

# Introduction to Artificial Neural Networks



Filippo Gatti<sup>†</sup>  1

<sup>†</sup>Laboratoire de Mécanique Paris-Saclay UMR 9026  
Université Paris-Saclay, CentraleSupélec, ENS Paris-Saclay, CNRS



# Outline

- ① Learning framework
- ② Optimize the statistical model
- ③ Create a statistical model
- ④ Train a statistical model
- ⑤ Neural networks for time series

# Learning framework

# General framework

Supervised learning : labelled database

$$\mathcal{D}_{XY} = \{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^N, (\mathbf{x}_i, \mathbf{y}_i) \sim p \text{ and i.i.d. } p(\mathcal{D}_{XY}) = \prod_{i=1}^N p(\mathbf{x}_i, \mathbf{y}_i)$$

Statistical model : parametric probability distribution

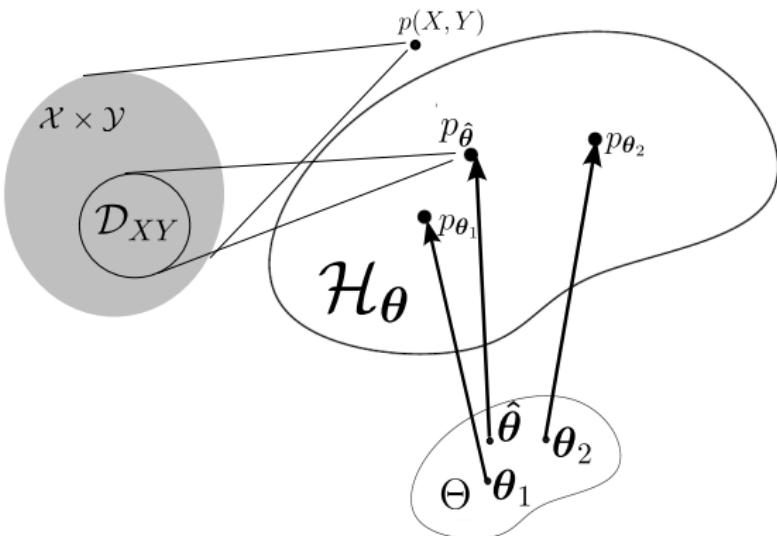
$$\mathcal{H}_\Theta := \{p_\theta, \theta \in \Theta \subset \mathbb{R}^{d_\Theta}\}$$

- Injective :  $\exists! \theta^* \iff \theta_1 = \theta_2 \implies p_{\theta_1} = p_{\theta_2}$
- Regular :  $\exists g \in L^1_{loc}(\mathcal{X}) : \left| \frac{\partial p_\theta}{\partial \theta_k} \right| \leq g$   $\mu$ -a.e. in open neighborhood  $U \in N(\mathbf{x})$
- ?  $\exists \theta^* \in \Theta : p_{\theta^*} = p$

Idea : find  $\hat{\theta} \approx \theta^*$  minimizing the *Negative Log-Likelihood*  $NLL$

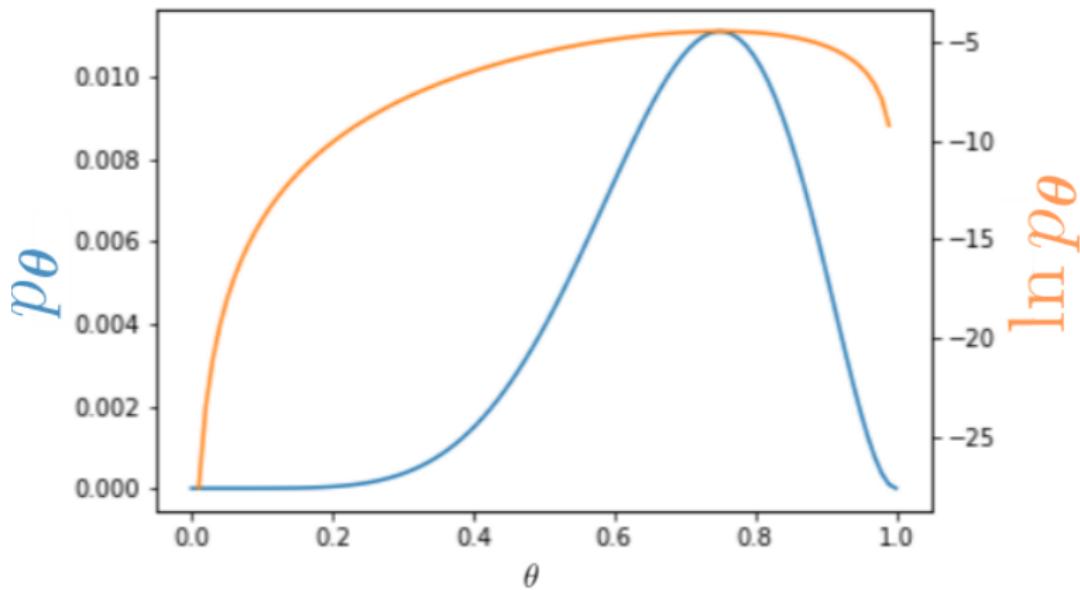
$$\hat{\theta}(\mathcal{D}_{XY}) = \arg \min_{\theta \in \Theta} -\ln p_\theta(\mathbf{y}|\mathbf{x}) = \arg \min_{\theta \in \Theta} -\sum_{i=1}^N \ln p_\theta(\mathbf{y}_i|\mathbf{x}_i), \quad p_{\hat{\theta}} \in \mathcal{H}_{\hat{\theta}}$$

# Visualize the learning framework



Minimizing  $\mathbb{E} [NLL] \implies$  minimizing Kullback-Leibler divergence  $D_{KL}(p||p_\theta) = \mathbb{E}_{x \sim p} \left[ \ln \frac{p}{p_\theta} \right] < +\infty$

$$\min_{\theta \in \Theta} \mathbb{E}_{x \sim p} \left[ -\ln p_\theta (\mathbf{x}) \right] = \min_{\theta \in \Theta} \mathbb{E}_{x \sim p} \left[ -\ln \frac{p_\theta}{p} \right] \underbrace{- \mathbb{E}_{x \sim p} [\ln p]}_{\geq 0} \geq \min_{\theta \in \Theta} D_{KL}(p||p_\theta)$$



**Figure** – Bernoulli distribution  $p_\theta(\mathcal{D}_X) = \prod_{i=1}^N \theta^{1_{x_i=1}} (1 - \theta)^{1_{x_i=0}}$

[stats.stackexchange.com/questions/486007/measuring-predictive-uncertainty-with-negative-log-likelihood-nll](https://stats.stackexchange.com/questions/486007/measuring-predictive-uncertainty-with-negative-log-likelihood-nll)

# Negative Log-Likelihood

For a “true” Gaussian distribution  $p(Y|X) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{y_i|x_i-\mu|^2}{2\sigma^2}}$  with  $\theta = (\mu, \sigma)$ :

$$\frac{\mathcal{NLL}(\mathcal{D}_{XY})}{N} = \frac{1}{2} \ln(2\pi) + \ln \sigma + \frac{1}{2\sigma^2 N} \sum_{i=1}^N (y_i|x_i - \mu|^2)$$

Idea : use a model  $y_i|x_i = h_\theta(x_i)$

$$\frac{\mathcal{NLL}(\theta; \mathcal{D}_{XY})}{N} = \frac{1}{2} \ln(2\pi) + \ln \sigma + \frac{1}{2\sigma^2 N} \sum_{i=1}^N (h_\theta(x_i) - \mu)^2$$

with  $\frac{1}{N} \sum_{i=1}^N (h_\theta(x_i) - \mu)^2 = \text{MSE}$

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \frac{\mathcal{NLL}(\theta; \mathcal{D}_{XY})}{N} \rightarrow h_{\hat{\theta}} = \mu$$

# The Shannon's approach

- Teacher (=source) sends a message to a pupil (=receiver)
- The message is meaningful if the receiver had no a priori knowledge of it
- Deterministic messages bare zero information
- Shannon's information = uncertainty or surprise associated to a random variable

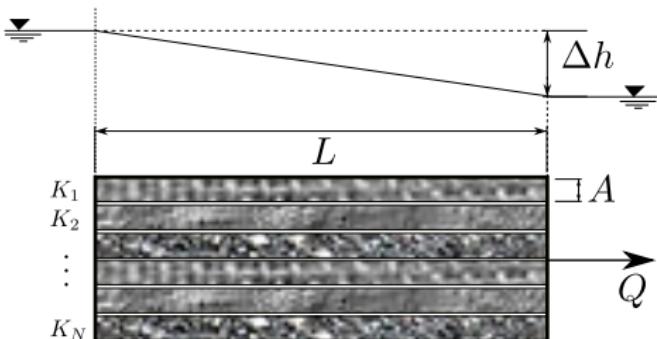
$$\mathbb{I}(\omega) = f(P(\omega)) = -\log_b P(\omega) = \frac{-\ln P(\omega)}{\ln b} \geq 0, \quad b > 1$$

quantifying the “surprise”

- $\mathbb{I}(\omega) = 0$  if  $P(\omega) = 1$
- $\mathbb{I}(\omega) = \infty$  if  $P(\omega) = 0$
- $\mathbb{I}(\omega_1 \cap \omega_2) = \mathbb{I}(\omega_1) + \mathbb{I}(\omega_2)$  with  $P(\omega_1 \cap \omega_2) = P(\omega_1) \cdot P(\omega_2)$

$$\mathbb{H}(p(\omega_1), \dots, p(\omega_N)) = \cdot \sum_{\omega_i \in \omega}^n p(\omega_i) \cdot \mathbb{I}(\omega_i) = \mathbb{E}_{\omega_i \sim p_i} [\mathbb{I}(\omega)]$$

# Example of Shannon's entropy

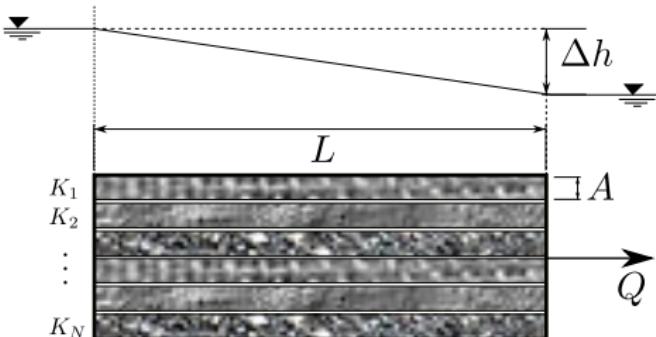


**Figure** – Darcy's law in a set of  $N$  parallel channels with random hydraulic conductivity  $K_i$ .

- Total discharge  $Q = \left( \sum_{i=1}^N K_i \right) \frac{\Delta h \cdot A}{L}$
- Collection of values  $\mathcal{D}_K = \{K_1, \dots, K_b\}$ ,  $\text{card}(\mathcal{D}_K) = b$  and  $p(\mathcal{D}_K) = (p_1, \dots, p_b)$
- Pick  $N$  independent channels with  $p_i = \frac{1}{b} : p\left(\{K_i\}_{i=1}^N\right) = b^{-N}$
- Entropy = average number of channel that deliver  $Q$

$$\mathbb{E}_{n \sim p}[N] = -\mathbb{E}_{n \sim p} \left[ \log_b p\left(\{K_i\}_{i=1}^N\right) \right] = H(p_1, \dots, p_N)$$

# Example of Shannon's entropy



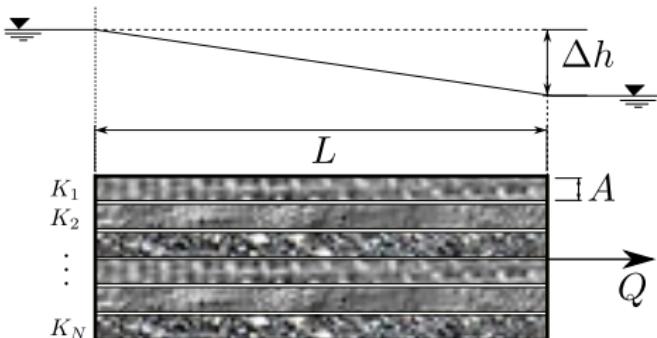
**Figure** – Darcy's law in a set of  $N$  parallel channels with random hydraulic conductivity  $K_i$ .

- Pick  $N$  independent channels with empirical frequency  $f_i = \frac{b_i}{N}$ ,  $\sum_{i=1}^N f_i = 1$

$$p\left(\{K_i\}_{i=1}^N\right) = \prod_{i=1}^N p_i^{b_i} = b^{N\left(\sum_{i=1}^N f_i \log_b p_i\right)} = e^{-N \cdot (\mathbb{D}_{KL}((f_1, \dots, f_N) || (p_1, \dots, p_N)) + \mathbb{H}(f_1, \dots, f_N) \ln b)}$$

- $b_i \sim \mathcal{B}(N, p_i) = \binom{N}{b_i} p_i^{b_i} (1 - p_i)^{N-b_i}$ ,  $\mathbb{E}_{b_i \sim \mathcal{B}(N, p_i)} [B_i] = N \cdot p_i \implies \mathbb{E}_{f_i} [f_i] = p_i$
- $f_i \xrightarrow[N \rightarrow +\infty]{} p_i \implies p\left(\{K_i\}_{i=1}^N\right) \xrightarrow[N \rightarrow +\infty]{} e^{-N \cdot \mathbb{H}(p(K)) \ln b}$

# Example of Shannon's entropy



**Figure** – Darcy's law in a set of  $N$  parallel channels with random hydraulic conductivity  $K_i$ .

- possible combinations (with repetitions) of values of hydraulic conductivity  

$$N_b = \frac{(N+b-1)!}{b! \cdot (N-1)!} \implies Q = \left( \sum_{i=1}^N f_i K_i \right) \frac{\Delta h \cdot A}{L}$$
- $Q$  is attained  $N_Q = \frac{N!}{\prod_{i=1}^N b_i!} \approx \frac{e^{\mathbb{H}(f_1, \dots, f_N)}}{\sqrt{(2\pi N)^{N-1} \prod_{i=1}^N f_i}}$  unique times with probability  

$$p(Q) = N_Q \cdot p\left(\{K_i\}_{i=1}^N\right) \approx \frac{e^{-N \cdot \mathbb{D}KL((f_1, \dots, f_N) \parallel (p_1, \dots, p_N))}}{\sqrt{(2\pi N)^{N-1} \prod_{i=1}^N f_i}} \xrightarrow[N \rightarrow +\infty]{} 1 \approx N_Q \cdot e^{-N \cdot \mathbb{H}(p(K)) \ln b}$$
- $$\frac{\ln N_Q}{N \ln b} \approx \mathbb{H}(p(K))$$

# MaxEnt principle

## Gibbs-Boltzmann theorem

$$\begin{cases} \max_{q \in \mathcal{H}} \mathbb{H}_d(q) = \max_{q \in \mathcal{H}} \left( - \int_{\mathcal{X}} \ln q(\mathbf{x}) \cdot q(\mathbf{x}) d\mathbf{x} \right) \\ c_k(q) = \mu_{y_k} - \int_{\mathcal{X}} y_k(\mathbf{x}) \cdot q(\mathbf{x}) d\mathbf{x} = 0, \quad c_k : \mathbb{R}^n \rightarrow \mathbb{R}, 1 \leq k \leq K \end{cases}$$

if it exists it is defined as  $p_{\theta} \in \mathcal{H}$ , which reads :

$$p_{\theta}(\mathbf{x}) = \arg \max_{q \in \mathcal{H}, \mathbf{c}: \mathbf{c} = \mathbf{0}} \mathbb{H}_d(q) = \frac{e^{-\sum_{k=1}^K y_k(\mathbf{x}) \cdot \theta_k}}{Z}, \quad \mathbb{H}_d(p_{\theta}) \geq \mathbb{H}_d(p)$$

$$\text{with } Z = \int_{\mathcal{X}} e^{-\sum_{k=1}^K \theta_k \cdot y_k(\mathbf{x})} \cdot \mu(d\mathbf{x})$$

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \max_{p_{\theta} \in \mathcal{H}_{\Theta}} \mathbb{H}_d(p_{\theta} | \mathbf{Y} | \mathbf{X})$$

$$\mathcal{H}_{\theta} \ni p_{\theta}(x) = \exp(-\langle y(x), \theta \rangle) / Z$$

## Log-Likelihood

$$\ell(\theta; x) = \ln p_{\theta}(x) = -\ln Z - \langle y(x), \theta \rangle$$

$$s_{\theta}(x) = \nabla_{\theta} \ln p_{\theta} = \mathbb{E}_{x \sim p_{\theta}} [y(X)] - y(x) = \mu_y - y(x)$$

