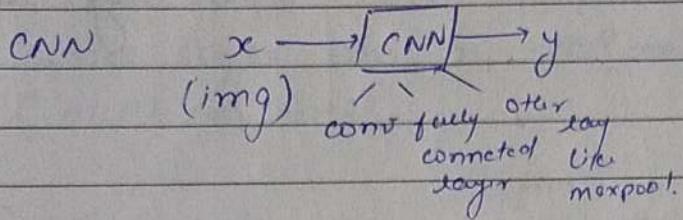
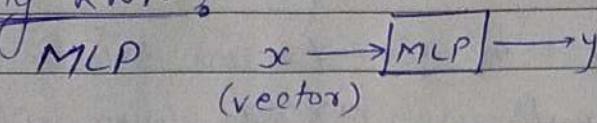


* Recurrent Neural Networks (RNN)

* Why RNN?



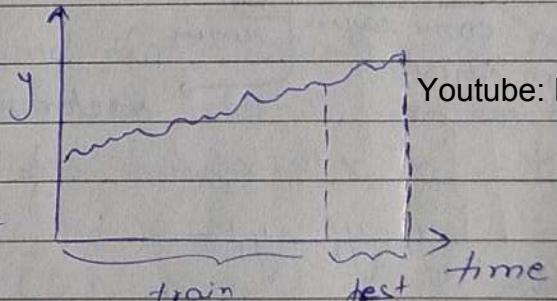
* x is a Sequence we use RNN.

Sequence like amazon reviews (sequence of words)

NOTE:- When we use Bow, TfIdf, w2v we don't care/ completely discard the sequence of the words in a sentence. But Sequence/ orientation of the words are important.

* Time-Series data

- windowed data
- fourier transform to decompose in constituent frequencies.



* Language translation

English to Italian

* Speech recognition

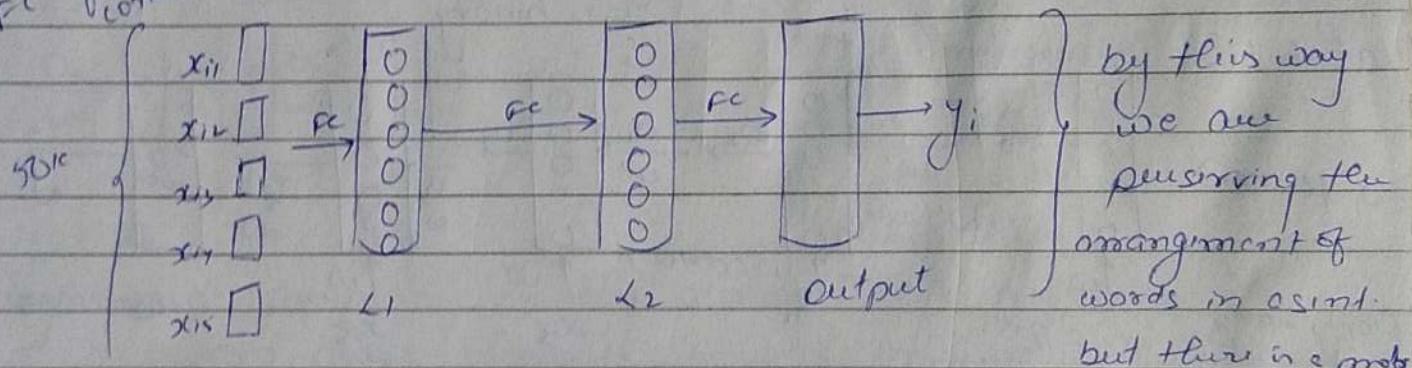
audio input \rightarrow text

* Input image \rightarrow output caption

* Core Idea: If output depends on the seq. or arrangement of the inputs we need new types of NN which retains and leverages the seq. info.

④ If we try to train MLP for amazon review polarity.

$\delta_1 \rightarrow x_i \rightarrow$ the phone is very fast
 (at 10k unique words (BOW))
 $x_{i1} \downarrow x_{i2} \downarrow x_{i3} \downarrow x_{i4} \downarrow x_{i5} \downarrow$



→ set them up 2 reviews Youtube: Programming Cradle

5 words $x_1 \rightarrow$ This phone is very fast $\rightarrow 50k$
 4 words $x_2 \rightarrow$ This phone is good. $\rightarrow 40k$
 :
 20 words $\rightarrow x_i \rightarrow 200k$
}

Sentences can be of
 diff lengths hence
 input ~~vector~~ size changes

→ one solⁿ to above prob can be we can fix the input size to ~~max~~ length of max length string let say 50 so for ~~word~~ sent. with 5 word will have input filled till 50K and rest 450K will be zero vectors (zero padding)

- Now what if in test date we have sent. of length 1000000
- Also we are getting huge vectors.

So please can the problems of using MLP for text data

* Recurrent Neural Network (RNN)

↳ Repeating.

eg) Amazon food review → binary classify $\{0, 1\}$ {size of vocab is d'}

$O = \{x_i, y_i\}$

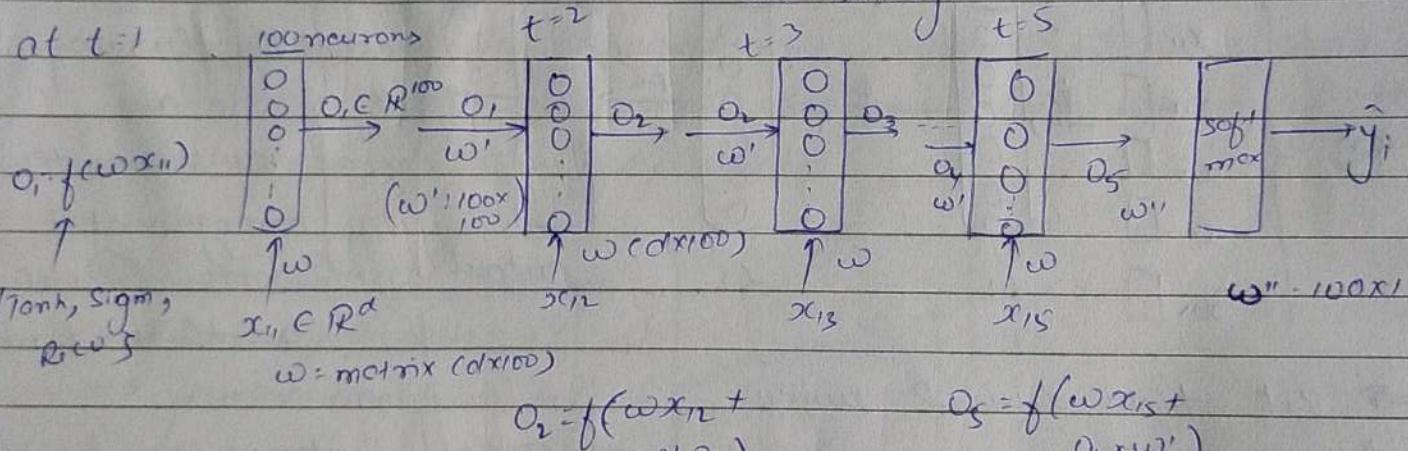
↳ binary
Sentences
(Seq. of words)

$y_i \leftarrow x_i \rightarrow (x_{11}, x_{12}, x_{13}, x_{14}, x_{15})$ words: one hot encoded
 $x_i \rightarrow x_{11}, x_{12}, x_{13}, \dots, x_m$
 m : length of sentence

