

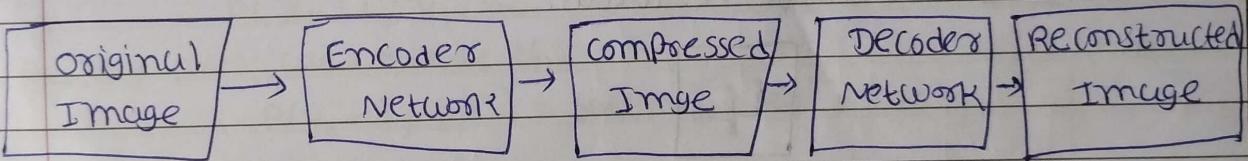
convolution

→ For anomaly detection train data on non-anomaly data and then anything fall out of the norm that automatically labeled as an anomaly.

Autoencoder in CNN

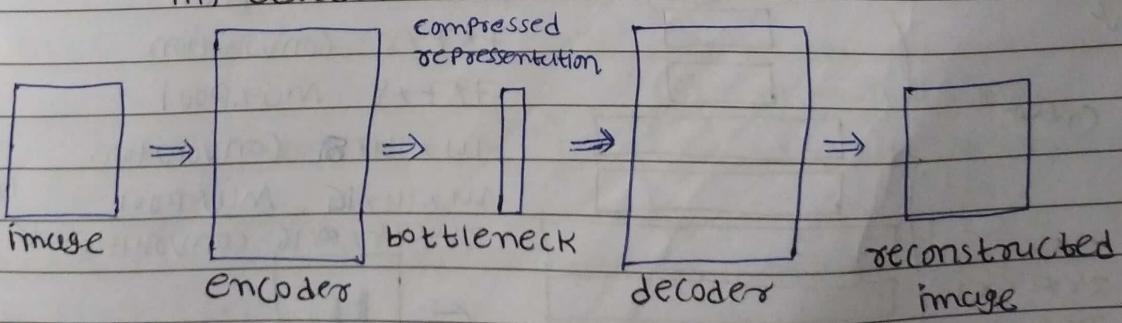
→ Reconstruction error is high for anomaly data

- ⇒ A convolution ~~into~~ en: autoencoder is a neural network (a special case of unsupervised learning model) that is train to reproduce its input image in output layer.
- ⇒ Image is pass through an encoder which is a convolution network that produce the low dimensional representation of Image.
- ⇒ The decoder which is another sample of convolution network (convnet) takes this compressed image and ~~recon~~ reconstructs the original image.
- ⇒ Encoder is use to compress the data and decoder is used to reproduce the original image.
- ⇒ compression logic is data-specific meaning it is learned from data rather than predefined compression algorithms such as JPEG, MP3, and so on.



- ⇒ autoencoders are deep learning architecture capable of reconstructing data instance from their feature vector.
- ⇒ They work on all sort of data but here we see application their application on image data.
- ⇒ Autoencoders is made with 3 main components

- i) encoder
- ii) bottleneck
- iii) decoder



\Rightarrow In autoencoder model will be train in a such a way so that we can get original data from compressed data that many compressed data contain all important information not just data
 \Rightarrow Auto encoder has analogy with PCA

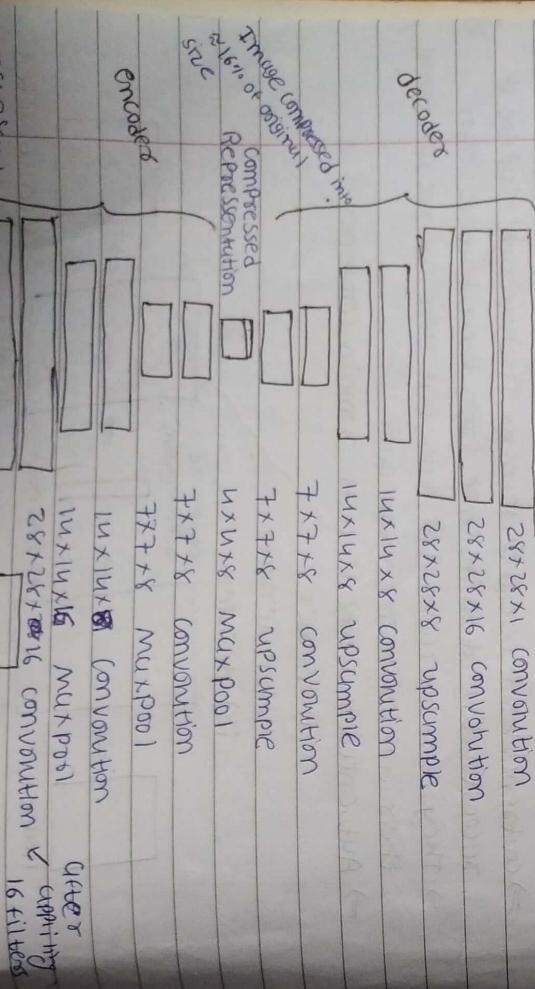
i) Encoder : The first section of auto encoders. encodes

extract the most salient(main) features from images and return them as a vector.

ii) Bottleneck: Bottleneck also called a code layer, serves as an extra layer which help to compress the extracted feature into a small vector representation. This is done in a bit to make it more difficult for the decoder to make sense of the features and force it to learn more complex mappings

iii) Decoders: decoders attempts to make sense of the features coming from the encoder, which have been subsequently compressed in the bottleneck so as to reconstruct the original image as it was.

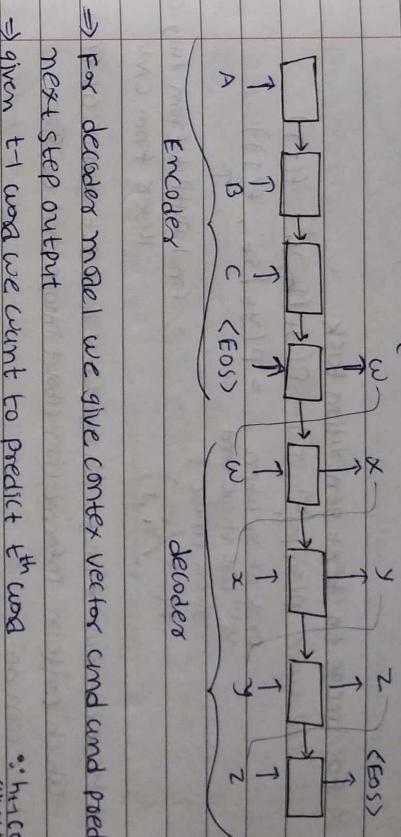
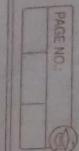
4



Encoder-decoder

NPTEL - NOC - IITM

context vector

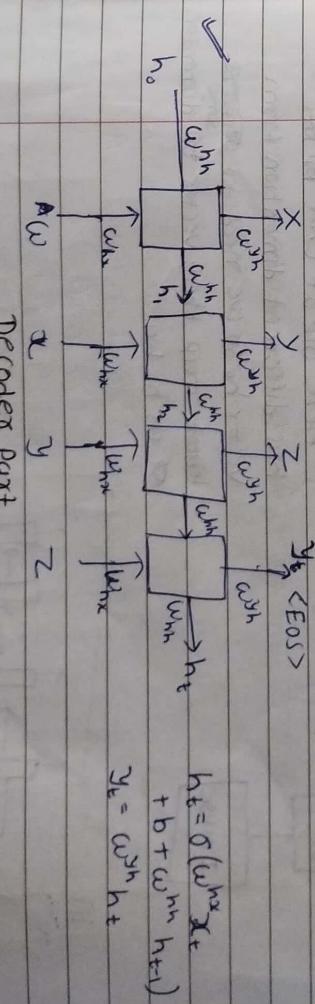


$$P(y_t = j | y_1^{t-1}) = \text{softmax}(\mathbf{w}^y h_t + \mathbf{c}_j)$$

$$\mathbf{y}_t = \mathbf{w}^y h_t + \mathbf{c}_j$$

$$h_t = \sigma(\mathbf{w}^h x_t + b + \mathbf{w}^h h_{t-1})$$

$$\mathbf{w}^y \quad \mathbf{x} \quad \mathbf{y} \quad \mathbf{z}$$



$$L_t(\theta) = -\log P(y_t = l_t | y_1^{t-1})$$

so that $L(\theta)$ is minimize

\Rightarrow and we do back propagation further and adjust weight

so that $L(\theta)$ is minimize

$\Rightarrow h_0$ is not computed but just randomly initialized

* For image to text generation task

$$\text{we take } p(y_t | y_{t-1}, I) = P(y_t | h_t, y_{t-1}, f(I))$$

$$\begin{aligned} \text{prob of } y_t \text{ when } &= P(y_t | s_t, f(I)) \\ y_{t-1} \text{ use these cond} & \\ \text{given image } I & \end{aligned}$$

$$p(y_t | y_{t-1}, I)$$

final output from fully connected
layer from CNN

\Rightarrow here we pass image through CNN and get output as $f(I)$

which contain information about image that we given.

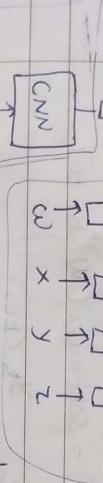
\Rightarrow we can do

i) connect op of CNN to cell as h_t .

OR ii) connect op of CNN to cell for every time stamp

i) connect op of CNN to cell for every time stamp

$h_0 = f(I)$



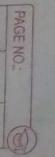
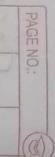
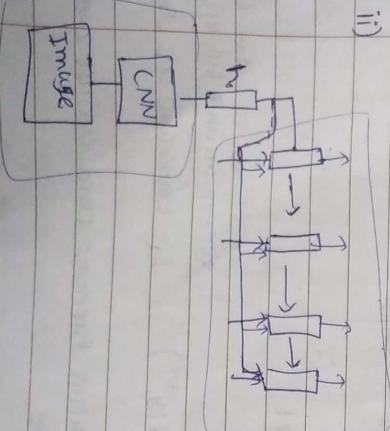
Here op of CNN has

vector h_0 so we need to ~~map~~ some

map that vectors to dimension

of h_0 here.

i)



Code for i)

TASK : factual entailment

Data : $\{x_i\}$ = premise, y_i = hypothesis, y_{i-1}^N

Encoder :

$$h_t = RNN(h_{t-1}, x_{it})$$

Decoder :

$$h_t = RNN(h_{t-1}, \hat{e}(y_{t-1}))$$

$h_0 = h_T$ (T is length of input)

$$y_t = RNN(s_t, e(y_{t-1}))$$

embedding

* Task: Image captioning

Data: x_i = image, y_i = caption, y_{i-1}^N

image & corresponding caption is given for training

Model:

Encoder: $h_0 = CNN(x_i)$ embedding of previous output it can be onehot, wordvec

Decoder: $h_t = RNN(h_{t-1}, \hat{e}(y_{t-1}))$

$P(y_t | y_{t-1}, I) = \text{Softmax}(\omega^{yt} h_t + b)$

parameters: $\omega^{yt}, \omega^{ht}, \omega^{xt}, \omega^{ht}$ biases

$$\text{Loss: } L(\theta) = \sum_{i=1}^T \log p(y_k = 1 | y_{k-1}^T, I)$$

Back Propagation gradient descent

Back Propagation gradient descent



★ Textual entailment

TASK: Textual entailment

Model: $\{x_i = \text{Premise}, y_i = \text{hypothesis}\}_{i=1}^N$

(option 1) Encoder: $h_t^{\text{Encoder}} = \text{RNN}(h_{t-1}^{\text{Encoder}}, x_{it})$

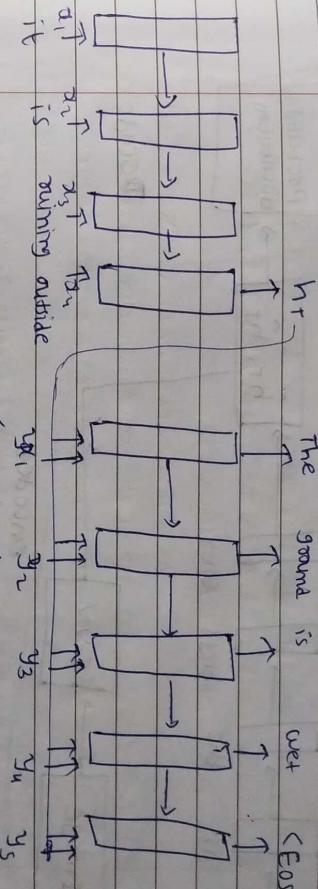
Decoder: $h_0 = h_T$ (T is length of input)

$$h_t = \text{RNN}(h_{t-1}, e(y_{t-1}))$$

$$P(y_t | y_{t-1}, x) = \text{softmax}(w^{yx} h_t + b)$$

parameters: w^{hh} , w^{hx} , w^{hy} , w^{by} , w^{bh} , w^{bx} , b
 Loss: $L(\theta) = -\sum_{i=t}^T \log P(y_i = l_i | y_{i-1}, x)$

Algo: gradient descent with back propagation



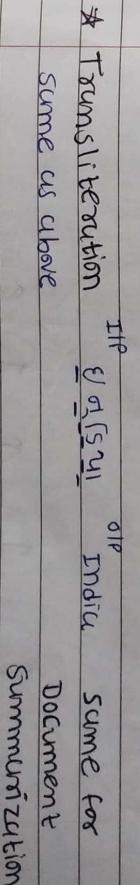
★ machine translation

Input: I am going home

TASK: Machine translation

Model: $\{x_i = \text{Source}, y_i = \text{target}\}_{i=1}^N$

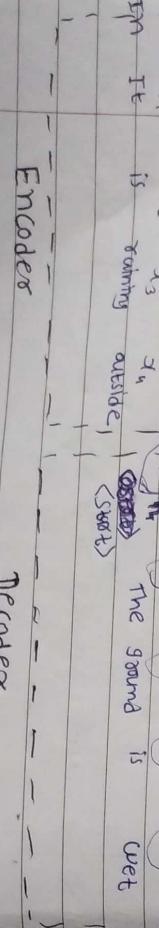
Other sume us above There is ~~two~~ two model



★ Image Question Answering

Model: $\{x_i = \{I, q_i\}, y_i = \text{Answer}\}_{i=1}^N$

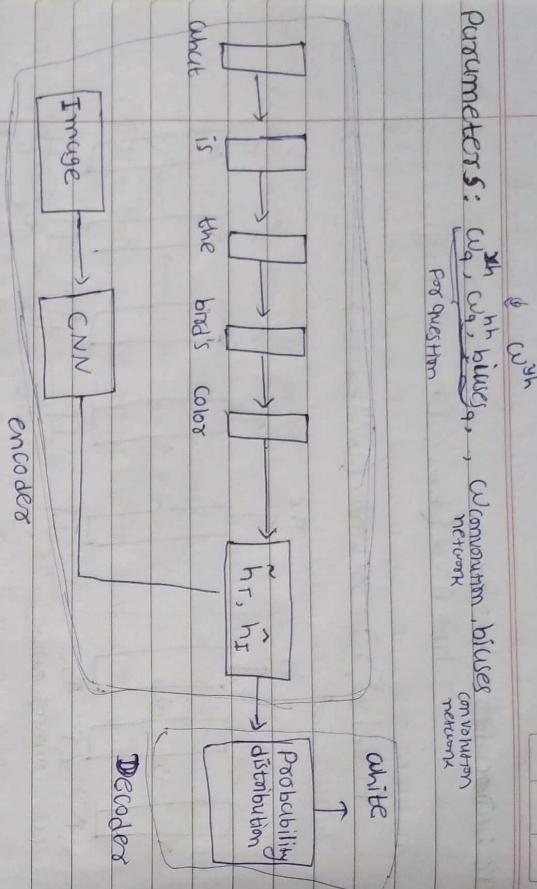
Image question (right)



Decoder: $P(y_i | I) = \text{softmax}(w^{yI} h_t + b)$

Prob distribution

Parameters: w_q, w_k, w_v, b_{ques} , W convolution, biases
For question network



encoder

Decoder

white

Probability distribution

* Video captioning

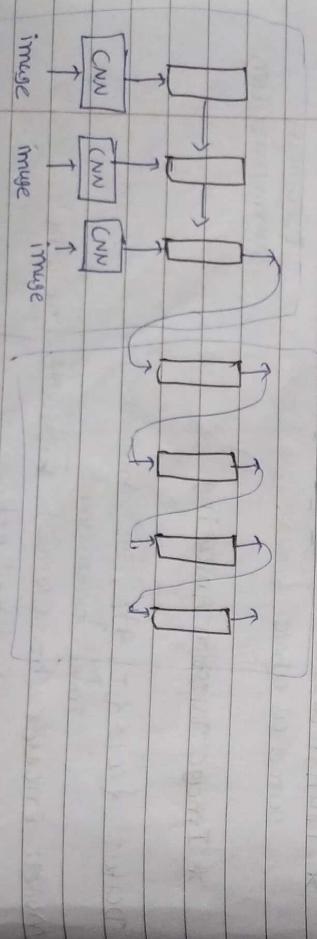
Data: x_i = video, y_i = desc $_i$; \hat{y}_{t+1}
↑
description

Model: Encoder: $h_o = h_T$

$$h_t = RNN(h_{t-1}, CNN(x_{it}))$$

$$P(y_t | y^{t-1}, x) = \text{softmax}(w_{yh} h_t + b)$$

All other same

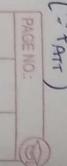


* Video classification

Encoder: same as video captioning

Decoder: Probability distribution

⇒ Attention model is better due to better modeling choice
⇒ have linear transformation followed by non linearity ("tanh")
Attention model



Attention Mechanism NLPFL-NOC-IIIM

⇒ Encoder feeds words the sentences only once and encode it

⇒ At each timestep the decoder uses this embedding to produce a new word.

⇒ Human try to produce each word in the output by focusing only on certain word in the input
⇒ Essentially at each time step we come up with a distribution on the input words.

Q1P : I am going home ⇒ This distribution tell us how much attention to give to each input word
t₁: [1 0 0 0 0] t₂: [0 0 0 0 1] t₃: [0 0 0.5 0.5 0] t₄: [0 1 0 0 0]

1P Main ghar ja rahi hoon

⇒ we could just take a weighted average of the corresponding word representations and feed it to decoder.

⇒ For E₂ timestep 3 we can just take weighted avg of 'ju' and 'ruhi'

⇒ This help to ignore irrelevant information (which is not important at that particular time step)

⇒ To do that we define function input jth word is how important for tth timestep

$E_{jt} = f_{ATT}(h_{t-1}, h_j)$ important for tth timestep
fth time step how important is jth word
decoder states of all the words from encoder

⇒ so E_{jt} is depend on what is happen in decoder so far but

and what is current word actually look like.

⇒ E_{jt} is captures the importance of the jth input word for decoding tth output word

⇒ we can normalize those weights by using softmax

$$d_{jt} = \frac{\exp(e_{jt})}{\sum_{j=1}^n \exp(e_{jt})}$$

⇒ d_{jt} is prob of focusing on the jth word to produce the tth output word

⇒ h_t^{dec} = decoder RNN's state at tth time step
 h_j^{enc} = Encoder RNN's state at jth time step

⇒ one possible choice for function f_{ATT} is (There are many possi)

$$\underline{e_{jt}} = \underline{v_{att}^T tanh(\underline{U_{att} h_{t-1}} + \underline{U_{att} h_j})}$$

⇒ $V_{att} \in \mathbb{R}^d$, $U_{att} \in \mathbb{R}^{d \times d}$, $W_{att} \in \mathbb{R}^{d \times d}$ are additional parameters of the model

⇒ These parameters will be learned along with the other parameters of the encoder and decoder

⇒ This model would make lot of sense if were given the true d_{jt} at training time

$$\text{let } d_{jt}^{\text{true}} = [0, 0, 0.5, 0.5, 0]$$

by minimizing loss by adjusting weights gives us appropriate values of weight

$E_{jt} = f_{ATT}(h_{t-1}, h_j)$ important for tth timestep

But, in practice it is very hard to get a true d_{jt}^{true}

⇒ Then how would this model work in the absence of such data?

⇒ It's work because of better modeling choice (we give attention on only relevant words)

⇒ This is a more informed model

⇒ we are essentially asking the model to approach the problem in a better way (natural way)
(like human can rise cycle if helped out of house on rainy then ...)

\Rightarrow given enough data it should be able to learn these attention weight just as humans do (E.g. code) \Rightarrow and in practice indeed this models work better than the vanilla encoder-decoder models

TASK: Machine Translation

DATASET: $\{x_i = \text{source}, y_i = \text{target}\}_{i=1}^N$

only input and output sentence are given. And which work uses attention to every time is not given

Encoder:

$$h_t^{\text{enc}} = RNN(h_{t-1}^{\text{enc}}, x_t)$$

$$\boxed{h_t^{\text{enc}} = h_{t-1}^{\text{enc}} + h_t}$$

Op from Encoder and input for Decoder

$$\text{Decoder: } e_{jt} = V_{\text{attn}}^T \tanh(W_{\text{attn}} h_j + W_{\text{attn}} h_t)$$

from decoder

$$c_t = \sum_{j=1}^J a_{jt} h_j^{\text{enc}}$$

Op from Decoder

$$\text{prob distri } d_{jt} = \text{softmax}(e_{jt})$$

parameters: $W_{\text{attn}}, W_{\text{dec}}, W_{\text{dec}}, b$, W_{enc}, b , W_{enc}, b .

