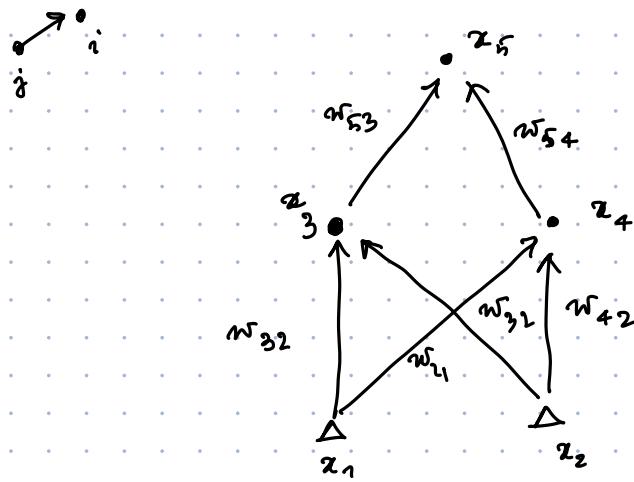
$$|\{j : pa(j) = \emptyset\}| = p$$

$$|\{j : ch(j) = \emptyset\}| = d$$



2. You have a map that associates to each vertex $i \xrightarrow{\psi} z_i \in \mathbb{R}$
 z_i is called "output of neuron i "

$$A \ni (j \rightarrow i) \mapsto w_{ij} \in \mathbb{R}$$



$$\begin{aligned} pa(1) &= \emptyset \\ pa(2) &= \emptyset \\ pa(3) &= \{1, 2\} \\ pa(4) &= \{1, 2\} \\ pa(5) &= \{3, 4\} \end{aligned}$$

3. I give the following rule on how to compute the outputs

$$z_i = \sigma \left(\sum_{j \in pa(i)} w_{ij} z_j \right) \quad \forall i \in V : pa(i) \text{ is non-empty.}$$

$$pa(i) = \{ j \in V : j \rightarrow i \in A \}$$

$$z_i = u_i \quad \text{if } pa(i) = \emptyset.$$

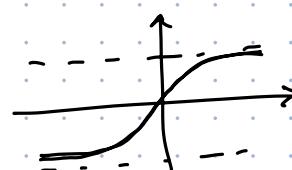
We call $u = (u_1, \dots, u_p)$ the input to the neural network.

$\sigma: \mathbb{R} \rightarrow \mathbb{R}$ is called "activation function" and:

1. Historically σ is a bounded increasing function

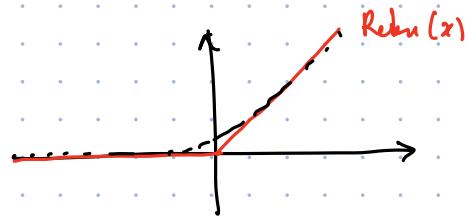
(sigmoidal-like) $\sigma' > 0$

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \text{sigmoid}$$



$$\sigma(x) = \tanh(x)$$

2. More recently another common activation function is the ReLU positive part



$$\text{ReLU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

$$(x)^+ , x \cdot [x \geq 0]$$

Example 1



$$x_1 = u_1$$

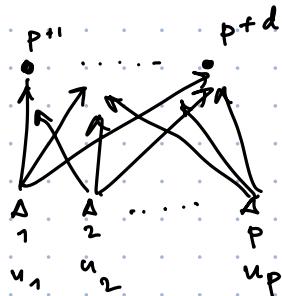
$$x_2 = \sigma(w x_1) = \sigma(w u_1)$$

$$f(u; w) = \sigma(w u)$$

$$u \mapsto f(u; w)$$

parameter or weight of the network.

Example 2



$$\begin{aligned} \text{if } j \geq p+1 \\ x_j &= \sigma\left(\sum_{i=1}^p w_{ij} u_i\right) \quad j \in \{1, \dots, d\} \\ &= \sigma\left(\sum_{i=1}^p w_{ij} u_i\right) \quad (*) \end{aligned}$$

$$W \in \mathbb{R}^{d \times p}$$

$$\sum_{j=1}^d w_{ij} u_i = W u$$

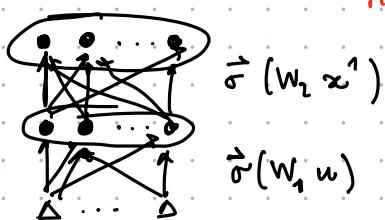
Therefore (*) is equivalent to

$$x = \vec{\sigma}(W u) \quad \vec{\sigma}(x_1, \dots, x_n) = (\sigma(x_1), \dots, \sigma(x_n)).$$

$$f(u; W) = \vec{\sigma}(W u) \quad f(\cdot; W) : \mathbb{R}^p \rightarrow \mathbb{R}^d$$

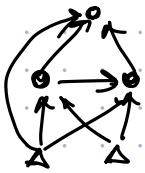
Example 3.