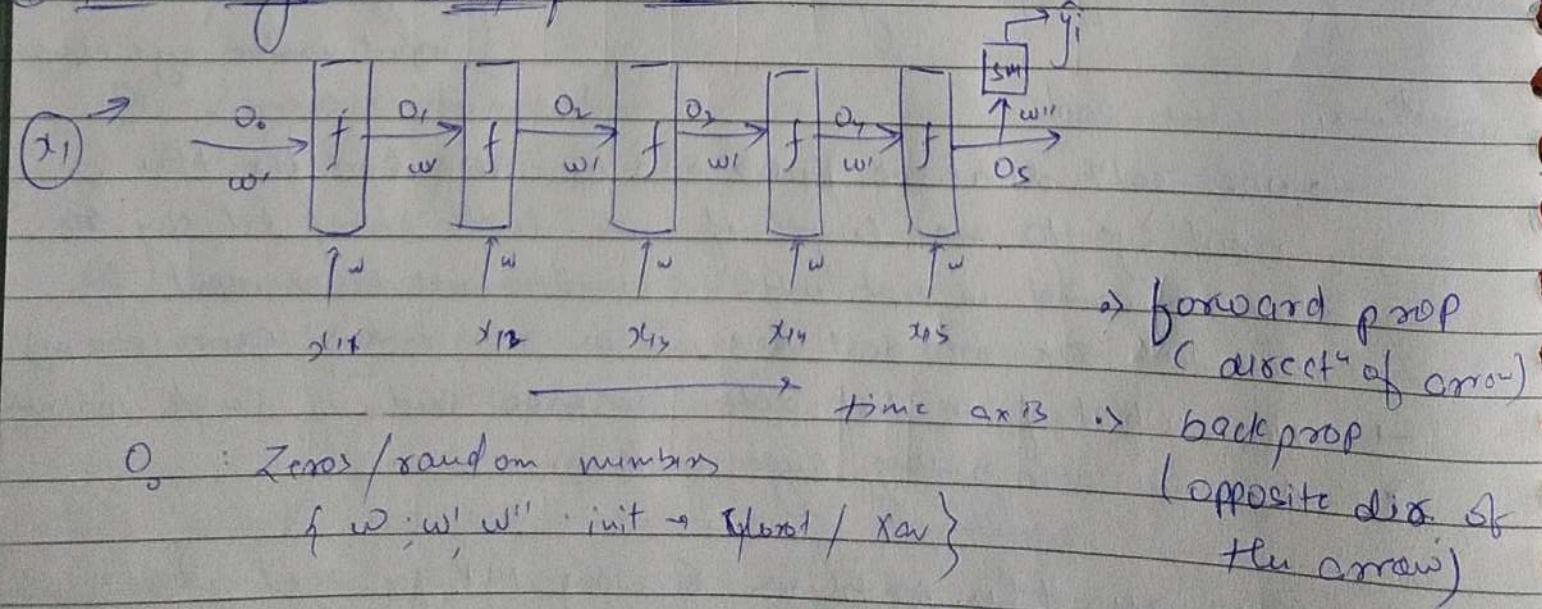
Task: Seq. of words given; predict the polarity.

$$x_i \rightarrow \{x_{11}, x_{12}, x_{13}, x_{14}, x_{15}\} \quad y_i \quad \{0, 1\}$$

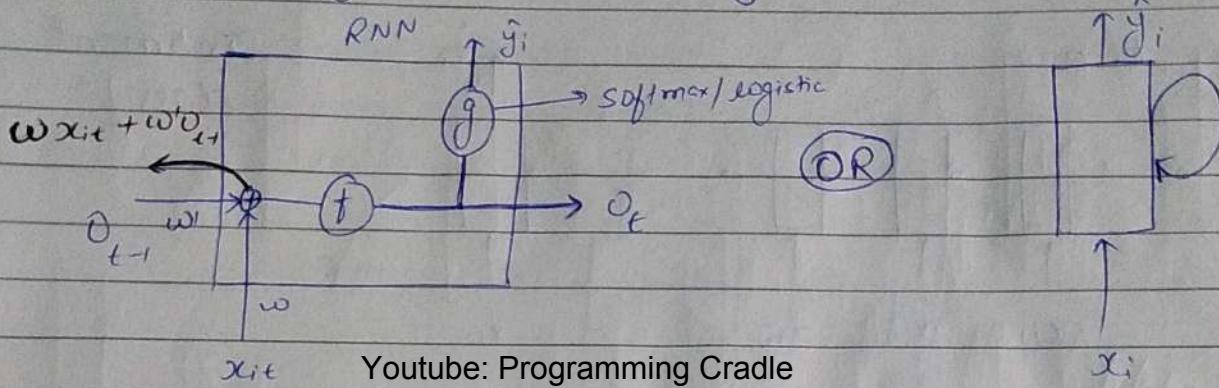
↳ binary



* Above diagram Simplified below

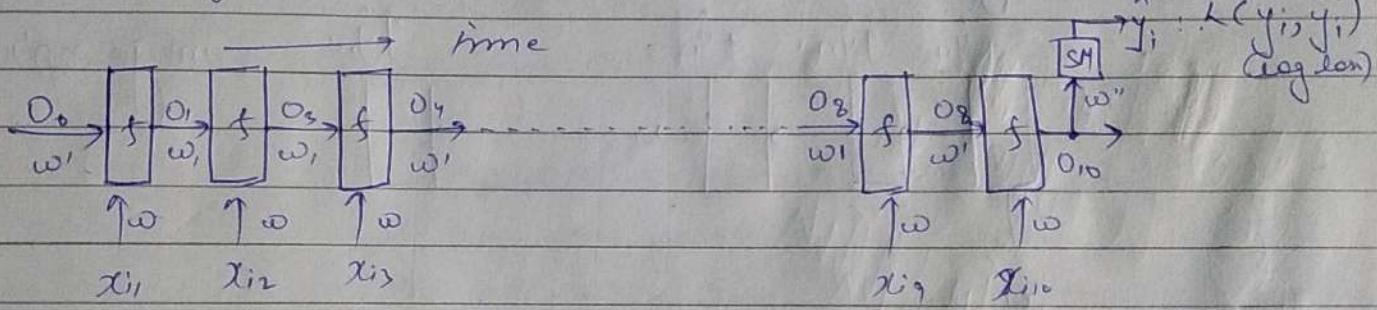


✳ Other way of representing RNN (Box equivalent)



Youtube: Programming Cradle

✳ Training RNN: Backprop



⇒ forward pass :- In the direction of arrows.

⇒ backward pass :- Opposite to the arrows.

$$\frac{\partial L}{\partial O_{10}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} ; \quad \frac{\partial L}{\partial w''} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w''} ; \quad \frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial w} ; \quad \frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial O_{10}} \cdot \frac{\partial O_{10}}{\partial w}$$

✳ Problem with this backprop.

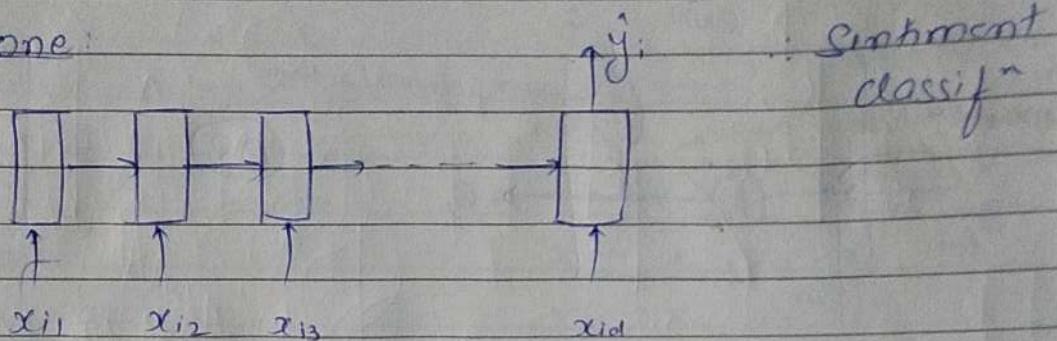
1st word in a sent : $\frac{\partial L}{\partial w}$: lots of multiplications of partial derivatives

NOTE we are getting this error ↳ all are less than 1 hence not bcz of many layers but vanishing grad. problem.
 bcz of repetition of the task ↳ possibility of exploding grad (performing forward and back prop over time)
 when derivat are more than 1