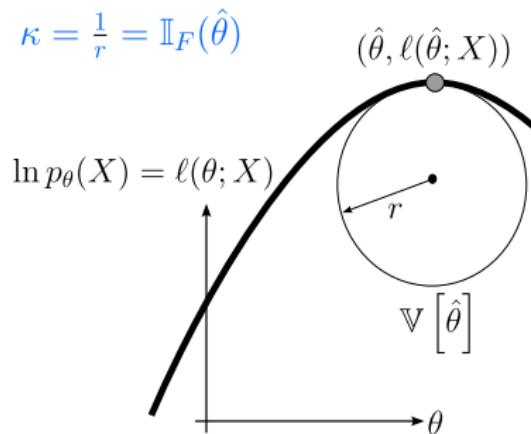
$$\mathbf{H}_{\ell}(\theta; X) = -\mathbb{C}_{x \sim p_{\theta}} (y(X))$$

- Unbiased score :  $\mathbb{E}_{x \sim p_{\theta}} [s_{\theta}(X)] = \mathbf{0} \longrightarrow$  stationarity
- Fisher Information Matrix (under regularity conditions)  
 $\mathbb{I}_F(\theta; X) = -\mathbb{E}_{x \sim p_{\theta^*}} [\mathbf{H}_{\ell}(\theta; X)] = \mathbb{C}_{x \sim p_{\theta}} (s_{\theta}(X))$
- $\exists (\hat{\theta}_n)_{n \in \mathbb{N}}$  convergent in  $\mathcal{N}(\mathbf{0}, \mathbb{I}_F^{-1}(\theta^*; X)) \in \mathcal{H}_{\theta}$

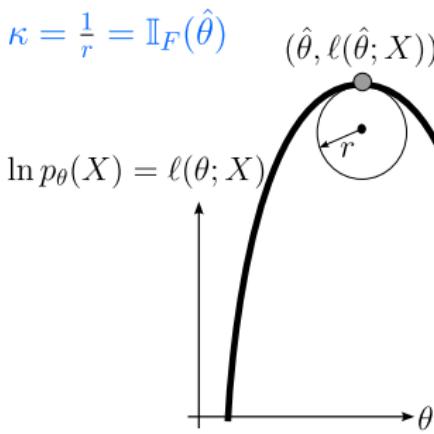
$$\theta^* = \lim_{n \rightarrow \infty} \hat{\theta}_n = \hat{\theta} = \arg \min_{\theta \in \Theta} \max_{p_{\theta} \in \mathcal{H}_{\Theta}} \mathbb{H}_d(p_{\theta}(Y|X))$$

# Sketching the Fisher Information

$$\kappa = \frac{1}{r} = \mathbb{I}_F(\hat{\theta})$$

**a**

$$\kappa = \frac{1}{r} = \mathbb{I}_F(\hat{\theta})$$

**b**

**Figure –** (a) Small FI/curvature (flat peak)  $\rightarrow$  large variance (small accuracy); (b) Large FI/curvature (sharp peak)  $\rightarrow$  small variance (high accuracy). Inspired by F. Nielsen presentation on Introduction to Geometry of Information

$$\mathbb{E}_{\mathbf{x} \sim p_{\boldsymbol{\theta}}} [\ell(\boldsymbol{\theta}; \mathbf{X})] = \mathbb{E}_{\mathbf{x} \sim p_{\hat{\boldsymbol{\theta}}}} [\ell(\hat{\boldsymbol{\theta}}; \mathbf{X})] + \left\langle \mathbb{E}_{\mathbf{x} \sim p_{\boldsymbol{\theta}}} [s_{\boldsymbol{\theta}}^0(\mathbf{X})], (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right\rangle - \frac{1}{2} \left\langle \boldsymbol{\theta} - \hat{\boldsymbol{\theta}}, \mathbb{I}_F(\boldsymbol{\theta}^*; \mathbf{X}) \cdot (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}) \right\rangle$$

# Empirical loss minimization

## Empirical Loss Minimization ( $\mathcal{P}$ )

Given a labeled dataset  $\mathcal{D}_{XY} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \in \mathcal{X} \times \mathcal{Y} \subset \mathbb{R}^{d_X+d_Y}$

Find  $\mathbf{h}_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ ,  $\mathbf{h}_\theta \in \mathcal{H}_\theta$  such that :  $\mathbf{h}_{\hat{\theta}}(\mathcal{D}_{XY}) = \arg \min_{\mathbf{h}_\theta \in \mathcal{H}_\theta} L_{\mathcal{D}_{XY}}(\mathbf{h}_\theta)$

$$L_{\mathcal{D}_{XY}}(\mathbf{h}_\theta) = \frac{1}{N} \sum_{i=1}^N \ell(\mathbf{h}_\theta(\mathbf{x}_k), \mathbf{y}_k), \quad \ell(\mathbf{y}, \mathbf{y}) = 0$$

## Empirical Loss

- $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^+$  : “distance” between “true” label and prediction
- $\ell(y, h_\theta(x)) = -\ln p_\theta(y|x) = \mathcal{NLL}(y, h_\theta(x))$
- $L_{\mathcal{D}_{XY}}$  is an approximation of the “true” loss  $L = \mathbb{E}_{(x,y) \sim p} [\ell(\mathbf{h}_\theta(\mathbf{X}), \mathbf{Y})]$ ,  
but  $p(\mathbf{X}, \mathbf{Y})$  is not known *a priori*

Loss landscape : <https://losslandscape.com/explorer>

# Optimize the statistical model

Existence of minimizing sequence (Gibbs-Boltzmann theorem)

$$\exists (\hat{\theta}_n)_{n \in \mathbb{N}} \text{ convergent in } \mathcal{N}(\mathbf{0}, \mathbb{I}_F^{-1}(\boldsymbol{\theta}^*; \mathbf{X})) \in \mathcal{H}_{\boldsymbol{\theta}}$$

“Delta rule”

The quest for  $\hat{\boldsymbol{\theta}}$  is iteratively conducted, following the direction of  $-\nabla_{\boldsymbol{\theta}} L_{\mathcal{D}_{XY}}$ , according to the following update rule from iteration  $i$  to iteration  $+i+1$  (the so called *delta rule* :

$$\boldsymbol{\theta}^{(i+1)} = \boldsymbol{\theta}^{(i)} - \eta^{(i)} \nabla_{\boldsymbol{\theta}} L_{\mathcal{D}_{XY}}(\boldsymbol{\theta}^{(i)}), \quad \eta^{(i)} \in \mathbb{R}^+$$

with  $\eta^{(i)} \in \mathbb{R}^+$  being the *learning rate*, and

$$L_{\mathcal{D}_{XY}}(\boldsymbol{\theta}^{(i)}) \geq L_{\mathcal{D}_{XY}}(\hat{\boldsymbol{\theta}})$$

# Theoretical background

Euler's inequality : Necessary condition for minimizer on convex sets

Given a function  $f: K \rightarrow \mathbb{R}$  defined over a non-empty convex set  $K \subset H$   
(Hilbert),  $f$  proper on  $K$  and  $f \in C^1(K)$ ,

$$\langle \nabla_x f(\hat{x}), y - \hat{x} \rangle \geq 0, \quad \forall y \in K \implies \hat{x} \text{ local minimizer of } f \text{ on } K$$

## Why Gradient Descent ?

- Taylor expansion :  $T_{\hat{\theta}} L_{\mathcal{D}_{XY}}(\theta) = L_{\mathcal{D}_{XY}}(\hat{\theta}) + \langle \nabla_{\theta} L_{\mathcal{D}_{XY}}(\hat{\theta}), \theta - \hat{\theta} \rangle + o(\|\theta - \hat{\theta}\|)$
- Euler's inequality holds if  $\langle \nabla_{\theta} L_{\mathcal{D}_{XY}}(\hat{\theta}), \theta - \hat{\theta} \rangle \geq 0, \quad \forall \theta \in \Theta$
- $\theta = \hat{\theta} + \eta \nabla_{\theta} L_{\mathcal{D}_{XY}}(\hat{\theta})$ 
  - $T_{\hat{\theta}} L_{\mathcal{D}_{XY}}(\theta(\eta)) \geq L_{\mathcal{D}_{XY}}(\hat{\theta}), \quad \eta > 0$
  - $\xi: (\eta, \theta) \mapsto L_{\mathcal{D}_{XY}}(\theta - \eta \cdot \nabla_{\theta} L_{\mathcal{D}_{XY}}(\theta))$
  - $T_0 \xi(\eta; \theta) = L_{\mathcal{D}_{XY}}(\theta) - \eta \|\nabla_x L_{\mathcal{D}_{XY}}(\theta)\|^2 + o(\eta)$
- $$-\frac{\nabla_{\theta} L_{\mathcal{D}_{XY}}(\theta)}{\|\nabla_{\theta} L_{\mathcal{D}_{XY}}(\theta)\|} = \frac{1}{r} \arg \min_{\|\delta \theta\|=r} L_{\mathcal{D}_{XY}}(\theta + \delta \theta)$$

# Learning rate

## Convergence analysis of functions with Lipschitz gradients

A function  $f: K \rightarrow \mathbb{R}$ , non-empty convex set  $K \subset H$  (Hilbert),  $f$  proper and strictly convex on  $K$ ,  $f \in C^1(K)$  and a gradient  $\nabla_x f \in \text{Lip}(K)$  of constant  $\beta > 0$ .  
 $\exists (\mathbf{x}_k)_{k \in \mathbb{N}} \in K$  such that

$$\eta_k \in \left[ \frac{1}{\beta}, \frac{2}{\beta} \right] \quad \mathbf{x}_{k+1} = \mathbf{x}_k - \eta_k \nabla_x f(\mathbf{x}_k)$$

$\exists A > 0$  such that :  $f(\mathbf{x}_k) - f(\hat{\mathbf{x}}) \leq \frac{A}{k+1} \rightarrow \lim_{k \rightarrow \infty} \mathbf{x}_k = \hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in K} f(\mathbf{x})$

If  $f$  is strongly convex of coefficient  $\alpha$ ,  $\exists \rho = \sqrt{\frac{\beta}{\alpha+\beta}} \leq \sqrt{\frac{1}{1+\kappa(\mathbf{H}_f)}} :$

$$\|\mathbf{x}_k - \hat{\mathbf{x}}\| \leq \rho^k \|\mathbf{x}_0 - \hat{\mathbf{x}}\|, \quad \forall k \in \mathbb{N}$$

$$\kappa(\mathbf{H}_f(\mathbf{x})) = \frac{\max_{1 \leq i \leq d_K} \lambda_i(\mathbf{H}_f(\mathbf{x}))}{\min_{1 \leq i \leq d_K} \lambda_i(\mathbf{H}_f(\mathbf{x}))} \text{ and } 0 < \kappa(\mathbf{H}_f(\mathbf{x})) \leq \frac{\beta}{\alpha}$$

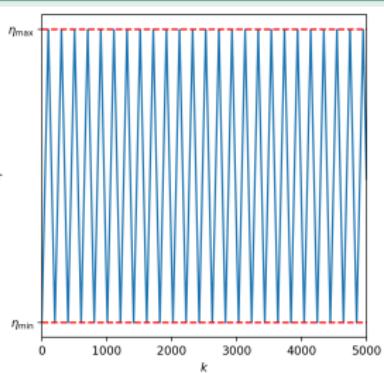
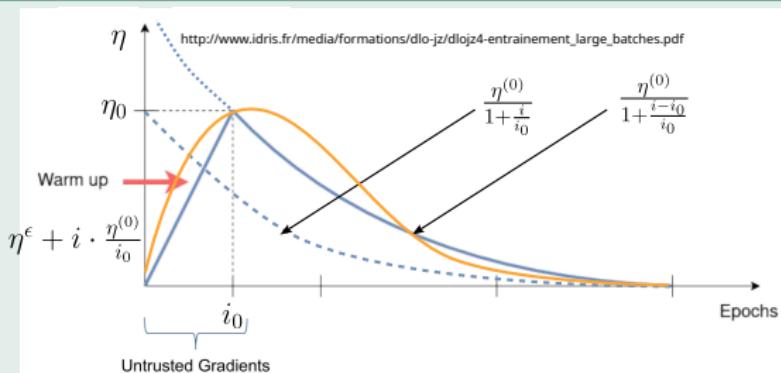
# Learning rate scheduler

## Classical scheduling

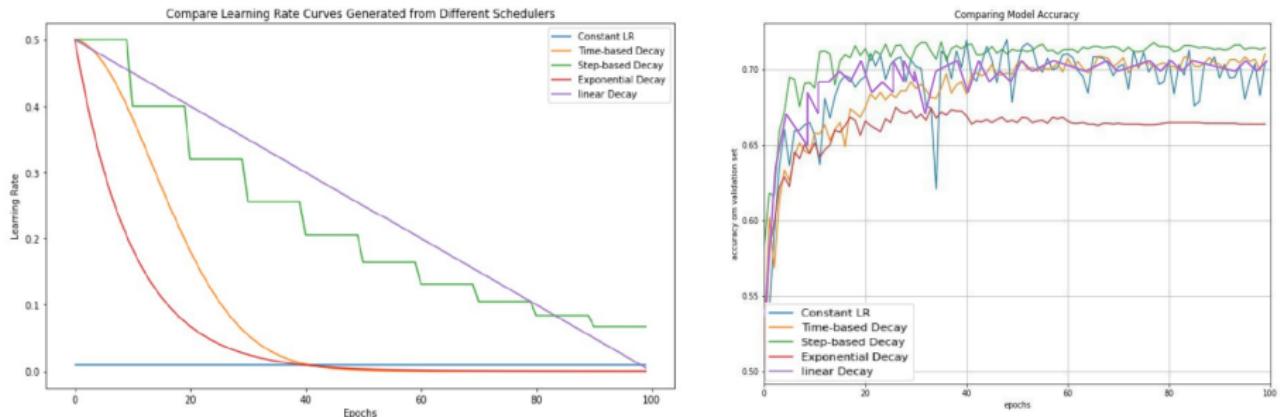
$$\eta^{(i)} = \frac{\eta^{(0)}}{1 + i \cdot \eta_d} \leq \eta^{(0)}$$

$\eta_d$  learning rate decay (often  $\eta_d = \frac{1}{i_0}$ ,  $i_0 > 0$ )

## Warm-up and cyclic scheduler



# Different learning rate schedulers



**Figure – Courtesy of IDRIS : “Deep Learning Optimisé - Jean Zay”**

# Stochastic Gradient Descent

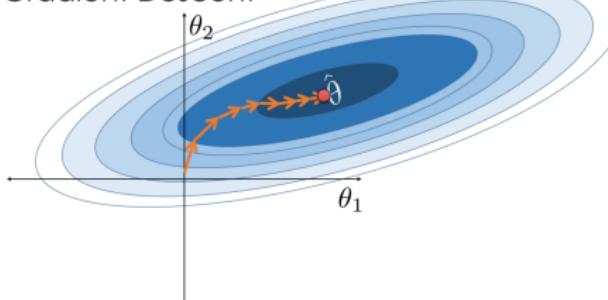
## Computing the gradient

$$\nabla_{\theta} L_{\mathcal{D}_{XY}} \left( \theta^{(i)} \right) = \frac{1}{N} \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}_{XY}} \nabla_{\theta} \ell \left( \mathbf{h}_{\theta} (\mathbf{x}_i), \mathbf{y}_i \right)$$

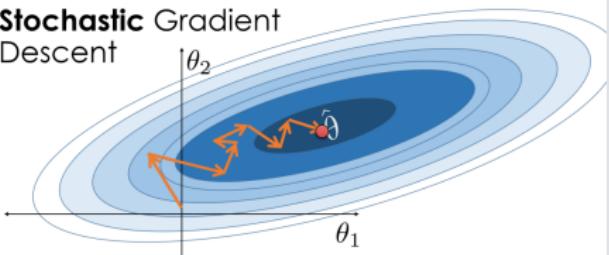
$$\theta^{(i+1)} = \theta^{(i)} - \frac{\eta^{(i)}}{N} \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}_{XY}} \nabla_{\theta} \ell \left( \mathbf{h}_{\theta} (\mathbf{x}_i), \mathbf{y}_i \right)$$

$$\theta^{(i+1)} = \theta^{(i)} - \eta^{(i)} \nabla_{\theta} \ell \left( \mathbf{h}_{\theta} (\mathbf{x}_i), \mathbf{y}_i \right)$$

