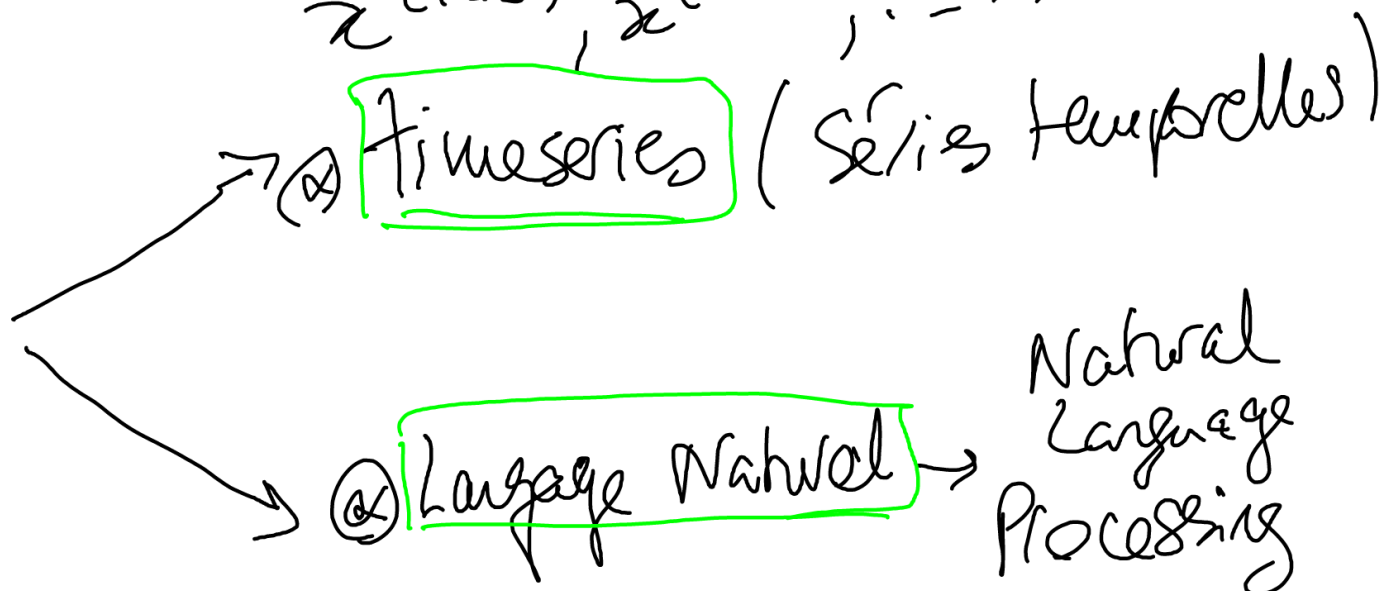# Day 4 : Recurrent Neural Networks (RNNs)
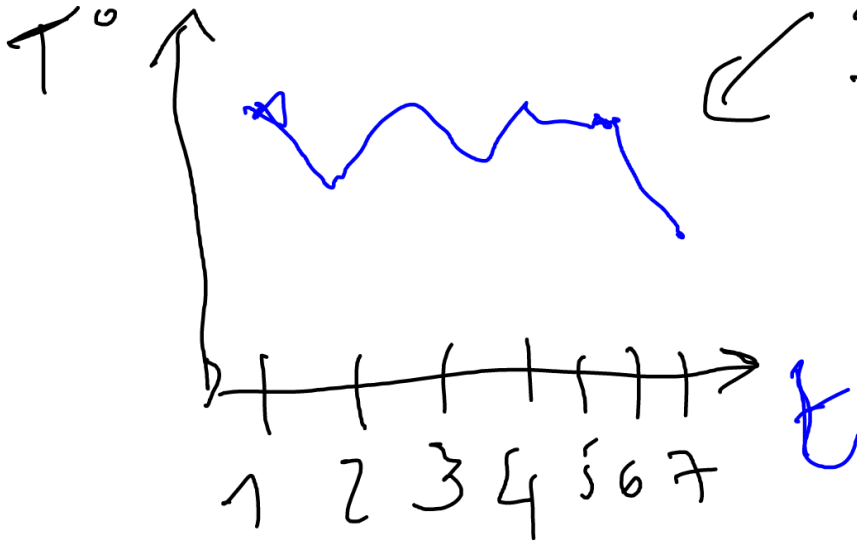## & Application to Sequence Modelling

RNN = special kind of NN for processing sequential data.

$$X = [x^{(1)} - - - - x^{(T)}]$$

each point $x^{(\iota)}$ of the sequence is dependant of the rest of the sequence points $x^{(i)}$

$x^{(t)}$ dependant of the $x^{(t-1)}, x^{(t-2)}, - - - -$

→ ⓐ timeseries ( Séries temporelles )

→ ⓐ Langage Natural → Natural Language Processing

$\otimes$ time series



T° [graph with blue curve]

1 2 3 4 5 6 7 → t

Série temporelle 1D.

