

① Course 5 :- Week 1

PAGE: 11
DATE: ..

Examples of Sequence data

- (1) Speech Recognition
- (2) Music generation
- (3) Sentiment classification

DNA analysis

Machine translation

Video Activity Recognition

Name Entity recognition.

Another Example:

Name Recognition in a sentence

Ex: Harry Potter and Hermione Granger invented a new spell

$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad \dots \quad x^{<9>}$

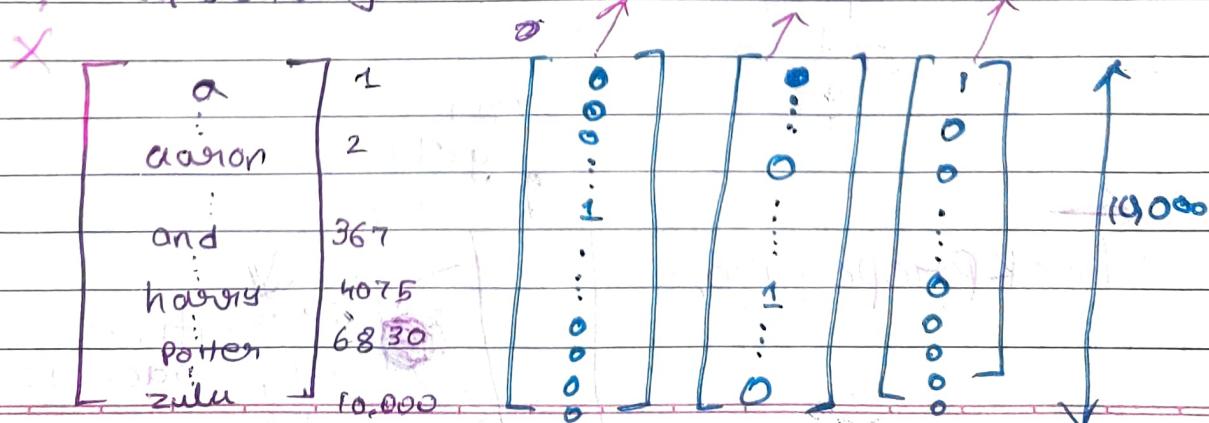
$Tx = 9$ bcz 9 words.

y

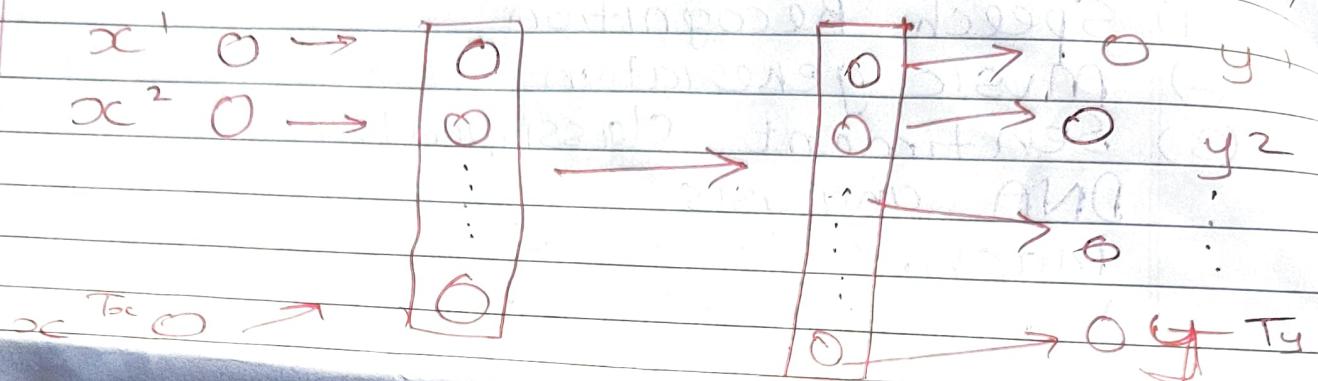
$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad \dots \quad y^{<9>}$

Representing words

Vocabulary Ex:- Harry Potter a



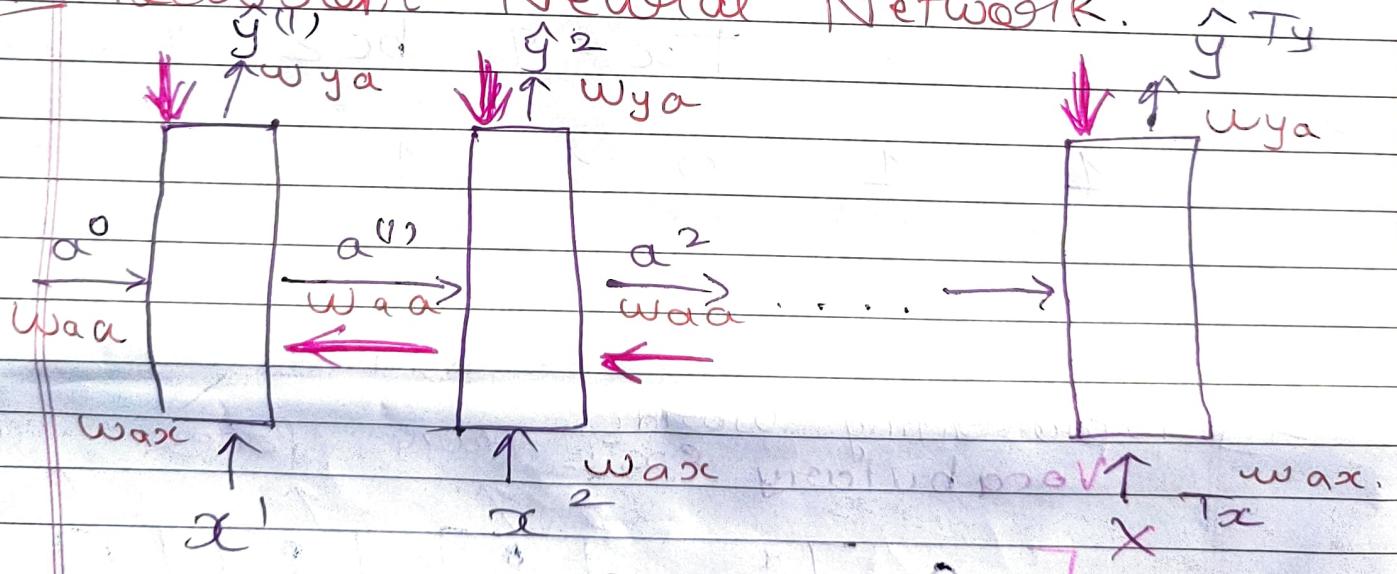
Why don't we use a standard neural network.



