Sequence Modelling $\begin{cases}\end{cases}$ classification

regression

seq 2 seg $\begin{cases} \text{Machine Translation} \\ \text{Health Monitor} \end{cases}$

Sequence generation $\rightarrow \begin{cases} \text{write poem} \\ \text{compose music} \end{cases}$ (creative)
(no ground-truth)

sequence prediction
(have ground-truth)

Today's topic

---

① Idea: Parameter-sharing across time

E.g. $\begin{cases} \text{I went Nepal in 2007.} \\ \text{In 2007 I went to Nepal.} \end{cases} \longrightarrow$ Same meaning
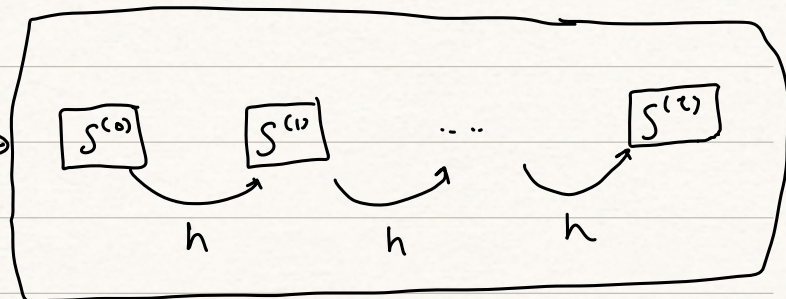
② Dynamic System

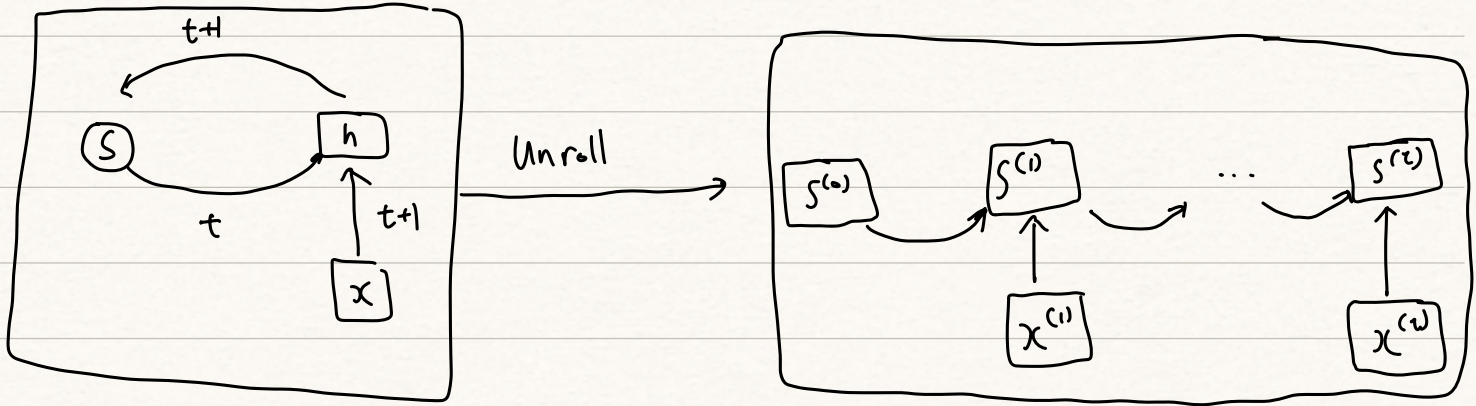a) $$S^{(t+1)} = h(S^{(t)} : \theta)$$



Unroll

b) $\quad S^{(t+1)} = h(S^{(t)}, x^{(t+1)}; \theta)$



Unroll

③ <u>Neural Network formulation</u>

<u>Recap</u>: $\begin{cases} S^{(t+1)} = h(S^{(t)}, x^{(t+1)}; \theta) \\[2mm] y^{(t)} = g(S^{(t)}; \varphi) := O^{(t)} \end{cases}$

time-t memory

current (t+1) observation

$\boxed{NN}$ $\begin{cases} S^{(t+1)} = \sigma_r(W\boxed{S^{(t)}} + U\boxed{x^{(t+1)}} + b) \\[3mm] y^{(t)} = \sigma_o(V S^{(t)} + c) \end{cases}$

$\boxed{\text{Elman Variant}}$

tanh usually (not ReLu)

