

Benefit of Self Attention

Parallel Training

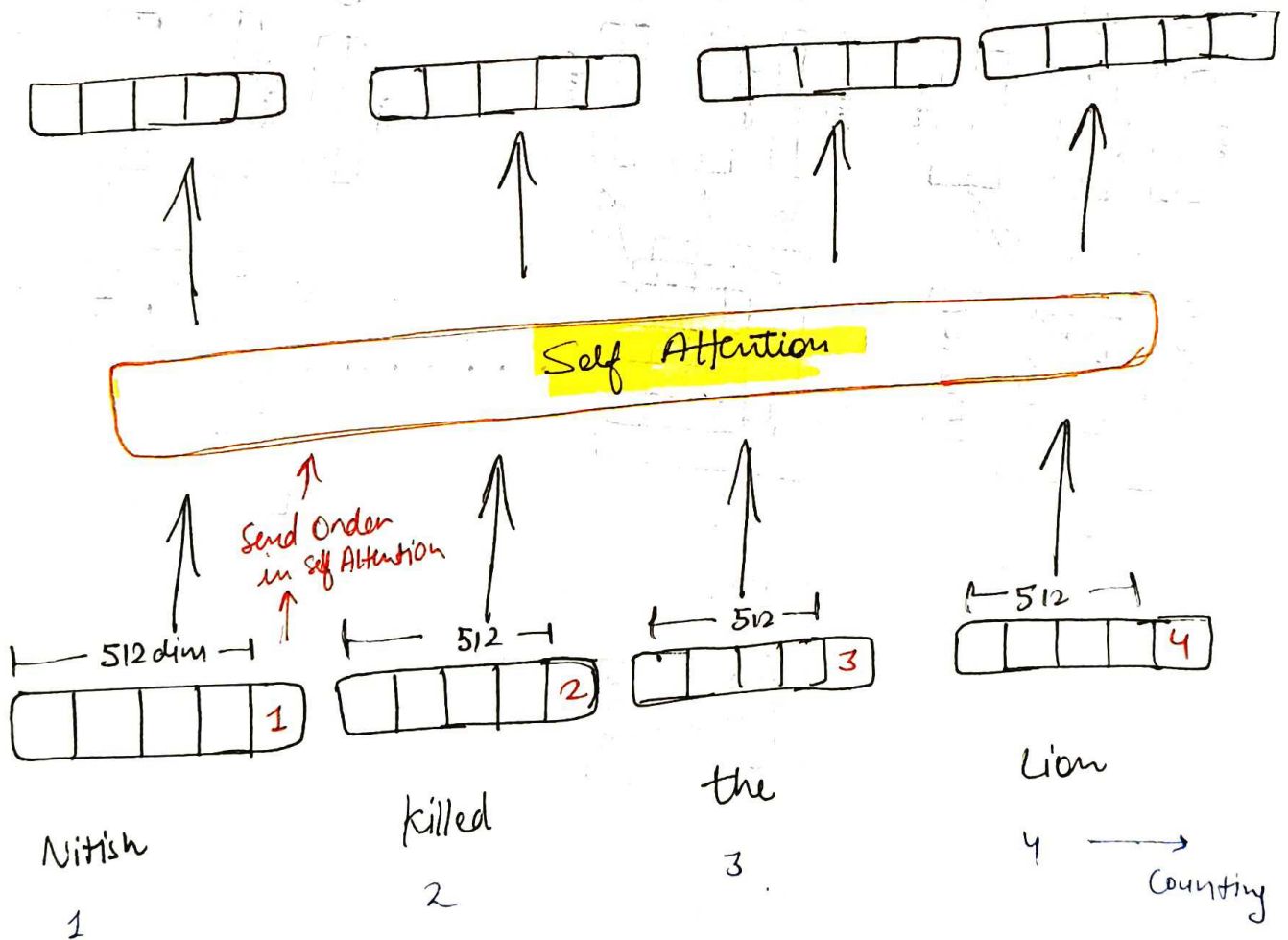
Drawback of Self Attention

word order

eg:- "Nirish killed lion" is same "lion killed Nirish"