Gradient Descent



Stochastic Gradient  
Descent



**Figure** – Lau, S., Gonzalez, J., Nolan, D. (2020)

[https://www.textbook.ds100.org/ch/11/gradient\\_stochastic.html](https://www.textbook.ds100.org/ch/11/gradient_stochastic.html),

<https://optimization.cbe.cornell.edu>

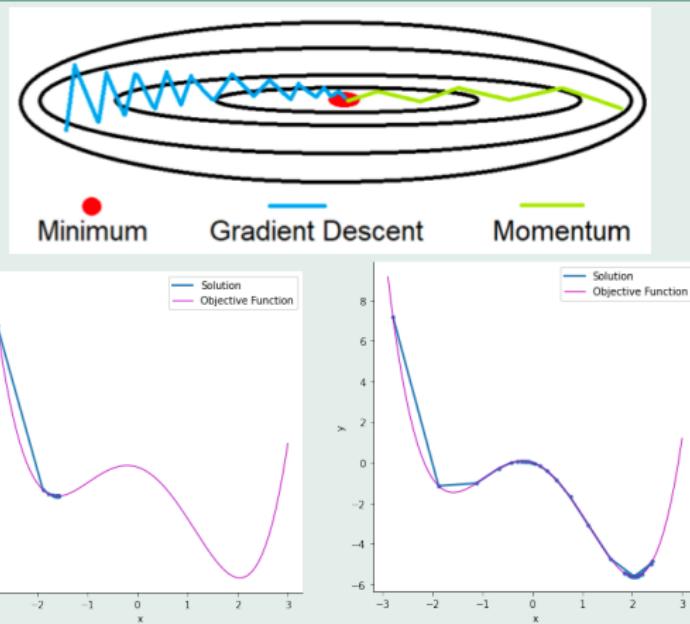
# SGD algorithm

## Use mini-batches

### Algorithm – Stochastic Gradient Descent on mini-batches

```
1 :  $i = 0$ 
2 : Initialize  $\theta^{(0)}$ 
3 :  $\eta^{(0)} = \eta_0$ 
4 : while  $i < n_e$  do
5 :     for  $j \sim \mathcal{U}\left(\{i\}_{i=1}^{\frac{N}{N_b}}\right)$  do
6 :          $\mathbf{g}_\theta = 0$ 
7 :         for  $t \sim \mathcal{U}\left(\{j_i\}_{i=1}^{N_b}\right)$  do
8 :              $\mathbf{g}_\theta += \nabla_{\theta} \ell\left(\mathbf{h}_\theta\left(\mathbf{x}_i\right), \mathbf{y}^{(t)} ; \theta^{(i)}\right)$ 
9 :         end for
10 :         $\theta^{(i+1)} = \theta^{(i)} - \frac{\eta^{(i)}}{N_b} \mathbf{g}_\theta, \quad \eta^{(i)} \in \mathbb{R}^+$ 
11 :    end for
12 : end while
```

The role of momentum : “pushing a ball down a hill”



**Figure** – Srihari, S. (n.d.). Basic Optimization Algorithms. Deep learning.

<https://cedar.buffalo.edu/~srihari/CSE676/8.3%20BasicOptimzn.pdf>,

<https://optimization.cbe.cornell.edu/index.php?title=Momentum>

# SGD algorithm

## The role of momentum

### Algorithm – SGD on mini-batches with Classical Momentum

```

1 :  $i = 0$ 
2 : Initialize  $\theta^{(0)}$ 
3 :  $\eta^{(0)} = \eta_0$ 
4 :  $v^{(0)} = 0$ 
5 : while  $i < n_e$  do
6 :   for  $j \sim \mathcal{U}\left(\{i\}_{i=1}^{\frac{N}{N_b}}\right)$  do
7 :      $\mathbf{g}_\theta = 0$ 
8 :     for  $t \sim \mathcal{U}\left(\{j_i\}_{i=1}^{N_b}\right)$  do
9 :        $\mathbf{g}_\theta += \nabla_{\theta} \ell\left(\mathbf{h}_\theta\left(\mathbf{x}_i\right), \mathbf{y}^{(t)} ; \theta^{(i)}\right)$ 
10 :      end for
11 :       $\mathbf{m}^{(i+1)} = \left(1 - \langle \tau^{(i)} - 1 \rangle\right) \mathbf{g}_\theta + \gamma^{(i)} \mathbf{m}^{(i)}$ 
12 :       $\theta^{(i+1)} = \theta^{(i)} - \frac{\eta^{(i)}}{N_b} \mathbf{m}^{(i+1)}$ 
13 :    end for
14 : end while

```

$\eta^{(i)} \in \mathbb{R}^+, \gamma^{(i)} \in [0, 1]$  (usually  $\gamma^{(i)} = 0.85$  or  $0.95$ )

$-\langle \tau^{(i)} - 1 \rangle$  dampens the gradient only if  $\tau^{(i)} > 1$  (“air resistance”)

## AdaGrad and RMSprop (Hinton's Lecture 6a-Coursera)

- Newton's method :  $\theta^{(i+1)} = \theta^{(i)} - \mathbf{H}_{L_{\mathcal{D}_{XY}}}^{-1}(\theta^{(i)}) \cdot \nabla_{\theta} L_{\mathcal{D}_{XY}}(\theta^{(i)})$
- $\text{diag}(\mathbb{I}_F(\theta; \mathbf{X})) = \text{diag}\left(\mathbb{V}_{x \sim p_{\theta^*}}(s_{\theta, k}(\mathbf{X}))\right) = \text{diag}(\mathbf{H}_{L_{\mathcal{D}_{XY}}})$
- Gradient memory :  $\mathbb{G}^{(i)} = \sum_{k=0}^{i-1} \nabla_{\theta} L_{\mathcal{D}_{XY}}(\theta^{(k)}) \otimes \nabla_{\theta} L_{\mathcal{D}_{XY}}(\theta^{(k)})$
- $\text{Tr}(\mathbb{G}^{(i)}) = \sum_{k=1}^{i-1} \|\nabla_{\theta} L_{\mathcal{D}_{XY}}(\theta^{(k)})\|_2^2 \approx \mathbb{E}\left[\|\nabla_{\theta} L_{\mathcal{D}_{XY}}(\theta)\|^2\right]$

$$\text{AdaGrad} : \theta^{(i+1)} = \theta^{(i)} - \eta \left( \text{diag} \mathbb{G}^{(i)} + \epsilon \mathbb{I} \right)^{-\frac{1}{2}} \nabla_{\theta} L_{\mathcal{D}_{XY}}$$

- Average gradient memory on short window :

$$\text{diag} \tilde{\mathbb{G}}^{(i)} = \alpha_r \text{diag} \mathbb{G}^{(i-1)} + (1 - \alpha_r) \left( \text{diag} \mathbb{G}^{(i)} - \text{diag} \mathbb{G}^{(i-1)} \right), \alpha_r = 0.9$$

$$\text{RMSprop} : \theta^{(i+1)} = \theta^{(i)} - \eta \left( \text{diag} \tilde{\mathbb{G}}^{(i)} + \epsilon \mathbb{I} \right)^{-\frac{1}{2}} \nabla_{\theta} L_{\mathcal{D}_{XY}}$$

## ADADELTA

- AdaGrad's  $\eta$  decay too aggressive !
- Introduce momentum vectors  $\mathbf{U}^{(i+1)} = \alpha_u \mathbf{U}^{(i)} + (1 - \alpha_u) \mathbf{u}^{(i)}$ ,  
$$\mathbb{U}^{(i)} = \left( \mathbf{U}^{(i)} \otimes \mathbf{U}^{(i)} \right)^{\frac{1}{2}}$$
- Compute vectors  
$$\mathbf{u}^{(i+1)} = \text{diag}^{\frac{1}{2}} \left( \mathbb{U}^{(i)} + \epsilon \mathbb{I} \right) \text{diag}^{-\frac{1}{2}} \left( \tilde{\mathbb{G}}^{(i)} + \epsilon \mathbb{I} \right) \cdot \nabla_{\theta} L_{\mathcal{D}_{XY}} \left( \boldsymbol{\theta}^{(i)} \right)$$
- Approximate Hessian :  $\mathbf{H}_{L_{\mathcal{D}_{XY}}} \left( \boldsymbol{\theta}^{(i)} \right) \approx \text{diag}^{\frac{1}{2}} \left( \mathbb{U}^{(i)} + \epsilon \mathbb{I} \right)$

$$\text{ADADELTA} : \boldsymbol{\theta}^{(i+1)} = \boldsymbol{\theta}^{(i)} - \mathbf{u}^{(i)}$$

- No need to set a default learning rate

Adam : “heavy ball with friction” [Heu+]

- Decay average of past gradients :

$$\mathbf{m}^{(i)} = \beta_1 \mathbf{m}^{(i-1)} + (1 - \beta_1) \nabla_{\theta} L_{\mathcal{D}_{XY}}(\boldsymbol{\theta}^{(i)})$$

- Decay average of past squared gradients :

$$\text{diag} \mathbb{V}^{(i)} = \beta_2 \mathbb{V}^{(i-1)} + (1 - \beta_2) \text{diag} \mathbb{G}^{(i)}$$

- Normalization :  $\hat{\mathbf{m}}^{(i)} = \frac{\mathbf{m}^{(i)}}{1-\beta_1}, \quad \hat{\mathbb{V}}^{(i)} = \frac{\mathbb{V}^{(i)}}{1-\beta_2}$

$$\text{Adam} : \boldsymbol{\theta}^{(i+1)} = \boldsymbol{\theta}^{(i)} - \eta \cdot \text{diag}^{-\frac{1}{2}} \left( \hat{\mathbb{V}}^{(i)} + \epsilon \mathbb{I} \right) \hat{\mathbf{m}}^{(i)}$$

- $\beta_1 = 0.9, \beta_2 = 0.999$

# Create a statistical model

# The artificial neuron

## Empirical loss minimization

Find  $\mathbf{h}_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ ,  $\mathbf{h}_\theta \in \mathcal{H}_\theta$  such that :

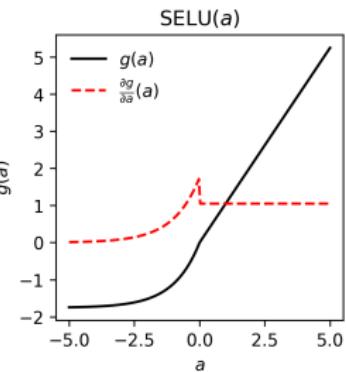
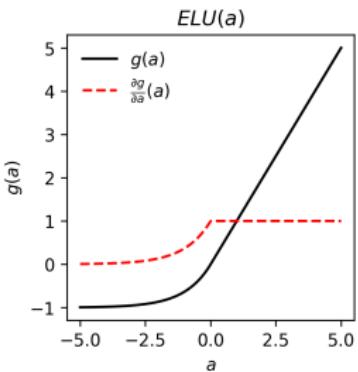
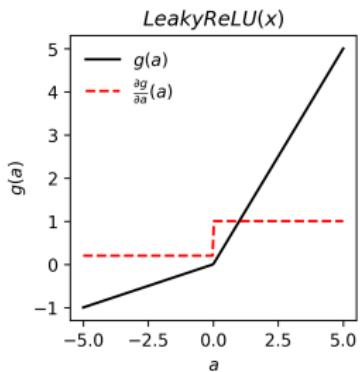
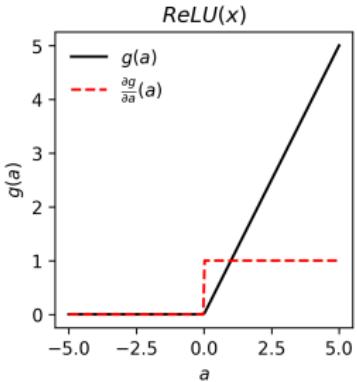
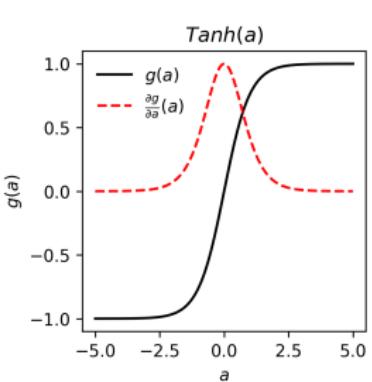
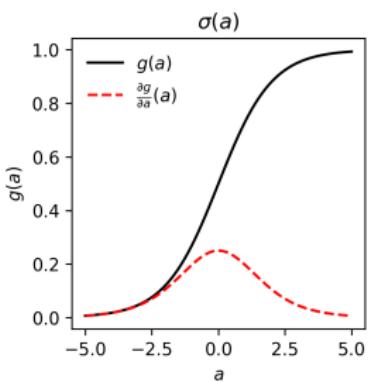
$$\mathbf{h}_{\hat{\theta}}(\mathcal{D}_{XY}) = \arg \min_{\mathbf{h}_\theta \in \mathcal{H}_\theta} L_{\mathcal{D}_{XY}}(\mathbf{h}_\theta)$$

$$\ell(y, h_\theta(\mathbf{x})) = -\ln p_\theta(y|\mathbf{x}) = \mathcal{NLL}(y, h_\theta(\mathbf{x}))$$

## The simplest statistical model [MP43 ; Ros57]

$$h_\theta(\mathbf{x}) = g(a(\mathbf{x})) = g(\langle \mathbf{w}, \mathbf{x} \rangle + b) = g\left(\sum_{c=1}^{d_X} w_c \cdot x_c + b\right); \quad \theta := \{\mathbf{w}; b\}$$

# Activation functions

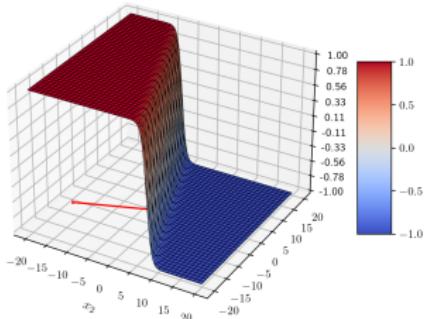
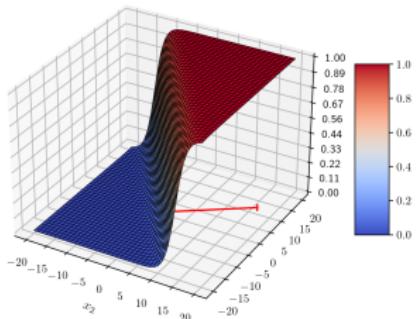


# Decision boundary

- Binary classification  $y \in \mathcal{Y} := \{0, 1\}$  (if  $g(\cdot) = \sigma(\cdot)$ ) or  $y \in \mathcal{Y} := \{-1, 1\}$  (if  $g(\cdot) = \tanh(\cdot)$ )
- hyperplane that separates the samples of class  $y = 0$  (or  $y = -1$ ) if  $h_{\theta}(\mathbf{x}) > 0.5$  from those labeled as  $y = 1$ , if  $h_{\theta}(\mathbf{x}) \leq 0.5$

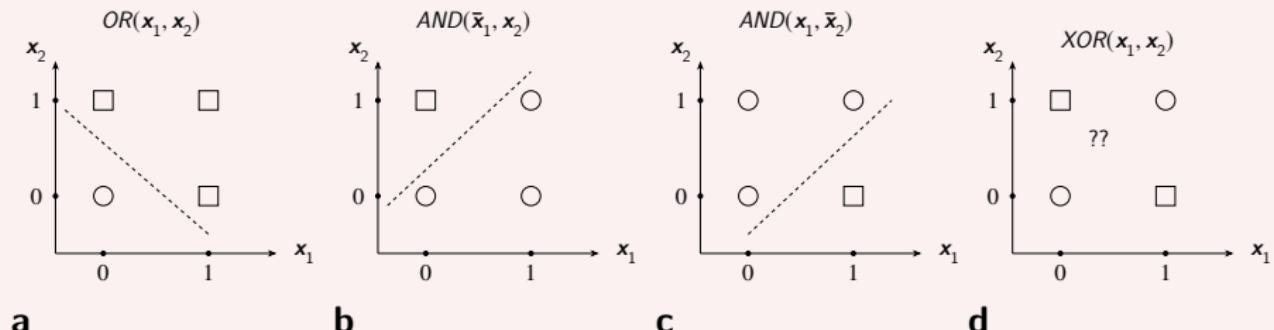
$$g(a) = \frac{1}{1+e^{-a}}$$

$$g(a) = \tanh(a)$$



# The role of hidden layers

## Unsolved decision boundary

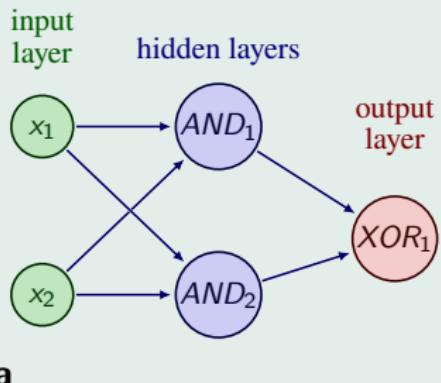
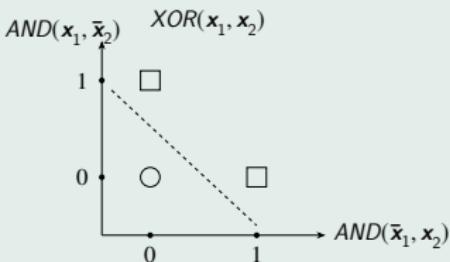


**Figure** – Solved and unsolved decision boundaries for different binary operations. (a)  $f_{OR}(x) = \{x_1\} \cup \{x_2\}$ , (b)  $f_{AND1}(x) = \{x_1\} \cap \{\{0, 1\} / \{x_2\}\}$ , (c)  $f_{AND2}(x) = \{\{0, 1\} / (x_1)\} \cap \{x_2\}$ , (d)  $f_{XOR}(x) = \{\{x_1\} \cup \{x_2\}\} / (\{\{x_1\} \cap \{x_2\}\})$  (unsolved). Reprinted from the video-lecture notes by Hugo Larochelle (<https://www.youtube.com/watch?v=iT5P4z6Fzj8&list=PL6Xpj9I5qXYEcOhn7TqghAJ6NAPrNmUBH&index=3>).

