Recurrent Neural Networks
( RNN )

Define a dynamical system

$$s^{(t)} = f(s^{(t-1)}; \Theta)$$

↑
output/
state

↑ set of
parameters

$$\frac{ds}{dt} = g(s, t)$$

Euler's method

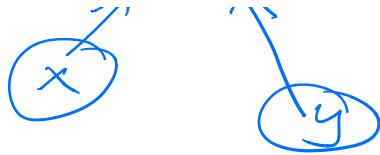$$s^{(t+1)} = s^{(t)} + h \; g(s^{(t)}, t)$$

↑
stepsize

recurrent equation since
s at time t+1 refers back
to the same definition at
time —t—

computational graphs

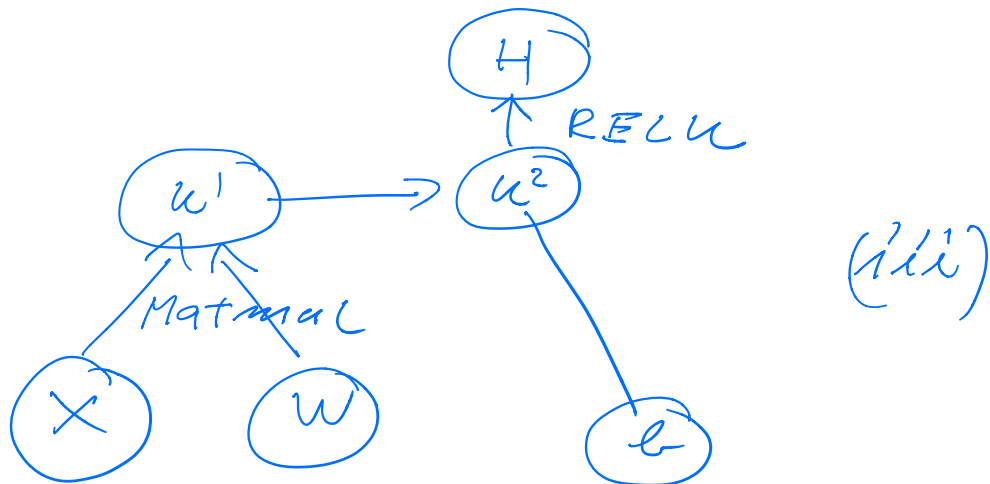$$z = x \cdot y$$



(i)

$$y = \sigma(X^T W + b)$$


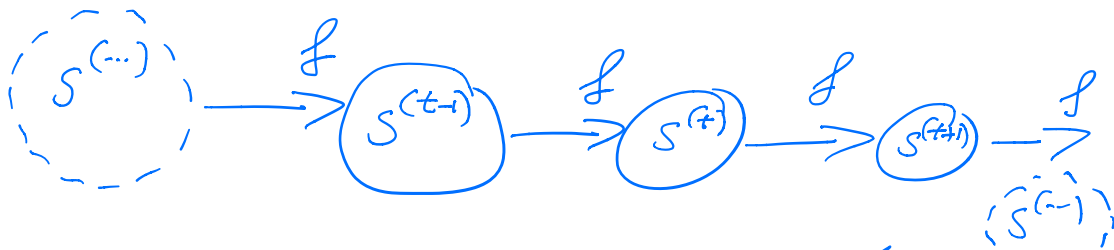
(ii)

$$H = \max\{0, XW + b\}$$



(iii)