Positional Encoding in Transformer

Proposing a Simple Solution



Problem: 1) unbound \rightarrow upper limit

eg:- 1 lakh word document \rightarrow last box is 100000

And we use NN architecture \rightarrow backpropagation

* create instability \leftarrow Hate big numbers \nwarrow

* gradient unstable

Sol. counting Number divide with Total Number
then range convert into 0 to 1.

for: $\frac{1}{100000} \approx 0$, $\frac{2}{100000}$, ... $\frac{100000}{100000} = 1$.

But this solution is not good.

Sentence 1 \rightarrow Thank

You
$\frac{2}{2}$

 $\xrightarrow{\text{position}}$ 2nd value is 1.

Sentence 2 \rightarrow Nitish

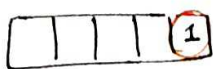
Killed
$\frac{2}{4}$

 then Lien $\frac{3}{4}$ $\frac{4}{4}$
 \rightarrow 2nd position value is 0.5

So, there is no consistency. Position value will be change. NN confused \rightarrow is 2nd position is 1 value or 0.5 value.

* Normalize Not Work

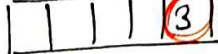
2)



Nitish



killed



the



lion

→ This number is discrete and it is not good for NN.

→ NN generally prefer smooth ^{Transition} ~~Number~~ like continuous Number.

→ Numerical Stability Problem.

3) Can't relative positions

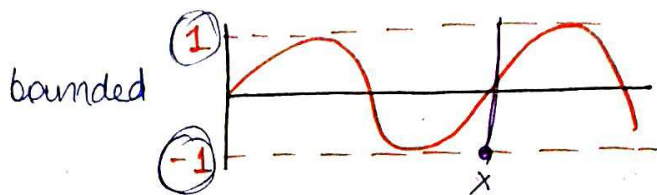


(Can't we find the distance with discrete function. (relative position))
 $(3 - 1 = 2)$

Summary of problem

- unbounded $\xrightarrow{\text{we want}}$ (bounded)
- discrete \longrightarrow (continuous)
- relative position \longrightarrow (periodical)

Sol Use Sin function (trigonometric)

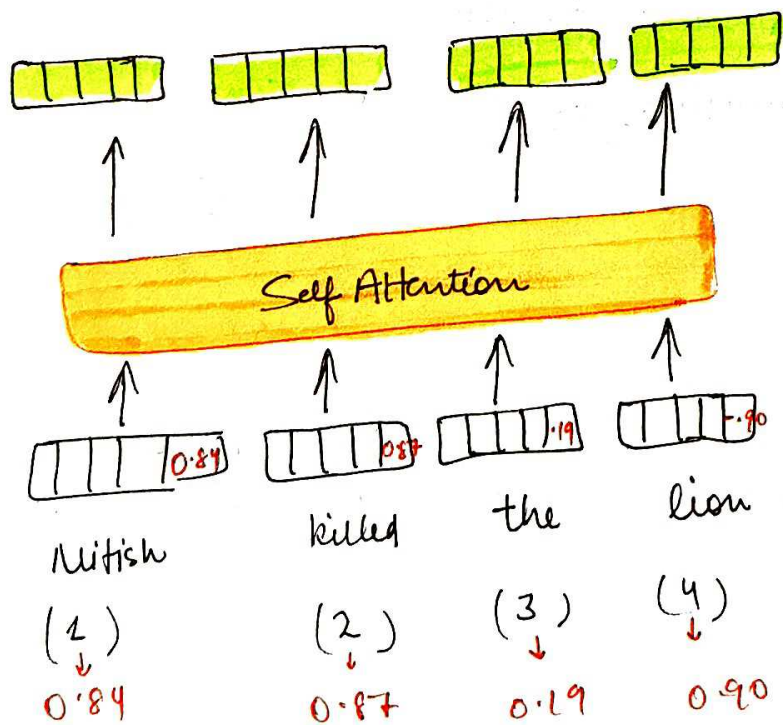
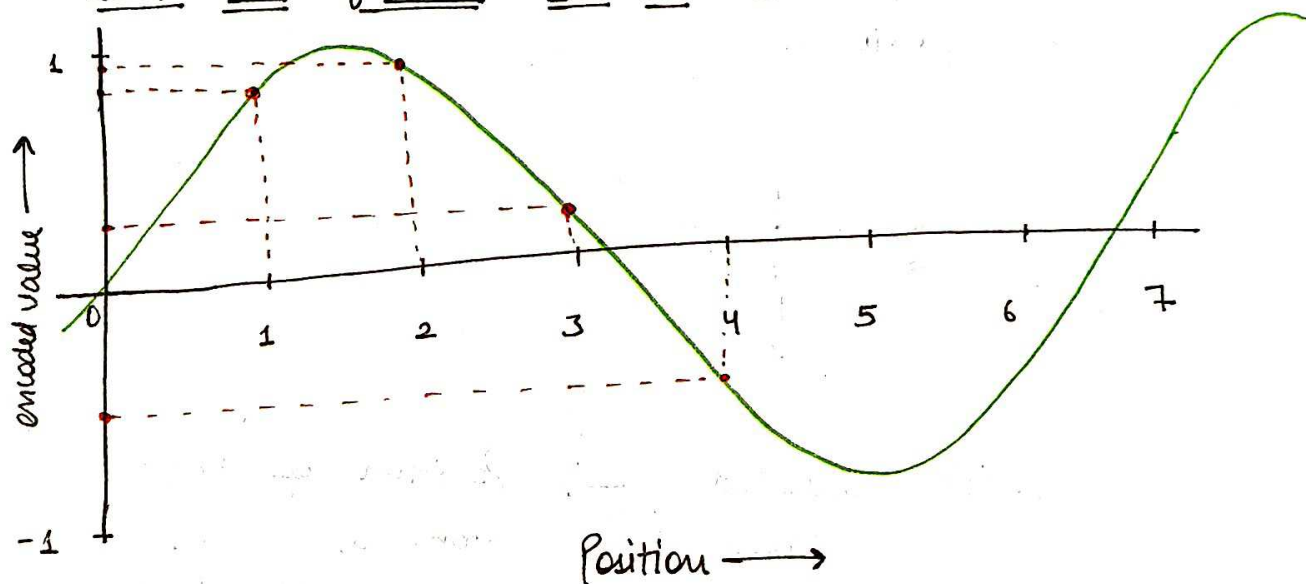


→ continuous → every x has y value

→ value repeats after action

Positional encoder → Sin → better solution.

The sine function as a solution



→ encoded value

$$y = \sin(\text{position})$$

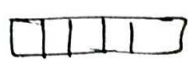
$$y = \sin(1) = 0.84$$

$$y = \sin(2) = 0.87$$

$$y = \sin(3) = 0.19$$

$$y = \sin(4) = -0.90$$

Big Problem with sin function



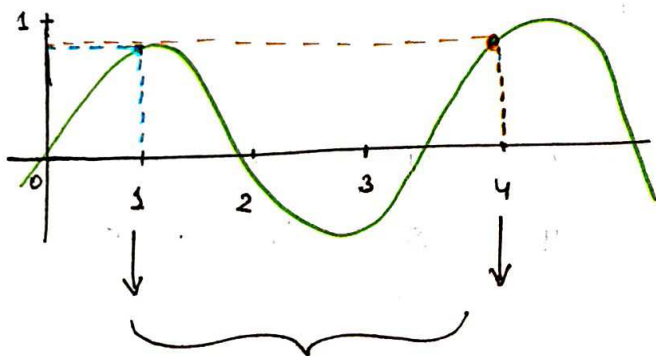
→ encoded value = 0.11

Nitish

* Koi aur word ka positional encoding same (0.11) nhi aana chahiye.

Of positional encoded value is same then position of both word is same. It will create issue.

Sin curve → periodic



* Model will confuse.

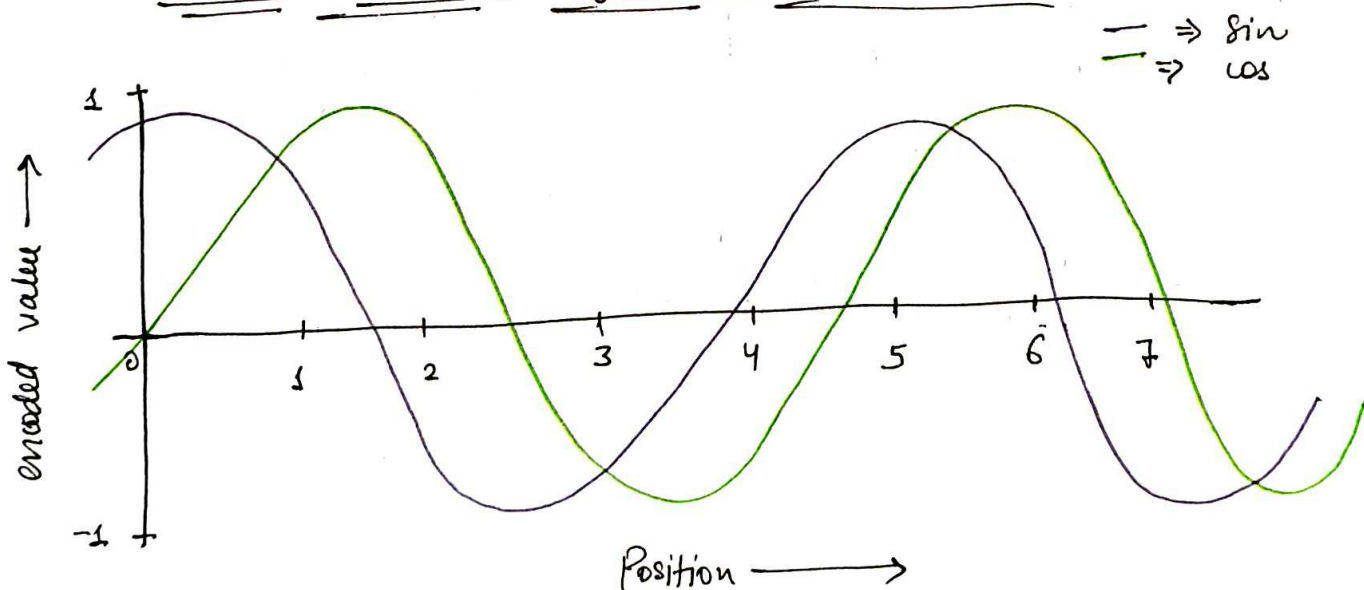


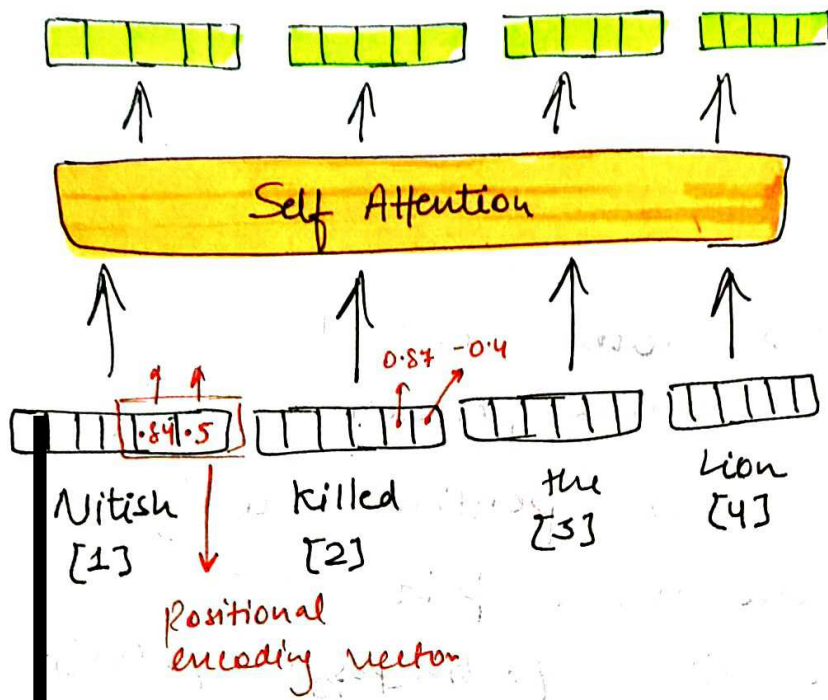
position encoded value is same

→ Position of both word is same according to sin function

Sol. Use double Trigonometric function

The sin and cos function as a solution





$$y_1 = \sin(\text{position})$$

$$y_2 = \cos(\text{position})$$

$$[y_1, y_2] \rightarrow \text{vector}$$

$$\text{Nitish} = [0.84, 0.5]$$

$$\text{killed} = [0.87, -0.4]$$

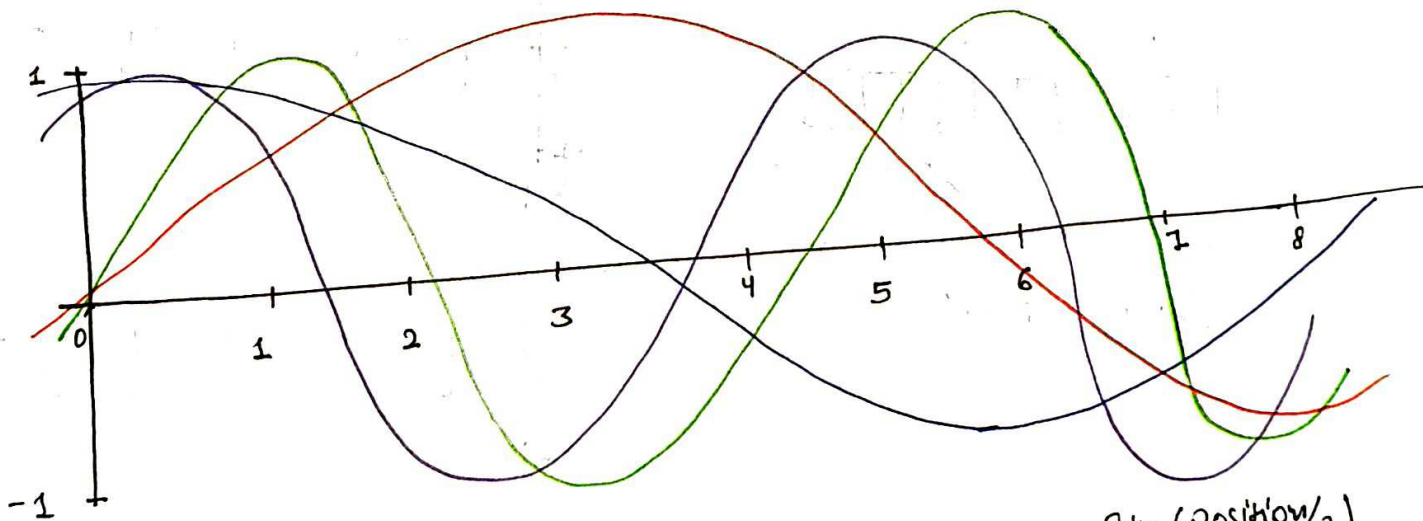
In this sol, also chance to get same vector of two different position word.

So, we use 4 different trigonometric function

$$y^1 = \sin(\text{pos}) \quad y^2 = \cos(\text{pos})$$

$$y^3 = \sin(\text{pos}/2) \quad y^4 = \cos(\text{pos}/2)$$

$$[y^1, y^2, y^3, y^4] \rightarrow \text{Vector}$$



$$\text{green line} = \sin(\text{position})$$

$$\text{black line} = \cos(\text{position})$$

$$\text{red line} = \sin(\text{position}/2)$$

$$\text{blue line} = \cos(\text{position}/2)$$



Nitish

[1]

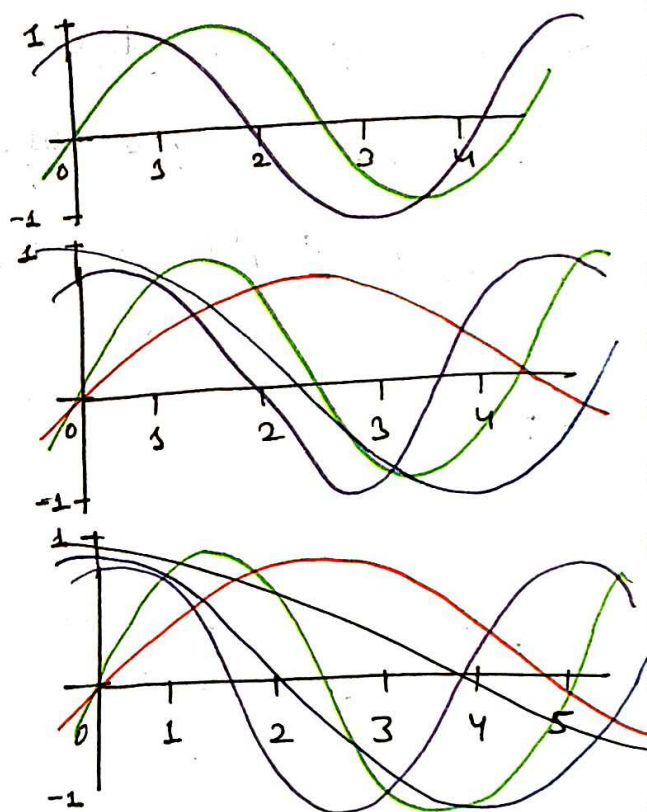
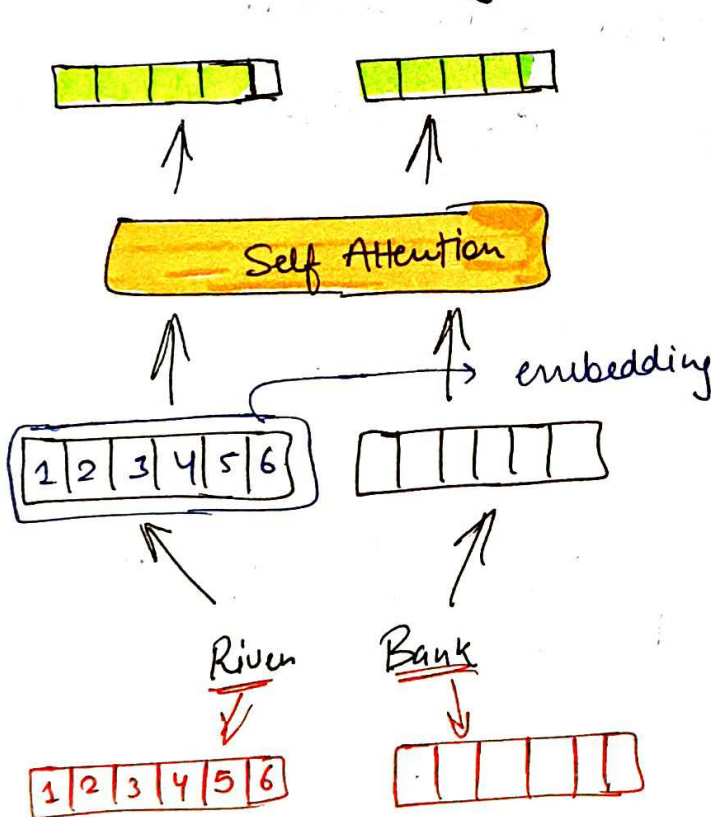
position encoding vector

$[-9 -9 -9 -]$

Very less probability \rightarrow two word have same position encoding

But still if you see same position encoding vector then add one more pair of trigonometric function like $y = \sin(pos/3)$ $y = \cos(pos/3)$

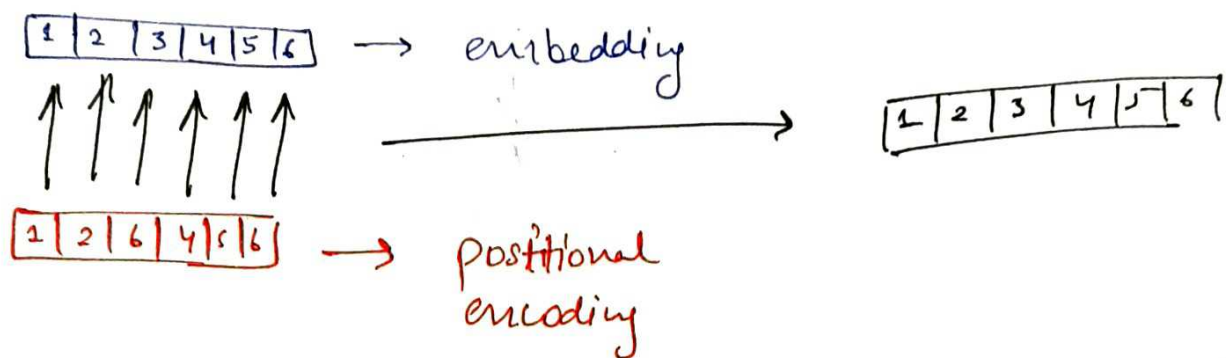
Positional Encoding in "Attention all you need"



* embedding is equal to vector size

\rightarrow Vector (positional encoding)

Embedding (6 dim) + Pos encoding (6 dim) = 6 dim Vector



If concatenate → $[6 \text{ dim}] [6 \text{ dim}] \rightarrow [12 \text{ dim}]$
 Training double up.

* Size of vector increase → frequency of sin and cos function is decrease.

So, how decide which frequency is good?

In the Research Paper → $PE_{(pos, 2i)} = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$

$$PE_{(pos, 2i+1)} = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

pos → Position of word

d_{model} → Dimensionality of embedding

$i \rightarrow [0] \text{ to } \left[\frac{d_{model}-1}{2}\right] \Rightarrow i=1, i=2, i=3 \dots$

Pos = 0

River

1 2 3 4 5 6



Pos = 1

Bank



for $i = 0$ → pos = 0 (River)

$$PE(0, 0) = \sin(0 / 10000^\circ) = 0$$

$$PE(0, 1) = \sin(0 / 10000^\circ) = 1$$

(2i+1)

for $i = 1$ → pos = 0 (River)

$$PE(0, 2) = \sin(0 / 10000^{1/3}) = 0$$

$$PE(0, 3) = \sin(0 / 10000^{1/3}) = 1$$

for $i = 2$ → pos = 0 (River)

$$PE(0, 4) = \sin(0 / 10000^{2/3}) = 0$$

$$PE(0, 5) = \cos(0 / 10000^{2/3}) = 1$$

Riven →

0	1	0	1	0	1
---	---	---	---	---	---

Bank word → Positional Encoding

for i = 0 → pos = 1 (Bank)

$$PE(1, 0) = \sin(1/10000^\circ) = \boxed{0.84}$$

$$PE(1, 1) = \cos(1/10000^\circ) = \boxed{0.54}$$

for i = 1 → pos = 1 (Bank)

$$PE(1, 2) = \sin(1/10000^{1/3}) = \boxed{0.04}$$

$$PE(1, 3) = \cos(1/10000^{1/3}) = \boxed{0.99}$$

for i = 2 → pos = 1 (Bank)

$$PE(1, 4) = \sin(1/10000^{2/3}) = 0.00$$

$$PE(1, 5) = \cos(1/10000^{2/3}) = 0.99$$

Bank →

0.84	0.54	0.04	0.99	0.00	0.99
------	------	------	------	------	------

↪ Positional Encoding vectors