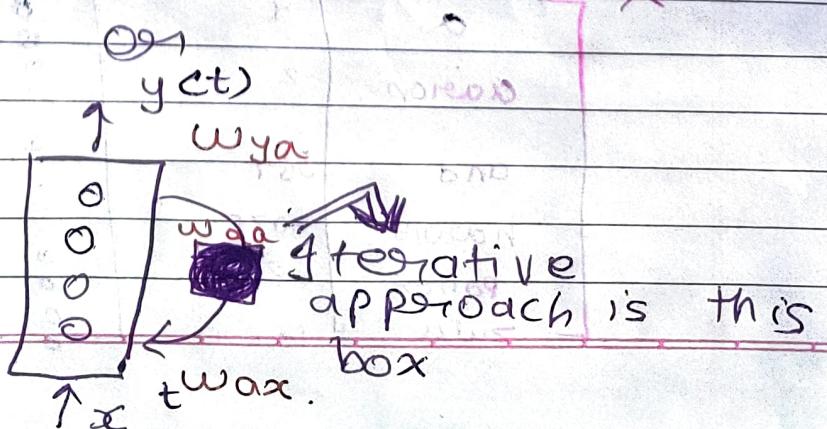
Problems:-

- (1) Inputs and outputs can be different length in different examples
- (2) Doesn't share features learned across different position of text

Recurrent Neural Network.



backpropagation.



In above diagram

tanh/ Relu

$$a^0 = \vec{0} \quad a' = g(w_{aa}a^0 + w_{ax}x' + b_a)$$

$$\hat{y}' = g(w_{ya}a' + b_y)$$

↑
Sigmoid.

General Notation.

$$a^t = g(w_{aa}a^{t-1} + w_{ax}x^t + b_a)$$

$$\hat{y}^t = g(w_{ya}a^t + b_y)$$

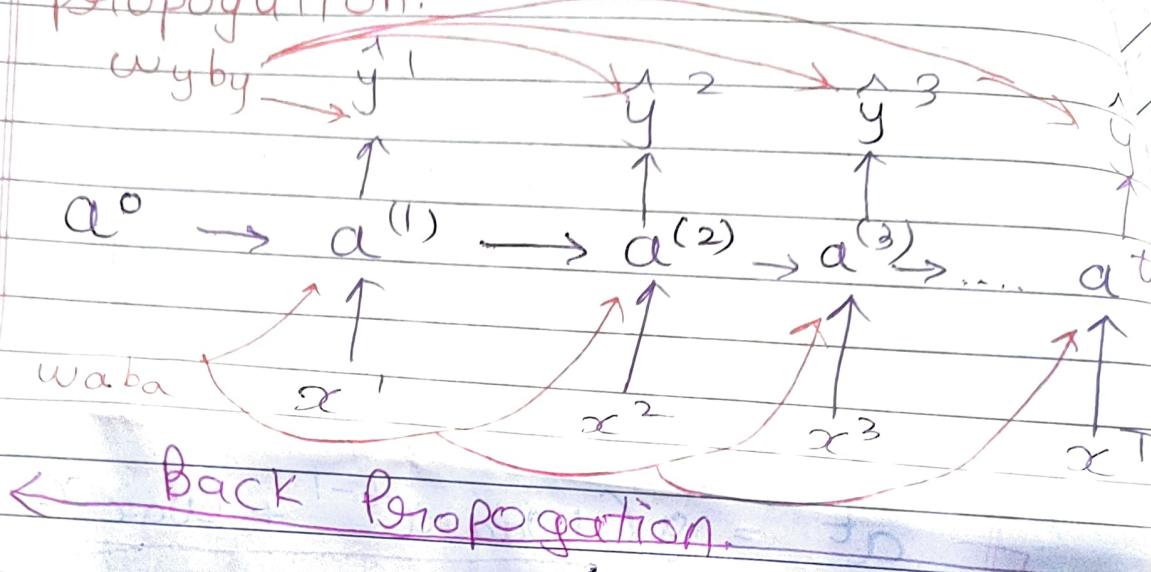
$$a^{<t>} = g(w_a[a^{t-1}, x^t] + b_a)$$

$$W_a = \begin{bmatrix} W_{aa} & W_{ax} \end{bmatrix} \quad [a^{t-1}, x^t] = \begin{bmatrix} a^{t-1} \\ x^t \end{bmatrix}$$

Dimensions: $W_{aa} \times 100$, $W_{ax} \times 100$, $a^{t-1} \times 100$, $x^t \times 100$

$(10100, 1)$

Forward Propagation and back propagation.



