## Course 5 (RNN)

RNN have transformed NLP and speech recognition.
↓
a sequence model

In some sequence problems, both i/p X and o/p Y are sequences and in some cases either of them is a sequence.

e.g. Speech Recognition (Input: Audio      Output: text)

Music Generation (Input: Nothing or genre      Output: Audio)
                      name of music
                      to be generated.

Sentiment classifica^n (Input: Text   Output: Number )
                                         may be

Machine translation ( "   "   "      Output: Text)

Video activity recogni^n (I/p: frames of images   Output: Activity identified
                                                          in the i/p)
, etc ... like named entity Recognition (i.e identifying names in an i/p)

e.g.

x: Harry Potter and Hermoine Granger invented a new spell.
   $x^{(1)}$   $x^{(2)}$   $x^{(3)}$   $x^{(4)}$       $x^{(5)}$      $x^{(6)}$   $x^{(7)}$   $x^{(8)}$   $x^{(9)}$

y:   1      1      0      1         1         0      0      0      0
   $y^{(1)}$   $y^{(2)}$   $y^{(3)}$   $y^{(4)}$      $y^{(5)}$      $y^{(6)}$   $y^{(7)}$   $y^{(8)}$   $y^{(9)}$

$T_x = 9$   (length of i/p sequence)

$T_y = 9$   { In this e.g. $y^{(i)}$ is if $x^{(i)}$ is a name or not }

$y^{(i)<t>}$ = $t^{th}$ element in the o/p sequence of $i^{th}$ training example.

$T_x^{(i)}$ = input sequence length for training example i.

A vocabulary is next made. One way to build this dictionary is to look through the training sets and find the top occurring words or look at some online dictionary to know the most common words in the English language saved.

Suppose, vocabulary is

$$
\begin{bmatrix}
a \\
aaron \\
and \\
harry \\
potter \\
\vdots \\
zulu
\end{bmatrix}
\begin{matrix}
1 \\
2 \\
367 \\
4075 \\
6830 \\
\\
\end{matrix}
$$

(if 10,000 words)

So, $x^{(1)}$  $x^{(2)}$

$$
\begin{bmatrix}
0 \\
0 \\
\vdots \\
0 \\
1 \\
0 \\
\vdots \\
0
\end{bmatrix} \leftarrow 4075
\qquad
\begin{bmatrix}
0 \\
0 \\
\vdots \\
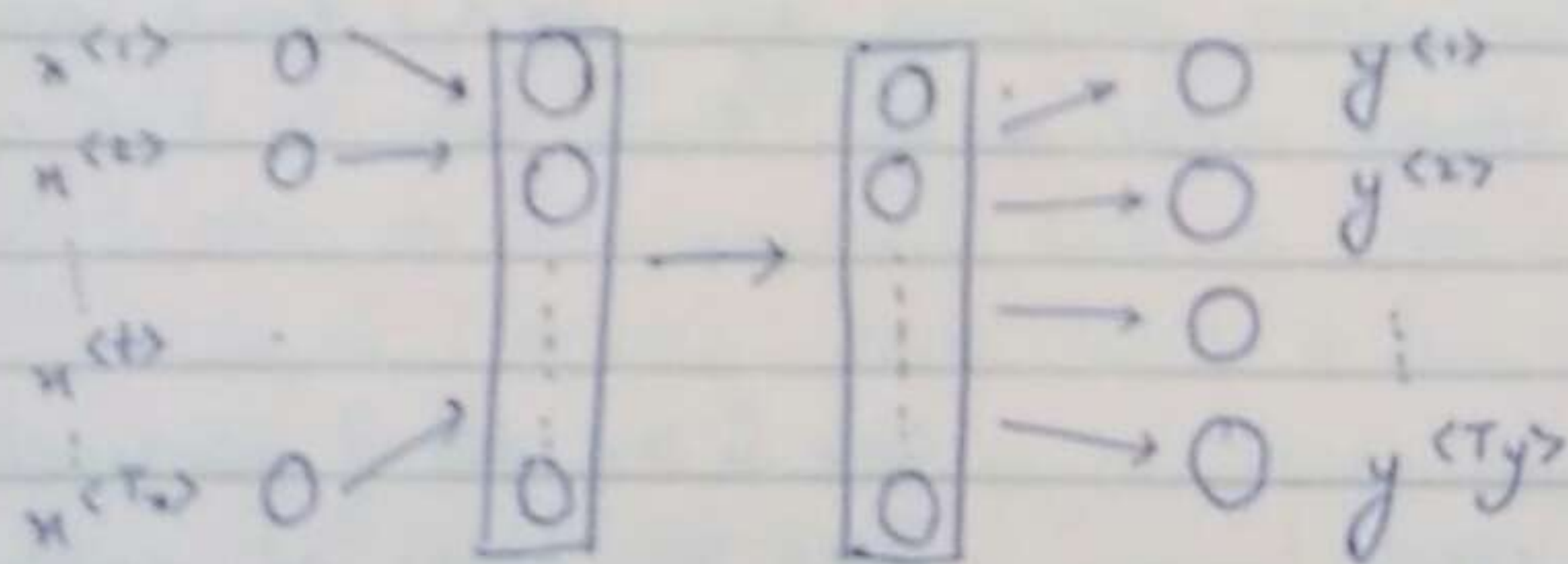\vdots \\
1 \\
\vdots
\end{bmatrix} \leftarrow 6830
$$

and so on.

All are 10,000 dimensional one-hot vector.

If some work that doesn't exist in the vocabulary is encountered, then a new token 'unknown' is added to the dictionary.

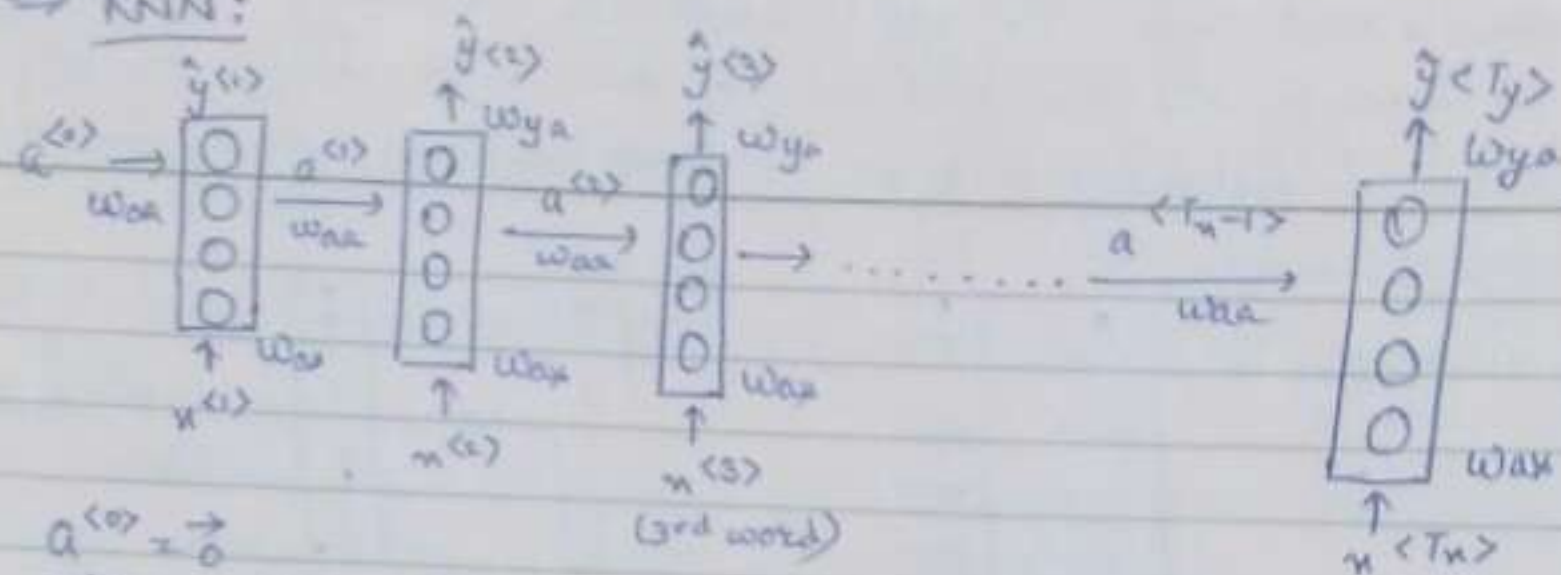→ Why not a standard network for such problems?



$x^{(t)}$ is a one hot vector of 10,000 dimension.

Disadvantages of using this model for such data:

padding can't always be done

- Inputs, Outputs can be different lengths in different examples and .
- Doesn't share features learned across different positions of text.
- Too large input size [$10,000 \times T_x$] for just one input sentence. which would lead to an enormous no. of parameters in the weight matrix.

RNN does not have these disadvantages.

# ⇒ RNN:



$$a^{<0>} = \vec{0}$$

$$a^{<1>} = g\left(W_{aa}\, a^{<0>} + W_{ax}\, x^{<1>} + b_a\right) \quad \leftarrow \text{tanh/Relu}$$
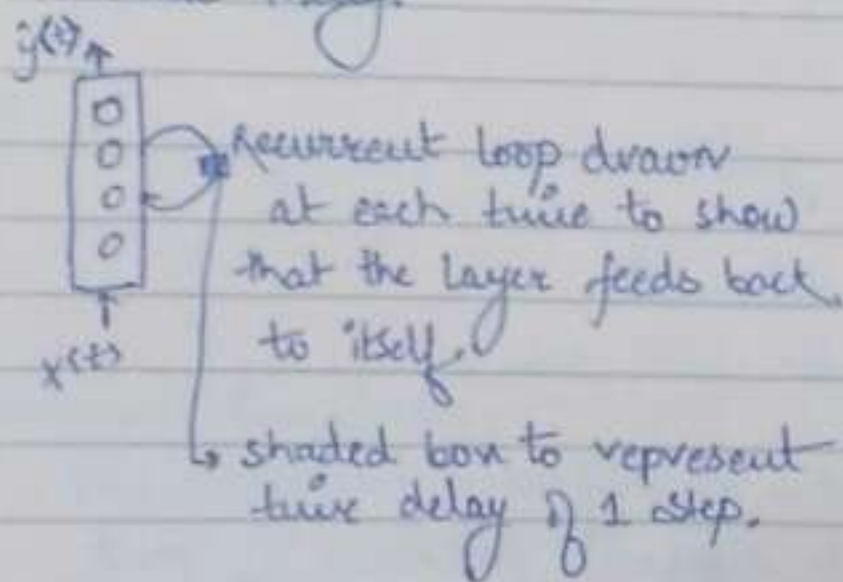
$$\hat{y}^{<1>} = g'\left(W_{ya}\, a^{<1>} + b_y\right) \quad \leftarrow \text{Sigmoid in case of classifica}^n \text{ problem}$$

$$a^{<t>} = g\left(W_{aa}\, a^{<t-1>} + W_{ax}\, x^{<t>} + b_a\right)$$

$$\hat{y}^{<t>} = g'\left(W_{yc}\, a^{<t>} + b_y\right)$$

$(W_{pq})$ means some $q$ type quantity is used to calculate some $p$ type quantity.

At every step RNN passes activa$^n$ from previous layer to the next layer to use. $a^{<0>}$ is a vector of zeros made up at time zero to kick up the whole thing.



Recurrent loop drawn at each time to show that the layer feeds back to itself.

↳ shaded box to represent time delay of 1 step.

$W_{ax}$ is the set of parameters used by RNN at each time step.

Horizontal connections are governed by the set of parameters $W_{aa}$ used at every time step.

Similarly $W_{ya}$ governs the output predictions.

\* Limita$^n$ of this network: It uses only the information from previous layers i.e only the info earlier in the sequence to

make a prediction, i.e while predicting $\hat{y}^{<3>}$, it doesn't disc info about the words $n^{<4>}$, $n^{<5>}$ or any later word in the sequence.

E.g.

i) He said, "Teddy Gupta was a great man".

ii) He said, "Teddy bears are on Sale".

In i) and ii), to know whether Teddy is the name of a person, knowing just 1st. 3 words is not sufficient. Info about the later words is necessary to decide if Teddy is a name or not. In i) it is and in ii) its not, this is decided on the basis of info from later words only.

For this we will use __B-RNN__ (Bidirectional RNN).

Also, $a^{<t>} = g(W_{aa} a^{<t-1>} + W_{ax} x^{<t>} + b_a)$ —①

$$\boxed{a^{<t>} = g(W_a [a^{<t-1>}, x^{<t>}] + b_a)}$$ —②

where $W_a = [W_{aa} \vdots W_{ax}] = (100, 10100)$ dimensional if I

If $a$ was 100 dimensional and $n$ 1000 dimensional, then $W_{aa}$ is $(100, 100)$ and $W_{an}$ is $(100, 10,000)$

$[a^{<t-1>}, x^{<t>}]$ means $= \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \updownarrow 10100$ justifying ① = ②.

