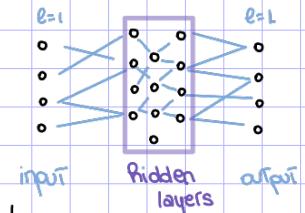


8 HOW TO COMPUTE THE WEIGHT?

8.1 FFNN recall



All the layers are fully connected with σ^e the activation function of the layer l

$$y_i^e = \sum_{j=1}^{N_{l-1}} \sigma^e(w_{ij}^e x_j^{e-1}) \quad \forall i \in \{1, \dots, N_l\} \quad \forall e \in \{2, \dots, L\} \Rightarrow \text{The non-linearity is contained in } \sigma^e$$

↓
neuron i
↓
 $y := \text{output}$

↓ Rewrite in compact form $y^e = f_e(w^e, y^{e-1})$ $y^{e-1} = z^e$ in which we consider the input of layer j is the output of layer $j-1$

Imagine to know the weight w (you already solve the optimization problem), the algorithm to the evaluation process is

ALGORITHM
to
compute
the
weight
 $y' = z$ → input
for $e=2, \dots, L$
 $z^e = w^e y^{e-1}$
 $y^e = \sigma^e(z^e)$
end
 $y = y^L$ → output

↳ The optimization problem is the offline phase that is NOT fast

↳ This is similar to the online phase that is very fast

How to compute the weight? TRAINING PHASE: we have different strategies

- ① Supervised learning: we have known snapshot $\{(x^n, y^n)\}_{n=1, \dots, N}$ that minimize the loss over input, n-layer
↳ the network is trained using example of input & output given by a teacher \Rightarrow goal: learn the rule to map $x \mapsto y$
- ② Unsupervised learning: we do NOT have snapshot \rightarrow PINN: Physics Informed Neural Network \rightarrow See next lessons
↳ the network is trained using only inputs, no outputs are given \Rightarrow goal: discover hidden patterns in the data
- ③ Reinforcement learning: You give information and reward and the net learn
↳ the network is trained in a dynamic environment to perform a fixed goal: while training, the environment gives feedback (positive & negative reward) and the network try to maximizes

8.2 Supervised learning: we want to approximate the target function

Consider the function $f: K \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ that is the function in the hidden layer that produce the output
↳ unknown $x \mapsto y$

Problem - Given FFNN fixed $(L, N_l, \sigma^e)_{e=2, \dots, L}$

Find w^e st.

$y = f(w^e, x) \rightarrow$ where f is an approximation of f

IDEA: we want to find w so the f approximate f (the hidden function) that we DO NOT have, but we have the model (FFNN) so we know \hookrightarrow

We suppose to know some snapshots $(x^n, y^n)_{n=1, \dots, N}$ and we select the loss function n -th loss error

$$E_n = \frac{1}{2} \|y^n - \bar{y}^n\|_2^2 = \frac{1}{2} \|f(w^e, x^n) - \bar{y}^n\|_2^2$$

↓
predicted value ↓ real value

The total loss error is $E = \sum_{n=1}^N E_n(w)$ and we want to find a local minimum

we try to get a minimization problem
we can have a local minimum
could be NOT unique

We need a minimization method \rightarrow to solve the minimization problem

a) GRADIENT DESCENT METHOD: based on the gradient of the loss function

Find $\min_{w \in \mathbb{R}^B} E(w)$ where $\Sigma_a := \{w \in \mathbb{R}^B : E(w) = a\}$ $\forall a \in \mathbb{R}$ are the level curve

To go in a curve with lower value we use the gradient \Rightarrow Iterative step: $w_i = w^{i-1} - \alpha \nabla E(w^{i-1})$ $\alpha > 0 \in \mathbb{R}^+$

BUT this works well iff the dimension is NOT so high, because we have $\nabla E(w) = \sum_{n=1}^N \nabla E_n(w)$ so it depends on N = number of snapshots that has to be high to learn properly

↓
Garantees only to reach a local minimum, not a global one

b) SOLUTION? STOCHASTIC GRADIENT DESCENT ALGORITHM: compute the gradient for only some snapshots

$$w_i = w_{i-1} - \alpha \frac{1}{M} \sum_{m=1}^M \nabla E_m(w_{i-1}) \quad M \ll N$$

I take $M \neq$ snapshots, chosen in a random way, M is called mini-batch

Example - \rightarrow Real gradient $\nabla E = \sum_{n=1}^N \nabla E_n$

\rightarrow Stochastic gradient set $S_M = \{n_1, \dots, n_M\} \quad n_i \in \{1, \dots, N\}$ yields $\tilde{\nabla} E = \frac{1}{M} \sum_{n \in S_M} \nabla E_n$

mini-batch: selected randomly at each epoch, is a subset of all the snapshots

Each iteration of the stochastic gradient descent is called epoch

↓ Adaptive Moment Estimation

ADAM ALGORITHM: improvement of the stochastic in which exploit the previous computation
all the previous gradient are called momentum of E

This is the standard used

L-BFGS: Second order algorithm \rightarrow exploit the Hessian, so it is more precise \rightarrow following quasi-Newton method

↓ Limited Memory Broyden-Fletcher-Goldfarb-Shanno

c) How to compute $\nabla E(w)$? [BACKPROPAGATION]

Consider $E(w) = \frac{1}{M} \sum_{m=1}^M \frac{1}{2} \|y^m(w) - \bar{y}^m\|^2$, where i.e. $w = \begin{pmatrix} w_{11} \\ w_{12} \\ \vdots \\ w_{24} \end{pmatrix}$ if $W^1 \in \mathbb{R}^{3,2}$ $W^2 \in \mathbb{R}^{2,4}$ perception

From definition, we have $\nabla E(w) = \frac{1}{2} \cdot 2 \left(y^m(w) - \bar{y}^m \right) \frac{\partial y^m(w)}{\partial w}$, where $y^m(w) = f(w, x^m)$
we know $y^m(w) = f(w, x^m)$
we do not know $\frac{\partial y^m(w)}{\partial w}$

HOW TO COMPUTE THAT?