$$\delta(\hat{y}, y) = -y(t) \log \hat{y}^t - (1-y(t)) \log(1-\hat{y}^t)$$

$$= \sum_{t=1}^{T_x} \delta^t(\hat{y}^t; y^t)$$

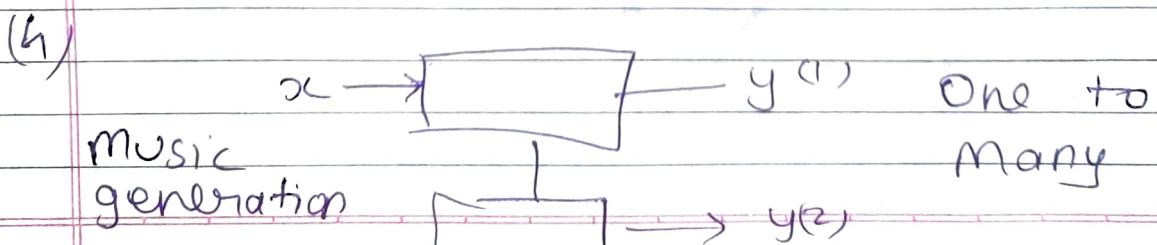
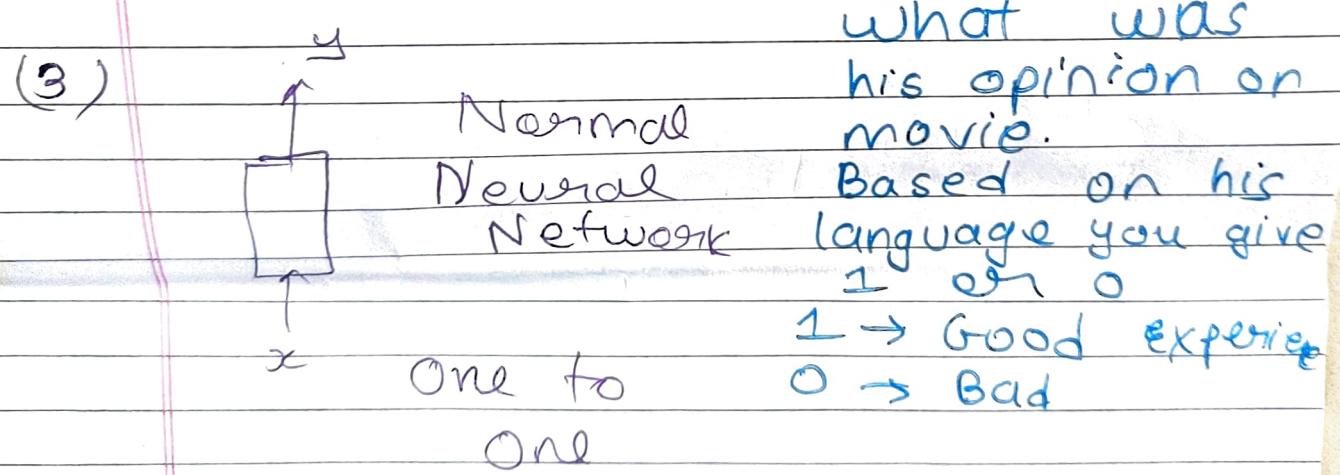
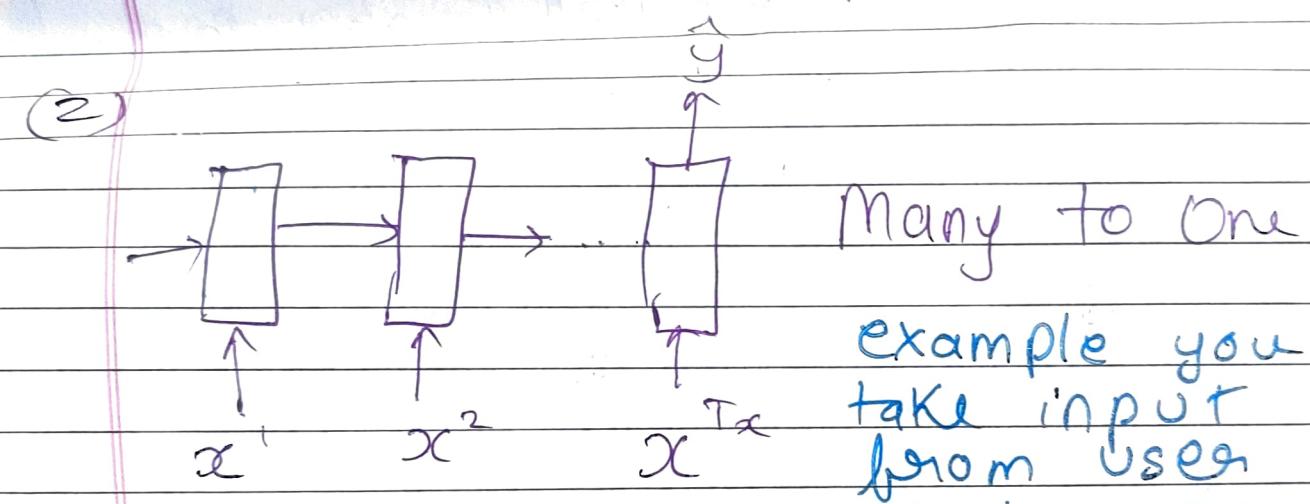
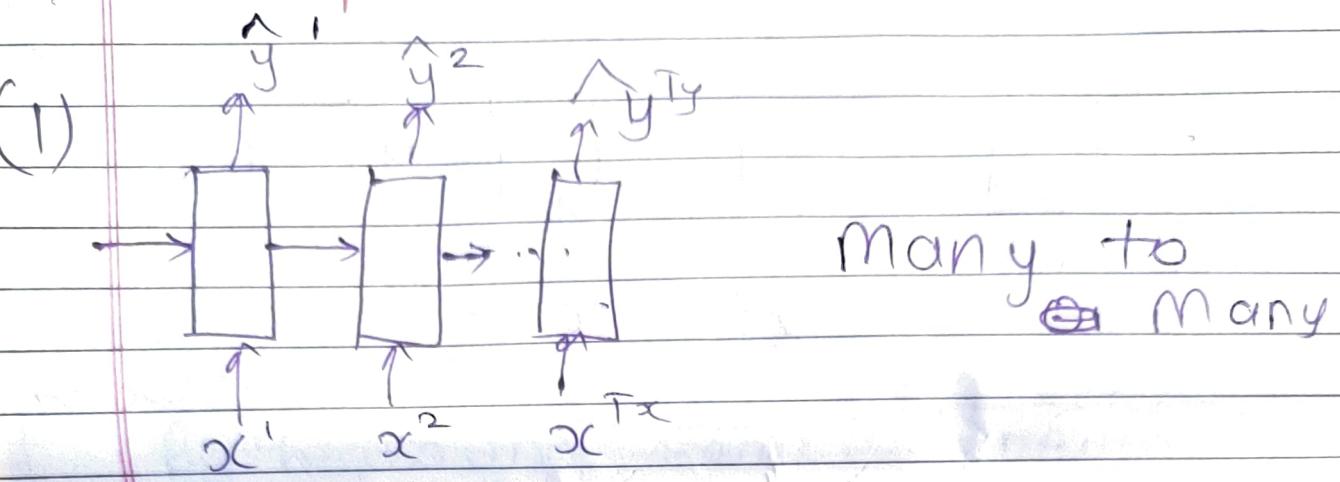
Note:- It's not necessary that T_x and T_y are similar in length or data type.

T_x can be voice.

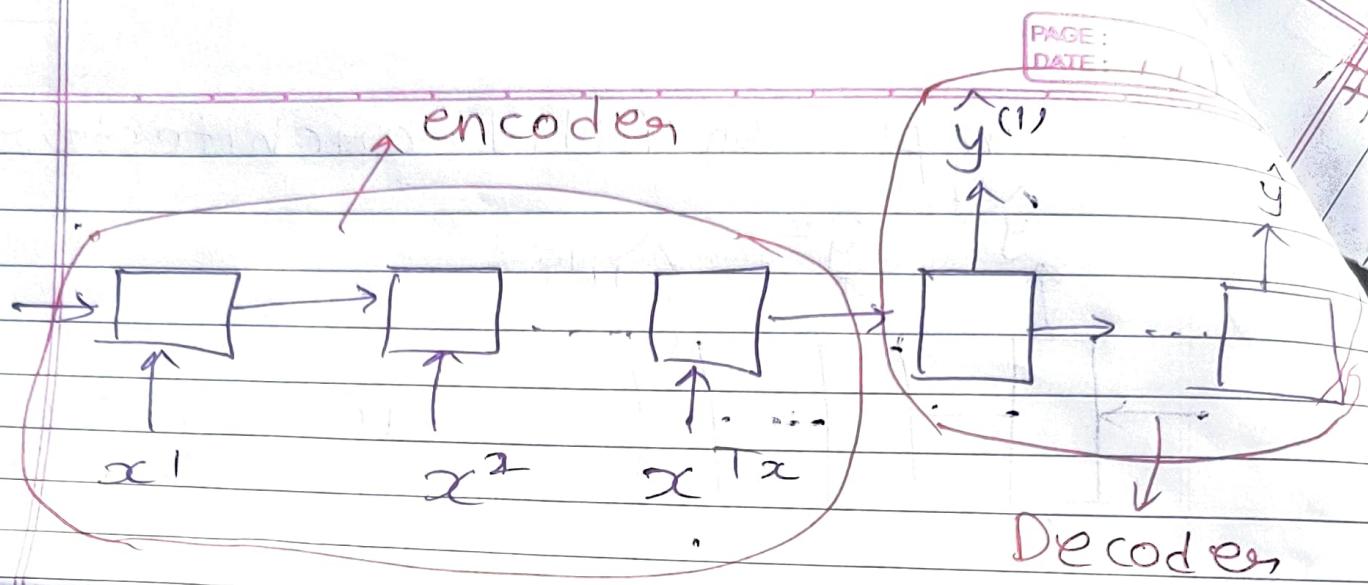
T_y can be speech recognition.

