

Query matrix

q _{money}
q _{bank}
q _{grows}

dot →

S ₁₁	S ₁₂	S ₁₃
S ₂₁	S ₂₂	S ₂₃
S ₃₁	S ₃₂	S ₃₃

k _{money}	k _{bank}	k _{grows}
--------------------	-------------------	--------------------

key^T matrix

Softmax ↓

Y _{money}
Y _{bank}
Y _{grows}

← dot

w ₁₁	w ₁₂	w ₁₃
w ₂₁	w ₂₂	w ₂₃
w ₃₁	w ₃₂	w ₃₃

Contextual Embeddings =

(A)

$$\text{Attention}(Q, K, V) = \text{softmax}(Q \cdot K^T) V$$

v _{money}
v _{bank}
v _{grows}

V matrix

↑
Is the mathematical representation of above process to calculate contextual embeddings.

Scaled Dot-Product :-

* Equation (A) is developed from scratch as you noticed guys! But when we headed to the research paper "Attention All you need". The formulation mentioned is:

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

1. Why they divided with $\sqrt{d_k}$?
2. What is d_k ?
3. How d_k is calculated?

QA:- d_k is the dimension of key vector

1A:- To solve the "unstable gradient" problem we divide with $\sqrt{d_k}$

I have two questions running in my mind!

1. What is unstable gradient?
2. Why we have divided with only $\sqrt{d_k}$, why not with d_k , d_q and d_v ?

1. Let me answer what is unstable gradient and vanishing gradient!

Ans: What If I say:

- * Low Dimensional Vector \propto Low Variance
- * High Dimensional Vector \propto High Variance

I think your face might be : 😞

Let me explain.....

For example:- you have 3 sets of vectors

In 2D

In 3D

1. $[1, 2] \cdot [1, 1] = 3$

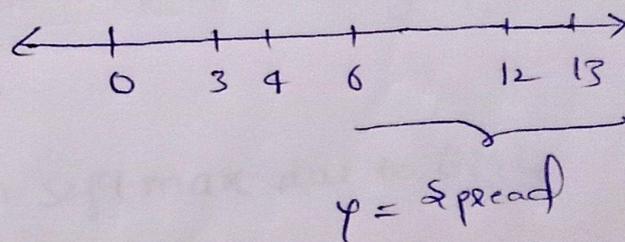
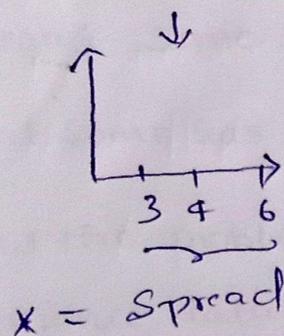
1. $[1, 2, 3] \cdot [1, 1, 1] = 6$

2. $[2, 1] \cdot [1, 2] = 2 + 2 = 4$

2. $[2, 1, 3] \cdot [1, 2, 3] = 13$

3. $[2, 2] \cdot [1, 2] = 2 + 4 = 6$

3. $[2, 2, 2] \cdot [1, 2, 3] = 12$



clearly, $y > x$

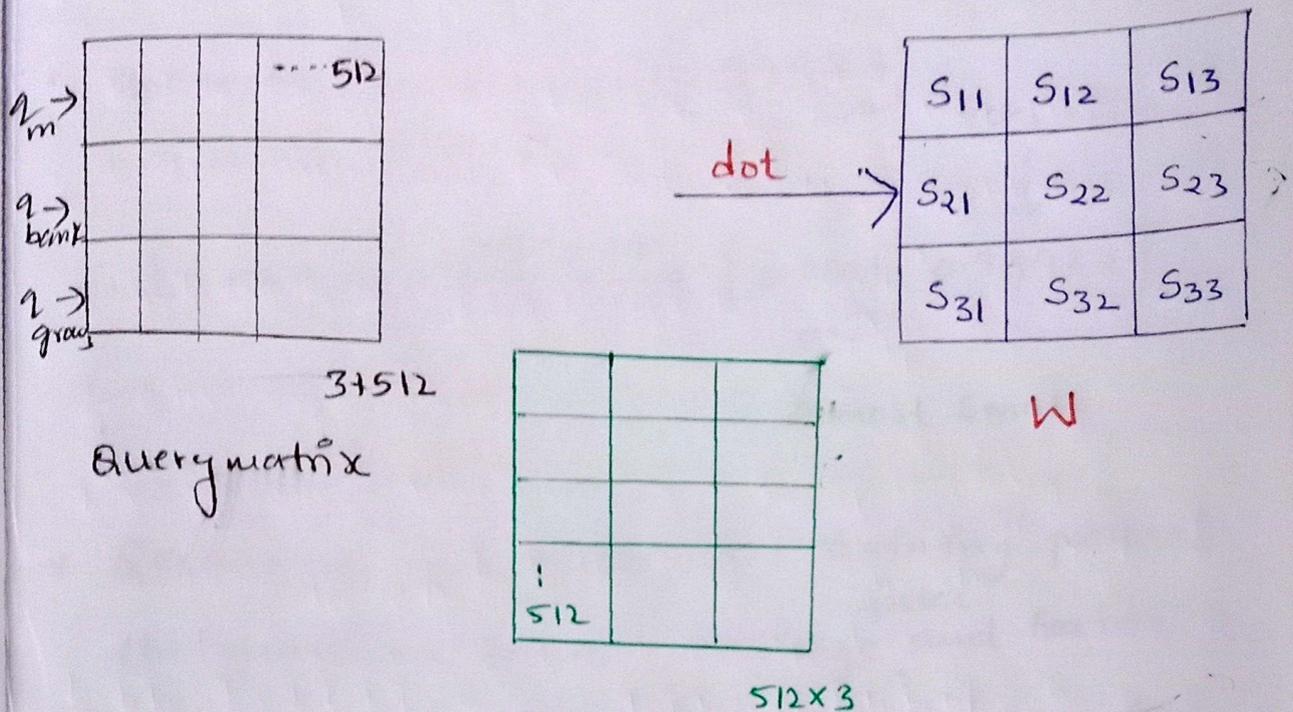
we know Variance is proportional to spread

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

observe $\Rightarrow \dim(y\text{-case}) = 3 > \dim(x\text{-case}) = 2$

* Think **dimension** as no. of elements in a vector

* Let's say the dimension of query and key vectors is 512



* Since over dimension is huge (512), then the values S_{11} ... S_{33} in W has more variance. It means some among them in W have large value and some has low value.

* Now the problem is with softmax due to high variance data.

What is softmax?

input (x) \longrightarrow $e^{\boxed{x}}$ \longrightarrow probabilities

let say we have $[1, 10]$ \rightarrow we have to convert into probabilities [we have to use softmax]

$$1 \longrightarrow \frac{e^1 + e^{10}}{e^1} = \frac{e^1}{e^1 + e^{10}} = 0.0001$$

$$10 \longrightarrow \frac{e^{10}}{e^1 + e^{10}} = 0.99987$$

very high prob

$$\therefore [1, 10] \xrightarrow{\text{Softmax}} [0.0001, 0.99987]$$

↑
Almost small

* Because of this during the training process the gradient will not converge ^{faster} and training becomes very slow! for larger values. And for the smaller values the gradient almost vanishes and smaller values will not update at all!

* Why we are having this problem?

It is because of high variance data! which is the consequence of high dimension (size).

But we can't avoid high dimension because embeddings having high dimension capture

Semantic meaning with other words. Now what is our goal?

Goal: Irrespective of any dimension, variance should be almost similar!

Solⁿ: what I am trying say is:

If it is 60-dimension and variance is 'x'

For 128 also $\approx x$

512 also $\approx x$

How can we achieve this?

The obvious answer is **normalization!**

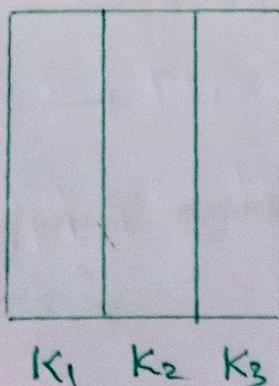
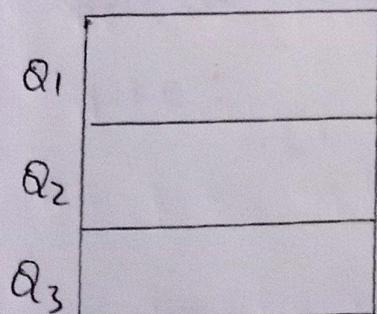
Let's take example:-

our problem with is with W matrix only right which contains high variance data.

S_{11}	S_{12}	S_{13}
S_{21}	S_{22}	S_{23}
S_{31}	S_{32}	S_{33}

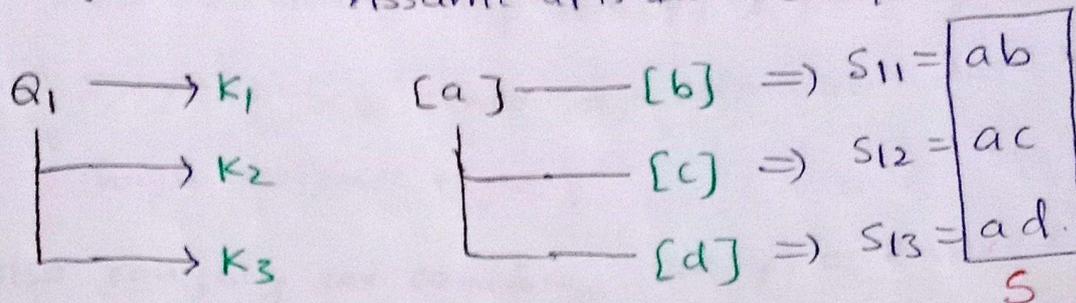
* Let's focus on only 1st row, then we scale up our concept to complete W.

[W]



* Now we are focusing on variance of 1st row of W only.

Assume Q_1 is 1dim & K_i also 1dim



* Think S_{11}, S_{12}, S_{13} are sample and W is population we want $\text{Var}(W)$ (predict) based on sample!

Let assume a, b, c, d are coming from a random variable 'x'. Lets consider the variance of S is $\text{Var}(x)$

* What if the vectors dim = 2?

$$\left. \begin{array}{l} [a \ b] \text{ --- } [c \ d] \Rightarrow S_{11} = ac + bd \\ \text{--- } [e \ f] \Rightarrow S_{12} = ae + bf \\ \text{--- } [g \ h] \Rightarrow S_{13} = ag + bh \end{array} \right\}$$

Here we want for a general vector variance, not

like:

$$\left. \begin{array}{l} [1 \ 1] \text{ --- } [1, 2] \Rightarrow 1+2=3 \\ \text{--- } [1, 3] \Rightarrow 1+3=4 \\ \text{--- } [2, 1] \Rightarrow 2+1=3 \end{array} \right\} \text{ This is Sample Variance}$$

Say $ac+bd, ae+bf$ and $ag+bh$ are coming from random variance 'y'.

* Now variance of 1st row is $\text{Var}(y)$

Is this relationship $\text{Var}(y) > \text{Var}(x)$ is True?

Ans! Is yes! bcz since high dimension data have high variance right!

Also roughly we can say $\text{Var}(y) \approx 2 \text{Var}(x)$

IF 3d?

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

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

* If d-dimension?

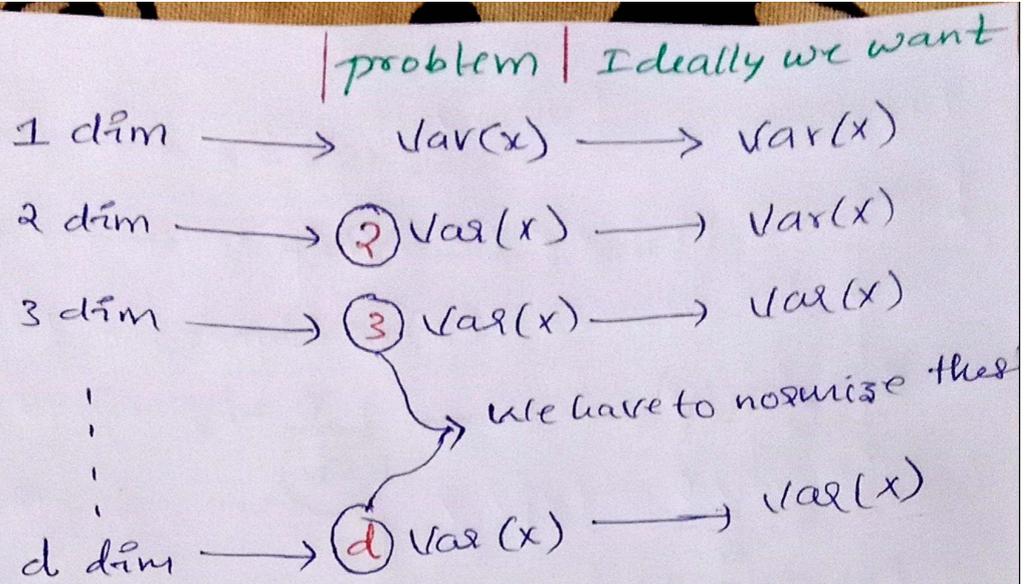
$$\text{Variance} \approx d \frac{\text{Var}(x)}{\text{dim}=1}$$

\therefore We figured out a linear relationship between variance and dimension

Do you agree !! 😊

\therefore We have mathematically quantified the relation: If $\text{dim} \uparrow \text{es} \rightarrow \text{Variance} \uparrow \text{es}$,,

Summary :-



* I hope you understood the problem! 😊

Mathematical Rule:- If you have variable X with a variance of $\text{Var}(x)$, and you create a new variable ' Y ' by scaling ' x ' with a constant ' c ', so that $Y = cX$, the variance:

$$\text{Var}(Y) = c^2 \text{Var}(X)$$

It says...

If you have a random variable ' x ' whose variance is $\text{Var}(x)$.

$$x \longrightarrow \text{Var}(x)$$

You defined a new variable $Y = cX$

$$\text{Then } \text{Var}(Y) \longrightarrow c^2 \text{Var}(x)$$

$$1 \text{ dim} \longrightarrow x = \text{Var}(x)$$

$$\textcircled{2} \text{ dim} \longrightarrow y = \sqrt{2} x \longrightarrow \text{Var}(y) =$$

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

$$y' = \frac{y}{\sqrt{2}} \longrightarrow \text{Var}(y') = \left(\frac{1}{\sqrt{2}}\right)^2 y$$

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

$$= \text{Var}(x)$$

\therefore If we normalize $\frac{y}{\sqrt{2}}$ then variance be
— same $\text{Var}(x)$

$$\textcircled{3} \text{ dim} \longrightarrow \frac{3 \text{Var}(x)}{\sqrt{3}} \longrightarrow \frac{1}{3} (3) \text{Var}(x)$$

$$\textcircled{4} \text{ dim} \longrightarrow \frac{4 \text{Var}(x)}{\sqrt{4}} \longrightarrow \frac{1}{4} (4) \text{Var}(x)$$

⋮

$$\textcircled{d} \text{ dim} \longrightarrow \frac{d \text{Var}(x)}{\sqrt{d}} \longrightarrow \text{Var}(x)$$

\therefore Irrespective of dimension variance
become roughly similar !!

\therefore I hope you understood why we are normalizing
with square-root and if you are clearly

observe the red circles were divideding with dimension itself.

** And that dimension should be dimension of k -vector because we are doing the dot-product with k -vector so we have to take d_k .

* ∴ we updated our first-principle approach by scaling with $\frac{1}{\text{sqrt}(d_k)}$

Let me summarize...

next page

Summary - Self-Attention

Q (3,n)

dot →

S ₁₁	S ₁₂	S ₁₃
S ₂₁	S ₂₂	S ₂₃
S ₃₁	S ₃₂	S ₃₃

(3x3)

K (n,3)

scale $\frac{1}{\sqrt{d_k}}$

↓

W ₁₁	W ₁₂	W ₁₃
W ₂₁	W ₂₂	W ₂₃
W ₃₁	W ₃₂	W ₃₃

(3x3)

← softmax

S' ₁₁	S' ₁₂	S' ₁₃
S' ₂₁	S' ₂₂	S' ₂₃
S' ₃₁	S' ₃₂	S' ₃₃

(3x3)

↓

W ₁₁	W ₁₂	W ₁₃
W ₂₁	W ₂₂	W ₂₃
W ₃₁	W ₃₂	W ₃₃

3x3

dot →

V-matrix 3xn

Y _{money}
Y _{bank}
Y _{grows}

I am Very happy

Congrats!!

