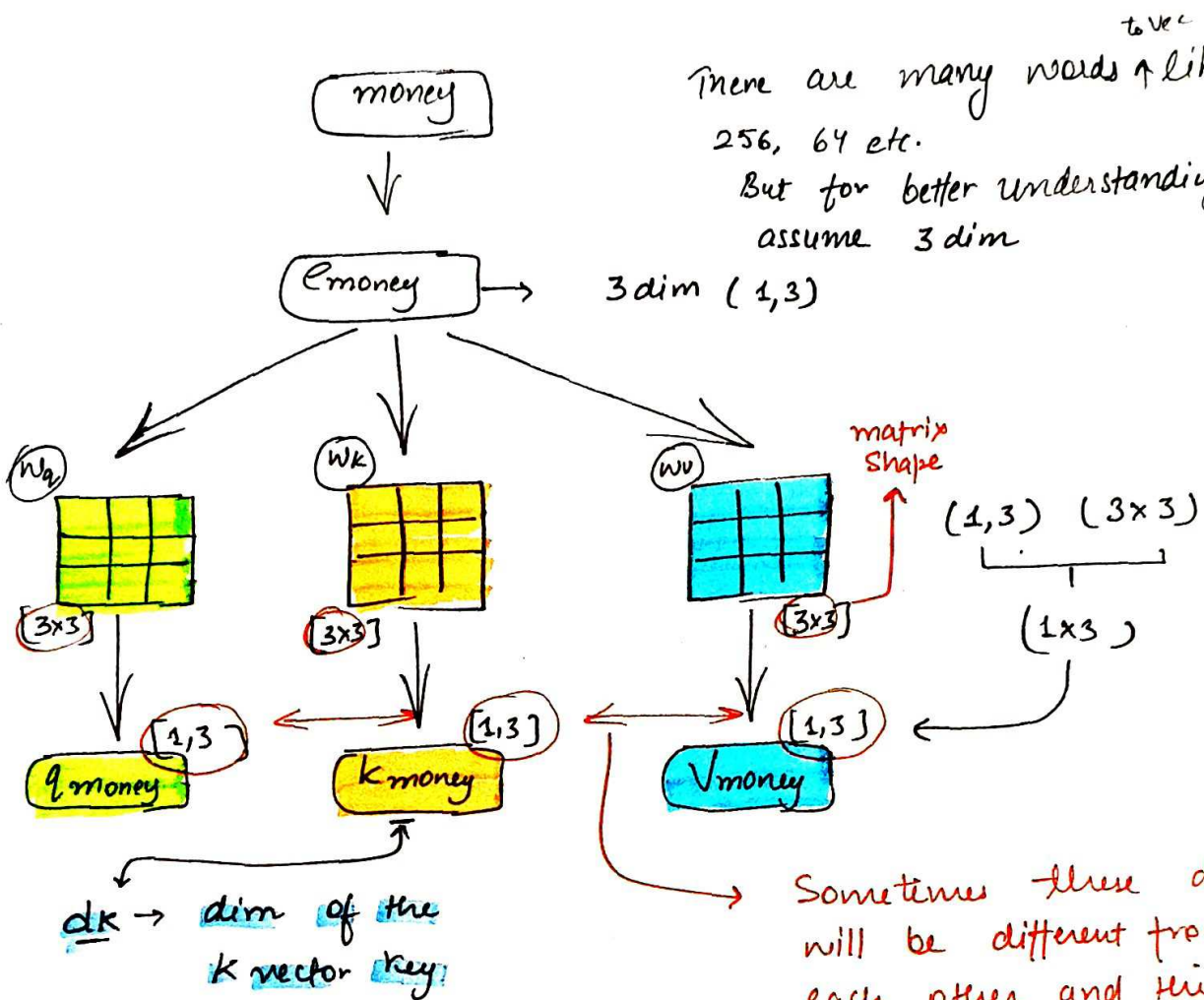


# Scalar Dot Product Attention



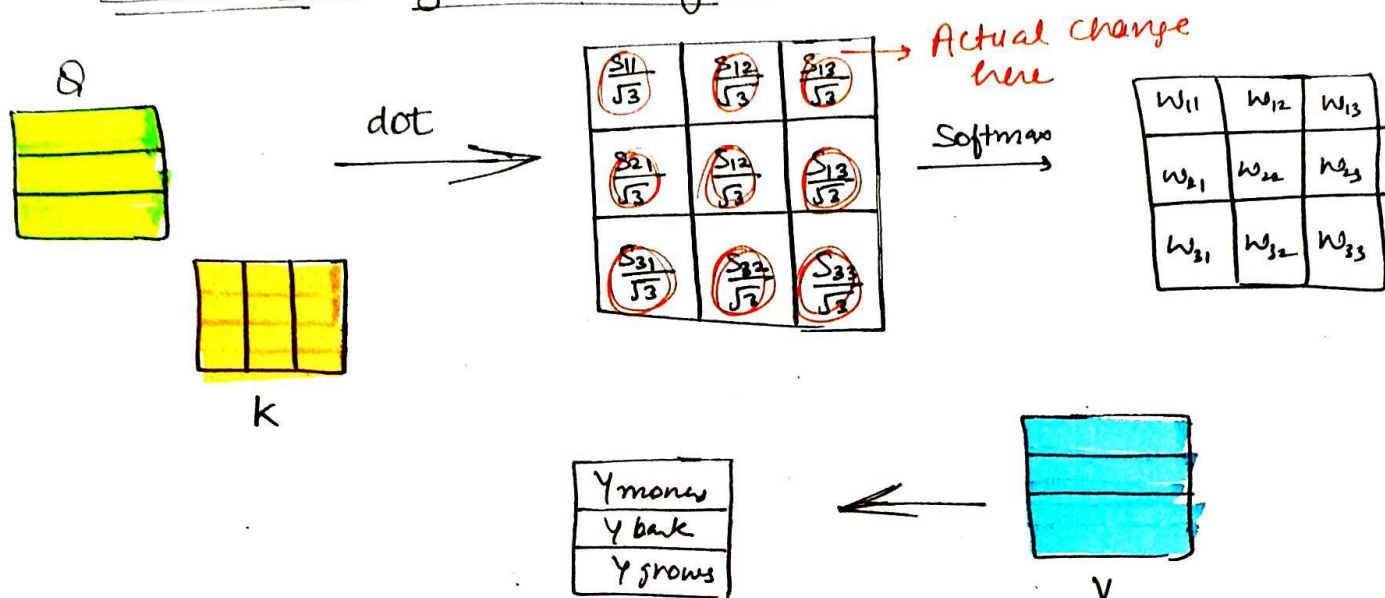
Sometimes these dim will be different from each other and this is depend on matrix shape

In our case,  
 $d_k = d_q = d_v = 3$

$$\text{Attention}(Q, k, V) = \text{Softmax}\left(\frac{Qk^T}{\sqrt{d_k}}\right)V \quad \rightarrow \text{Scaled formula}$$

$$\text{Attention}(Q, k, V) = \text{Softmax}\left(\frac{Qk^T}{\sqrt{3}}\right)V \quad \text{in our case}$$