We know everything, but not the partial derivative $f: \mathbb{R}^S \rightarrow \mathbb{R}^M \quad \frac{\partial f}{\partial w} = \left\{ \frac{\partial f_j}{\partial w_{pq}} \right\}$ \rightarrow we need to compute the Jacobian

Imagine to fix the layers $l-1, l$ that are fully connected, and select neuron p and q
 $y^l = \phi_{l-1}(w^l \cdot y^{l-1})$



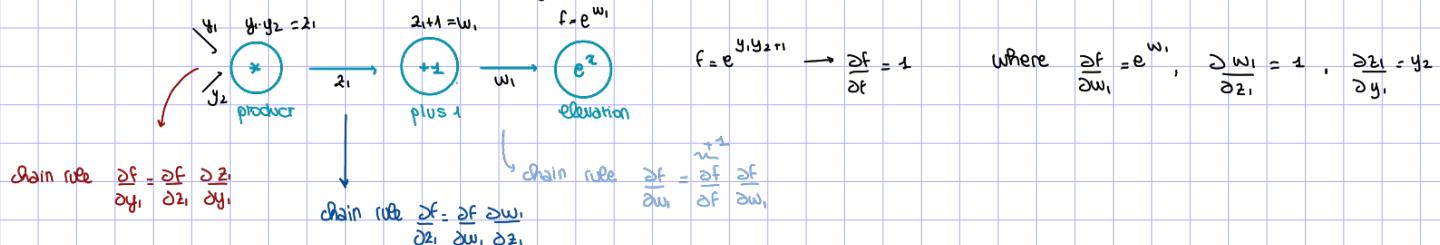
So we can use the chain rule to compute the partial derivative

$$\frac{\partial f_j}{\partial w_{pq}} = \frac{\partial f_j}{\partial y^p} \cdot \frac{\partial y^p}{\partial w^l_{qp}}$$

because I derive wrt w^l_{qp} so all the other elements becomes ϕ

we do not know i.e. if $f_j = e^{y_1 y_2 + 1} \Rightarrow \frac{\partial f_j}{\partial y_1} = e^{y_1 y_2 + 1} \cdot y_2 \quad \frac{\partial f_j}{\partial y_2} = e^{y_1 y_2 + 1} \cdot y_1$

BUT how to compute that in NN logic?



Hence, go back we have all the element and we can compute everything

$$\frac{\partial f}{\partial y_1} = \frac{\partial f}{\partial z_1} \cdot \frac{\partial z_1}{\partial y_1} \quad \frac{\partial f}{\partial y_2} = \frac{\partial f}{\partial z_2} \cdot \frac{\partial z_2}{\partial y_2}$$

Starting to the end, we know everything so we can compute the derivative!

Now, we can apply this IDEA to compute $\frac{\partial f_j}{\partial y^p} = \frac{\partial f_j}{\partial y_h^H} \cdot \frac{\partial y_h^H}{\partial y^p}$, where by definition

$$y_h^H = \text{softmax} \left(\sum_{q=1}^{N^H} w_{hq}^H y_q^H \right)$$

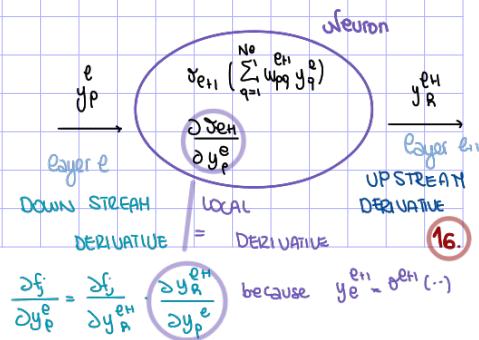
$$f(w, x^m) = y^m$$

BUT! Know we do NOT know this element

We can put in the case $l+1=L \Rightarrow \frac{\partial f}{\partial y^L} = \frac{\partial y^L}{\partial y^L} = I$ identity

So we can do all the steps: $\frac{\partial f}{\partial y^{l-1}} \rightarrow \frac{\partial f}{\partial y^{l-2}} \rightarrow \dots \rightarrow \frac{\partial f}{\partial y^1} \rightarrow \frac{\partial y^1}{\partial x} \Rightarrow \text{zoom}$

$$\text{And } \frac{\partial y^1}{\partial x} = \frac{\partial}{\partial x} (O_1(w \cdot x)) = O_1'(w \cdot x) w^1$$



$$\frac{\partial f}{\partial y^p} = \frac{\partial f}{\partial y^H} \cdot \frac{\partial y^H}{\partial y^p} \quad \text{because } y^H = \text{softmax}(\dots)$$

To see more detail on UNIVERSAL APPROXIMATION THEOREMS see lecture notes Lesson 8