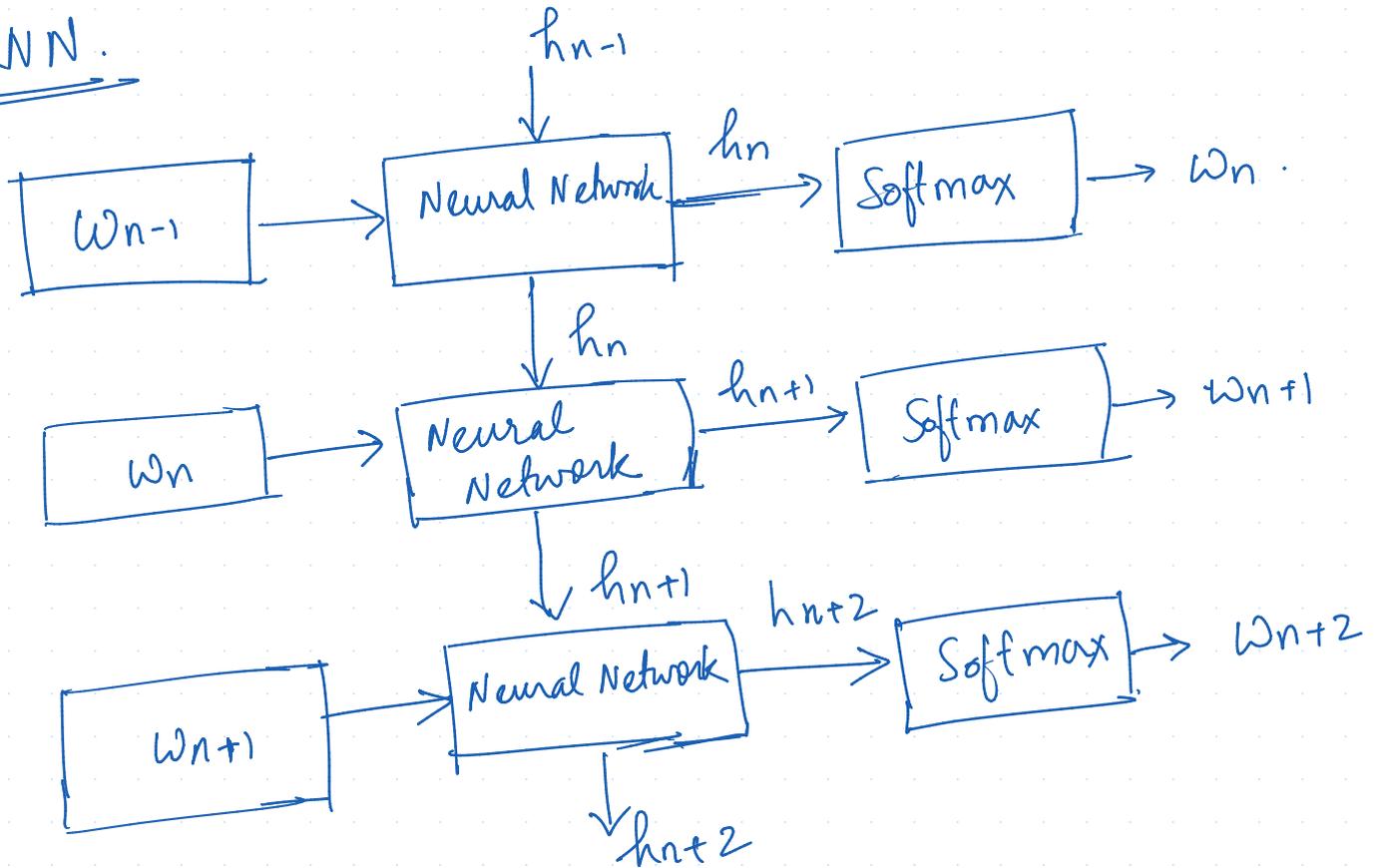# Lecture - 9

## RNN.



$$h_n = \tanh(W\, x_{n-1} + b)$$

$$f(w_n \mid w_{n-1}, h_{n-1}) = \text{Softmax}(U h_n + \beta).$$
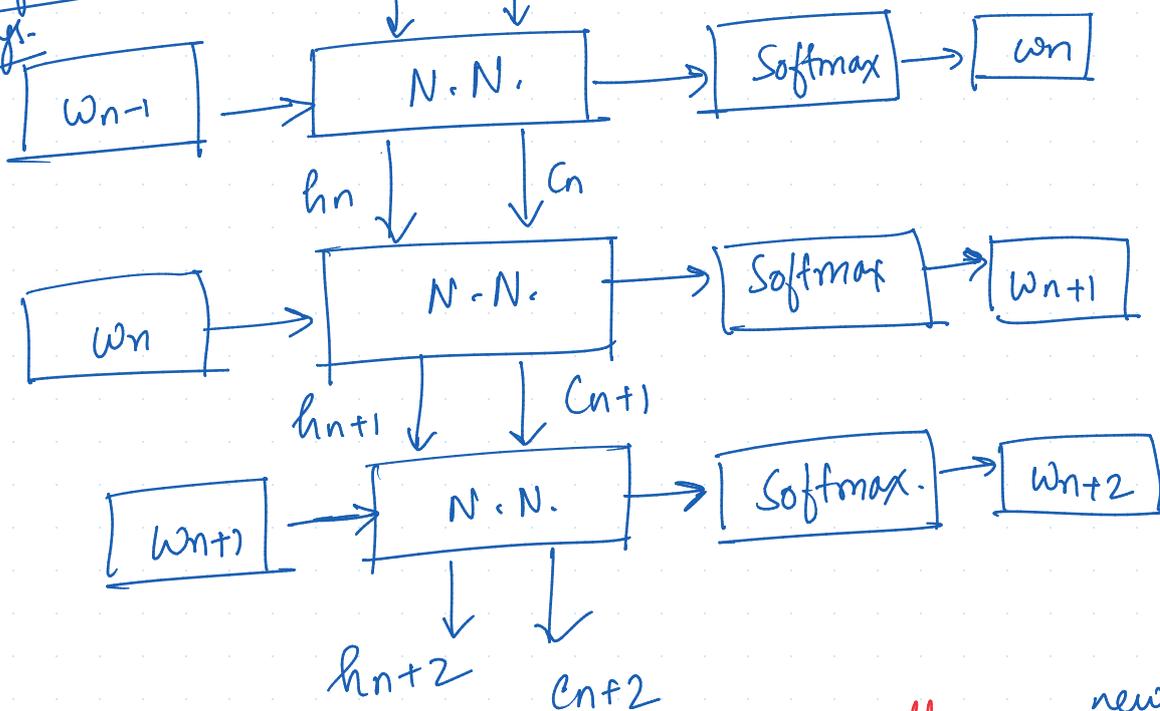
$$w_1, w_2, \ldots, \quad w_{n-1}, \quad w_n$$

# Long - short - term Memory

model predicted.
next word.

sequence of word embeddings.

$h_{n-1}$    $C_{n-1}$ memory cells.

| $W_{n-1}$ | → | N.N. | → | Softmax | → | $W_n$ |

$h_n$ ↓    $C_n$ ↓

| $W_n$ | → | N.N. | → | Softmax | → | $W_{n+1}$ |

$h_{n+1}$ ↓    $C_{n+1}$ ↓

| $W_{n+1}$ | → | N.N. | → | Softmax. | → | $W_{n+2}$ |

$h_{n+2}$    $C_{n+2}$

update   $C_n = f_n \odot C_{n-1} + i_n \odot \tilde{C}_n$

memory cell previous word

$f_n \odot C_n$   forgotten info.

memory cell   $C_n$ → (-1, 1).   new updated memory cell.

$i_n \odot \tilde{C}_n$   new input.   $C_n$

| $C_{n-1}$ |

$O_n \odot \tanh(c_n)$

forget gate   (0-1)   input gate   (0-1)

$f_n = \sigma(W_f x_{n-1} + b_f)$

$i_n = \sigma(W_i x_{n-1} + b_i)$

$O_n = \sigma(W_o x_{n-1} + b_o)$   (0-1).

$\tilde{C}_n = \tanh(W_c x_{n-1} + b_c)$   (-1, 1).

output gate

| $x_{n-1} = [W_{n-1}, h_{n-1}]$ |

input previous word

$h_n$ hidden vectors

Control networks
↳ $f_n$, $i_n$, $O_n$.
↳ controls the amount of information within the models.

$$[ x_1, x_2, \cdots \cdots, x_n ].$$

weighted
multiplication $[ \underset{\times}{0.5}, \underset{\times}{0.1}, \cdots 1 \cdots \underset{\times}{x \cdot 3} ] \rightarrow$

↓      ↓      ↓      ↓

higher    lower - full weightage   less weightage
weightage

## forget gate      (0-1).

- If we want to forget some information, we multiply
  it with a very low value. $(0 \cdot 1, 0 \cdot 001)$

- If we want to retain some information we multiply
  it with some high value $(1, 0 \cdot 99)$.

- Input gate controls the amount of new data to
  the neural network.

- Output gate — controls the degree to which the
  memory cell goes to the o/p of the hidden vector.

$\underline{I/p}$ :    $\underline{x_t}, \underline{h_{t-1}}, c_{t-1}$ (memory cell).    $\underline{(LSTM)}$

$$g_t = \tanh(W_c x_t + V_c h_{t-1})$$

$$f_t = \sigma(W_f x_t + V_f h_{t-1}) \quad \rightarrow \text{forget gate.}$$

$$i_t = \sigma(W_i x_t + V_i h_{t-1}) \quad \text{i/p gate}$$

$$c_t = f_t c_{t-1} + i_t g_t$$

$$o_t = \sigma(W_o x_t + V_o h_{t-1}) \quad \text{o/p gate } -$$

$$h_t = o_t \tanh(c_t)$$

$\underline{GRU}$    Input : $x_t, h_{t-1}.$    : no memory cell.

$$z_t = \sigma(W_z x_t + U_z h_{t-1}) \quad \rightarrow \text{update gate}$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1}) \quad \rightarrow \text{reset gate.}$$

$$\tilde{h}_t = \tanh(W x_t + U(r_t \odot h_{t-1}))$$

$$h_t = (1 - z_t) h_{t-1} + z_t \tilde{h}_t$$

# Different types of attention Mechanism :-

## Additive attention

$$a(q,k) = w_v^T \tanh(w_q \cdot \underline{q} + w_k \cdot \underline{k}).$$

## General

$$a(q,k) = q^T \underline{W_a} k$$

## Dot-product attention

$$a(q,k) = q^T k.$$

## Scaled - dot product attention

$$a(q,k) = \boxed{\frac{q^T k}{\sqrt{d_k}}} \quad \rightarrow \text{softmax in transformers.}$$

to lie b/w $(0-1)$

## Transformer.

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{Q K^T}{\sqrt{d_k}}\right) V.$$

Hard attention : non-differentiable/ gradient estimation $(RL)$

Soft Attention : differentiable / softmax based.