$n \rightarrow$ number of filters

$m \rightarrow$ number of words in a sentence

$k \Rightarrow$ Kernel size of filters

$emb\_dim \longrightarrow$ embedding dimension

# How does ConvID work with word embedding?

Suppose we have a Sentence consisting of 'm' words where each word has been represented using word embedding :

| Sample Sentence | feature (i.e. word embedding) |
|---|---|
| word-1 | |
| word-2 | |
| word-3 | |
| word-4 | |
| word-5 | |
| word-6 | |
| . | |
| . | |
| . | |
| . | |
| . | |
| . | |
| word_m-1 | |
| word_m | |

जब दुःख अपनी चरम सीमा पर होता है, तब सुख ज्यादा दूर नहीं होता।

Now we would like to apply
1D convolution layer consisting
of 'n' different filters with kernel
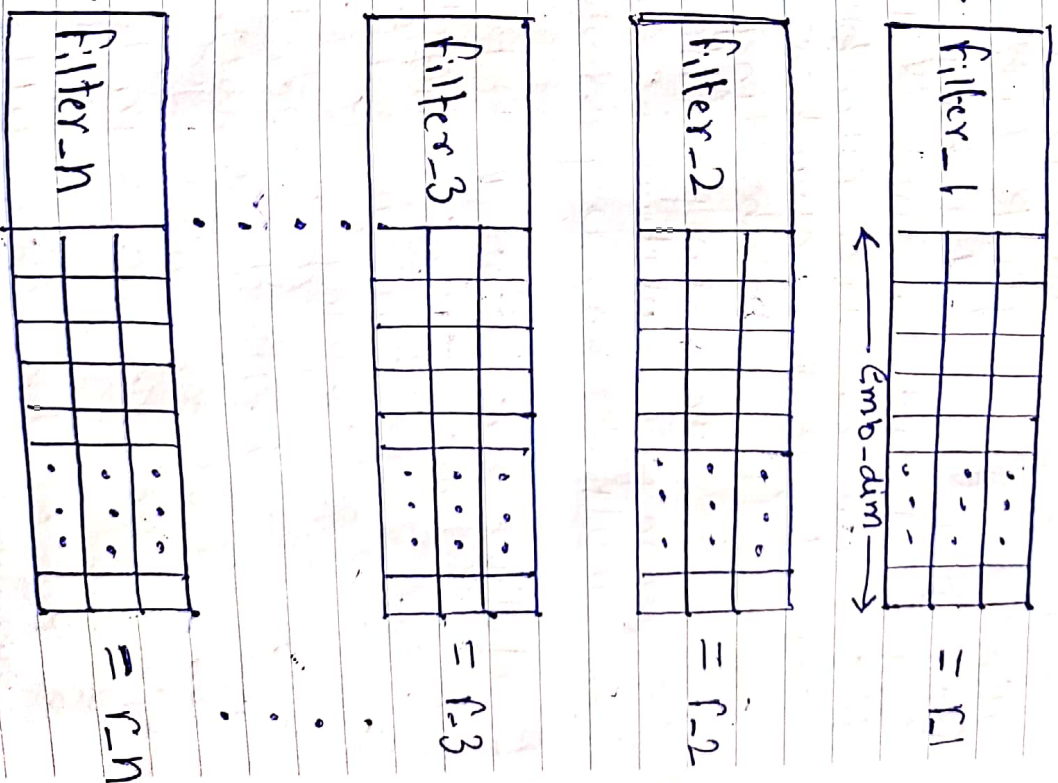size of 'K' on this data. To do so,
Sliding windows of length 'K'
are extracted from the data and
then each filter is applied on each
of those extracted windows.
Here is an illustration of what
happens.
Here, I have assumed $K=3$ and
removed the bais parameter
of each filter for simplicity.

Window of length K (=3)

filter with kernal size K (=3)

filter response (convolution result)



word-i
word-(i+1)
word-(i+2)

← emb-dim →

\* 

Element-wise multiplication (convolution)

filter-1 ← emb-dim → = f-1

filter-2 = f-2

filter-3 = f-3

filter-n = f-n

As we can see in the figure above, the response of each filter is equivalent to the result of its convolution (i.e element-wise multiplication and then summing all the results) with the extracted window of length K (i.e ith to (i+K-1)-th words in the given sentence). further, note that each filter has the same number of channels as the number of feature (i.e word-embedding dimension) of the training sample (hence performing convolution, i.e element-wise multiplication, is possible).

Essentially, each filter is detecting the presence of a particular feature of pattern in a LOCAL window of training data (eg whether a couple of specific words exists in this window or not). After all the filters have been applied on all the window of length 'K' we would have an Output of like this which is the result of convolution.

| Filter response | Filter-1 | Filter-2 | . . . . . | Filter-n |
|---|---|---|---|---|
| window-1 | | | | |
| window-2 | | | | |
| window-3 | | | | |
| window-4 | | | | |
| window-5 | | | | |
| window-6 | | | | |
| window-7 | | | | |
| window-8 | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| | | | | |
| window-(m-k+1) | | | | |

As we can see, there are $m-k+1$ windows in the figure since we have assumed that the padding = 'Valid' and Stride=1 (default behavior of Con1D layer in Keras). The Stride argument determines how much the window should Slide (i.e shift) to extract the next window (e.g in our ~~example~~ example above, a Stride of 2 would extract windows of words: (1,2,3), (3,4,5), (5,6,7)..... instead). The padding argument determines whether the window should entirely consists of the words in training Sample or there Should be padding at the beginning and at the end: this way, the Convolution ~~respone~~ response may have the Same length (i.e. $m$ and not $m-k+1$) as the training Sample (e.g. In our example above, padding = 'Same' would extract window of words: (PAD, 1, 2), (2, 3, 4),.....,(m-2, m-1, PAD)).

number of parameters in the

1-D convolution layer is equal to :

num_filters * (Kernal_size * n_features)
+ one_bias_per_filter

$$= n * (K * emb\_dim) + n$$

for example, if, $n=32$, $m=20$, $k=3$,

$$emb\_dim = 100$$

$$= n * (K * emb\_dim) + n$$

$$= 32 * (3 * 100) + 32$$

$$= \underline{\underline{9632}}$$