

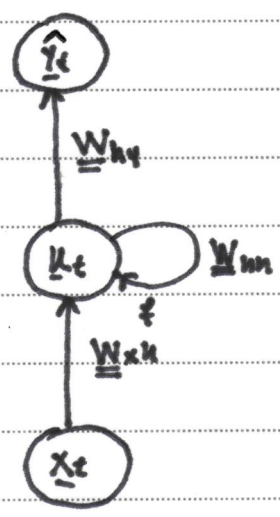
RNN Training with Back Propagation Through Time [BBTT]

- Gradient Descent:
- (a): Take the derivative of the loss (gradient) with respect to each parameter
 - (b): Shift parameters to the opposite direction in order to minimize loss.

Let $X = \{x_1, x_2, \dots, x_z\}$ a sequence of l -dimensional vectors ^{input}
 $x_t \in \mathbb{R}^l$ for $1 \leq t \leq z$.

Let $Y = \{y_1, y_2, \dots, y_z\}$ a sequence of m -dimensional output vectors
 $y_t \in \mathbb{R}^m$ for $1 \leq t \leq z$.

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



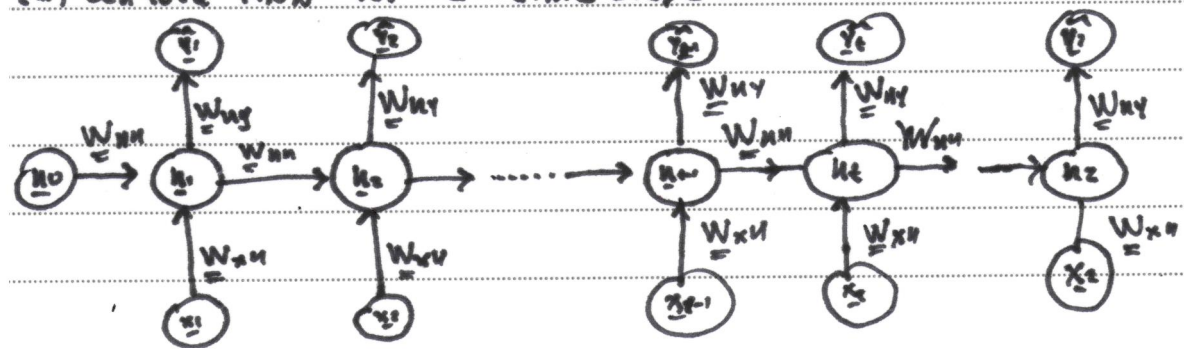
(A): Update Hidden State:

$$h_t = f \left(\sum_{1 \leq t \leq z} W_{xh} \cdot x_t + W_{hh} \cdot h_{t-1} \right) \quad (1)$$

(B): Update Output Vector:

$$\hat{y}_t = W_{hy} \cdot h_t \quad 1 \leq t \leq z \quad (2)$$

(*) Unfold RNN for z time steps:

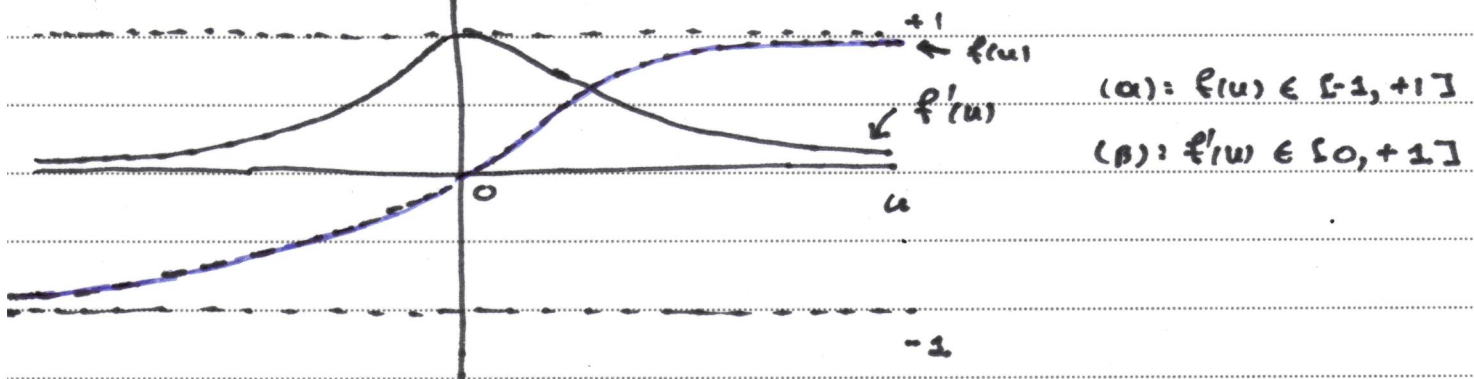


h_0 : Some Random Initialization

► The activation function utilized to update the hidden state is given by:

$$f(u) = \tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}} \quad (3)$$

SIGMOID
FORM



► A very helpful technique involves expressing the derivative of a given activation function as a function of itself.

► Try expressing $f'(x)$ as a function of $f(x)$:

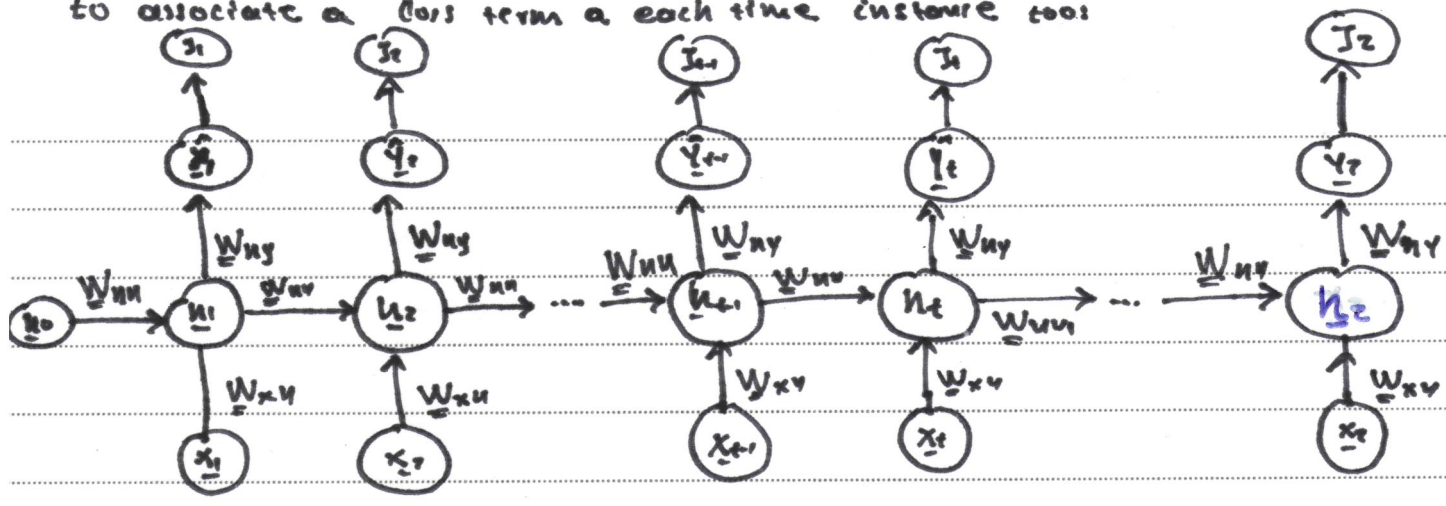
(i): let $p(x) = e^x - e^{-x}$ and $q(x) = e^x + e^{-x}$ so that (ii): $f(x) = \frac{p(x)}{q(x)}$

► We may write that: (iv): $\begin{cases} p'(x) = e^x + e^{-x} \\ p'(x) = q(x) \end{cases}$ and $\begin{cases} q'(x) = e^x - e^{-x} \\ q'(x) = p(x) \end{cases}$

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

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

*) Since we are making a prediction at each time step we have to associate a loss term at each time instance too:



*) Updating Equations: ($1 \leq t \leq z$):

$$\begin{cases} h_t = f(W_{xh} \cdot x_t + W_{hh} \cdot h_{t-1}) & (5) \\ \hat{y}_t = W_{hy} \cdot h_t & (6) \end{cases}$$

*) Form the expression of the total loss for the RNN throughout the complete sequence X :

$$J(\theta) = \sum_t J_t(\theta) \quad (7)$$

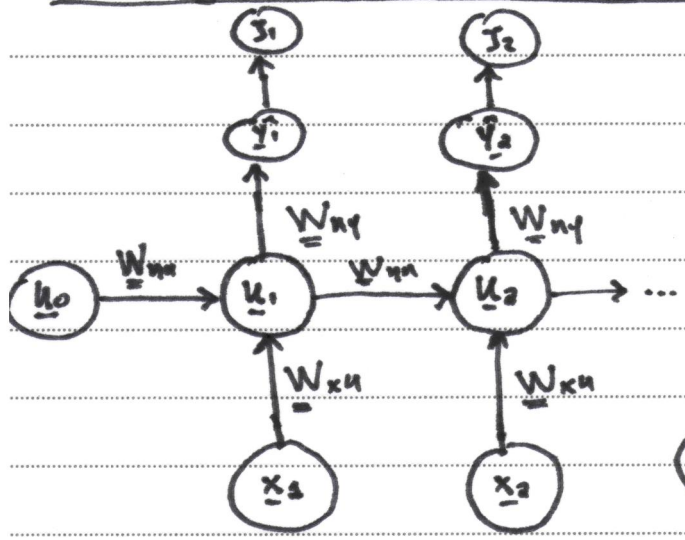
*) Compute the total gradient by summing the partial derivative with respect to each parameter θ of the RNN:

$$\frac{\partial J}{\partial \theta} = \sum_t \frac{\partial J_t(\theta)}{\partial \theta} \quad (8)$$

Let's initially focus on the weight-matrix \underline{W}_{xh} which is used to update the internal state of the cell according to the current input signal at each time-step.

$$\frac{\partial J}{\partial \underline{W}_{xh}} = \sum_t \frac{\partial J_t}{\partial \underline{W}_{xh}} \quad (9)$$

We may focus on time-step $t=2$:



CHAIN RULE FOR DERIVATIVES:

$$\frac{\partial J_2}{\partial \underline{W}_{xh}} = \frac{\partial J_2}{\partial \hat{y}_2} \cdot \frac{\partial \hat{y}_2}{\partial h_2} \cdot \frac{\partial h_2}{\partial \underline{W}_{xh}} \quad (10)$$

But, h_2 may be written as:
 $h_2 = f(\underline{W}_{xh} \cdot \underline{x}_2 + \underline{W}_{hh} \cdot h_1)$ (11)
where h_1 also depends on \underline{W}_{xh} . Therefore, the term $\frac{\partial h_2}{\partial \underline{W}_{xh}}$ cannot be treated as a constant.

We need to know how exactly the cell-state at $t=2$ depends on the weight-matrix \underline{W}_{xh} .

We need to expand the term $\frac{\partial h_2}{\partial \underline{W}_{xh}}$ by considering/ substituting the corresponding expression for the total derivative w.r.t. \underline{W}_{xh} , $\frac{d h_2}{d \underline{W}_{xh}}$.

Equation (5) suggests that:
$$\underline{h}_2 = \underline{h}_2(\underline{h}_1(\underline{w}_{xv}), \underline{h}_0(\underline{w}_{xv}); \underline{w}_{xu}) \quad (12)$$

Thus, the total derivative of \underline{h}_2 w.r.t \underline{w}_{xu} may be written according to the chain rule as:

$$\frac{d\underline{h}_2}{d\underline{w}_{xu}} = \frac{\partial \underline{h}_2}{\partial \underline{w}_{xu}} + \frac{\partial \underline{h}_2}{\partial \underline{h}_1} \cdot \frac{\partial \underline{h}_1}{\partial \underline{w}_{xu}} + \frac{\partial \underline{h}_2}{\partial \underline{h}_0} \cdot \frac{\partial \underline{h}_0}{\partial \underline{w}_{xu}} \quad (13)$$

Eq. (13) suggest the following compact form:

$$\frac{d\underline{h}_2}{d\underline{w}_{xu}} = \sum_{k=0}^2 \frac{\partial \underline{h}_2}{\partial \underline{h}_k} \cdot \frac{\partial \underline{h}_k}{\partial \underline{w}_{xu}} \quad (14)$$

* Apparently, the summation term for $k=2$ yields:

$$\frac{\partial \underline{h}_2}{\partial \underline{h}_2} \cdot \frac{\partial \underline{h}_2}{\partial \underline{w}_{xu}} = \frac{\partial \underline{h}_2}{\partial \underline{w}_{xu}} \quad (a)$$

Finally, Eq. (10) may be written as:

$$\frac{\partial J_2}{\partial \underline{w}_{xu}} = \sum_{k=0}^2 \frac{\partial J_2}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial \underline{h}_2} \cdot \frac{\partial \underline{h}_2}{\partial \underline{h}_k} \cdot \frac{\partial \underline{h}_k}{\partial \underline{w}_{xu}} \quad (15)$$

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

⊙ In this setting, the general form of Eq. (15) may be written as:

$$\frac{\partial J_t}{\partial \underline{W}_{xu}} = \sum_{k=0}^t \frac{\partial J_t}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial h_k} \cdot \frac{\partial h_k}{\partial \underline{W}_{xu}} \cdot \frac{\partial h_k}{\partial \underline{W}_{xu}} \quad (16)$$

↓ Contributions of \underline{W}_{xu} in previous timesteps to the error at timestamp t .

⊙ What about $\frac{\partial J}{\partial \underline{W}_{uu}}$ and $\frac{\partial J}{\partial \underline{W}_{uv}}$?

(*) Keep in mind that:

- (a): h_t does depend on \underline{W}_{uu}
- (b): h_t does not depend on \underline{W}_{uv} .

* Gradient of the loss function at time t :

$$\frac{\partial J_t}{\partial W_{xh}} = \sum_{k=0}^t \frac{\partial J_t}{\partial \hat{y}_k} \cdot \frac{\partial \hat{y}_k}{\partial h_k} \cdot \left[\frac{\partial h_t}{\partial h_k} \right] \cdot \frac{\partial h_k}{\partial W_{xh}} \quad (16)$$

requires the computation of the term:

$$\frac{\partial h_t}{\partial h_k} \quad (17)$$

* For example at $t=2$ and for $k=0$, we would have to compute the contribution of the term $\frac{\partial h_2}{\partial h_0}$ which according to the interdependence of the previous 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, if we are focusing on a timestep very far away in the future?

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

* The general form of Eq. (19) can be written 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 reveals the necessity to compute the term:

$$\frac{\partial h_t}{\partial h_{t-1}} \quad (21)$$

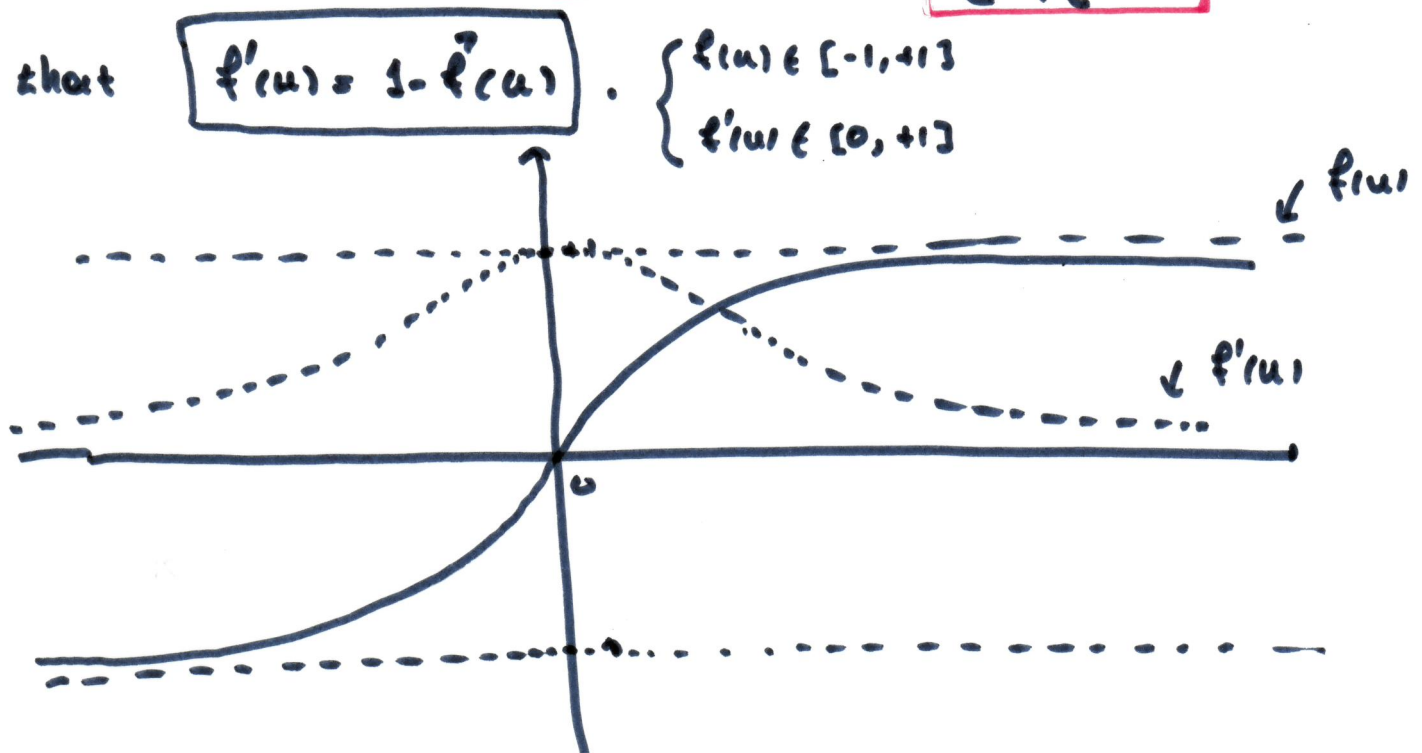
* For the case of a trivial (vanilla) RNN with a single neuron, we may assume that:

- (1): $\underline{x}_t \in \mathbb{R}^{n \times 1}$, $1 \leq t \leq 2$
- (2): $y_t \in \mathbb{R}$, $1 \leq t \leq 2$
- (3): $\underline{h}_t \in \mathbb{R}$, $1 \leq t \leq 2$ ($\underline{h}_t \equiv h_t$)
- (4): $\underline{W}_{xn} \in \mathbb{R}^{n \times 1}$, $1 \leq t \leq 2$ ($\underline{W}_{xn} \equiv \underline{W}_{xn}$)
- (5): $\underline{W}_{nn} \in \mathbb{R}$, $1 \leq t \leq 2$ ($\underline{W}_{nn} \equiv \underline{W}_{nn}$)
- (6): $\underline{W}_{ny} \in \mathbb{R}$, $1 \leq t \leq 2$ ($\underline{W}_{ny} \equiv W_{ny}$)

* According to the previous definitions, we may write for the forward pass of the information within the network that:

$$\begin{cases} h_t = f(\underline{W}_{xn}^T \cdot \underline{x}_t + \underline{W}_{nn} \cdot h_{t-1}) & [27] \\ y_t = W_{ny} \cdot h_t & [28] \end{cases}$$

* At this point mention (Page #2) of previous notes in order to prove that for: $f(u) = \tanh(u) = \frac{e^u - e^{-u}}{e^u + e^{-u}}$ we have



⊛ In 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:

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

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

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

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

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

$$\frac{\partial h_t}{\partial h_{t-1}} = [1 - f'(\underline{W}_{xn}^T \cdot \underline{x}_t + W_{nn} \cdot h_{t-1})^2] \cdot W_{nn} \quad [27]$$

⊛ Eq. (26) suggest that:

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

(ii): W_{nn} are sampled from the standardized normal distribution (for the i.i.d case, $\mu = 0, \sigma^2 = 1$) so that $W_{nn} < 1$.

⊛ We are multiplying a lot small numbers together.

⊛ Errors due to further back timesteps have increasingly smaller gradients.

⊛ Weight-parameters become biased towards capturing shorter-term dependencies.