# Recurrent Neural Networks Notes   #1
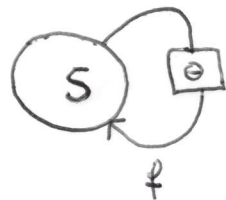
**I:** Modelling sequences requires:

(α): deal with variable-length sequences

(β): maintain sequence order

(γ): keep track of long-term dependencies

(δ): share parameters across the sequence

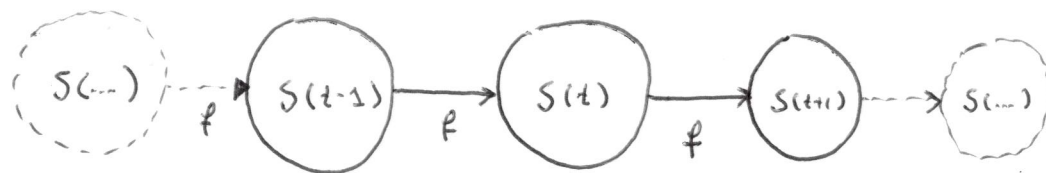**II:** Classical form of an autonomous dynamical system:

$$\underline{S}(t) = f(\underline{S}(t-1); \underline{\theta})$$
$$t \geq 1$$

where $\underline{S}(t)$ is the state-vector of the system and $\underline{\theta}$ is a vector of internal constant parameters

**III:** Computational Graph of Classical Dynamical System:
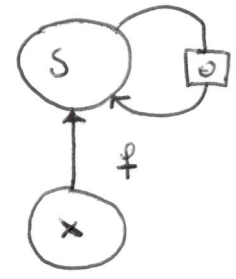


**IV:** Unfolded Computational Graph of Classical Dynamical System:

## V: Autonomous Dynamical System Driven by External Signal $\underline{x}(t)$:

$$\underline{S}(t) = f\left(\underline{S}(t-1), \underline{x}(t); \underline{\theta}\right)$$
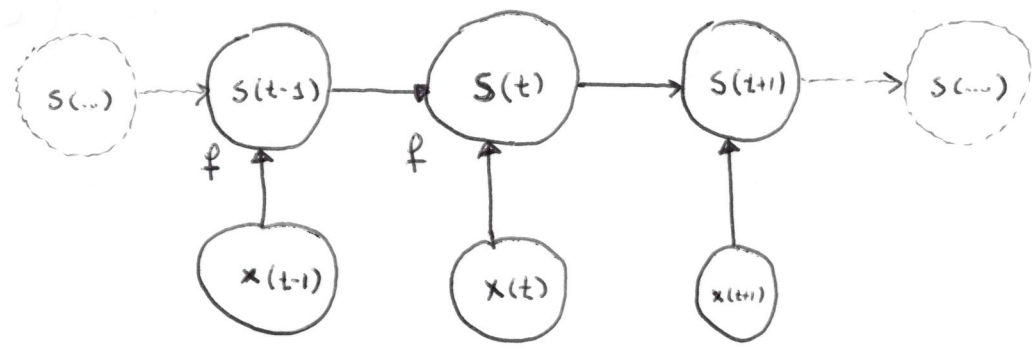$$t \geq 1$$

## VI: Computational Graph of an Autonomous Dynamical System Driven by External Signal $\underline{x}(t)$:
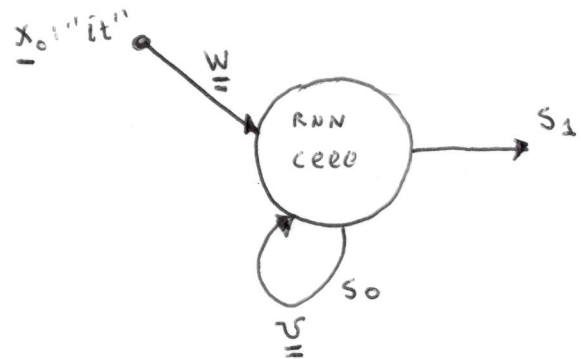


## VIII: Many RNNs can be described by a similar equation, such as:

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

## VII: Unfolded Computational Graph:

**IX:** RNNs are capable of remembering their previous state:

$x_0$: "it"



W

RNN cell

$S_1$

$S_0$

U

$x_0$ : vector representing first word

$S_0$ : cell state at $t=0$ (some initialization is required)

$S_1$ : cell state at $t=1$

$W$, $U$ : are weight matrices

⊛ Update Equation for next cell state:

$$S_1 = \tanh(W x_0 + U S_0)$$

$x_1$ : "was"



W

RNN cell

$S_2$

$S_1$

U

$x_1$ : vector representing second word

$S_1$ : cell state at $t=1$

$S_2$ : cell state at $t=2$

$$S_2 = \tanh(W x_1 + U S_1)$$

$W$, $U$ : weight matrices

**X :** Unfolding the RNN across time;

(★) In practise, one may have many
hidden units and many layers of hidden units.