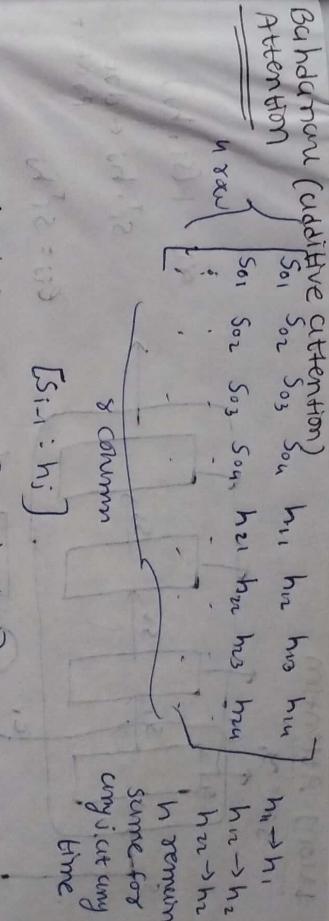
Encoder, V_{attn}

LOSS and ALGORITHM: sum of $-\log(d_{j,y_t})$ others

entropies

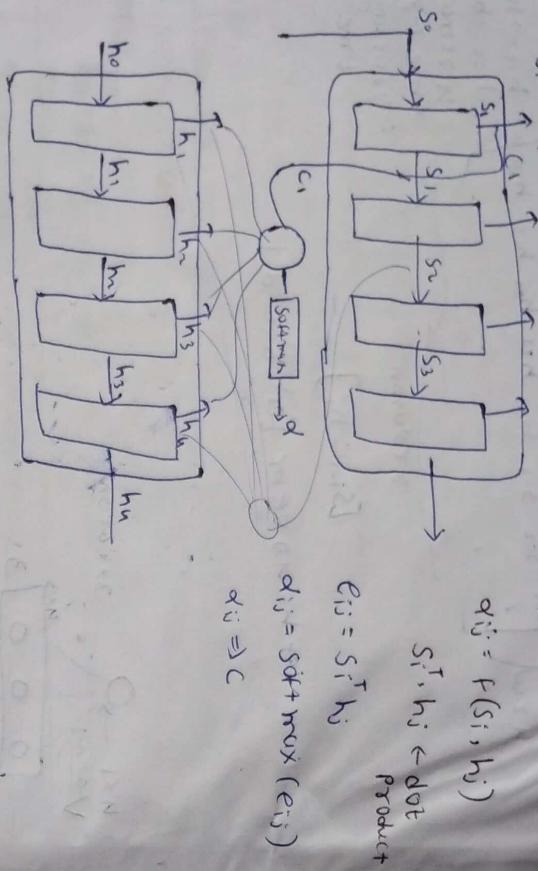
EN with attention

Batch gradient



Wrong Attention

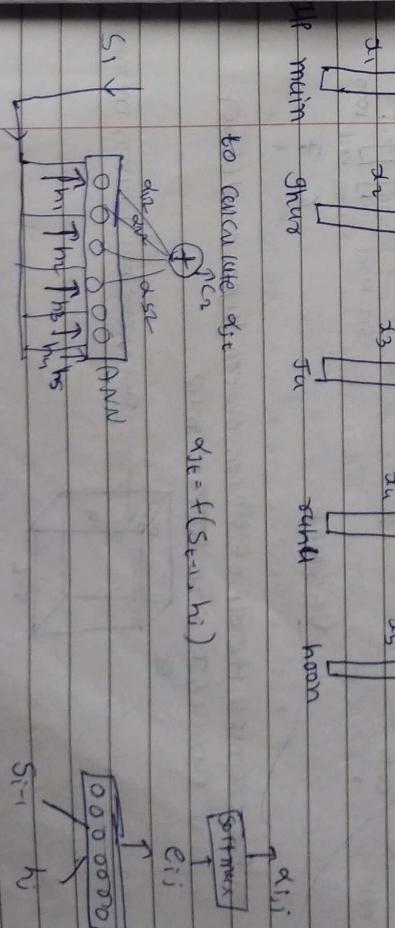
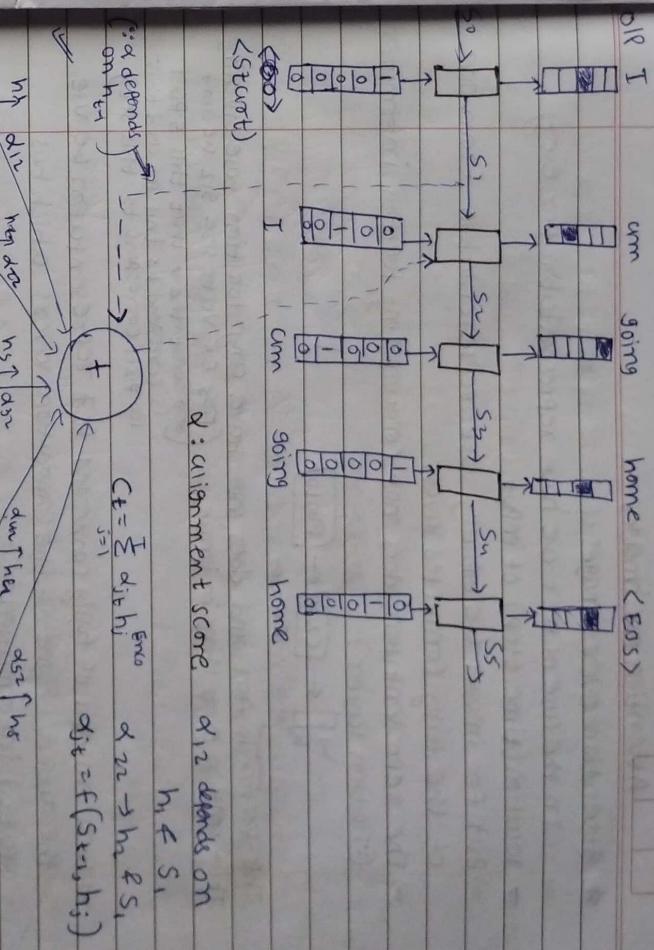
$$s_i \cdot k_i = s_i$$



Deco
Prat

\Rightarrow Here we use s_i . So we get recent information and get more accurate c value.

$$[c_1, c_2, c_3, c_{in}] \rightarrow c_i$$



Run
Loss
on

Attention on Images

* Attention over images

\Rightarrow In the case of text we have representation for every location (time step) of input

\Rightarrow But for images we typically use representation from one of the fully connected layers

\Rightarrow The representation does not contain any location information

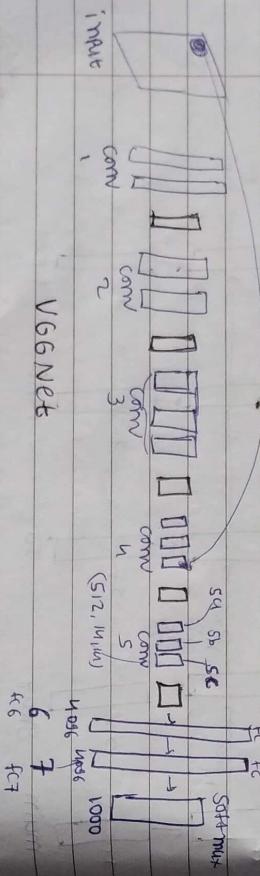
(\therefore location layers)



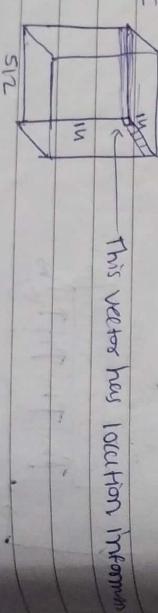
This vector h is flat and does not have any location bias encoding that contains everything in image

(For ex. vectors size 512 we can't say exactly that this 5 pixel is correspond to this set and that 100 pixel correspond to this set)

\Rightarrow Well instead of fully connected FC7 representation we use the output of one of the convolution layers which has spatial information

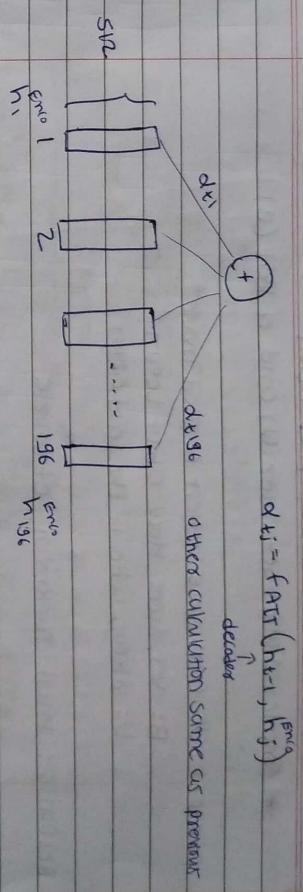


\Rightarrow For example the 5th convolution layer of VGGNet is a 1x1x512 size feature map



This vector has location information

\Rightarrow we could think of this as 196 locations (each having 512 dimension representation)



Hierarchical Attention

→ consider dialog: bet user (U) and a bot (B)

context U: Can you suggest a good movie?

B: Yes, sure how about Logan?

U: Okay, who is the lead actor?

Response: Hugh Jackman, Of course

⇒ The dialog contains a sequence of utterance bet' the user and bot

⇒ Each utterance in turn is a sequence of words

⇒ Thus what we have here is a "sequence of sequence" of words

⇒ Can you think of an encoder for such a sequence of sequences? ⇒ RNN or RNN

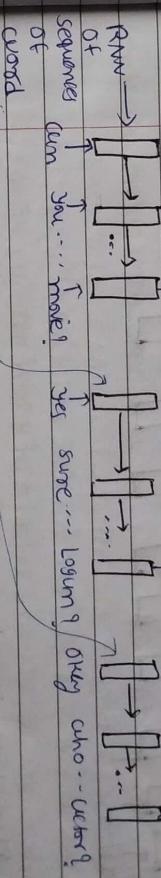
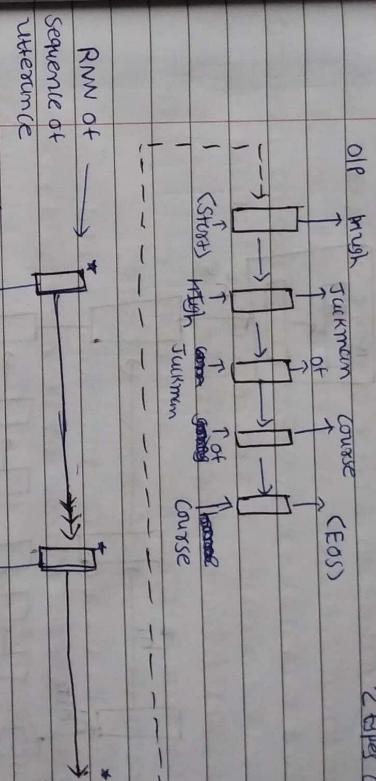
two level hierarchical RNN encoder

⇒ The first level RNN operates on the sequence of words in each utterance and give us a representation

⇒ we have known sequence or utterance representations

⇒ so another RNN encodes this sequence and gives a single representation for sequences of utterances

⇒ decoder can then produce an output sequence conditioned on this utterance



⇒ Another example document classification or summarization

⇒ A doc. is seq. of sentences

each sentence in turn is a sequence of words

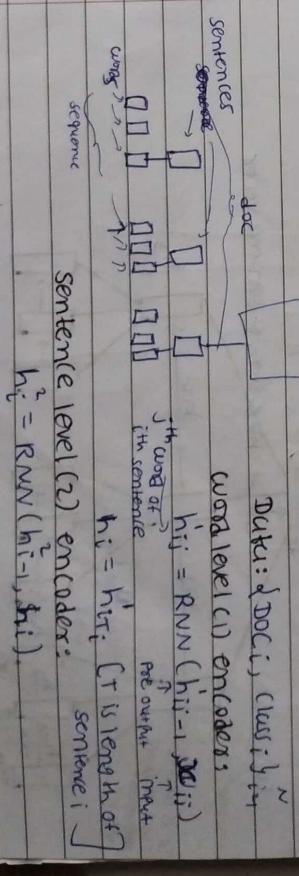
⇒ we can again use hierarchical RNN to model this

⇒ Another example document classification or summarization

⇒ A doc. is seq. of sentences

each sentence in turn is a sequence of words

⇒ we can again use hierarchical RNN to model this



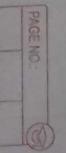
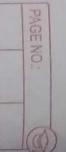
$$h = h_K \quad (K \text{ is number of sentences})$$

$$\text{Decoder: } p(y|d) = \text{softmax}(w^T h + b) - \frac{\alpha}{\text{softmax}(Vh + b)}$$

Parameters: w^T, V, b

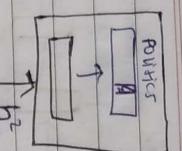
Others same as previous

2 types of models prevs



Data: $\{x_i\}_{i=1}^N$, Doc, Class: y_i

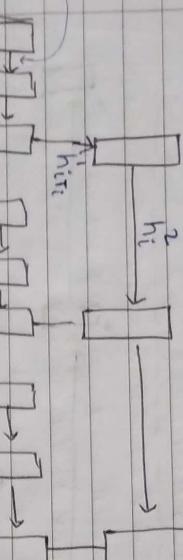
word level encoders: $h_{ij} = RNN(h_{i-1}, x_{ji})$



$$e_{ij} = \text{tanh}(w_{ui} + b_{ui})$$

(:: op is one
only
no previous time step
decodes one word at a time)

$$h_{ij}$$



$$x_{ji}$$

★ Attention in Hierarchical encoder-decoder model

- First we need to attend to important words in sentence.
- Then we need to attend to important sentences in doc.

Output:

$$S = \sum_i d_i h_i$$

$$h_i = RNN(h_{i-1}, S_i)$$

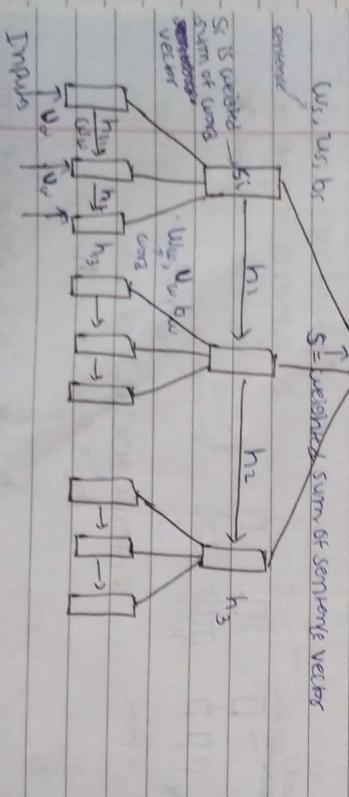
$$e_i = v_i^T \tanh(w_i h_i + b_i)$$

$$d_i = \text{softmax}(e_i)$$

$$S_i = \sum_j d_{ij} h_i$$

\$S_i\$ is representation of sentence
in term of vector

sentence level (w)-encoder:



LSTM & Gated Recurrent Units (GRUs)

NPEL PARAGRAPH

\Rightarrow In ML when ever we want to learnt something always introduce parametric form of these quantity and the learnt the parameter of that function

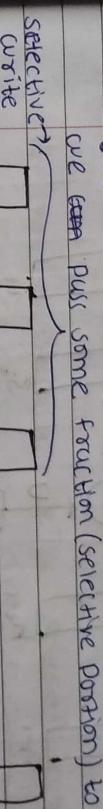
\Rightarrow Remember that there is no explicit supervision here it's just that we have a better modeling choice.

\Rightarrow Recall that RNNs we used s_{t-1} to compute s_t

$$s_t = \sigma(w s_{t-1} + u x_t) \quad \text{ignoring bias}$$

\Rightarrow but now we ~~pass~~^(current) only some important portion to next in this case ~~staircase~~ decision is binary (will given information pass 100% or not, pass 100%) 0 or 1

\Rightarrow but more sensible way of doing this is instead of 0 or 1 we ~~pass~~ pass some fraction (selective portion) to next state



\Rightarrow so we want to do

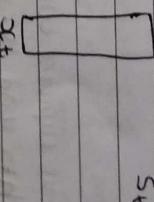
~~Forget~~: Forget the information add by stop words (a, the, etc.)

~~Selective read~~: The information added by previous

Sentiment bearing words (awesome, amazing, etc.) selectively write: selectively write new information from current word of the state

\Rightarrow So what we want to do here is we have computed a state s_{t-1} at time step $t-1$ and now we want to overload it with new information (x_t) and compute a new state (s_t)

$O_{t-1}^{\text{output gate}}$: decides what fraction of s_{t-1} should be pass



\Rightarrow each element of O_{t-1} multiplied with each element of the corresponding element of s_{t-1} and each ele of O_{t-1} is less than 1

\Rightarrow Hence problem is how we compute O_{t-1} ? does RNN knows what function of the state to pass on?

so we introduce parametric form for O_{t-1}

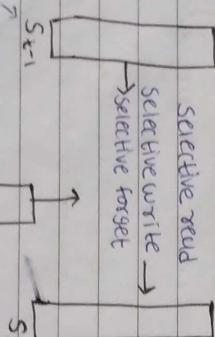
we compute O_{t-1} and h_{t-1} as

$$O_{t-1} = \sigma(w o h_{t-1} + u_o x_{t-1} + b_o)$$

σ is used because we want O_{t-1} between 0 to 1

$$h_{t-1} = O_{t-1} \odot s_{t-1}$$

Op from pre time step



S_t output from t^{th} time step

x_t : Input for t^{th} time step

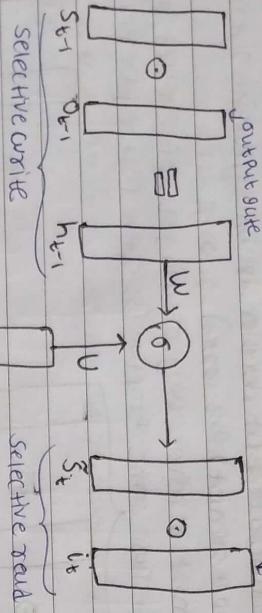


Selective read

\Rightarrow we use h_{t-1} to compute new state at the next time step t
 \Rightarrow we will also use i_t which is input at time step t

$$\tilde{S}_t = \sigma(W h_{t-1} + U i_t + b)$$

$$Tanh(\tilde{S}_t)$$



$\Rightarrow \tilde{S}_t$ thus captures all the information from the previous

state h_{t-1} and the current input i_t

\Rightarrow However, we may not want to use all info and selectively

read from it before constructing the new cell state S_t

\Rightarrow to do that we introduce new gate called input gate

$$i_t = \sigma(W h_{t-1} + U i_t + b_i)$$

\Rightarrow so far we have

previous state : S_{t-1}

output gate : $O_{t-1} = \sigma(W_0 h_{t-1} + U_0 i_t + b_0)$

selective write : $h_{t-1} = O_{t-1} \odot S_{t-1}$

current state : $\tilde{S}_t = \sigma(W h_{t-1} + U i_t + b)$
 (temporary)

input gate : $i_t = \sigma(W h_{t-1} + U i_t + b_i)$

selective read : $i_t \odot \tilde{S}_t$

Selective forget

\Rightarrow how do we combine S_{t-1} and \tilde{S}_t to get the new state
 \Rightarrow here simple way to doing this

$$S_t = S_{t-1} \odot f_t \oplus \tilde{S}_t$$

\Rightarrow But we may not want to use the whole of S_{t-1} but

forget some parts of it

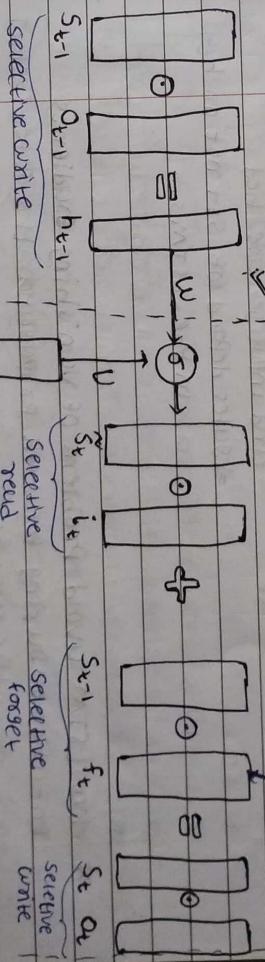
$\star \Rightarrow$ Here you should have some intuition why we use S_{t-1}
 Why not use h_{t-1} and why should we do selectively forget
 For this ~~why~~ that we have already done in selective write

so the answer is selective write is used to
 compute how to read the information but now once you have

compute how to read the new into you want to see how to
 assimilate (adopt) it back with old into that you had

~~forget that why we introduce separate gate~~

forget gate



\Rightarrow we have full set of eqn for LSTM

States:

$$O_t = \sigma(W_0 h_{t-1} + U_0 i_t + b_0), \quad \tilde{S}_t = \sigma(W h_{t-1} + U i_t + b)$$

$$i_t = \sigma(W h_{t-1} + U i_t + b_i) \quad S_t = f_t \odot S_{t-1} + i_t \odot \tilde{S}_t$$

$$f_t = \sigma(W h_{t-1} + U i_t + b_f) \quad h_t = O_t \odot S_t$$

GRU

\Rightarrow the full set of eqn for GRUs \Rightarrow GRU is type of RNN

Gates: $O_t = \sigma(w_o s_{t-1} + w_{oI} i_t + b_o)$ but in certain cases, $i_t = \sigma(w_i s_{t-1} + w_{iI} i_t + b_i)$ advantages over LSTM

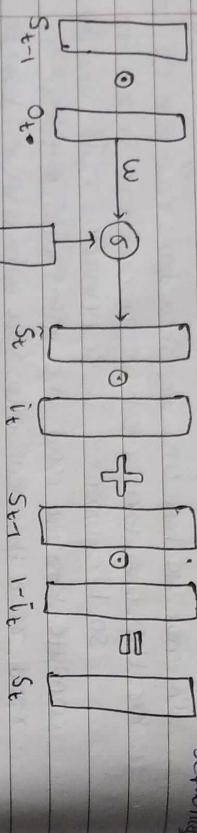
States: $\tilde{S}_t = \sigma(w_o (O_t \odot S_{t-1}) + w_{iI} i_t + b)$ it uses less memory and LSTM better accurate when using dataset with longer sequences

$S_t = (1 - i_t) \odot S_{t-1} + i_t \odot \tilde{S}_t$ LSTM better accurate when using dataset with longer sequences

it uses less memory and LSTM better accurate when using dataset with longer sequences

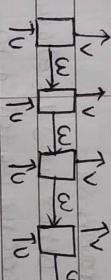
\Rightarrow It state at time $t-1$ did not contribute much to the state at time t (i.e. $I_t || h_{t-1} \rightarrow 0$ and $|h_{t-1}| \rightarrow 0$) then during backprop the gradient flowing into s_{t-1} will vanish \Rightarrow this kind of vanishing is OK because if s_{t-1} didn't contribute to s_t we don't want to hold it responsible for crimes of s_t

\Rightarrow so in LSTM key difference from RNN vs Vanilla RNN is that flow of flow of info and gradient controlled by the gates which ensure that gradients vanish only when they should. (i.e. when s_{t-1} didn't contribute much to s_t)



\Rightarrow no explicit forget gate the forget i_t and input get intertwined

\Rightarrow gates depend on s_{t-1} not on h_{t-1}



\Rightarrow recall RNNs had multiplication terms which limited the gradient to vanish

$$\frac{\partial L(\theta)}{\partial w} = \frac{\partial L_t(\theta)}{\partial s_t} \sum_{j=1}^{t-1} \frac{\partial s_{t+j}}{\partial s_j} \frac{\partial s_j}{\partial w}$$

\Rightarrow in particular, if the loss at $L_t(\theta)$ is high because s_t was not good enough to compute s_{t+1} correctly then into will not be propagated

* How LSTM avoid problem of vanishing gradients

Intuition: During forward propagation the gates control

the flow of information and they prevent information from being written to state

\Rightarrow similarly during backpropagation they control the flow of gradients

$S_t = f_t \odot S_{t-1} + i_t \odot \tilde{S}_t$

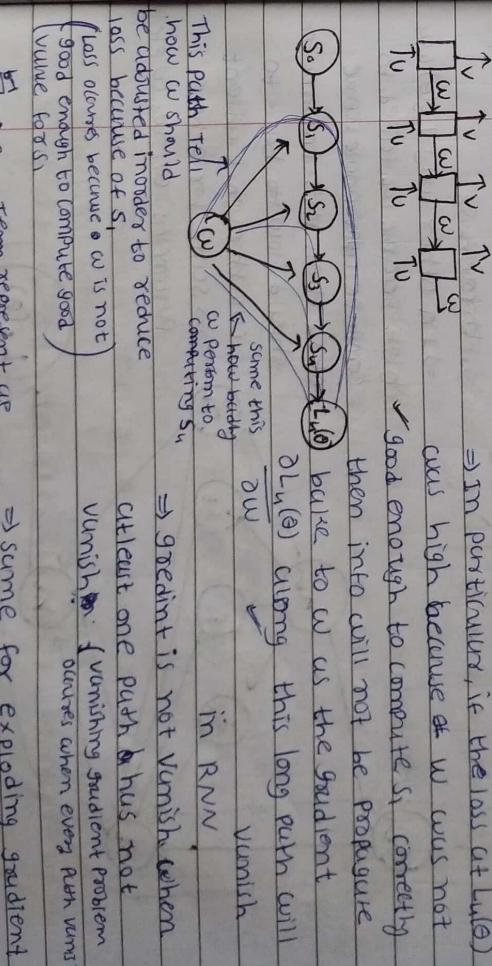
Here f_t decides how much goes on forward direction

and for backward prop $\frac{\partial S_t}{\partial S_{t-1}} = f_t$ (let neglect w)

so it also decides how much $\frac{\partial S_t}{\partial S_{t-1}}$ goes on backward direction.