# The role of hidden layers

Add intermediate feature=hidden layers

$$h_{\theta}(\mathbf{x}) = g(a(\phi(\mathbf{x}))) = g\left(\sum_{c=1}^{d_{\phi}} w_c \cdot \phi_c(\mathbf{x}) + b\right); \quad \theta := \{\mathbf{w}; b\}$$

**a****b**

**Figure** – Simple example of Multi-Layer Perceptron (1 hidden layers  $h^{(1)}$  and 1 output  $y \in \mathbb{R}$ ). layer. The figure was drawn with tikz.net.

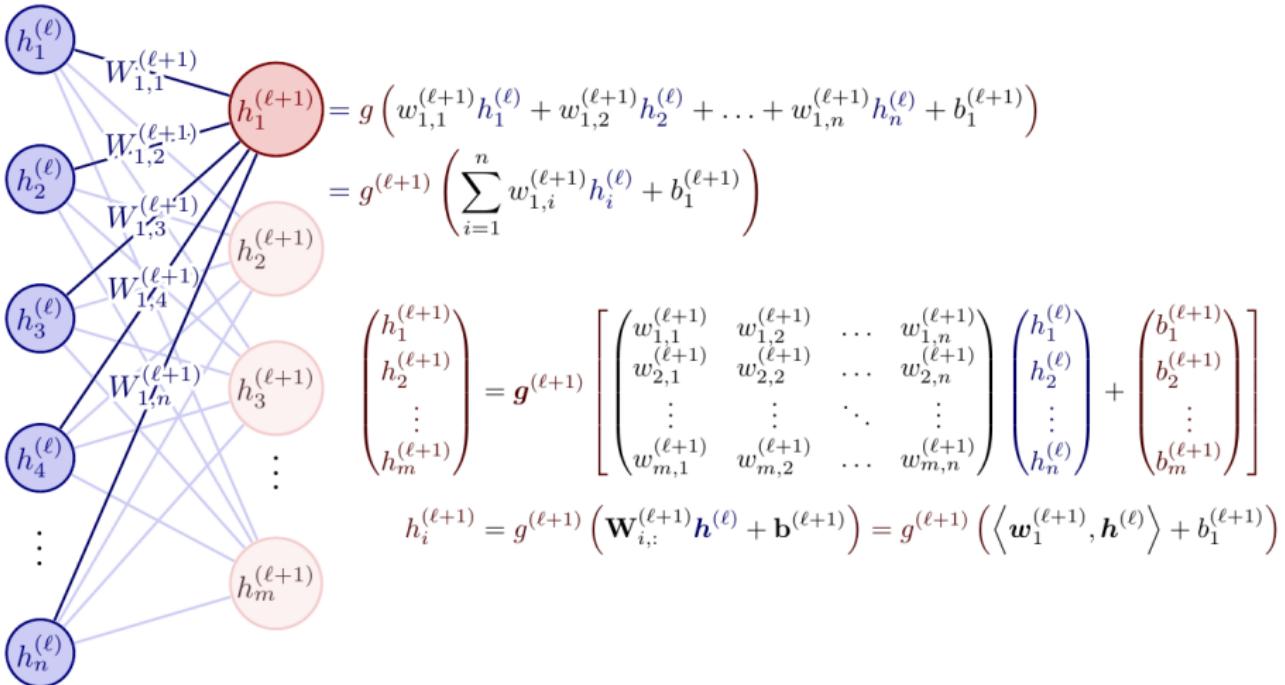
# Multi Layer Perceptrons

## Stacking layers of neurons

$$\begin{aligned} h_{\theta}(\mathbf{x}) &= h^{(o)} \circ \mathbf{h}^{(1)}(\mathbf{x}) = g^{(o)}(a^{(o)}(\mathbf{h}^{(1)}(\mathbf{x}))) = g\left(\langle \mathbf{w}^{(o)}, \mathbf{h}^{(1)}(\mathbf{x}) \rangle + b^{(o)}\right) = \\ &= g^{(o)}\left(\sum_{c=1}^{d_{h_1}} w_c^{(o)} \cdot h_c^{(1)}(\mathbf{x}) + b^{(o)}\right); \quad \theta := \{\mathbf{w}^{(1)}, \mathbf{w}^{(o)}; b^{(1)}, b^{(o)}\} \end{aligned}$$

with

$$h_c^{(1)} = g^{(1)}(a_c^{(1)}(\mathbf{x})) = g^{(1)}\left(\langle \mathbf{w}^{(1)}, \mathbf{x} \rangle + b^{(1)}\right) = g^{(1)}\left(\sum_{c=1}^{d_X} w_c^{(1)} \cdot x_c + b^{(1)}\right)$$



# How many layers and neurons ?

Lower bounds for approximation by  $\mathcal{MLP}$  neural networks [Mai99]

Given a function  $f: [-T, T]^{d_X} \in \mathcal{X} \rightarrow \mathbb{R}$ ,  $T > 0$ , with  $f \in \mathcal{W}^{k,2}$ , then it exists a 1-hidden-layer  $\mathcal{MLP}$   $h_\theta$ , with sigmoid activation function and  $N_K$  hidden neurons, such that the  $L^2$ -error  $e = \|f - h_\theta\|_{L^2([-T, T]^{d_X})}$  is lower than a tolerance  $\varepsilon$  if :

$$N_K \approx \varepsilon^{\frac{1-d_X}{k}} \quad (1)$$

$$d_X = \dim(\mathcal{X})$$

For  $f \in L^2([-T, T]^{d_X})$ , to reduce the error by an order or magnitude, the number of neurons must be multiplied by  $10^{\frac{d_X-1}{2}}$  !

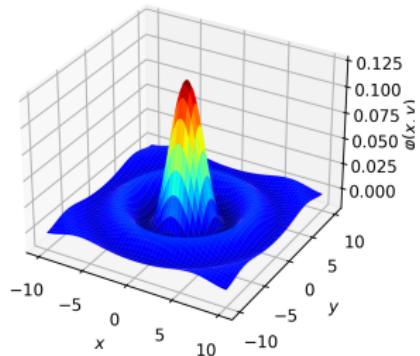
# Hands-on 1

## Approximate radial functions [ES16]

$$\varphi(\mathbf{x}) = \left( \frac{R_{d_X}}{\|\mathbf{x}\|} \right)^{\frac{d_X}{2}} J_{\frac{d_X}{2}}(2\pi R_{d_X} \|\mathbf{x}\|) = \int_{\mathbf{w}: \|\mathbf{w}\| \leq R_d} e^{-2\pi i \langle \mathbf{x}, \mathbf{w} \rangle} d\mathbf{w}$$

$J_{\frac{d_X}{2}}(2\pi R_{d_X} \|\mathbf{x}\|)$  : Bessel function of first kind of order  $\frac{d_X}{2}$

- Is 1 hidden layer enough ?
- How many neurons ?
- Approximation error ?
- Try with two hidden layers !



**Figure** –  $d_X = 2$  and  $R_d = 0.2$ .

# Design a $\mathcal{MLP}$

## Main strategies

- Reduce the dimensionality :  $\phi : \mathbf{x} \mapsto \mathbf{z} \in S$ ,  $\dim(S) < \dim(\mathcal{X})$  (PCA, ROM)
- Separate the interactions (disentanglement) :

$$\exists \{f_n\}_{n \in I}, f_n : \mathcal{X}_n \rightarrow \mathbb{R}, \mathcal{X}_n \subset \mathcal{X} \quad f(\mathbf{x}) = \sum_{n \in I} f_n(\mathbf{x})|_{\mathbf{x} \in \mathcal{X}_n}, \forall \mathbf{x} \in \mathcal{X}$$

$$\varepsilon = \sum_{n \in I} \varepsilon_n \approx \text{card}(I) \cdot N_K^{\frac{k}{1-d_S}} \quad N_K \approx \left( \frac{\varepsilon}{\text{card}(I)} \right)^{\frac{1-d_S}{k}}$$

$$\exists \{f_n\}_{n=1}^{d_X}, f_n : \mathbb{R} \rightarrow \mathbb{R} \quad f(\mathbf{x}) = \sum_{n=1}^{d_X} f_n(w_n \cdot x_n), \forall \mathbf{x} \in \mathcal{X}$$

- Promote sparsity : dictionary learning  $D = \{(k, v(\mathbf{x}))_n\}_{n \leq N_D}$  ( $k$  is the key and  $v$  the corresponding value)

$$h_{\theta}(\mathbf{x}) = \sum_{n \in I} \theta_n [k_n] v_n(\mathbf{x}), \quad \text{card}(I) \leq N_D$$

# Train a statistical model

# Automatic differentiation

## Chain rule

$$\nabla_{\theta} L_{\mathcal{D}_{XY}} = \frac{1}{N} \sum_{k=1}^N \nabla_{\theta} \ell(h_{\theta}(x_k), y_k)$$

$$\nabla_{\theta} \ell(h_{\theta}(x)) = \mathbb{J}_h^T \cdot \nabla_{h_{\theta}} \ell(h_{\theta}) \quad \mathbb{J}_h^{(k)} = \left[ \begin{array}{c} \frac{\partial h_{\theta}^{(k)}}{\partial \theta_1^{(k)}} \frac{\partial h_{\theta}^{(k)}}{\partial \theta_2^{(k)}} \cdots \frac{\partial h_{\theta}^{(k)}}{\partial \theta_{u^{(k+1)}}^{(k)}} \end{array} \right]$$

Example : Autodiff of a classifier  $h_{\theta}(x) = \sigma_y \left( \sum_{c=1}^{N_y} h_c^{(N_{\ell})}(x) e_c \right)$

- $h^{(o)} = g^{(o)}(a^{(o)}) = \sigma_y(a^{(o)}) = \text{softmax}_y(a^{(o)}) = \frac{e^{a_y^{(o)}(x)}}{\sum_{c=1}^{N_y} e^{a_c^{(o)}(x)}}$
- $a_c^{(o)}(x) = h_c^{(N_{\ell})} = g_c^{(N_{\ell})}(a_c^{(N_{\ell})}(x))$ ,  $1 \leq c \leq N_y$
- $\ell(h_{\theta}(x), y) = - \sum_{c=1}^{N_y} \ln(\sigma_c(a^{(o)}(x))) \chi_{(y=c)}$

1 Derivative at output activation :  $\frac{\partial \ell(h_{\theta}(x), y)}{\partial \sigma_c} = -\frac{\chi_{(y=c)}}{\sigma_c(a^{(o)}(x))}$

2 Gradient at output activation :

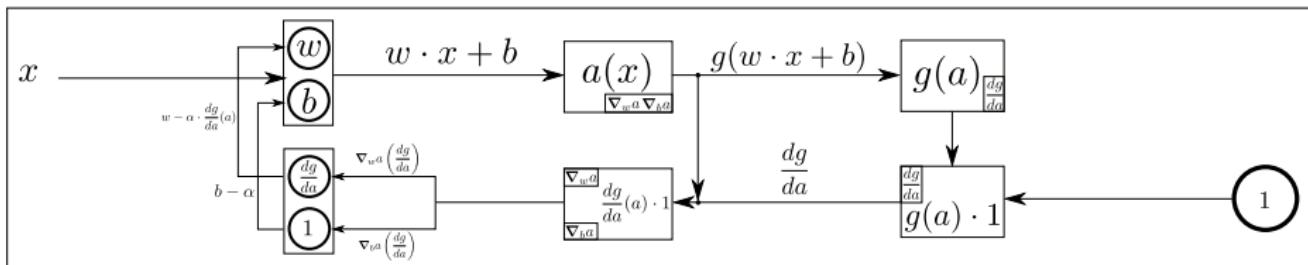
$$\nabla_{\sigma} \ell(h_{\theta}(x), y) = \sum_{c=1}^{N_y} \frac{\partial \ell(h_{\theta}(x), y)}{\partial \sigma_c} e_c = - \sum_{c=1}^{N_y} \frac{\chi_{(y=c)}}{\sigma_c(a^{(o)}(x))} e_c = - \frac{e_y}{\sigma_y(a^{(o)}(x))}$$

3 Derivative at  $N_{\ell}$ -layer activation :  $\frac{\partial \ell(h_{\theta}(x), y)}{\partial h_c^{(N_{\ell})}} = \sum_{c'=1}^{N_y} \frac{\partial \ell(h_{\theta}(x), y)}{\partial \sigma_{c'}'} \cdot \frac{\partial \sigma_{c'}'}{\partial h_c^{(N_{\ell})}} \Big|_{a^{(o)}(x)}$

with :  $\frac{\partial \sigma_{c'}'}{\partial h_c^{(N_{\ell})}} \Big|_{a^{(o)}(x)} = \frac{\partial \sigma_{c'}'}{\partial a_c^{(o)}} = \sigma_{c'}'(a^{(o)}(x)) \cdot (\chi_{(c'=c)} - \sigma_c(a^{(o)}(x)))$

4 Gradient at  $N_{\ell}$ -layer activation :

$$\begin{aligned} \nabla_{h^{(N_{\ell})}} \ell(h_{\theta}(x), y) &= \sum_{c=1}^{N_y} \frac{\partial \ell(h_{\theta}(x), y)}{\partial h_c^{(N_{\ell})}} e_c = \sum_{c=1}^{N_y} \sum_{c'=1}^{N_y} \frac{\partial \ell(h_{\theta}(x), y)}{\partial \sigma_{c'}'} \cdot \frac{\partial \sigma_{c'}'}{\partial h_c^{(N_{\ell})}} \Big|_{a^{(o)}(x)} = \\ &= - \sum_{c=1}^{N_y} (\chi_{(y=c)} - \sigma_c) e_c = - (e_y - \sigma(a^{(o)}(x))) \end{aligned}$$



5 Derivative at  $N_\ell$ -layer pre-activation :  $\frac{\partial \ell(h_{\theta}(x), y)}{\partial a_c^{(N_\ell)}} = \sum_{c'=1}^{N_y} \frac{\partial \ell(h_{\theta}(x), y)}{\partial h_{c'}^{(N_\ell)}} \cdot \frac{\partial h_{c'}^{(N_\ell)}}{\partial a_c^{(N_\ell)}}$  with :

$$\frac{\partial h_{c'}^{(N_\ell)}(x)}{\partial a_c^{(N_\ell)}} = \frac{\partial g_{c'}^{(N_\ell)}}{\partial a_c^{(N_\ell)}} \left( a_{c'}^{(N_\ell)}(x) \right) \chi_{(c'=c)} \quad \frac{\partial \ell(h_{\theta}(x), y)}{\partial a_c} = \\ - \sum_{c'=1}^{N_y} \left( \chi_{(y=c')} - \sigma_{c'} \right) \frac{\partial g_{c'}^{(N_\ell)}}{\partial a_c^{(N_\ell)}} \left( a_{c'}^{(N_\ell)}(x) \right) \chi_{(c'=c)} = - \left( \chi_{(y=c)} - \sigma_c \right) \frac{\partial g_c^{(N_\ell)}}{\partial a_c^{(N_\ell)}} \left( a_c^{(N_\ell)}(x) \right)$$