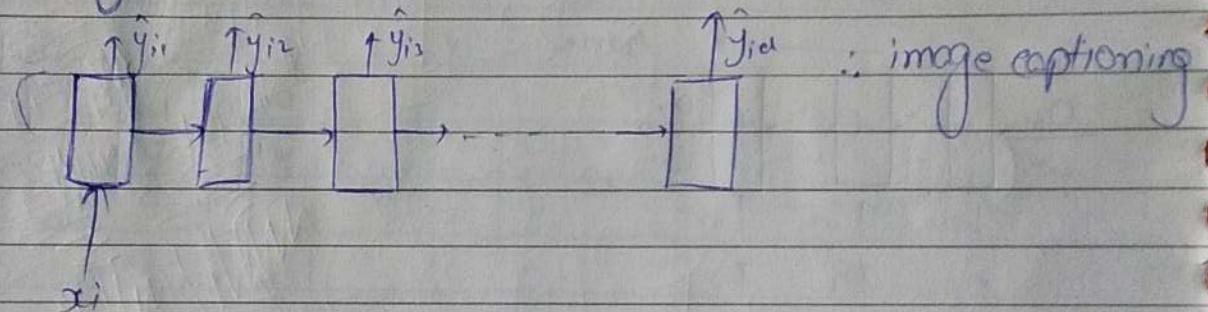
) Also Simple RNN have short term memory.

* Types of RNN :-

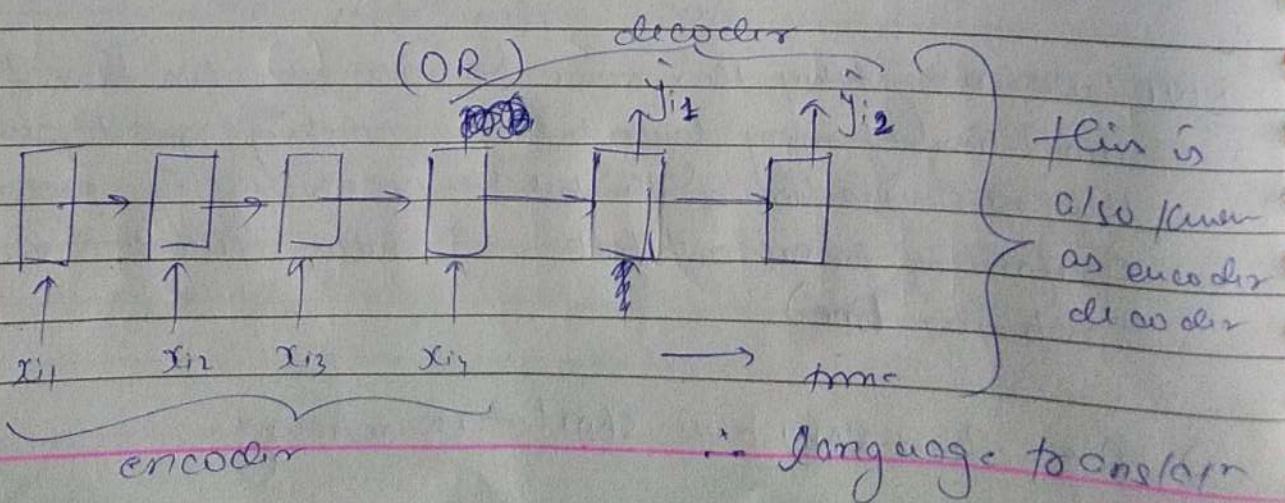
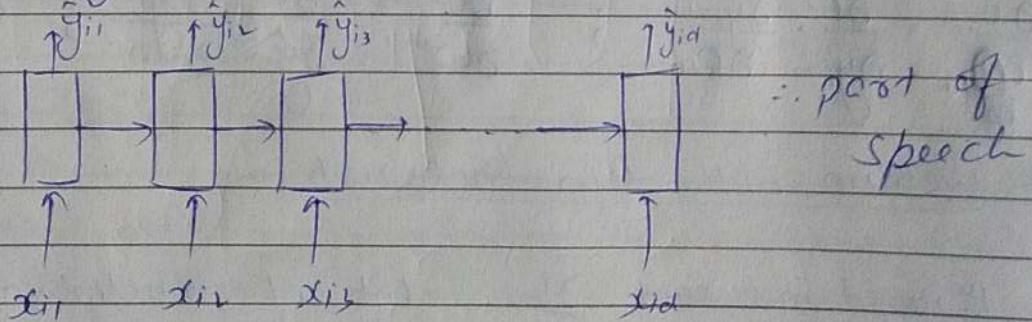
* many to one:



* One to many: Youtube: Programming Cradle



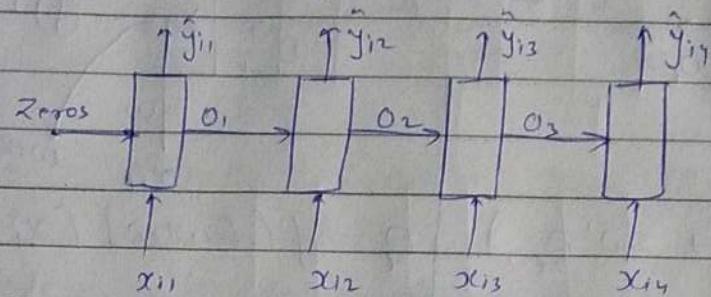
* Many to Many:



① Need for LSTM / GRU

② Problem with Simple RNN

→ many-many (same length)

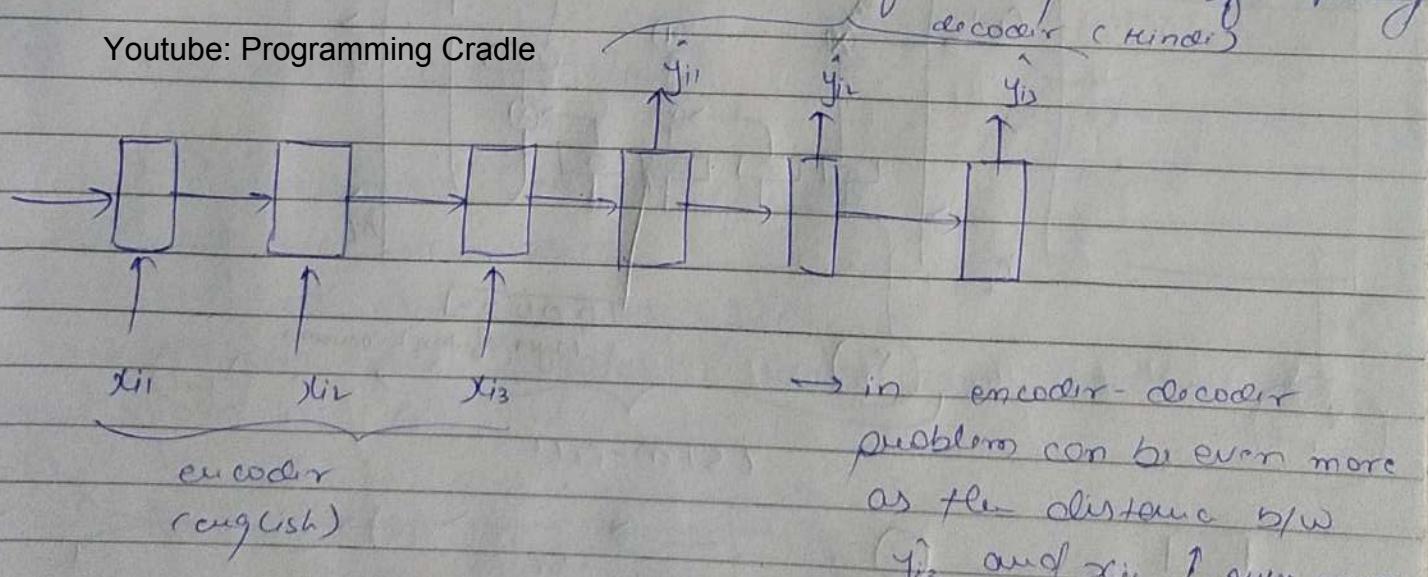


\hat{y}_{i_4} → depends a lot on x_{i_4} and O_3
 \hat{y}_{i_4} → depends less on x_{i_1} and O_1

} we can think
of it as short term
memory.