\Rightarrow here vanishing occurs in forward us well as backward prop so this types of vanishing is OK. It in forward pass S_t is not contribute to s_t than

us well as backward prop so this types of vanishing is OK. It in forward pass S_t is not contribute to s_t than



\Rightarrow gradient is not vanish when it can't one path has not

\Rightarrow gradient is not vanish when every path vanish

(loss doesn't because w is not good enough to compute good value to s_n)

\Rightarrow some for exploding gradient

in the case

→ LSTMs & there exists at least one path through which the gradients can flow (and hence no vanishing gradients)

⇒ starting from h_{k-1} and s_{k-1} we have reached h_k and s_k and the recursion will now continue till the last timestep

⇒ Hence we are showing only few parameters.

not showing U, U_b, U_t, U_i

⇒ we use now interest

in knowing if the

gradient from $L_t(\theta)$

flow back unvanishing

time step k

⇒ for example we are

interested in knowing

if gradient flows to w_t

through s_k

⇒ in other way $L_t(\theta)$ was

high because w_t failed to

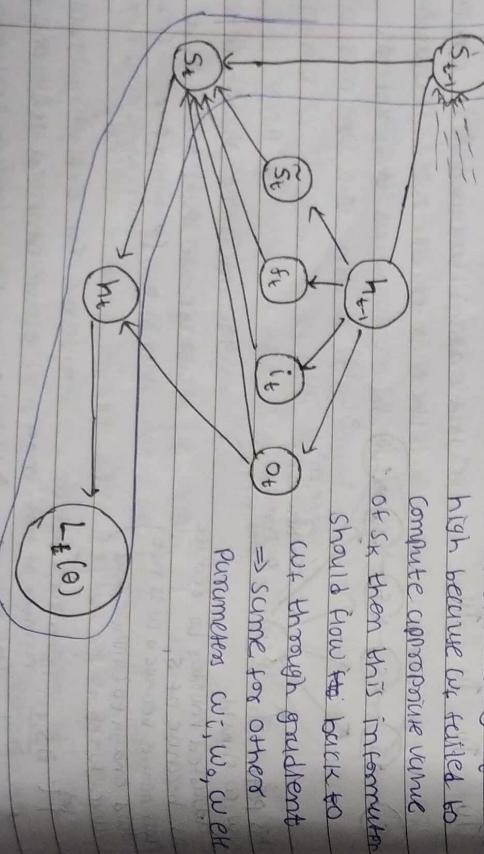
compute appropriate value

of s_k then this information

should flow back to

w_t through gradient

parameters w_i, w_o, w_h



→ These are multiple paths (many many) from $L_t(\theta)$ to s_k (you just need to reverse the direction of arrows for back prop)

⇒ many path like s_{k-1} another path through h_k and further there is multiple path to reach h_k itself.

→ let consider one path and let denote gradient along this path to,

$$\frac{\partial \ell}{\partial h_t} = \frac{\partial L_t(\theta)}{\partial h_t} \frac{\partial h_t}{\partial s_k} \frac{\partial s_k}{\partial s_{k-1}} \dots \frac{\partial s_k}{\partial s_1}$$

This will not vanish because h_t directly connected to $L_t(\theta)$ and there is no intermediate node which can cause gradient vanish

$$h_t = O_t \circ S_t$$

$h_t = [h_{t1} \ h_{t2} \ \dots \ h_{tn}]$

$$O_t = [O_{t1} \ O_{t2} \ \dots \ O_{tn}]$$

$$S_t = [S_{t1} \ S_{t2} \ \dots \ S_{tn}]$$

h_t only depend on one S_t

$$\frac{\partial h_t}{\partial s_t}$$

and only diagonal term remain where $i=j$

$$\frac{\partial h_t}{\partial s_t} = D(O_t \circ S_t^*)$$

diagonal matrix $\in \mathbb{R}^d$ representation of diagonal matrix which diagonal vector inside D

$s_k = f_t \circ S_{t-1} + i_t \circ \tilde{S}_t \leftarrow \tilde{S}_t$ depends on s_{t-1} so we

can't treat as const $\frac{\partial s_t}{\partial s_{t-1}}$

→ we dealing with ordered network so thus $\frac{\partial s_t}{\partial s_{t-1}}$ will be a sum of an explicit term and implicit term

(treating \tilde{S}_t as const)

⇒ now making const (use assumption implicit term vanish)

so explicit term f_t

⇒ gradient from explicit term is given us $D(f_t)$

$$\begin{aligned}
 t_0 &= \frac{\partial L(\theta)}{\partial h_t} \frac{\partial h_t}{\partial s_t} \frac{\partial s_t}{\partial s_{t-1}} \dots \frac{\partial s_t}{\partial s_k} \\
 &= L'(h_t) D(\overbrace{a_t \odot s_t}) D(f_t) D(f_{t-1}) \dots D(f_{k+1}) \\
 &= L'(h_t) D(a_t \odot s_t) D(f_t) \underset{i=k+1}{\underbrace{D(a_i \odot s_i)}_{f_i}}
 \end{aligned}$$

- This term does not vanish → contain multiple of forget gates
- The forget gate thus regulate the gradient flow depending on explicit contribution of a state (s_t) to the next State (s_{t+1})

→ It during forward pass it didn't contribute much to s_{t+1} (because $f_t \rightarrow 0$) then during backprop also gradient will not reach s_t so it did not contribute much to s_{t+1} because if it did not contribute much to s_{t+1} then during back propagation (f_t does the same regulation during forward pass and back pass which fails)

⇒ so in LSTM gradient will vanish only when required otherwise it will not vanish like RNNs.

* why LSTM do not solve exploding gradient

- ⇒ we will see the path through which the gradient can explode
- ⇒ let's start from θ
- gradient along this path is $t_0 = \frac{\partial L(\theta)}{\partial h_{k-1}}$ gradient along that which contains this $\frac{\partial h_{k-1}}{\partial h_{k-1}}$ gradient

path

$$t_1 = \frac{\partial L(\theta)}{\partial h_t} \left(\frac{\partial h_t}{\partial s_t} \frac{\partial s_t}{\partial s_{t-1}} \right) \dots \left(\frac{\partial h_t}{\partial s_k} \frac{\partial s_k}{\partial s_{k-1}} \right)$$

$$= L'(h_t) [P(s_{t-1} \odot s_t) \cdot w] \dots [P(s_{k-1} \odot s_k) \cdot w]$$

$|w| \leq \|L'(h_t)\| \|w\| \|w\| \dots \|w\|$ t^{k-1} this is diagonal matrix

⇒ so dependency on norm of w , the gradient may explode

→ so to solve this problem popular trick is to use gradient clipping

⇒ while back prop if the norm is exceed a certain value it is scaled to keep its norm within an acceptable threshold.

⇒ This is fine because gradient care about direction not for a magnitude, when gradients are not manage manageable in terms of their magnitude the we clip it to some manageable while be useful of the direction so for exploding gradient is easy but in case of vanishing gradient we does not have direction because either gradient become zero.



Intuition behind self-Attention

→ Attention is most important part of an input
 ⇒ If we want to find video on youtube about deep learning then we have type in search bar

[Q] deep Learning

Machine Learning

Key (K_3)

deep learning
key (K_2)
value (V)
dict [u] = 2

dict [u]: 2, b: 3, c: v

word

sentence

um →

bus?

word

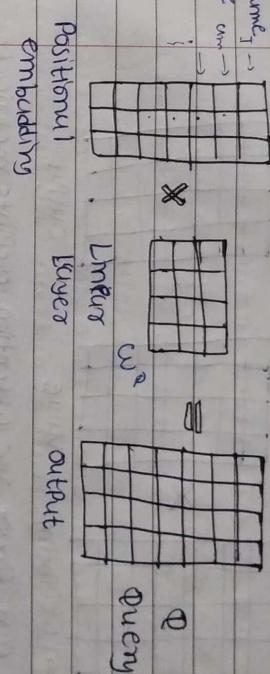
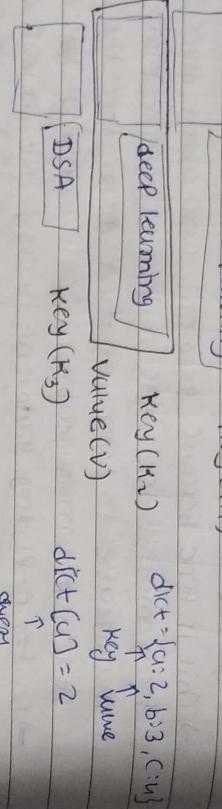
int

key

value

query

output



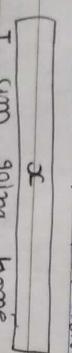
→ for every video in database there some key information related to that video its called **key**
 → and now • Algorithm will compare how keys and query use related to each other. This algo do similarity computation

→ Final step now we know which key is relevant

so extracting the information what we want to pay attention to.(video)

⇒ Learning self attention with neural networks
 Goal: identify and attend to most important features in input

2) encode positional information



embedding

I am going home

Datu is fed in all at once!

Positional info P_0, P_1, P_2, P_3 into to understand order

Need to encode position!

softmax $(Q \cdot K^T / \text{Scaling})$

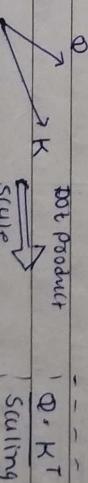
Positional-aware encoding

attention weighting

Attention weighting: where to attend to!
 How similar is the key to the query?

I am going home ⇒ This operation gives us the score which define how the component of input data related to each other.

similarity metric (cosine similarity)

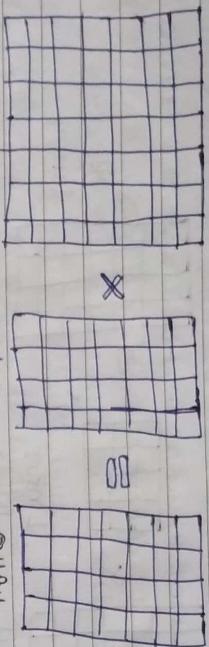


Attention score: compute pairwise similarity between each query and key

$$\begin{matrix} Q & \xrightarrow{\text{dot product}} & Q \cdot K^T \\ \downarrow \text{scale} & & \downarrow \text{scaling} \\ \text{similarity metric (cosine similarity)} \end{matrix}$$



v) self-attend to extract features



Attention weighting

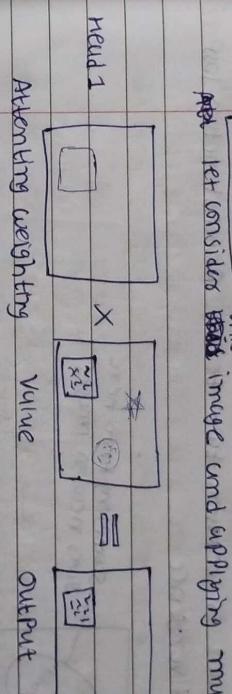
value

Output

now we can finally use attention weights metric to extract features that use decreasing or high attention

$$\text{softmax} \left(\frac{Q \cdot K^T}{\text{scaling}} \right) \cdot V = A(Q, K, V)$$

let consider this image and applying multiple self-attention heads



others and picking out others part of the data ★ and ☺ etc.

→ ★ why do we use scale self-attention? $\text{softmax} \left(\frac{Q \cdot K^T}{\sqrt{d_k}} \right) \cdot V$

low dimension \rightarrow dot product \rightarrow low variance
high dimension \rightarrow dot product \rightarrow high variance \rightarrow dim of weight \rightarrow mean distribution of numbers is different \rightarrow drops to zero

(large range) for high dimension so dot product \rightarrow max at steps it will neglect small values and give more value to big numbers

$$\text{softmax} = \frac{e^{x_i}}{\sum_j e^{x_j}} \cdot \frac{v_i}{\|v_i\|} \rightarrow \underbrace{\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}}_{\text{softmax}} \cdot \underbrace{\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}}_{\text{Input}} = \underbrace{\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}}_{\text{Output}}$$

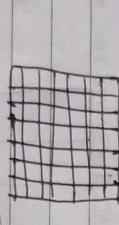
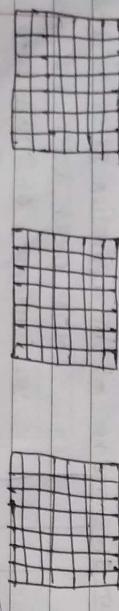
→ And vanishing gradient occurs so to reduce high variance to low softmax prob become comparable so we need to scale.

$1D \rightarrow \text{Value}$, $X = \text{random variable}$

$2D \approx 2 \text{Var}(x)$ we need variance const $\text{Var}(x)$ for any dimension
 $3D \approx 3 \text{Var}(x)$ $y \rightarrow \text{Var}(y) \rightarrow 2 \text{Var}(x)$
 $dD \approx d \text{Var}(x)$ $\frac{y}{d} \rightarrow \frac{1}{d} \text{Var}(y) \rightarrow \frac{1}{d} \text{Var}(x)$
If $y = cx$
 $\text{Var}(y) = c^2 \text{Var}(x)$

Positional Encoding

single self-attention head



money bank grows

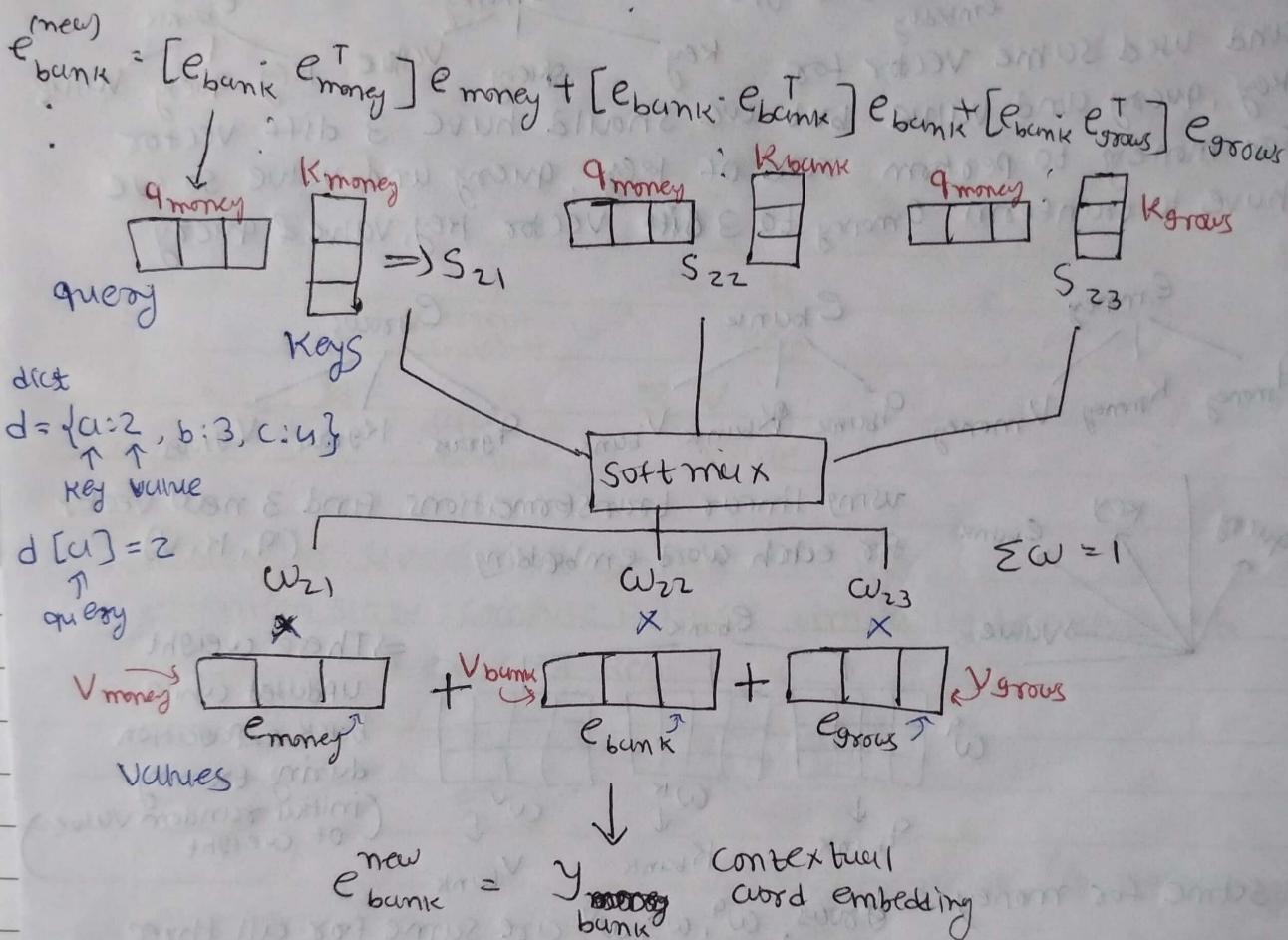
$$\text{money} = 0.7 \text{ money} + 0.2 \text{ bank} + 0.1 \text{ grows}$$

$$\text{bank} = 0.25 \text{ money} + 0.7 \text{ bank} + 0.05 \text{ grows}$$

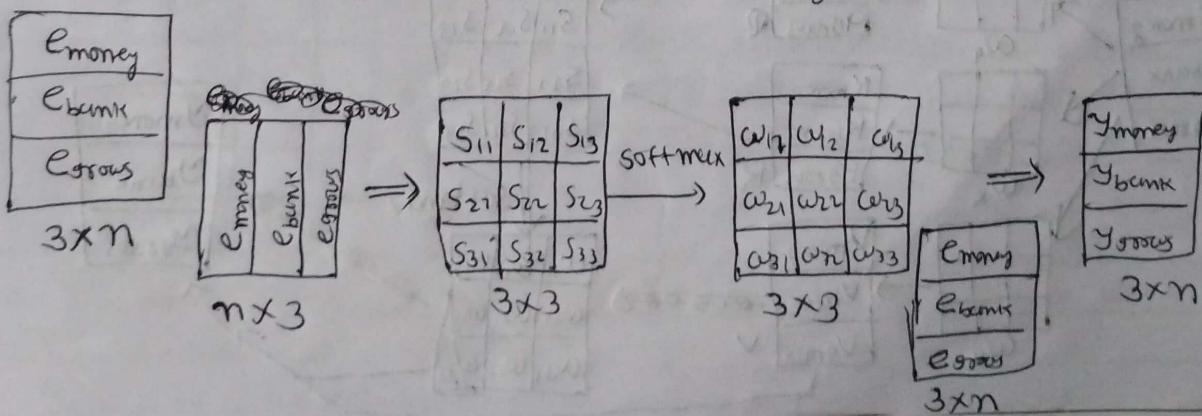
$$\text{grows} = 0.1 \text{ money} + 0.2 \text{ bank} + 0.7 \text{ grows}$$

\Rightarrow In mathematics form

$$e_{\text{money}} = 0.7 e_{\text{money}} + 0.2 e_{\text{bank}} + 0.1 e_{\text{grows}}$$



same way can find y_{money} & y_{grows}

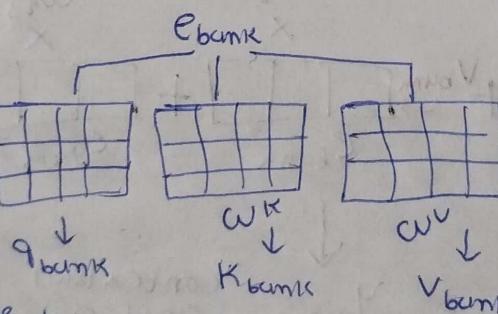
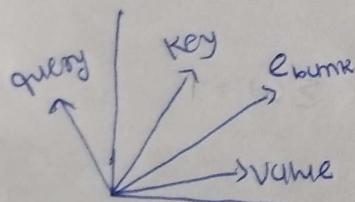
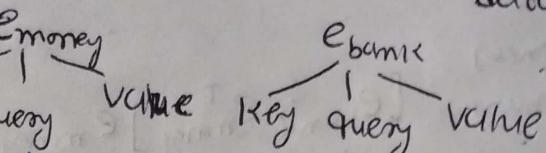
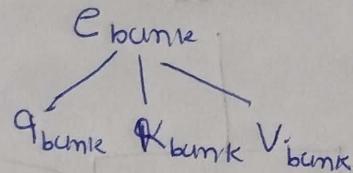
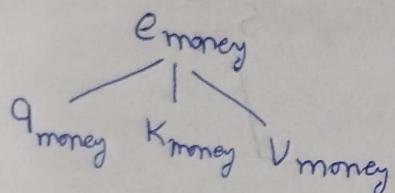


- ★ Points to consider
 - This operation is a parallel operation
 - There are no learning parameters involved (generated new embedding is not learnt from data that we have)
 - So this simple attention model give general contextual embedding but we need task-specific embedding for better result.

So we introduce W & b so self attention model can learn from data.

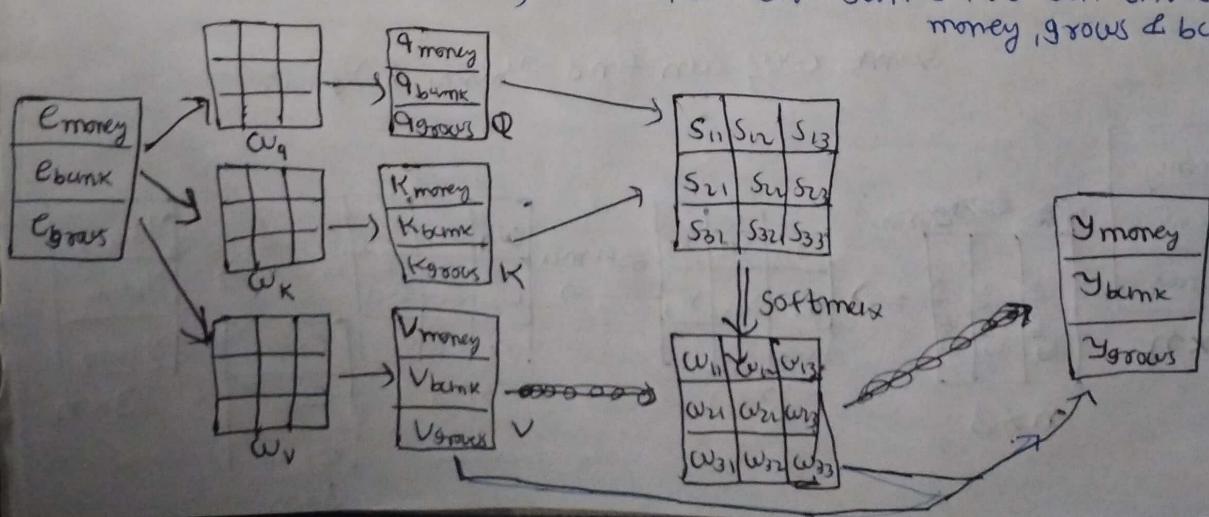
⇒ We use each word as 3 type embedding

and used same vector for key, query and value but we should have 3 diff vector for money to perform role of key, query and value so we have to transform e_{money} to 3 diff vector key, value & query.