6 Gradient at  $N_\ell$ -layer pre-activation (provided that  $g_c^{(N_\ell)} = g \quad \forall c$ )

$$\nabla_{a^{(N_\ell)}} \ell(h_{\theta}(x), y) = \sum_{c=1}^{N_y} \frac{\partial \ell(h_{\theta}(x), y)}{\partial a_c^{(N_\ell)}} e_c = \nabla_h \ell(h_{\theta}(x), y) \odot \nabla_a g = - \left( e_y - \sigma(a^{(o)}(x)) \right) \odot \nabla_a g$$

7 Derivative at  $N_\ell$ -layer weights :  $a_c^{(N_\ell)} = \sum_{c'=1}^{N_y} W_{cc'}^{(N_\ell)} h_{c'}^{(N_{\ell-1})} + b_c^{(N_\ell)}$

$$\frac{\partial \ell(h_{\theta}(x), y)}{\partial W_{cc'}^{(N_\ell)}} = \sum_{n=1}^{N_y} \frac{\partial \ell(h_{\theta}(x), y)}{\partial a_n^{(N_\ell)}} \frac{\partial a_n^{(N_\ell)}}{\partial W_{cc'}^{(N_\ell)}} = \frac{\partial \ell(h_{\theta}(x), y)}{\partial a_c^{(N_\ell)}} h_{c'}^{(N_{\ell-1})}$$

8 Derivative at  $N_\ell$ -layer bias :  $\frac{\partial \ell(h_{\theta}(x), y)}{\partial b_c^{(N_\ell)}} = \sum_{n=1}^{N_y} \frac{\partial \ell(h_{\theta}(x), y)}{\partial a_n^{(N_\ell)}} \frac{\partial a_n^{(N_\ell)}}{\partial b_c^{(N_\ell)}} = \frac{\partial \ell(h_{\theta}(x), y)}{\partial a_c^{(N_\ell)}} \cdot 1$

9 Gradient at  $N_\ell$ -layer weight  $W_{cc'}^{(N_\ell)}$  :  $\nabla_{W^{(N_\ell)}} \ell(h_{\theta}(x), y) = \nabla_a \ell(h_{\theta}(x), y) \otimes h^{(N_{\ell-1})}$

10 Gradient at  $N_\ell$ -layer bias :  $\nabla_{b^{(N_\ell)}} \ell(h_{\theta}(x), y) = \nabla_a \ell(h_{\theta}(x), y)$

⑪ Derivative at hidden layer  $N_{\ell-1}$  activation :

$$\frac{\partial \ell(h_{\theta}(x), y)}{\partial h_c^{(N_{\ell-1})}} = \sum_{n=1}^{N_y} \frac{\partial \ell(h_{\theta}(x), y)}{\partial a_n^{(N_{\ell})}} \frac{\partial a_n^{(N_{\ell})}}{\partial h_c^{(N_{\ell-1})}} = \sum_{n=1}^{N_y} \frac{\partial \ell(h_{\theta}(x), y)}{\partial a_n^{(N_{\ell})}} W_{nc}^{(N_{\ell})}$$

⑫ Gradient at hidden layer  $N_{\ell-1}$  activation  $\nabla_{h^{(N_{\ell-1})}} \ell(h_{\theta}(x), y) = W^{(N_{\ell})T} \cdot \nabla_a \ell(h_{\theta}(x), y)$

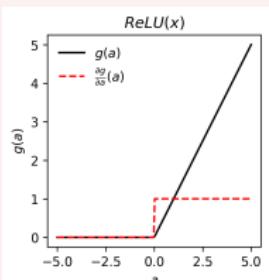
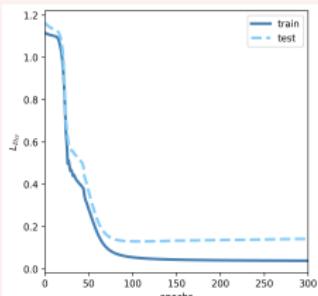
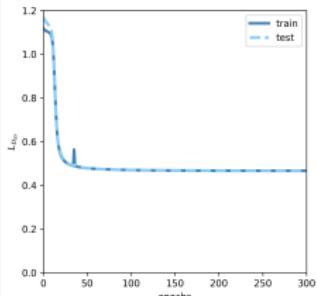
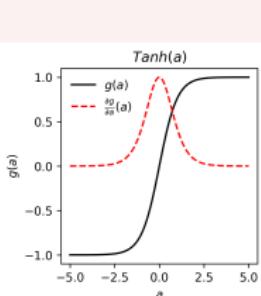
⑬ Reiterate  $\ell = 1 \dots N_{\ell} - 1$

## Vanishing gradients

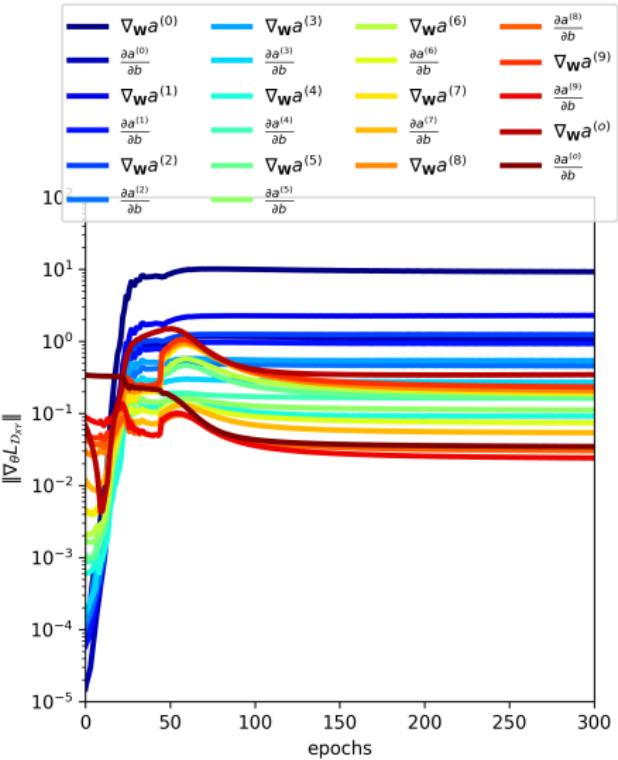
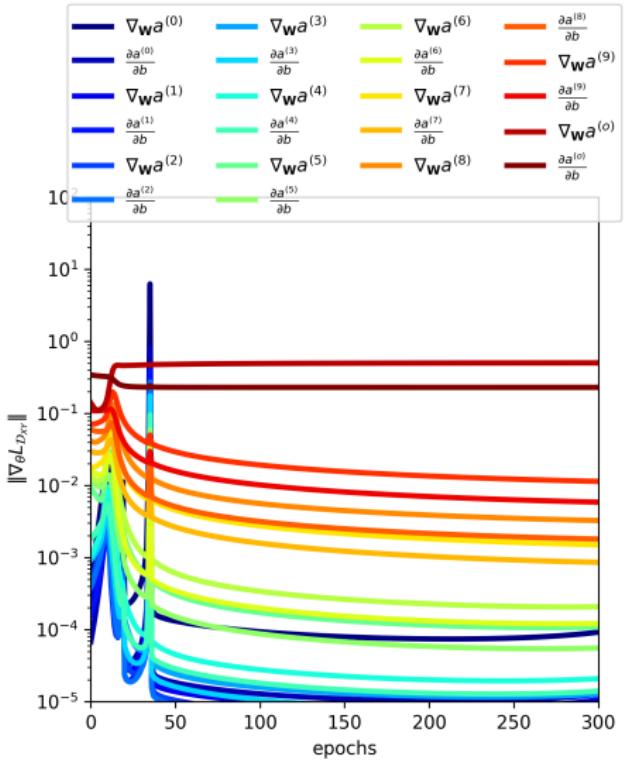
$$h_{\theta}(x) = g^{(o)} \circ a^{(o)} \circ g^{(N_\ell)} \circ a^{(N_\ell)} \circ \dots \circ g^{(\ell)} \circ a^{(\ell)} \circ \dots g^{(1)} \circ a^{(1)}(x)$$

$$\frac{\partial L_{\mathcal{D}_{XY}}}{\partial w_c^{(1)}} = \frac{\partial g^{(o)}}{\partial a^{(o)}} \cdot \frac{\partial a^{(o)}}{\partial h^{(N_\ell)}} \left( \prod_{\substack{\ell=1 \\ N_\ell > 1}}^{N_\ell-1} \frac{\partial g^{(N_\ell+1-\ell)}}{\partial a^{(N_\ell+1-\ell)}} \cdot \frac{\partial a^{(N_\ell+1-\ell)}}{\partial h^{(N_\ell-\ell)}} \right) \frac{\partial g^{(1)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial w_c^{(1)}}$$

$$\nabla_{w^{(1)}} L_{\mathcal{D}_{XY}} \approx 0 \dots \frac{\partial g^{(N_\ell+1-\ell)}}{\partial a^{(N_\ell+1-\ell)}} \approx 0$$



## Vanishing gradients



# Initialize weights

Example : 2-hidden layers  $\mathcal{MLP}$

$$\frac{\partial L_{\mathcal{D}XY}}{\partial w_{1c}^{(1)}} = \frac{\partial g^{(o)}}{\partial a^{(o)}} \cdot \frac{\partial a^{(o)}}{\partial h_1^{(1)}} \frac{\partial h_1^{(1)}}{\partial a_1^{(1)}} \cdot x_c \quad \frac{\partial L_{\mathcal{D}XY}}{\partial w_{2c}^{(1)}} = \frac{\partial g^{(o)}}{\partial a^{(o)}} \cdot \frac{\partial a^{(o)}}{\partial h_2^{(1)}} \frac{\partial h_2^{(1)}}{\partial a_2^{(1)}} \cdot x_c$$

$$\mathbb{V}(a_i^{(k)}) = \mathbb{V}\left(\sum_{j=1}^{u^{(k)}} W_{ij}^{(k)} h_j^{(k-1)}\right) + \mathbb{V}\left(b_i^{(k)}\right)$$

$$\mathbb{V}\left(\sum_{j=1}^{u^{(k)}} W_{ij}^{(k)} h_j^{(k-1)}\right) = \sum_{j=1}^{u^{(k)}} \mathbb{V}\left(W_{ij}^{(k)}\right) \cdot \left(\mathbb{E}\left[h_j^{(k-1)}\right]\right)^2 + \sum_{j=1}^{u^{(k)}} \left(\mathbb{E}\left[W_{ij}^{(k)}\right]\right)^2 \cdot \mathbb{V}\left(h_j^{(k-1)}\right) + \sum_{j=1}^{u^{(k)}} \mathbb{V}\left(W_{ij}^{(k)}\right) \cdot \mathbb{V}\left(h_j^{(k-1)}\right)$$

## How to initialize weights ?

- $\theta^{(0)} = cst.$ 
  - If  $N_\ell = 1$  : no vanishing gradients
  - Symmetric evolution :
 
$$a_1^{(1)} = a_2^{(1)} = 0 \rightarrow h_1^{(1)} = g(a_1^{(1)}) = h_1^{(1)} = g(a_1^{(1)}) \rightarrow w_{1c}^{(1)(i)} = w_{2c}^{(1)(i)}$$

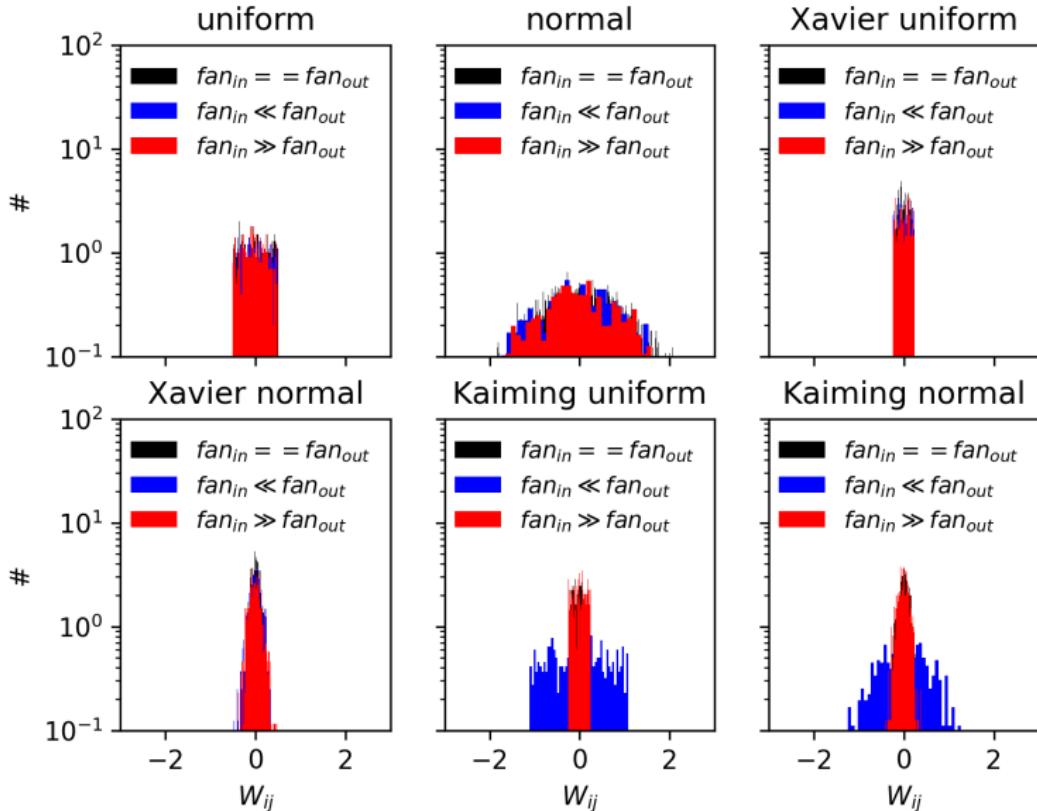
# Initialize weights

Example : 2-hidden layers  $\mathcal{MLP}$

$$\frac{\partial L_{\mathcal{D}XY}}{\partial w_{1c}^{(1)}} = \frac{\partial g^{(o)}}{\partial a^{(o)}} \cdot \frac{\partial a^{(o)}}{\partial h_1^{(1)}} \frac{\partial h_1^{(1)}}{\partial a_1^{(1)}} \cdot x_c \quad \frac{\partial L_{\mathcal{D}XY}}{\partial w_{2c}^{(1)}} = \frac{\partial g^{(o)}}{\partial a^{(o)}} \cdot \frac{\partial a^{(o)}}{\partial h_2^{(1)}} \frac{\partial h_2^{(1)}}{\partial a_2^{(1)}} \cdot x_c$$

## How to initialize weights ?

- $\theta_i \sim \mathcal{N}(\mathbf{0}, \sigma_\theta \mathbb{I})$  or  $\theta \sim \mathcal{U}(-\sqrt{3\sigma_\theta}, \sqrt{3\sigma_\theta})$ 
  - which  $\sigma_\theta$  ?
  - $\mathbb{V}(a_i^{(k)}) = (u^{(k)} \cdot \mathbb{V}(W^{(k)}))^i \cdot \mathbb{V}(h^{(k-1)}) + \mathbb{V}(b_i^{(k)})$
  - $(u^{(k)} \cdot \mathbb{V}(W^{(k)}))^i$  can vanish
  - Decorrelate :  $\mathbb{V}(W^{(k)}) = \frac{1}{u^{(k)}} \Rightarrow \mathbb{V}(a_i^{(k)}) = \mathbb{V}(h^{(k-1)}) + \mathbb{V}(b_i^{(k)})$
  - Xavier/Glorot :  $\mathbb{V}(W^{(k)}) = \frac{2}{u^{(k)} + u^{(k+1)}}$
  - He/Kaiming :  $\mathbb{V}(W^{(k)}) = \frac{2}{u^{(k)}}$