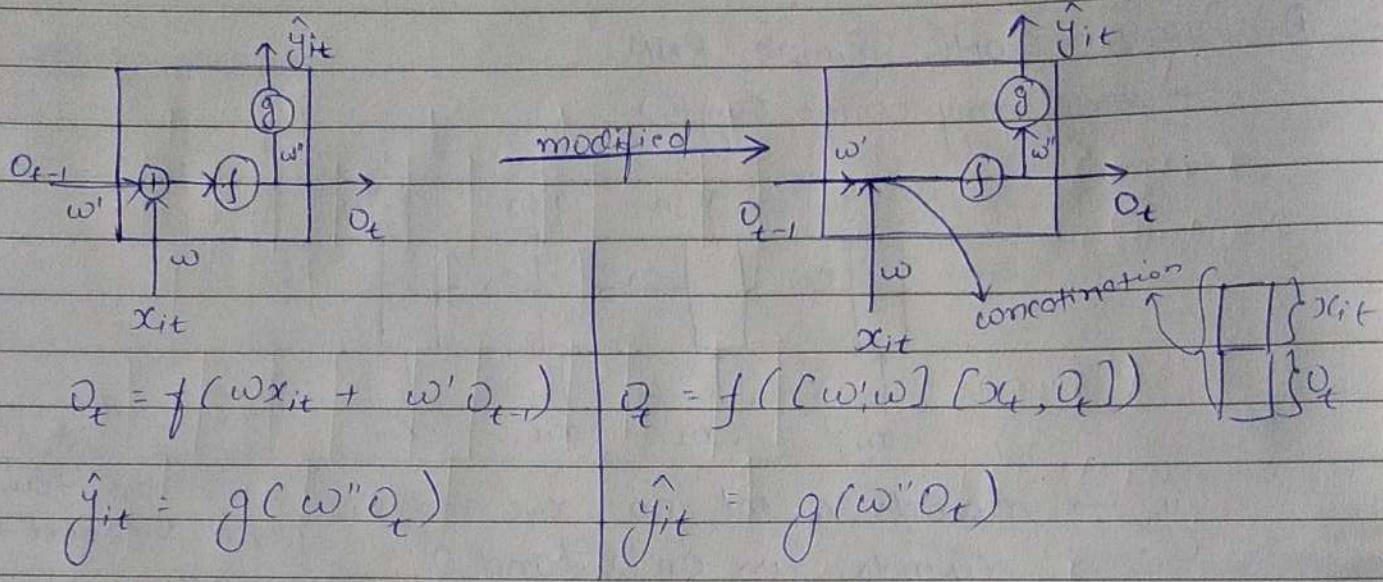
∴ Because of back prop/forward prop, contribution of x_i and O_i will reduce. They are far from \hat{y}_i .

→ Simple RNN can not take care of long term of dependency

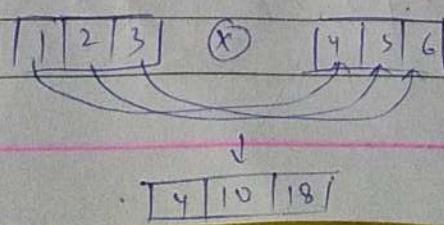
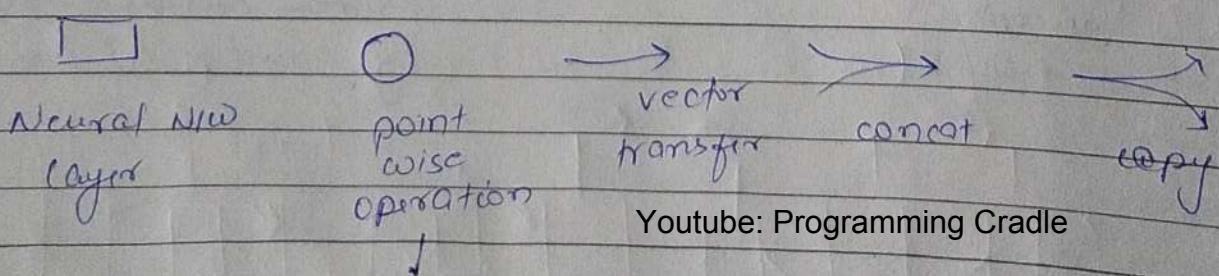
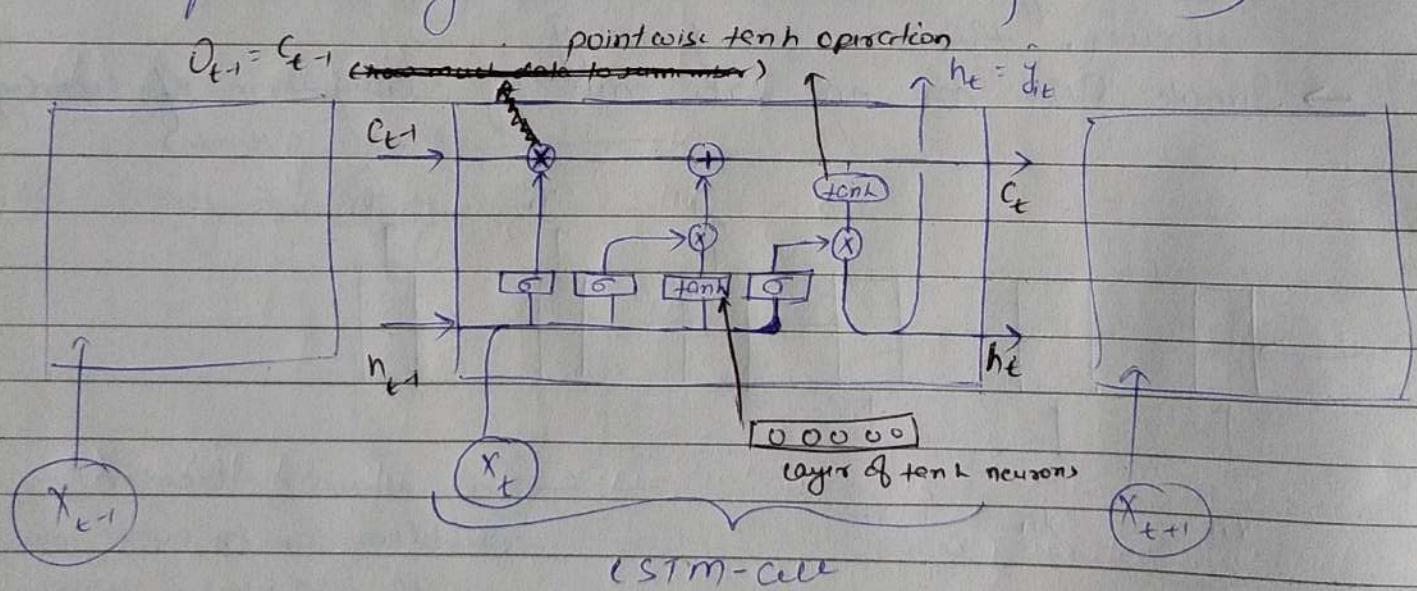