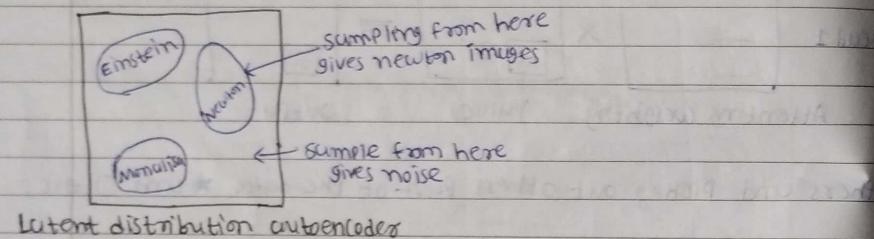
⇒ These weight updated while back propagation during training (initial random values of weight)

Same for money & grows, w^q, w^k, w^v are same for all three money, grows & bank



Variational Autoencoder

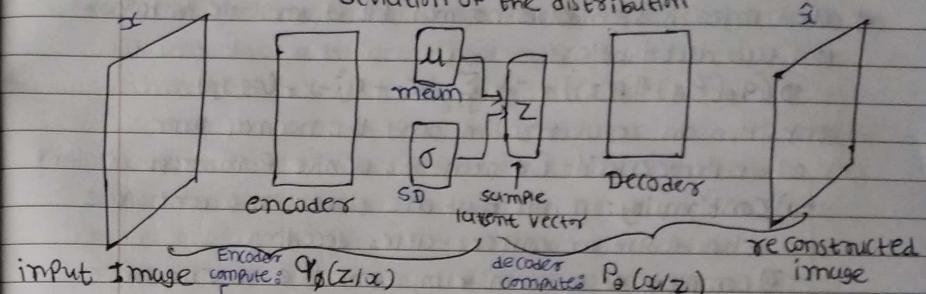
- ⇒ idea behind VAE is that instead of mapping any input to a fix vector we map our input into a distribution
- ⇒ The difference in VAE is that normal bottleneck vector z replace by two separate vectors one representing the mean of distribution and other representing standard deviation of that distribution
- ⇒ And know fed a sample from distribution to decoder
- ⇒ Loss function is contain two terms Reconstruction loss and KL divergence



How to define the latent variable space?

- ⇒ Instead of mapping input to a fixed latent vector, we map it to a distribution
- ⇒ Force latent variable to be normally distributed so its easy for us to define normal distribution using mean and variance
- ⇒ instead of passing encoder output to the decoder, use mean and standard describing the distribution.
- ⇒ now we have only two variable that need to be trained not entire distribution
- ⇒ Quantify the distance betⁿ learned distri and standard normal distri using Kullback-Leibler (KL) divergence
- ⇒ so while training we are going to force a normal distribution as close as possible to standard normal distribution minimize one of the loss which is given by KL divergence. and there is other loss is due to reconstruction loss

Mean and standard deviation of the distribution



⇒ Here we use reparameterization trick for back propagation

Learned parameters
during back propagation

$$Z = \mu + \sigma \theta \epsilon \quad \begin{array}{l} \text{mean} \\ \text{standard deviation} \end{array} \quad \begin{array}{l} \rightarrow \text{standard normal distribution} \\ \rightarrow \text{Fixed space where we randomly sample} \\ \rightarrow \text{not learned during BP.} \end{array}$$

$$L(\phi, \theta, x) = \text{reconstruction loss} + \text{regularization term}$$

$$\rightarrow \text{like log likelihood, } \|x - \hat{x}\|^2$$

$\mathcal{D}(q_{\phi}(z|x) || p(z))$
inferred latent fixed prior on latent distri (some initial hypothesis about what that latent variable look like)

Regularization term
capturing dis betⁿ encoding of latent variable and prior hypothesis

- ⇒ So during training we are trying that each of latent variables adopts a probability distri. that similar to that prior
- ⇒ common choice of prior is Normal Gaussian

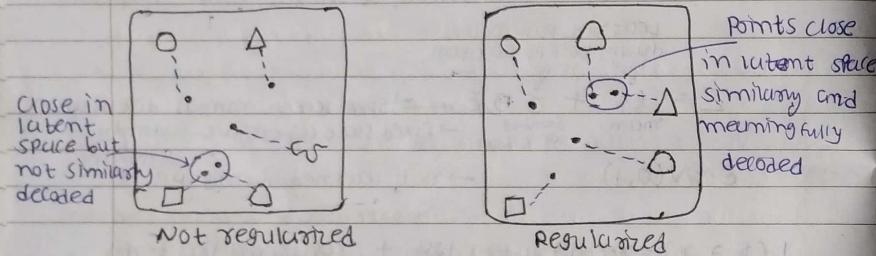
$$p(z) = N(\mu=0, \sigma^2=1)$$

Encourages encoding to distribute encoding evenly around the center of the latent space.

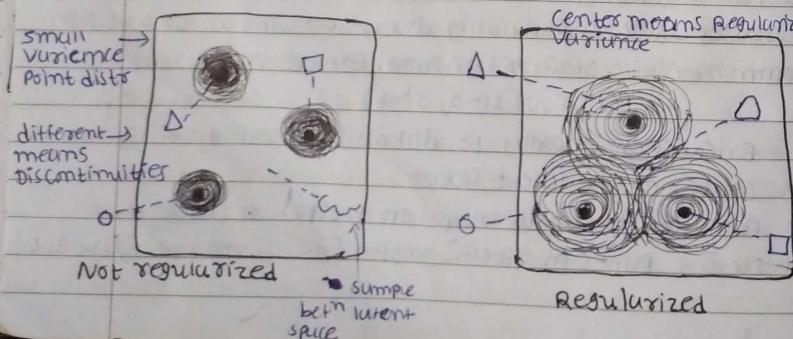
Penalize the network when it tries to cheat by clustering points in specific region (i.e. by memorizing the data).



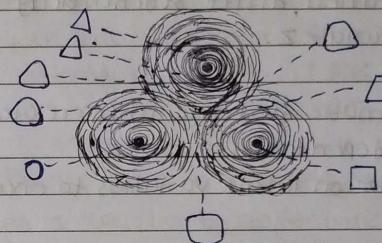
- ⇒ When Prior is standard normal, KL divergence between the two distribution
- $$D(q_\phi(z|x) || P(z)) = -\frac{1}{2} \sum_{j=0}^{K-1} (\sigma_j + \epsilon_j^2 - 1 - \log \sigma_j)$$
- ⇒ Intuition on regularization and the normal prior what properties do we want to achieve from regularization?
- 1) Continuity: points that are closer in latent space should similar content after decoding
 - 2) Completeness: sampling from the latent space should have meaningful content after decoding



- ⇒ with regularization we are able to achieve that by ~~time~~ training to minimize regularization term its not sufficient to just employ the reconstruction loss alone to achieve continuity and this completeness
- ⇒ without regularization just encoding and reconstructing does not guarantee the properties of continuity and completeness.



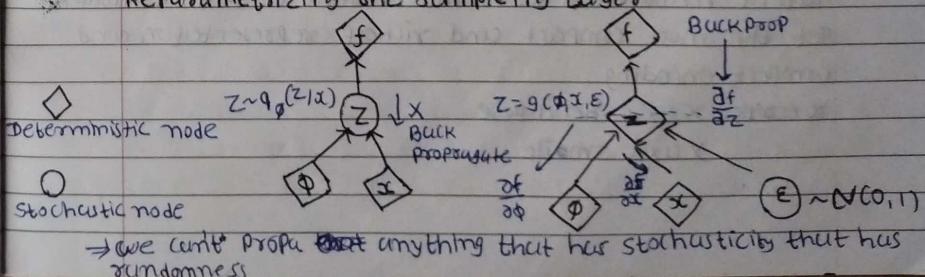
- ⇒ Normal Prior allows us for learned dists. of those latent variables to effectively overlap in latent space because everything is regularized to have according to this prior mean zero standard deviation one and that centers the mean regularizes the variances for each of the independent latent variable distributions.
- ⇒ by adding effect of regularization we can achieve continuity and completeness in the latent space



- ⇒ when we have notion of randomness of probability we we can't sample directly through that layer instead with reparametrization what we do is we redefine how a latent variable vector is sampled as a sum of the fix deterministic mean μ & a fixed vector of SD σ .
- ⇒ and now we divert all the randomness on the sampling to a random constant ϵ drawn from a normal distribution

$$z = \mu + \sigma \epsilon \quad \epsilon \sim \mathcal{N}(0, 1)$$

Reparametrizing the sampling layer



→ reparametrization allows us to shift this diagram where now we completely diverted that sampling operation off to the side to this const E which is drawn from normal prior.

→ now when we look back at our latent variable it is deterministic w.r.t that sampling operation what this mean is that we can back propagate to update neural network weights completely end to end without worry about direct randomness, direct stochasticity within those latent variables.

Latent perturbation:

→ we want latent variable that are uncorrelated with each other and independent

→ Enforce diagonal prior on latent variables to encourage independence

→ for example slowly increase or decrease a single latent variable keep all other variable fixed we can see Head pose is shifting and this is due to perturbation of single latent variable

→ turning the single variable and we can see how that particular latent variable affects the decoded reconstruction

→ different dimensions of z different interpretable latent features

→ we can even multiple latent variable simultaneously compare one against the others and ideally we want those latent features to be as independent as possible in order to get the most compact and richest representation and compact encoding

X axis head pose

Y axis smile

Latent space disentanglement with β -VAEs

β VAE loss:

$$L(\theta, \phi; x, z, \beta) = \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - \beta D_{KL}[q_{\phi}(z|x) || p(z)]$$

Recon. term

→ weighting constant β controls how powerful that regularization term is in the overall loss of VAE

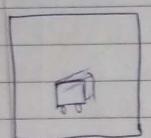
→ By increasing value of β you can try to encourage greater disentanglement more efficient encoding to enforce this latent variable to be uncorrelated with each other.

→ let look into face reconstruction as a task of interest with standard VAE ($\beta=1$): pose and rotation of the head correlated with smile as the head pose is changing apparent smile is also changing

but with β VAE imposing β value much much greater than 1 we can try to enforce greater disentanglement now we can consider only single latent variable head pose and the smile is constant compared to standard VAE.

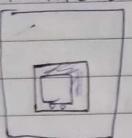
Object Detection / Localization

Image classification



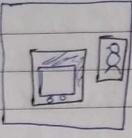
'car'

Classification with localization

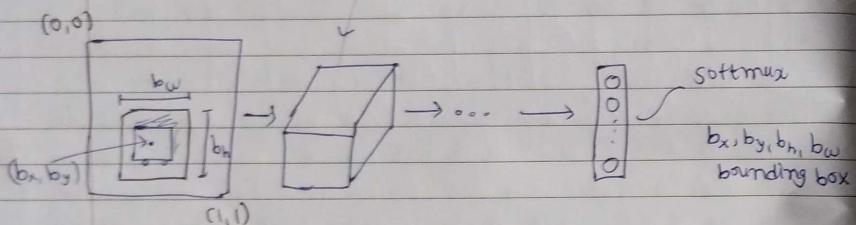


1 object

Detection



multiple objects



Define the target label

y_1	y_2	y_3
b_x	b_x	b_x
b_y	b_y	b_y
b_h	b_h	b_h
c_1	c_2	c_3
y_1	y_2	y_3

If there is object class 1, 2 or 3

1 - Pedestrian
2 - Car
3 - Motorcycle
4 - background & nothing

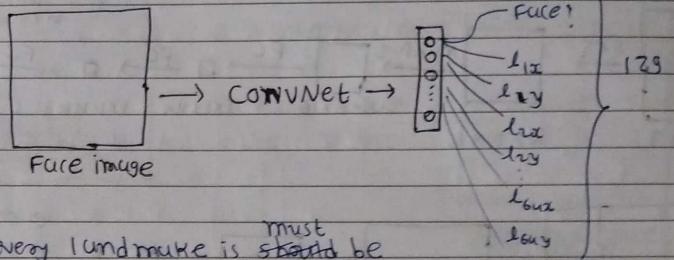
For example

$$x \quad y = \begin{bmatrix} b_x \\ b_y \\ b_h \\ 0 \\ 0 \end{bmatrix}$$

$$x \quad y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{LOSS} = L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_r - y_r)^2, & \text{if } y \neq 0 \\ (\hat{y}_1 - y_1)^2, & \text{If } y = 0 \end{cases}$$

⇒ Landmark detection



⇒ every landmark is ~~should~~ must be consistent in all images

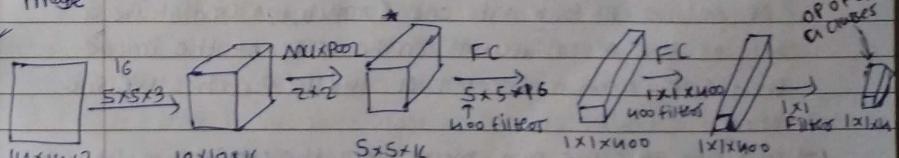
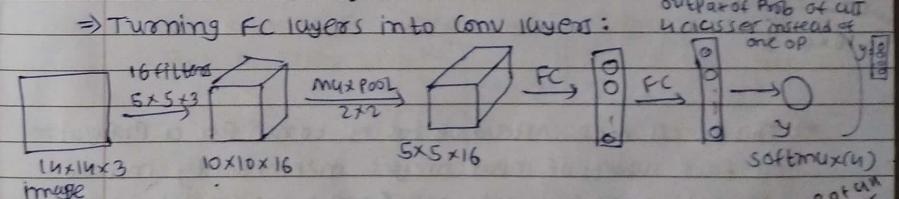
Ex (l_{1x}, l_{1y}) contain left side of eyes (corner landmark)
Table must be consistent

⇒ Car detection example

In this we have to give many labeled data for train
So we take different size of window and try to detect
car in every position of given image (here we use
some stride). This also called Sliding window detection.

⇒ disadvantage of Sliding window algo is its take much
more computation time

⇒ Turning FC layers into Conv layers:



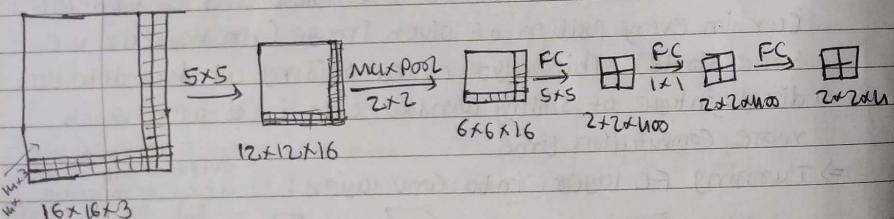
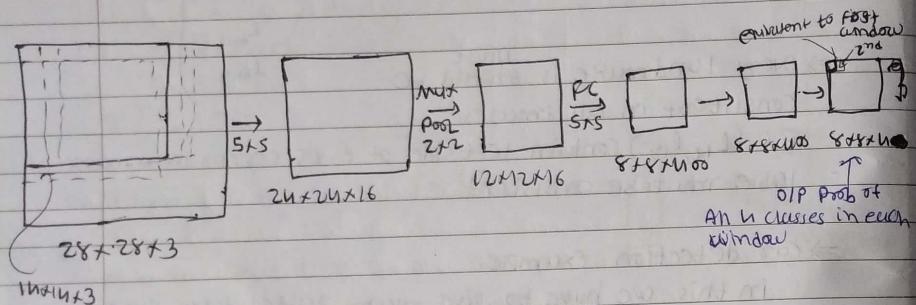
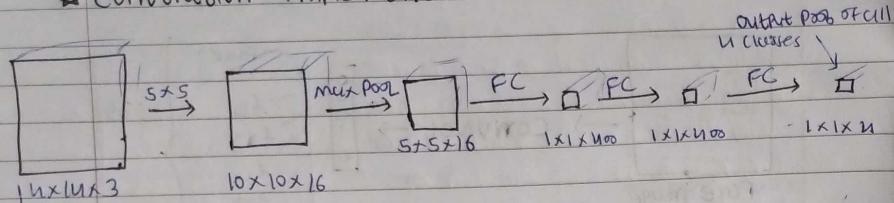
each of 100 values is
some arbitrary linear function
of this 5x5x16 activations from
the previous layer *

$\Rightarrow ?$ means some numbers since $P_c = 0$ so no need to worry about that number that number does not signify anything

PAGE NO.:

PAGE NO.:

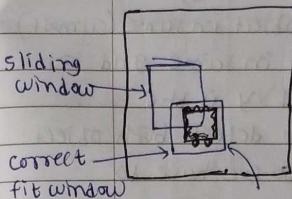
* Convolution implementation of sliding window



\Rightarrow This process does instead of focusing to run for a propagation on four subsets of input images independently instead it combines all four into one computation and share a lot of the computation in the regions of the image that are common for all 4 16x16 patches we see here.

\Rightarrow so using this we can make all the computation at same times and no need to compute same things many times

* Intersection over union (IOU)



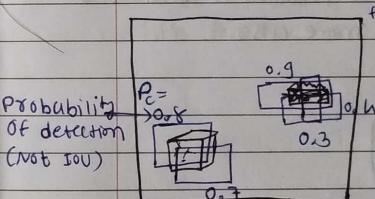
Intersection over area = $\frac{\text{size of intersection}}{\text{size of union}}$

Sliding window is correct for given object If $\text{IOU} \geq 0.5$

Threshold value can take according to your choice

* Non-max suppression

\Rightarrow algo may find multiple detection of same object non-max suppression is a way to make sure that your algo detect the object only ones



\Rightarrow let take example for only one object detection: car detection

$$\text{Each Prediction is } \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_w \\ b_h \end{bmatrix}$$

Algorithm i) discard all boxes with $P_c < 0.6$ (low prob of box)

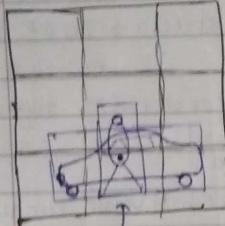
ii) while there is a any remaining boxes:

- pick the box with largest P_c output that us prediction
- discard any remaining box with $\text{IOU} > 0.5$ with the box output in the previous step.

algo for some multiple obj ex: multiply cars
only car not covering both



★ Overlapping Objects



This grid contains mid point of both

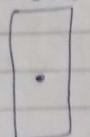
- C_1 = pedestrian
- C_2 = car
- C_3 = background (nothing)
- ⇒ mid point of both object in same (almost) location and both in same grid cell
- For that grid cell y is the
- not be able to detect both object

$y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$

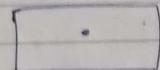
so we have to pick one of the two detections detection as output

- With the idea of Anchor boxes we going to predefine two different shapes called Anchor boxes.
- And now we able to associate two predictions with the two anchor boxes and can use more also

Anchor box 1



Anchor box 2



$y =$

Anchor box 1

$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$

Anchor box 2

$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$

Previously: Each object in training is assigned to grid cell that contains that object's midpoint

$$OIP = 3 \times 3 \times 8$$

(here

$$\text{box grid each } y = \begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

with two anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint and anchor boxes

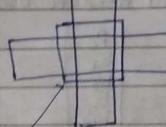
for grid cell with highest IOU

$$OIP = 3 \times 3 \times 16$$

Anchor 2×8 x size

✓

Anchor boxes



grid

(anchor box)

We see which has higher IOU with respect to ground truth bounding box and we assign that object to grid size and anchor box. For training labeled output.

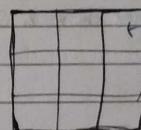
$$OIP = 3 \times 3 \times 16$$

size

⇒ two anchor boxes with 3 object in same cell: this case is not handle by only two anchor box

⇒ two object in same cell and have same anchor box size: this case also does not handle by this algorithm

with Anchor box

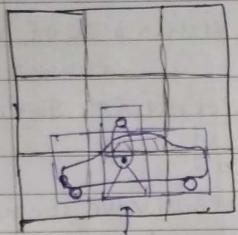


every box has OIP

$$y = \begin{bmatrix} P_c \\ C_1 \\ C_2 \\ C_3 \\ P_c \\ C_1 \\ C_2 \\ C_3 \\ P_c \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

↑ anchor box 1
↑ anchor box 2

★ Overlapping Objects



This grid contains mid point of both

C_1 = pedestrian

C_2 = car

C_3 = background (nothing)

⇒ mid point of both object in same (almost) location and both in same grid cell

→ For that grid cell y is the not be able to detect both object

$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}$

so we have to pick one of the two detections detection as output

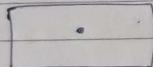
→ with the idea of Anchor boxes we going to predefine two different shapes called Anchor boxes.

→ And now we able to associate two predictions with the two anchor boxes and can use more also

Anchor box 1



Anchor box 2



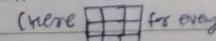
$y =$

Anchor box 1

Anchor box 2

Previously: Each object in training is assigned to grid cell that contains that object's midpoint

OIP $y = 3 \times 3 \times 8$



Box grid over $y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ c_1 \\ c_2 \\ c_3 \end{bmatrix}^8$

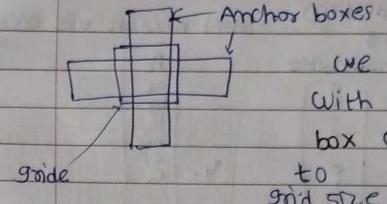
with two anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint and anchor boxes

for grid cell with highest IOU

OIP $y = 3 \times 3 \times 16$

Anchor box $2 \times 8 \times 8$

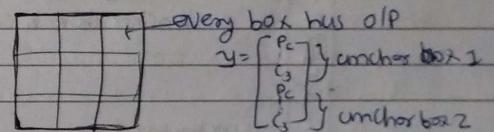


OIP $y = 3 \times 3 \times 16$

⇒ two anchor boxes with 3 object in same cell: this case is not handle by only two anchor box

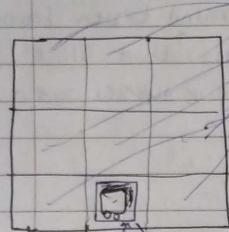
⇒ two object in same cell and have same anchor box size: this case also does not handle by this algorithm

with Anchor box



* YOLO Algorithm

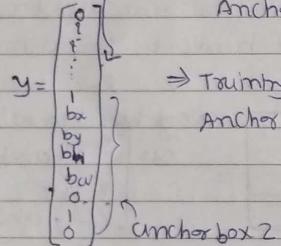
C_1 = Pedestrian
 C_2 = car
 C_3 = motorbike



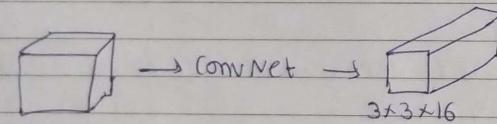
Training labeled
(given in training
OP table)

Anchor box 1

Anchor box 2



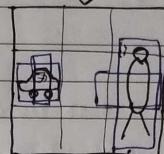
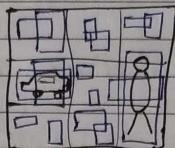
⇒ Training labeled has higher IOU with
Anchor box 2. So we choose anchor box 2.



don't need to use background
class when

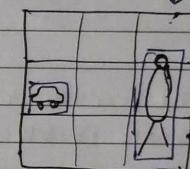
$$y = \begin{bmatrix} 0 \\ 1 \\ bx \\ by \\ bw \\ bh \\ C_1 \\ C_2 \\ C_3 \\ P_C \\ bx \\ by \\ bw \\ bh \\ C_1 \\ C_2 \\ C_3 \\ P_C \end{bmatrix}$$

* Outputting the non-max suppressed output



⇒ For each ~~for~~ grid cell, get 2
Predicted bounding boxes
⇒ get rid of low prob predictions.
For object

⇒ For each class (pedestrian, car,
motorbike) use non max suppression
to generate final prediction



Reinforcement Learning

MIT 2020
MIT 6.S191 Reinforcement learning
57.33 min

Supervised Learning unsupervised Learning Reinforcement Learning

Data: (x, y)

Data: x

Data: State-Action pair

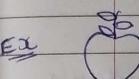
x is data, y is label

no labels

Goal: Learn function
to map $x \rightarrow y$

Goal: Learn underlying
structure

Goal: maximize future
rewards over many
time step



This is apple



This thing is like
other things



Eat this thing because it
will keep you alive

Key Concepts:

Agent: take actions

Environment: The world in which the agent exists and

Actions: a move that agent can
make in environment

Action Space A: The set of all possible actions that agent
can make in environment

Observation: of environment after taking actions.

State: a situation which agent perceives

Reward: feedback that measures the success or failure
of the agent's action.

Total Reward
(Return)

$$R_t = \sum_{i=t}^{\infty} r_i = r_t + r_{t+1} + r_{t+2} + \dots$$

Discounted

Total Reward
(Return)

$$R_t = \sum_{i=t}^{\infty} r^i r_i = r^t r_t + r^{t+1} r_{t+1} + \dots + r^{t+n} r_{t+n} + \dots$$

r : discount factor

$$0 < r < 1$$

greedy ness

→ γ is designed for make future rewards essentially worth less than rewards that we might see this instance

→ Total reward R_t is discounted sum of all rewards obtains from time t

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$$

state action

$$Q = (s_t, a_t) = E [R_t | s_t, a_t]$$

→ Q function captures the expected total future reward an agent in state s_t , can receive by executing a_t certain action a_t

→ Ultimately the agent needs a policy $\pi(s)$, to infer the best action to take at its state s

Strategy: the policy should choose an action that maximizes future reward

$$\pi^*(s) = \arg \max_a Q(s, a)$$

Value Learning

Find $Q(s, a)$

$$u = \arg \max_a Q(s, a)$$

Policy Learning

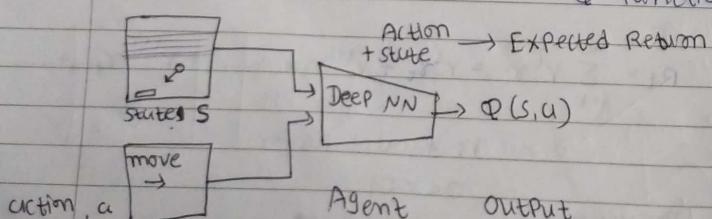
Find $\pi(s)$

sample $a \sim \pi(s)$

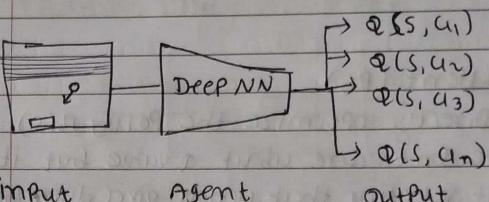
Value Learning algo
How can we use deep neural networks to model input

Deep Q Networks (DQN) Policy $\pi(s)$
use to infer the optimal Q function

Q -functions?



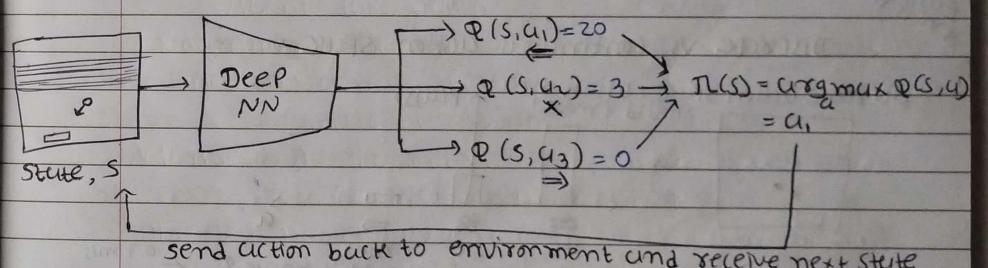
State \rightarrow Expected return for Each action



$$L = E \left[\left\| \underbrace{(\gamma + \max_a Q(s', a'))}_{\text{target}} - \underbrace{Q(s, a)}_{\text{predicted}} \right\|^2 \right]$$

Q-Loss

target: initial reward + γ times max Q value for next state
predicted: Q value for current state and action



send action back to environment and receive next state

★ Downsides of Q -Learning

take only those action which give best result at that time so it repeat same action and don't learn diff actions

Complexity: can model scenarios where the action space is discrete (and small)
cannot handle continuous action spaces

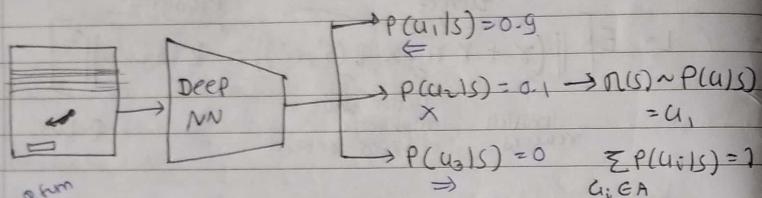
Flexibility: Policy is deterministically computed from the Q function by maximizing the reward
→ cannot learn stochastic policies

⇒ To address these consider a new class of RL training algorithms: policy gradient methods

Policy learning algorithm

Policy Gradient (PG)

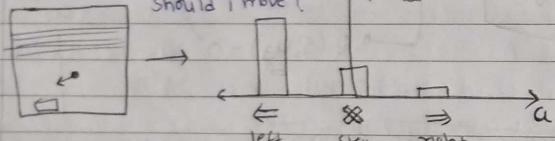
Policy gradients: Directly optimize the Policy $\pi(s)$
 ⇒ in this approach we don't care about Q value but it's give Probability of selecting that action gives desirable output



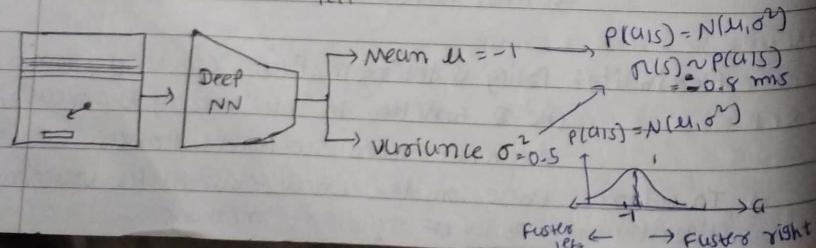
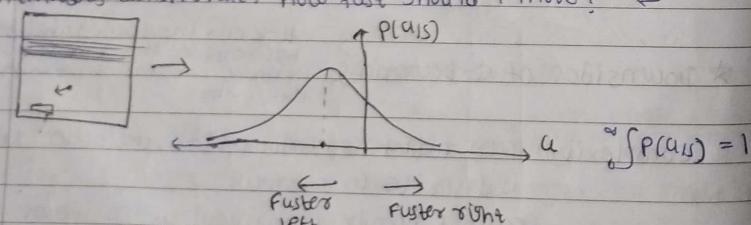
Easy to train than RL

Discrete vs continuous action spaces

discrete action space: which direction should i move!

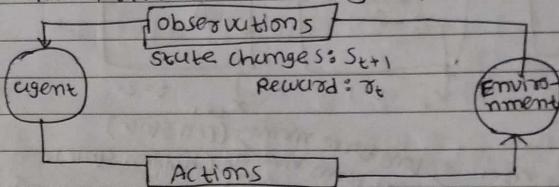


continuous action space: how fast should i move? $\leftarrow 0.7 \text{ m/s}$



⇒ Training policy gradients

Reinforcement learning loop



case-study: self driving cars a_t

Agent: vehicle

state: camera, lidar, etc

Action: steering wheel angle

Reward: distance traveled

Training algorithm:

- 1) initialize the agent
- 2) Run a policy until termination
- 3) Record all state, actions, rewards
- 4) decrease Prob of action that resulted in low reward
- 5) increase Prob of " " " " " high "

$$\text{loss} = -\log p(a_t|s_t) R_t$$

↑
log-likelihood of action ↓
reward

gradient decent update

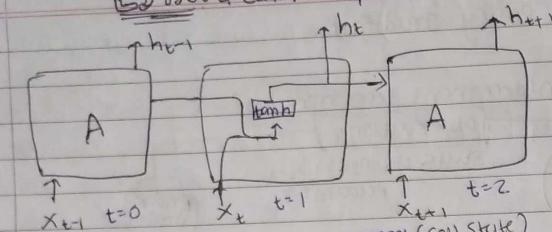
$$w' = w - \eta \text{loss}$$

$$w' = w + \nabla \log p(a_t|s_t) R_t$$

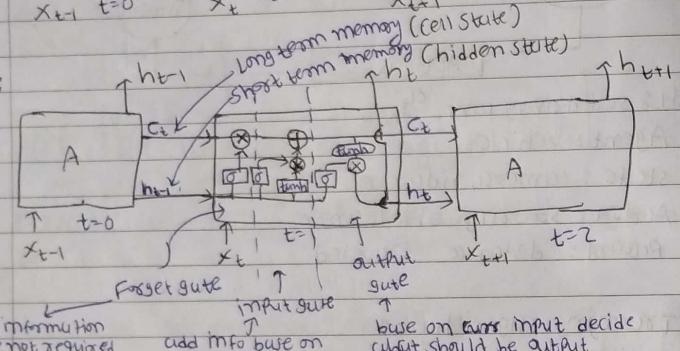


LSTM & GRU

RNN:



LSTM:



Remove information which is not required base of current input

add info base on current input

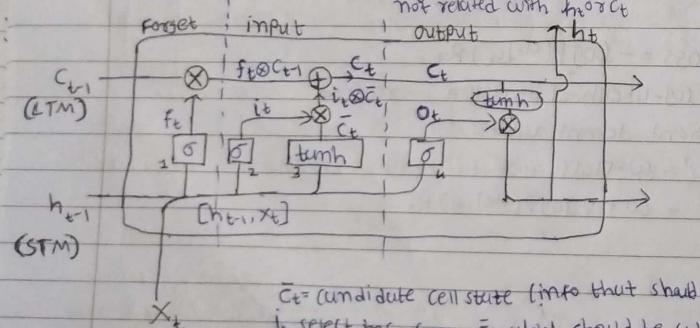
base on curr input decide what should be output

$\rightarrow h_t \text{ & } c_t$ is vector

dimension of $h_t \text{ & } c_t$ is same & also $h_t, c_t, f_t, i_t, \bar{c}_t, o_t$ has same dimension

$\Rightarrow x_t$ is input at time t (any dimension vector)

not related with h_{t-1} or c_t



\bar{c}_t = candidate cell state (info that shall be added in LTM)
select info from \bar{c}_t which should be add to LTM

\Rightarrow [] neural network layers with activation function σ / tanh
[] NO of node for every neural network \mapsto sum
NN has only one HL (neural network layer)

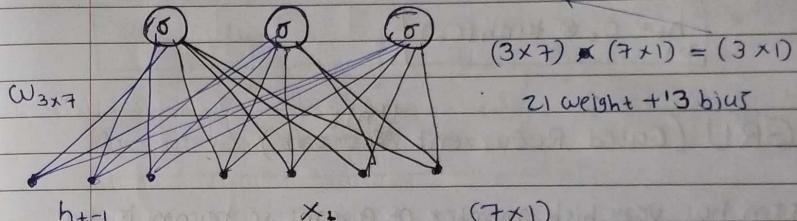
$$\begin{aligned} \otimes &\Rightarrow [u \ 5 \ 6] \otimes [1 \ 2 \ 3] \Rightarrow [u \ 10 \ 18] \\ \oplus &\Rightarrow " \oplus " \Rightarrow [5 \ 7 \ 5] \\ \text{tanh} &\Rightarrow [1 \ 2 \ 3] \Rightarrow [\text{tanh}_1 \ \text{tanh}_2 \ \text{tanh}_3] \end{aligned}$$

\Rightarrow there are three pointwise operation \otimes , \oplus , tanh

\Rightarrow calculation for forget gate: (let assume node = 3 & input dim = 4)

$$\begin{array}{c} f_1 \quad f_2 \quad f_3 \\ | \quad | \quad | \\ b_1 \quad b_2 \quad b_3 \end{array} \rightarrow [f_1 \ f_2 \ f_3]$$

$$(3 \times 1) + (3 \times 1) = (3 \times 1)$$



$$f_t = \sigma \left(W_f [h_{t-1}, x_t] + b_f \right)$$

$$(3 \times 1) + (3 \times 1) \rightarrow (3 \times 1) \sigma(3 \times 1)$$

$$C_{t-1} \otimes f_t \rightarrow C_t$$

Forgot some info

(add info)

\Rightarrow Calculation for input gate:

\bar{c}_t is calculate same as above but use tanh activation fun

$$\bar{c}_t = \text{tanh} \left(W_i [h_{t-1}, x_t] + b_i \right)$$

i_t same w/ above

$$i_t = \sigma \left(W_i [h_{t-1}, x_t] + b_i \right)$$

Pointwise operation of \bar{c}_t & $i_t \Rightarrow \bar{c}_t \otimes i_t$

Now for c_t do \oplus for $f_t \otimes C_{t-1} \oplus i_t \otimes \bar{c}_t$

$$C_t = f_t \otimes C_{t-1} \oplus i_t \otimes \bar{c}_t$$

⇒ calculation for output gate (OP for current timestamp h_t)

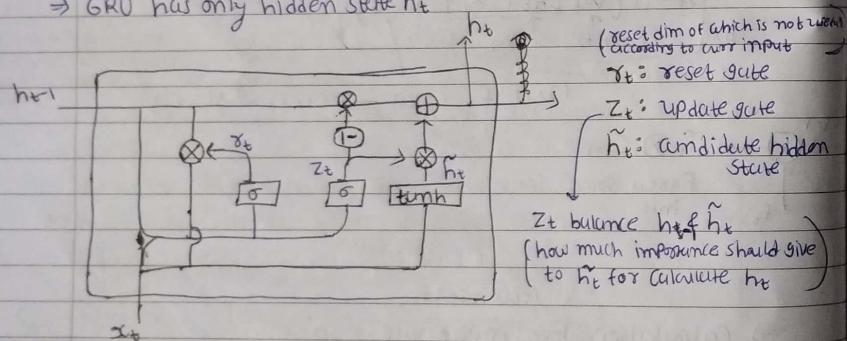
$$O_t = \sigma [W_o [h_{t-1}, x_t] + b_o]$$

$$h_t = O_t \otimes \text{tanh}(C_t)$$

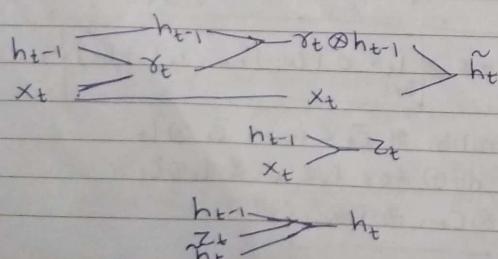
★ GRU (Gated Recurrent Unit)

⇒ LSTM has very high number of parameters so training is slow
so GRU is used for fast training.

⇒ GRU has only hidden state h_t



⇒ $h_{t-1}, h_t, r_t, z_t, \tilde{h}_t$ has same dim



⇒ choice betw GRU & LSTM is depends on dataset and task
try both and find best one

$$r_t = \sigma (W_r [h_{t-1}, x_t] + b_r)$$

$$\tilde{h}_t = \text{tanh} (W_c [h_{t-1} \otimes r_t, x_t] + b_c)$$

$$z_t = \sigma (W_z [h_{t-1}, x_t] + b_z)$$

z_t balance h_{t-1} & \tilde{h}_t

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

LSTM

⇒ 3 gates: input, output, forget \Rightarrow 2 gates: reset, update

⇒ cell state (C_t)

hidden state (h_t)

⇒ more parameters

for input size d & hidden size h
LSTM has $3 \times ((d \times h) + (h \times h) + h)$
parameters: $3 \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

⇒ due to extra gates

computationally intensive

⇒ LSTM give better result
than the GRU

GRU

vs

GRU

⇒ hidden state (h_t)

⇒ less parameters

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$

$h \times ((d \times h) + (h \times h) + h)$



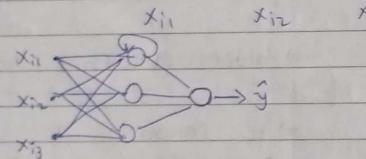
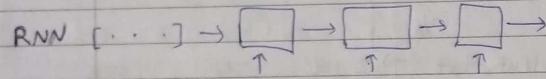
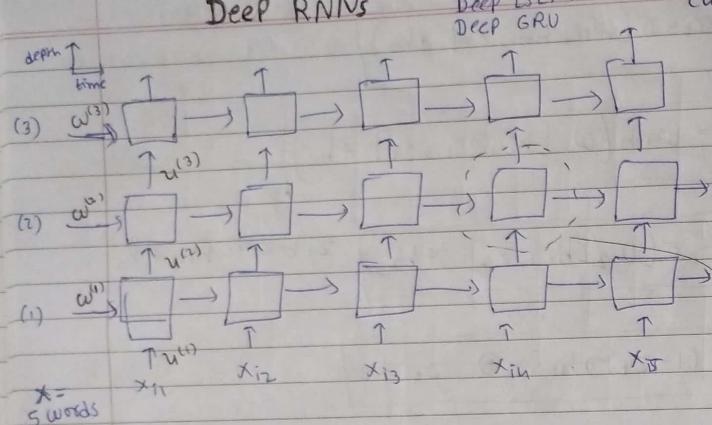
Deep RNNs

DEEP LSTM
DEEP GRU

PAGE NO.:

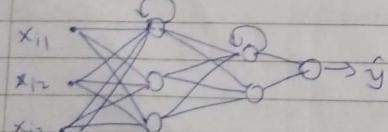
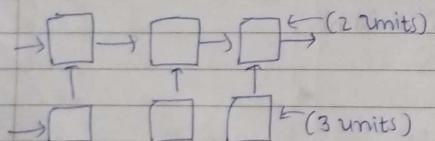


CampusX



Paragraph level
Sentence level
Word level

Deep RNN

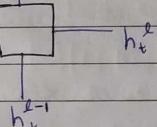


$t=1$

$t=2$

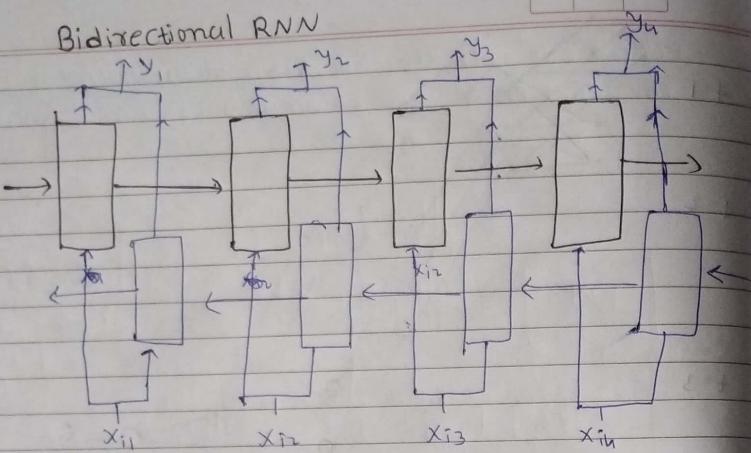
$t=3$

$$h_t^2 = h_3$$



$$h_t^l = \text{tanh}(w^{(2)} h_{t-1}^{(1)} + u^{(2)} h_t^{(l-1)} + b)$$





$$\vec{h}_t = \text{tanh}(\vec{W} \vec{h}_{t-1} + \vec{U} \vec{x}_t + \vec{b})$$

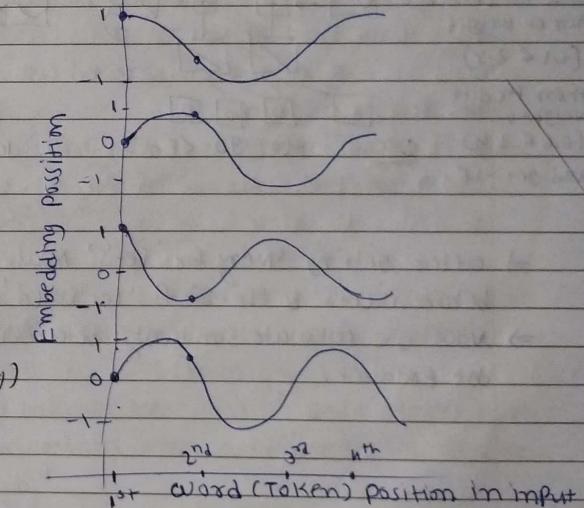
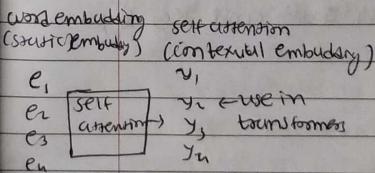
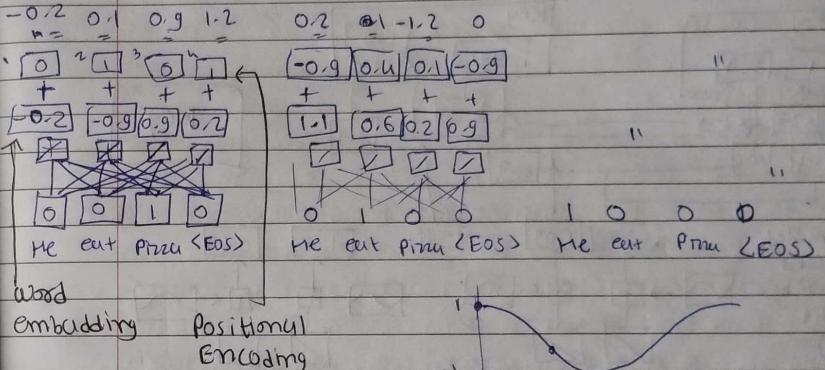
$$\overleftarrow{h}_t = \text{tanh}(\overleftarrow{W} \overleftarrow{h}_{t+1} + \overleftarrow{U} \overleftarrow{x}_t + \overleftarrow{b})$$

$$y_t = \sigma(V[\vec{h}_t, \overleftarrow{h}_t] + b)$$

↑ ↑ ↑
weight concatenated bias

Transformers

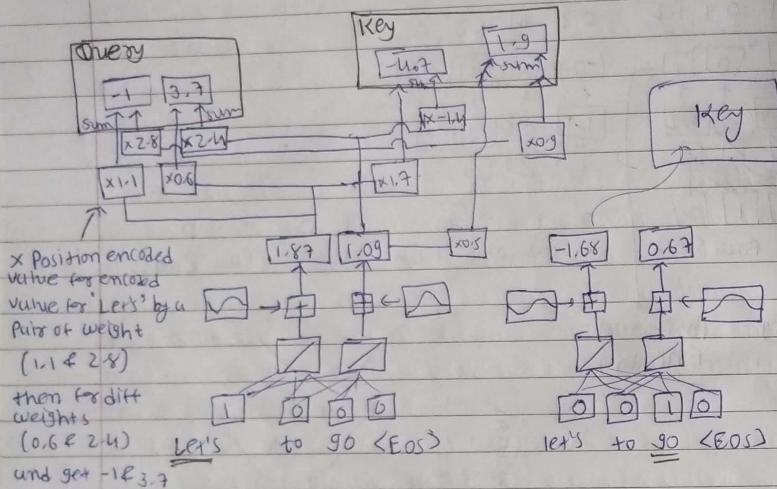
Position Encoding: position encoding allows a Transformer to keep track of word orders.



Self Attention: self attention works by seeing how similar each word is to all of the words in the sentence including itself.

The pizza came out of the oven and it tasted good!

→ self attention calculate similarity of the 'The' full of the words in sentence. sum for 'Pizza' of all the words



- after getting query for 'let's' now we calculate similarity between 'let's' to 'let's' & 'let's' to 'go'
- Now we calculate similarity between Query and Key by dot product.

DeepLearning.AI

Sequence models

Word to Vec

- ★ why Problem with softmax classification

target word
embedding vector

$$P(t|c) = \frac{e^{\Theta_t^T e_c}}{\sum_{j=1}^{vocab_size} e^{\Theta_j^T e_c}}$$

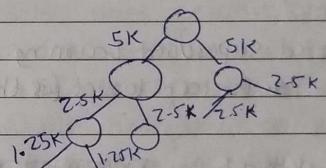
deeplearning.ai andrew ng
20 Nov 23

WEEK 2
w2vec

here problem is: we need to calculate this complex calculation every time

⇒ There are many solution one is Hierarchical softmax classifiers

⇒ This hierarchical softmax classifier instead of trying to categorize something into all 10000 categories on one go image you have one classifier which tell target word is in first 5000 or second 5000 and second classifier tell word is in 2500 or 2500 & third 2500 or 2500 ... until leaf node tell you get down to what word is exactly it is.

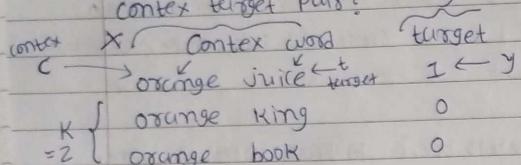


so each step it is binary classifier so we do not need to sum all of 10000 at denominator and complexity become vocab size to log (vocab size)

⇒ it is not perfectly balanced tree in practice tree has top leaf node as common words and deeper leafs are less common words.

Negative Sampling

In this algo we create new supervised learning problem given a pair of word we are going to predict is this a context target pair?



I want a glass of orange juice to go along with my cereal. So for create this dataset we choose word within + window of word orange so for that target is 1.

- For target 0 we choose random word from dictionary and so probability that is not in windows and target for this is 0.
- and its OK if random selected word is from dictionary but we target 0.
- so now we have created supervised Learning problem for given context & word → Predict target for those words (can put together?)

$K = 5-20$ for small dataset

$K = 2-5$ large dataset

for large dataset for

context

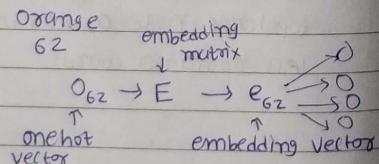
$K = \frac{\text{positive negative to Positive example ratio}}$

Model: softmax

$$P(t|c) = e^{\theta_t^T e_c}$$

$$\sum_j^{10000} e^{\theta_j^T e_c}$$

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



so by creating this now we are not predicting all
For given word orange
target word is juice or not? and so on for whole
dictionary 10,000

So that have 10000 binary regression classifier but instead of all 10,000 of them every iteration we are only going to train K of them (here K=2) + target 1

⇒ so instead have 10,000 way giant softmax which is very expensive to compute we are instead turn it into 10000 binary classification problem each of very cheap to compute and every iteration we are only going to train K+1 of them. which is relatively cheap to do negative example ↑ positive example. on every iteration we are updating 10,000 very softmax classifier

This technique called Negative sampling.

→ How to Select ~~most~~ negative word?

~~most~~ one thing we can do is sample according to empirical frequency of the word in corpus but problem is with that is that you end up with a very high representation of the, of, a, an...

→ one other thing we can do is $\frac{1}{|\text{Vocab size}|}$ sample negative example uniformly random But that also very non representative of distribution

→ so researcher found some heuristic value which is little bit in between the two extremes of sampling from the empirical frequencies. meaning from whatever is the observed distribution in english text to the uniform distribution. sample proportional to their frequencies of the word to the power of 3/4.

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10000} f(w_j)^{3/4}}$$

$w_i = \text{freq of word in corpus}$

Glove (global vector for word representation)

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

$$x_{ij} = \text{number of times word } j \text{ appears in context of } i$$

$c \uparrow$ \uparrow \uparrow
 t t i

x_{ij} is count captures how often do words i & j appear with each other

Glove model:

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

\uparrow \uparrow \uparrow \uparrow
 weighting term θ_i^T e_j b_i
 $= 0 \text{ if } x_{ij} = 0$

Solve for parameters θ_i & e_j with using gradient descent to minimize above eqn

$$f(x_{ij}) = 0 \text{ when } x_{ij} = 0$$

also $f(x_{ij}) \geq 0$ for stop words, this, is, a, am, the... so this weighting function is function which give less weight to more freqⁿ word and more to less freqⁿ word.

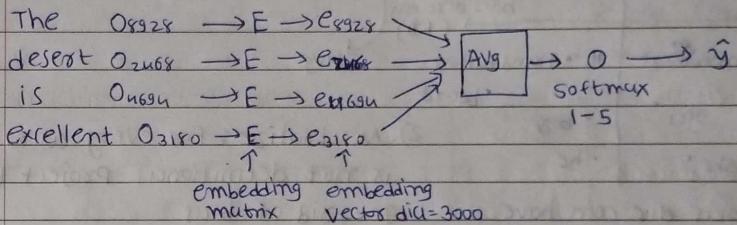
⇒ In this algo role of θ_i and e_j are completely symmetric

Sentiment classification model

21 Nov 23

The desert is excellent

8828 2468 1694 3180



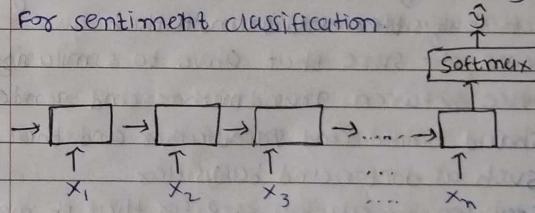
⇒ This algo avg all the word

⇒ problem with algo is ignore order of the word

⇒ Not good sentence like: Completely lacking in good taste

good service and good ambience

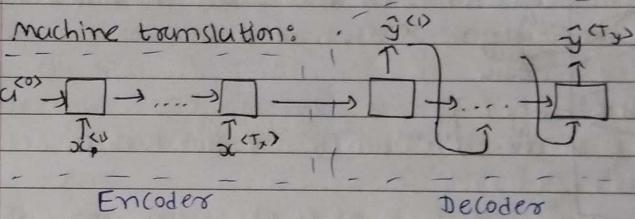
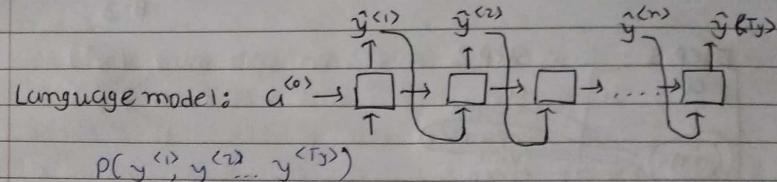
⇒ So we use RNN instead of summing all of word embedding for sentiment classification.



Week 3

Addressing bias in word embeddings

- non-bias direction (2D dimension)
- 1) identify bias direction
-
- ⇒ bias dir can have more than 1D dim by using PCA & SVD to handle it.
- 2). Neutralize: for every word that is not definitional project to get rid of bias
→ project ~~grandmother~~ word like doctor and baby sitter onto axis element to ~~reduce~~ their component on the bias direction
- 3). Equalize pairs.
⇒ This example similarity betⁿ grandmother and baby sitter is smaller than babysitter and grandfather.
⇒ This step we make ~~the~~ sure that above similarity is same. distance between grandmother and gender neutral word should same and grandfather and their gender neutral word such as doctor and babysitter
⇒ so there are few linear algebra step for that it move grandfather and grandmother from the axis in the middle!



conditional language model
 $p(y^{(1)}, \dots, y^{(T)}, | x^{(1)}, \dots, x^{(T)})$

why not use greedy search?

$p(y^{(1)}|x)$
First op word
French sentence

Find Prob of first op word given x
choose most prob word the go for
 $p(y^{(2)}|x)$ choose most prob word
For second word and so on...

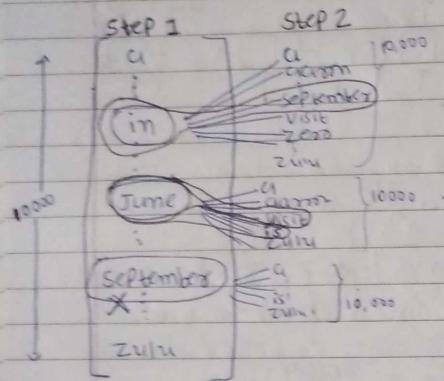
greedy op: June is going to be visiting Africa in september
→ Here 'going' is most common after 'is' in English dictionary
so greedy algo choose above segment for French to eng conversion of sentence X = June visite l'Afrique en septembre

→ but actual accurate op should be
June is going to be visiting Africa in september.

⇒ For greedy algo space complexity for 10,000 vocab and 10 word long prediction take $10,000^{10}$ space so we use approx search algo and do crg max ($p(y^{(1)}, y^{(2)}, \dots, y^{(T)})|x$)

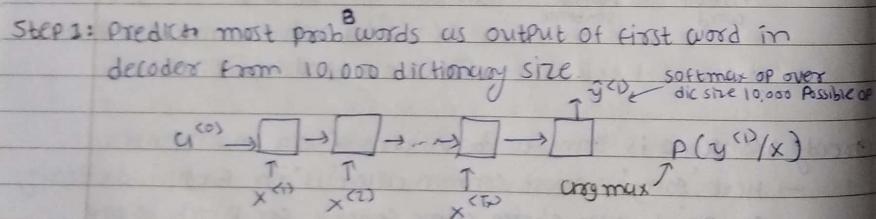
Beam search algorithm

Here Beam width
B=3



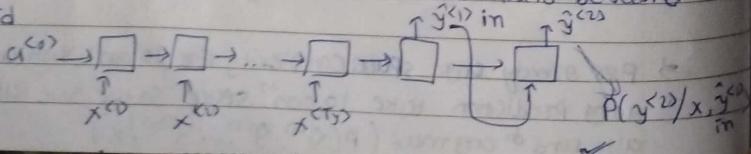
⇒ large B better result, slower
small B worse result, faster

⇒ Unlike BFS & DFS Beam search runs faster but is not guaranteed to find exact maximum for arg max $p(y|x)$



Let in June & September is highest first 3 prob words

Step 2: For each of this 3 choices consider what should be second word



do same for month september

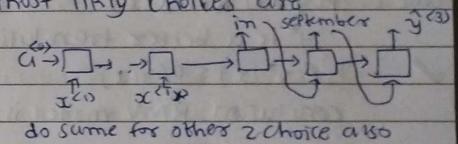
and finally we select most prob first 3 For m. JUNE SEPTEMBER

$$p(y^{(1)}, y^{(2)} | x) = p(y^{(1)} | x) p(y^{(2)} | x, y^{(1)})$$

⇒ do same for june and september and we get total 9 most Prob op out of those select most 3 likely op

⇒ let out of 30,000 choices most likely choices are

in, september
June, is
June, visit



do same for other 2 choice also

⇒ B=1, is greedy search algo

★ Length normalization

$$\text{argmax } \prod_{t=1}^{T_y} p(y^{(t)} | x, y^{(1)}, y^{(2)}, \dots, y^{(t-1)})$$

$$\text{argmax } \sum_{t=1}^{T_y} \log p(y^{(t)} | x, y^{(1)}, y^{(2)}, \dots, y^{(t-1)})$$

⇒ First eqn \prod has long length and all prob is very 0 to 1 so first eqn become very small so we take log and make 2nd eqn

⇒ we can normalize by dividing by T_y = total word m op

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

generally divide by T_y^α where $0 < \alpha \leq 1$

$\alpha = 0.7$
generally



Error Analysis in Beam Search

Error is due to one of those two
RNN or Beam search

$\xrightarrow{\text{better translation}}$ Human: June visits Africa in September. (y^*)

Algo: June visited Africa last September. (\hat{y})

$$\text{case 1: } P(y^*|x) > P(\hat{y}|x)$$

Beam search choose \hat{y} But y^* attains higher $P(y|x)$

✓ conclusion: Beam search is at fault

$$\text{case 2: } P(y^*|x) \leq P(\hat{y}|x)$$

y^* is a better translation than \hat{y} . But RNN predicted

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

✓ conclusion: RNN model is at fault.

\Rightarrow RNN computes $P(y|x)$

so here we compute $P(y^*|x)$ & $P(\hat{y}|x)$ with RNN
and do ~~error~~ above two case

Error analysis Process

Human	Algo	$P(y^* x)$	$P(\hat{y} x)$	At fault?
-	-	-	-	B
-	-	-	-	RNN
-	-	-	-	B
⋮	⋮	⋮	⋮	⋮

\Rightarrow And we can find out what fraction of error is due to
RNN & Beam search

\Rightarrow RNN fault \Rightarrow deeper network
increases delta ...

\Rightarrow Beam search \Rightarrow increases beam width
fault

Bleu Score (Used for translation)

Bleu score compute a score how good that machine
translation.

French: Le chat est sur le tapis

Reference 1: The cat is on the mat.

Ref. 2: There is a cat on the mat

MT OP: the the the the the the the (g)

$$\text{Precision} = \frac{7}{7} (\because \text{every words in MT OP appears on ref 1 or 2})$$

$$\text{Modified precision} = \frac{2}{7} \leftarrow \begin{array}{l} \text{count "the"} \\ \text{count clip "the"} \end{array}$$

bleu score on bigrams:

count_{clip} is that input
appears on ref 1 or ref 2

MT OP: The cat the cat on the mat (g) count count_{clip}

the cat 2 1

cat the 1 0

cat on 1 1

on the 1 1

the mat 1 1
 $\frac{1}{6}$ $\frac{1}{6}$

$$P_1 = P_2 = \dots = P_n = 1$$

when MT output is
summe as one of the ref

$$P_1 = \frac{\sum_{\text{unigrams}} \text{Count clip (unigram)}}{\sum_{\text{unigrams}} \text{Count (unigram)}}$$

$$P_m = \frac{\sum_{n\text{-grams}} \text{Count clip (n-gram)}}{\sum_{n\text{-grams}} \text{Count (n-gram)}}$$



$$P_n = \text{Blew score of } n \text{ grams only}$$

$$\text{combine Bleu score} = \text{BP} \exp\left(\frac{1}{n} \sum_{m=1}^n P_m\right)$$

BP = brevity penalty

\Rightarrow when MTOP is very short translation then precision goes high so BP penalizes translation system that output translation are too short.

$$\text{BP} = \begin{cases} 1, & \text{if MTOP} > \text{reference OP length} \\ \exp(1 - \frac{\text{ref OP len}}{\text{MTOP len}}), & \text{otherwise} \end{cases}$$

Rouge (use for summarization)

\Rightarrow compares a summary to one or more reference summaries

Reference (human): It is cold outside

Generated op: It is not cold outside

$$\text{Rouge-1} = \frac{\text{unigram matches}}{\text{unigram in reference}} = \frac{4}{6} = 0.67$$

$$\text{Precision} = \frac{\text{unigram matches}}{\text{unigram in op}} = \frac{4}{5} = 0.8$$

$$\text{Rouge-1} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}} = 2 \times \frac{0.8 \times 0.67}{0.8 + 0.67} = 0.89$$

For Rouge-2 bigram instead of unigram

Rouge-1 unigram based on LCS (longest common subsequence)

\Rightarrow Rouge clipping

$$\text{humans: It is cold outside} \quad \text{Rouge-1 precision} = \frac{4}{6} = 0.67$$

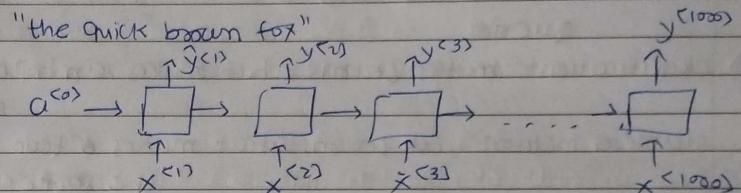
$$\text{OP: It It It It} \quad \text{modified precision} = \frac{\text{clip(unigram matches)}}{\text{unigram in op}} = \frac{1}{6} = 0.17$$

P₁, P₂, P₃, P_n

Speech Recognition - Audio Data

1 march 2024

★ CTC cost for speech recognition
(connectionist temporal classification)



\Rightarrow number of time step is very large for audio input
time step is much bigger than output

\exists For 10 sec of audio comes at 100 Hz so we have 100 sample per second and 10sec audio clip has 1000 timestep
But OP might not have 1000 alphabets it would be much lesser.

\Rightarrow CTC cost function allows RNNs to generate ~~multiple~~ output like this

\downarrow
the 9
ttt - h - eee -- - l - - 999 - - \Rightarrow the 9
 \uparrow blank
space

★ Trigger word detection (trigger words like: Alexa, okay, Google)

\Rightarrow So to detect trigger word we give input audio to train model which has table as below

For training: when trigger words comes in input audio mark output of model 1 otherwise 0 for that time interval and train model



Scanned with OKEN Scanner

Deep Learning. AI

NLP with sequence models

Siamese Networks

: identify similarity between things

⇒ siamese network in NLP

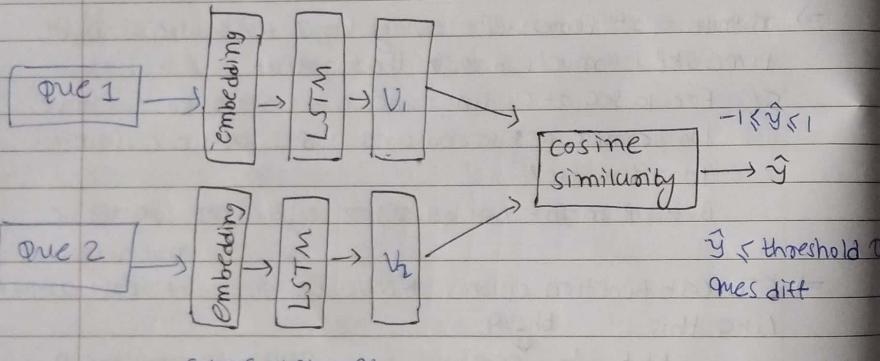
Handwritten checks ~~✓~~ ~~✓~~ ~~✓~~

question duplicate

queries

⇒ classification: network learns what makes an input about it is

siamese network: network learns about makes or two inputs the same



⇒ Not all siamese network use LSTMs.

Loss Function

How old are you? Anchor

What is your age? Positive

Where are you from? Negative

$$\cos(V_1, V_2) = \frac{V_1 \cdot V_2}{\|V_1\| \|V_2\|}$$

$$\cos(A, P) \approx 1$$

$$\cos(A, N) \approx -1$$

$$-(\cos(A, P) + \cos(A, N)) \leq 0$$

PAGE NO.:
WEEK 3
1 march 2014

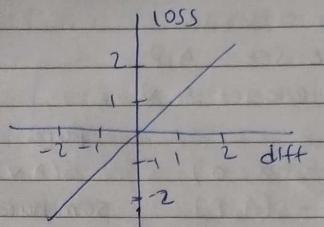
PAGE NO.:
PAGE 3

Simple loss:

$$diff = S(A, N) - S(A, P)$$

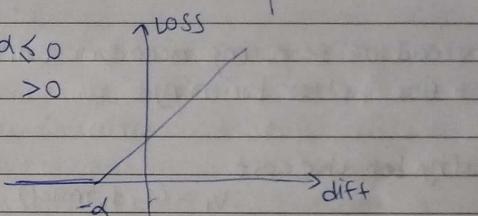
with non-linearity

$$L = \begin{cases} 0 & \text{if } diff \leq 0 \\ diff & \text{if } diff > 0 \end{cases}$$



with a margin

$$L = \begin{cases} 0 & \text{if } diff + td \leq 0 \\ diff + td & \text{if } diff + td > 0 \end{cases}$$



⇒ you need your model confident enough and distinguish between negative and positive values with ease so you need your model to learn until the model achieves a specific value for the difference between similarities so we need to use margin α . It's ensuring that learning goes on until the difference between similarity is equal or close to the margin value.

$$\text{so Triple loss } L = \begin{cases} 0 & \text{if } diff + td < 0 \\ diff + td & \text{otherwise} \end{cases}$$

$$L(A, P, N) = \max(\text{diff} + td, 0)$$

From neural network

You can use any similarity function or distance metric

Triplet selection

Triplet A, P, N
 / duplicate set A, P
 non duplicate set A, N

Random $L = \max(\text{diff} + \alpha, 0)$

$$\text{diff} = S(A, N) - S(A, P)$$

Easy to satisfy Little to learn $S(A, N) < S(A, P)$
 $< S(A, N) + \alpha$

Hard $S(A, N) \approx S(A, P)$

Harder to train more to learn $S(A, N) > S(A, P)$

Easy:

$$S(A, N) < S(A, P)$$

semihard

⇒ we can use off-diagonal information to make some modification to the loss function & can improve model performance

⇒ to do that we used to concepts

mean negative: mean (avg) of all off-diag value in each row

closest negative: off diagonal value closest to (but less than) the value on diagonal in each row

mean negative: mean of off-diagonal values.

$$\text{Longman} = \max(\underbrace{S(A, N) - S(A, P)}_{\text{diff}} + \alpha, 0)$$

$$L_1 = \max(\text{mean_neg} - S(A, P) + \alpha, 0)$$

closest negative: closest off diagonal value

$$L_2 = \max(\text{closest_neg} - S(A, P) + \alpha, 0)$$

$$L_{\text{full}}(A, P, N) = L_1 + L_2$$

$$J = \sum_{i=1}^m \cdot L_{\text{full}}(A^{(i)}, P^{(i)}, N^{(i)})$$

		V ₁			
		V ₁₁	V ₁₂	V ₁₃	V ₁₄
Batch 1	1	1	-1	1	-1
	2	-1	1	-1	1
	3	1	-1	1	-1
	4	-1	1	-1	1

		V ₂			
		V ₂₁	V ₂₂	V ₂₃	V ₂₄
Batch 2	1	1	-1	1	-1
	2	-1	1	-1	1
	3	1	-1	1	-1
	4	-1	1	-1	1

$S(V_1, V_2)$

		V ₁			
		-1	2	3	-4
		-1	0.9	0.3	-0.5
		-2	-0.8	-0.5	0.1
		-3	0.3	0.1	0.7
		-4	-0.5	0.2	0.1

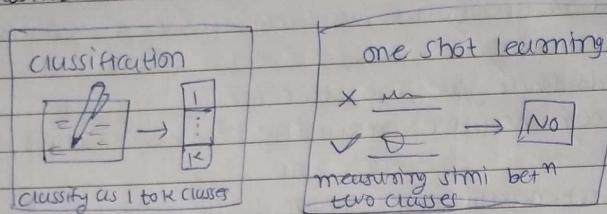
Other than this : similarities for non duplicate

Similarity
bet^n duplicate question

For well train model this value should be greater than off diagonal similarities.

Deep learning. AI

Classification vs one shot Learning



→ problem change which class → measuring similarity bet^n two classes

→ one shot learning no need to retrain when new data with new op comes to model. we need to just check similarity score bet^n signs.

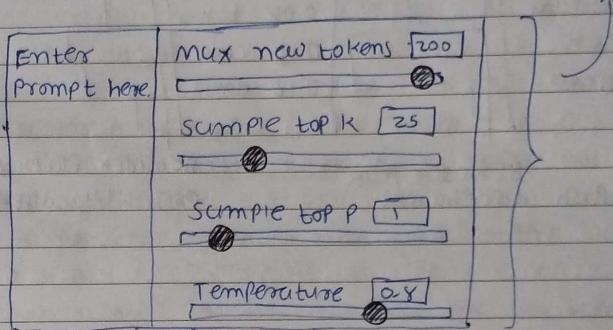
Generative AI with Large Language Models.

PAGE NO. _____

1 march 24

week 1

Generative configuration - inference parameters



greedy vs random sampling

greedy: The wordToken with the highest prob is selected

random(-weighted) sampling: select token using random-weighted strategy across the probabilities of all tokens

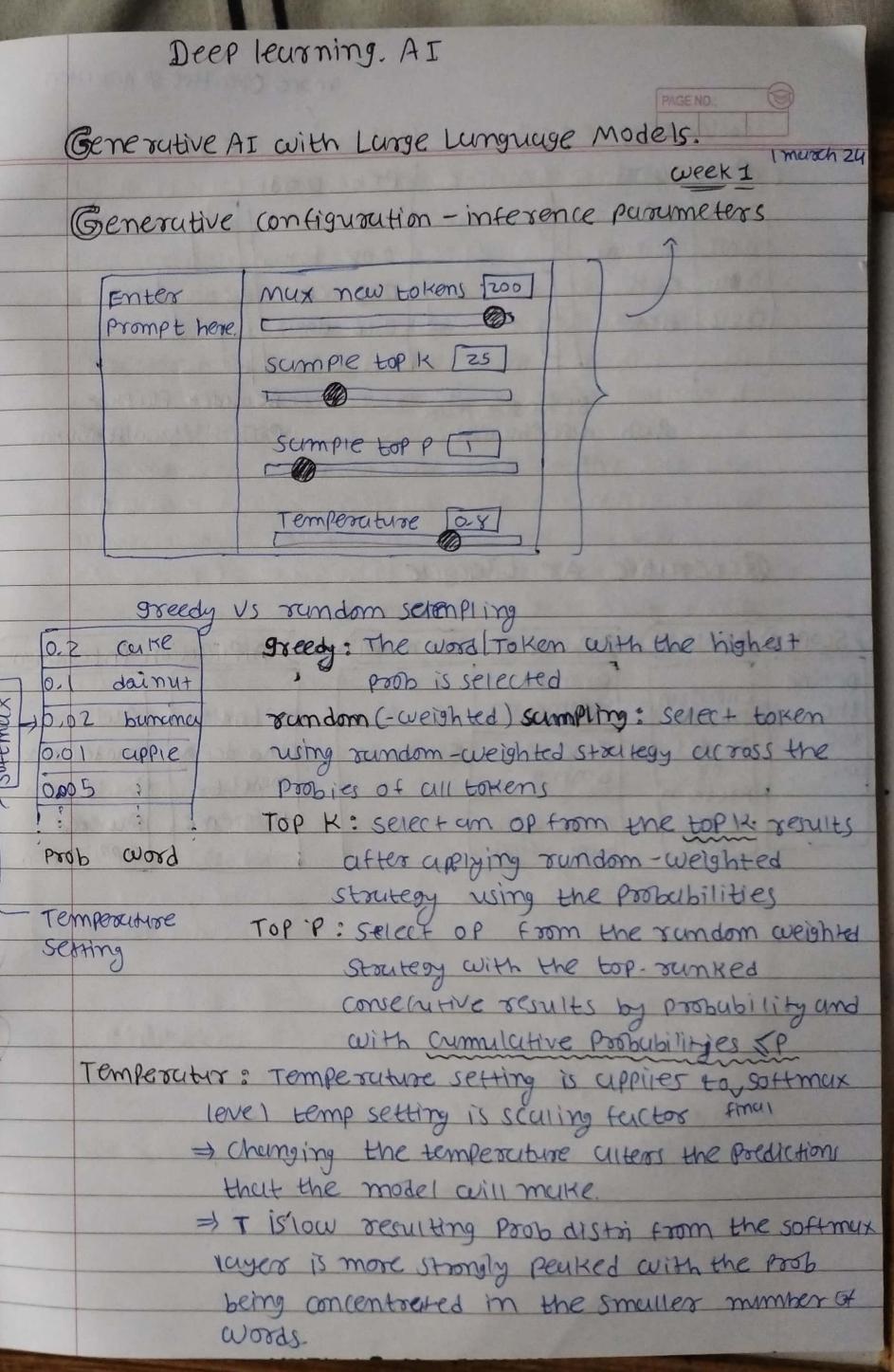
Top K: selection op from the top K results after applying random-weighted strategy using the probabilities

Top P: Select op from the random weighted strategy with the top-ranked conservative results by probability and with Cumulative Probabilities $\leq P$

Temperature: Temperature setting is applied to softmax level temp setting is scaling factor final

→ changing the temperature alters the predictions that the model will make.

→ T is low resulting Prob distn from the softmax layers is more strongly peaked with the prob being concentrated in the smaller number of words.



Cooler Temp (e.g <1)

Prob word

0.001	apple
0.002	banana
0.004	cake
0.012	donut
...	...

Higher Temp (e.g >1)

Prob word

0.01	apple
0.08	banana
0.15	cake
0.12	donut

Stronger peak distribution

Encoder flatter probability distribution

more creative generation

PAGE NO.

Types of Transformers

→ There are three variants of Transformer model

i) Encoder only models

ii) Decoder only models

iii) Encoder-decoder models

i) Encoder only models (Autoencoding models)

This models pre-training using ~~the~~ Masked Language modeling (MLM)

→ This models post-training using ~~the~~ MLM

The original text: teacher teaches the student

→ Here tokens of input sequences are randomly masked and the training objective is to predict the masked token in order to reconstruct original sentence.

→ It's we bidirectional representations of the input sequence that means model ^{has} understand full context of the tokens not ~~the~~ just of the words come before mask language modeling (MLM)

The teacher <MASK> the student

Encoder only LLM

Objective: Reconstruct text ("denoising")

The teacher teaches the student

← bidirectional context

use cases: sentiment analysis

name entity recognition

word classification

Ex BERT
ROBERTA

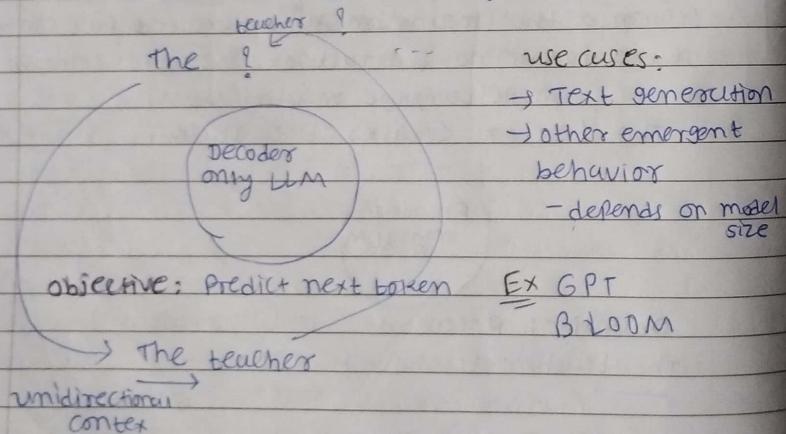


ii) Decoder only (Auto-regressive models)

Decoder only models pretrain using causal language modeling (CLM)

- Here the training objective is to predict next token based on the previous sequence of tokens.
- Predicting the next token is sometimes called full language modeling by researchers.
- This model masks the input sequence and can only see the input tokens leading up to the token in question. Model has no knowledge of end of the sentence.
- The model iterates over the input sequence one by one to predict the following token.
- Decoder only model has only unidirectional context.

Causal Language Modelling (CLM)

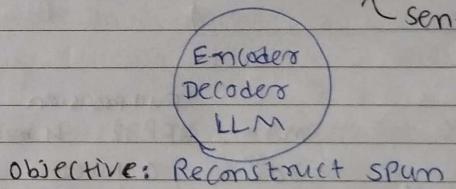


iii) Encoder Decoder Models (sequence-to-sequence models)

- Pretraining details ~~are~~ are varies from model to model.
- Popular seq-to-seq model T5 pretrains ~~with~~ the encoder using SPM corruption which mask random sequences of input tokens and those random sequences are replaced with the unique sentinel token.
- Sentinel tokens are special tokens added to vocab but do not correspond to any actual word from the vocab.
- Decoder is then tasked with reconstructing the masked token sequences auto-regressively.
- The output ~~is~~ is the sentinel token followed by predicted tokens.

SPM Corruption

The teacher <mask> <mask> student
The teacher <x> student
Sentinel token



use cases: Translation
Summarization Text
Q&A answering

Ex:
T5
BART (nabert)

Computational challenges of training LLMs

1 Parameter = 4 bytes (32-bit float)
(real number)

1B Parameters = 4×10^9 bytes = 4GB

4GB @ 32-bit
Full precision

Additional GPU RAM needed for train 1B Parameters

- Adam optimizer (2 states) + 8 bytes
- Gradients + 4 bytes
- Activation and temp memory + 8 bytes

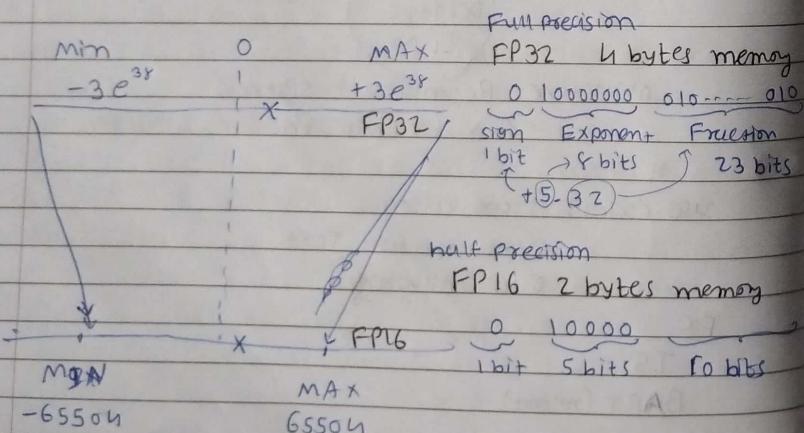
= 4 bytes per parameter

+ 20 extra bytes per parameter

So GPU needed is not 4GB

but 24 GB @ 32-bit
Full precision

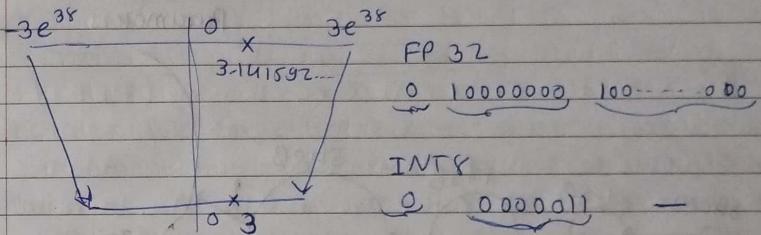
★ Quantization



BFLOAT 16 / BF16 (not suitable for integer calculation)

1 10000000 1001001
1 bit 8 bits 7 bit
Sign Exponent Fraction
 $\approx 3e^{38}$ $\approx +3e^{38}$

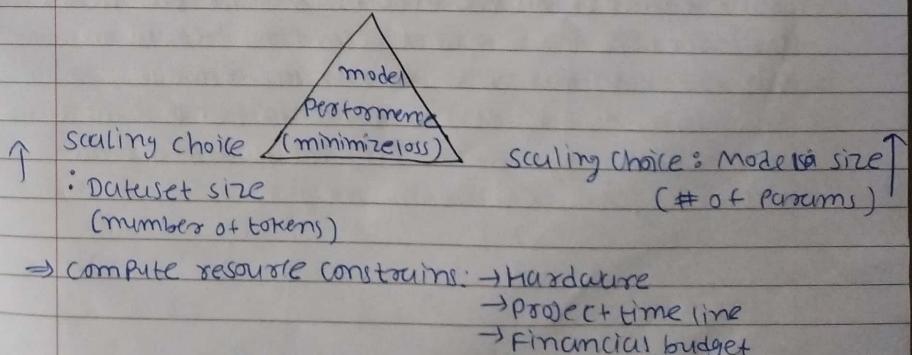
⇒ Let convert pi from FP32 to 8 bit space

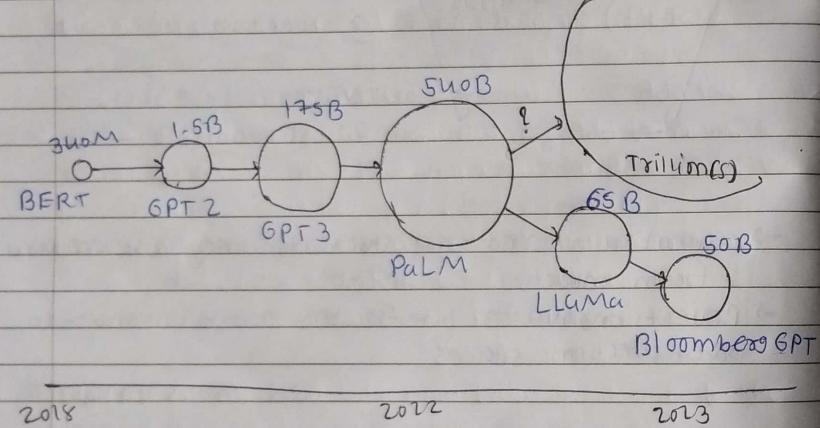
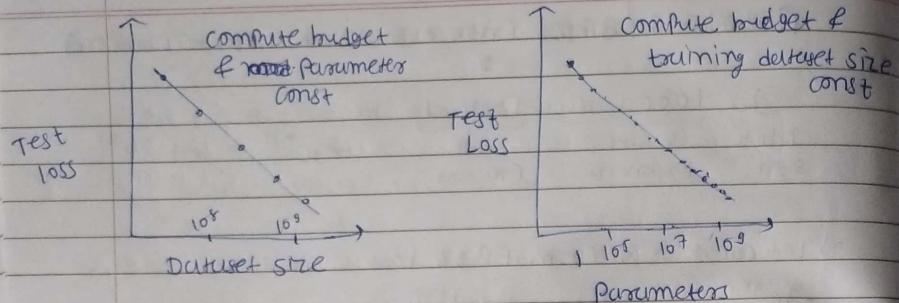


⇒ Quantization reduce required memory to store and train models

⇒ Projects original 32 bit floating point number to lower precision spaces

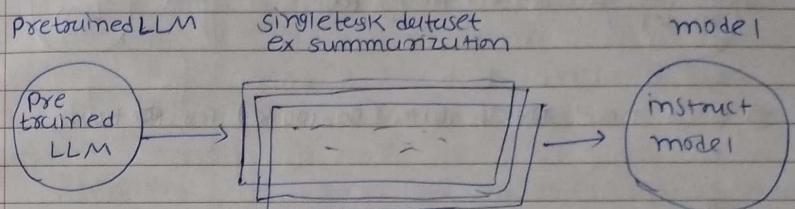
Scaling choices for pre-training Goal: maximize model performance





Fine-tuning LLMs with instruct

★ Fine-tuning on single task



⇒ Catastrophic forgetting happens here. because of full fine tuning process modifies the ~~the~~ weight's of original LLM while this leads to great performance on single fine tuning task but it can degrade performance on other task.

Ex before fine tuning model give best performance on name entity recognition but after pretraining it will give bad result.

⇒ How to avoid catastrophic forgetting.

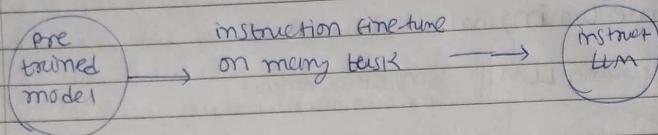
→ If you need only ~~the~~ single task model and do not need multi task then no need to do anything.

→ Fine tune on multi task at the same time

→ consider Parameter Efficient Fine-tuning (PEFT)

↳ PEFT is set of techniques that preserves the weights of original LLM & trains only small number of task specific adapter layers and params.

* multi-task, instruction fine-tuning



Parameter efficient fine-tuning (PEFT)

	selective	Reparameterization	Additive
initial LLM	select subset of parameters to fine-tune	repurpose existing model weight using a low-rank representation	Add trainable layers or parameters to model

Two main Approaches:
 i) Adapters
 ii) soft Prompts

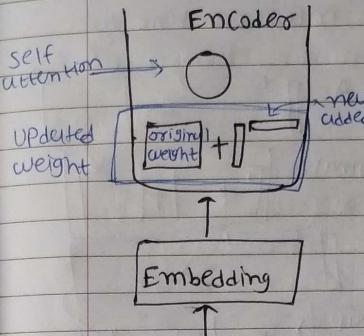
LORA

= **Adapters**: methods add new trainable layers to the architecture of the model typically inside the encoder or decoder component after the attention or feed-forward layers.

= **soft Prompts**: on the other hand keep the model architecture fixed and frozen and focus of manipulating the input to achieve the better performance.

This can be done by adding trainable ~~parameters~~ prompts to the prompt embeddings or keeping input fixed and retraining the embedding weight.

* PEFT techniques 1: LORA (Low Rank Adaptation)



1. Freeze most of original LLM weights
2. inject 2 rank decomposition matrices
3. train the weights of the smaller matrices

Steps to update model for inference:
 1. matrix multiply the low rank matrices

$$B \boxed{\quad} * \boxed{A} = \boxed{B \times A}$$

2. Add to original weight

$$\boxed{\text{original weight}} + \boxed{B \times A}$$

⇒ the model has same number of parameters as the original so there is little to no impact on inference latency.

⇒ Researchers have found that applying LORA to just self attention layers of the model is often enough to fine-tune a task to achieve performance gains.

⇒ However in principle you can also use LORA on other components like a feed forward layer, but since most of the parameters of the LLM are in the attention layers you get a biggest savings in trainable parameters by applying LORA to these weight matrices.

⇒ base on paper attention is all you need

Transformer weights $d \times K = 512 \times 64 = 32768$ trainable parameters in LORA with $\text{rank } r = 8$:

$$A = r \times K = 8 \times 64 = 512 \text{ parameters}$$

$$B = d \times r = 512 \times 8 = 4096 \text{ parameters}$$

86% reduction in parameters to train

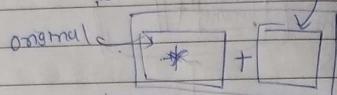
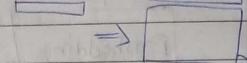
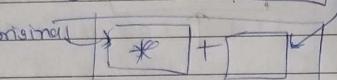
$$32768 \longrightarrow 4096$$



Prompt tuning is not prompt engineering

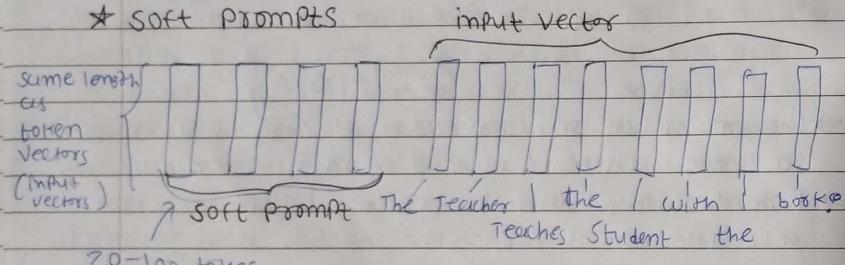
↑
not change weights but giving
Prompt in such a way so
Performance will increase

⇒ we can fine-tune to different set for each task
and then switch them out at inference time by
updating the weights.

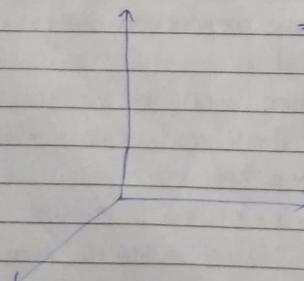
1. Train different tasks Task A:  
2. Update weights before inference Task B:  

PEFT Techniques 2:

★ Soft Prompts

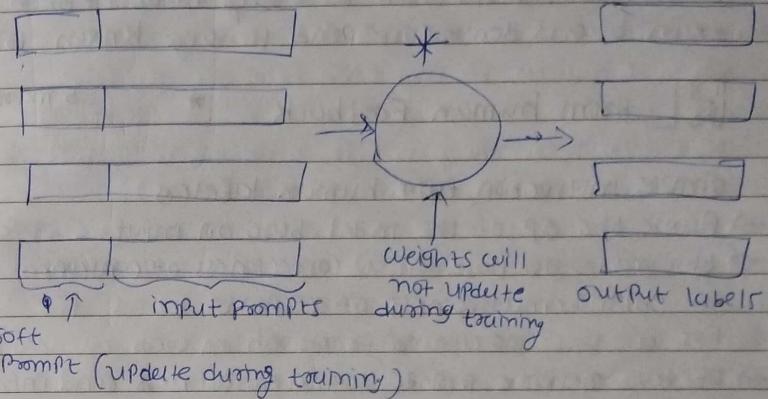


⇒ with prompt tuning you add additional trainable tokens to your prompt and leave it up to the supervised learning process to determine their optimal values.



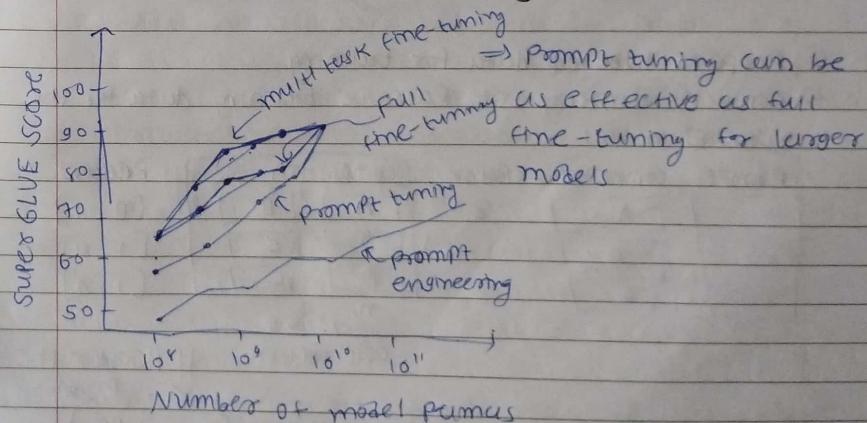
⇒ Soft Prompt can take any value in embedding space to maximize performance

weights of model frozen and soft prompt trained



⇒ Similar to LoRA can train different soft prompt for different tasks and switch out soft prompt at inference time to change task!

Performance of Prompt tuning



Interpretability of soft prompts

- because of prompt can take any value in token space
- train tokens don't correspond to any known token.



from human feedback

★

5 march 2h

Simple instruction for human labelers

- Rank the op of the model base on input.
- this rank is base on a) correctness of answer
 - b) informativeness of response
- For (a) you are allowed to search on web
- If two response provide same correctness and informative ness and there is no clear winner you may rank them as a same but please only use this sparingly
- If ans for given response is nonsensible, irrelevant and highly confusing or does not clearly respond to given prompt then label it with 'F' (for fail) rather than its rank. so poor quality ans can be removed.

Prepare labeled data for training

- convert ranking into pairwise training data for the reward model

Prompt	Completion	Rank	Completion	Reward
A	A	2	A	[0, 1]
B	B	1	A	[1, 0]
C	C	3	B	[1, 0]

$$\binom{n}{C_2} \binom{3}{C_2}$$

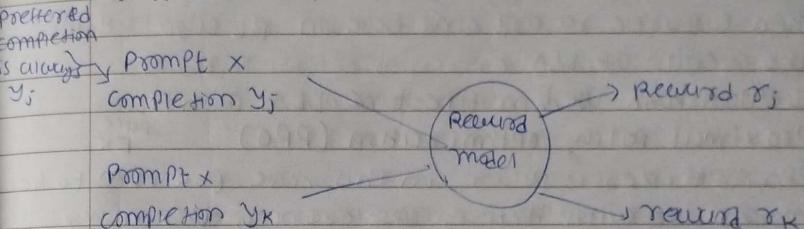
combinations

↓ Reorder the
prompt

completion		reward
B	A	[1, 0]
B	C	[1, 0]
A	B	[1, 0]

Train reward model

train model to predict preferred completion from $\{y_j, y_k\}$ for prompt x



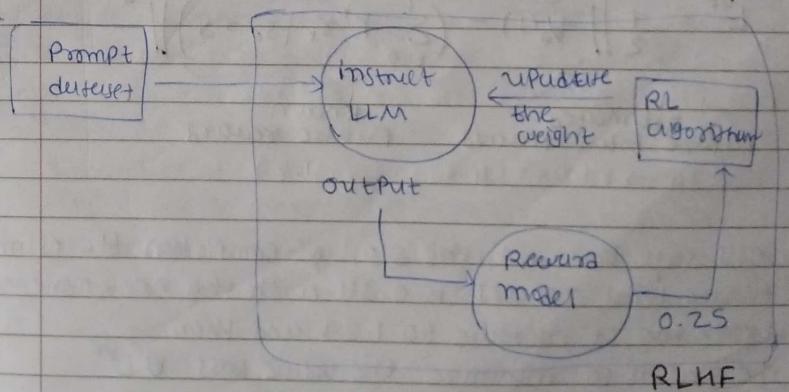
$$\text{loss} = \log(\sigma(r_j - r_k))$$

Use the reward model

- Use reward model as a binary classifier to provide reward value for each prompt completion pair

Prompt completion	Reward model	Logits	Probabilities
		+ve class 3.1753	0.99
		-ve class -2.6943	0.003
		+ve class -0.53215	0.3377
		-ve class 0.8314	0.664

Use reward model to fine-tune LLM with RL



- we are repeating this process until we get threshold of user's human need or till number of epochs.
- In each epoch's input, data feed to LLM and then calculate reward and base on RL algo update the weight of LLM.

RL algorithm used in fine-tune LLM with RL

Proximal Policy Optimization (PPO) ^{Poof}_{ER}

- ⇒ PPO optimizes a policy, in this case the LLM to be more aligned with human preferences over many iteration. PPO make updates to LLM ~~the~~

- ⇒ The updates are small and within a bounded region, resulting in an updated LLM that is close to the previous version.

- ⇒ Keeping the changes within this small region result in a more stable learning.

- ⇒ PPO ensures to keep the model updates within a certain small region called the trust region. This is where the proximal aspect of PPO comes into play. Ideally this series of the small update move the model towards higher rewards.

$$L^{VF} = \frac{1}{2} \| V_b(s) - \left(\sum_{t=0}^T r_t \gamma_t | s_0 = s \right) \|_2^2$$

Estimated future total reward Known Future reward
 $0.34 \rightarrow 1.23 \rightarrow 1.4 \rightarrow 1.87$

- ⇒ Let's say that at this step of completion, the estimated future total reward is 0.34 with the next generated token ~~the~~. It increases to 1.23 and then ...

- ⇒ The goal is to minimize the value loss L^{VF}

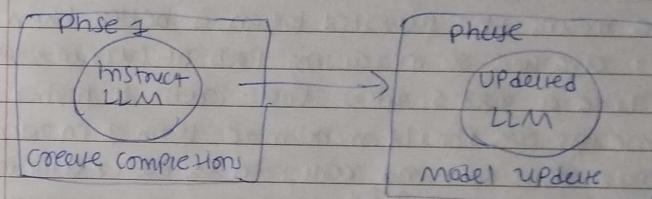
- ⇒ The value function then used in advantage estimation in Phase 2.

- ⇒ This is similar to when you start writing a passage and you have a rough idea of its final form even before you write it.

PPO Phase 2: Model update

In Phase 2 you make small updates to a model and evaluate the impact of those updates on your alignment goal for the model.

- ⇒ The model weight updates are guided by prompt+completion losses and rewards.



PPO phase 2 : calculate Policy loss

$$L_{POLICY} = \min \left(\frac{\pi_{old}(a|s)}{\pi_{new}(a|s)} \cdot \hat{A}_t, \text{clip} \left(\frac{\pi_{old}(a|s)}{\pi_{new}(a|s)}, 1-\epsilon, 1+\epsilon \right) \cdot \hat{A}_t \right)$$

trust region

Prob of the next token with initial LLM (which is frozen)

Prob of the next token with updated LLM (which we can change for a better reward)

there is other technique like PPO is Q-Learning

At: estimated advantage terms

estimates how better or worse the current action is compare to all possible actions at state.

⇒ Notice that second expression defines a region where two policies are near each other.

Second term is

Guarantees: Keeping the policy in the trust region

why here we have choose min of 1st & 2nd term:

- At +ve means the suggested token is better than the current average so increasing prob of current token seems like a good strategy that leads to higher reward.
- This means we should maximize term 1 in order to get max reward. and result in a better aligned LLM
- directly maximizing L^{POLICY} would lead into problems because our calculations are reliable under the assumption that our advantage estimation are valid
- ★ → This advantage estimation are valid only when old and new policy are close to each other

PPO phase 2: calculate entropy loss

$$L^{\text{ENT}} = \text{entropy}(\pi_{\theta}(\cdot | s_t))$$

⇒ while the Policy loss moves the model toward the alignment goal, entropy allows the model to maintain creativity.

→ If you kept entropy low you might end up always completing the prompt in the same way.

→ higher entropy guide the LLM towards more creativity. This is similar to temperature setting of LLM.

→ the diff that ~~then~~ temperature influences model creativity at the inference time while the entropy is during training.

$$L^{\text{PPO}} = \underbrace{L^{\text{POLICY}}}_{\text{Policy loss}} + \underbrace{c_1 L^{\text{VF}}}_{\text{Value loss}} + \underbrace{c_2 L^{\text{ENT}}}_{\text{Entropy loss}}$$

c_1 & c_2 : hyperparameters

→ ~~every~~ The PPO objective updates the model weights through back propagation over several steps.

→ Once the model weights are updated PPO starts a new cycle. For the next iteration LLM is replaced with updated LLM and new PPO cycle starts.

→ after many iterations you arrive at the human-aligned LLM

★ Reward Hacking

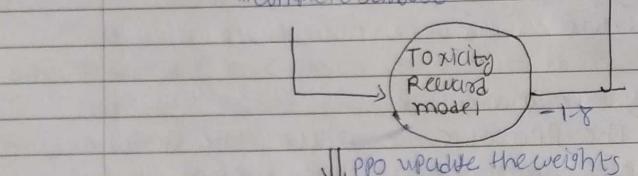
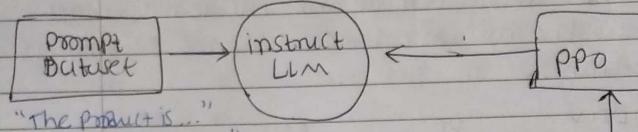
An interesting problem that can emerge in RL is known as reward hacking where the agent learns to cheat the system by favoring actions that maximize the reward received even if those actions don't align well with the original objective.

Eg: suppose we are using RLHF to detoxify and instruct model. You have reward model that can tell if output sentiment analysis and classify model completion is toxic or non toxic.

You select prompt from toxicity data and pass it to the instruct LLM which generates the completion and reward model



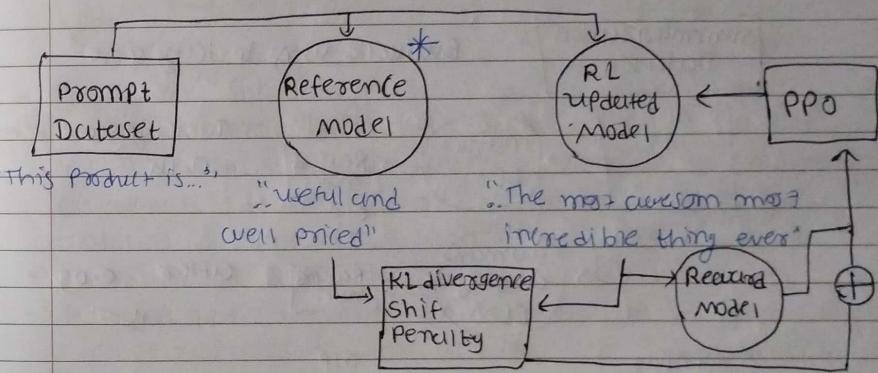
gives a bad reward and PPO update the model weight to create less toxic responses
 → as the policy tries to optimize the reward it can diverge too much from the initial language model



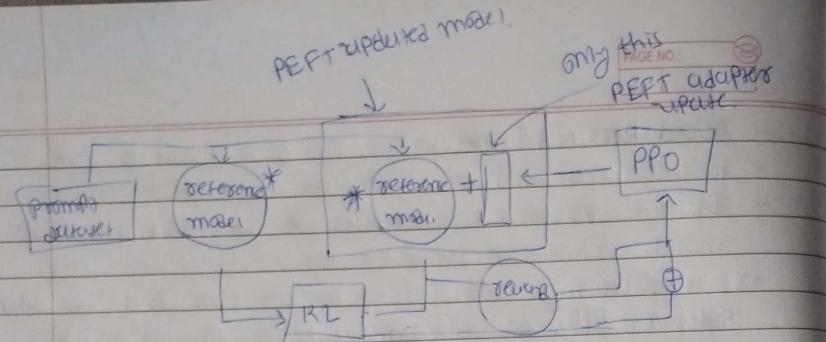
→ The model also can generate nonsensical, grammatically incorrect text that just happens to maximize the reward in a similar way.

Avoid Reward hacking

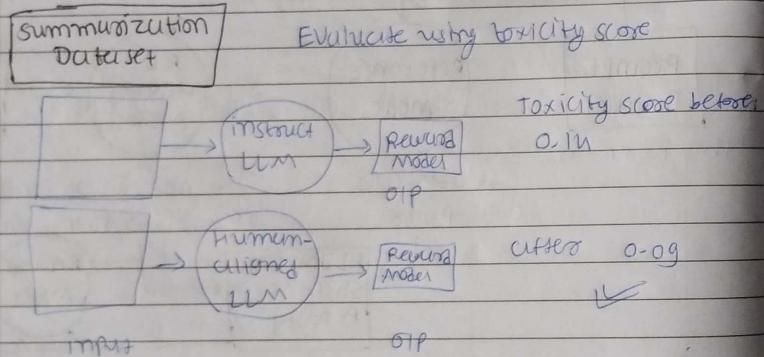
- ⇒ to avoid reward hacking we can use the initial instruct LLM as performance reference lets call it reference model. weight of it's are frozen.
- ⇒ During training each prompt is pass through both model reference LLM and intermediate LLM updated model at this point you can compare the completion and calculate KL divergence
- ⇒ KL divergence are statistical measure of how diff two prob distri are.



- ⇒ KL divergence is calculated for each generated token across the whole vocabulary of the LLM.
- ⇒ once KL divergence betw two model is calculated add said term to the reward calculation
- ⇒ This penalize RL updated model if it shift too far from the reference LLM.
- ⇒ By the way you can benefit from combining our relationship with PEFT in this case only update the weight of PEFT adapters



Evaluate the human-aligned LLM



Constitutional AI

Constitutional AI is a method for training model using a set of rules and principles that govern the model's behavior.

Principles like → choose helpful, honest and harmless responses
→ choose ethical and moral ...

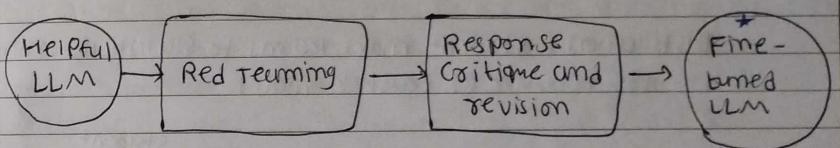
⇒ training model will done in two phases.

i) supervised learning stage

ii) RL stage

i) Supervised Learning Stage

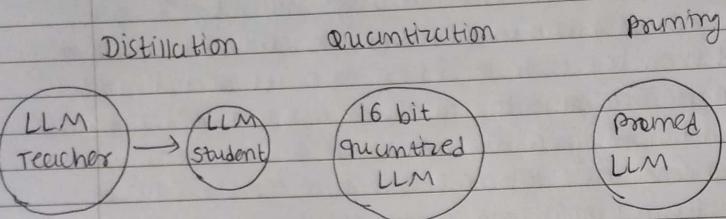
- Prompt the model in a ways that try to get it to generate harmful responses. This is called red teaming and then tell to model to critique its own harmful responses according to constitutional principle and revise them to comply with those rules.
- Once done you finetune the model using the pairs of red team prompt and revised constitutional responses.



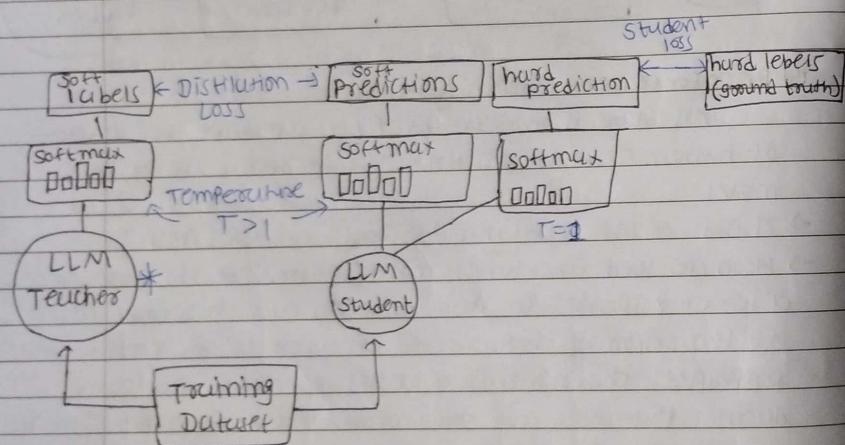
ii) Reinforcement Learning Stage

- This stage is similar to RLHF, except that instead of human feedback we know use feedback generated by model.
- This is called RL from AI Feedback (RLAIF).
- Here we used fine-tuned model * from (i) to generate a set of responses to prompt. Then ask the model which of the responses is preferred according to the constitutional principles. The result is a model generated preference dataset that you can use to train your reward model.
- Using this reward model can now fine-tune model further using RL algo like PPO.

Model optimizations for deployment



Distillation: student model learns to statistically mimic the behavior of teacher model



→ The teacher model is already fine tuned on training data so prob distn likely closely matches the ground truth data and won't much variation in tokens \Rightarrow that's why distillation require Temperature parameter to the softmax function.

→ Combining distillation & student losses are used to update the weights of student model through a back prop.

- In practice distillation not as effective for generative decoder models. It typically more effective for encoder only model like ~~GPT~~ bert.
- In distillation we train new smaller model we are not reducing model size of initial LLM in any way

Pruning: element the weights ~~are~~ that are not contributing much to overall model performance these are the weights with value very close to zero, or equal to zero.

→ however there may not much impact on the size and performance if only a small percentage of the model weights are close to zero.

LLM Powered applications

8-3-24

* Programming-aided language (PAL) models

PAL Example

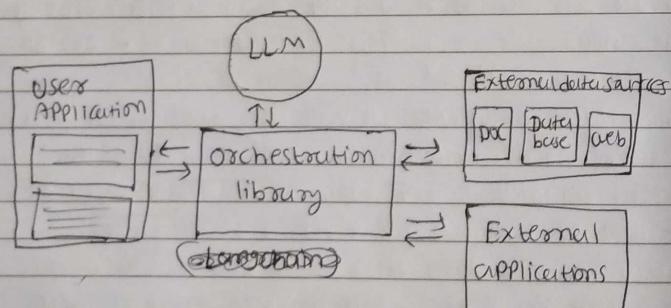
Ques: Roger has 5 balls. He buys 2 more cans of tennis ball. Each can has 3 balls. How many tennis balls does he have now?

Ans: # Roger started with 5 ball
balls = 5

2 cans of balls each is bought - balls = 2 × 3

Tennis balls. cans is
answer = balls + bought - ball

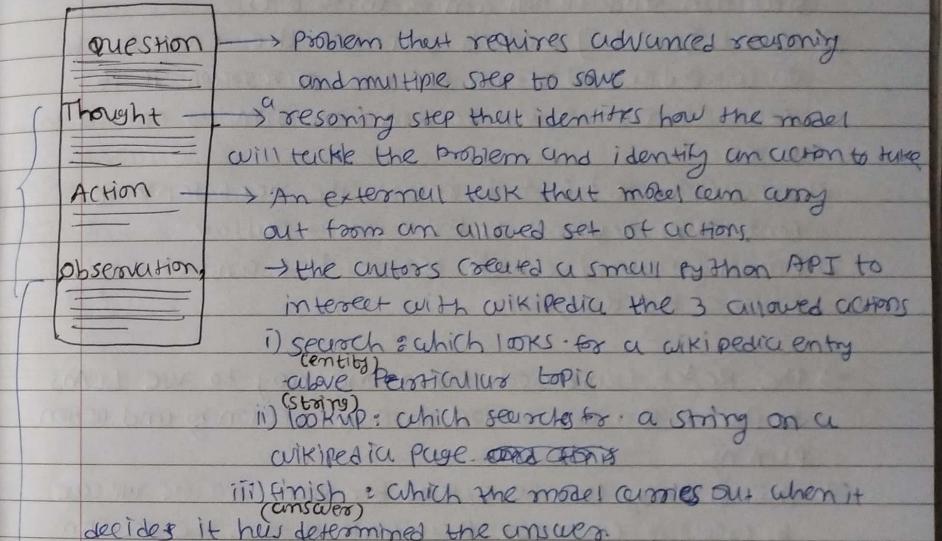
- ⇒ Step of each reasoning step starts with # so line can be skip as a comment by python interpreter
- ⇒ we can give this example as one-shot.



- ⇒ Orchestrator manages the flow of info and initiation of calls to external data sources or applications.
- ⇒ It can also decide what actions to take based on the info contained in the output of LLM

* REACT : Combining reasoning and action.

REACT is a prompting strategy that combines chain of thought reasoning with action planning



→ Observation : the result of carrying out the action.

→ This is where the new information provided by the external search is brought into the context of the prompt for the model to interpret the prompt.

→ The prompt then repeats the cycle { as many time as it necessary to obtain the final answer.

→ Give a question answering task with inter-leaving thought Action observation steps.

→ Thought can be reason about the current situation and action can be three types:

1) Search [entity] : which searches the exact entity on Wikipedia and return the first paragraph if it exist. If not it will return some similar entities to search.

2) Lookup [keyword] : which returns the next sentence containing keyword in the current passage.

3) Finish [answer] : which returns the answer and finishes the task.

→ The React frame work shows one way to use LLMs to power an application through reasoning and action planning.

→ This strategy can be extend for your specific use case by creating examples that work through the decision and craft action that take place in your application.

Ex: Langchain

Langchain framework provides you with modular pieces that contain component to work with LLMs.

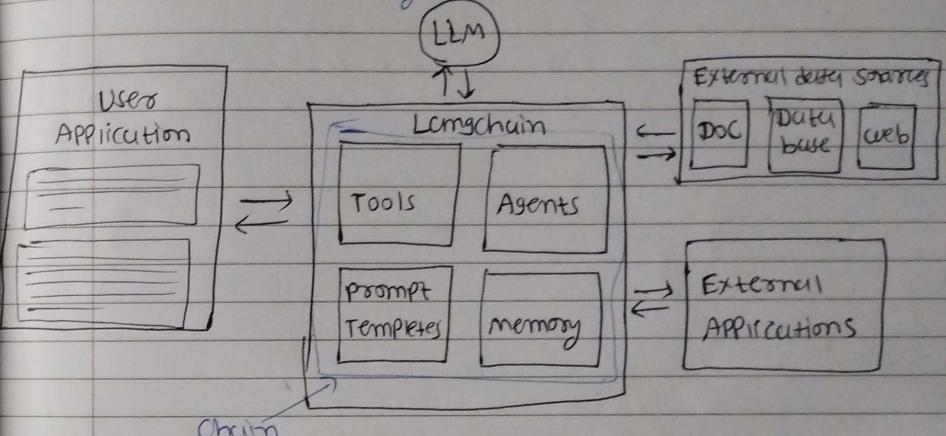
→ These component include Prompt templates for many different use cases that you can use to format both input examples and model completions and memory that you can use to store interaction with LLM.

The framework also include pre-built tools that enable you to carry out a wide variety of tasks, including calls to external datasets and various APIs.

→ Langchain has some pre-defined chains that have been optimized for diff use cases.

→ It also includes agent which interpret the input from the user and determine which tools to use to complete the task.

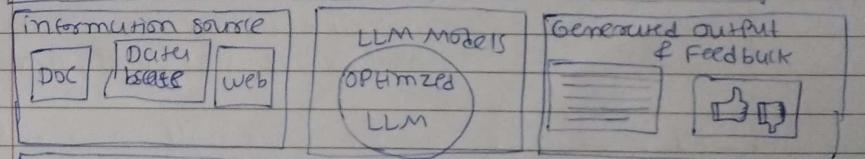
→ Langchain currently includes agent for both PAL and REACT among others.



LLM application architectures

users ex consumers
systems
Application interface ex website, mobile APP, API etc.

LLM Tools & Frameworks Ex: Langchain, Model hubs



Application components Ex training | Fine tuning, serving,