(α): Notice that the same parameter matrices $\underline{W}$, $\underline{U}$ are used.

(β): $S_n$ can contain information from all past timesteps.

**XI :** Training RNN with Back Propagation Through Time ( BBTT )

(α): Take the derivative of the loss (gradient) with respect to each parameter.

(β): shift parameters to the opposite direction in order to minimize the loss.

Gradient Descend !!!

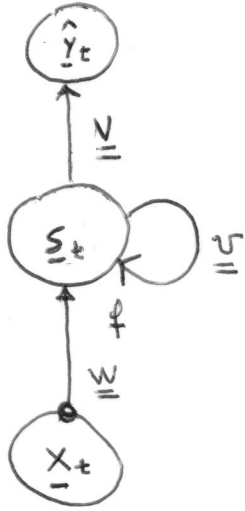$\underline{\underline{XII}}$: Let $\Upsilon = \{ \underline{x}_1, \underline{x}_2, ..., \underline{x}_z \}$ a sequence of $\ell$-dimensional input vectors ($\underline{x}_t \in \mathbb{R}^\ell$, $\forall t \in [z]$)  #5

Let $y = \{ \underline{y}_1, \underline{y}_2, ..., \underline{y}_z \}$ the corresponding sequence of target output vectors ($y_t \in \mathbb{R}^m$, $\forall t \in [z]$)

$\circledast$ The estimated output of the RNN cell at each time step may be computed as:



(A): Update Hidden State:

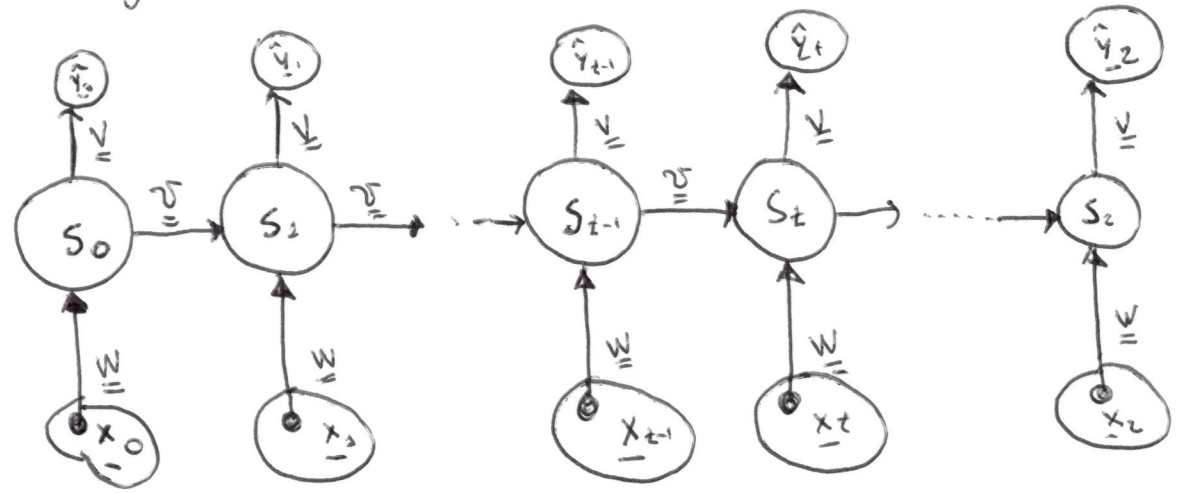$$\underline{S}_t = f( \underline{\underline{W}} \circ \underline{x}_t + \underline{\underline{U}} \circ \underline{S}_{t-1}) \quad (A)$$

where $f(u) = \tanh(u) = \dfrac{e^u - e^{-u}}{e^u + e^{-u}}$

(B): Output Vector:

$$\hat{\underline{y}}_t = \underline{\underline{V}} \circ \underline{S}_t \quad (B)$$

$\circledast$ By unfolding the RNN cell for $z$ time steps we get:

Unfolded Description of the RNN Cell!!!