$\boxed{\begin{cases} \theta = (W, U, b; \sigma_r) \\[2mm] \varphi = (V, c; \sigma_o) \end{cases}}$ $\quad$ <u>parametrization</u>

④ <u>Loss Function</u> : $\begin{cases} \underline{\text{classification}} \longrightarrow \underline{\text{cross-entropy}} \\[3mm] \underline{\text{sequence - prediction}} \longrightarrow \underline{\text{sum across time}} \end{cases}$

⑤ <u>Jordan Variant of RNN</u>



ELMAN

JORDAN

---

⑥ <u>Example</u> of [Purpose of Hidden Layer] !

consider a scaler <u>Time Series</u>   $\{ x^{(t)} : t = 1, 2, \dots \}$

output   $\{ y^{(t)} : t = 1, 2, \dots \}$

Here,   $y^{(t)} = x^{(t)} + x^{(t-1)} + x^{(t-2)}$   for   <u>$t \geq 1$</u>

| $t$ | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| Then. if $x^{(t)}$ | 1 | 2 | 3 | 4 | 5 |
| then $y^{(t)}$ | 1 | 3 | 6 | 9 | 12 |

<u>Goal</u> : build model to predict   $\boxed{\hat{y}^{(t)} = y^{(t)}}$   $\forall t$

Model : 1. $\hat{y}^{(t)} = FCNN(x^{(t)})$

Issue : No memory

2. $\hat{y}^{(t)} = $ General linear function on $\{x^{(1)}, ..., x^{(t)}\}$

$$= \sum_{s=1}^{t} a^{(s)} x^{(s)}$$

fit $a^{(s)}$ to data

Issue : cannot make inference (have memory)

3. (RNN)  →  3 unit of memory

$$h_1^{(t)} = x^{(t)}$$
$$h_2^{(t)} = x^{(t+1)}$$
$$h_3^{(t)} = x^{(t-2)}$$

→ $\hat{y}^{(t)} = h_1^{(t)} + h_2^{(t)} + h_3^{(t)}$

$\Rightarrow$
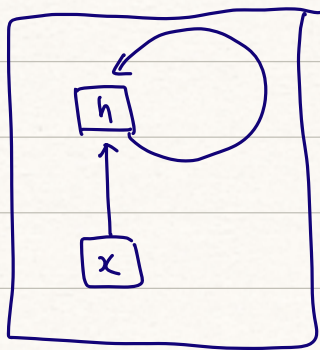$$h^{(t)} = W h^{(t-1)} + U \cdot x^{(t)} + b$$ → identical activation
$$y^{(t)} = V h^{(t)} + c$$

where
$$W = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} \quad U = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$
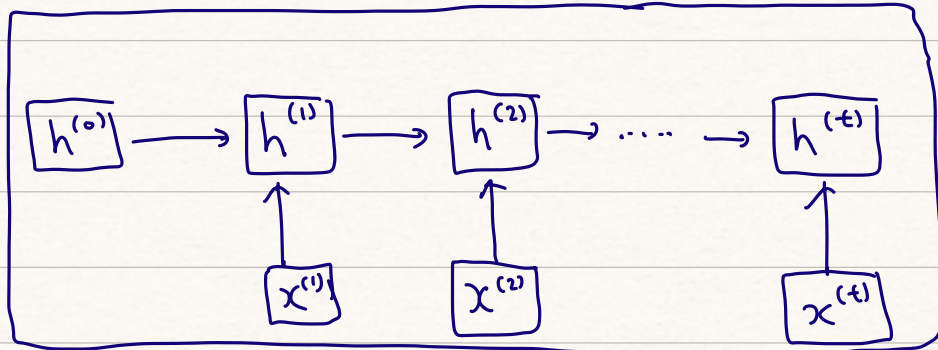
$$V = \mathbb{1}^T \qquad b = c = 0$$

⑦ How to train RNN ?