Let's see on Next Page

Example of RNN architecture.



(5) x^t can be \emptyset



This is popularly used in
Machine Translation

What is Language Modelling

(1) Ex :- In speech Recognition,

1st sent The Apple and pair salad

2nd sent The apple and pear salad.

Both sound same however
2nd sentence is preferred bcz it
has more probability and
more logical sense in it.

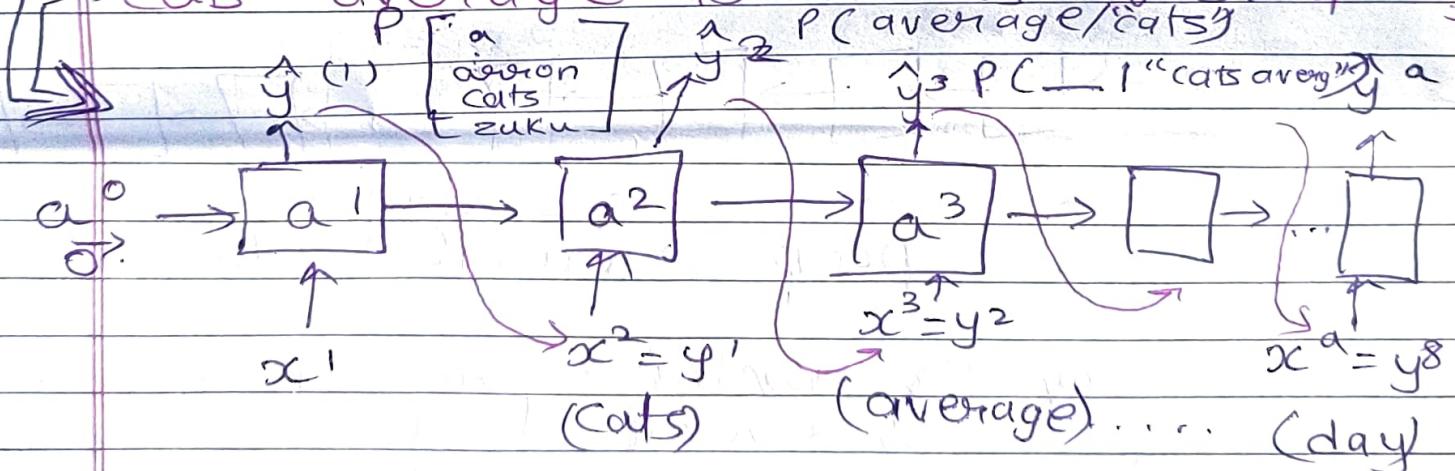
Language modelling mainly deals
with Probability of words in
a sentence.

Ex:- $P(1^{\text{st}} \text{ sentence}) = 3.2 \times 10^{-10}$

$$P(2^{\text{nd}} \text{ sentence}) = 5.7 \times 10^{-10}$$

This sentence is first broken into tokens

Cats average is hours of sleep a



$$P(y_1^{(1)}) \rightarrow P(y^2 | y^1) \rightarrow P(y^3 | y^1 y^2)$$

$$\dots \leftarrow P(y^4 | y^1 y^2 y^3)$$

$\langle \text{UNK} \rangle$ = Unknown token !!

$\langle \text{EOS} \rangle$ = end of Sentence

$$L(\hat{y}^t, y^t) = -\sum_i y_i^t \log \hat{y}_i^t$$

$$L = \sum_t L^t(\hat{y}^t, y^t)$$

$y^{(1)} = \text{cats}$

$y^{(2)} = \text{cats Average}$

$y^{(3)} = \text{Average is}$

$y_8 = \text{sleepy day}$ Cats Average is hours of

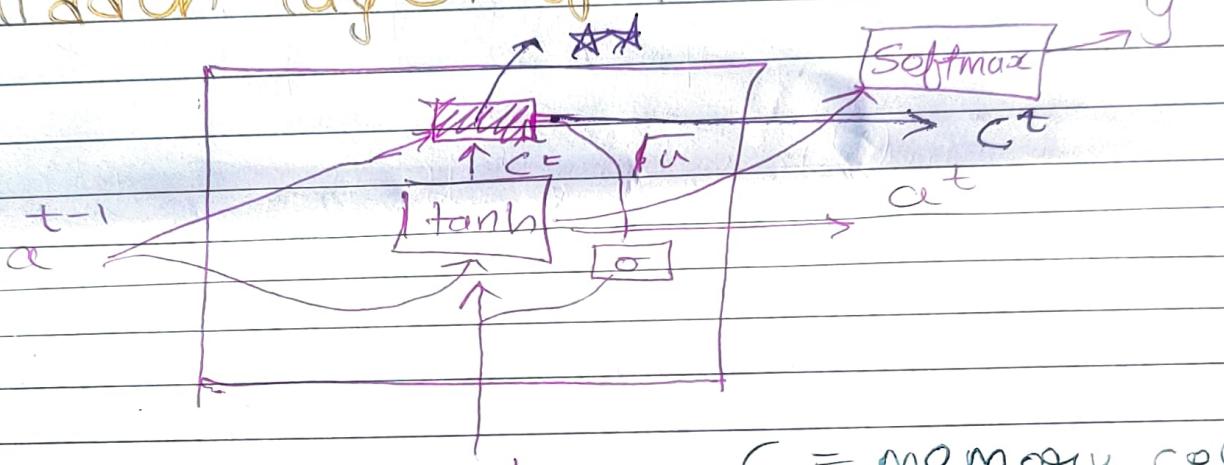
Problem with RNNs is vanishing gradient problem



when the sentence proceeds further it is observed that new characters produced are less dependent and often overloop the meaning of words used previously thereby creating errors in sentence grammatical or logical

when we train Deep Neural Network we often develop vanishing and exploding gradient problem

Hidden layer of RNN Unit.



c = memory cell

Gated Recurrent Unit which remembers that a_t is singular or plural so gives like that which are grammatically correct.

$$c_t = a_t$$

$$\tilde{c}^t = \tanh(w_c [c^{t-1}, x^t] + b_c)$$

$$\Gamma_u = \sigma(w_u [c^{t-1}, x^t] + b_u)$$

Gamma

u = for update

$$\star \text{ is } c^t = \Gamma_u * \tilde{c}^t + (1 - \Gamma_u) * c^{t-1}$$

$$\{(1+1)+2\}[3]\{(0+1)*2\}+3\}$$

Course 5:-

Week 2:-

WORLD STAR™

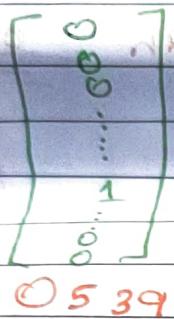
Date : _____

Page : _____

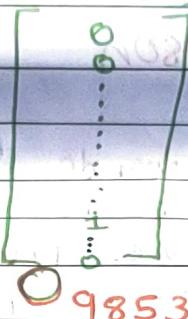
Word representation.

1-hot representation

Man
(5391)



Woman
(9853)



Ex :-

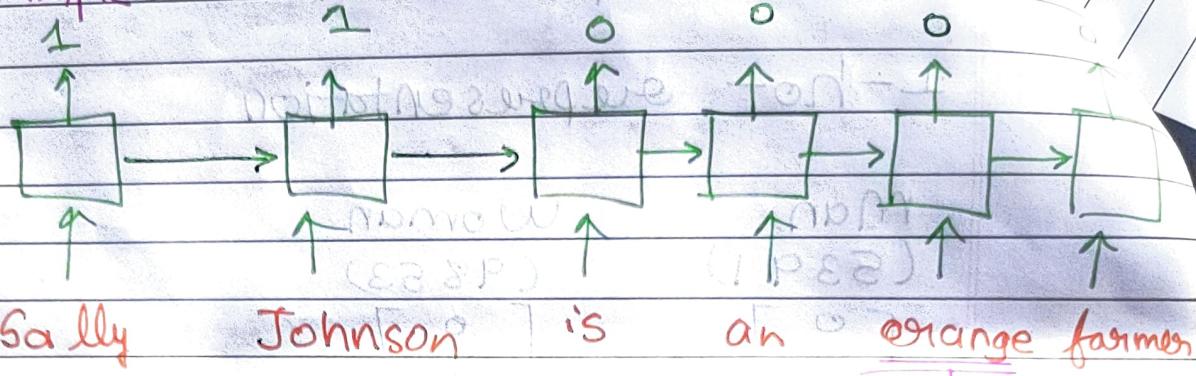
	Man	Woman	King	Queen	Apple	Orange
Gender	5391	9853	4916	7157	456	0
Royal	0.01	0.02	0.93	0.95	-0.01	0.0
Age	0.03	0.02	0.7	0.69	0.03	-0.0
Food	0.09	0.01	0.02	0.01	0.95	0.97
Parameters	300	300	300	300	300	300
	e ₅₃₉₁	e ₉₈₅₃				

Visual word Embeddings

man	woman	dog
woman	man	cat
dog	cat	fish
king	queen	grape
queen	king	apple
apple	grape	orange
orange	apple	grape

Example of Named entity recognition.

example



1B to 100B words
of data set to
train above model.

Based on this sentence we
can understand that orange can be
replaced by another fruit.

Learning and word Embeddings.

(1) Learning words from large text based
corpus (1-100B) words.

→ (2) Transfer embedding to new task with
smaller training set ex. look words

→ We Reduce the words set and other
things even more based on need and
also

ANALOGIES

Man and woman
King and ?

(1) Analogy using matrix

$$\text{eman} - \text{ewoman} \approx$$

-2	→ Gender
0	→ Royal
0	→ Age
0	→ Food

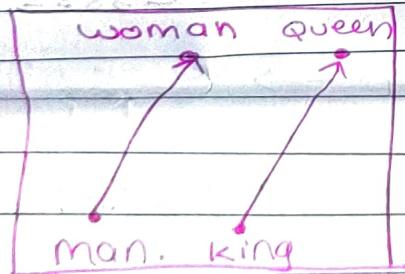
$$\text{eKing} - \text{eQueen} \approx$$

-2	→ Gender
0	→ Royal
0	→ Age
6	→ Food

$$\therefore \text{eman} - \text{ewoman} \approx \text{eKing} - \text{e}?$$

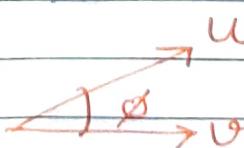
∴ ? = Queen based on prediction

(2) Analogy using word vectors.

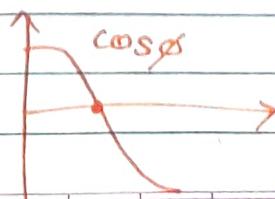


$$\text{Sim}(\text{e}_w, \text{e}_{\text{King}} - \text{eman} + \text{e}_w)$$

$$\text{Sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

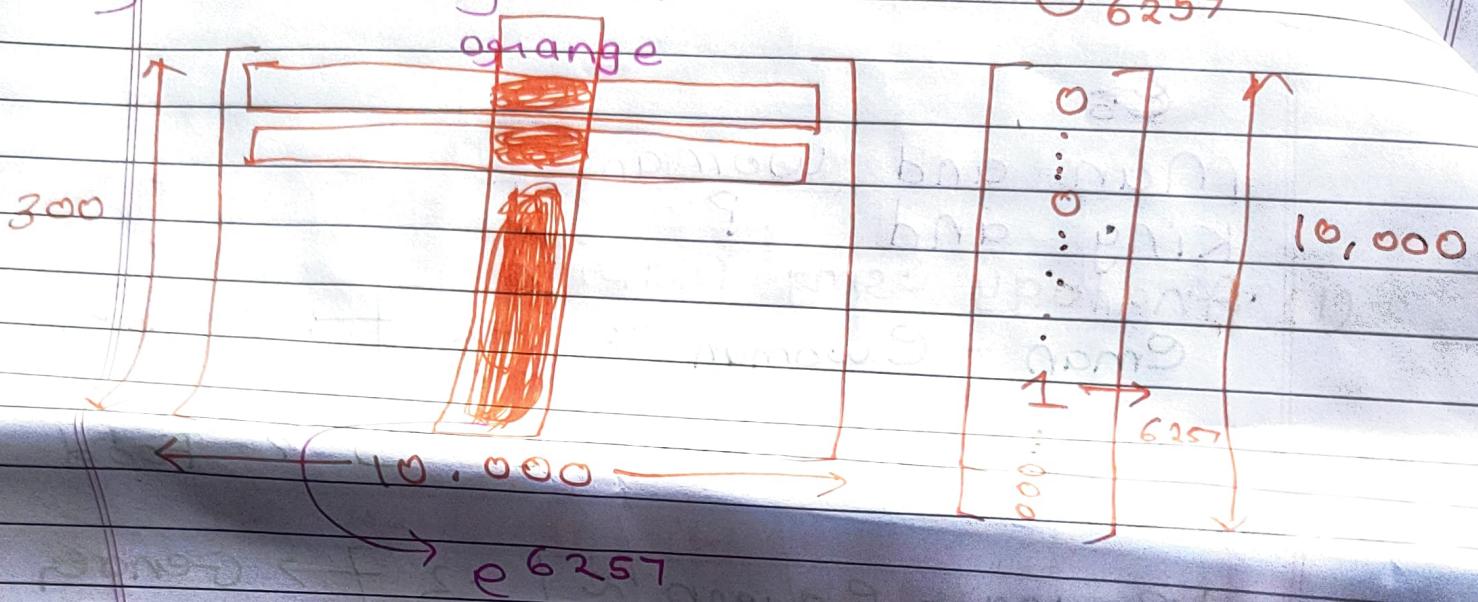


$$\|u - v\|^2$$



Embedding matrix

① 6257



Neutral Language modelling.

g

o₄₃₄₃ → E → e₄₃₄₃want b₉₆₆₅ → E → e₉₆₆₅

a

o₁ → E → e₁

glass

o₃₈₅₂ → E → e₃₈₅₂

ef

o₆₁₆₃ → E → e₆₁₆₃

Orange

o₆₂₅₇ → E → e₆₂₅₇

Juice

e₆₂₅₇ ←

softmax

softmax

#] Other contexts / Target Pairs

I want a glass of orange juice
to go along with my cereal.

Context last 4 words

a glass of orange — to go alg
with.

last 1 word.

orange —

#] Model

Vocab size = 10000

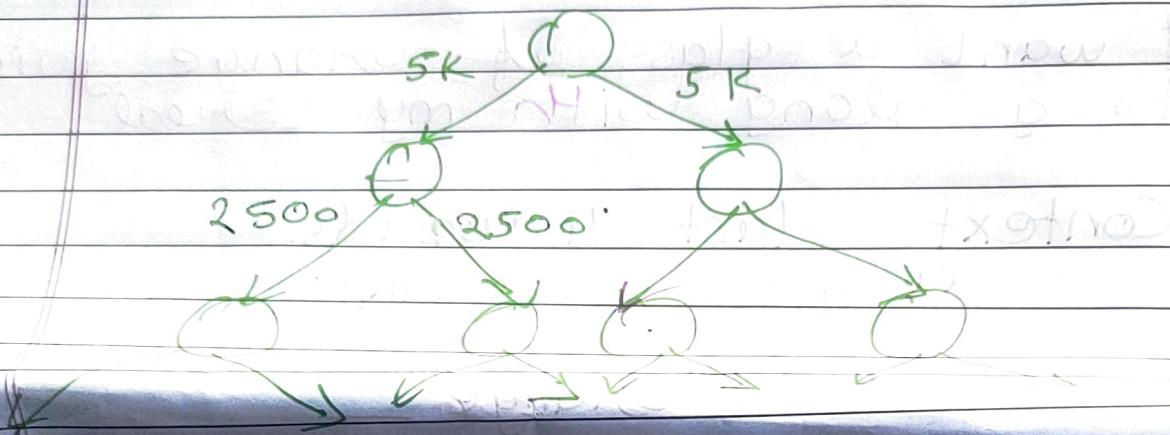
$x_c \rightarrow y$
context target
(c)

orange glass
orange juice
orange my

$$\text{Salmon} \quad P(t|c) = \frac{\sum_{j=1}^{10000} e^{\theta_t^T \mathbf{x}_c}}{\sum_{j=1}^{10000} e^{\theta_j^T \mathbf{x}_c}}$$

$$L(\hat{y}, y) = -\sum_{i=1}^{10000} y_i \log \hat{y}_i$$

Hierarchical softmax



How to sample the context 'c'
 → The, of, a, and to, ...
 → orange, apple, division, book

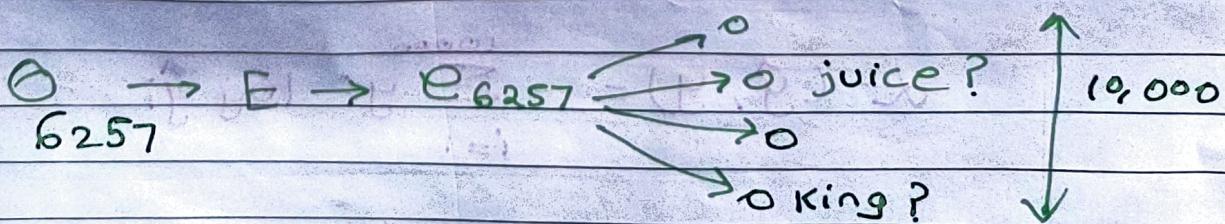
Negative Sampling

Context	word	target
orange	juice	target no 1
orange	King	(2) 0
orange	book	approx 0
orange	in the	approx 0

↑ ↑ ↑

C t y

$$P(y=1 | c, t) = \sigma(\theta_t^T e_c)$$



GLOVe (Global Vectors for Word Representation).

$x_{ij} = i$ appears in context of j

$\begin{matrix} \uparrow & \uparrow & \uparrow \\ c & t & t \\ \end{matrix}$
 \uparrow
 c

$$x_{ij} = x_{ji}$$

$$\text{minimize } \sum_{i=1}^{10000} \sum_{j=1}^{10000} f(x_{ij}) \cdot (\Theta^T e_j + b_i + b_j - \log x_{ij})^2$$

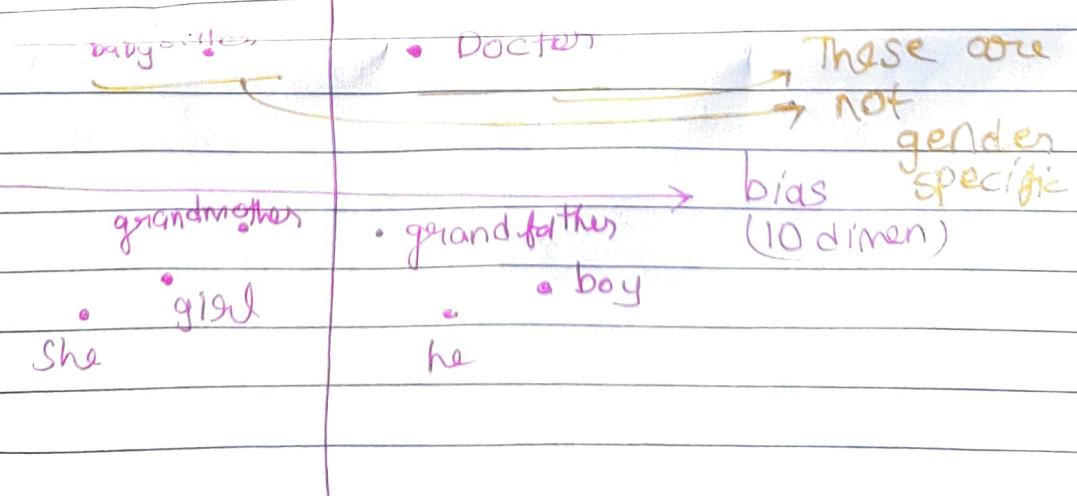
Problem of bias in word embeddings

Man x woman as King x Queen ✓

Father x Doctor as Mother x Nurse \Rightarrow ✗

↓
This problem is arising here
To avoid this we use

non bias direction θ_{99D}



① Course 5 Week 3:-

WORLDSTAR™

Date :

Page :

Why not Greedy Search

Greedy Search means printing the words in a line which are most connected to each other i.e. with highest probability however this is not at all useful bcz we need to make sentence based on conditions, constraints and input.

Beam search

Length Normalization

$$\arg \max_y \prod_{t=1}^{T_y} P(y^t | x, y^1, \dots, y^{t-1})$$

$$\arg \max_y \sum_{t=1}^{T_y} \log P(y^t | x, y^1, \dots, y^{t-1})$$

$$\frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^t | x, y^1, \dots, y^{t-1})$$

Beam Width B

large B: better result, slower

smaller B: worse result, faster.

Beam search runs faster but is not guaranteed to find exact

minimum for arg max $P(y|x)$.

Error Analysis Process.

Human

Jane visits
Africa in
September

Algorithm $P(y^*(x)) P(\hat{y}|x)$

Jane visited
Africa last
September

2×10^{-10} 2×10^{-10}
Beam
at
fault
 $\therefore y^* > \hat{y}$

If $y^* < \hat{y}$
then RNN at
fault.

Attention Model Intuition.

COURSE 5 WEEK 4

In RNNs GRUs, LSTMs the complexity of Transformer model keeps on increasing.

Hence we developed Transformer which works on Attention + CNN

- ↳ Self - Attention
- ↳ Multi - Head Attention

$A(Q, K, V) = \text{Attention based Vector}$
 ↳ representation of a word
 ↳ calculate for each word.

$A^1, A^2, \dots, A^n \Rightarrow n$ words

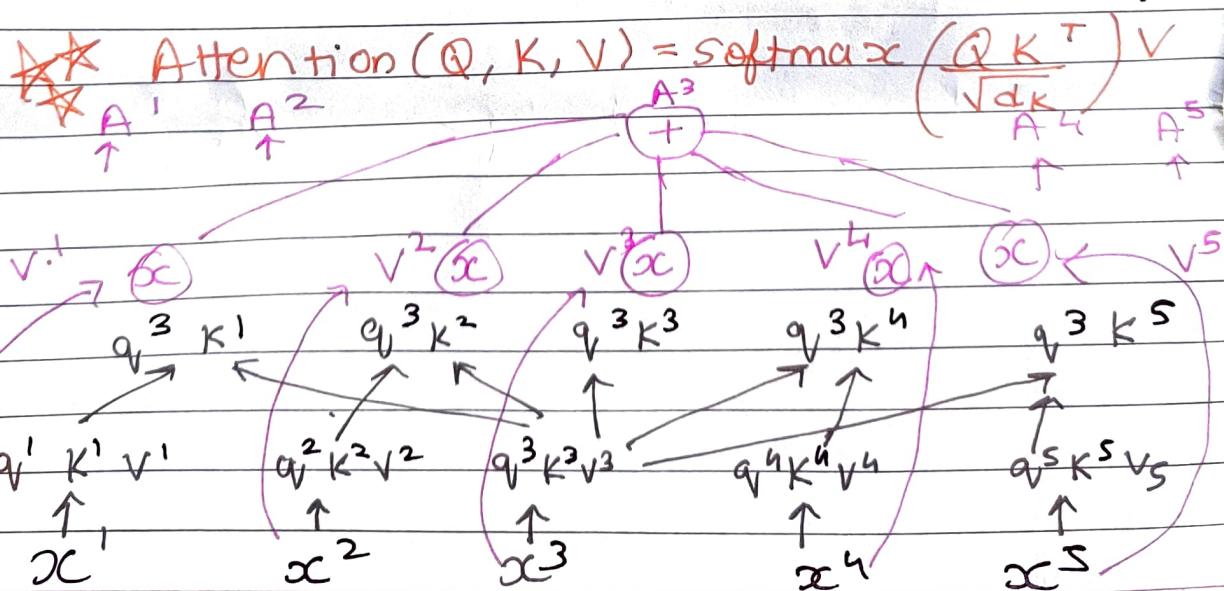
have 1 Attention each

RNN Attention.

$$\alpha^{t, t'} = \frac{\exp(e^{t, t'})}{\sum_{t'=1}^T \exp(e^{t, t'})}$$

Transformer's Attention

$$A(Q, K, V) = \underbrace{\sum_i \exp(q_i K_i)}_{\sum_j \exp(q_j K_j)} \times v_i$$



Jane visite l' Afrique en septembre

Attention \rightarrow queries, Key and Values

\hookrightarrow queries	q_1, q_2, \dots, q_T	$q_i \in \mathbb{R}^d$
\hookrightarrow Keys	k_1, k_2, \dots, k_T	$k_i \in \mathbb{R}^d$
\hookrightarrow values	v_1, v_2, \dots, v_T	$v_i \in \mathbb{R}^d$

Vectors in a given dimension.

Number of Keys, queries can differ from the number of keys and values in practice.

In self-Attention q, k, v are drawn from same source-

The output of previous layer is x_1, \dots, x_T
(One vector per word)

we can let $q_i = k_i = v_i = x_i$ use same
vector for query, key and value [Can we use same
vector for query, key and value then?]

$\checkmark e_{ij} = q_i^T k_j$
 $T \times T$ matrix with scalar values not
compute key - bounded by size
query affinities

$$\checkmark a_{ij} = \frac{\exp(e_{ij})}{\sum_j \exp(e_{ij})}$$

Normalisation

constant [Compute attention weights from
(softmax) affinities]

$\text{Output}_i = \sum_j a_{ij} v_j$ going to weigh
 Compute outputs as weighted sum of
 values.

Attention treats each words representation
 as a query to access and incorporate
 information from set of values.

* Why cannot we have self-attention
 itself as a building block i.e. a
 replacement of RNNs and LSTM why
 need of Transformer - ~~for n~~

↓
 we don't use bcz the self Attention
 model does not accept input in
 particular order like

The Chef makes food
 and the food Chef makes
 is the same for it even though
 there is lot of difference. it does
 not give incidies to q_i, k_i, v_i
 no importance of value of a_{ij}

To solve above problem we use
 $p_i \in \mathbb{R}^d$ for $i \in 1, 2, \dots, T$
 → position vectors

• $\tilde{q}_i, \tilde{v}_i, \tilde{k}_i$ are our old values
 $q_i = \tilde{q}_i + p_i$ }
 $v_i = \tilde{v}_i + p_i$ }
 $k_i = \tilde{k}_i + p_i$ } you can do
 many other things.

* Position representation done through sinusoidal waveforms.

$$\begin{aligned} \mathbf{P}_i &= \begin{pmatrix} \sin(i/10000^{2\pi/\lambda}) \\ \cos(i/10000^{2\pi/\lambda}) \end{pmatrix} \\ &\quad \text{dimension: } (d) \quad \text{index in: } (l) \\ &\quad \text{search: } (r) \\ &\quad \boxed{\sin(i/10000^{(3*\frac{2}{2})/\lambda})} \\ &\quad \boxed{\cos(i/10000^{(2*\frac{2}{2})/\lambda})} \end{aligned}$$



Pros :-

- (1) Periodicity indicates that maybe absolute position isn't as important
- (2) Can extrapolate to longer sequences as period starts.

Cons :-

- (1) Not learnable and the extrapolation doesn't really work.

We have modified \mathbf{P}_i to make it learnable

→ $\mathbf{P} \in \mathbb{R}^{d \times T}$ let each \mathbf{P}_i be a row
Most systems use column d that this one matrix.

Pros :-

- (1) gets to be learned to fit the data
- (2) Cannot extrapolate to indices outside $1, \dots, T$

Problem 2 :-

Note :- In self Attention - there are no element-wise non-linearities in self-Attention; we have more self-Attention layers just to give averages value vectors.

This can be easily resolved by a feed-forward network to post-process each output vector.

Problem 3 :-

We cannot make sentence based on the knowledge what we are going to produce in future. This is hypothetical or impossible. This happens then no need to train network.

So to inhibit it from predicting future based on present.

i.e. to enable parallelization we mask out attention to future words by setting attention scores to $-\infty$.

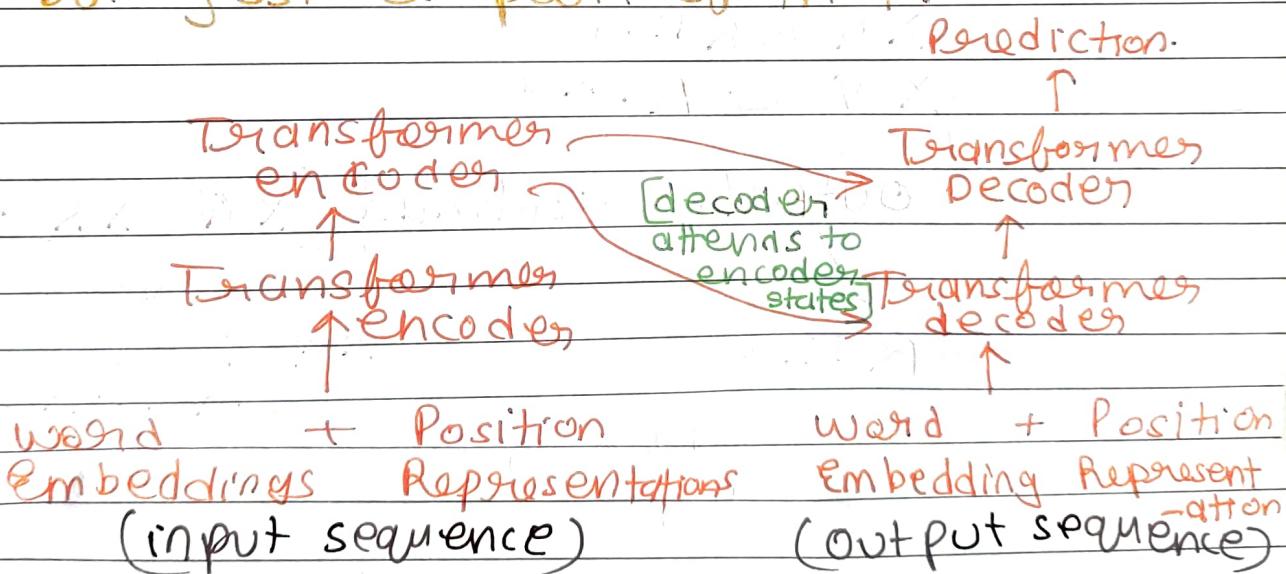
$$e_{ij} = \begin{cases} q_i^T K_j, & k < j \\ -\infty, & k > j \end{cases}$$

	start	The	chef	who
start	$-\infty$	$-\infty$	$-\infty$	$-\infty$
The			$-\infty$	$-\infty$
chef				$-\infty$
who				$-\infty$

Necessities for building a self-
Attention building block.

- (1) Self - Attention.
 - (2) Position Representations
 - (3) Non-linearities with simple feed forward network. (will be at output of self attention block)
 - (4) Masking
 - In order to parallelize operations while not looking at the future.
 - Keeps the information about the future from "leaking" to the past.

Above model is not a transformer
but just a part of that.



K, q, V come from same source in transformer

Let x_1, \dots, x_T be input vectors to the transformer encoder $x_i \in \mathbb{R}^d$

$K_i = K \cdot x_i$, where $K \in R^{d \times d}$ is key
 $q_i = Q \cdot x_i$, where $Q, V \in R^{d \times d}$ matrix
 $v_i = V \cdot x_i$, where $V \in R^{d \times d}$ is query and value matrix.

Transformer Encoder: Key - Query - Value Attention

Let's see how attention is computed in matrices

Let $X = [x_1, \dots, x_T] \in R^{T \times d}$ be the concatenation of input vectors

$$XK \in R^{T \times d}$$

$$XQ \in R^{T \times d}$$

$$XV \in R^{T \times d}$$

$$\text{OUTPUT} = \text{softmax}(XQ(XK)^T) * XV$$

$$E \in R^{T \times d}$$

$$\begin{aligned} XQ &= K + Q^T \\ XQ &\quad KT \in R^{T \times T} \end{aligned}$$