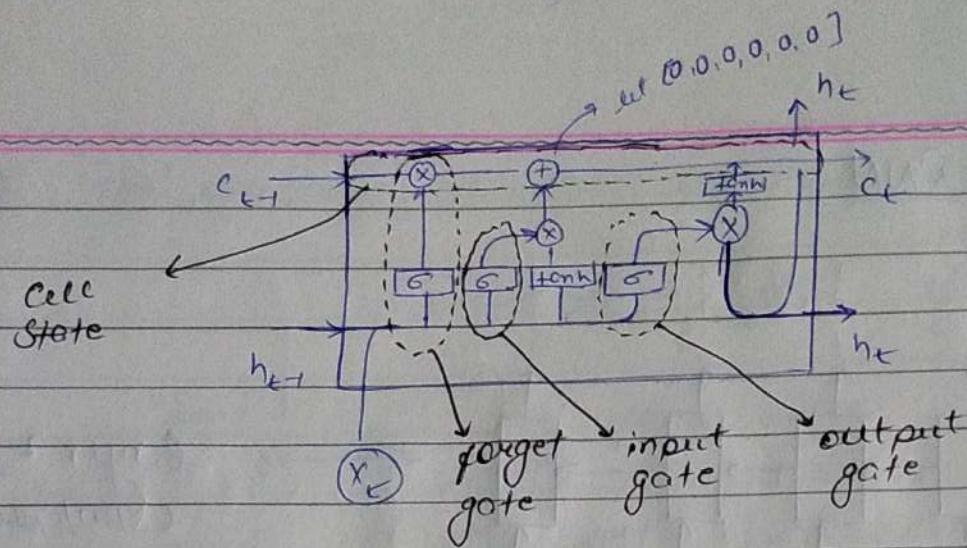
→ in encoder-decoder problem can be even more as the distance b/w y_{i_3} and x_{i_1} is even more

* LSTM Long Short term Memory RNN



(refer colah.github.io understanding - LSTMs)





- Cell State : Works similar to skip connection in ResNet.
- forget gate : Decides ~~how much date to be retained or~~ how much date to be forgotten

example :- $h_t \quad c_{t-1} : (2, 8, 9, 4, 0, 6)$
 forget gate : $(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})$
 $c_t = [(1, 4, 4.5, 2, 0, 3), \oplus (0, 0, 0, 0, 0)]$

Youtube: Programming Cradle

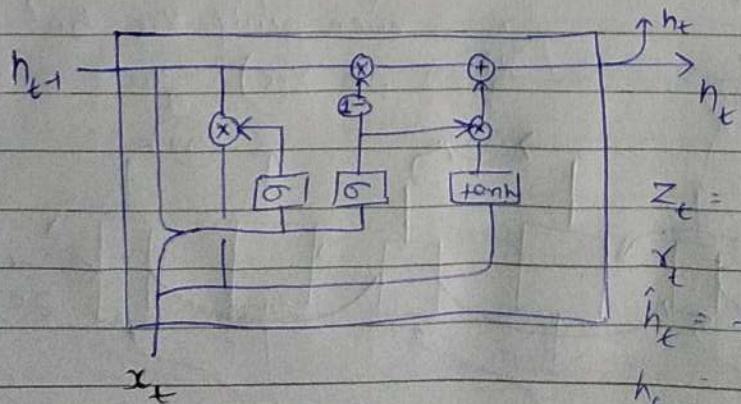
⊗ GRU (Gated Recurrent Unit)

↳ Simplified Version of LSTM

↳ faster to train

↳ As powerful as LSTM

↳ has only 2 gates (reset and update gate)



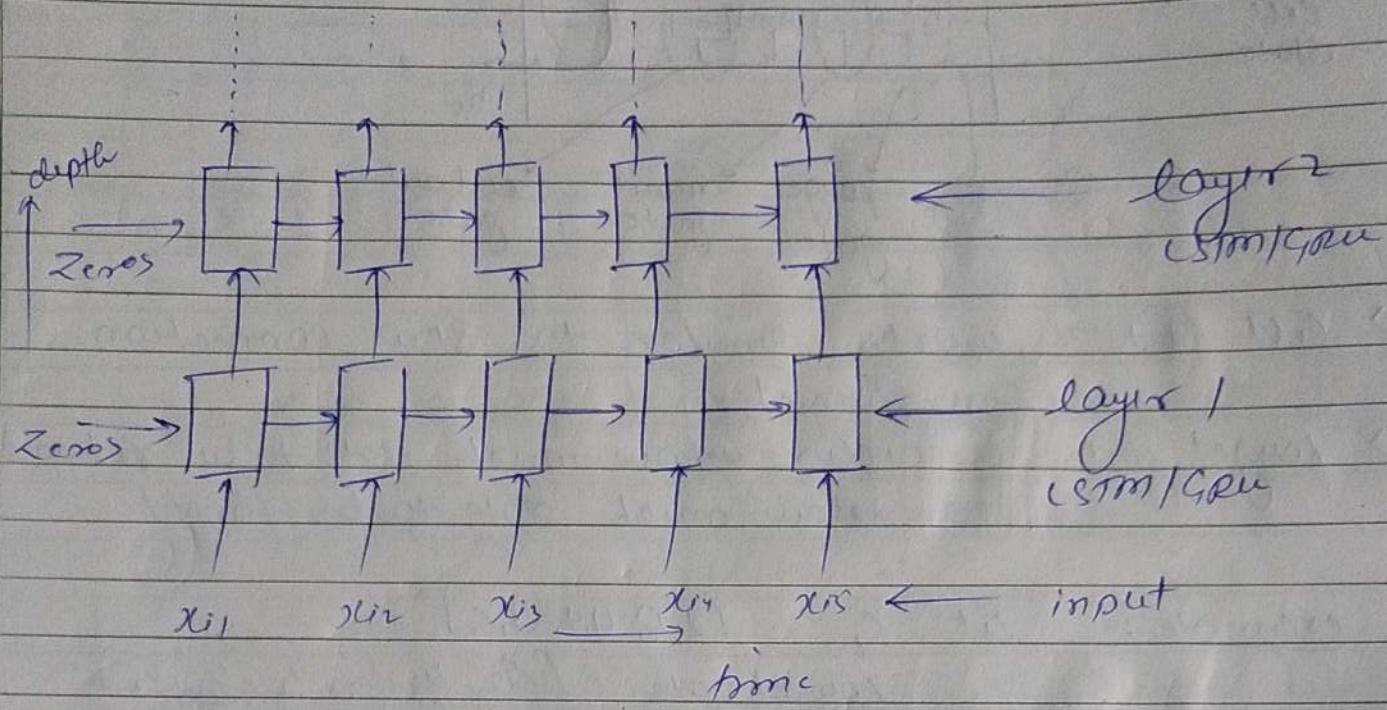
$$z_t = \sigma(\omega_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(\omega_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(\omega \cdot [r_t * h_{t-1}, x_t])$$

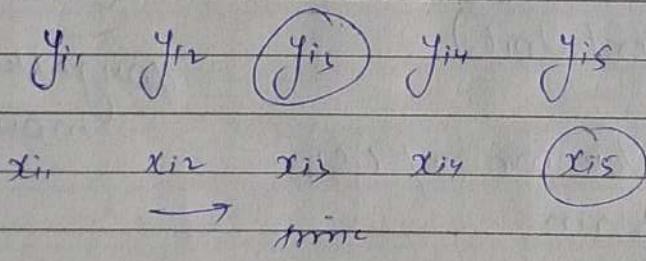
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