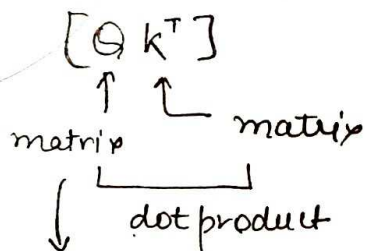
# Actual meaning in diagram



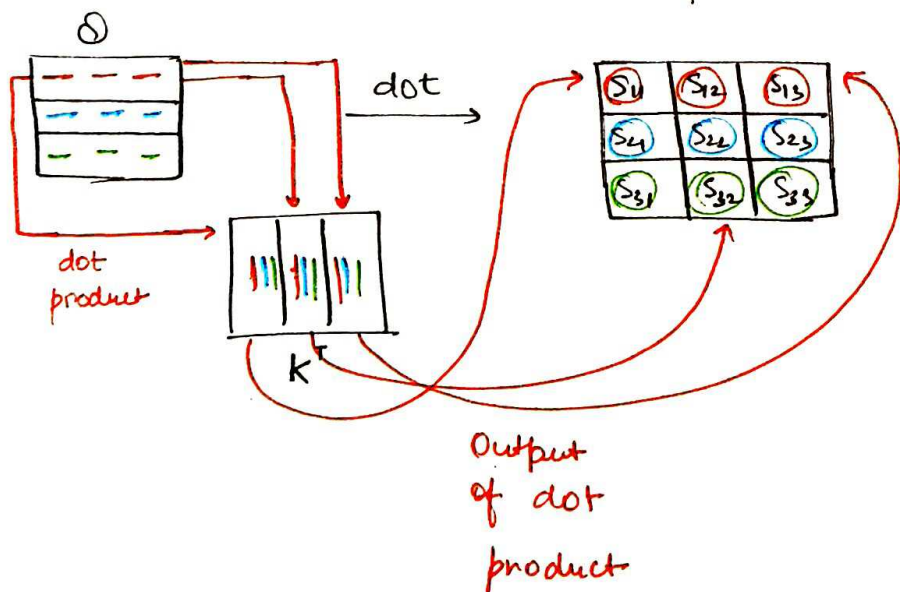
Why?

↳ dot-product ka Nature

In this formula



How many vectors in this dot product?



Total 9 vector-vector dot products

which means we get 9 Number.

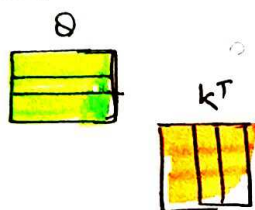
If we have Number then we can easily calculate mean and variance.

Now, what is the meaning of "Dot-product ka Nature"?

(i) low dimension vector  $\rightarrow$  dot product of  $\xrightarrow{\text{is}}$  low dimension vector  $\rightarrow$  low variance

(ii) High dimension vector  $\rightarrow$  dot product of  $\xrightarrow{\text{is}}$  high dimension vector  $\rightarrow$  high variance

Example



let say,  
we have ③ dim vector  
 $\theta(1,3)$   $\leftarrow$  vector  
 $(1,3)$   
 $k^T$

$\rightarrow$  ⑨ dot product values  $\rightarrow$  3 dim vector is very low dim

then variance of 9 dot product value is also low.

let say,

we have ⑤⑫ dim vector  $\rightarrow$  still we get 9 dot product value

$\hookrightarrow$  But 512 dim vector is very high so, dot product variance is very high.

Another example

$\rightarrow$  2-D

1>	[1,2]	[3,2]
2>	[ ]	[ ]
3>	[ ]	[ ]
4>	[ ]	[ ]
5>	[ ]	[ ]