Finite # of time steps $\tau$
from $t = \underline{1}$ to $t = \tau$

$\tau = 3$

$$S^{(3)} = f(S^{(2)}; \theta)$$

$$= f\left(f(s^{?}; \Theta); \Theta\right)$$

$$S^{(...)} \xrightarrow{f} S^{(t-1)} \xrightarrow{f} S^{(t)} \xrightarrow{f} S^{(t+1)} \xrightarrow{f} S^{(...)}$$
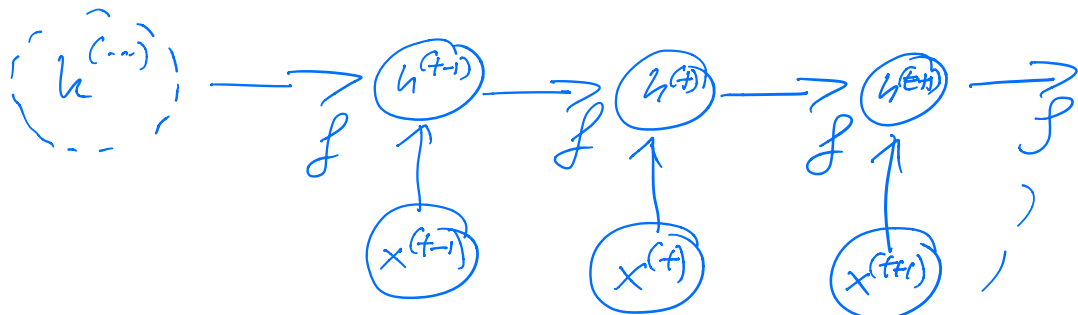
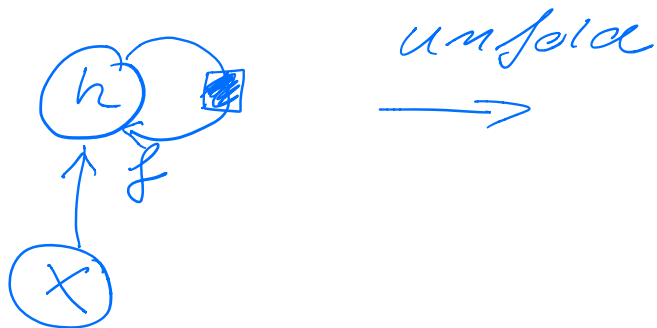Dynamical system driven
by an <mark>external signal</mark> $x^{(t)}$

$$S^{(t+1)} = f\left(S^{(t)}, x^{(t+1)}; \Theta\right)$$

Define hidden layer

$$h^{(t+1)} = f\left(h^{(t)}, x^{(t+1)}; \Theta\right)$$

unfold

$\longrightarrow$

$h$ $f$

$x$

$$h^{(...)} \xrightarrow{f} h^{(t-1)} \xrightarrow{f} h^{(t)} \xrightarrow{f} h^{(t+1)} \xrightarrow{f}$$

$$x^{(t-1)} \qquad x^{(t)} \qquad x^{(t+1)}$$

$$h^{(t)} = f\left(x^{(t)} U + S^{(t-1)} W + b\right)$$

$\left(\widehat{h^{(\text{...})}}\right)$

Example

$$m \frac{d^2 s}{dt^2} = -ks + A\cos(t)$$

$$\begin{cases} \dfrac{ds}{dt} = v(s, t) \\[2mm] \dfrac{dv}{dt} = -\dfrac{K}{m} s + \dfrac{A}{m} \cos(t) \end{cases}$$