So, $\tilde{y}^{<t>} = g(W_y a^{<t>} + b_y)$
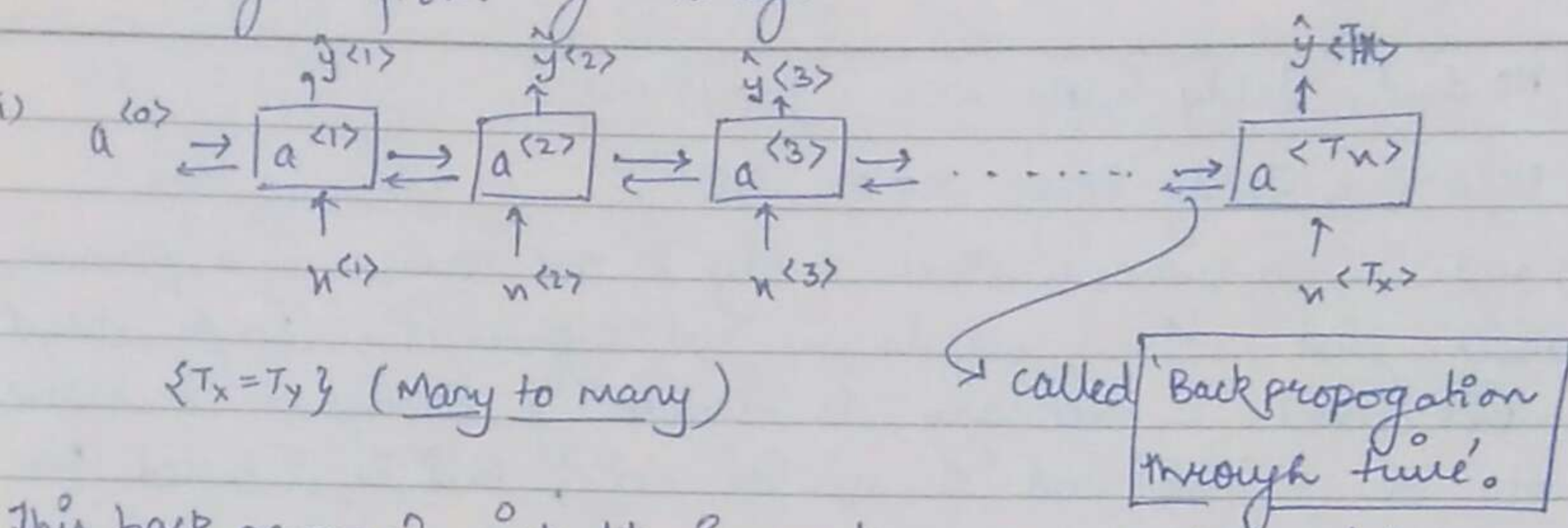
$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) = -y^{<t>} \log \hat{y}^{<t>} - (1-y^{<t>}) \log(1-\hat{y}^{<t>})$

↳ loss associated with single word of a single sentence.

$$\boxed{L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{<t>}(\hat{y}^{<t>}, y^{<t>})}$$ {loss for entire sequence}

In the back propogaⁿ procedure, the significant calculaⁿ is the one that goes from right to left

i)

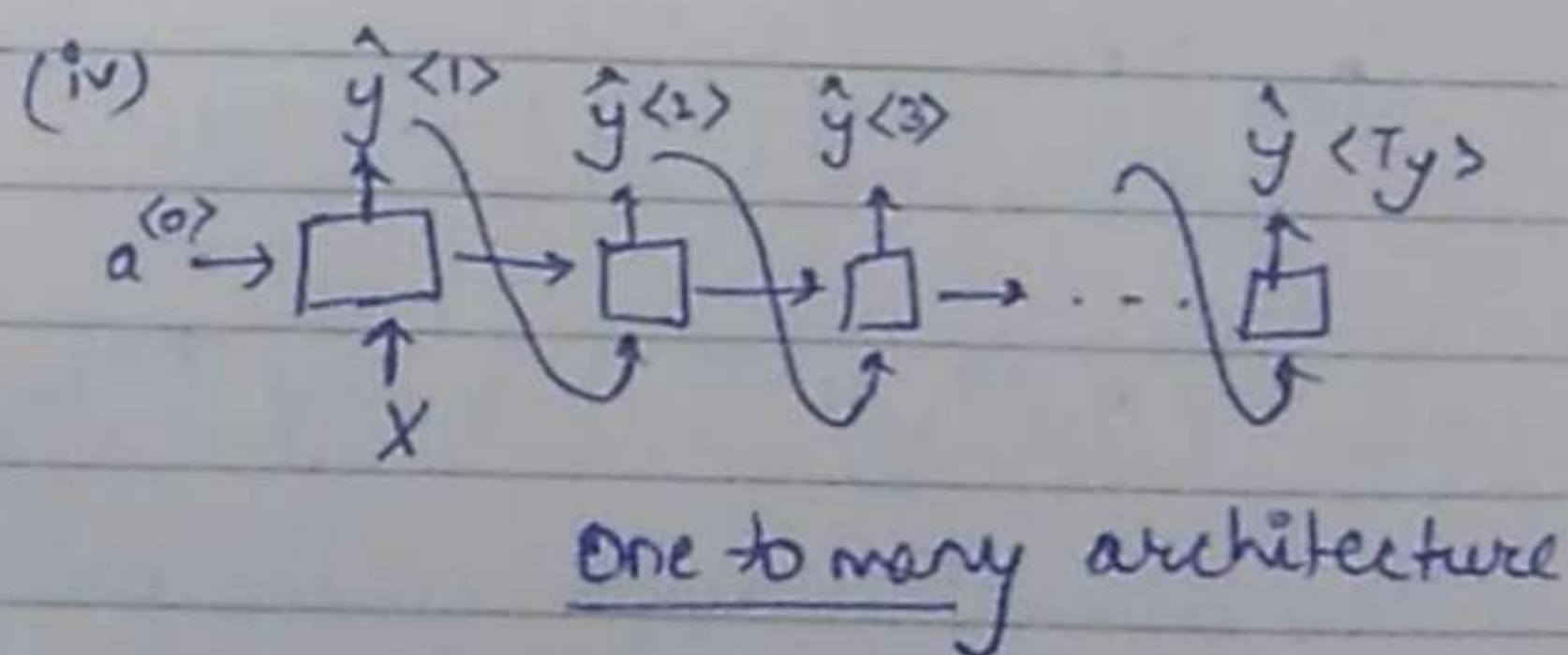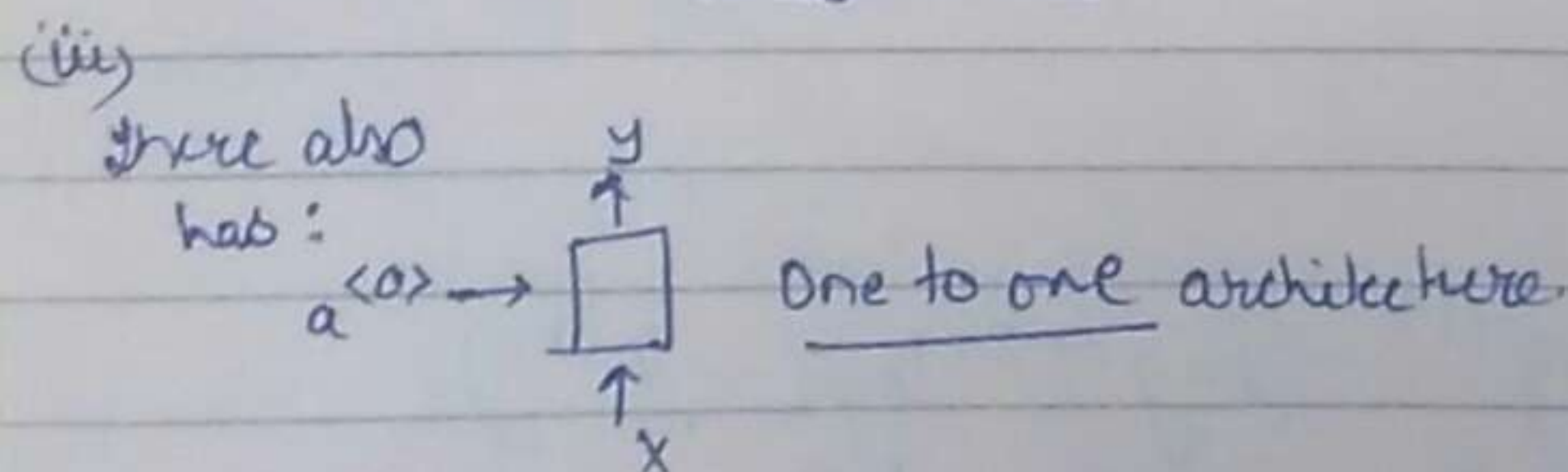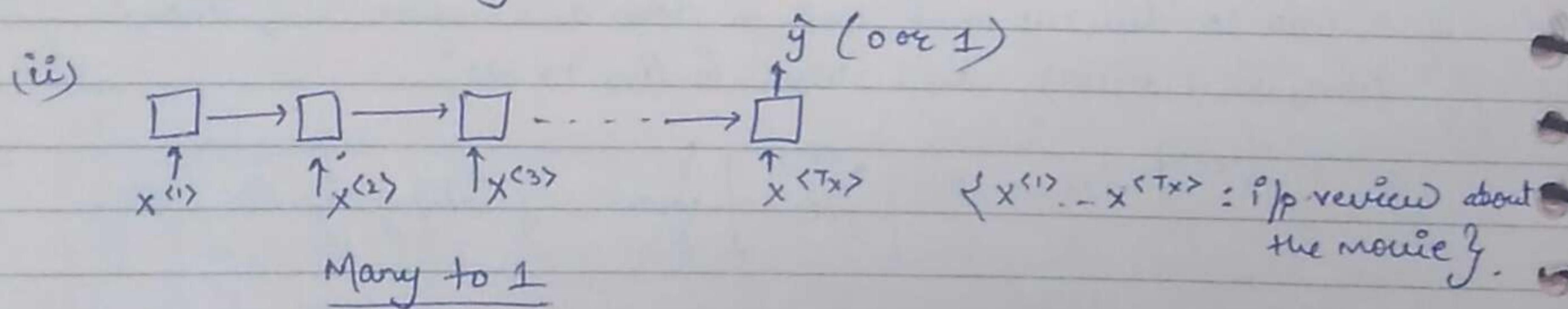$$a^{<0>} \rightleftarrows \boxed{a^{<1>}} \rightleftarrows \boxed{a^{<2>}} \rightleftarrows \boxed{a^{<3>}} \rightleftarrows \cdots \cdots \rightleftarrows \boxed{a^{<T_n>}}$$

with outputs $\hat{y}^{<1>}$, $\hat{y}^{<2>}$, $\hat{y}^{<3>}$, $\hat{y}^{<T_n>}$ and inputs $x^{<1>}$, $x^{<2>}$, $x^{<3>}$, $x^{<T_x>}$

$\{T_x = T_y\}$ (Many to many)

$\rightarrow$ called Back propogation through time.

This back propogaⁿ just like in earlier case, helps update parameters through gradient descent.

$\rightarrow$ what if $T_x \neq T_y$ i·e no. of i/p $\neq$ no. of o/p?
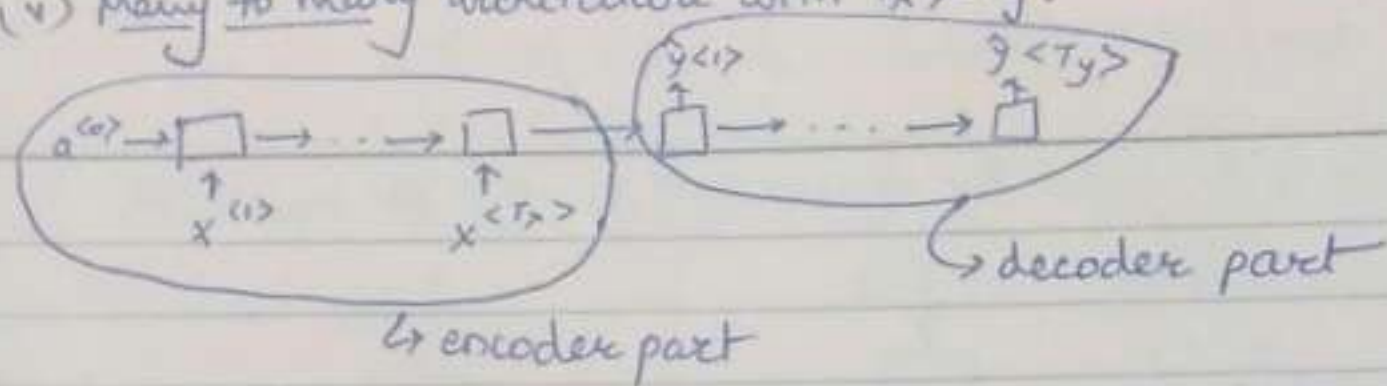
e·g· sentiment classificaⁿ.

i/p is a sequence of words but o/p is 0/1 i·e either +ve or -ve review for this we use <u>many to 1</u> architecture.

(ii)

$$\boxed{} \rightarrow \boxed{} \rightarrow \boxed{} \cdots \cdots \rightarrow \boxed{}$$

with inputs $x^{<1>}$, $x^{<2>}$, $x^{<3>}$, $x^{<T_x>}$ and output $\hat{y}$ (0 or 1)

$\{x^{<1>} - x^{<T_x>} : $ i/p review about the movie $\}$.

<u>Many to 1</u>

(iii)
there also has:
$$a^{<0>} \rightarrow \boxed{} \rightarrow \hat{y}$$
with input $x$

<u>One to one</u> architecture.

(iv)

$$a^{<0>} \rightarrow \boxed{} \rightarrow \boxed{} \rightarrow \boxed{} \rightarrow \cdots \rightarrow \boxed{}$$

with outputs $\hat{y}^{<1>}$, $\hat{y}^{<2>}$, $\hat{y}^{<3>}$, $\hat{y}^{<T_y>}$ and input $x$

<u>One to many</u> architecture

e·g· music generaⁿ or sequence generaⁿ.

(v) <u>Many to many</u> architecture with $T_x \neq T_y$.



&rarr; decoder part

&rarr; encoder part

e.g. French sentence translated to English. i/p and o/p may have diff length of sequences.

&rArr; <u>Building a language model using RNN</u> : (e.g. to generate Shakespeare like texting).

Speech recognition:

$P(\text{The apple and pair salad}) = 3.2 \times 10^{-13}$

$P(\text{The apple and pear salad}) = 5.7 \times 10^{-10}$

{ language model learns that which of the two sentences has higher probability of being spoken depending on its learning. this is used both in speech recognition & machine/language translation }

&rarr; How to build a language model?

Training set : large <u>corpus</u> of text in whatever language we want to build a model. &rarr; chunk / large body
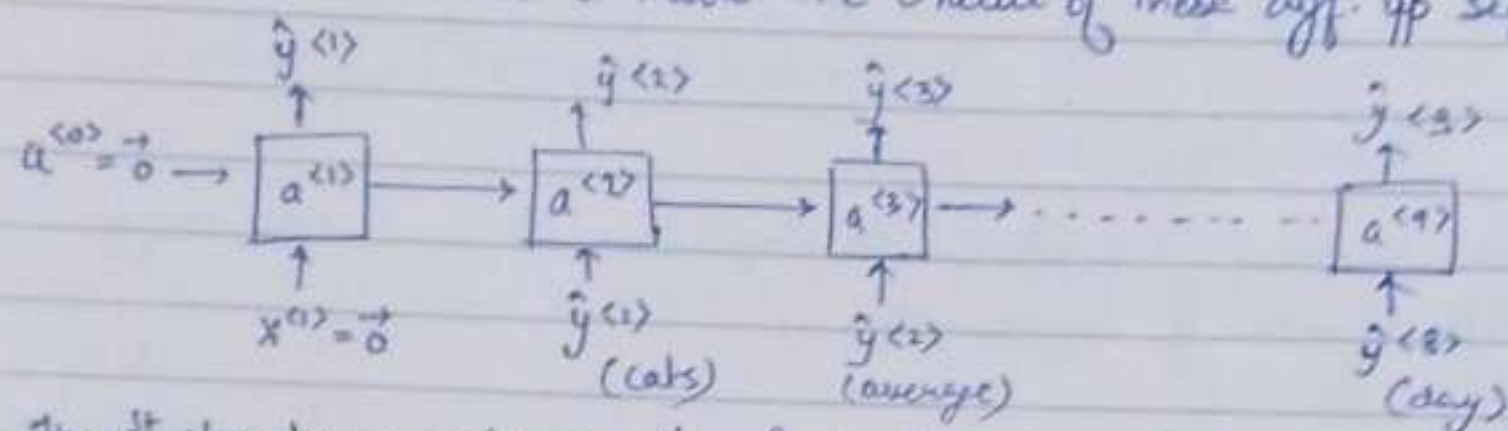
Step1 : Tokenize.

e.g. if i/p is Cats average 15 hrs of sleep a day.
$y^{<1>}$ $y^{<2>}$ $y^{<3>}$ . . . . . . $y^{<8>}$  &lt;EOS&gt;

&rarr; one-hot vectors to represent them in dictionary.

additional token to tell end of sentence, if needed.