→ Keypoint: Unroll the computational graph

**Unroll**

$$h^{(0)} \rightarrow h^{(1)} \rightarrow h^{(2)} \rightarrow \cdots \rightarrow h^{(t)}$$

with inputs $x^{(1)}$, $x^{(2)}$, $x^{(t)}$

**Issue:** Gradient
- **Vanishing** $\rightsquigarrow$ for far-away nodes
- **Explosion**

hard to learn long-term dependency

idea: $a^k$
$$\begin{cases} +\infty & a > 1 \quad \text{as } k \to \infty \quad (\text{explosion}) \\ 0 & a < 1 \quad \text{as } k \to \infty \quad (\text{vanish}) \end{cases}$$

$\boxed{\text{Modification}} \rightarrow \underline{GRU} \ (LSTM)$

idea: use gates to control the
accumulation of knowledge

**Diagram**

==RNN Cell==



RNN Cell

$h_{t-1}$ — $\tanh$ — $h_t$

$h_t$

transform

$x_t$
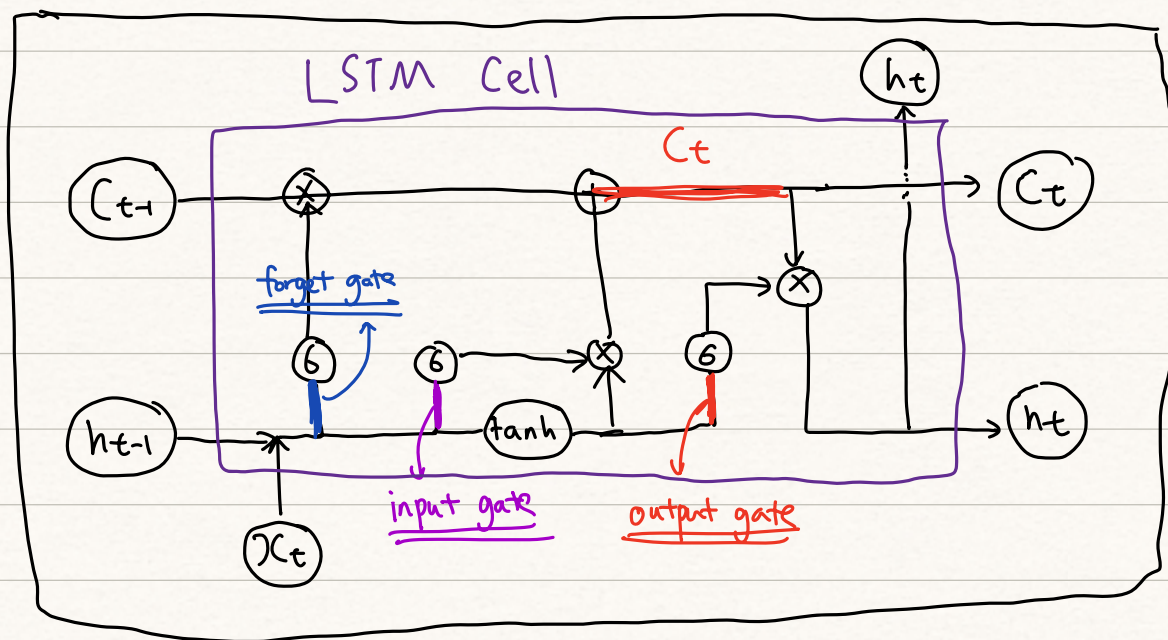
# LSTM Cell

$\bigstar$

## Self-understanding :

(adaptive)

→ Gate Unit is trying to learn from $\begin{cases} h_{t-1} \\ x_t \end{cases}$

in order to determine $\begin{cases} 1.\ \text{how much memory should be left} \\ 2.\ \text{how much memory should be updated} \\ 3.\ \text{how much hidden unit should be transformed} \end{cases}$

⑧ Deep RNN

→ Shallow RNN $\begin{cases} h^{(t)} = \sigma_r ( W h^{(t-1)} + U x^{(t)} + b) \\ \hat{y}^{(t)} = \sigma_o ( V h^{(t)} + c) \end{cases}$

→ Deep RNN $\begin{cases} h^{(t)} = \sigma_r ( W_1 h^{(t-1)} + U_1 x^{(t)} + b_1) \\ z^{(t)} = \sigma_r ( W_2 z^{(t-1)} + U_2 h^{(t)} + b_2) \\ \hat{y}^{(t)} = \sigma_o ( V z^{(t)} + c) \end{cases}$

⑦ Other Variants

→ 1. Bi-directional RNN ( Translation )

→ 2. Seq 2 Seq ( Encoder – Decoder Architecture )
⇓

MODEL { input : sequence
        output : sequence