Multi Layer Perceptron



$$f(u, w) = \vec{\sigma}\left(W_2 \vec{\sigma}(W_1 u)\right)$$

$$\underbrace{W_2, W_1, u}_M$$



Backpropagation

PD book

Goeffrey Hinton

Suppose you want to compute the gradient of a function $\varphi \in C^1(\mathbb{R}^p; \mathbb{R})$

$$\nabla \varphi \approx \left(\frac{\varphi(w_1 + h, w_2, \dots, w_p) - \varphi(w_1, w_2, \dots, w_p)}{h}, \frac{\varphi(w_1, w_2 + h, \dots, w_p) - \varphi(w_1, w_2, \dots, w_p)}{h}, \dots, \frac{\varphi(w_1, w_2, \dots, w_p + h) - \varphi(w_1, \dots, w_p)}{h} \right)$$

$\varphi: \mathbb{R}^p \rightarrow \mathbb{R}$

Computing φ is $O(p)$

In general computing $\nabla \varphi$ numerically requires $O(p^2)$ which is a problem

Relatively modern NN has 10^8 parameters, so to compute gradients we would require to do $O(10^{16})$ operations, which is computationally unfeasible.

(RNN)

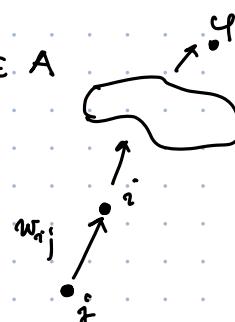
It turns out that, instead if φ is not a generic map but it is a FNN you can improve the $O(p^2)$ bound and reduce it to $O(p)$ which is optimal. $[\Omega(p)]$

Suppose for simplicity that $\varphi: \mathbb{R}^p \rightarrow \mathbb{R}$

The parameters of φ are w_{ij} for $i \leftarrow i \in A$

$$\frac{\partial \varphi}{\partial w_{ij}} = \boxed{\frac{\partial \varphi}{\partial a_i}} \frac{\partial a_i}{\partial w_{ij}}$$

δ_i - error of neuron i



$$\frac{\partial a_i}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_{k \in \text{pa}(i)} w_{ik} x_k$$

$$= \sum_{k \in \text{pa}(i)} \delta_{kj} x_k = x_j$$

$$a_i = \sum_{j \in \text{pa}(i)} w_{ij} x_j$$

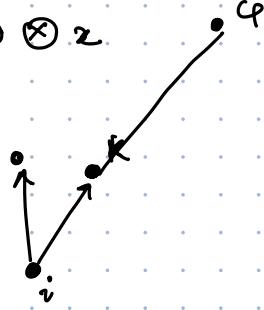
$$\delta_{kj} = \begin{cases} 1 & \text{if } k=j \\ 0 & \text{otherwise} \end{cases}$$

Kronecker δ symbol

So we found that

$$\frac{\partial \varphi}{\partial w_{ij}} = \delta_i z_j$$

$$\nabla \varphi = \delta \otimes z$$



$$\begin{aligned} \delta_i &= \frac{\partial \varphi}{\partial a_i} = \sum_{k \in \text{ch}(i)} \frac{\partial \varphi}{\partial a_k} \frac{\partial a_k}{\partial a_i} \\ &= \sum_{k \in \text{ch}(i)} \delta_k \frac{\partial a_k}{\partial a_i} \end{aligned}$$

$$z_i = \sigma \left(\sum_{j \in \text{pa}(i)} w_{ij} z_j \right)$$

$$x_i = \sigma(a_i)$$

$$a_k = \sum_{m \in \text{pa}(k)} w_{km} x_m = \sum_{m \in \text{pa}(k)} w_{km} \sigma(a_m)$$

$$\frac{\partial a_k}{\partial a_i} = \sum_{m \in \text{pa}(k)} w_{km} \sigma'(a_m) \delta_{mi} = w_{ki} \sigma'(a_i)$$

$$\delta_i = \sum_{k \in \text{ch}(i)} \delta_k w_{ki} \sigma'(a_i)$$

$$\delta_i = \sigma'(a_i) \sum_{k \in \text{ch}(i)} \delta_k w_{ki}$$

when i is an output neuron

$$\delta_i = \frac{\partial \varphi}{\partial a_i} = \sigma'(a_i)$$

$$x_i = \sigma \left(\sum w_{ij} x_j \right) \rightarrow x = \sigma(Wz) \quad \downarrow O(p)$$

$$\delta_i = \sigma'(a_i) \sum \delta_k w_{ki} \rightarrow \delta \propto W^T \delta \quad O(p) \quad W^T \text{ is transpose}$$