If some word appears in the training set which is not not in the dictionary, then token UNK (unknown) is used for it.

Step2: build RNN model to model the chance of these diff. i/p sequences.



The 1st step has a softmax trying to know the 1st word depending on probabilities of the 1st word being any one from the dictionary of 10,000 words. Hence, this softmax o/p is 10,000 dimensional vector. Which ever word gets the highest probability is considered to be the 1st word of sentence.

In 2nd step, o/p is again predicted by softmax with probability of 2nd word being any one from dictionary given 1st words is cats.

$$P(\underline{\quad}|\text{cats}) = \hat{y}^{<2>}.$$

$$\hat{y}^{<3>} = P(\underline{\quad}|\text{"cats average"})$$

$$\hat{y}^{<9>} = P(\underline{\quad}|\text{"cats average 15 hrs of sleep a day"}).$$
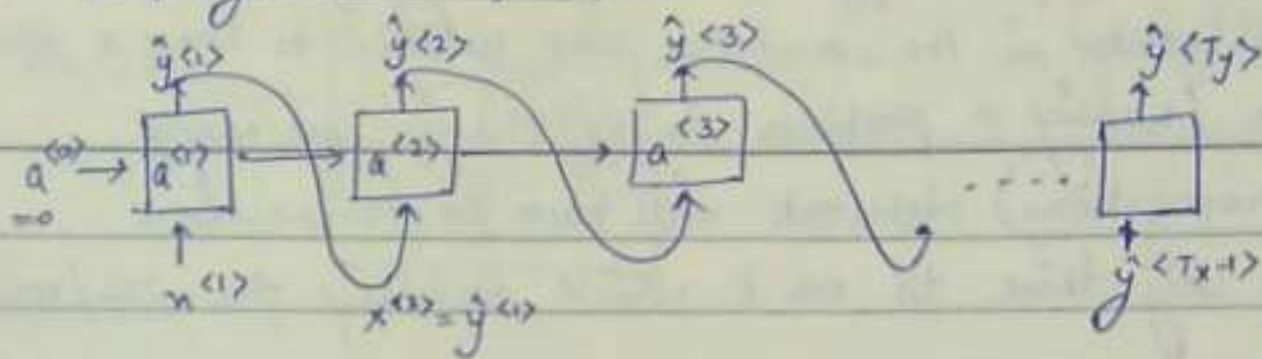
Note: In this model $\underline{x^{<t>} = y^{<t-1>}}$

This RNN learns to predict one word at a time given from left to right.

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

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

The diff. time stamps in RNN are nothing but layers.

⇒ Sampling novel sequences:



$\hat{y}^{<1>}$ i·e the probabilities help sample the first words. Then the 2nd timestamp helps sample the second words based on $\hat{y}^{<1>}$ i·e the 1st word sampled, and so on.

where $\hat{y}^{<T_y>}$ = EOS, we know end of sentence is reached and we can stop. This can also be done by setting a word limit.

Also, to avoid UNK from getting sampled, any samples containing it are rejected and resampling is done from remaining vocabulary. This all will generate a random sentence.

**Note** Depending on application character level RNNs can also be built, with $n^{<t>}$, $y^{<t>}$ as characters and not words. Vocabulary also consists of characters (alphabets, digits, space, punctua⁰) in this case.

Pros: no worry about UNK tokens, any word can be generated.

Cons: much longer sequences in character level model and not able to capture long range dependencies (i·e how the earlier parts of sentence also affect the later part of sentence) like the word level model does.
Also, they are computationally expensive to train.

→ Vanishing gradients with RNN:

eg of long range dependencies: The cat which already ate ··········, was full.
            The cats    "    "    "    ···---··, were full.

The RNNs we have seen so far suffer from vanishing gradient descent and as a result later words in the sentence are unable to have a effect on the earlier words (mainly a problem in upgrading so many weights during backpropogation.) Network will have to remember cat/cats for a very long time to use it while deciding for was/were

Hence, in basic RNN value of $\hat{y}^{<t>}$ is affected mainly by $\hat{y}$ which are close to this $\hat{y}^{<t>}$, hence not capturing long term dependencies.

It is easy to spot exploding gradients than vanishing gradients as the parameters just blow up resulting in lot of NaNs (not a number).
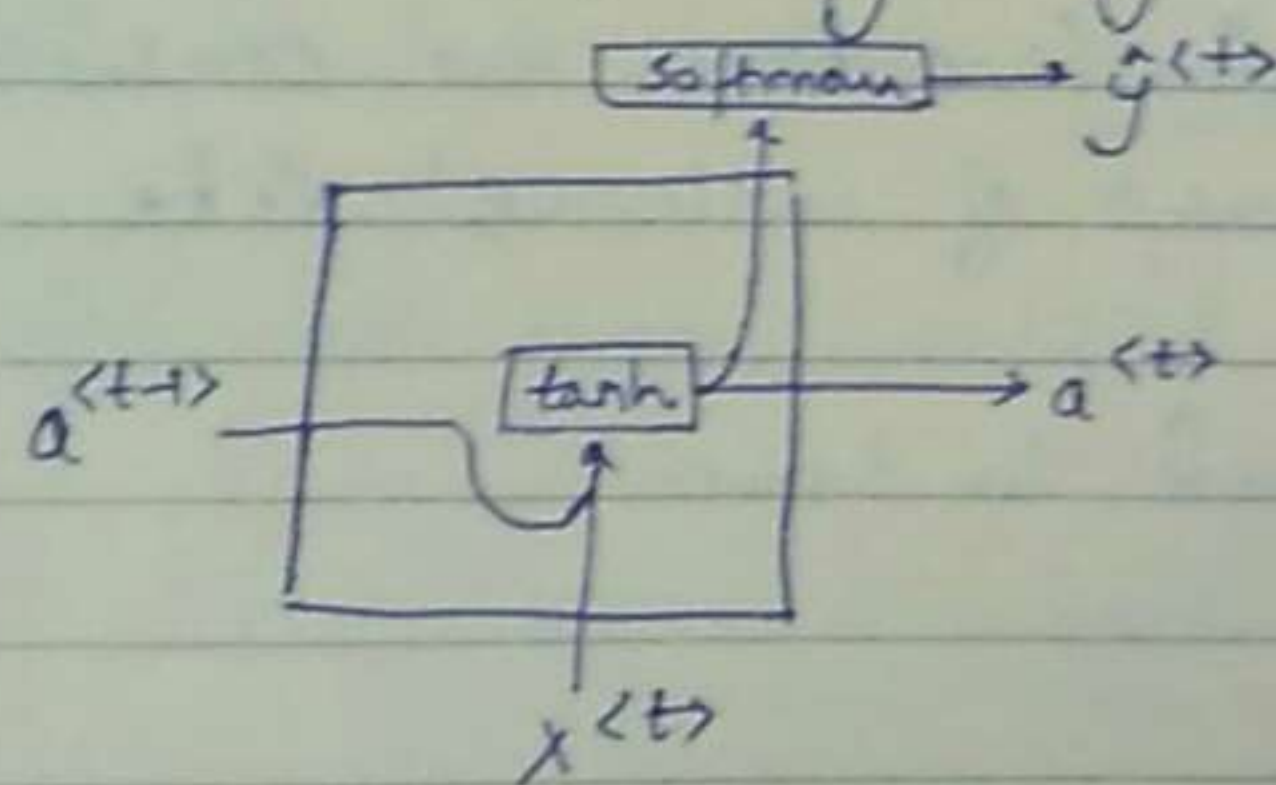
Note Solution to exploding gradient descent : Gradient clipping

i.e clipping the gradient vectors whose dimensions have grown beyond a threshold.

* Difficult to solve vanishing gradient descent problem.

⟹ [Gated Recurrent Unit (GRU)] {solu" to vanishing gradient descent}

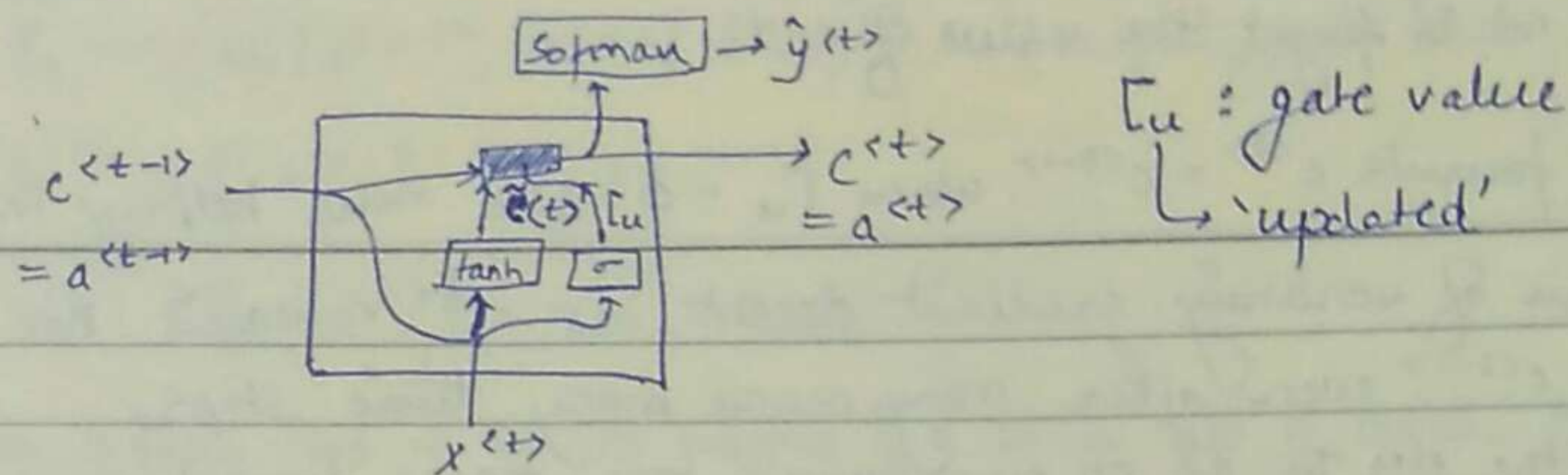It helps in capturing long distance dependencies.



Basic RNN

Note GRU has a c = memory cell to remember whether cat was singular or plural.

for GRU ; $c^{<t>} = a^{<t>}$

memory cell value at timeslot 't' = activa" output at timeslot 't'.

GRU picture

At every time step, memory cell value is overwritten by $\tilde{c}^{<t>}$

$$\tilde{c}^{<t>} = \tanh\left(W_c\left[\underset{\underset{c^{<t-1>}}{\downarrow}}{c^{<t>}}, x^{<t>}\right] + b_c\right)$$

$$\Gamma_u = \sigma\left(W_u\left[c^{<t-1>}, x^{<t>}\right] + b_u\right) \qquad \{\text{this is either 0 or 1}$$
(gate)                                                         most of the time.}

$u \rightarrow$ update

(value b/w 0 and 1) $\longrightarrow$ Suppose if we set $c^{<t>}$ as 1 when cat is singular.

e.g.    $c^{<t>} = 1$ ........ - - - - - ... - ... $c^{<t'>} = 1$

The <u>cat</u>, which already ate . . . . . . . , <u>was</u> full.

$c^{<t>}$ remains to be 1 in $c^{<t'>}$ as well representing that singular is there implying 'chose was'.

<u>Job of $\Gamma_u$</u> is <u>to decide when to update these values of $c^{<t>}$</u>. When we came across cat (the subject of the sentence), it was considered the right time to update $c^{<t>}$ and then maybe when it was done being used, no need to memorize it any longer!

$\longrightarrow$ element wise multiplica⁀

So, $c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$

If $\Gamma_u = 1 \Rightarrow$ set $c^{<t>}$ to a new value $\tilde{c}^{<t>}$

$\ast \Rightarrow \Gamma_u$ was 0 every where b/w cat and was, but got = 1 at 'was'. indicating to update $c^{<t>}$ at 'was'. $\Gamma_u = 0$ indicated not to

update and not to forget the value of $c^{<t>}$.

from formula $c^{<t>} = c^{<t-1>}$ when $\Gamma_u = 0$ and thus helping in the problem of vanishing gradient decent as $c^{<t>}$ remains the same as $c^{<t-1>}$ even after many many many time steps, allowing the NN to go on even very long range dependencies.

$c^{<t>}$ can be a n-dimensional vector, so, $\tilde{c}^{<t>}$ is also of same dimension. And so will be $\Gamma$, telling which of the n-dimensional bits in $\tilde{c}^{<t>}$ to be updated. 1 bit might be to remember that cat is singular and other might be remembering that food is being talked about.

## Full GRU: (commonly used version)

$$\tilde{c}^{<t>} = \tanh\left(W_c\left[\Gamma_n \times c^{<t-1>}, x^{<t>}\right] + b_c\right)$$

tells the relevance of $c^{<t-1>}$ in computing $c^{<t>}$.

$$\Gamma_r = \sigma\left(W_r\left[c^{<t-1>}, x^{<t>}\right] + b_r\right)$$

(relevance gate)

→ **LSTM (long short Term Memory)**

- same advantages as that of GRU but even more powerful & general version of GRU.
- In LSTM $a^{<t>} \neq c^{<t>}$

eqn of LSTM:

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

$$\Gamma_u = \sigma\left(w_u\left[a^{<t-1>}, x^{<t>}\right] + b_u\right) \qquad \{\text{update gate}\}$$

$$\Gamma_f = \sigma\left(w_f\left[a^{<t+1>}, x^{<t>} + b_f\right]\right) \qquad \{\text{forget gate}\}$$

$$\Gamma_0 = \sigma \left( W_0 \left[ a^{<t-1>}, x^{<t>} \right] + b_0 \right) \quad \{\text{output gate}\}$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_0 * \left( c^{<t>} \right) = \Gamma_0 * \tanh c^{<t>}$$

So, 3 gates used in LSTM instead of 1 as in case of GRU. Hence, it's a bit more complicated and places the gates at different places.



LSTM picture.