* Deep RNN

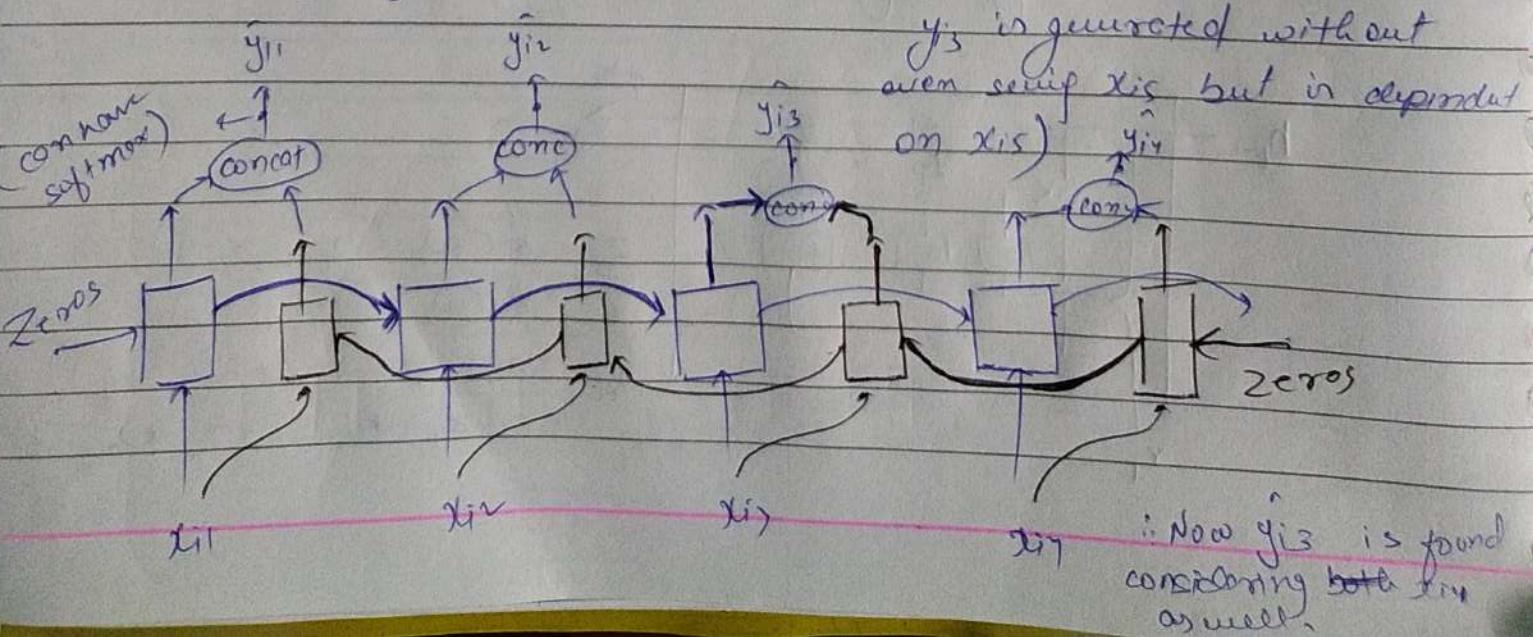


* 'Bi'-directional RNN

Youtube: Programming Cradle



Q) What if y_3 is dependent on x_5 ? (in unidirection)

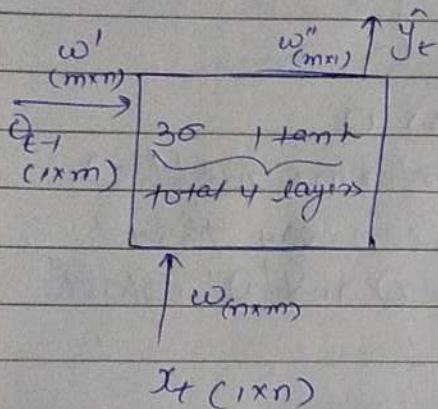


* Total number of trainable params:-

$$y(m^2 + nm + m) \text{ or } y(m^2 + nm)$$

↑
when bias (b)
is considered

when bias (b)
is not considered



m : no. of artificial neurons in each layer
 n : dimensions of input x_t

* Generative Adversarial Networks :-

Task :- Given MNIST images $\rightarrow D \Rightarrow$ Create/generate new MNIST imgs. very similar to images in D but not same.

Youtube: Programming Cradle

* At first solve simpler problem:-

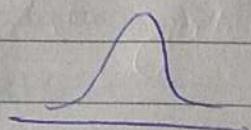
given D = set of heights of students

$$= \{h_1, h_2, h_3, \dots, h_m\}$$

generate $D' = \{h'_1, h'_2, \dots, h'_m\}$

Step 1 find the distribution of the D ; lets say it is normally distributed $\sim N(150, 30)$

$$\text{PDF: } P_D(H) = N(\theta = \{\mu, \sigma\} = [150, 30])$$



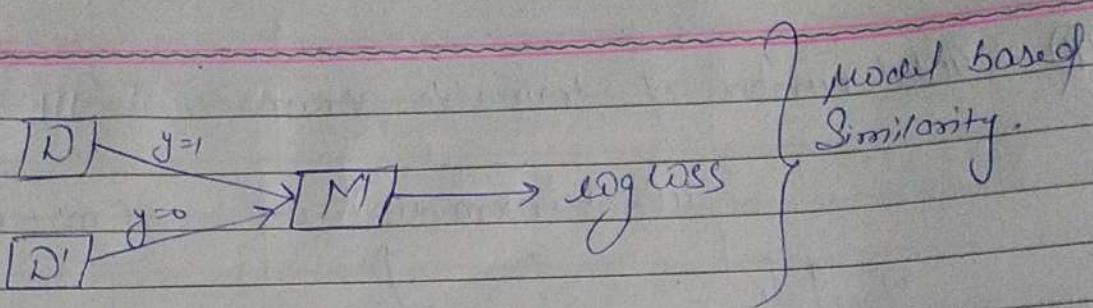
Step 2 Generate random samples using parameters (μ, σ)

(μ) found in previous step.

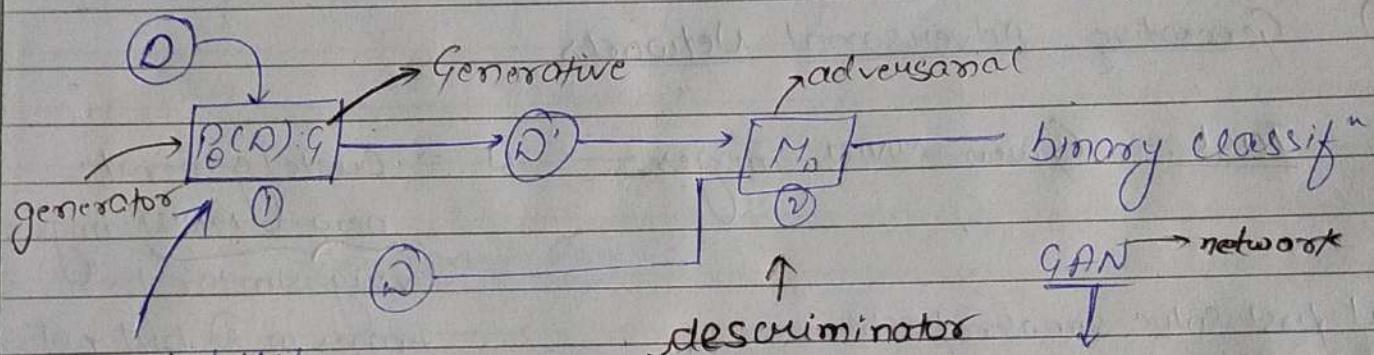
(it can be
- normal
- uniform
etc.)

Step 3 measure / test how similar $D \approx D'$ (how similar dist. of D and D' are)

\rightarrow we can use $\& \&$ plot, KS test, KL divergence but there is another way.... (Model based Similitiy)



- Label D as class 1 and D' as class 0
- if model M (can be any model) performs well => log loss will be less and hence we can say that model is able to distinguish b/w D and D' which means $D \neq D'$



probability dist. f" used
to create new data

(JS-dist) won't work
(KL-div) very well in case of high dim.