# Neural networks for time series

# Time-forward prediction

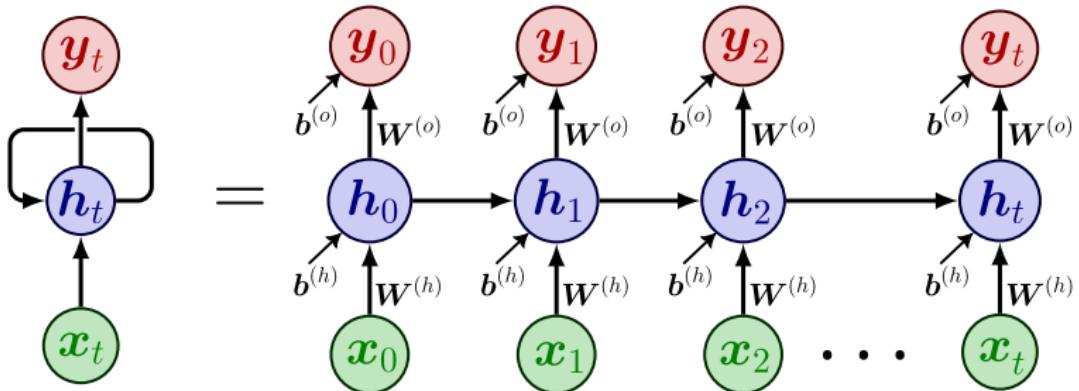
## Time-histories

$\mathcal{D}_{XY} = \{(\mathbb{X}_i, \mathbb{Y}_i)\}_{i=1}^N$ : i.i.d. discrete input  $\mathbb{X} = \{\mathbf{x}[t]\}_{t=1}^{N_t}$  and output signals  $\mathbb{Y} = \{\mathbf{y}[t]\}_{t=1}^{N_t}$ , with  $\mathbf{x}[t] \in \mathbb{R}^{d_X}$  and  $\mathbf{y}[t] \in \mathbb{R}^{d_Y}$ , both of length  $N_t$

## Recurrent Neural Networks $\mathcal{RNN}$

$$\begin{cases} \mathbf{a}_t(\mathbf{h}, \mathbf{x}) = \mathbf{W}^{(h)} \mathbf{h}_{t-1} + \mathbf{W}^{(x)} \mathbf{x}[t] + \mathbf{b}^{(h)} \\ \mathbf{h}_t(\mathbf{a}) = \mathbf{g}^{(h)}(\mathbf{a}_t) \\ \mathbf{a}_t^{(o)} = \mathbf{W}^{(o)} \mathbf{h}_t + \mathbf{b}^{(o)} \\ \mathbf{z}_t(\mathbf{a}^{(o)}) = \mathbf{g}^{(o)}(\mathbf{a}_t^{(o)}) \\ \mathbf{h}_\theta(\mathbf{x}[t]) = \mathbf{z}_t \circ \mathbf{h}_t \circ \mathbf{a}_t(\cdot, \mathbf{h}_{t-1}) \circ \mathbf{x}[t] \end{cases}$$

hidden state  $\mathbf{h}[t] \in \mathbb{R}^{d_h}$ : “memory” variable that keeps track of past states



**Figure – Sketch of RNN scheme according to ??**

- Feedback :  $\mathbf{a}_t (\mathbf{h}, \mathbf{x}) = \mathbf{W}^{(h)} \mathbf{h}_{t-1} + \mathbf{W}^{(x)} \mathbf{x}[t] + \mathbf{b}^{(h)}$
- Biases  $\sim$  exogenous variables
- Loss function  $L_{\mathcal{D}_{XY}}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{t=1}^{N_t} \sum_{i=1}^N \ell(\mathbf{y}_i[t], \mathbf{h}_{\boldsymbol{\theta}}(\mathbf{x}_i[t]))$

# RNN training

- Training instability because of long-term time-dynamics

$$\frac{\partial \ell}{\partial W_{ij}^{(x)}} \left( \mathbf{y}_i[t], \mathbf{h}_{\theta}(\mathbf{x}_i[t]) \right) = \sum_{k,l,m,n} \frac{\partial \ell}{\partial z_{t,k}} \frac{\partial g^{(o)}}{\partial h_{t,l}} \frac{\partial g^{(h)}}{\partial a_{t,m}} \left( \frac{\partial a_{t,m}}{\partial h_{t-1,n}} \frac{\partial h_{t-1,n}}{\partial W_{ij}^{(x)}} + \frac{\partial a_{t,m}}{\partial W_{ij}^{(x)}} \right)$$

$$\begin{cases} \frac{\partial g^{(o)}}{\partial h_{t,l}} = \sum_p \frac{\partial g^{(o)}}{\partial a_{t,p}^{(o)}} W_{lp}^{(o)} \\ \frac{\partial a_{t,m}}{\partial h_{t-1,n}} = W_{mn}^{(h)} \\ \frac{\partial a_{t,m}}{\partial W_{ij}^{(x)}} = \delta_{mi} \mathbf{x}_j[t-1] \\ \frac{\partial h_{t-1,n}}{\partial W_{ij}^{(x)}} = \sum_{q,r} \frac{\partial g^{(h)}}{\partial a_{t,q}} W_{q,r}^{(h)} \frac{\partial h_{t-2,r}}{\partial W_{ij}^{(x)}} + \frac{\partial g^{(h)}}{\partial a_{t,i}} \mathbf{x}_j[t-1] \end{cases}$$

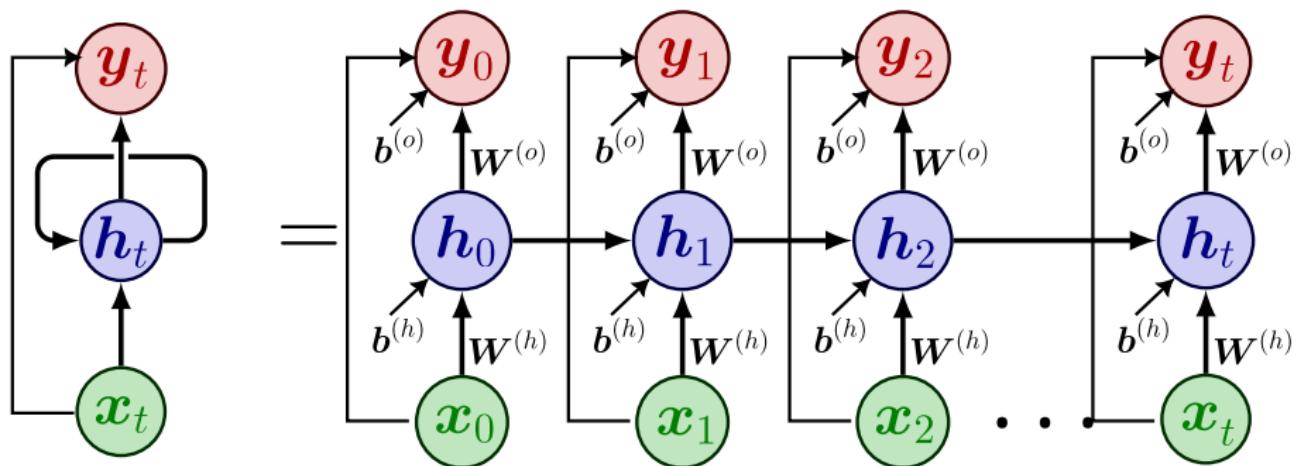
Non-local gradient graph.

- Back-Propagation Through Time (BPTT [Sal+18])

$$\left\{ \begin{array}{l} \nabla_{W^{(x)}} \ell(y_i[t], h_\theta(x_i[t])) = \nabla_{h_t} \ell \cdot \left( \sum_{\tau=1}^t (h_t \otimes \nabla_{h_\tau}) \otimes \nabla_{W^{(x)}} h_\tau \right) \\ \nabla_{h_\tau} h_t = \prod_{s=\tau+1}^t \frac{\partial h_s}{\partial h_{s-1}} \\ \frac{\partial h_s}{\partial h_{s-1}} = W^{(h)T} \left( \sum_k \frac{\partial g^{(o)}}{\partial h_{s-1,k}} e_k \otimes e_k \right) \end{array} \right.$$

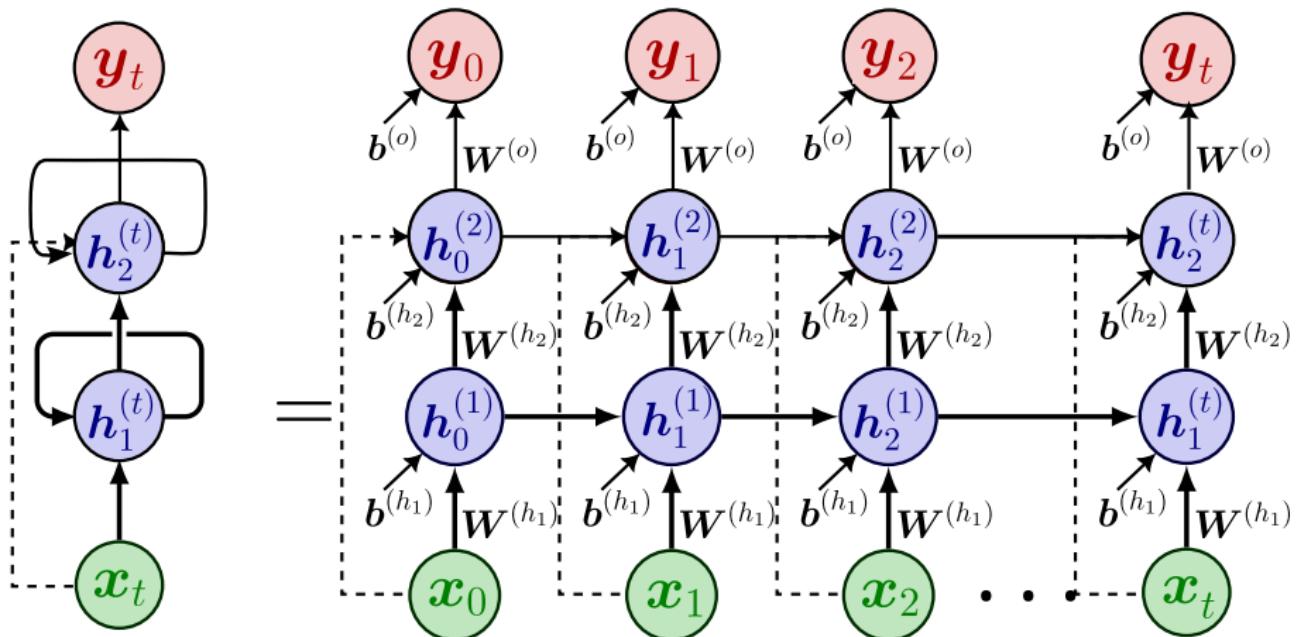
- “immediate”  $\nabla_{W^{(x)}} h_\tau$
- “short-term”  $W^{(h)T} \left( \sum_k \frac{\partial g^{(o)}}{\partial h_{s-1,k}} e_k \otimes e_k \right)$ ,  $\left\| \sum_k \frac{\partial g^{(o)}}{\partial h_{s-1,k}} e_k \otimes e_k \right\| \leq \gamma$
- $\rho(W^{(h)}) = \max_i |\lambda_i| < \frac{1}{\delta}$ ,  $\lambda_i \in \text{Spec}(W^{(h)})$
- “long-term”  $\prod_{s=\tau+1}^t \frac{\partial h_s}{\partial h_{s-1}} \rightarrow \left\| \frac{\partial h_s}{\partial h_{s-1}} \right\| \leq \frac{\gamma}{\delta}$
- $\left\| \nabla_{W^{(x)}} \ell(y_i[t], h_\theta(x_i[t])) \right\| \leq \left(\frac{\gamma}{\delta}\right)^{t-\tau} \left\| \nabla_{h_t} \ell \right\|$  Vanishing gradients !

# Other types of $\mathcal{RNN}$ : input-to-hidden



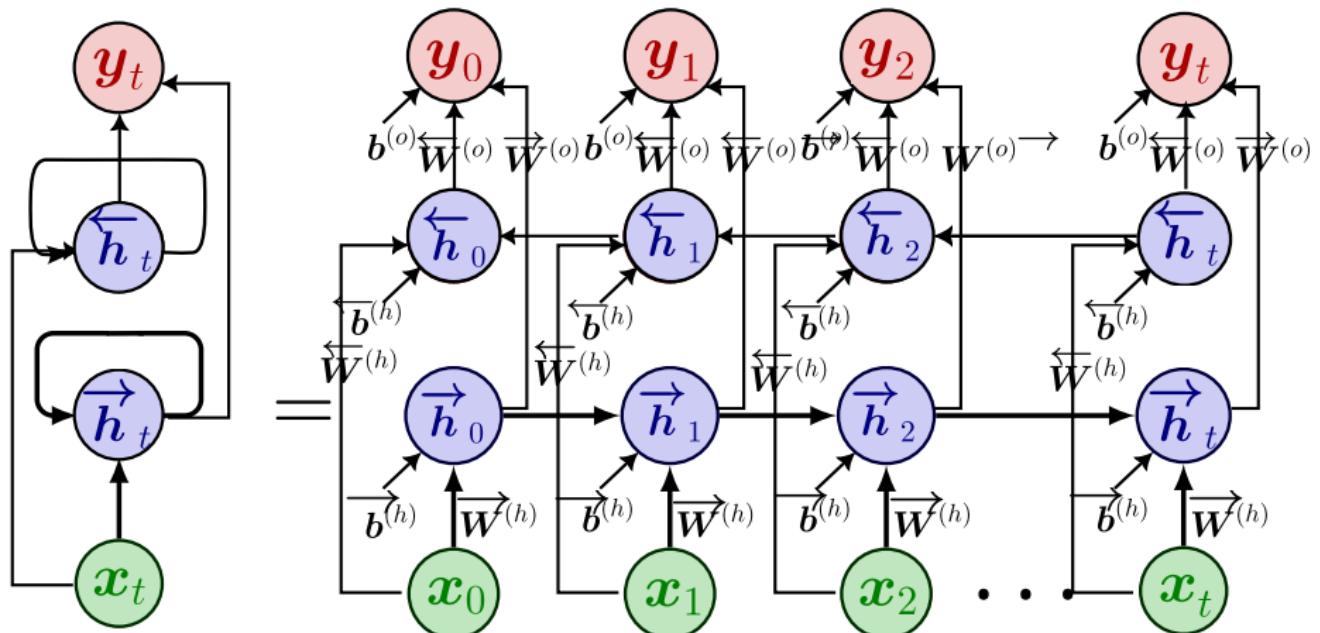
**Figure** – Sketch of the  $\mathcal{RNN}$  version with input-to-hidden connection.

# Other types of $\mathcal{RNN}$ : hidden-to-hidden

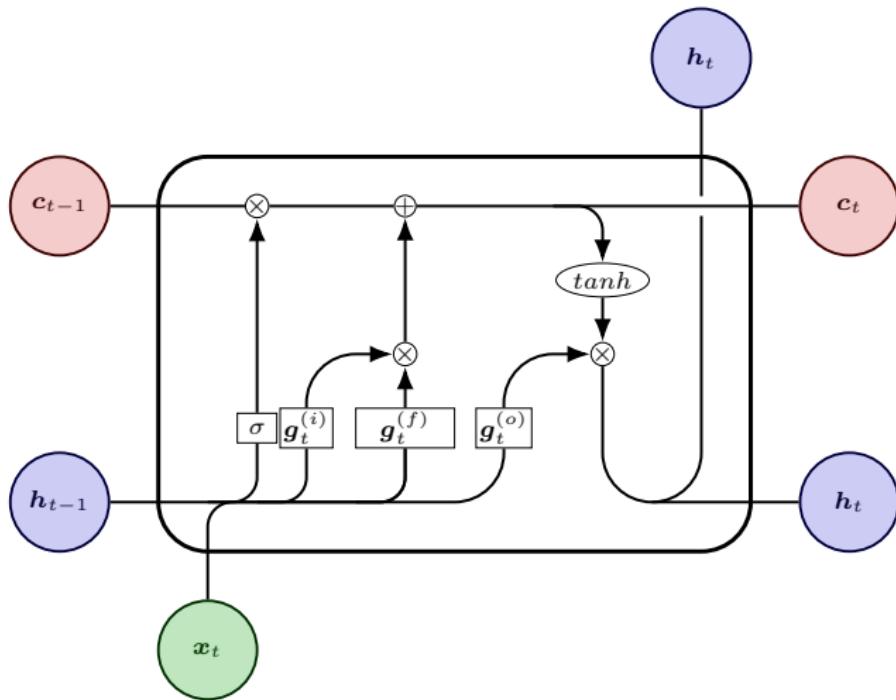


**Figure – Sketch of the  $\mathcal{RNN}$  version with hidden-to-hidden connection.**

# Other types of $\mathcal{RNN}$ : bidirectional



**Figure** – Sketch of the  $\mathcal{BRNN}$  version with hidden-to-hidden connections.



**Figure** – Sketch of the standard  $\mathcal{LSTM}$  architecture.

## Ideas

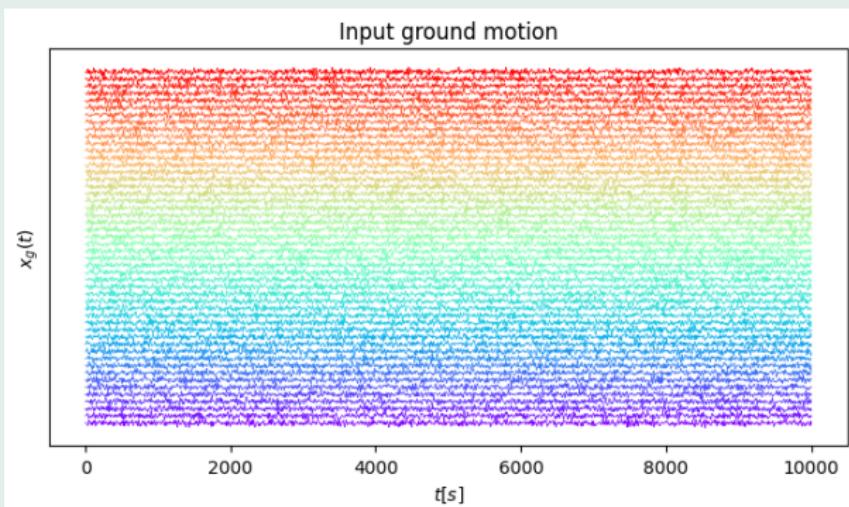
- Add “memory cells”  $c_t$  with “gates”
- Introduce a forget gate  $g^{(f)}$  for past hidden states

$$\begin{cases} z_t^{(i)} = \mathbf{g}^{(i)} \left( \mathbf{W}^{(ii)} \mathbf{x}[t] + \mathbf{W}^{(ih)} \mathbf{h}_{t-1} + \mathbf{W}^{(ic)} c_{t-1} + \mathbf{b}^{(i)} \right) \\ z_t^{(f)} = \mathbf{g}^{(f)} \left( \mathbf{W}^{(fi)} \mathbf{x}[t] + \mathbf{W}^{(fh)} \mathbf{h}_{t-1} + \mathbf{W}^{(fc)} c_{t-1} + \mathbf{b}^{(f)} \right) \\ c_t = \sum_j z_{t,j}^{(i)} \cdot \tanh \left( \langle \mathbf{w}_j^{(ci)}, \mathbf{x}[t] \rangle + \langle \mathbf{w}_j^{(ch)}, \mathbf{h}_{t-1} \rangle + b_j^{(c)} \right) \mathbf{e}_j + \sum_j z_{t,j}^{(f)} \cdot c_{t-1,j} \mathbf{e}_j \\ z_t^{(o)} = \mathbf{g}^{(o)} \left( \mathbf{W}^{(oi)} \mathbf{x}[t] + \mathbf{W}^{(oh)} \mathbf{h}_{t-1} + \mathbf{W}^{(oc)} c_t + \mathbf{b}^{(o)} \right) \\ \mathbf{h}_t = \sum_j z_{t,j}^{(o)} \cdot \tanh \left( z_{t,j}^{(c)} \right) \mathbf{e}_j \end{cases}$$

## Predict the non-linear response of a MDOF with $\mathcal{LSTM}$

$$\ddot{\mathbf{M}}\ddot{\mathbf{y}}(t) + \mathbf{C}\dot{\mathbf{y}}(t) + \mathbf{F}(\mathbf{y}(t), \dot{\mathbf{y}}(t)) = -\ddot{\mathbf{M}}\boldsymbol{\Gamma}\ddot{x}_G(t)$$

$$\dot{f}_i(t) = k_i(\dot{y}_{i+1}(t) - \dot{y}_i(t)) - \alpha_i |\dot{y}_{i+1}(t) - \dot{y}_i(t)| \cdot |f_i(t)|^{n_i-1} \cdot f_i(t) - \beta_i (\dot{y}_{i+1}(t) - \dot{y}_i(t)) \cdot |f_i(t)|^{n_i-1}$$





# Recap of basic probability

## Probability measure

$\mathbf{X} : (\Omega, \mathcal{E}, \mathbb{P}) \rightarrow (X, \Xi)$ , with a  $\sigma$ -algebra  $\mathcal{E} \subset \mathcal{B}(\mathbb{R})$  on  $\Omega$

$$\forall A \in \mathcal{E} \quad P_X(A) = \mathbb{P}(\mathbf{X}^{-1}(A)) = \mathbb{P}(\mathbf{X} \in A)$$

## Absolute continuity

$P_X$  is  $\sigma$ -finite measure dominated by a  $\sigma$ -finite measure  $\mu$ , i.e. iff  
 $\forall \varepsilon > 0, \exists \delta(\varepsilon)$  such that  $\forall A \in \Xi$  such that  $\mu(A) < \delta(\varepsilon)$

$$P_X(A) < \varepsilon$$

## Theorem

### **Radon-Nikodym theorem**

For two  $\sigma$ -finite measures on a  $\sigma$ -algebra  $\Xi$ , namely  $\mu$  and  $P_X$  such that  $P_X$  is absolutely continuous with respect to  $\mu$ ,  $\exists p_X \in L^1(\mu)$  such that

$$P_X(A) = \int_A p_X(x) \cdot \mu(dx), \quad \forall A \in \Xi$$

$p_X$  is called probability density of  $P_X$  and it corresponds to the Radon-Nikodym derivative  $p_X = \frac{dP_X}{d\mu}$ .

## Theorem

### **Strong law of large numbers (Kolmogorov)**

Given a set of random variables  $(X_n)_{n=1}^N$  independent identically distributed (i.i.d.) and Lebesgue integrable, with expected value  $\mu_X < +\infty$ , and a sample average  $\bar{X}_N = \frac{1}{N} \sum_{n=1}^N X_n$ , then :

$$\mathbb{P} \left[ \lim_{N \rightarrow +\infty} \bar{X}_n - \mu_X \right] = 1$$

which means that the sample average  $\bar{X}_N$  converges almost surely to the expected value  $\mu$ .

## Theorem

### **Weak law of large numbers (Khintchine)**

Given a sequence of random variables  $(X_i)_{i=1}^N$  independent identically distributed (i.i.d.) and Lebesgue integrable, with expected value  $\mu_X < +\infty$ , and a sample average  $\bar{X}_N = \frac{1}{N} \sum_{i=1}^N X_i$ , then :

$$\forall \varepsilon > 0, \quad \lim_{N \rightarrow +\infty} \mathbb{P} [|\bar{X}_n - \mu_X| \leq \varepsilon] = 1$$

which means that the sample average  $\bar{X}_N$  converges in probability to  $\mu$

## Theorem

Given a sequence of random variables  $(X_n)_{n=1}^N$  that converges in probability to  $X$ , then any continuous function  $f: \mathcal{X} \rightarrow \mathbb{R}$  converges in probability too.

## Consistent estimator

An estimator  $\theta_N = \theta(X_1, \dots, X_N)$  is said to be *consistent* if it converges in probability to its limit  $\theta$ .

## Remark

$\bar{X}_N = \frac{1}{N} \sum_{i=1}^N X_i$  is a consistent estimator of  $\mu_X$ .  $\bar{X}_N^2$  converges in probability to  $\mu_X^2$  and the estimator  $\bar{Y}_N = \frac{1}{N} \sum_{i=1}^N X_i^2$  converges in probability to  $\mathbb{E}[X_i^2]$ . Therefore, the variance estimator of each  $X_i$   $s_N^2$  that reads :

$$s_N^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X}_N)^2$$

converges in probability to the variance of  $X_i$ ,  $\sigma_{X_i}^2 = \mathbb{V}[X_i]$

## Theorem

### **Central limit theorem**

Given a sequence of i.i.d. random variables  $X = \{X_i\}_{i=1}^N$  with expected value  $\mathbb{E}[X_i] = \mu$  and variance  $\mathbb{V}[X_i] = \sigma^2 < +\infty$  and a random variable  $\bar{X}_N = \frac{1}{N} \sum_{i=1}^N X_i$ , the random variable

$$Z_N = \sqrt{N} \frac{\bar{X}_N - \mu}{\sigma}, \quad \mathbb{E}[Z_N] = 0, \mathbb{V}[Z_N] = 1$$

$$p(Z_N) \xrightarrow{n \rightarrow \infty} \mathcal{N}(0, 1)$$

## Statistical model

$\mathcal{H}_\Theta := \{P_\theta, \theta \in \Theta\}$ , if  $P \in \mathcal{H}_\Theta$ , there exist  $\theta^*$  such that  $P_{\theta^*} = P$ .  
 $\theta^*$  is unique if and only if  $\mathcal{H}_\Theta$  is *identifiable*, i.e.,  $\theta \mapsto P_\theta$  is injective.

## Likelihood maximization

Idea : find  $\hat{\theta} \approx \theta^*$  maximizing the *likelihood*

$$\hat{\theta}(\mathcal{D}_X) = \arg \max_{\theta \in \Theta} p_\theta(\mathbf{x} | \mathbf{x} \in \mathcal{D}_X) = \arg \max_{\theta \in \Theta} = \prod_{i=1}^N p_\theta(\mathbf{x}_i)$$

- $\hat{\theta}(\mathcal{D}_X)$  consistent (i.e. it must converge in probability to  $\theta^*$ )
- $\mathcal{D}_X$  i.i.d. and *exhaustive*, i.e., sufficient to characterize  $P = P_{\theta^*}$

C1  $\Theta$  is an open set and

$$\forall \mathbf{x} \in \mathcal{X}, \forall (\boldsymbol{\theta}, \boldsymbol{\theta}') \in \Theta^2 \quad p_{\boldsymbol{\theta}}(\mathbf{x}) > 0 \iff p_{\boldsymbol{\theta}'}(\mathbf{x}) > 0$$

$p_{\boldsymbol{\theta}} \in \mathcal{H}_{\Theta}$  have the same support  $\text{supp}(p_{\boldsymbol{\theta}})$ ,  $\forall \boldsymbol{\theta} \in \Theta$ .

C2  $\forall \boldsymbol{\theta} \in \Theta$ ,  $p_{\boldsymbol{\theta}}$  can be differentiated under the integral<sup>2</sup>:

$$\nabla_{\boldsymbol{\theta}} \int_{\text{supp}(p_{\boldsymbol{\theta}})} p_{\boldsymbol{\theta}} \cdot \mu(d\mathbf{x}) = \int_{\text{supp}(p_{\boldsymbol{\theta}})} \nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}} \cdot \mu(d\mathbf{x}) \quad (2)$$

---

2. According to [Bil95], C2 holds if the gradient is locally dominated by an integrable function  $g$ , i.e., it exist a neighborhood  $\square$  of  $\boldsymbol{\theta}$  and  $g$  such that  $\int_{\text{supp}(p_{\boldsymbol{\theta}})} g(x)\mu(dx) < +\infty$  such that almost everywhere on a neighborhood  $V$  of  $\mathbf{x}$ :

$$\left| \frac{\partial p_{\boldsymbol{\theta}}}{\partial \theta_k} \right| \leq g$$

## Theorem

### **Maximum Likelihood**

Given a sequence of i.i.d. random variables  $\mathcal{D}_X = \{\mathbf{x}_i \in \mathcal{X}\}_{i=1}^N$ , then :

$$\forall (\theta, \theta^*) \in \Theta, \theta \neq \theta^*, \forall \varepsilon > 0 \quad \lim_{N \rightarrow +\infty} \mathbb{P} [p_\theta(\mathcal{D}_X) - p_{\theta^*}(\mathcal{D}_X) < \varepsilon] = 1$$

- if  $p \in \mathcal{H}_\Theta$ ,  $\exists \theta^* \in \Theta$  such that  $p = p_{\theta^*} = \arg \max_{\theta \in \Theta} p_\theta(\mathbf{X})$
- $\hat{\theta}(\mathcal{D}_X) = \arg \max_{\theta \in \Theta} p_\theta(\theta; \mathcal{D}_X) = \arg \max_{\theta \in \Theta} \sum_{i=1}^N p_\theta(\mathbf{x}_i)$

# Maximum log-likelihood

## Jensen inequality

$$\mathbb{E}_{\mathbf{x} \sim p_{\theta^*}} \left[ \ln \left( \frac{p_{\theta} (\mathbf{X})}{p_{\theta^*} (\mathbf{X})} \right) \right] < 0$$

## Theorem

**Maximum Likelihood of bijection  $g(\theta)$ .** [Bil95]

Given a sequence of i.i.d. observations  $\mathcal{D}_X = \{\mathbf{x}_i \in X\}_{i=1}^N$  and  $\theta$  the of parameters that maximizes the likelihood  $p_{\theta}$ . Then  $g(\hat{\theta})$  is the maximum likelihood estimator of  $g(\theta)$  ( $g$  bijection).

# Maximum Log-Likelihood

## Theorem

**Maximum Log-Likelihood** Given a sequence of i.i.d. random variables

$\mathcal{D}_X = \{\mathbf{x}_i \in \mathcal{X}\}_{i=1}^N$ , then :

$$\forall \theta \neq \theta^*, \forall \varepsilon > 0 \quad \lim_{N \rightarrow +\infty} \mathbb{P} \left[ \frac{1}{N} \ln \left( \frac{p_\theta(\mathcal{D}_X)}{p_{\theta^*}(\mathcal{D}_X)} \right) < \varepsilon \right] = 1$$

- if  $p \in \mathcal{H}_\Theta$ ,  $\exists \theta^* \in \Theta$  such that  $p = p_{\theta^*} = \arg \max_{\theta \in \Theta} \ln p_\theta(\mathbf{X})$
- $\hat{\theta}(\mathcal{D}_X) = \arg \max_{\theta \in \Theta} \mathcal{L}(\theta; \mathcal{D}_X) = \arg \max_{\theta \in \Theta} \sum_{i=1}^N \ln p_\theta(\mathbf{x}_i)$

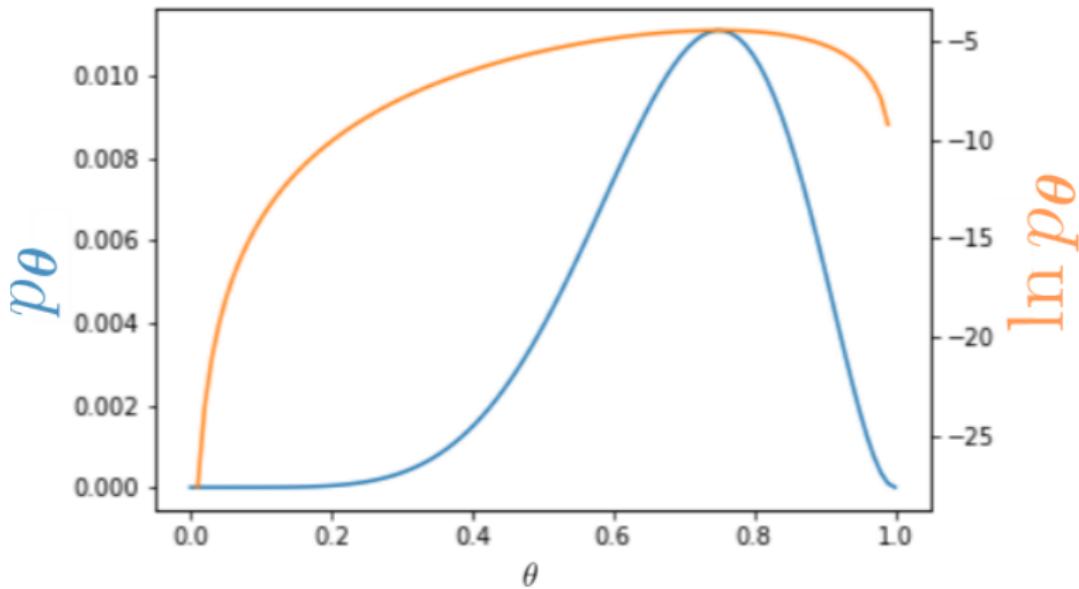
# “True” probability distribution

Converging sequence to maximum log-likelihood

- Open neighborhood  $O_\varepsilon(\theta^*)$
- Sequence  $\theta_N := \theta(\mathcal{D}_X)$ ,  $\mathcal{D}_X = \{\mathbf{x}_i \in X\}_{i=1}^N$

$$S_N := \left\{ \mathbf{x}_i \in \mathcal{D}_X \mid \ln p_{\theta^*}(\mathbf{x}_i) > \sup_{\theta \in \partial O_\varepsilon(\theta^*)} \ln p_\theta(\mathbf{x}) \right\}$$

$$\forall \varepsilon > 0 \quad \lim_{N \rightarrow +\infty} \mathbb{P}[|\mathbb{P}(S_N) - 1| < \varepsilon] = 1 \implies \theta^* \in \text{Int}(\Theta)$$



**Figure** – Bernoulli distribution  $p_\theta(\mathcal{D}_X) = \prod_{i=1}^N \theta^{1_{x_i=1}} (1 - \theta)^{1_{x_i=0}}$ .

[stats.stackexchange.com/questions/486007/measuring-predictive-uncertainty-with-negative-log-likelihood-nll](https://stats.stackexchange.com/questions/486007/measuring-predictive-uncertainty-with-negative-log-likelihood-nll)

## Converging sequence to maximum log-likelihood

- Open neighborhood  $O_\varepsilon(\theta^*)$ ,  $\theta_N := \theta(\mathcal{D}_X)$
- Regularity C2

$$\nabla_{\theta} \int_{\text{supp}(p_{\theta})} p_{\theta} \cdot \mu(dx) = \int_{\text{supp}(p_{\theta})} \nabla_{\theta} p_{\theta} \cdot \mu(dx)$$

- Rolle’s theorem :  $\exists \hat{\theta}_N \in \bar{O}_\varepsilon(\theta^*)$  such that  $\nabla_{\theta} \ln p_{\hat{\theta}_N}(x) = \mathbf{0}$
- Convergence

$$S_N \subset \tilde{S}_N := \left\{ x_i \in \mathcal{D}_X \mid \exists \hat{\theta}_N \in \bar{O}_\varepsilon(\theta^*) \text{ such that } \nabla_{\theta} \ln p_{\hat{\theta}_N}(x_i) = \mathbf{0} \right\}$$

$$\mathbb{P}(S_N) \leq \mathbb{P}(\tilde{S}_N) \leq 1 \quad \forall \varepsilon > 0 \quad \lim_{N \rightarrow +\infty} \mathbb{P} \left[ |\mathbb{P}(\tilde{S}_N) - 1| < \varepsilon \right] = 1$$

## Recap

- $\exists (\hat{\theta}_N)_{N \in \mathbb{N}^*} = \hat{\theta}(\mathcal{D}_X) \xrightarrow[N \rightarrow \infty]{\mathbb{P}} \theta^*$
- If C2 holds :  $s_{\hat{\theta}_N}(\mathcal{D}_X) = \nabla_{\theta} p_{\hat{\theta}_N} = \mathbf{0}$
- Under C1 and C2 holds :

$$\mathbb{E}_{x \sim p_{\theta^*}} \nabla_{\theta} \ln p_{\theta^*}(\mathbf{X}) = \mathbb{E}_{x \sim p_{\theta^*}} s_{\theta^*}(\mathbf{X}) = \mathbf{0}$$

- $\ell(\theta; \mathcal{D}_X) = \ln p_{\theta}(\mathcal{D}_X)$

Fisher Information Matrix at  $\theta^*$  : auto-covariance of the score

$$\mathbb{I}_F(\theta^*; \mathbf{X}) = \mathbb{E}_{x \sim p_{\theta^*}} \left[ \nabla_{\theta} \ln p_{\theta^*}(\mathbf{X}) \otimes \nabla_{\theta} \ln p_{\theta^*}(\mathbf{X}) \right] \geq 0$$

If  $x_i$  i.i.d. :  $\mathbb{I}_F(\theta^*; \mathcal{D}_X) = N \mathbb{I}_F(\theta^*; X_1)$

Meaning of the FIM (under regularity conditions of  $p_{\theta}$ )

$$\mathbb{I}_F(\boldsymbol{\theta}^*; \mathcal{D}_X) = -\mathbb{E}_{\mathbf{x} \sim p_{\boldsymbol{\theta}^*}} [\mathbf{H}_\ell(\boldsymbol{\theta}^*; \mathcal{D}_X)] \quad (p_{\boldsymbol{\theta}} \in C^2(\Theta))$$

### Theorem

**Cramér-Rao bound** : statistical algorithm to estimate  $\boldsymbol{\theta}^*$ , that reads  $\boldsymbol{\theta} = \mathbf{A}(\mathbf{x})$  :

$$\bar{\boldsymbol{\theta}} = \mathbb{E}_{\mathbf{x} \sim p_{\boldsymbol{\theta}^*}} [\mathbf{A}(\mathbf{X})]$$

and the variance of the estimator is bounded by below as follows :

$$\mathbb{V}_{\mathbf{x} \sim p_{\boldsymbol{\theta}^*}} (\mathbf{A}(\mathbf{X})) \geq \frac{\|\nabla_{\boldsymbol{\theta}^*} \bar{\boldsymbol{\theta}}\|^2}{N \cdot \mathbb{I}_F(\boldsymbol{\theta}^*; X_1)}$$

$$\mathbf{A} \text{ is unbiased if : } \mathbb{V}_{\mathbf{x} \sim p_{\boldsymbol{\theta}^*}} (\mathbf{A}(\mathbf{X})) = \mathbb{E}_{\mathbf{x} \sim p_{\boldsymbol{\theta}^*}} \left[ (\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}^*)^2 \right] \geq \frac{1}{N \cdot \mathbb{I}_F(\boldsymbol{\theta}^*; X_1)}$$

# CRB as an optimum solution

## Theorem

**MLE convergence in standard normal probability [Bil95].**

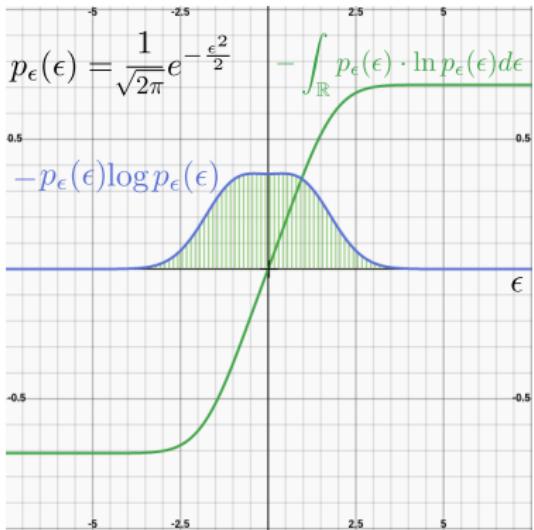
$$\forall \varepsilon > 0 \quad \lim_{N \rightarrow +\infty} \mathbb{P} \left[ |p_{\sqrt{N}(\hat{\theta}_N - \theta^*)} - \mathcal{N} \left( \mathbf{0}, \mathbb{I}_F^{-1} (\theta^*; \mathcal{D}_X) \right)| < \varepsilon \right] = 1$$

*The MLE is an asymptotically “optimum” estimator :*

$$\forall \varepsilon > 0 \quad \lim_{N \rightarrow \infty} \mathbb{P} \left[ \left| \frac{1}{N \cdot \mathbb{I}_F(\theta^*; X_1) \cdot \mathbb{V}_{x \sim p_{\theta^*}} (A(X))} - 1 \right| < \varepsilon \right] = 1$$

# Shannon's differential entropy (Jaynes correction)

$$\mathbb{H}_d(p(\mathbf{X})) = - \int_{\mathcal{X}} p(\mathbf{x}) \cdot \ln\left(\frac{p(\mathbf{x})}{m(\mathbf{x})}\right) \cdot \mu(d\mathbf{x})$$



$$\mathbb{H}_d(\mathcal{N}_{(\mu, \Sigma)}(\mathbf{X})) = \frac{d_X}{2} \ln 2\pi + \frac{1}{2} \text{Tr}(\Sigma)$$

$$(d_X = 1, \Sigma = 1 \rightarrow \mathbb{H}_d = 1.418938533204673)$$

# Some properties of Shannon's entropy

- The conditional entropy of two discrete random variables  $X$  and  $Y$  reads :

$$\mathbb{H}_d(p(\mathbf{Y}|\mathbf{X} = \mathbf{x}_i)) = -\mathbb{E}_{\mathbf{y} \sim p(y|x)} [\ln(p(\mathbf{Y}|X = x_i))]$$

and

$$\mathbb{H}_d(p(\mathbf{Y}|\mathbf{X})) = -\mathbb{E}_{\mathbf{x} \sim p(x)} [\mathbb{H}_d(p(\mathbf{Y}|\mathbf{X} = \mathbf{x}))] = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y}_i)} [\ln p(\mathbf{Y}|\mathbf{X})]$$

- The joint entropy of two discrete random variables  $X$  and  $Y$  reads :

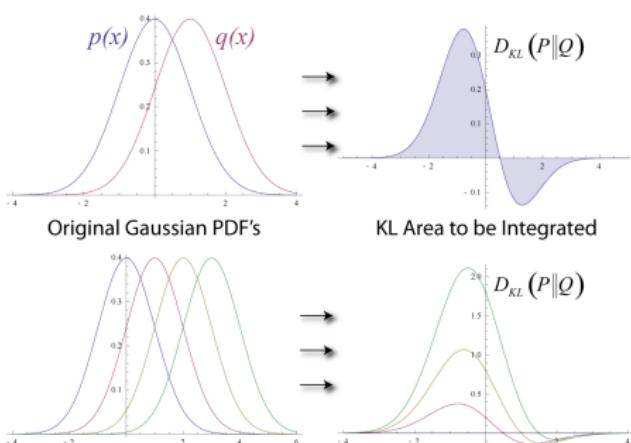
$$\mathbb{H}_d(p(\mathbf{X}, \mathbf{Y})) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y})} [\ln(p(\mathbf{X}, \mathbf{Y}))]$$

$$\mathbb{H}_d(p(\mathbf{X}, \mathbf{Y})) = \mathbb{H}_d(p(\mathbf{Y}|\mathbf{X})) + \mathbb{H}_d(p(\mathbf{X})) = \mathbb{H}_d(p(\mathbf{X}|\mathbf{Y})) + \mathbb{H}_d(p(\mathbf{Y}))$$

# Some properties of Shannon's entropy

- Kullback-Leibler divergence :

$$0 \leq \mathbb{D}_{KL}(p\|q) = \mathbb{E}_{x \sim p} \left[ \ln \frac{p}{q} \right] < +\infty \quad \mathbb{D}_{KL}(p\|q) \geq \mathbb{E}_{x \sim p} \left[ \frac{p - q}{p} \right]$$



**Figure –**  $\mathbb{D}_{KL}(p\|q) = 0$  implies that  $p = qa.e.$  if and only if  $\text{supp}(p) \cap \text{supp}(q) \neq \emptyset$

- Mutual information :

$$\mathbb{I}(p(\mathbf{X}, \mathbf{Y})) = \mathbb{D}_{KL} (p(\mathbf{X}, \mathbf{Y}) \| p(\mathbf{X}) \cdot p(\mathbf{Y})) \geq 0$$

$$\mathbb{H}(p(\mathbf{Y})) - \mathbb{H}(p(\mathbf{Y}|\mathbf{X})) = \mathbb{I}(p(\mathbf{X}, \mathbf{Y})) = \mathbb{H}(p(\mathbf{X})) - \mathbb{H}(p(\mathbf{X}|\mathbf{Y}))$$

- $\sup_{\mathbf{Y}} \mathbb{I}(p(\mathbf{X}, \mathbf{Y})) = \mathbb{H}_d(p(\mathbf{X})) - \inf_{\mathbf{Y}} \mathbb{H}_d(p(\mathbf{Y}|\mathbf{X})) \rightarrow 0$ , i.e., forcing the *disentanglement* between  $\mathbf{X}$  and  $\mathbf{Y}$  [Che+16].

# MaxEnt principle

## Theorem

**Gibbs-Boltzmann theorem :** available observations  $(\mathbf{x}_k, y_k)$

$$\max_{q \in \mathcal{H}} \mathbb{H}_d(q) = \max_{q \in \mathcal{H}} \left( - \int_{\mathcal{X}} \ln q(\mathbf{x}) \cdot q(\mathbf{x}) d\mathbf{x} \right) \quad (3)$$

$$\left\{ \begin{array}{l} c_k(q) = \mu_{y_k} - \int_{\mathcal{X}} y_k(\mathbf{x}) \cdot q(\mathbf{x}) d\mathbf{x} = 0, \quad c_k : \mathbb{R}^n \rightarrow \mathbb{R}, 1 \leq k \leq K \end{array} \right. \quad (4)$$

if it exists it is defined as  $p_\theta \in \mathcal{H}$ , which reads :

$$p_\theta(\mathbf{x}) = \arg \max_{q \in \mathcal{H}, \mathbf{c} : \mathbf{c} = \mathbf{0}} \mathbb{H}_d(q) = \frac{e^{-\sum_{k=1}^K y_k(\mathbf{x}) \cdot \theta_k}}{Z}$$

with  $Z = \int_{\mathcal{X}} e^{-\sum_{k=1}^K \theta_k \cdot y_k(\mathbf{x})} \cdot \mu(d\mathbf{x})$  a normalization constant. Moreover,  
 $\mathbb{H}_d(p_\theta) \geq \mathbb{H}_d(p)$ .

$$p_{\theta}(\mathbf{x}) = \frac{e^{-\sum_{k=1}^K y_k(\mathbf{x}) \cdot \theta_k}}{Z}$$

## Log-Likelihood

$$\ln p_{\theta}(\mathbf{x}) = -\ln Z - \langle \mathbf{y}(\mathbf{x}), \boldsymbol{\theta} \rangle$$

$$\mathbf{s}_{\theta}(\mathbf{x}) = \mathbb{E}_{\mathbf{x} \sim p_{\theta}} [\mathbf{y}(\mathbf{X})] - \mathbf{y}(\mathbf{x}) = \boldsymbol{\mu}_{\mathbf{y}} - \mathbf{y}(\mathbf{x})$$

$$\mathbf{H}_{\ell}(\boldsymbol{\theta}; \mathbf{X}) = -\mathbb{C}_{\mathbf{x} \sim p_{\theta}} (\mathbf{y}(\mathbf{X}))$$

- Unbiased score :  $\mathbb{E}_{\mathbf{x} \sim p_{\theta}} [\mathbf{s}_{\theta}(\mathbf{X})] = \mathbf{0}$
- FIM is the covariance of the observations  $\mathbf{y}(\mathbf{X})$
- $\boldsymbol{\theta}_N$  convergence in  $\mathcal{N}(\mathbf{0}, \mathbb{I}_F^{-1}(\boldsymbol{\theta}^*; \mathcal{D}_X)) \in \mathcal{H}_{\theta}$

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta} \in \Theta} \max_{p_{\theta} \in \mathcal{H}_{\Theta}} \mathbb{H}_d(p_{\theta}(Y|\mathbf{X}))$$

# Bibliography

- [Bil95] **BILLINGSLEY, Patrick.** *Measure and Probability*. John Wiley et Sons : New York, 1995.
- [Che+16] **CHEN, Xi et al.** “InfoGAN : Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets”. en. In : arXiv :1606.03657 (juin 2016). arXiv :1606.03657 [cs, stat]. URL : <http://arxiv.org/abs/1606.03657>.
- [ES16] **ELDAN, Ronen et SHAMIR, Ohad.** “The power of depth for feedforward neural networks”. In : *Workshop and Conference Proceedings*. T. 49. PMLR, 2016, p. 1-34.
- [Heu+] **HEUSEL, Martin et al.** “GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium”. en. In : () .

- [Mai99] **MAIOROV, V.E.** “On Best Approximation by Ridge Functions”. en. In : *Journal of Approximation Theory* 99.1 (juill. 1999), p. 68-94. ISSN : 00219045. DOI : 10.1006/jath.1998.3304. URL : <https://linkinghub.elsevier.com/retrieve/pii/S0021904598933044>.
- [MP43] **McCULLOCH, Warren S. et PITTS, Walter.** “A logical calculus of the ideas immanent in nervous activity”. en. In : *The Bulletin of Mathematical Biophysics* 5.4 (déc. 1943), p. 115-133. ISSN : 0007-4985, 1522-9602. DOI : 10.1007/BF02478259. URL : <http://link.springer.com/10.1007/BF02478259>.

- [Ros57] **ROSENBLATT, Murray.** “Some Purely Deterministic Processes”. In : *Journal of Mathematics and Mechanics* 6.6 (1957), p. 801-810. ISSN : 00959057, 19435274. URL : <http://www.jstor.org/stable/24900623>.
- [Sal+18] **SALEHINEJAD, Hojjat et al.** “Recent Advances in Recurrent Neural Networks”. en. In : arXiv :1801.01078 (fév. 2018). arXiv :1801.01078 [cs]. URL : <http://arxiv.org/abs/1801.01078>.