$\hookrightarrow$ Multivariate time-series (séries temporelles multivariées)

variables :

(T°, $\rho_{air\ pressure}$, humidity rate, % de précipitation)

nombre de pas de temps

→ X input data (N, T, F)

nombre d'échantillons

nombre de variables.

F = 4 for example in the above example.

$\otimes$ langage Naturel (NLP)

ex: Sentences

Sample 1: 'The cat is on the mat' 6 words

Sample 2: 'Time flies like a arrow' 5 words

...

Sample n: 'I am fine today' 4 words

mapping (w_i : idx)

**Raw data**

└ "data cleaning"

→ ① Tokenization → vocabulaire are les mots

Sample 1:
↳ ['The', 'cat', 'is', 'on', 'the', 'mat']

└→ ② words → ids

Sample 1 → (1, 6, 3, 15, 63, 26)

└→ ③ padding

max sentence length = 10
Samples → (1,6,3,15,63,26,0,0,0,0)

④ Embedding

list_idx

$\boxed{embedding}$ Emb (list) $\in \mathbb{R}^D$.

text representation

① encode similarities between words meanings

② encode grammar

→ INPUT X (TENSOR --)

(N , T , D)
samples   d
        ex: 10 words          embedding size

✗ other sequential data types

↳ 'Speech'

↳ 'Music'

RNNs → represent the notion of sequentiality by having parameter sharing

↳ a RNN shares the same weights across several time-steps.

↳ Apply these models to data with ≠ input sequence lengths.

$$s^{(t)} = f\left(s^{(t-1)}, \theta\right)$$

↖ parameter θ

dynamical system driven by an external signal $x^{(t)}$

$$s^{(t)} = f\left(s^{(t-1)}, x^{(t)}, \theta\right)$$

# Typical equation of a RNN

$$h^{(t)} = f\left(h^{(t-1)}, x^{(t)}, \theta\right)$$

↑ Neural Network parameter

State of RNN

→ hidden representation of the RNN

↳ output of the RNN hidden layer.

↳ $h^{(t)} =$ 'lossy' summary of the input data until timestep $t$

$$\left(x_1 \cdots x^{(t)}\right)$$

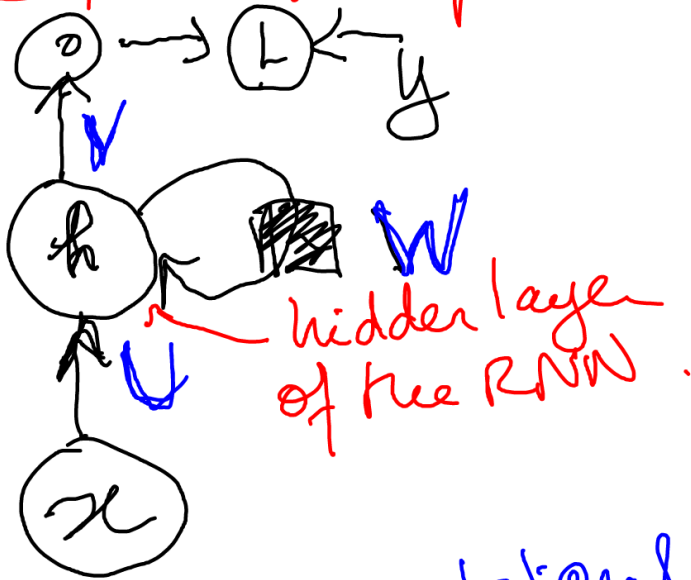→ Map an arbitrary 'long' sequence in a fixed-length vector.

↳ unfolded dynamical system.