\therefore PDF will not work well the data is high dimensional
since we use DNN or CNN
in D box, G box is always some model (can be NN also)

for training GAN

- fix M_D model's param and train G
- find G and train M_D using gradient ascent

Encoder - Decoder Model

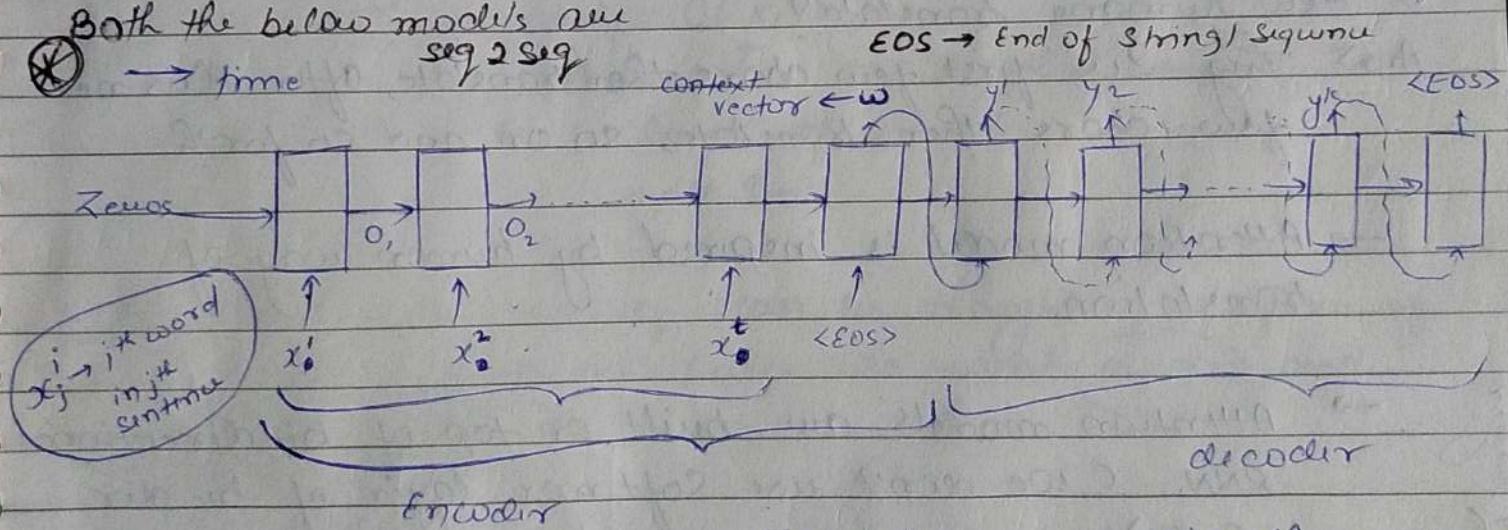
* Machine - Translation :-

NOTE:- t need not equal to k

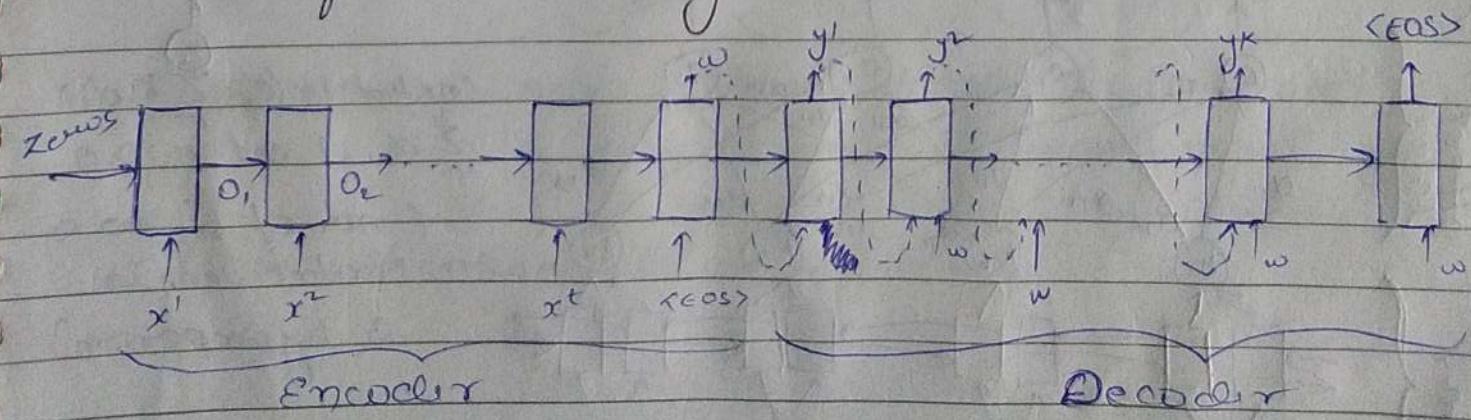
⑥ Image-captioning / descriptions: Youtube: Programming Cradle

Image input $\xrightarrow{M} y^1 y^2 y^3 \dots y^K$ output

 Both the below models are seq 2 seq
→ time



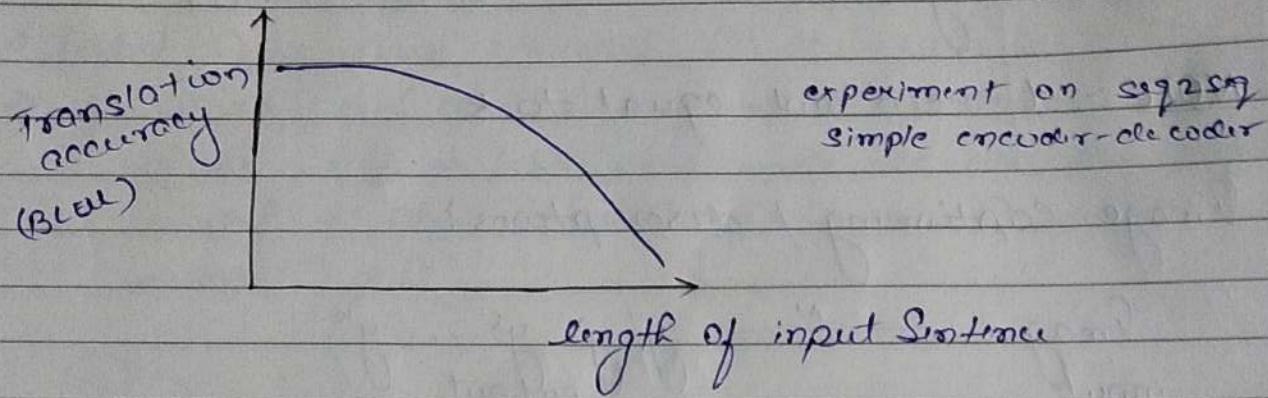
* Above model works really well for short sentences. But what if sent. are long? we can use below model.





Attention Models

- * Biggest problem with Seq2Seq models is they do not perform very well with long sentences.