✸ A very helpfull result relates to expressing the derivative of a given activation function
as a function of itself. Let's consider the hyperbolic tangent function for example.

$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad [\text{SIGMOID FORM}]$$



● Express the derivative $f'(x)$ as a function of $f(x)$?

Let $p(x) = e^x - e^{-x}$ and $q(x) = e^x + e^{-x}$, so that $f(x) = \frac{p(x)}{q(x)}$.

We may observe that:  (i): $p'(x) = e^x + e^{-x} = q(x)$

(ii): $q'(x) = e^x - e^{-x} = p(x)$

Thus, we have that: $f'(x) = \frac{p'(x)q(x) - p(x)q'(x)}{q^2(x)} \Rightarrow$

$$f'(x) = \frac{q^2(x) - p^2(x)}{q^2(x)} \Rightarrow f'(x) = 1 - \frac{p^2(x)}{q^2(x)} \Rightarrow$$

$$f'(x) = 1 - \left(\frac{p(x)}{q(x)}\right)^2 \Rightarrow \boxed{f'(x) = 1 - f^2(x)}$$

(★) Since we are making a prediction at each time step, we have to associate a loss at each timestep too:



Unfolded RNN for each time-step $0 \le t \le 2$.

{ Updating Equations for each time-step $1 \le t \le 2$.

$$S(t) = f(\underline{\underline{W}} \circ x(t) + \underline{\underline{U}} \circ \underline{S}(t-1) + \underline{b})$$

$$\hat{y}(t) = \underline{\underline{V}} \circ \underline{S}(t) + \underline{c}$$

(★) Bias-terms may be associated with each updating equation.

(★) Form the expression for the total loss of the RNN throughout the complete sequence $J$:

$$J(\underline{\Theta}) = \sum_{t=0}^{2} J_t(\underline{\Theta})$$

Total Loss

(★) Compute the total gradient by summing the partial derivative w.r.t each parameter of the RNN foreach time-step:
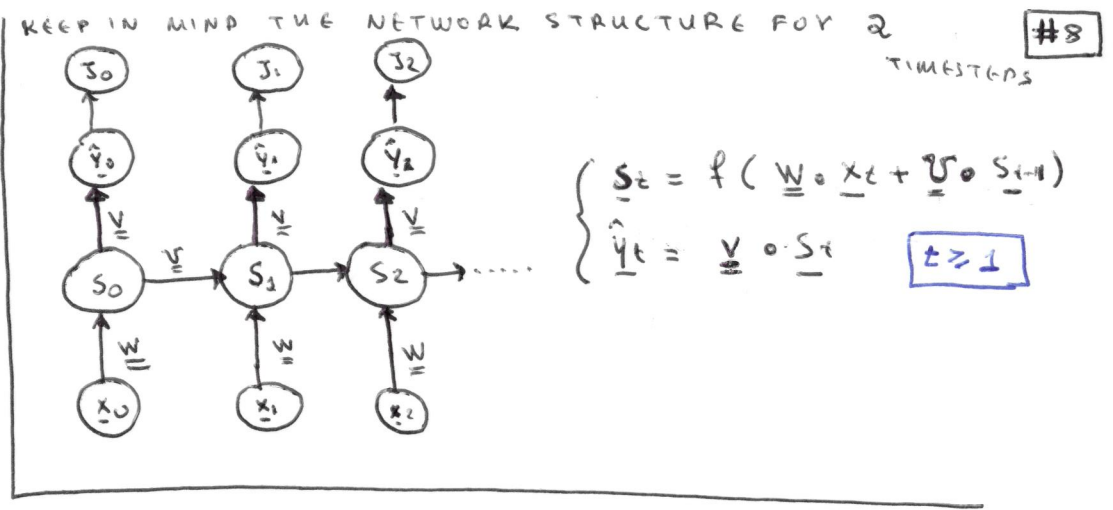
$$\frac{\partial J}{\partial P} = \sum_t \frac{\partial J_t}{\partial P} \quad (1)$$



$$\begin{cases} \underline{S_t} = f(\underline{\underline{W}} \cdot \underline{x_t} + \underline{\underline{U}} \cdot \underline{S_{t-1}}) \\ \underline{\hat{y}_t} = \underline{\underline{V}} \cdot \underline{S_t} \end{cases} \quad \boxed{t \geq 1}$$

(★) Let's initially focus on the weight-matrix $\underline{\underline{W}}$:

$$\frac{\partial J}{\partial \underline{\underline{W}}} = \sum_t \frac{\partial J_t}{\partial \underline{\underline{W}}} \quad (2)$$

(A) Let's focus on a single time step where $t=2$:

$$\frac{\partial J_2}{\partial \underline{\underline{W}}} = \frac{\partial J_2}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial \underline{S_2}} \cdot \boxed{\frac{\partial \underline{S_2}}{\partial \underline{\underline{W}}}} \quad (3)$$

CHAIN RULE OF DERIVATIVES.