$\begin{bmatrix} a \\ b \\ c \\ d \\ e \end{bmatrix} \rightarrow \sigma_1^2$

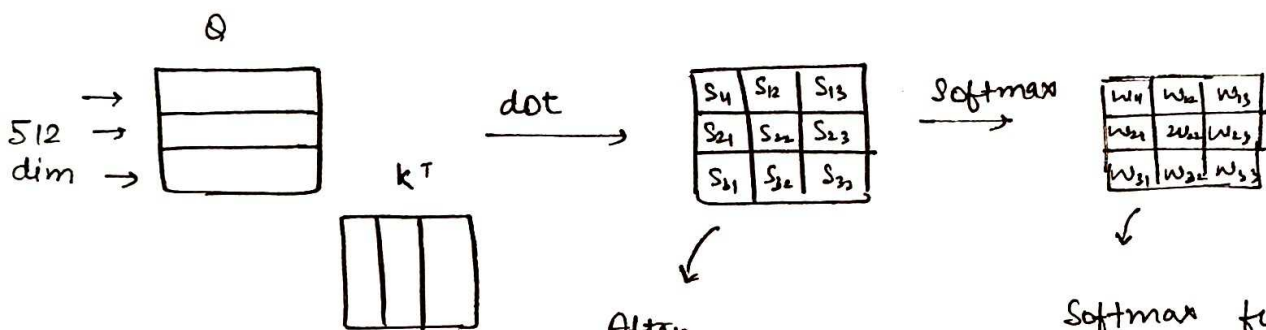
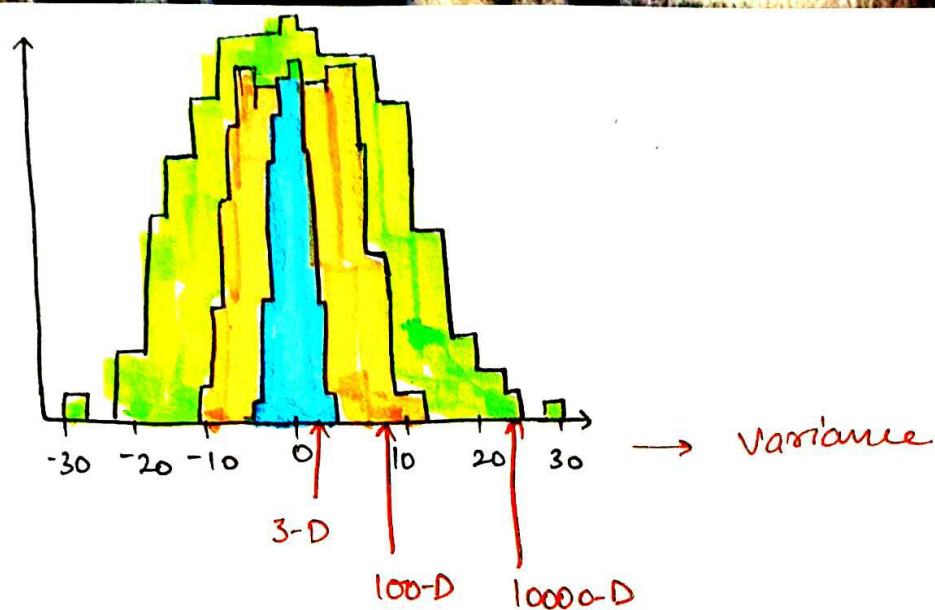
$\rightarrow$  3-D

1>	[1,2,3]	[3,2,1]
2>	[ ]	[ ]
3>	[ ]	[ ]
4>	[ ]	[ ]
5>	[ ]	[ ]

$\begin{bmatrix} f \\ g \\ h \\ i \\ j \end{bmatrix} \rightarrow \sigma_2^2$

$\sigma_2^2$  variance is greater than  $\sigma_1^2$





After dot product variance will be higher.  
(like some values are very high and some values are very low)

Softmax function give small value to small probability and large value to large prob.

during backpropagation gradient descent only focus on large value after Softmax probab and neglect small value. This will cause vanishing gradient problem.

Training process  $\rightarrow$  kharab ho gya.

Softmax: [4, 5] low variance

Small probab [0.45, 0.55]

[1, 10] large variance

[0.001, 0.999]

Small prob  $\hookrightarrow$  large probab.

Real life example :- large variance height student in class. Teacher ask to hand raise for doubt. Because of large height student small height student's hand not see. Only large height student doubt solve and small height student's doubt is ignore. So, Training on the basis of large student.

Now, let's imagine. Student's is almost same in the class. And Hand of every student is easily seen by Teacher. Doubt is solved of every student. So, Training on the basis of all student. Variety of question is increase.

Problem :-

$S_{p1}$	$S_{p2}$	$S_{p3}$
$S_{21}$	$S_{22}$	$S_{23}$
$S_{31}$	$S_{32}$	$S_{33}$

} Variance is very high

One solution is to use low dim vector from startly but in low vector we cannot get good information.

Softmax

$w_{11}$	$w_{12}$	$w_{13}$
$w_{21}$	$w_{22}$	$w_{23}$
$w_{31}$	$w_{32}$	$w_{33}$

max probab

Probability is

(99%)

(1%)

min probab

during backpropagation gradient descent focus on (99%) and ignore (1%).  
Vanishing gradient problem.

another solution is  $\rightarrow$  (i) [10, 20, 30, 40, 50, 60, 70]  
variance is 400.0

Scaled the value of

doted producted which variance is very high.

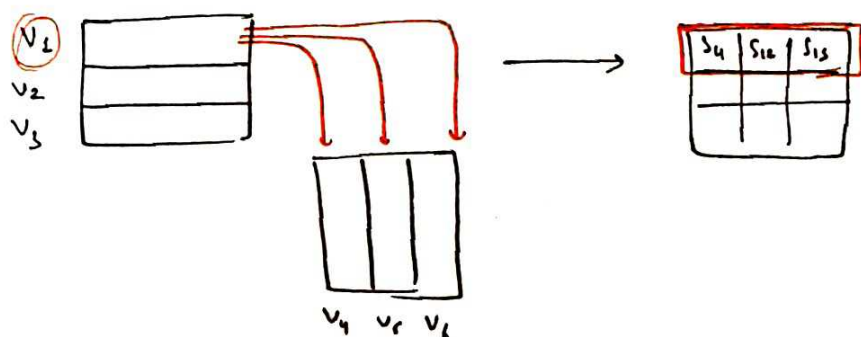
(ii) [1, 2, 3, 4, 5, 6, 7]

variance is 4.0

After scaling variance will be low

Now, we have to find that Scaling factor or number which is divide with value of dot product.

If dimension  $\uparrow\uparrow$  then how much variance  $\uparrow\uparrow$   
In math quantify:



let say all these vector are 1-D

$v_1 \rightarrow v_4$   
 $\quad \quad \quad \rightarrow v_5$   
 $\quad \quad \quad \rightarrow v_6$

1-D  
 $[a] \rightarrow [b]$   
 $\quad \quad \rightarrow [c]$   
 $\quad \quad \rightarrow [d]$

here we don't want Sample variance  
 we want expected Variance (population Variance) ( $\text{Var } x$ )

Now Increase dim  
from 1-D to 2-D

$[a \ b]$

$[c \ d] \rightarrow ac+bd$

$[e \ f] \rightarrow ae+bf$

$[g \ h] \rightarrow ag+bh$

$y \rightarrow \text{Var}(y)$   
 $\text{Var}(y) > \text{Var}(x)$

Now Increase dim  
from 2-D to 3-D

$[a \ b \ c]$

$[d \ e \ f] \rightarrow ad+be+cf$

$[g \ h \ i] \rightarrow ag+bh+ci$

$[k \ l \ m] \rightarrow ak+bl+cm$

$z \rightarrow \text{Var}(z)$

$$\text{Var}(z) > \text{Var}(y) > \text{Var}(x)$$

$$\text{Var}(z) \approx 3 \text{Var}(x)$$

If  $d$  dim then  $dn \rightarrow d \text{Var}(x)$

\* Somehow  $\begin{bmatrix} a_1 + b_1 \\ a_2 + b_2 \\ \vdots \end{bmatrix}$  use number se use divide karne hai ki every time  $\text{Var}(x)$  hi aaye

$$\left. \begin{array}{c} [a \ b] \quad [c \ d] \\ \quad \quad \quad \sqrt{2} \\ [e \ f] \\ \quad \quad \quad \sqrt{2} \\ [g \ h] \\ \quad \quad \quad \sqrt{2} \end{array} \right\} 2 \text{Var}(x)$$

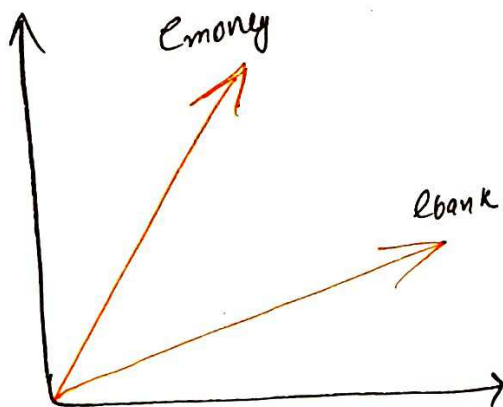
dim  $\swarrow$

$$\frac{1}{2} \text{Var}(y) \rightarrow \frac{1}{2} 2 \text{Var}(x) = \text{Var}(x)$$

$\sqrt{2} \rightarrow \sqrt{dk}$

$dk \rightarrow$  dimension

## Self Attention Geometric Intuition



money bank

$$\begin{array}{cc} \boxed{E_{\text{money}}} & \boxed{E_{\text{bank}}} \\ \hookrightarrow [2, 7] & [9, 3] \leftarrow \end{array}$$