$$h^{(t)} = g\left(x^{(t)}, x^{(t-1)}, x^{(t-2)}, \cdots x^{(1)}, h^0\right)$$

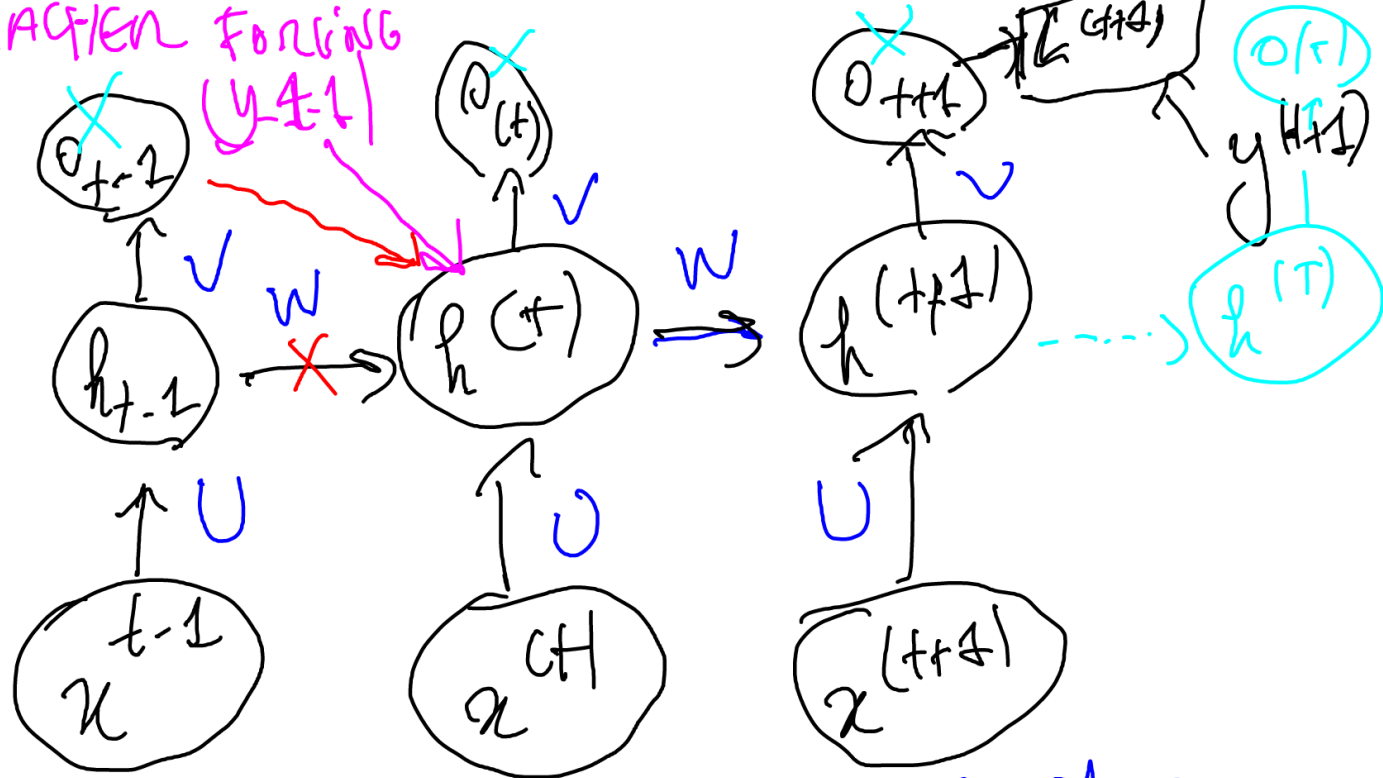↳ Generally, we use the same transition function $f$ with the same $\theta$ at every parameters timestep

# IB) Graph of Recurrent Neural Networks

$o \longrightarrow L \longleftarrow y$

$V$

$h \quad\quad W$

hidden layer of the RNN.

$U$

$x$

'Compact computational graph for RNN'

TEACHER FORCING $(y_{t-1})$

Unfolded computational graph of a RNN

Some ≠ types of RNNs

— RNN producing an output at each timestep & have recurrent connections between hidden units

① produce an output at each timestep but have only recurrent connections only from the output at timestep at the hidden units of the timestep

③ RNN with recurrent connections between hidden units → 'reads' an entire sequence and outputs a single output $o(c)$

## forward popagation equations for the RNN

input-to-hidden connection

⊗ RNN layer :
$$a^{(t)} = \boxed{U} x^{(t)} + \widehat{W} h^{(t-1)} + b$$
$$h^{(t)} = \tanh\left(a^{(t)}\right)$$

hidden-to-hidden connection.

⊗ dense layer : 
$$\underbrace{o^{(t)}}_{logits} = \boxed{V} \times h^{(t)} + c$$

↖ hidden-to-output connection

$$\left[\hat{y}^{(t)} = softmax\left(o^{(t)}\right)\right]$$

⊗ Associated Loss / Cost function

$$\mathcal{L}\left(\left\{x^{(1)} \ldots x^{(T)}\right\}, \left\{y^{(1)} \ldots y^{(T)}\right\}\right)$$
$$= \sum_{t \in \{1 \cdots T\}} \mathcal{L}^{(t)} = -\sum_{t \in \{1 \cdots T\}} \log P_{model}\left(y^{(t)} \mid \{x^{(1)} \cdots x^{(t)}\}\right)$$

