Let us consider a statement – 'reason with logic.'
[Source statement in English]

Translated statement in Hindi – 'तर्क के साथ कारण'

Here $T_x = 3$ & $T_y = 4$

We associate a hidden state for each word of English statement

| Word | hidden state | Encoding |
|---|---|---|
| reason | $h_1^E$ | $[0.1, 0.7, 0.5, 0.3]$ |
| with | $h_2^E$ | $[-0.1, -0.7, 0.4, 0.6]$ |
| logic | $h_3^E$ | $[0.7, 0.2, -0.3, 0.4]$ |

Here $h_i^E$ is a vector of 4 dimensions.

We know we can get $h_i^E$ as below.

$$h_i^E = f(W_E h_{i-1}^E + U x_i + b)$$

Here $W_E$ is of order $4 \times 4$
$h_{i-1}^E$ is of order $4 \times 1$
$U$ is of order $4 \times V_{in}$   $\Big\}$ $V_{in}$ is input vocab
$x_i$ is of order $V_{in} \times 1$
$b$ is of order $4 \times 1$
$f$ is an activation function [like tanh, sigmoid etc]

$30,000$   $\begin{pmatrix} 0 \\ 0 \\ \hline 1 \\ \vdots \\ 0 \end{pmatrix} \rightarrow$ reason

$h_t = g(W h_{t-1} + U_{x_t} + b)$

[Note:- In paper, they used $h_t = f(x_t, h_{t-1})$ [Page 2]]
$\quad \hookrightarrow h_i \in \mathbb{R}^n$

... to predict the word 'तर्क'

$\langle 300 \rangle$

We suppose that we need to predict the word 'तक'.

The hidden state be $s_D = [0.3, 0.4, 0.1]$

[Note :- We used an hidden state with dimension = 3 ]

The hidden state from encoder was of size $4 \times 1$

We will consider a matrix of size $3 \times 4$ to convert the $4 \times 1$ vector $h_i$ to $3 \times 1$

The transformed vectors are.

$$\hat{h}_1^E = W_c \, h_1^E$$

$$\hat{h}_2^E = W_c \, h_1^E$$

$$\hat{h}_3^E = W_c \, h_3^E$$

[Here $\boxed{W_c}$ is of order $\boxed{3 \times 4}$
Wc will be a trainable matrix]

$$(W_c) \, h_1^E \rightarrow \underline{4 \times 1}$$
$$\phantom{(W_c)}_{3 \times 4}$$

$$h_1^E \rightarrow 3 \times 1 \qquad h_1^E, \; h_2^E, \; h_3^E$$
$$\underbrace{\phantom{h_1^E, \; h_2^E, \; h_3^E}}_{\underline{\underline{s_0}}}$$

We consider $W_{dec}$ & $W_{enc}$ as two matrices such that

$$\boxed{\begin{array}{l} W_{enc} \, \hat{h}_1^E + W_{dec} \, s_{0} = \underline{z_1} \\[2ex] W_{enc} \, \hat{h}_2^E + W_{dec} \, s_0 = \underline{z_2} \\[2ex] W_{enc} \, \hat{h}_3^E + W_{dec} \, s_0 = \underline{z_3} \end{array}}$$

[Here $W_{enc}$ &
$W_{dec}$ are of order
$1 \times 1$ ]

We have

$$e_{11} = V_a \, \text{Tanh}(z_1)$$
$$e_{12} = \overline{V_a} \, \text{Tanh}(z_2)$$
$$e_{13} = V_a \, \text{Tanh}(z_3)$$

$$\left[ \text{Here } V_a \text{ is } |X| \text{ vector} \right]$$

Now $a_{11} = \dfrac{\exp(e_{11})}{\sum\limits_{j=1}^{T_x} a_{ij} \exp(e_{1j})} = \dfrac{\exp(e_{11})}{A} = 0.7$

where $A = \sum\limits_{j=1}^{T_x} \exp(e_{1j})$

$$a_{12} = \exp(e_{12})/A = 0.7$$
$$a_{13} = \exp(e_{13})/A = 0.1$$

Now $c_1 = a_{11} \, h_1^E + a_{12} \, h_2^E + a_{13} \, h_3^E$

$$= 0.7 \,(0.1, 0.7, 0.5, 0.3) +$$
$$0.7 \,(0.1, -0.7, 0.4, 0.6) +$$
$$0.1 \,(0.7, 0.2, -0.3, 0.4)$$

$$\boxed{c_1 = (0.12, 0.37, 0.4, 0.37)}$$

The decoder predictions are done in following way:

$$s_1 \, \overset{y}{_1} = g(y_0, s_0, c_1)$$

$$\langle DOS \rangle,$$
$$s_0 = \begin{pmatrix} 0.3 \\ 0.4 \\ 0.5 \end{pmatrix} \, 3 \times 1$$

$$\begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \langle SOD \rangle$$
$$20\,000$$

$$s_1 = \overset{y}{_1}$$

$s_1 =$   `'6,5 3×1`

$$g\left(W_D \begin{bmatrix} s_0 \end{bmatrix}_D + U_D \begin{pmatrix} y_0 \end{pmatrix} + \begin{pmatrix} W_a & c_1 \end{pmatrix} + b_D\right)$$

$$\begin{bmatrix} \text{Here; } & s_0 \text{ is of order } 3\times 1 \\ & W_D \text{ is of order } \underline{3\times 3} \\ & U_D \text{ is of order } 3\times V_{out} \\ & y_0 \text{ is of order } \underline{V_{out}\times 1} \\ & c_1 \text{ was of order } \underline{4\times 1} \\ & W_a \text{ is of order } 3\times 4 \\ & b_D \text{ is of order } 3\times 1 \end{bmatrix}$$

$W_a$   $\underline{\underline{4\times 1}}$
$3\times 4$

The prediction $y_1$ will be done as → $3\times 1$

$$y_1 = softmax\left(V\left(s_1\right) + c\right) \qquad \begin{array}{l} \to b_{1\,\infty} \\ \to V_{out}\times 3 \end{array}$$

$\underline{\underline{20000\times 3}}$

$$\begin{bmatrix} \text{Here } V \text{ is of order } V_{out}\times 3 \\ s_1 \text{ is of order } 3\times 1 \\ c \text{ is of order } \underline{V_{out}\times 1} \end{bmatrix}$$

$$\begin{pmatrix} 0.001 \\ 0.00 \\ \boxed{0.87} \\ 0 \\ 0 \\ 0.07 \end{pmatrix}$$

$c$ तक ←

Now, suppose argmax $y_1 = \dfrac{c}{\text{तक}}$ ←

We will then proceed to find $y_2$.

| Reason | truth | logic |

$$h_1^E = f(W_E h_0^E + Ux_1 + b)$$

$$h_2^E = f(W_E h_1^E + Ux_2 + b)$$

$$= h_3^E = f(W_E h_2^E + Ux_3 + b)$$

$$h_1^E \qquad h_2^E \qquad h_3^E$$

$$\alpha_{11} = 0.7 \qquad \alpha_{12} = 0.2 \qquad \alpha_{13} = 0.1$$

$$c_1 = [0.12, 0.37, 0.4, 0.37]$$

Decoder: $s_1 = g(W_D s_0 + U_D y_0 + W_a C_1 + b_D)$

$y_1 = \text{softmax}(Vs_1 + c)$

$$\langle start \rangle$$
$$y_0$$
$$=$$

$$तर्क$$
$$y_1$$