<u>Note</u> Peephole connection : means that the gate values may depend not just on $a^{<t-1>}$ and $x^{<t>}$ but also on the previous memory cell value i·e $c^{<t-1>}$.
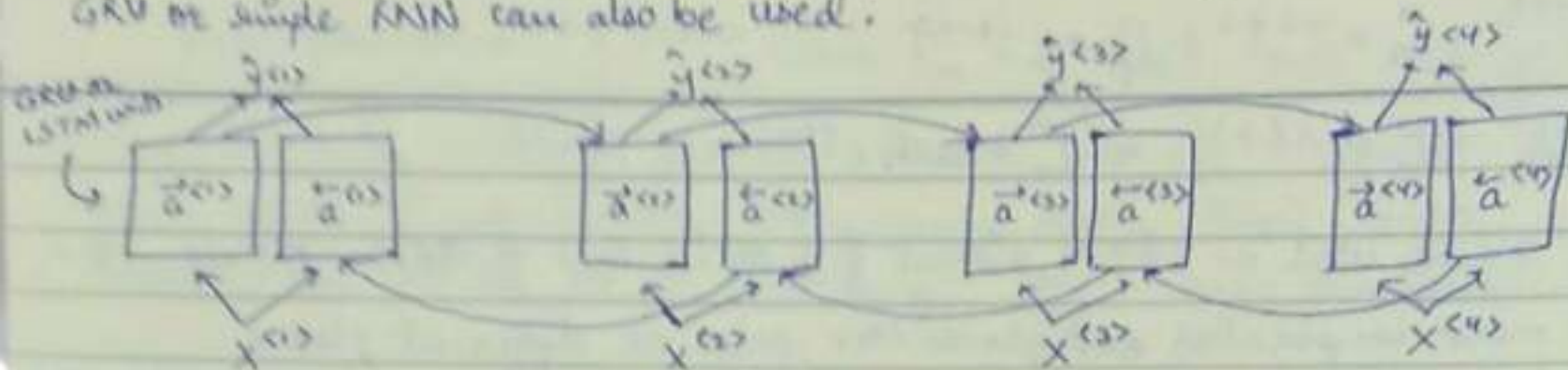
* GRUs came later than LSTMs as a simplified version of the much complex LSTMs.

⇒ Bidirectional RNN

GRUS were much simpler than LSTMs and had less computational cost due to tnce of 2 gates but LSTMs are more powerful & effective becoz of tnce of 3 gates. Due to simplicity GRU is easier to be implemented in big deep networks.

- BRNN allows to take info both from earlier and later sequences, thus helping us with that 'Teddy Bear' example taken earlier.

Bidirectional RNN with a LSTM appears to be commonly used, but GRU or single RNN can also be used.



$$\hat{y}^{<t>} = g\left(W_y[\overrightarrow{a}^{<t>}, \overleftarrow{a}^{<t>}] + b_y\right)$$

BRNN picture

In this BRNN, the sentence is being read both from left to right and right to left. $\overrightarrow{a}^{<1>}$ feeds into $\overrightarrow{a}^{<2>}$, $\overrightarrow{a}^{<2>}$ feeds into $\overrightarrow{a}^{<3>}$ and so on. $\overleftarrow{a}^{<4>}$ feeds into $\overleftarrow{a}^{<3>}$, $\overleftarrow{a}^{<3>}$ into $\overleftarrow{a}^{<2>}$ and so on. And both of them feed into $y^{<t>}$ thus taking into consider both earlier and later sequence while deciding if Teddy is a name or not for example.

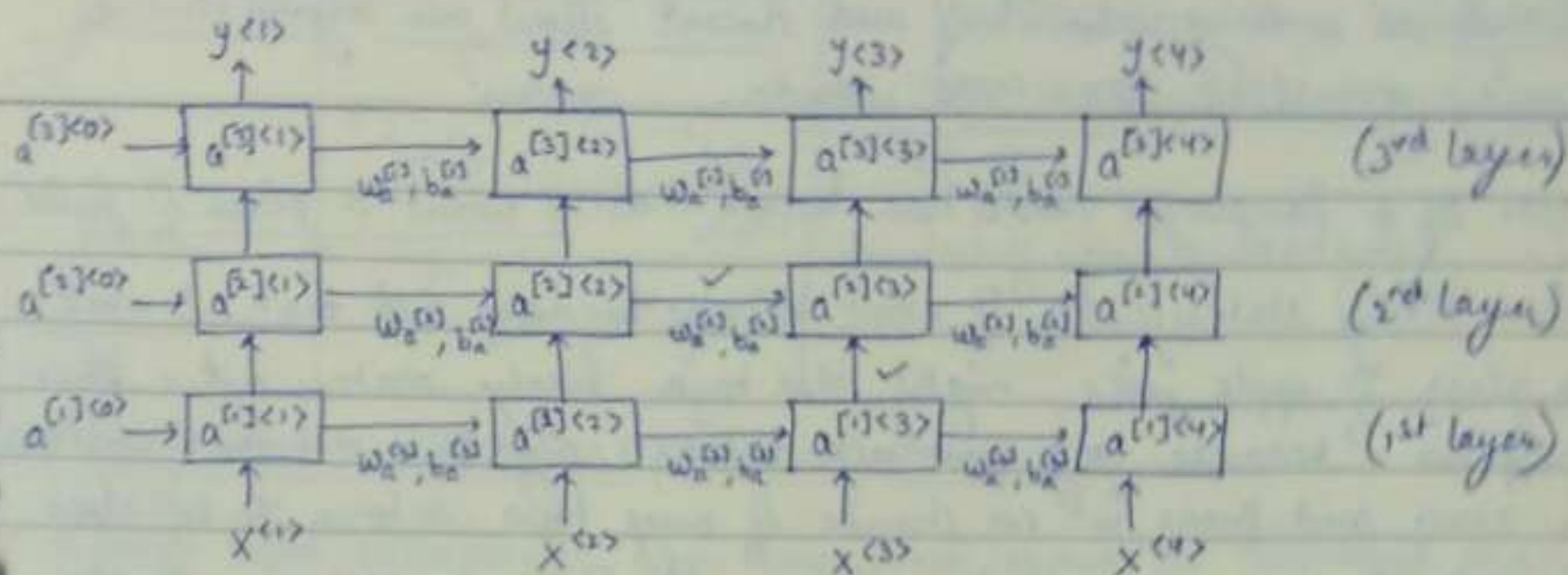Hence, it uses info potentially from the entire sentence.

Disadv: entire sequence of data is needed before making any decision

eg. If a speech recognition system is made using BRNN, then we will have to wait for the person to utter the complete sentence before any processing can be done to reach a conclusion, hence it won't be possible to do real time speech recognition.

Adv: when we have entire sequences available like in NLP problem, BRNN is quite effective.

To get better results we can stack multiple layers of RNN, and get deep RNN versions.

## ⟹ Deep RNN



| | | | | |
|---|---|---|---|---|
| $a^{[3]\langle 0\rangle}$ → $a^{[3]\langle 1\rangle}$ | $a^{[3]\langle 2\rangle}$ | $a^{[3]\langle 3\rangle}$ | $a^{[3]\langle 4\rangle}$ | (3rd layer) |

3 hidden layer RNN

$a^{[\ell]\langle t\rangle}$ : $t^{th}$ activation of $\ell^{th}$ layer of deep RNN.

$$a^{[2]\langle 3\rangle} = g\left(w_a^{[2]}\left[a^{[2]\langle 2\rangle}, a^{[1]\langle 3\rangle}\right] + b_a^{[2]}\right)$$

for RNNs, having 3 layers is quite deep bcoz of the temporal dimension, that already makes the network quite big. More layers can be stacked up to predict ŷ but without horizontal connections.

Deep RNNs are computationally very expensive, hence not much hidden layer used as are used in case of CNNs.
Bidirectional RNN can also be converted to deep RNN.

## ⟹ Word embeddings :
- a way of representing words.

Vocabulary = [a, aaron, ....., zulu, <UNK>]

every word is then represents' as one hot encoding with 1 at the posi" where it exists in the vocabulary, Represented by $O_{zulu}$ or $O_{aaron}$

In this one-hot representation, weakness is that it treats each word as a thing unto itself and doesn't allow an algorithm to easily generalise the cross words.

e.g. If a language model has learnt that 3 want a glass of juice is a likely statement, it won't be able to make out that 3 want a glass of apple juice might also be a likely statement. This happens because any product between any two one-hot vectors is zero and there is no chance of any two vectors to be close to each other as compared to others. So, it can't recognize that apple and orange are more similar than any other word like king.

Hence, a need for featurised representation rather than one-hot vector for each of the word

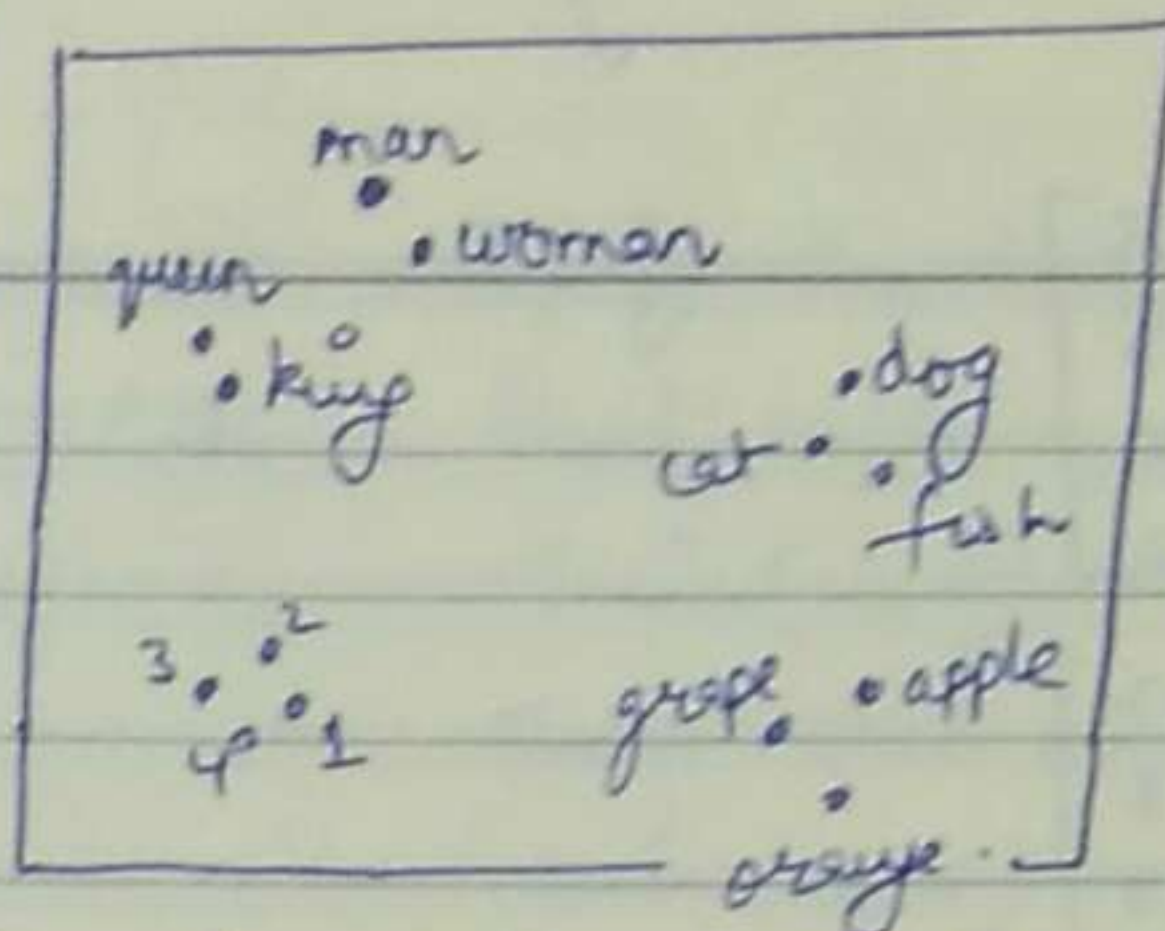so as to learn a set of features and values for each of them?

e.g.

| Features | Man | Woman | King | Queen | Apple | Orange | some words in the vocabulary |
|---|---|---|---|---|---|---|---|
| Gender | -1 | 1 | -0.95 | 0.97 | 0.00 | 0.01 | ↑ |
| Royal | 0.01 | 0.02 | 0.93 | 0.95 | -0.01 | 0.00 | the vocabulary −ary 3 |
| Age | 0.03 | 0.02 | 0.7 | 0.69 | 0.03 | -0.02 | |
| Food | 0.04 | 0.01 | 0.02 | 0.01 | 0.95 | 0.97 | |
| Noun | | | | | | | |

So, now each word is a n-dimensional vector and now its easier to notice the relationship b/w words like king and queen, man & woman, orange & apple, etc. So, it generalises better across different words. Though it is difficult to learn what exactly the word says or represent, but easy to tell which words are more related and which are less.

<u>t-SNE</u> algorithm to convert these n-dimensional vectors to 2D space to visualize them better. It a very complicated & non-linear mapping.

e.g.

```
┌─────────────────────────┐
│   man                   │
│ queen  •woman           │
│   •king        •dog     │
│           cat•  •        │
│                  fish    │
│  3  •2                   │
│  4  •1     grape• •apple │
│                          │
│        • orange          │
└─────────────────────────┘
```

{concept of word embeddings}

→ <u>Using word embeddings for transfer learning:</u>

1) learn word embeddings from large text corpus or download pre-trained word embeddings.

2) Transfer embedding to new task, with smaller training set.

<u>Note</u> One hot vectors are fast than the featurized vectors but give less info.

3) continue to fine tune the word embeddings with new data, if wanted [just a case of choice]

Word embeddings tend to make the biggest difference when the task, to be carried out has a relatively smaller training set.

Applica²: Named entity recogni", text summarize², co-reference resolution, parsing.

<u>Note</u> Transfer learning from A to B is useful when A has lot of data than B.

Word embeddings ~ Picture encoding (as seen in Siamese network)

Only difference is, in word embedding vocabulary is limited but pictures can be any number in case of picture encoding.

— <u>word embeddings help in knowing analogies</u> e.g. if man : woman ::
king : ?.

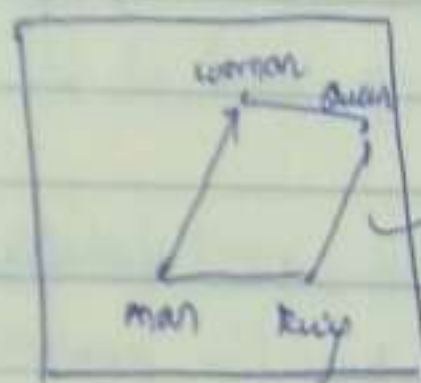To solve this we use the word embedding vectors as made on previous pages.

$$e_{man} = \begin{bmatrix} -1 \\ 0.01 \\ 0.03 \\ 0.09 \end{bmatrix} \qquad e_{king} = \begin{bmatrix} -0.95 \\ 0.93 \\ 0.70 \\ 0.02 \end{bmatrix}$$

$$e_{woman} = \begin{bmatrix} 1 \\ 0.02 \\ 0.02 \\ 0.01 \end{bmatrix} \qquad e_{queen} = \begin{bmatrix} 0.97 \\ 0.95 \\ 0.69 \\ 0.01 \end{bmatrix}$$

interesting property of these vectors:

$$e_{man} - e_{woman} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix} \approx e_{king} - e_{queen} \implies \text{main diff b/w the two is gender.}$$

used to calculate the analogy man : woman :: king : ?



this vector is basically the one representing gender and is approx. equal.

<u>119m analogy</u>

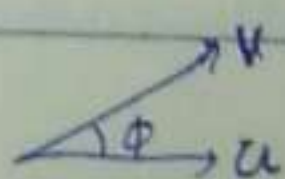$$\boxed{\text{find word } w \text{ such that : } \underset{w}{\arg\max} \; sim(e_w, \; e_{king} - e_{man} + e_{woman})}$$

i.e finding a word w that maximises its similarity with $e_{king} - e_{man} + e_{woman}$.