# C/ Teacher Forcing

From RNNs that have output to hidden
recurrent connections
→ can be trained with a techni-
que called **Teacher Forcing :)**

→ replacing the output-to hidden
connection
by label-to- hidden connection

→ $h^{(t)}$ depend not of $o^{(t-1)}$ but
of $y^{(t-1)}$
ground truth ← prediction of the RNN

↳ <u>Advantages</u> :

ⓧ Improve performance by giving the <u>ground-truth</u> from previous timestep

ⓧ Ease training of the Recurrent Neural Network by avoiding the '<u>Backpropagation Through Time</u>' (BPTT) which compute gradients across all time steps

<u>Maximum likelihood criterion</u>

<u>2 time steps</u>   $(x_1, x_2)$   $(y_1, y_2)$
Input data

$\log \ P(y_1, y_2 \mid x_1, x_2)$

$= \log \ P(y_2 \mid y_1, x_1, x_2) + \log P(y_1 \mid x_1, x_2)$

↳ Cons

↳ discrepancy between what is done during training & what is done at inference

# D) The challenge of long-term Dependencies

"long-term dependencies"
$\hookrightarrow$ Input data with long sequences.
$\Longrightarrow$ IN THE TRAINING PROCESS.

pb of $\left\{\begin{array}{l} \rightarrow \text{Vanishing gradients} \\ \quad \hookrightarrow \text{gradient très petits} \\ \\ \hookrightarrow \text{exploding gradient} \\ \quad \hookrightarrow \text{gradient très grands.} \end{array}\right.$

## Intuition derrière

$$h^{(t)} = W^{T} h^{(t-1)}$$

$\hookrightarrow$ parameter sharing
Same W for each timestep. $\rightsquigarrow > 1$

$$h^{(t)} = \boxed{(W^t)^{T}} h^{(0)}$$

Let's assume that (to simplify) $\boxed{W^t = U \Lambda^t U}$

$$W = U \Lambda U^{T}$$

$\curvearrowleft$ matrice de valeurs propres.

$\begin{pmatrix} d_1 & & \\ & \ddots & \\ & & d_N \end{pmatrix}$

$$\Lambda = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_N \end{bmatrix}$$

$\Lambda$, diagonal matrix

$W^t \approx \Lambda^t$

$$(\lambda_1 \cdots \lambda_N)^t$$

$> 1 : (\lambda_1 \cdots \lambda_N)^t$ large values.

$< 1 : (\lambda_1 \cdots \lambda_N)^t$ small values

exploding gradients

Vanishing gradients

→ 'gradient clipping'

└→ $\text{clip}(\min(\nabla_\Sigma))$ a epsilon value if too large.

gradient clipped to

⟹ Recurrent Neural Networks with Gated Units

└→ LSTM

Long Short Term Memory Networks.

# II / LSTM (and other Gated Recurrent Neural Networks)

LSTM : In practice, one of the RNN the most used to process sequential data.

↳ Introduce a cell memory (recurrent self-loop) where the gradient can flow for long durations (long sequences)

Internal recurrence

addition to the outer recurrence

$$h^{(t+1)} = f(h^{(t)}, x^{(t)})$$

Compute internal gates

→ forget gate $f^{(t)}$
↳ external input gate $g^{(t)}$
→ output gate $q^{(t)}$

# Recurrent equations in a LSTM

## ⊘ forget gate

$$f_i^{(t)} = \sigma\left( b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)} \right)$$

$U^f x^{(t)}$

## ⊗ Internal state of the LSTM cell

$$\boxed{s_i^{(t)}} = \underbrace{\boxed{f_i^{(t)}}}_{\text{forget gate}} s_i^{(t-1)} + \boxed{g_i^{(t)}} \times \sigma\left( b_i^s + \sum_j U_{i,j}^s x_j^{(t)} + \sum_j W_{i,j}^s h_j^{(t-1)} \right)$$

External Input gate

## ⊗ External input gate

$$g_i^{(t)} = \sigma\left( b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)} \right)$$

## ⊘ output gate $q_i^{(t)}$

$$\boxed{q_i^{(t)}} = \sigma\left( b_i^q + \sum_j U_{i,j}^q x_j^{(t)} + \sum_j W_{i,j}^q h_j^{(t-1)} \right)$$

output $h_i^{(t)}$ (hidden representation of the
of the RNN)                                              LSTM
                                                         cell

$$h_i^{(t)} = \tanh\left(s_i^{(t)}\right) \times q_i^{(t)}$$

internal
recurrent state                          output
                                         gate

→ still a dynamical system

$$h^{(t)} \rightarrow f_\theta\left(h^{(t-1)}, x^{(t)}\right)$$