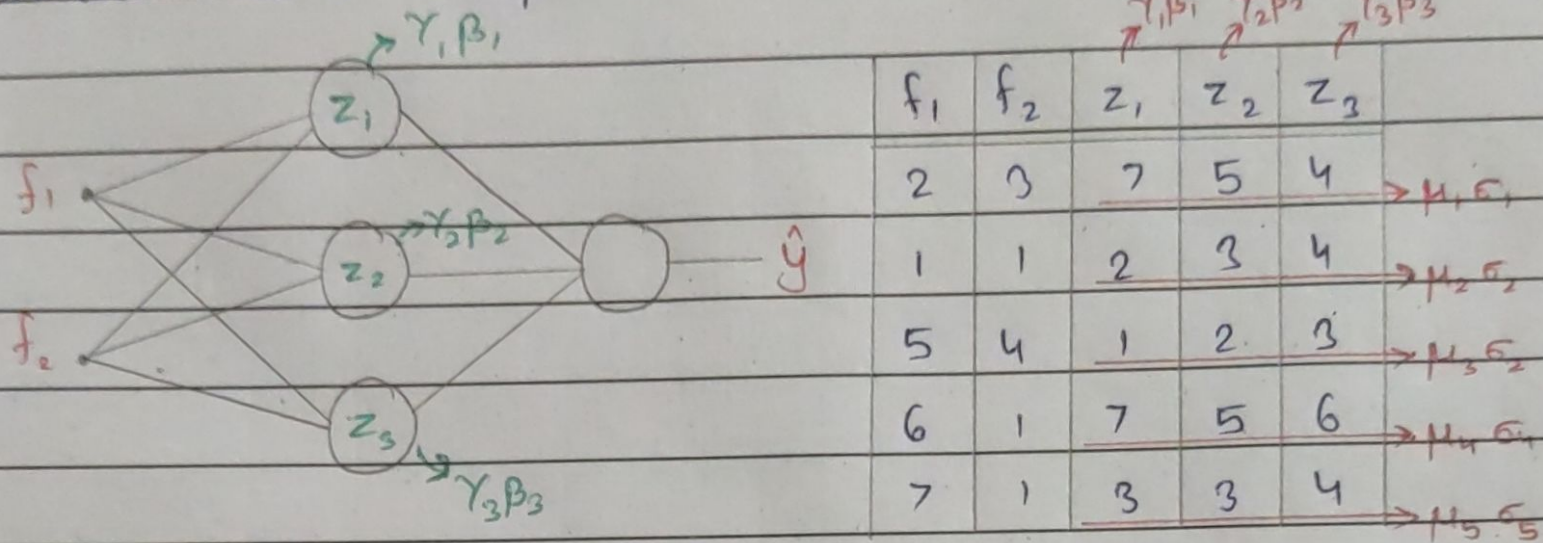


Consider same setup as of Batch Normalization



In layer normalization  $\mu$  &  $\sigma$  are calculated across features (row wise) & are readjusted using the  $\gamma$  &  $\beta$  values of the respective columns (values of the pre activation) ( $z_1 \rightarrow \gamma_1 \beta_1$ ,  $z_2 \rightarrow \gamma_2 \beta_2$ ,  $z_3 \rightarrow \gamma_3 \beta_3$ )

In same way for normalizing the values of text embedding matrix we use Layer Normalization

	$\gamma_1 \beta_1$	$\gamma_2 \beta_2$	$\gamma_3 \beta_3$	
Hi	6.5	2.41	3.21	$\rightarrow \mu_1 \sigma_1$
Probal	2.21	0.4	3.6	$\rightarrow \mu_2 \sigma_2$
<Pad>	0	0	0	$\rightarrow \mu_3 \sigma_3$
<Pad>	0	0	0	$\rightarrow \mu_4 \sigma_4$
How	7.5	9.2	1.5	$\rightarrow \mu_5 \sigma_5$
are	2.2	1.1	6.7	$\rightarrow \mu_6 \sigma_6$
you	2.9	6	9	$\rightarrow \mu_7 \sigma_7$
today	9.9	2.3	6.5	$\rightarrow \mu_8 \sigma_8$

Here, as we can see, each embedding vector has it's own statistical representation ( $\mu$  &  $\sigma$ ) & the zero padding embedding vectors are not affecting the true statistical representations of the embedding matrix.

Eg. For 1st dimension of embedding vector of word "How" -

$$\frac{7.5 - \mu_5}{\sigma_5} = 0.7 \rightarrow \text{Hypothetical}$$

$\left\{ \begin{array}{l} \text{Row wise } \mu_5 \text{ \& } \sigma_5 \\ \text{Column wise } \gamma_1 \text{ \& } \beta_1 \end{array} \right.$

Readjustment:  $0.7 \gamma_1 + \beta_1$