Many 119m analogy relationships will be broken by <u>t-SNE</u> algorithm. One must not count on it to find out $e_w$.

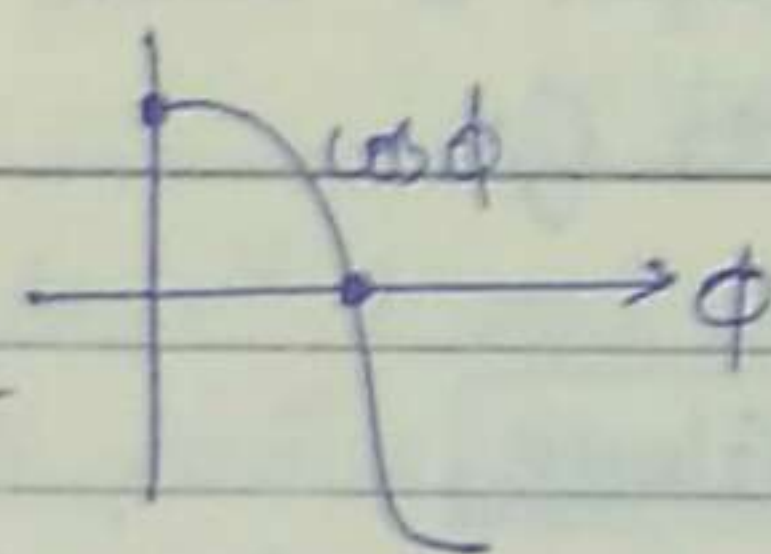<u>Cosine similarity:</u>

$$\boxed{\cos\phi = sim(u,v) = \frac{u^T v}{||u|| \, ||v||}}$$

— inner product b/w u and v.

If u and v are similar, their inner product will be a large value.



So, if the two words have $\phi = 0$ b/w them that means they are similar and so $\text{sim}(u,v) = \cos\phi = 1$. If they are anti then $\text{sim}(u,v) = \cos\phi = -1$ and if $\phi = 90°$, $\text{sim}(u,v) = \cos\phi = 0$.

Also, $\boxed{||u-v||^2}$ is measure of dissimilarity b/w u and v unlike $\text{sim}(u,v)$ which is a measure of similarity. $\text{sim}(u,v)$ is used more often. 'u' & 'v' are nothing but word embeddings. All kind of analogies can be learnt by just running a word embedding algorithm on the large text corpus.

⇒ <u>Embedding matrix</u> : (E)

$E = (300, 10,000)$ matrix

↓ no. of features     ↓ vocabulary size

$E \cdot O_{orange} = e_{orange}$ ⟶ In general $\boxed{E \cdot O_j = e_j}$ = embedding for word $j$ in vocabulary.

↓ embedding matrix   ↓ one-hot vector    embedding vector for orange.

a  aaron ... orange ... zulu



300

10,000

$E$

$O_{orange}$ {1 at pos^n where orange lies in the vocabulary}

$e_{orange}$

<u>Note</u> But multiplying one-hot vector as a matrix is not practised bcoz most of this matrix is 0. So, a specialised fn is used to just look up to a column in E.

In Keras, the embedding layer simply pulls out the column corresponding to the word (e.g. orange) from matrix E instead of doing matrix multiplication $\bar{e}$ the one-hot-vector $O_j$.
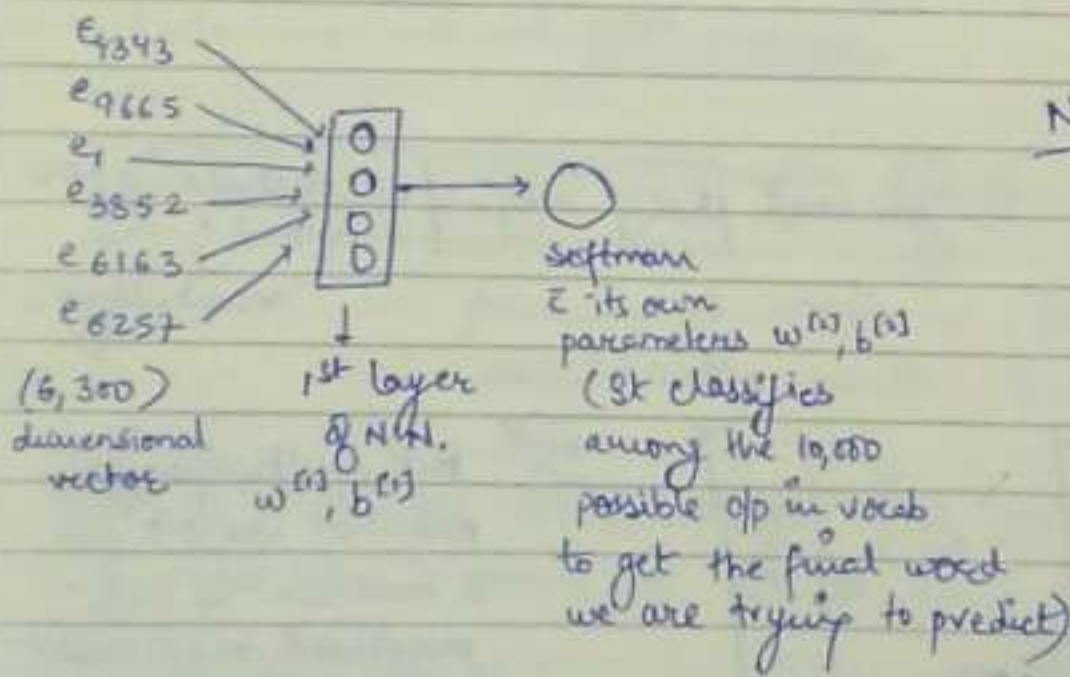
⟹ learning word embeddings : [word 2 Vec and Glove]

Building a neural language model is the small way to learn a set of embeddings. Here we will build a NN to predict the next word in the sentence    I want a glass of orange ____ .

4343   9665   1   3852   6163   6257

↳ posi^n of these words in the vocabulary.

I        $O_{4343} \longrightarrow E \longrightarrow \boxed{e_{4343}}$
want     $O_{9665} \longrightarrow E \longrightarrow e_{9665}$       → each of this is
a        $O_1 \longrightarrow E \longrightarrow e_1$                 a 300 dimensional
glass    $O_{3852} \longrightarrow E \longrightarrow e_{3852}$       embedding vector.
of       $O_{6163} \longrightarrow E \longrightarrow e_{6163}$
orange   $O_{6257} \longrightarrow E \longrightarrow e_{6257}$

$e_{4343}$
$e_{9665}$
$e_1$
$e_{3852}$
$e_{6163}$
$e_{6257}$

(6,300)
dimensional
vector

1st layer
of NN.
$w^{[1]}, b^{[1]}$

softmax
$\bar{e}$ its own
parameters $w^{[2]}, b^{[2]}$
(It classifies
among the 10,000
possible o/p in vocab
to get the final word
we are trying to predict)

Note  A history hyperparameter can be used. It says to predict juice & just need the last 4 words of the sequence and not the entire sequence, thus reducing the i/p size. Hence, arbitrarily long sequences can be dealt easily as i/p size is always fixed.

Perform backprop to apply gradient decent to learn decent word embeddings in this algorithm.

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

content          target word.      content

So, 1st 4 word and last 4 words are fed into the NN to learn to predict the target word juice.
1 word can also be used while feeding content into NN, this will help train the model in predicting nearby words.
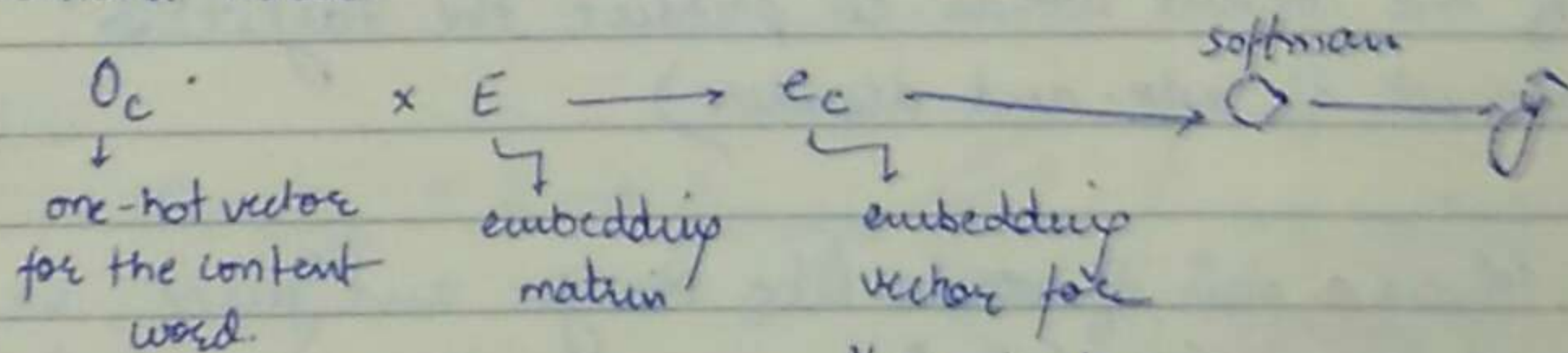All this helps us in learning good word embeddings.

➡️ [Word to Vec Algo] (simple & more efficient to learn embeddings)

- Skip gram model : Here, unlike the previous algo, the content is not taken to be the immediate words lying before or after the target word but any random word can be chosen as the content for the target word. e.g.

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

| Content | Target | |
|---------|--------|---|
| orange | Juice | {word chosen within a ± word window} |
| orange | glass | {  "  "  "  "  ±2  "  "} |
| orange | my | {  "  "  "  "  ±5  "  "} |

So, the aim of this supervised learning problem is set to predict a target word for the content within a given window size around the content word.

$O_c$      × E  ⟶  $e_c$  ⟶  softmax ◯ ⟶ $\hat{y}$

one-hot vector for the content word.     embedding matrix     embedding vector for

i/p content c.

$$\text{Softmax}: p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10K} e^{\theta_j^T e_c}}$$

target word   content word
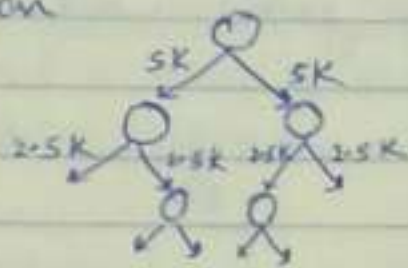
$\theta_t$ : parameter associated with output t.

$$L(\hat{y}, y) = -\sum_{i=1}^{10k} y_i \log \hat{y}_i \qquad \hat{y}, y = \text{one hot vector}.$$

Skip gram model called so boz it takes as i/p one word like orange and then tries to predict some words skipping a few words from the left or the right

Disadv: computational speed boz of $\sum_{j=1}^{10k} e^{\theta_j^T e_c}$ term in denominator.

Solu: Heirarchial softmax

classifier

$\downarrow$

(it is generally designed such that more common words are at the top & less common are buried much deeper in the tree)

this is classifier to classify y the word is in the 1st 5000 words or the 2nd 5k words and so on i.e binary classification.

This solves the problem of summing over entire vocab (size in the denominator.

Note we remove frequently occurring words like 'of','a','the','and','to','etc. boz otherwise much effort will be wasted in updating their $e_c$, so as to focus on the less common but important words.
(under skip-gram)?

→ [Negative Sampling] : (modified algo similar to skip gram but more efficient)

Note CBow: Continuous Backward model (it considers the surrounding words of the content words to predict the target & has its own set of adv. and disadv.)

learning problem: Given a pair of words like orange and juice, we need to predict is this a content-target pair?
in this e.g. orange juice is a +ve example.
but orange and king is a -ve example.

| content | word | target |
|---------|------|--------|
| orange | juice | 1 |
| orange | king | 0 |
| orange | book | 0 |
| orange | the | 0 |
| orange | of | 0 |

$\underbrace{\quad}_{t_c} \quad \underbrace{\quad}_{t_t} \quad \underbrace{\quad}_{y} \;\; \text{training set}$

to create such pairs, for +ve, target word is taken from a window of some size around the content word while for -ve example, target word is taken at random from the dictionary.

After taking 1 +ve e.g. take k random words from the dictionary for the same content word and label them 0. It's ok if by chance that random word appears to be in the window around the content word. Now this table created will act as X and Y for the supervised learning algorithm. Output is Y i.e 1 or 0 given a pair of words, i.e the algorithm predicts if it thinks the two i/p words were obtained by sampling two words close to each other or not.

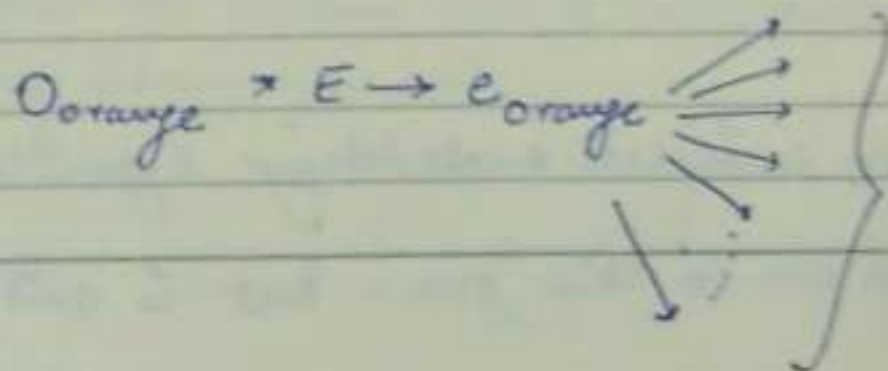k is 2-5 for larger datasets and 5-20 for smaller datasets.

Softmax: 
$$P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10k} e^{\theta_j^T e_c}}$$

→ embedding vector for each possible content word.

→ parameter $\theta$ for each possible target word.

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

NN:

$O_{orange} * E \to e_{orange}$ 

10k logistic regression classification problems.

but we will train the model for only k+1 of them.

Hence, this way instead of having one giant 10,000 way softmax, which is very expensive to compute, we've instead turned it into 10K binary classificaⁿ problems. each of which is quite cheap to compute and at every "itereⁿ" we'll train only 5 i.e K+1 of them (i.e k -ve and 1 +ve enample) hence further reducing the computaⁿ cost.

Hence, is called **Negative sampling**. But how to choose the -ve enamples ?

→ We can sample the words for a given content word depending on the empirical frequency i.e how often different words appear in the vocabulary, or use $\frac{1}{\text{Vocab size}}$ , sample the -ve enamples uniformly at random but that is also very non-represen tative of the distribuⁿ of English words.

So, the middle path of the above two solutions :

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10K} f(w_j)^{3/4}}$$

$f(w_i)$ is observed frequency of the word $w_i$ in English language or in the training corpus.

Note,
Hence, this approach is somewhere b/w the extreme of taking uniform distribuⁿ and the other extreme of just taking whatever was the observed distribution in the training set ( bcoz this approach sampled the, &, and, etc which are actually not unimportant words).

