

Lecture October 28

RNNs Basic definitions

Example ;

ODE :

$$m \frac{d^2 x}{dt^2} + x(t) = F(t)$$

$$v = \frac{dx}{dt}$$

$$v(t_0) = v_0$$

$$x(t_0) = x_0$$

$$a = \frac{dv}{dt}$$

$$m \frac{dv}{dt} + x(t) = F(t)$$

$$\frac{dx}{dt} = v$$

Euler's method (velocity)

$$v_{i+1} \approx v_i + \Delta t v'_i \Big|_{t=t_i}$$

$$\frac{dv}{dt} = -\frac{x}{m} + \frac{F}{m}$$

$$= -\tilde{x} + \tilde{F}$$

$$v_{i+1} \approx v_i + \Delta t (-\tilde{x}_i + \tilde{F}_i)$$

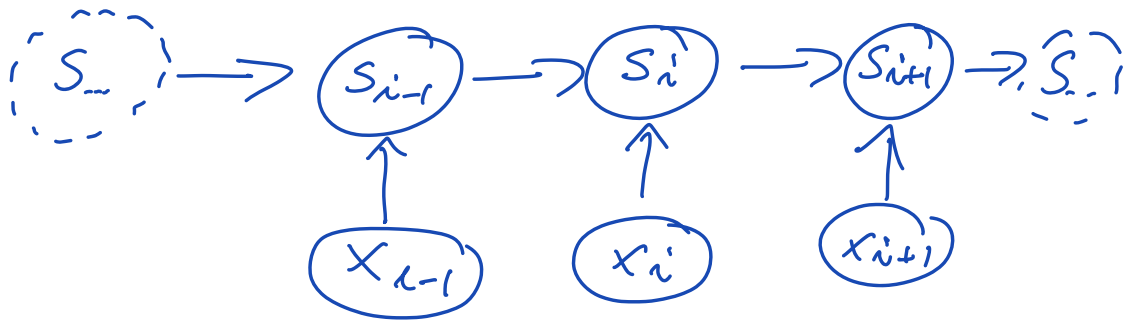
$$v_{i+1} = f(\tilde{x}_i, \tilde{F}_i, v_i)$$

$$= f(x_i, F_i, v_i)$$

$$v_{i+1} = f(f(x_{i-1}, F_{i-1}, v_{i-1}), F_i, x_i)$$

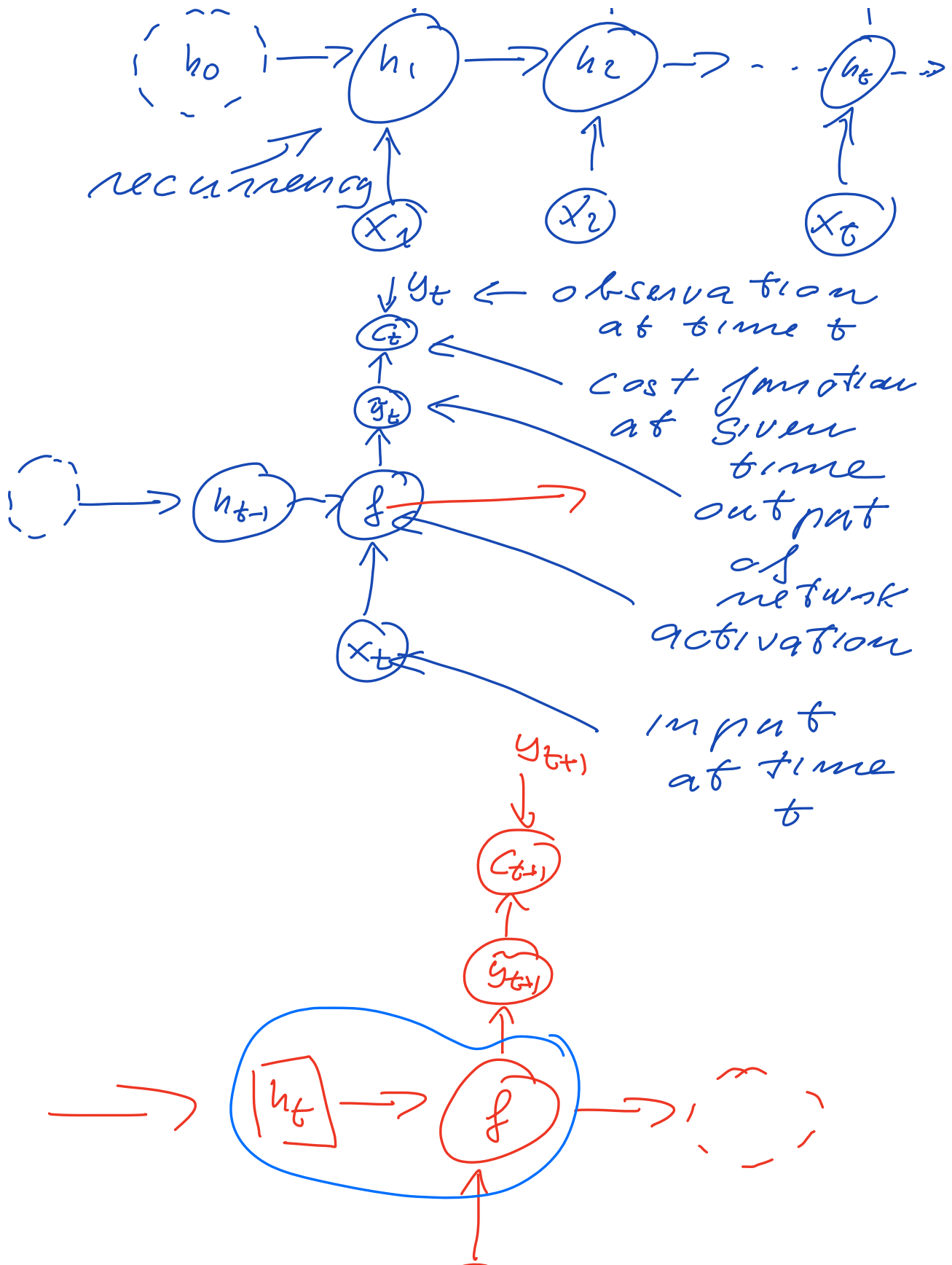
↓

$$s_{i+1} = h(s_i, x_i; \epsilon)$$



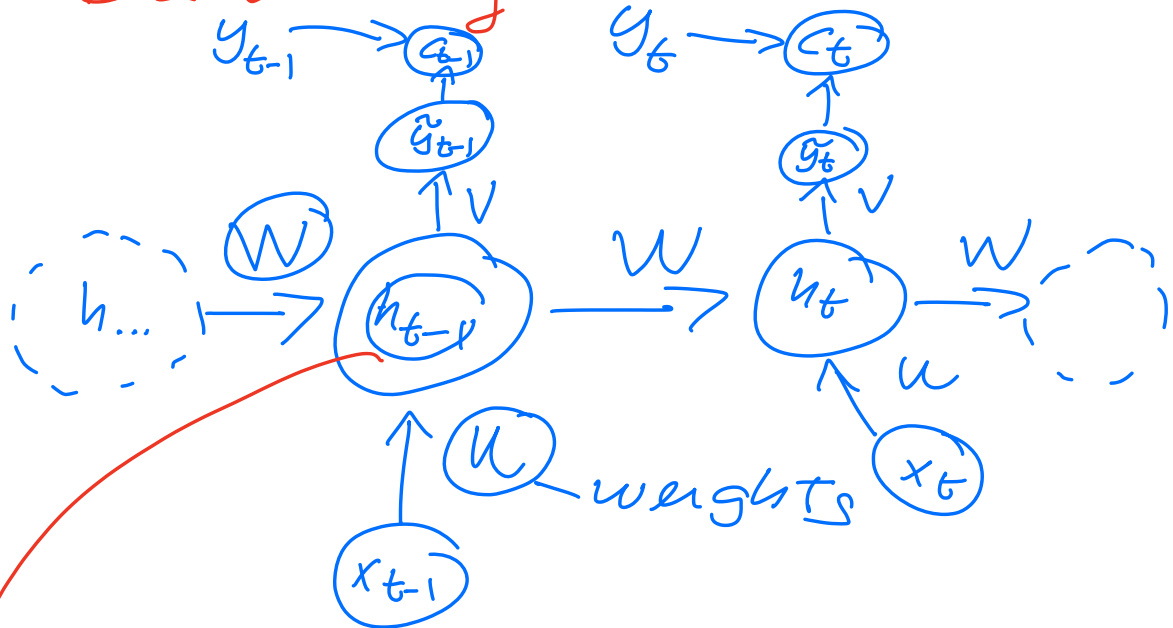
expand to include an
output y_t





x_{t+1}

Breaking it down



parameters:

W (weights)

$$h_{t-1} = f(z_{t-1})$$

$$z_{t-1} = b + Wh_{t-2} + Ux_{t-1}$$

$$\tilde{y}_{t-1} = g(h_{t-1}V + c)$$

the cost function at
time $t-1$

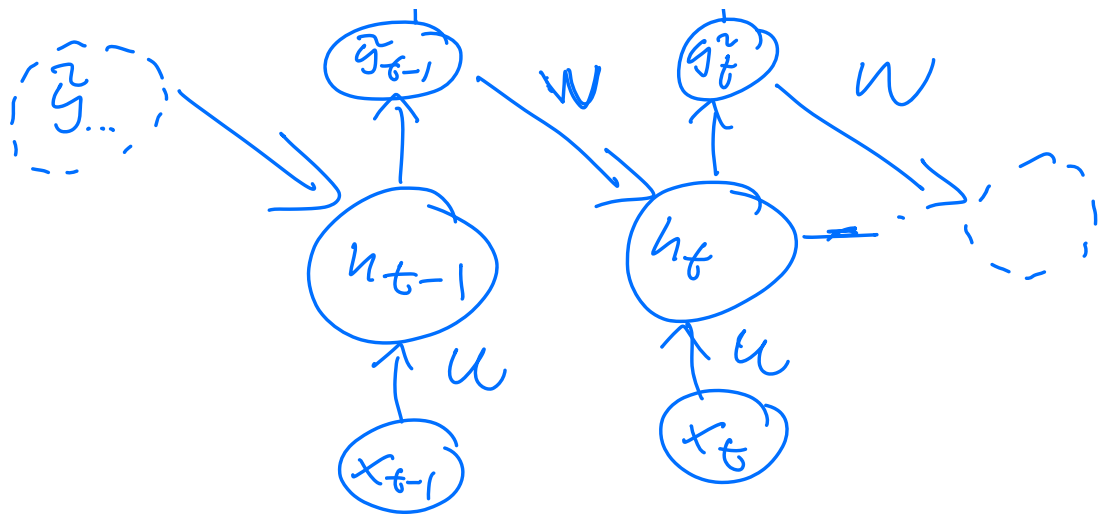
$$C_{t-1}(y_{t-1}, h_{t-1}, \tilde{y}_{t-1}, x_{t-1})$$

- 1) Feed Forward part with initialization of weights and biases at all t
- 2) Back propagation through time.
(BPTT)

Memory and CPU intensive and keeps long term memory.

- instead of connections between hidden layers, we can feed in output \tilde{y}_{t-1} into h_t





Gradient training problems

- vanishing gradients
less of a problem
- exploding gradients

Simple example

$$h_t = W h_{t-1}$$

\nearrow
 output from hidden nodes (no x_t no bias)

W is the same at all times - t -

$$h_1 = W h_0 \dots W h_{t-1}$$

$$h_t = W^t h_0$$

$$W w_i = \lambda_i w_i$$

$$h_0 = \sum_i \alpha_i w_i$$

$$W^t h_0 = \sum_i \alpha_i \lambda_i^t w_i$$

$$\lambda_0 > \lambda_1 > \lambda_2 \dots$$

$\lambda_0 > 1$, then $W^t h_0$ will diverge,

$\Rightarrow \nabla_{h_t}$ they will diverge

$\lambda_0 < 1$ the $W^t h_0$ will decrease \Rightarrow

$$\nabla_{h_t} \Rightarrow 0$$

To avoid exploding gradients:

gradient clipping

gradient g .

if $\|g\|_2 \geq \epsilon$

$$g \leftarrow \frac{\epsilon}{\|g\|_2} g$$

end if.