(A) We need to know how exactly the repp-state at time $t=2$ depends on the weight-matrix $\underline{\underline{W}}$:

(F) In other words, we need to disambiguate the term $\frac{\partial S_2}{\partial W}$ for which we need to express its total derivative with respect to $\underline{\underline{W}}$ as:

$$\frac{dS_2}{d\underline{\underline{W}}} = ?$$

But, $S_2$ may be written as

$$\underline{S_2} = f(\underline{\underline{W}} \cdot \underline{x_2} + \underline{\underline{U}} \cdot \underline{S_1}) \quad (u)$$

where $\underline{S_1}$ also depends on $\underline{\underline{W}}$. Therefore, the term $\frac{\partial S_2}{\partial \underline{\underline{W}}}$ cannot be treated as a constant.

⊛ Eq. (3) suggests that:

$$S_2 = S_2 \left( S_1(\underline{w}), S_0(\underline{w}) ; \underline{w} \right) \qquad (5)$$

▶ Thus, the total derivative of $S_2$ w.r.t. $\underline{w}$ may be written as: (chain rule for composite functions)

$$\frac{dS_2}{d\underline{w}} = \frac{\partial S_2}{\partial \underline{w}} + \frac{\partial S_2}{\partial S_1} \cdot \frac{\partial S_1}{\partial \underline{w}} + \frac{\partial S_2}{\partial S_0} \cdot \frac{\partial S_0}{\partial \underline{w}} \qquad (6)$$

▶ Eq. (6), suggests the following compact form:

$$\frac{dS_2}{d\underline{w}} = \sum_{K=0}^{2} \frac{\partial S_2}{\partial S_K} \cdot \frac{\partial S_K}{\partial \underline{w}} \qquad (7)$$

▶ Apparently, the summation term for $k=2$ yields $\boxed{\dfrac{\partial S_2}{\partial S_2} \cdot \dfrac{\partial S_2}{\partial \underline{w}} = \dfrac{\partial S_2}{\partial \underline{w}}}$

⊛ Finally, Eq. (3) becomes:

$$\frac{\partial S_2}{\partial \underline{w}} = \sum_{K=0}^{2} \frac{\partial J_2}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial S_2} \cdot \frac{\partial S_2}{\partial S_K} \cdot \frac{\partial S_K}{\partial \underline{w}} \qquad (8)$$

↓

Contributions of $\underline{w}$ in previous time-steps to the error at time-step $t=2$.

---

⊛ In this setting, the general form of Eq. (8) at any given time-step $t$ may be given as:

$$\frac{\partial J_t}{\partial \underline{w}} = \sum_{K=0}^{t} \frac{\partial J_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial S_t} \cdot \frac{\partial S_t}{\partial S_K} \cdot \frac{\partial S_K}{\partial \underline{w}} \qquad (9)$$

↓

Contribution of $\underline{w}$ in previous timesteps to the error at timestep $t$.

⊛ What about $\frac{\partial J}{\partial \underline{v}}$ and $\frac{\partial J}{\partial V}$ ?

▶ $$\frac{\partial J}{\partial \underline{v}} = \sum_{t} \frac{\partial S_t}{\partial \underline{v}} \qquad (10)$$

▼ $$\frac{\partial J}{\partial \underline{v}} = \sum_{t} \frac{\partial J_t}{\partial \underline{v}} \qquad (11)$$

⊛ We have to aknowledge that

(i): $S_t$ does depend on $\underline{v}$

(ii): $S_t$ does not depend on $\underline{V}$

⊛ Thus, we may write that :

$$\frac{\partial J_t}{\partial \underline{\underline{V}}} = \frac{\partial J_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial \underline{\underline{V}}} \quad (12)$$

⊛ For the case of the weight-matrix $\underline{\underline{U}}$ we may write that:

$$\frac{\partial J_t}{\partial \underline{\underline{U}}} = \frac{\partial J_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial s_t} \cdot \left(\boxed{\frac{\partial s_t}{\partial \underline{\underline{U}}}}\right) \quad (13)$$

⊛ We need to take into consideration the contribution of $\underline{\underline{U}}$ in the previous timesteps to the error at timestamp t.

⊛ To do so, we have to formulate the explicit dependence of each $s_t$ on $\underline{\underline{U}}$ alongside with the implicit dependence of each $s_t$ on the previous state vectors $s_{t-1}, s_{t-2}, \dots, s_1, s_0$.

⊛ Therefore, we may write that :

$$s_t = s_t \left( s_{t-1}(\underline{\underline{U}}), s_{t-2}(\underline{\underline{U}}) \dots, s_0(\underline{\underline{U}}); \underline{\underline{U}} \right) \quad (14)$$

⊛ We need the total derivative counterpart of the quantity, $\frac{d s_t}{d \underline{\underline{U}}}$

⊛ Eq.(14) suggests that we may write:

$$\frac{d s_t}{d \underline{\underline{U}}} = \sum_{k=0}^{t} \frac{\partial s_t}{\partial s_k} \cdot \frac{\partial s_k}{\partial \underline{\underline{U}}} \quad (15)$$

⊛ Finally, by combining Eqs.(13) and (15) we get :

$$\frac{\partial J_t}{\partial \underline{\underline{U}}} = \sum_{k=0}^{t} \frac{\partial J_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial s_t} \cdot \frac{\partial s_t}{\partial s_k} \cdot \frac{\partial s_k}{\partial \underline{\underline{U}}} \quad (16)$$

✪ Next Question : "Why RNNs are so hard to train ?"

In Eqs.(9) and (16) share a common factor $\frac{\partial s_t}{\partial s_k}$ which may be written as:

$$\frac{\partial s_t}{\partial s_k} = \prod_{m=k+1}^{t} \frac{\partial s_m}{\partial s_{m-1}} \quad (17)$$

⊛ The problem of vanishing gradient:

According to Eq. (16), the gradient of the loss function at time (t) requires the computation of the term

$$\frac{\partial h_t}{\partial h_k} \quad [17]$$

For example for $t=2$ and $k=0$, we would have to compute $\frac{\partial h_2}{\partial h_0}$ which according to the interdependence of the various hidden states of the RNN can be written as:

$$\frac{\partial h_2}{\partial h_0} = \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial h_0} \quad [18]$$

But, what if we are focusing on a timestep very far in the future?

$$\frac{\partial h_t}{\partial h_0} = \frac{\partial h_t}{\partial h_{t-1}} \cdot \frac{\partial h_{t-1}}{\partial h_{t-2}} \cdot \dots \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial h_0} \quad [19]$$

⊛ The general form of the term $\frac{\partial h_t}{\partial h_k}$ may be expressed as:

$$\frac{\partial h_t}{\partial h_k} = \prod_{m=k+1}^{m=t} \frac{\partial h_m}{\partial h_{m-1}} \quad [20]$$

which in turn reveals the necessity of computing the following term:

$$\frac{\partial h_t}{\partial h_{t-1}} \quad [21]$$

For the case of a trivial vanilla ANN with a single neuron, we may assume that:

(1): $x_t \in \mathbb{R}^{\ell \times 1}$, $1 \leq t \leq 2$

(2): $y_t \in \mathbb{R}$, $1 \leq t \leq 2$

(3): $h_t \in \mathbb{R}$, $1 \leq t \leq 2$ ($u_t \equiv h_t$)

(4): $W_{xh} \in \mathbb{R}^{\ell \times 1}$, $1 \leq t \leq 2$ ($W_{xh} \equiv W_{xh}$)

(5): $W_{hh} \in \mathbb{R}$, $1 \leq t \leq 2$ ($W_{hh} \equiv W_{hh}$)

(6): $W_{hy} \in \mathbb{R}$, $1 \leq t \leq 2$ ($W_{hy} \equiv W_{hy}$)

⊛ According to the previous definitions, we may re-write the updating equations for the forward pass of the information within the network as:

$$h_t = f(\underline{W}_{xu}^T \cdot \underline{x}_t + W_{nn} \cdot h_{t-1}) \quad [22]$$

$$y_t = W_{ny} \cdot h_t \quad [23]$$

⊛ At this point, we must mention the content of Lecture Page #6 concerning the hyperbolic tangent and its derivative.

⊛ In the light of the previous declarations we may write that:

(i): Let $u \in \mathbb{R}$ be the expression that forms the input argument for the transfer function $f(\cdot)$ as:

$$u = \underbrace{\underline{W}_{xu}^T \cdot \underline{x}_t}_{[1 \times e] \cdot [e \times 1]} + \underbrace{W_{nn} \cdot h_{t-1}}_{[1 \times 1] [1 \times 1]}$$
$$[1 \times 1] + [1 \times 1] = [1]$$

(ii): $\dfrac{\partial h_t}{\partial h_{t-1}} = \dfrac{\partial f}{\partial u} \cdot \dfrac{\partial u}{\partial h_{t-1}} \quad [24]$

(iii) $\dfrac{\partial u}{\partial h_{t-1}} = W_{nn} \quad [25]$

⊛ Combining (eqs. (24) and (25)) yields:

$$\boxed{\dfrac{\partial h_t}{\partial h_{t-1}} = f'(u) \cdot W_{nn} \quad [26]}$$ , which gives:

$$\dfrac{\partial h_t}{\partial h_{t-1}} = (1 - f^2(\underline{W}_{xu}^T \underline{x}_t + W_{nn} \cdot h_{t-1})) W_{nn} \quad [27]$$

⊛ According to Eq. (26), we may identify that:

(i): Since, $f(u) = \tanh(u) \Rightarrow f'(u) \in [0,1]$

(ii): $W_{nn}$ are sampled from the standard normal distribution, (for the 1-d case, $\mu = 0$, $\sigma^2 = 1$) so that $W_{nn} < 1$
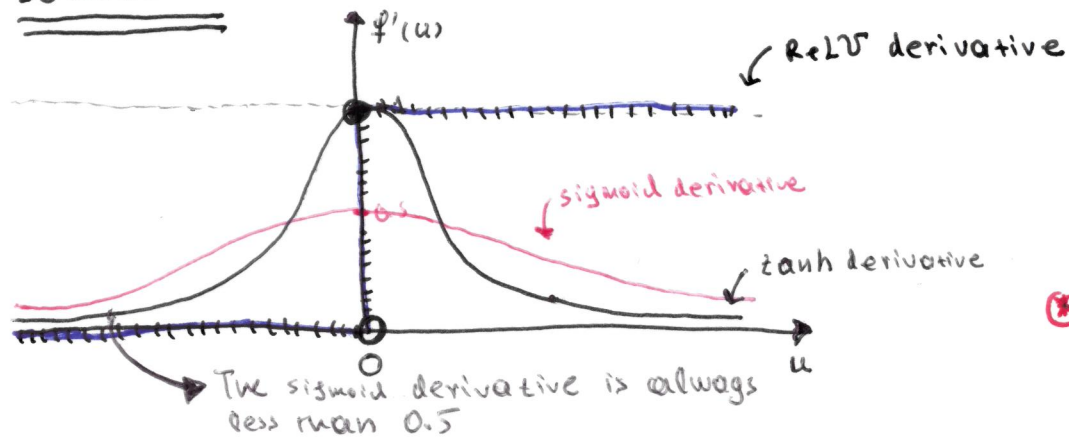
⊛ In other words, we are multiplying a lot of small numbers together.

(I): Errors due to further back timesteps have increasingly smaller gradients. (We are actually multiplying long sequences of decimals which are smaller than 1)

(II): Weight-parameters become biased towards capturing shorter-term dependencies.

**⊛ Addressing the Problem of Vanishing Gradients:**

**Solution I:** Choice of Activation Functions:



ReLU derivative

sigmoid derivative

tanh derivative

The sigmoid derivative is always less than 0.5

**⊛ ReLU Activation Function:**

$$f(x) = \max(0, x) = \begin{cases} x, & x > 0; \\ 0, & x \leq 0. \end{cases} \quad [28]$$

$$f'(x) = \begin{cases} 1, & x > 0; \\ 0, & x \leq 0. \quad [29] \end{cases}$$

**⊛ Thus, choosing ReLU as the activation function prevents $f'$ from shrinking !!!**

**Solution II:** Initialize Weight Matrices:

(i): Weight-matrices initialized to the Identity Matrix $(I)$, at least prevents vanishing the derivative at the first steps of the calculation

(ii): Biases-vectors could be initialized to zero.

**Solution III:**

Rather than each node being just a simple RNN cell, we could make each cell a more complex unit with gates controlling what information is passed through.

$$RNNs \quad vs \quad \{ LSTMs, GRUs, etc... \}$$

Long Short Term Memory cells are able to keep track of information throughout many timesteps.