→ | **GloVe word vectors** | (even simpler embedding algorithm)

Not used as much as word to vec or skip gram but is quite simple.

GloVe : Global vectors.

$X_{ij}$ = no-of times i appears in content of j $\approx$ $X_{ji}$ {acc-to defini

$\uparrow$ target $\qquad$ $\uparrow$ context

- Count that captures how often do words i and j appear z each other, or close to each other.

if content and target are being select as + nt within 10 or so word limit of each other

But $X_{ij} \neq X_{ji}$ if target lies before content always or something like that.

## GloVe model:

$$\text{minimize} \quad \sum_{i=1}^{10k} \sum_{j=1}^{10k} f(X_{ij})\left(\theta_i^T c_j + b_i + b_j' - \log X_{ij}\right)^2 \qquad —\text{①}$$

$\underbrace{}_{\text{weighting term}}$

$\theta^T c_j \equiv (A\theta_i)^T (A^{-T} c_j)$

(it does not give too much weight to stopwords and gives significant weight to the less frequent but imp words)

$f(X_{ij}) = 0$ if $X_{ij} = 0$ $\qquad$ {as log 0 is undefined}

$$e_w^{(final)} = \frac{e_w + \theta_w}{2} \qquad \text{\{as } \theta_i \text{ and } c_j \text{ are symmetric\}}$$

In ①, such kind of algo, it can't be guaranteed that the axis used to represent the features will be well-aligned with what might be easily humanly interpretable axis.

→ Sentiment classification (Application of these algorithms)

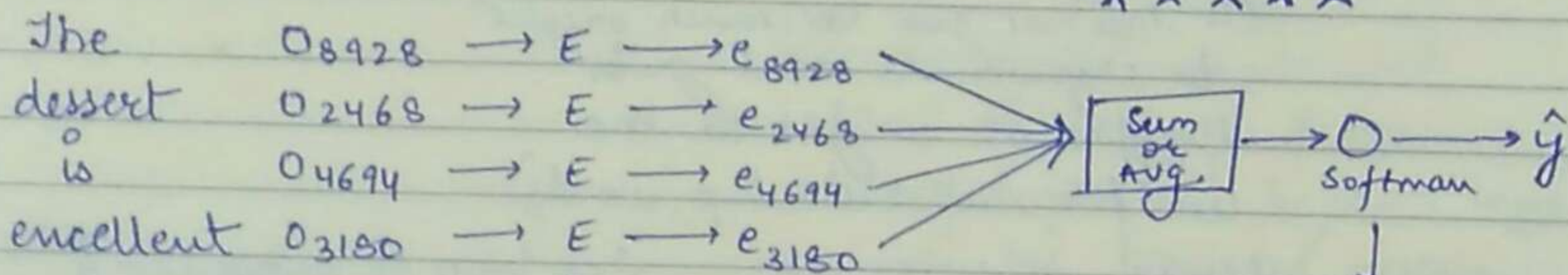| x (i/p text) | y (sentiment) |
|---|---|
| The dessert is excellent. | ★ ★ ★ ★ ☆ |
| Service was quite slow. | ★ ★ ☆ ☆ ☆ |
| Good for a quick meal, but nothing special. | ★ ★ ★ ☆ ☆ |
| Completely lacking in good taste, good service, and good ambiance. | ★ ☆ ☆ ☆ ☆ |

training dataset to figure out the sentiment of a review.

→ labelled data is not always available for training. But with word embeddings, good sentiment classification model can be built just by using only modest-size label training set.
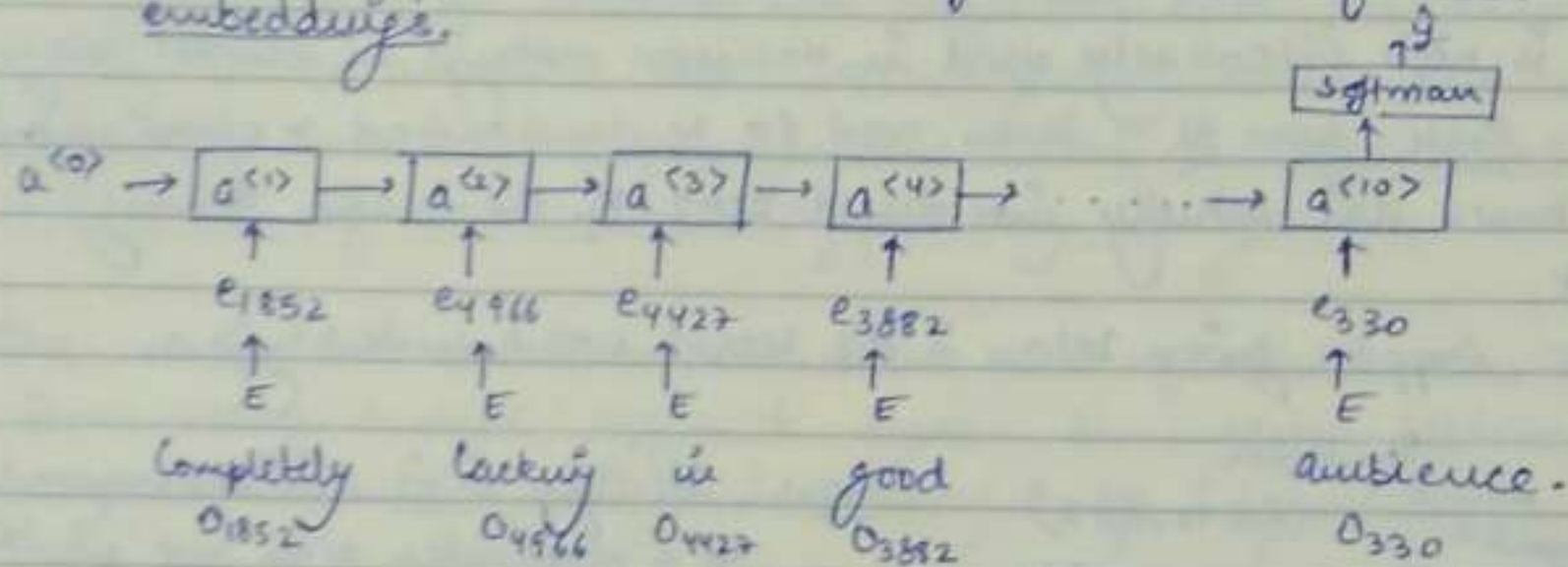
e.g.   The dessert is excellent.

$$y$$
★ ★ ★ ★ ☆

The        $O_{8928} \rightarrow E \rightarrow e_{8928}$
dessert   $O_{2468} \rightarrow E \rightarrow e_{2468}$
is           $O_{4694} \rightarrow E \rightarrow e_{4694}$
excellent $O_{3180} \rightarrow E \rightarrow e_{3180}$

→ [Sum or Avg.] → O softmax → $\hat{y}$

already learnt from a large training corpus, different from this problem's labelled training set, thus helps in using external info also.

to calculate the probability of getting the different ratings from 1 to 5, hence giving o/p $\hat{y}$.

So, this algo works for reviews of any word length as it averages up all the embedding vectors of the words. But this algo ignored word order, and simply took the average.

Disadv: In the review "Completely lacking in good taste, good service and good ambience" is a negative review but bcoz 'good' appears too many times, this algorithms will probably classify this review as good even though it is actually harsh. This was bcoz word order was not taken into account.

Sol^n: Use RNN for sentiment classifica^n instead of word embeddings.



Many to one RNN architecture

So, considering the order of sequence, this algo will help classify this review as -ve.

Training this algo gives a pretty decent sentiment classifica^n algo and also the use of E (learnt from a large corpus of text from somewhere else) helps the model take decision even on such words which are not +nt in the labelled training dataset.

⇒ Debiasing word embeddings

- means removing gender, ethnicity, sexual orientation bias, and not the technical bias used in ML.
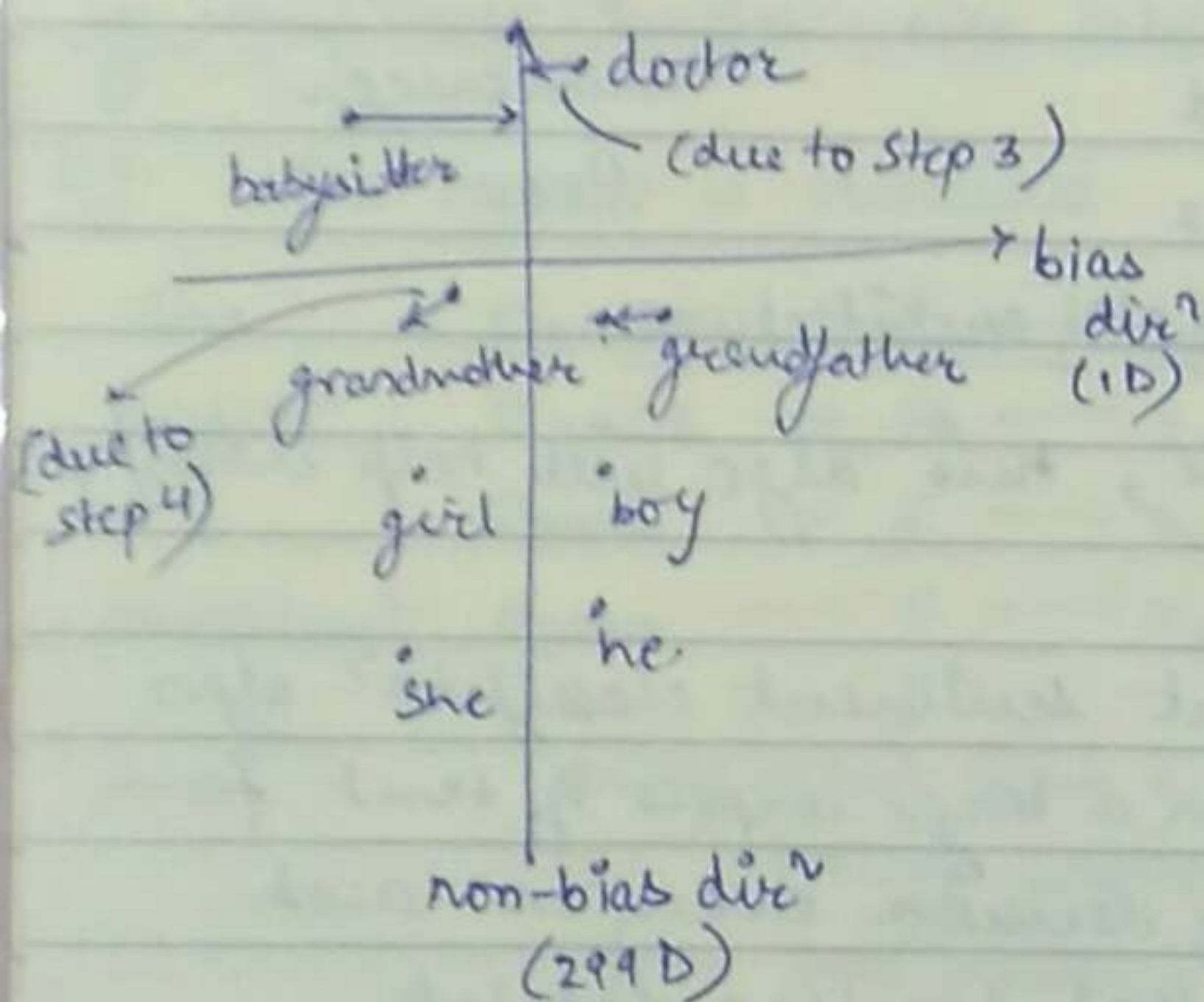
eg Man: Woman :: King: Queen

but | Man : Computer programmer :: Woman : ? |

It was seen that on learning from word embedding the result came out as 'homemaker' and that enforces a unhealthy gender stereotype and its wrong. An unbiased system would have outputted 'computer programmer' as answer.

Also, ~~DoFather~~ Father : Doctor :: Mother : Nurse were observed.

As ML is now extensively used in decision making in almost every field, these types of biases need to be diminished or eliminated. These biases are actually introduced bcoz of the biased training dataset.

Soln: Suppose given below is the learnt word embedding.



**Step 1:** Identify bias dirⁿ corresponding to a particular bias we want to reduce or eliminate. Take for example gender bias.

**Step 2:** $e_{he} - e_{she}$
$e_{male} - e_{female}$ } take their
⋮
aug or use SVD (singular value decomposition)

aug ⇒ horizontal dirⁿ is the bias direction.

**Step 3:** Neutralizeⁿ step - for every word that is not <u>definitional</u>, project to get rid of bias.

Note Neutralizeⁿ is done to eliminate their component in the bias direction. eg. of gender neutral words like doctor, nurse, etc.

(means words which do not represent a gender, eg. he, she, father, lady, etc)

Step 4: Equalization pairs - make sure that pairs like grandfather & grandmother are both exactly the same similarity or exactly the same distance from words that should be gender neutral, such as doctor or babysitter. This happens using linear algebra that will move grandmother & grandfather to a pair of points that are equidistant from the axis in the middle and its effect is that now the distance b/w babysitter and grandfather is same as b/w babysitter and grandmother unlike in the previous case where grandmothers end up babysitting more than grandfathers. This is done for all such pair of words which differ only in gender.

The words that need to be neutralized and equalized can be learn easily through different algos. In English dictionary, very few words are gender specific by definition.

* Sequence models & Attention mechanisms :
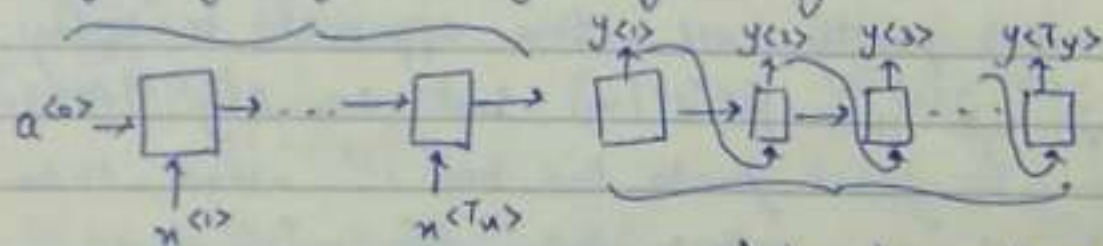
from machine transle? to speech recogni?.