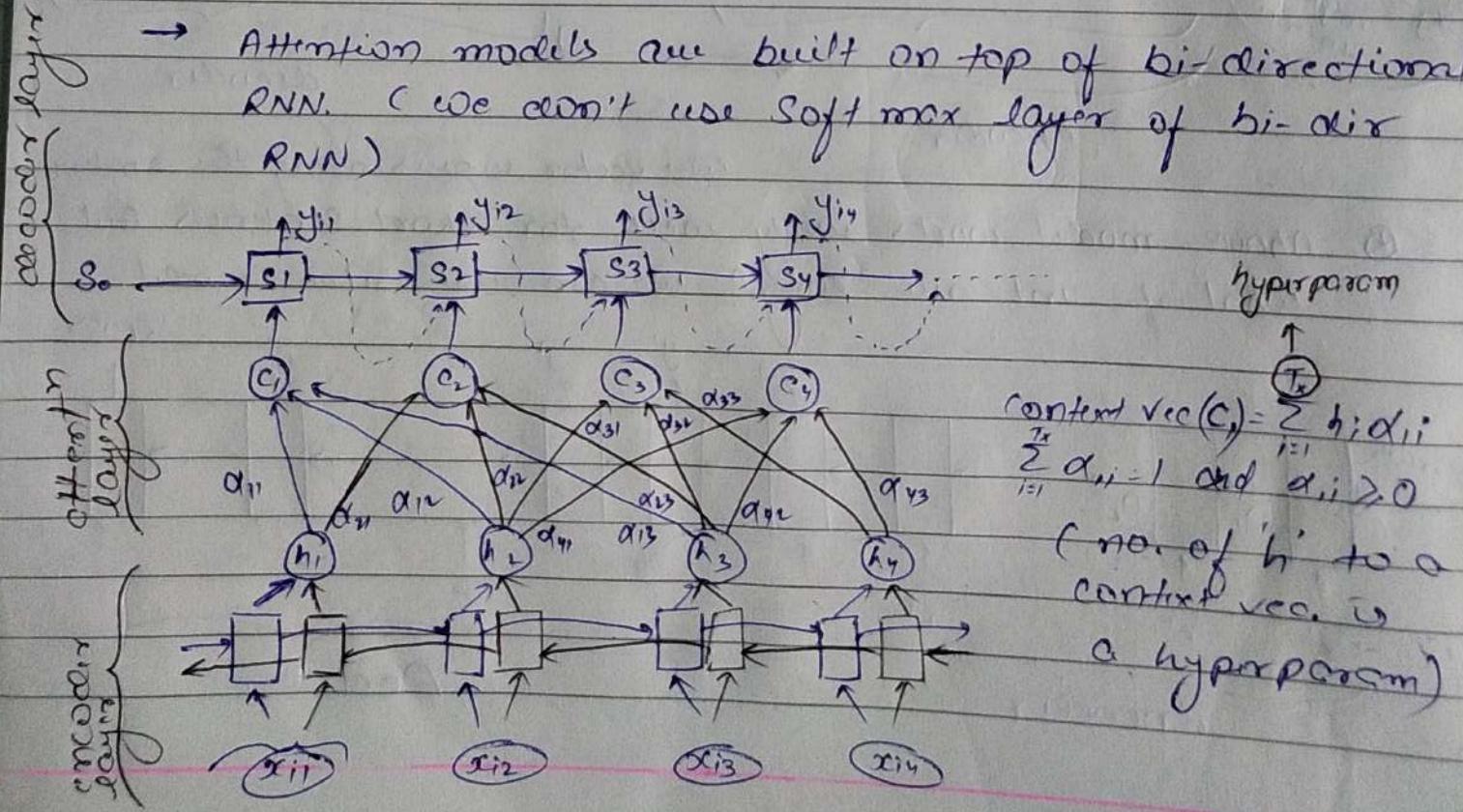


• How humans translate? Youtube: Programming Cradle

Ans> They see first few chars then translate, after their next few chars then translate, so on and so forth.

→ Attention model is inspired by human way of translation.

→ Attention models are built on top of bi-directional RNN. (We don't use Softmax layer of bi-dir RNN)



Neutral Machine ~~translates~~
translation by jointly learning to align and translate.

Q> How to compute α_{ij} 's?

$$c_i = \sum_{j=1}^{k_2} \alpha_{ij} h_j \quad (\text{eqn 5 in the research paper})$$

→ attention-model

$$\alpha_{ij} = \frac{\exp(c_{ij})}{\sum_{k=1}^{k_2} \exp(c_{ik})}$$

$$c_{ij} = a(s_{i-1}, h_j)$$

feed forward ANN (Small)

example :-

$$c_{12} = a(s_0, h_2)$$

⊗ Drawback Youtube: Programming Cradle

Time Complexity: $O(k_1 \cdot k_2)$

\uparrow , length of output sent.
 \downarrow , length of input sent.

NOTE :) constraint on α_{ij} can be thought of as a form of regularization. Since we want $\alpha_{ij} \geq 0$ and

$$\sum_{i=1}^{k_1} \alpha_{ij} = 1$$

) We are using $\exp(c^*)$ in calculation of α_{ij} to make sure $\alpha_{ij} \geq 0$ and when we sum all α_{ij} or d_{ij} or $d_{ij} - \dots$ we will get it as 1

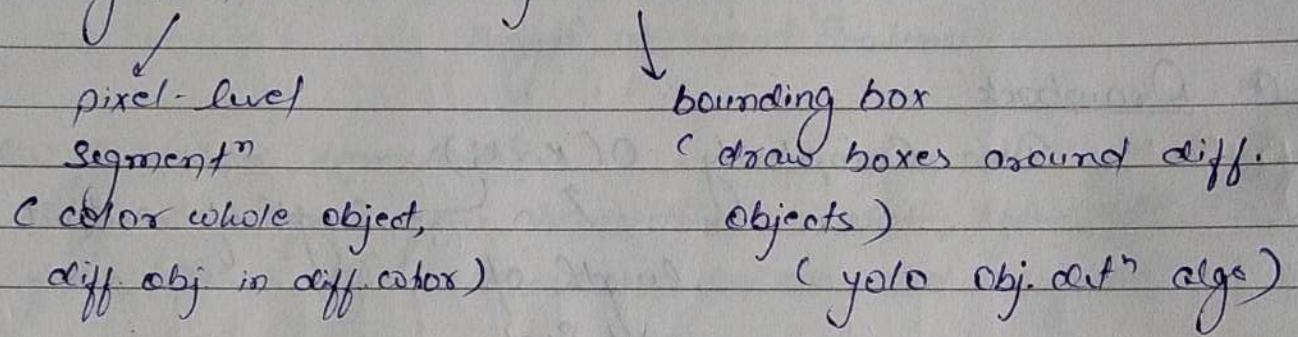
⊗ Transformers and BERT Blog: 'The illustrated transformer' by Jay Alammar

Blog: "Attention is all you need."

* Image Segmentation

- input is an image and output is also an img.
- Diff. subject in the image will be colored in diff. colors.
- coco dataset
- image segmentation traffic dataset ("cityscapes")
- google medical segmentation ("gg-mm-segment")

→ Segmentation vs Object detection :-



1 → Performance Metric :- Youtube: Programming Cradle

i) Average Precision and Recall for each class/ (coco ds) color

$$\text{i)} \text{ Jaccard Sim} = \frac{R_i \cap \hat{R}_i}{R_i \cup \hat{R}_i}$$

(intersection over union)

ideal value
= 1

The diagram shows a square image containing a blue circle. A smaller blue circle inside it is labeled \hat{R}_i . The area of intersection between the two circles is labeled $R_i \cap \hat{R}_i$. The entire area of the image is labeled $R_i \cup \hat{R}_i$. Arrows point from the labels to their corresponding regions in the image.

→ classical approaches :-

- ① clustering. { K-means, ... }
- ② Edge - detection
- ③ Graph - Theoretic (spectral clustering)

→ Paper: imb.informatik.uni-freiburg.de/people/mennebir/u-net/

→ Deep-learning based method.

- U-nets
- Fully convolutional Networks
- Mask R-CNN
- SegNet
- ...

④ Object Detection Youtube: Programming Cradle

- input : img / vid
- output : bounding box-1, obj. class
bounding box-2, obj. class
- :
- Dataset: coco (80-class-data)
- Trade off: Speed vs mAP (avg. precision)
 - ↓
 - ↳ medical

Self driving cars
- Diff algorithms :-

R-CNN

SSD (Single Shot detector)

YOLO v1

Fast R-CNN

Retina Net

YOLO v2 / YOLO 9000

Faster R-CNN

YOLO v3

YOLO v3 : feature extractor

→ fully convolutional Network

→ Conv-Batch Norm - Leaky ReLU

→ Conv. with stride - return sample

→ no padding

→ pre-trained model.

paper → pycroddicay/medical files/papers/YOLOv3.pdf.

9/10/10

Q) why Darknet-53 was developed?

- A) To match the accuracy of ResNet-152 and at the same time to match the "no. of opers" of Darknet-19. and also for higher fps.

* Per class Sigmoids

Youtube: Programming Cradle

→ for multiclass classification softmax can be used but it will give only one class which has the max. probability. But in case of situations like we have a class "person" and "woman" here a woman can be a person so 2 classes can be returned. BUT this can't be done using softmax. Hence logistic regression / sigmoids are build for each classes.

* Loss :

$$\begin{aligned}
 & \sum_{\text{No coordinate}} \text{No coordinate} * \text{Squared loss} + \sum_{\substack{\text{no. obj in boxes} \\ \text{on } \{x_w, y_h, tx \text{ and } ty\} \rightarrow (0 \text{ to } \infty)}}^N N \cdot \log\text{-loss for } p_i \\
 & + \text{bounding boxes} + \log\text{-loss for } p_0 \\
 & \text{containing obj in ground truth} + \log\text{-loss for each } p_i
 \end{aligned}$$

Yolov3 Network architecture bengt.damprof.com