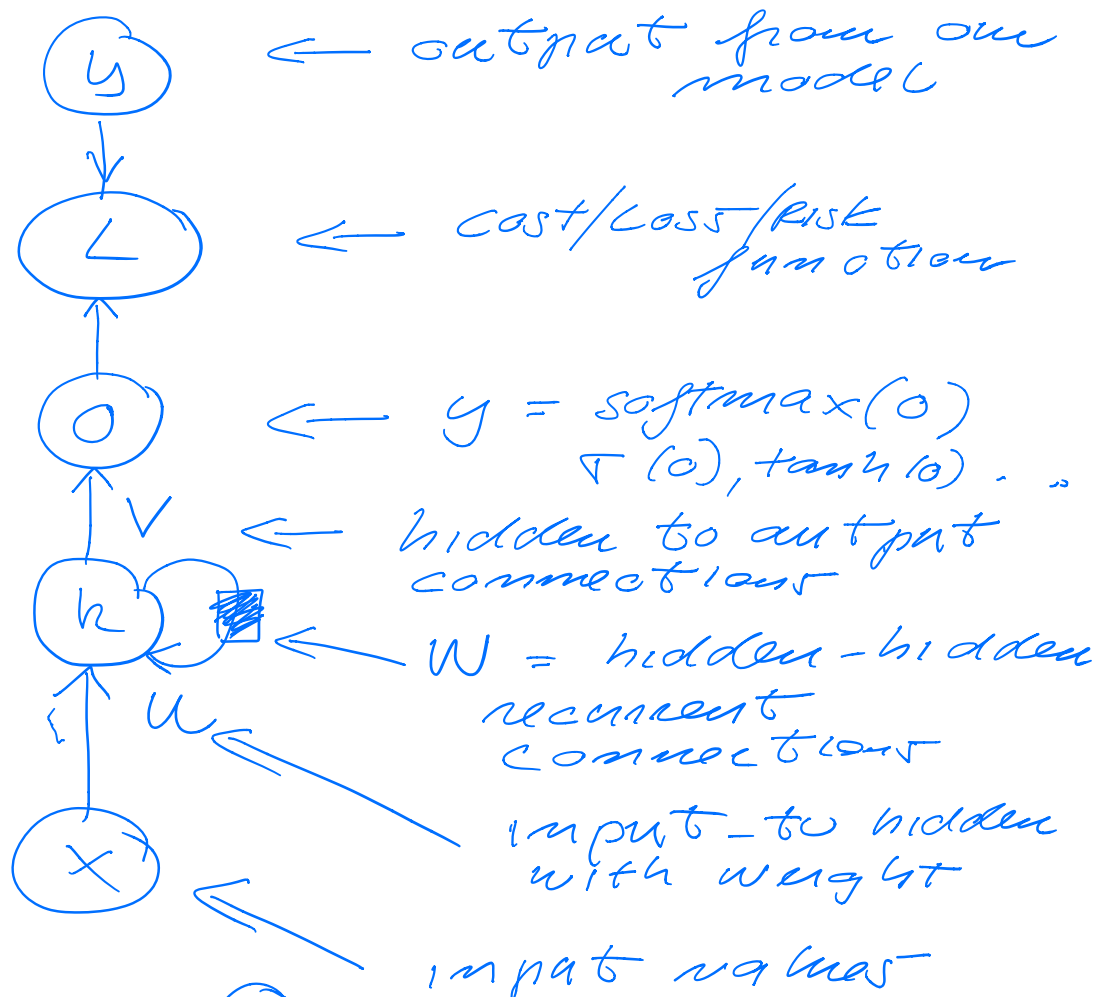Euler's method

$$S_{t+1} = S_t + h \cdot v_t$$

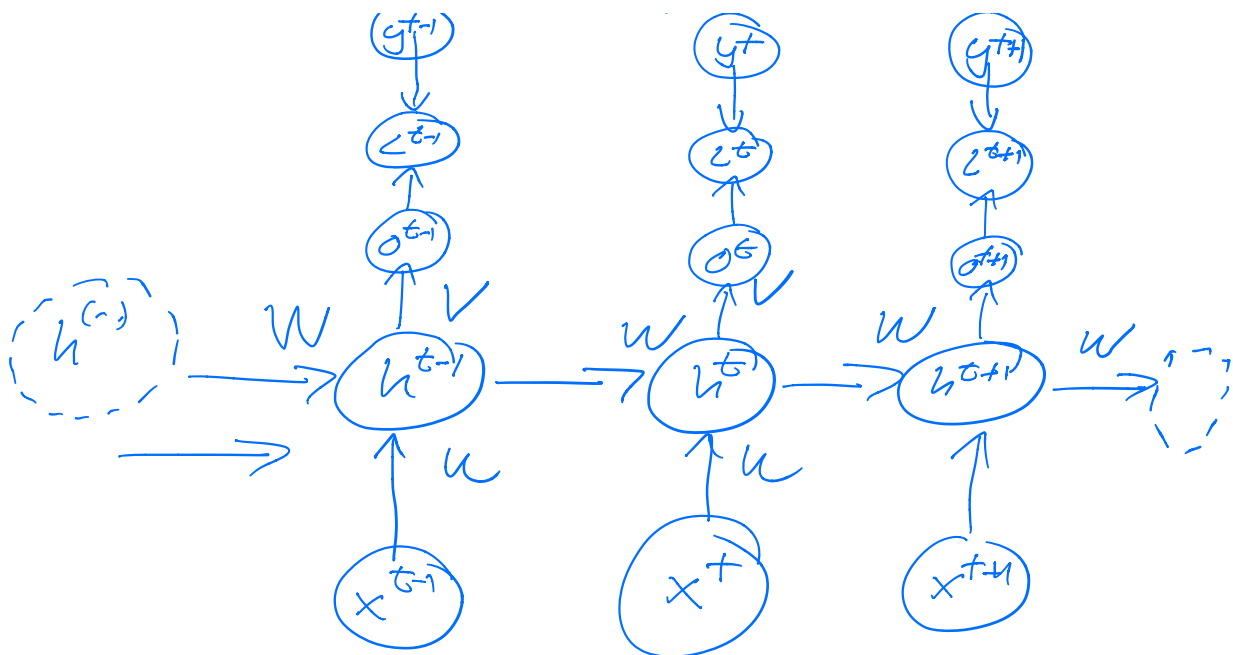$$v_{t+1} = v_t + h \left(\frac{A}{m} \cos(t) - \frac{K}{m} S_t\right)$$

# Recurrent NN (RNN)

1st model : RNN that produces an output $y$ at each time step and recurrent connections between hidden units

: RNN that produces an output at each time step and have recurrent connections from the output at one time step to the hidden unit at the next time step.

## 1st model

y  ← output from our model

L  ← cost/Loss/Risk function

O  ← $y = softmax(O)$
$\sigma(O), tanh(O) \ldots$

V  ← hidden to output connections

h  ← $W =$ hidden-hidden recurrent connections

U  ← input-to hidden with weight

X  ← input values

$$a^{(t)} = b + W h^{(t-1)} + u x^{(t)}$$

$$h^{(t)} = \tanh(a^{(t)}) \quad (\sigma(a^{(t)})$$
$$\text{ReLU}(a^{(t)}) ..)$$

$$y^{(t)} = \text{softmax}(o^{(t)}) / \sigma(o^{(t)})$$

$$o^{(t)} = c + V h^{(t)}$$

with $y^{(t)}$ + target$^{(t)}$ we
can calculate the Loss
function L.

Need: $\nabla_h L$, $\nabla_o L$, $\nabla_c L$

$\nabla_b L$, $\nabla_V L$, $\nabla_W L$, $\nabla_u L$
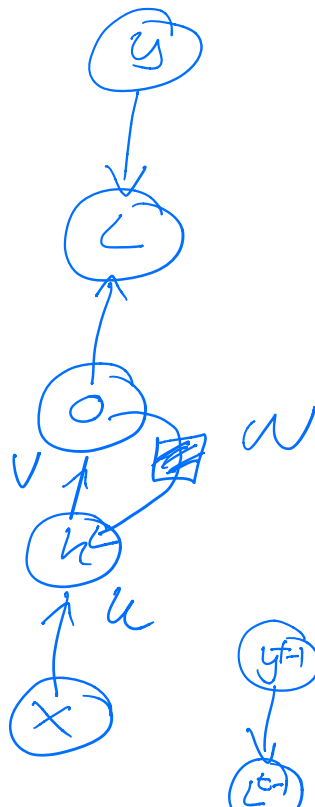
$X^{(t)}$ does not have any parameters.

- Training
  - Feed forward from left to right
  - Back prop in time from the right to the left ( BPTT )

Expensive to train and is difficult to parallelize.

2nd Model

$y$

$L$

$O$

$V$     $W$

$h$

$u$

$X$

open up

$y^{t-1}$

$L^{t-1}$

$y^{t}$

$L^{t}$

$y^{t+1}$

$L^{t+1}$

$O^{(\cdots)}$    $W$    $o^{t-1}$    $W$    $o^{t}$    $W$    $o^{t+1}$    $W$

$V$     $V$     $V$

$h^{t-1}$    $h^{t}$    $h^{t+1}$

$u$     $u$     $u$

$x^{t-1}$    $x^{t}$    $x^{t+1}$