① Sequence to Sequence model -
$$x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad x^{<4>} \quad x^{<5>}$$

eg  Jane visite l'Afrique en septembre.
→ Jane is visiting Africa in September.
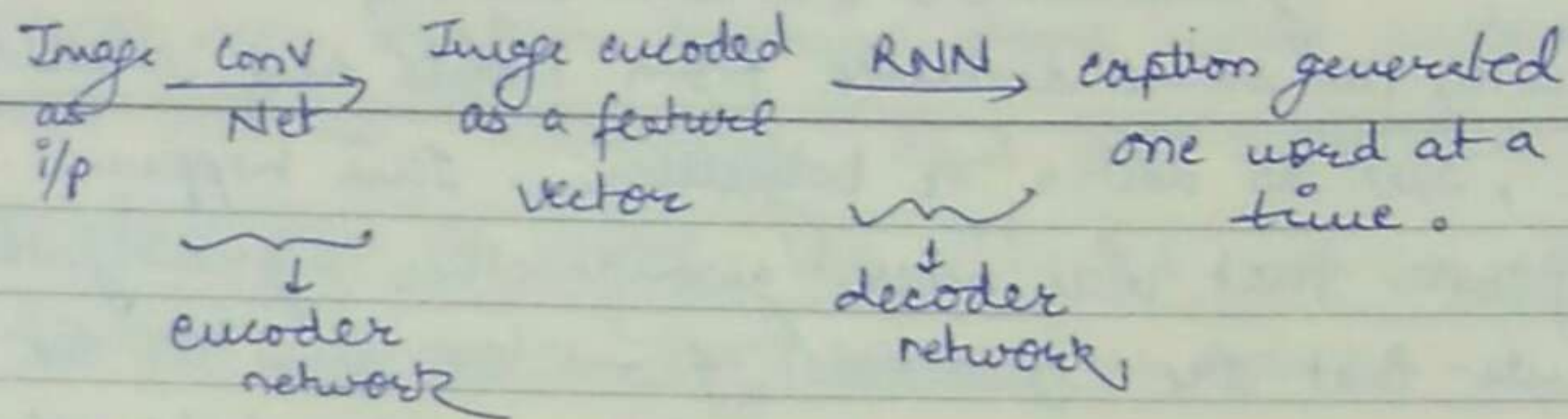$$y^{<1>} \quad y^{<2>} \quad y^{<3>} \quad y^{<4>} \quad y^{<5>} \quad \hat{y}^{<6>}$$



'encoder network'
(built as RNN)

'decoder network'

This model works well for given a enough pairs of French & English sentences.

eg. Image captioning

Image    Conv    Image encoded    RNN    caption generated
as    ──────→    as a feature    ──────→    one used at a
i/p    Net    vector    ⌣    time.
    ⌣         ⌣
    ↓              ↓
  encoder          decoder
  network          network

→ this works well especially if the caption is not too long.

The difference b/w this sequence to sequence model and the earlier tent generator model is that now we want the most likely sentence as o/p (like in transl^a or captioning) unlike the tent generator case of random o/p which were just gramatically correct.

The similarity are there in language model & Machine translation model except that instead of always starting along with the vector of all zeros, machine transl^a has an encoded network, that figures out some representa^n for the input sentence. to start off the decoded network. Hence, it is called 'conditional language model'. Now instead of modelling the probability of any sentence, it is now modelling the probability of an English translation as o/p given a French i/p. And the one with max probability is selected as the o/p.
For this Beam Search is used.

'Greedy Search' doesn't work pretty well. this approach basically doesn't focus on $P(\hat{y}^{<1>}, \hat{y}^{<2>}, \ldots \hat{y}^{<T_y>} | x)$ but on $P(\hat{y}^{<1>} | x)$ ie best individual words. It first finds the best 1st word, then the best 2nd word and so on rather than finding the probability of best overall sentence. e.g. in Greedy search:
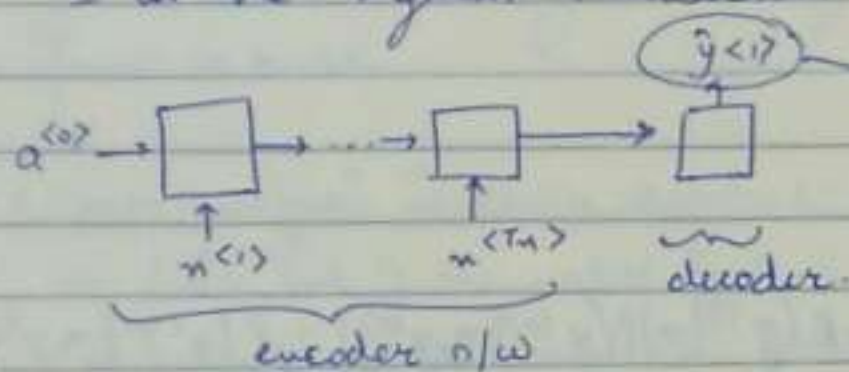
P(Jane is visiting Africa in September) > P(Jane is going to be visiting Africa in September.)

bcoz Occurrence of going > visiting in the vocabulary, and hence bcoz of the greedy approach, 2nd statement might be selected as the o/p which is good but not as good as the 1st one.

→ **Beam Search Algorithm:** (approx/heuristic search algo)

B = beam width (no. of guesses per word made by this algo)

It works similar to greedy but instead of finding just the best word, it finds 'B' words from the vocabulary for posi^n 1 in the English translation.
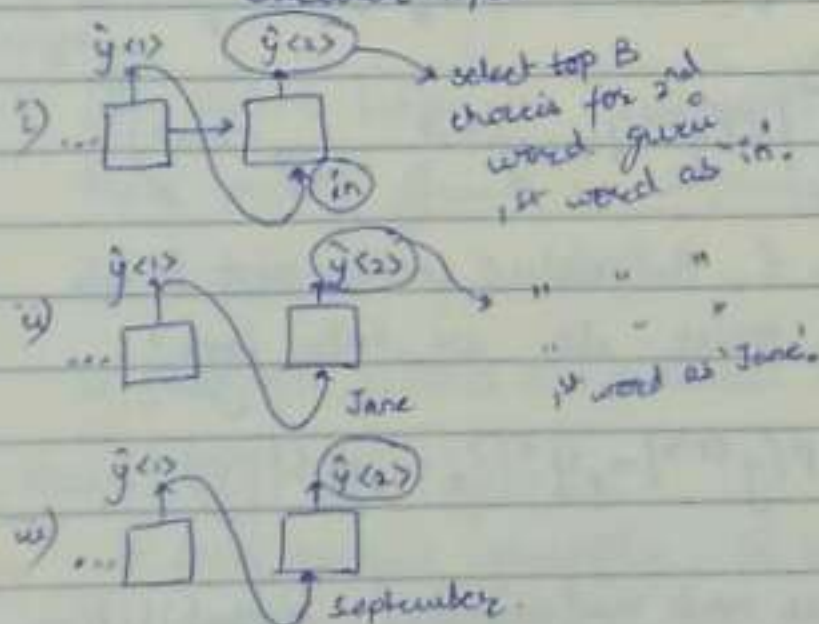


**Step1:**
Softmax o/p with overall 10,000 (vocab size) possibilities but only B of them which are the most probable are kept in memory.

encoder n/w

**Step2:**
For each of the 'B' choices for the 1st word, now the 2nd word is guessed, and this time not only $P(\hat{y}_2)$ is maximised but $P(\hat{y}_1, \hat{y}_2 | x)$ is maximised.
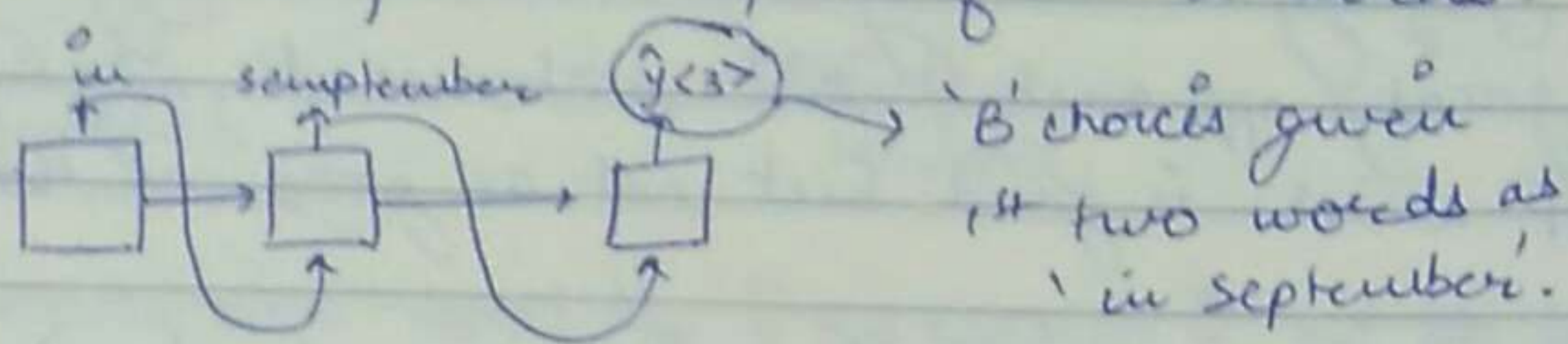
$$P(\hat{y}_1, \hat{y}_2 | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>})$$

Now the top 'B' of these possible combin^s of 1st 2 words is stored in memory.

→ select top B choices for 2nd word given 1st word as 'in'.

Jane

September.

3] the B choices for 1st word were in, Jane and September.

**Step 3:** If the top B choices for 1st two words were 'jane is',
'in september', 'jane visits', then the same procedure
as step 2 applied to predict the best possible 3 choices for
the 3rd word for each pair of 1st 2 words.



→ 'B' choices given
1st two words as
'in september'.

$$P(\hat{y}^{<3>} \mid x, \text{ "in september"})$$

and now similar to step 2 $P(\hat{y}^{<1>}, \hat{y}^{<2>}, \hat{y}^{<3>} \mid x)$ is maximised.
and 'B' choices (top) for the 1st 3 words are stored in memory.

**Note** Beam search with $B = 1 \Rightarrow$ greedy search algorithm.

• Refinements to beam search:

① Length Normalize^ –

In beam search, $P(y^{<1>}, y^{<2>} \dots y^{<T>} \mid x) = P(y^{<1>} \mid x) P(y^{<2>} \mid x, y^{<1>}) \dots P(y^{<T_y>} \mid x, y^{<1>}, \dots y^{<T_y-1>})$

$$\text{Beam search} = \underset{y}{\text{argmax}} \prod_{t=1}^{T_y} P(y^{<t>} \mid x, y^{<1>}, \dots, y^{<t-1>})$$

bcoz all probabilities are less than 1, the value comes out to be
very small and round off errors occur. So, we take log:

$$= \underset{y}{\text{argmax}} \sum_{y=1}^{T_y} \log P(y^{<t>} \mid x, y^{<1>}, \dots, y^{<t-1>})$$

so that rounding errors are reduced as now individual probabilities
are added and not multiplied.

Drawback is that even this formula works well for only short sentences,
so normalize^ is done. As all probabilities are less than 1, their log is -ve

and so including more terms (i.e longer sentences) makes it even more -ve. So,

$$\boxed{\frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{<t>} \mid x, y^{<1>}, \ldots, y^{<t-1>})}$$

$\alpha \approx 0.7$

in b/w full and zero normali- (hyperparameter) $-za$.

$\alpha$ can be varied acc. to what value gives the best results.

— How to decide B (beam length)?
Larger B means more possibilities and hence better results but higher computational overhead and memory requirement & hence slow.
there is high gain in results in going B from 1 to 3 to 10 to 100 but the gains are not that big when B goes from 1K to 3K

Note: Unlike BFS, DFS; beam search runs faster but is not guaranteed to find exact max for any max $P(y \mid x)$

• Error analysis in beam search:
to find out if the value of B is the problem or the RNN model used should be redesigned for better performance.

Jane visite l'Afrique en septembre.
Jane visits Africa in September. $(y^*)$ — Human translan
Jane visited " last " $(\hat{y})$ — Algorithm transla?

to find out where the problem lies, we need to know $P(\hat{y} \mid x)$ & $P(y^* \mid x)$ using RNN model.

Case 1: $P(y^* \mid x) > P(\hat{y} \mid x)$
⇒ RNN guessed correctly, so beam search is at fault bcz of which it didn't chose max P o/p & o/p was wrong.

Case 2: $P(y^* \mid x) < P(\hat{y} \mid x)$
⇒ RNN guess was wrong and it was not bcz of beam search that $\hat{y}$ was obtained as o/p.

Going through various such examples, we find out how many times beam search is at fault and how many times RNN is at fault.

Note RNN generates the objective function that beam search is supposed to maximise.
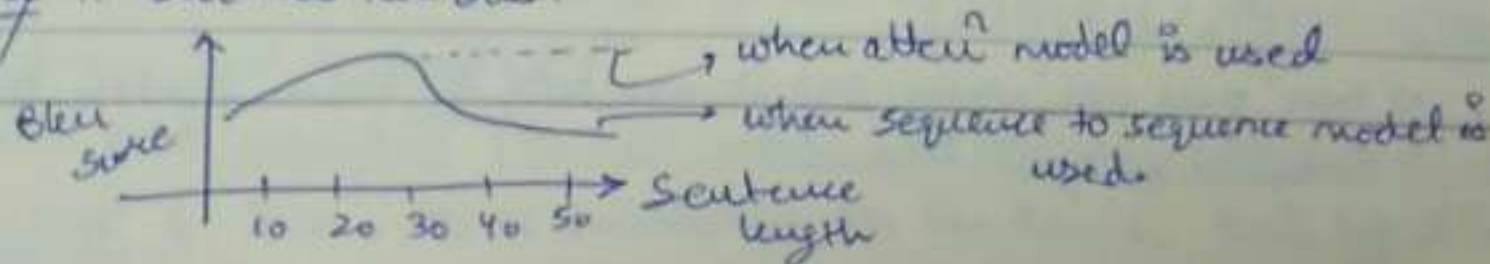
If RNN is found at fault, a deeper layer of analysis can be added, and regularizaⁿ , more training data or modifying N/w architecture can be done. Instead if beam search is found to be at fault, try Z greater B.

→ **BLEU score** (to help with the problem of having multiple correct
bilingual                      o/p or English translaⁿ to a given french sentence.)
evaluation

It evaluates the o/p of a machine translation systems instead of a human doing so. This is done by looking at each of the words in the o/p and see if it appears in the references given by humans (as test set) This is called **precision** of the machine translaⁿ output.

② **Attention model** —

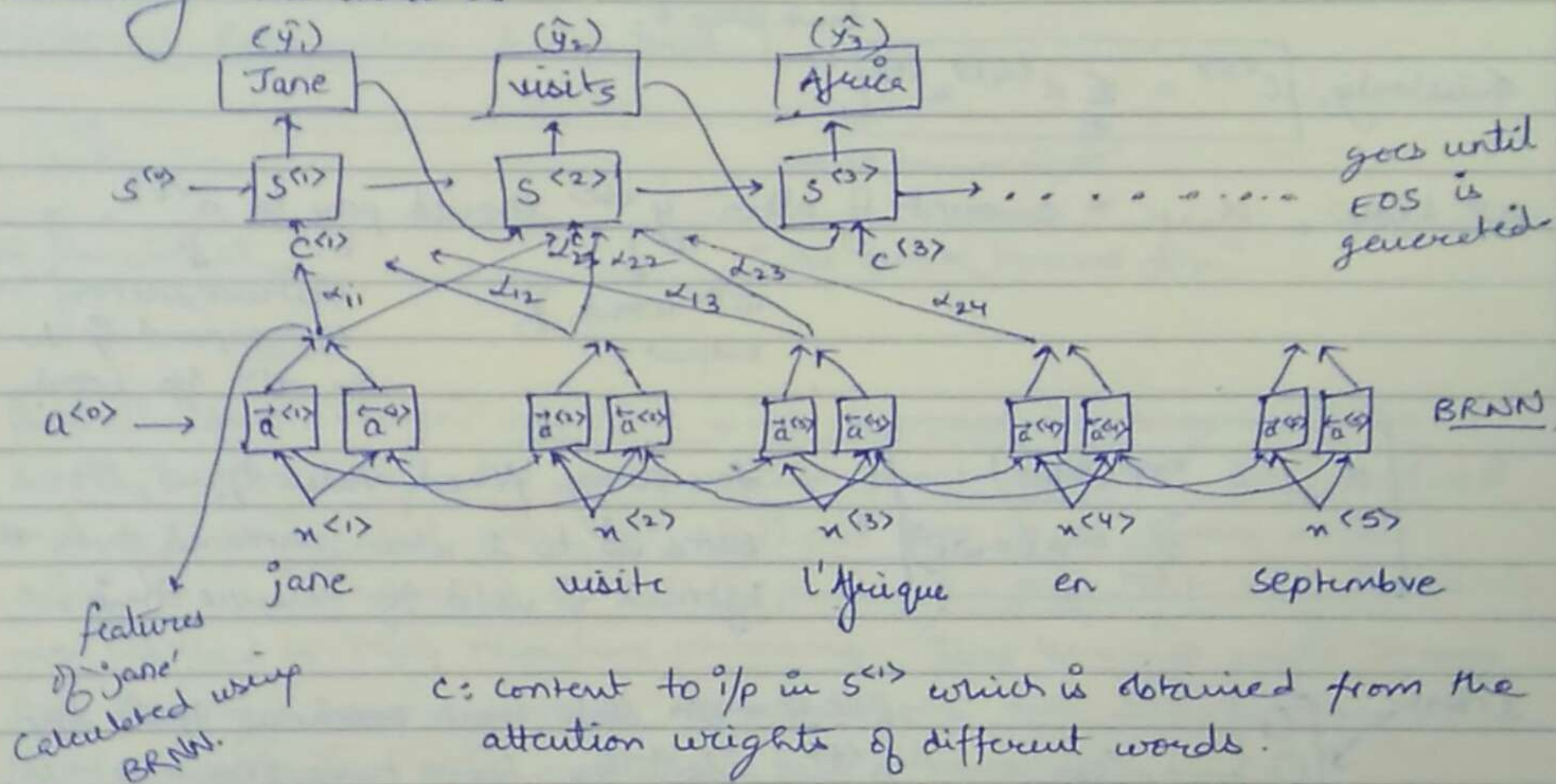The encoder-decoder mechanism works well for short sentences, as it involves memorizing the whole sentence, encoding it in vector form and then decoding it to give o/p. But in case of long sentences, memorizing the whole sentence at once and then translating it becomes tedious.



, when attenⁿ model is used

when sequence to sequence model is used.

BLEU score (y-axis) vs Sentence length (x-axis: 10 20 30 40 50)

for very short sentences, bleu score is less bcoz its difficult to predict the exact words and for very long sentences, bleu score is less bcoz its difficult to memorize them.

Hence, comes Attention models that work like human beings i.e reading the sentence part by part and hence able to translate long sentences.

Attention model actually uses <u>attention weights</u> $(\underline{\alpha_{ij}})$ where $\alpha_{ij}$ means how much attention is to be paid to word $j$ while predicting the word $i$.



features of 'jane' calculated using BRNN.

$c$: content to i/p in $s^{<i>}$ which is obtained from the attention weights of different words.

The fwd activa<sup>n</sup> $\underline{\vec{a}^{<t>}}$, the bwd activa<sup>n</sup> $\underline{\overleftarrow{a}^{<t>}}$, $\underline{s^{<t-1>}}$ <u>contribute</u> <u>in deciding</u> $\alpha$ wt i.e <u>attn<sup>n</sup> weight</u> of <u>word t</u> <u>in predicting word</u>.

This allows the model to look only within a local window of the French sentence to predict a particular word of the English translation.

$a^{\langle t \rangle} = (\vec{a}^{\langle t \rangle}, \overleftarrow{a}^{\langle t \rangle}) =$ feature vector for time step $t$.

$\sum_t \alpha^{\langle 1, t \rangle} = 1$    (None of the $\alpha$ is -ve)

$$c^{\langle 1 \rangle} = \sum_t \underset{\underset{\substack{\text{attention} \\ \text{weights}}}{\uparrow}}{\alpha^{\langle 1, t \rangle}} \, \underset{\underset{\substack{\text{feature vector} \\ \text{corresponding} \\ \text{to i/p word at} \\ \text{time step } t.}}{\uparrow}}{a^{\langle t \rangle}}$$

Similarly, $\boxed{c^{\langle 2 \rangle} = \sum_t \alpha^{\langle 2, t \rangle} a^{\langle t \rangle}}$

we know, $\alpha_{wt} =$ amount of atten$^n$ $\underset{\underset{\substack{w^{th} \text{ word of} \\ \text{output.}}}{\downarrow}}{y^{\langle w \rangle}}$ should pay to $\underset{\underset{\substack{\text{feature vector} \\ \text{corresponding to} \\ \text{the } t^{th} \text{ i/p word.}}}{\downarrow}}{a^{\langle t \rangle}}$.

$$\boxed{\alpha_{wt} = \frac{\exp(e_{wt})}{\sum_{t=1}^{T_x} \exp(e_{wt})}}$$

for every fixed value of $w$, it sums up to 1 when summed over $t$. softmax is used to ensure this.

$s^{\langle w-1 \rangle}$  $\rightarrow e_{wt}$

$a^{\langle t \rangle}$

Train this NN to get $e_{wt}$ corresponding to diff $s$ and $t$, using gradient descent. It is found that it works pretty well in deciding the correct atten$^n$ weights.

$s^{\langle w-1 \rangle}$: NN state from previous time step

$a^{\langle t \rangle}$: features from time step $t$ of i/p.

Obviously, $\boxed{e_{wt}}$ depends on $s^{\langle w-1 \rangle}$ & $a^{\langle t \rangle}$.

$\boxed{\text{Note} \quad \begin{array}{l} \text{total no. of atten}^n \\ \text{parameters} = T_x T_y \end{array}}$

(so, quadratic cost of this algo)

$\underline{a^{<t>}} = (\vec{a}^{<t>}, \overleftarrow{a}^{<t>}) \times$ feature vector for time step $t$.

$$\sum_t \alpha^{<1,t>} = 1 \qquad (\text{None of the } \alpha \text{ is } -ve)$$

$$c^{<1>} = \sum_t \underbrace{\alpha^{<1,t>}}_{\substack{\text{attention} \\ \text{weights}}} \underbrace{a^{<t>}}_{\substack{\text{feature vector} \\ \text{corresponding} \\ \text{to i/p word at} \\ \text{time step } t.}}$$

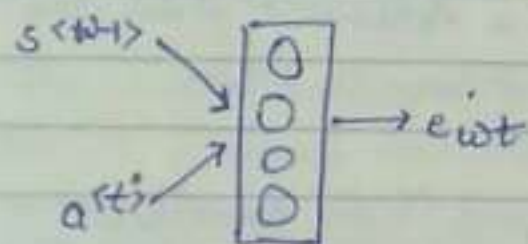Similarly, $\boxed{c^{<2>} = \sum_t \alpha^{<2,t>} a^{<t>}}$

we know, $\alpha_{wt} = $ amount of atten$^n$ $y^{<w>}$ should pay to $a^{<t>}$.

$w^{th}$ word of output.

feature vector corresponding to the $t^{th}$ i/p word.

$$\boxed{\alpha_{wt} = \frac{\exp(e_{wt})}{\sum\limits_{t=1}^{T_x} \exp(e_{wt})}}$$

for every fixed value of $w$, it sums up to 1 when summed over $t$. softmax is used to ensure this.

$s^{<w-1>}$

$a^{<t>}$

$\longrightarrow e_{wt}$

Train this NN to get $e_{wt}$ corresponding to diff $s$ and $t$, using gradient decent. It is found that it works pretty well in deciding the correct atten$^n$ weights.

$s^{<w-1>}$: NN state from previous time step
$a^{<t>}$: features from time step $t$ of i/p.

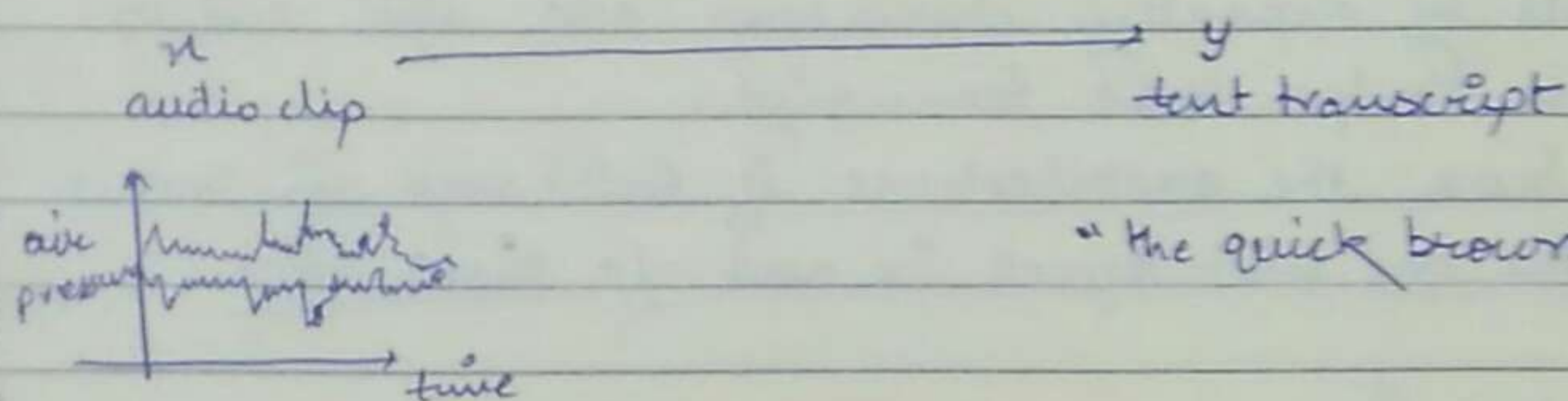Obviously, $\boxed{e_{wt}}$ depends on $s^{<w-1>}$ & $a^{<t>}$.

Note $\boxed{\text{total no. of atten}^n \text{ parameters} = T_x T_y}$

(so, quadratic cost of this algo)

- Attenⁿ models are also applied in __image captioning__, using the same kind of architecture as just discussed i.e "paying attenⁿ to only some parts of picture ~~to the~~ at a time while predicting the caption.

- It is also used to get __normalized dates__.

  __e.g.__  July 20th 1969 $\longrightarrow$ 1969-07-20

  23 April, 1564 $\longrightarrow$ 1564-04-23

$\Rightarrow$ | Speech Recognition - Audio Data. |

$x$ $\xrightarrow{\hspace{5cm}}$ $y$
audio clip                                    text transcript

air
pressure $\uparrow$ ⟨waveform⟩                    "the quick brown fox".
$\xrightarrow{}$ time

Initially speech recogniⁿ was done using __phonemes__ i.e recognizing the sound and outputting the corresponding phonemes to it. It involved hand engineering but were not efficient bcoz if a person said quick, it was outputted as kwik (phonem). Different kind of sounds were labelled to their respective phonems. Deep learning made it easy thus removing hand-engineered representations and simply training the n/w with audio clips and their text transcripts.

- __Attention Models__ can be used for speech recogniⁿ.

- __CTC cost for speech recognition__ is another method.
connectionist temporal classification

__Note:__ In speech recognition the number of output time steps is much lesser than the input time steps.

e.g. If in a 10sec video, features come at 100 samples per second, then a 10 sec audio will have 1K features on i/p. But the o/p might not have a thousand alphabets / characters. So, what CTC does is it allows the RNN to generate an output like:

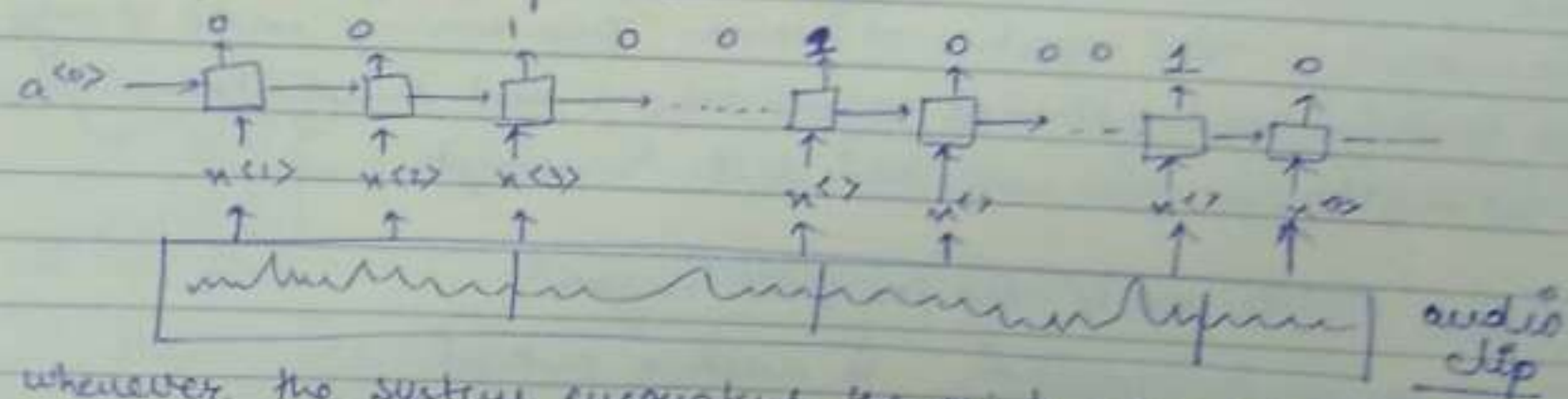$$(ttt\_h\_eee\_\_\_\text{---}\_\_\text{---}qqq\text{---}) \text{ and its considered}$$

blank ↓
space

the correct output for the audio i/p "the quick". the rule of CTC is to collapse output characters not separated by "blanks". so the sequence is collapsed into "the q" which allow the NN to have 1000 outputs by repeating characters and still end up with a much shorter output text transcript.

CTC was needed bcoz the architecture of RNN used in speech recognition was the one with equal i/p and o/p time steps.

→ **Trigger Word detection :**

Trigger word is actually like "Hey Siri" for Siri, "Ok Google" for Google, 'Mahinir' for me and likewise. this trigger word helps initiate a process.



whenever the system encounters the points where someone has just finished saying the trigger word, then in the training set target labels are set to be zero for everything before that point and 1 right at that point. later on if trigger word is said again,

then target label is again set to 1 right after the o's. The only disadvantage is that it creates an unbalance with lot of zeroes and very few 1s. So instead of setting just a single time step equal to 1, we can output 1 for a time dure^ before reverting back to zero. This evens out the ratio of 1s to 0s slightly to some extent.

DOUBTS: