

Why -

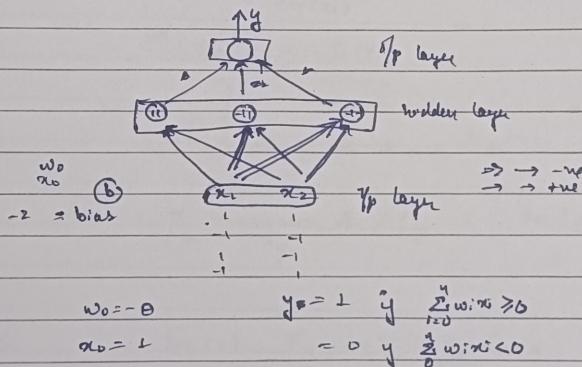
- Threshold value
- weights
- response is binary
- linear - non-linear

Perception

- Perceptron learning algorithm \rightarrow finding optimal weights
- \rightarrow still has problems in implementing XOR (non-linear)

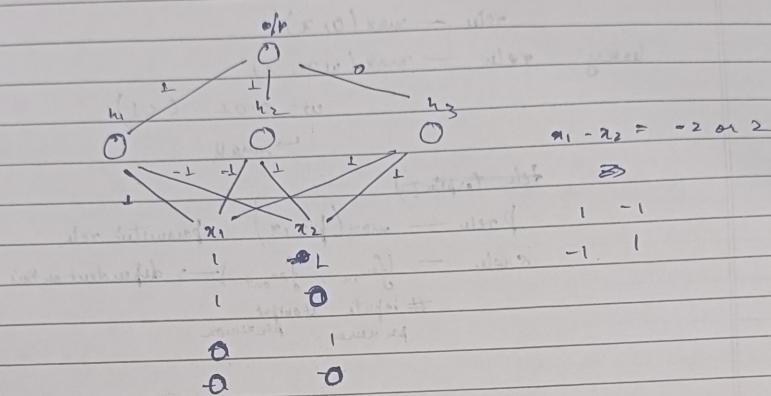
reticular theory

Multi-layer perceptron Architecture (MLP)



how to determine num layers and num neurons





$$x_1 - x_2 = -2 \text{ or } 2$$

\Rightarrow

$$1 - 1$$

$$-1 - 1$$

$$b = -2$$

$$h_1 : x_1 - x_2 + b = 0$$

$$h_2 : x_2 - x_1 + b = 0$$

$$h_3 : x_1 + x_2 + b = 0 \quad (?)$$

$$b = -L$$

$$x_1 + x_2 + b - x_3 + b = 0$$

$$x_1 + x_2 = L \quad x_1 - x_2 = -L \text{ or } L$$

$$x_1 + x_2 \neq 2 \text{ or } 0$$

$$h_1 : x_1 - x_2 - L = 0$$

$$h_2 : x_2 - x_1 - L = 0$$

$$h_3 : x_2 + x_1 - L = 0$$

$$y : \cancel{x_1 + x_2} + h_1 + h_2 + h_3 - 2 = 0$$

Relu $\rightarrow R \rightarrow [0, \infty)$

$$\begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

no upper bound

Not differentiable at 0

Other loss functions may take more time to reach goal. Make why we generally associate 2 loss function with y problem.

Page No.:
Date: / /

other class $\# \text{ class} = \# \text{ # layers neurons}$ (3 neurons) \rightarrow 3 neurons

3 neurons

$$\text{relu} = \max(0, z)$$

$$\text{leaky relu} = \max(wz, z)$$

$$w = 0.01 (< 1)$$

\rightarrow small

tanh (approx.)

$$\text{parametric tanh}$$

sigmoid \rightarrow fan in, fan out \rightarrow dependent on this

inputs \rightarrow output per neuron

$$\tanh = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

center $\rightarrow 0$

slope will decrease around

as $z \rightarrow 0$, the center more rapidly

\rightarrow next-gen sigmoid due to

better / clear distinction b/w 0, 1

Updates in Backpropagation

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a_i^3} \times \frac{\partial a_i^3}{\partial h_i^3} \times \frac{\partial h_i^3}{\partial w_i} = \frac{\partial L}{\partial a_i^3} \times \frac{\partial a_i^2}{\partial a_i^3} \times \frac{\partial a_i^2}{\partial h_i^2} \times \frac{\partial h_i^2}{\partial w_i} = \frac{\partial L}{\partial a_i^3} \times \frac{\partial a_i^2}{\partial a_i^3} \times \frac{\partial a_i^2}{\partial h_i^2} \times \frac{\partial h_i^2}{\partial w_i}$$

Input layer Hidden layer Output layer

$a_i^j \rightarrow$ activation func
 $h_i^j \rightarrow$ Z value
 $j \rightarrow$ layer no.
 $i \rightarrow$ neuron no.

Multilabel

Multiclass

Loss func depends on prob.

→ Hyperparameters

- Grid search

→ Number of Hidden layers, I/p layer, O/p layer
 relation b/w I/p \downarrow O/p \downarrow \rightarrow # clauses/labels

→ # Neurons / layer

→ Activation fn

- Vanishing gradient \rightarrow too small
- Exploding gradient \rightarrow too large

→ Weight Initialization

\rightarrow Xavier init or Glorot init

Glorot \rightarrow random number from known distribution

\rightarrow number of ~~neur~~ should be in a particular range

$$\rightarrow \sigma^2 = 1/\text{fan out}$$

$$\rightarrow \text{fan out} = \frac{(\text{fan in} + \text{fan out})}{2}$$

\rightarrow Activation - None, tanh, sigmoid, softmax

\Rightarrow taking max due to symmetric behaviour

He \rightarrow Activation \rightarrow ReLU and variants

$$\rightarrow \sigma^2 = 2/\text{fan in}$$

\rightarrow Based on obs that 50% neurons are switched off

Lecun \rightarrow Activation \rightarrow SELU

$$\rightarrow \sigma^2 = 1/\text{fan in}$$

\rightarrow orthogonal

$$\text{② Xavier Uniform} \\ b_{\text{bias}} = \sqrt{\frac{6}{N_{\text{in}} + N_{\text{out}}}}$$

init wt = 0

if we init with 0 \rightarrow no change in gradient \rightarrow no learning \rightarrow same o/p \rightarrow problem

init wt = 1

all 1/p treated equally \rightarrow no change \rightarrow not just for any constant \rightarrow problem

init wt = random

unstable gradients if variance not in check

GBLU

QR Decomposition

LLM wt init

Normalization & Standardization

Page No.:

Date: / /

If l/p distribution changes in test and train, we need normalization
initial wt (\rightarrow)

Or Marginal \rightarrow wt random wt with W

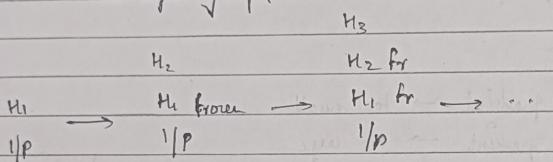
\rightarrow $W = QR$ (QR decomposition, SVD can also be used)

\rightarrow ~~WT~~ $W^T = I$ (constraint)

\rightarrow results exploding

\Rightarrow Unsupervised Pretraining

train layer by layer



\Rightarrow Gradient Clipping

\rightarrow Manage gradients by placing bounds.

If val is in range take it

If beyond threshold, take threshold \rightarrow this may change direction of movement
another approach is to normalize the weights so that

they are within the range

This keeps the direction same as before what it would be if we didn't clip.

\Rightarrow Batch Normalization

\rightarrow We want to maintain the same distribution for the inputs of all the hidden layers.

$$\hat{x} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$\mu \rightarrow$ mean

$\sigma^2 \rightarrow$ variance

$\gamma \rightarrow$ o/p scale parameter vector

$\beta \rightarrow$ output shift

$z \rightarrow$ o/p of BN also

$$z = \gamma \otimes \hat{x} + \beta$$

page
rank freq
 $R \propto f$
 $\log R = \text{page}$

ML: Intro to Non-differentiable / Reversible Reinforcement
Optimization

Normal/real
BN vs LN vs NIN vs Std
Page No.:
Date: 1-1

Batch v/s Layer Normalization | Normalization v/s Standardization

Gödel, Gödel, Power law

Normalization can be done on the basis of # neuron and layer by layer
also layer normalization
↳ # layer

BN handles both expl and vars. problem

LN handles only one (check)

Optimizer

→ Improvement on GD

Scheduler → manages/adjusts LR

Power, exp, c = power

Power scheduling: $\eta(t) = \eta_0 / (1 + t/\tau)^c$, t = iteration number, τ = total steps

Exponential scheduling: $\eta(t) = \eta_0 \cdot (0.1)^{t/\tau}$

Decay in constant scheduling: $\eta_0 = \alpha_i$ for K epochs, then a smaller LR for another number of epochs

Performance scheduling: measure val acc. every N steps and reduce LR by a factor λ (can be loss $f(x)$ or utility $f(x)$)

SGD with Momentum

initially: $w_{t+1} = w_t - \eta \Delta w_t$

with momentum: $w_{t+1} = w_t - \nu_t$

$$\nu_t = \beta \nu_{t-1} + \eta \Delta w_t$$

$$\beta \in (0, 1)$$

→ tackles — local minima

— saddle point

— High curvature

$$\beta(\nu_{t-1}) + \eta \Delta w_t + \dots$$

$$\beta(\beta \nu_{t-2} + \eta \Delta w_{t-1}) + \eta \Delta w_t$$

etc

$$= \eta \left(\Delta w_t + \sum_{i=1}^{\infty} \beta^i \Delta w_{t-i} \right) + \dots$$

$$= \eta \sum_{i=0}^{\infty} \beta^i \Delta w_{t-i}$$

→ Zeta step

⇒ Nesterov Accelerate Gradient
→ Look ahead

$$w_{t+1} = w_t - \beta v_t -$$

$$v_t = \beta v_{t-1} + \eta \Delta w_t = (w_t - w_{t-1}) + \eta \Delta w_{t-1}$$

$$w_{t+1} = w_t - v_t$$

→ may get stuck in local minima (check why?)

⇒ Adagrad (verify?)

$$w_{t+1} = w_t - \frac{\eta \Delta w_t}{\sqrt{v_t + \epsilon}} \rightarrow (?) \text{ check}$$

$$v_t = v_{t-1} + (\nabla w_t)^2$$

→ Not preferred in Deep NN

⇒ RMSProp
similar to Adagrad → just add decay (check)

$$v_t = \beta v_{t-1} + (1-\beta) (\nabla w_t)^2$$

⇒ Adam (check)

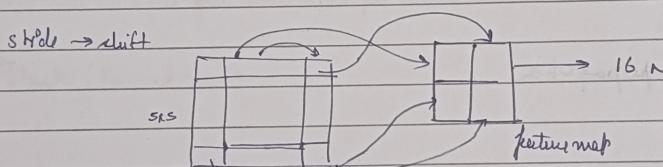
$$w_{t+1} = \frac{w_t - \eta \times m_t}{\sqrt{v_t}} \quad \begin{cases} m_t = \beta_1 m_{t-1} + (1-\beta_1) \nabla w_t \\ v_t = \beta_2 v_{t-1} + (1-\beta_2) (\nabla w_t)^2 \end{cases}$$

⇒ Regularization

- Dropout → n% neurons are switched off / disabled
- don't emit η_p or gradient
- sparsity in net
- explore multi. architectures using dropout → ensemble

CNN

→ If we process images ~~as~~ in the regular way, we need width nodes for each i.e. 1 neuron for each pixel



Kernel has ^{same} weights throughout traversal, each kernel for one has own wts.
Padding → To reduce bias of common area

→ dimensionality reduction of images

$$n_{\text{out}} = \left\lfloor \frac{n_{\text{in}} + 2p - k}{s} \right\rfloor + 1$$

n_{in} = # ip features

n_{out} = # op features

k = conv. kernel size

p = conv. padding size

s = conv. stride

→ 3 Filters Vertical Horizontal etc.

→ We can use n kernels for n different kernels and get n diff. feature maps to capture more features

$$Z_{ijk} = b_k + \sum_{u=0}^{f_h-1} \sum_{v=0}^{f_w-1} \sum_{l=0}^{f_{w'}-1} x_{i+u,j+k} \cdot w_{u,v,k,l} \quad \begin{cases} l' = l + S_w u \\ j' = j + S_w v \end{cases}$$

16×16 input 3 channels max filter \times filters
each filter gives $m \times n \rightarrow 3D \times m \times n$

→ Back propagation

Pooling

- Mathematical aggregation
- No learning
- kernel is traversed throughout
- No idea about loss and info. of each pixel involving kernel
- Adv. → dim. red.

↓ Translation invariance → Pattern / feature learnt can be identified at any location

→ Good for object detection

ARCHITECTURES

- AlexNet

- ResNet

- LeNet

... All but GoogleNet

Ch 14 HOML

Ques

Page No.:
Date: / /

Data Augmentation \Rightarrow $\downarrow \downarrow \downarrow \downarrow$ shift, brightness etc.

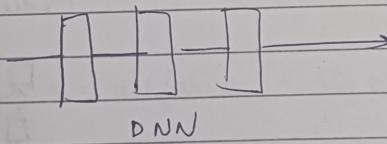
ANN
CNN
RNN
LSTM
GRU

Local Response Normalization (in AlexNet)
(formula from book)

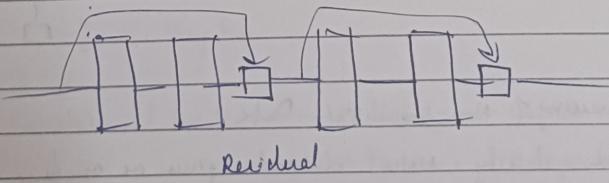
Where else can we use CNN

- Tabular data
- ID data
- But other techniques give better results for Tabular data.

Mark / Forward mark



DNN



Residual

Xception

No. of Parameteric neurons: 4F G

Each conv \Rightarrow Parameters = $(K_w \times K_h \times C_{in} + 1) \times C_{out}$

Fully connected : Parameters = $(\text{input units} \times \text{output units}) + \text{bias units}$

If BN : Parameters = $2 \times \text{num features}$

$\frac{1}{2}(K_w \times K_h \times C_{in} + 1) \times C_{out}$

EE-Project
Date: _____
Page No.: _____

Ph-3 — [14-20 Nov] — Lab - and an one eval — 404B ETL — 15-30 mins
Stat will be given
Report + Repo + Readme + code + what we did
(Template)
(20% phg)

Page No.: _____
Date: / /

$$67712 = 784$$

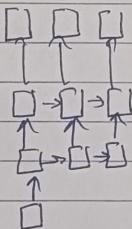
RNN Recurrent Neural Network

- Regular feed-forward networks can not capture sequential relations.
- Theoretically, CNNs can work, but it is not quite right.

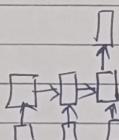
why CNN

- for spatial relations
- for lower input
- For for getting proper response for non-target containing inputs

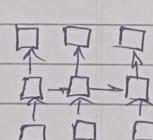
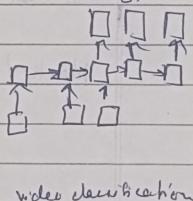
One to many
(long caption)



many to one
(video action prediction)



many to many
(Video captioning)

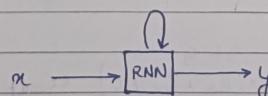


Sequential Processing of Non-sequential Data

→ RNN can classify MNIST dataset (given as a image)



RNN



RNNs have an internal state that is updated as a seq. goes on

Compare different outputs from different activation functions

Page No.: / /
Date: / /

Vanilla RNN

Hidden state update

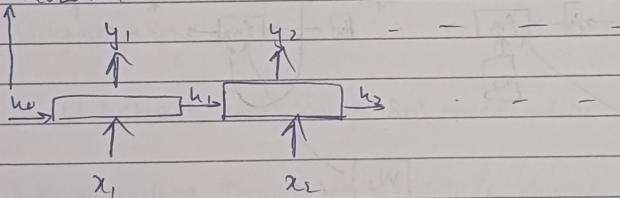
$$h_t = f_w(h_{t-1}, x_t)$$

new neuron ↓ old state input vector at some time step
 state some fat
 with parameters W

$$y_t = f_W y_t (h_t) \quad (\text{like activation } f)$$

output another for new state
 with parameters W_0

initialized value



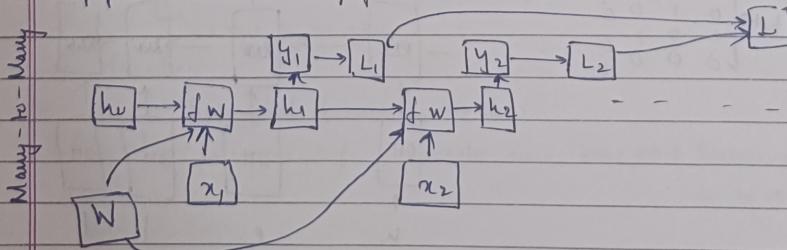
Elman RNN
in Vanilla RNN

$$\text{Eg: } h_t = \tanh(W_{hh} h_{t-1} + W_{hx} x_t)$$

$$y_t = W_{yy} h_t$$

Q. Why tanh
Why not ReLU
GEM
Sigmoid

Computational Graph \rightarrow Static \rightarrow complete required
 \downarrow \rightarrow Dynamic \rightarrow layer by layer - output
 Map of operations that will be performed



Similarly
 Many to One
 (2 types)
 One to Many

In one-to-many, what would be the input for the intermediate / non-initial nodes
 → replicate input (based on problem/modality)
 or
 → replicate output B

~~X~~ 1 → Garbage value (not preferred → garbage in garbage out)

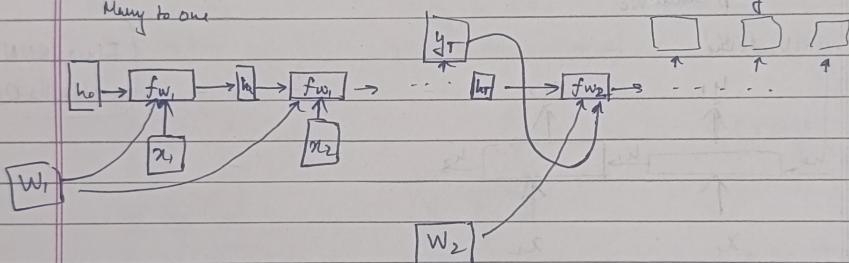
1 → 0s or 1s (Since W is same for all, it would not be beneficial)

Sequence to Sequence → Many to One + One to Many

Many to Many → character level language model

Many to one

one to many

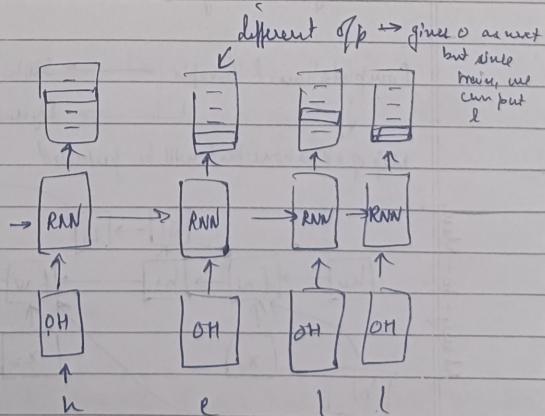


Here Char - LM

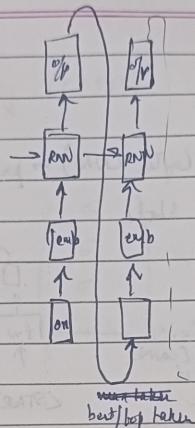
Train

hello → [h, e, l, o]
 one hot →

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1



Test



Backpropagation

In std. backprop \rightarrow contribution of first neuron reduces drastically \rightarrow vanishing

Better approach \rightarrow Truncate \rightarrow calculate in chunks

RNN Tradeoffs

Advantages

- \hookrightarrow can process any length $1/p$
- \hookrightarrow computation for step t uses info from prev. steps
- \hookrightarrow Model size doesn't increase for longer $1/p$
- \hookrightarrow same W on each

Disadv.

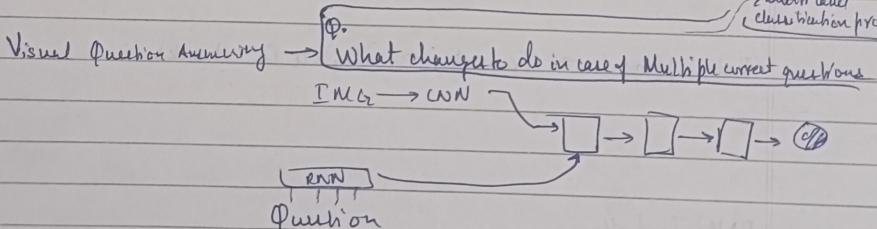
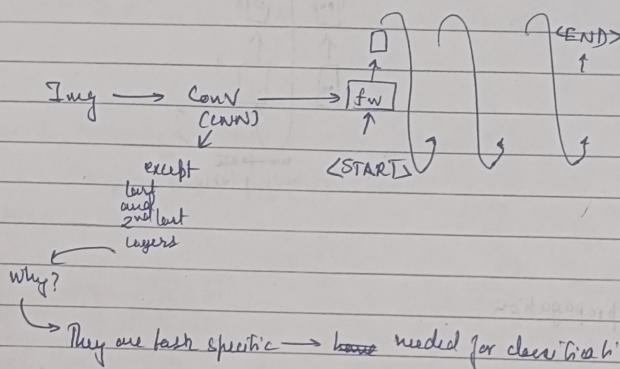
- \hookrightarrow Sequential processing \rightarrow slow
- \hookrightarrow previous info is lost after some time \Rightarrow + backprop is diff due to vanishing grad.

15-18 words can be worked with nicely in vanilla RNN

Page No.: / /
Date: / /

Image Captioning

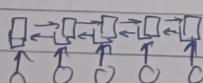
- Image as Input (after CNN) → problem: same ip for all
- Images as Hidden State



Vis QA → Dataset Bias

→ Generalization

Multilayer RNNs ~~task~~ (VERIFY)



GC

E.S. [Code may also be asked for End-Sem exam]

END - SEM

Page No.:

Date: / /

Long Short Term Memory (LSTM)

→ how much info
 $g \rightarrow$ gate gate → to write to cells
 $c \rightarrow$ cell state → long term

input $\begin{pmatrix} 1 \\ f \\ o \\ g \end{pmatrix}$ = $\begin{pmatrix} 6 \\ 6 \\ 6 \\ \tan h \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$

update ← forget
 info ← output
 o/p ← $g \odot g$

vanilla rnn

$c_t = f \odot c_{t-1} + i \odot g$

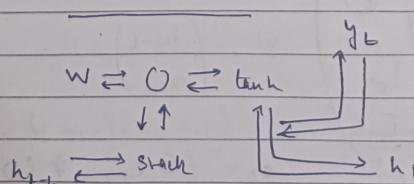
→ determining how much to write
 $h_t = o \odot \tan h(c_t)$

Vanilla RNN gradient flow

Backpropagation Through Time → calc. and backprop at each step (no need to wait for end)

$$\begin{aligned} h_t &= \tan h (W_{hh} h_{t-1} + W_{hx} x_t) \\ &= \tan h ((W_{hh} \quad W_{hx}) (\begin{matrix} h_{t-1} \\ x_t \end{matrix})) \\ &= \tan h (W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}) \end{aligned}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \tan h' (W_{hh} h_{t-1} + W_{hx} x_t) W_{hh}$$



$$\frac{\partial h_t}{\partial w} = \frac{\partial \tanh(W_{in}^T W_{out})}{\partial w} \left(\frac{w_{t-1}}{w_t} \right)$$

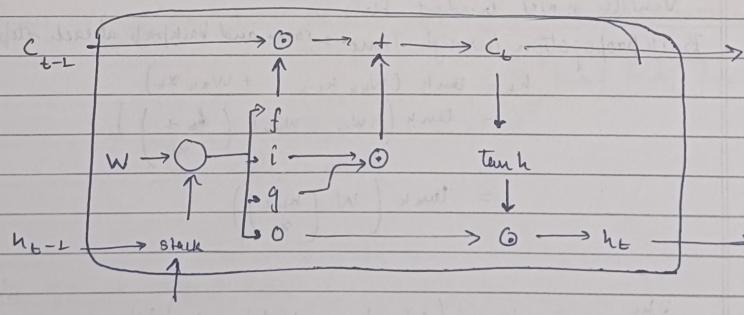
$$\Rightarrow \frac{\partial \tanh(W_{in}^T W_{out})}{\partial w} \cdot \frac{w_{t-1}}{w_t} = \tanh'(W_{in}^T W_{out}) \left(\frac{w_{t-1}}{w_t} \right)$$

$$\frac{\partial L}{\partial w} = \sum_{t=1}^T$$

Page No.: / /
Date: / /

- have to change RNN structure
- Vanishing gradients for tanh sigmoid
- Exploding for ReLU leaky ReLU
- Gradient Clipping

⇒ LSTM



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{tanh} \\ \sigma \\ \sigma \\ \sigma \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

↑ why → capture more → keep / del

LSTM - 43 words

Page No.:
Date: / /

But backprop has no W multiplication, so easy to calculate

Do LSTM solve the gradient problems?

→ LSTM makes RNN ^{it's easier for} to retain/reuse info

→ LSTM doesn't guarantee freedom from vanishing/exploding gradients.

→ LSTM Gradient Flow.

Backpropagation in LSTMs (VERIFY)

Get start

$$h_t = \sigma \odot \tanh(f \odot c_{t-1} + i \odot g)$$

$$\begin{array}{c} i \\ f \\ o \\ g \end{array} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \text{Tanh} \end{pmatrix} \left(\begin{array}{c} \begin{matrix} W_1 & W_2 \\ W_2 & W_3 \end{matrix} \\ \begin{matrix} W_1 & W_2 \\ W_2 & W_3 \end{matrix} \\ \begin{matrix} W_1 & W_2 \\ W_2 & W_3 \end{matrix} \\ \begin{matrix} W_1 & W_2 \\ W_2 & W_3 \end{matrix} \end{array} \right) \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$L \rightarrow \frac{\partial L}{\partial y}$

$$\begin{array}{c} \sigma'W_1h_{t-1} + \sigma'W_2x_t \\ \sigma'W_1h_{t-1} + \sigma'W_2x_t \\ \sigma'W_1h_{t-1} + \sigma'W_2x_t \\ \tanh'W_1h_{t-1} + \tanh'W_2x_t \end{array} \quad \begin{array}{c} i \\ f \\ o \\ g \end{array}$$

$$h_t = \sigma \odot \tanh \left(\left(\sigma'W_1h_{t-1} + \sigma'W_2x_t \right) \odot \tanh \left(\left(\sigma'W_1h_{t-1} + \sigma'W_2x_t \right) \odot C_{t-1} \right) \right)$$

$$\begin{aligned} \frac{\partial h_t}{\partial h_{t-1}} &= (\sigma'(W_1h_{t-1} + W_2x_t) \odot W_1 + 0) \\ &\quad + (\sigma'(W_1h_{t-1} + W_2x_t) \odot \tanh'(\sigma(W_1h_{t-1} + W_2x_t) \odot C_{t-1})) \\ &\quad \odot \tanh'(W_1h_{t-1} + W_2x_t) \\ &\Rightarrow \sigma'(\cdot)W_1 \end{aligned}$$

$$55/100 = \boxed{2.75}/5 \quad 2.75 \rightarrow 3$$

$$70 \rightarrow 3.5/5 \quad 34.75/40 \quad 43.75/45$$

Page No.:	51/3
Date:	1/1

$17 \quad 34.75/20$

$15+2.75 \quad 17.75/20$

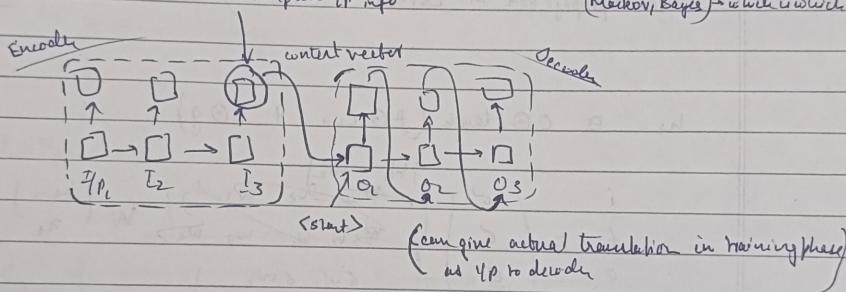
$34+2.75/45$

8.75

Other RNN variants

GRU \rightarrow forget gate and forget gate only
 ↴ ↵
 32/38 words

Encoder - Decoder Architecture



Why not put context vec. as hidden state to all?

↳ because prev. node ~~contains~~ captures all prev info, CN becomes un-updated if we keep going with that.

For long sentences context vec alone are not enough, so concatenate ~~will go on hidden state~~ (what?) and consider it as context vector.

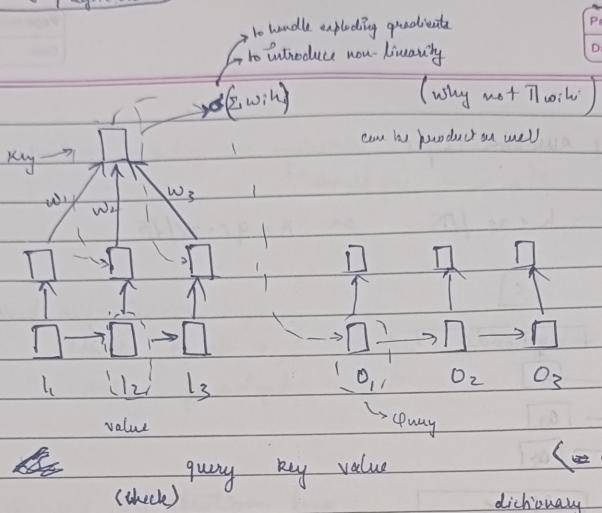
Problem

concatenation may lead to domination of some inputs. \hookrightarrow why?

Luang Attention
before self-att.

to handle exploding gradients
to introduce non-linearity

Page No.: / /
Date: / /



(why not $\prod w_i$)

can be product as well!

\hookrightarrow compare with attention

dictionary analogy

$$d = \{a: 1, b: 2, c: 3\}$$

Key $d[a]$
value
 \hookrightarrow query

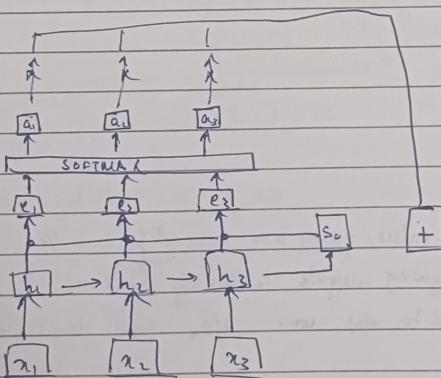
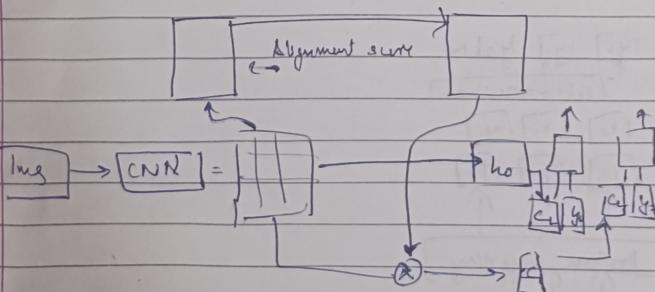
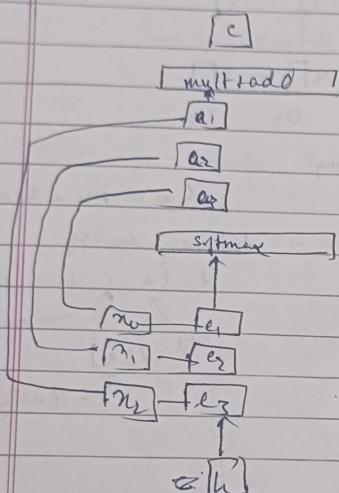


Image Captioning using Spatial features



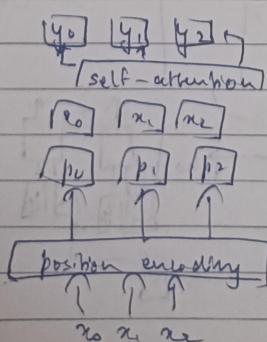
General Attention Layer

$$A = h \cdot \text{softmax} / \sqrt{D} \quad \text{or} \quad A = q \cdot p \cdot v / \sqrt{D}$$



Self Attention \rightarrow key, value, query are derived from 'p'
 \rightarrow Assigning weights to identify relations
and to stop some 'p's from dominating

Position Encoding



Q Relation b/w residual connections and layer normalization
Residual vs RNN

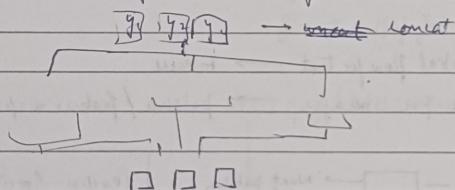
Page No.:
Date: / /

position encoding \rightarrow constraints
 \hookrightarrow should give unique vals

Masked self-attention layer

Mask & all $x_{i+k} \neq x_i$ because mask i/p has not been found yet.

Multi-head self attention layer

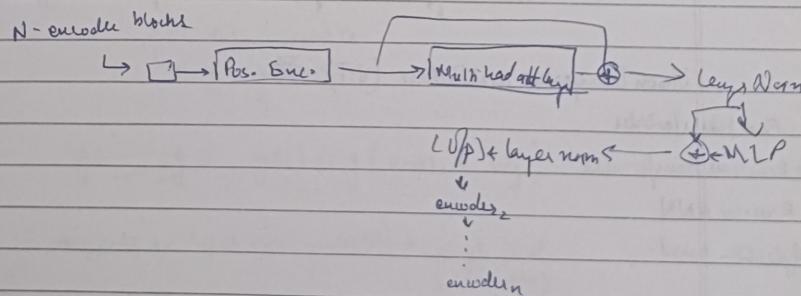


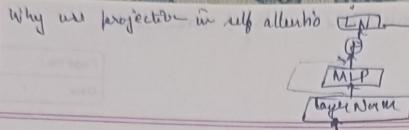
General att v seq Att

RNN v Transformer

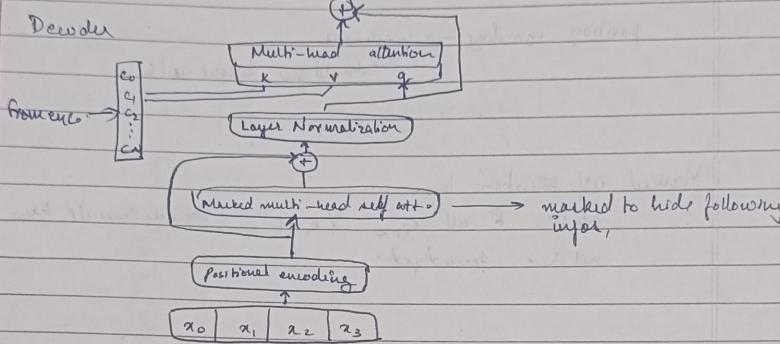
Normalization whenever residuals are added

Transformer



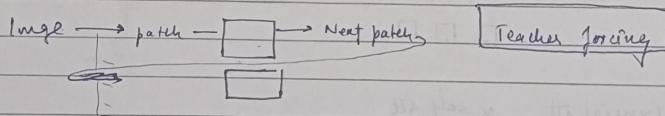


Page No.: / /
Date: / /



tokenize enforce sequential flow for text → tokens

convolution does same for images → patches / feature maps



Img → patch → [] → Next patch

Img → enc → [] → has info of entire space
How to handle different kinds of inputs (image, audio, text)

What off do we get at each step (what off and what type)

vision transformer architecture (ViT)

From papers/articles

↪ Recurrent Transformers

Recurrent RNN

Application based.

Complexity

Recursive NN → process structured data → can work with topologies
with graphs rather than feature based approach

initially only worked with acyclic str. under assumptions
structured domain has little math str.

$$g \rightarrow \text{output for } f \rightarrow \text{MLP}$$

$al(v) = f(al(\text{child}[v]), I(v), v, W_f)$

$$a(\text{ch}_i[v]) = [a(u_1[v]), a(\text{ch}_2[v]), \dots, a(u_{\ell_i}[v])]$$

$$a(nl) = a_0 \rightarrow \text{frontier state}$$

$$S \rightarrow \text{super source}$$

$$\hookrightarrow y = g(a(s))$$

$$y \cancel{c} = y(v) = g(a(v), v, w_g)$$

mapping of input str $\xrightarrow{\text{RECN}} \text{single node}$
(str)

max transfer \rightarrow an intelligent sys builds know by max transfer of prev know
max recov. \rightarrow an int. sys must be able to infer cause

Page No.:	1 / 1
-----------	-------

Recursion in existing architecture

Hierarchy → collects promised info at levels

info flow + logical net

→ sequential flow

B Improve state transition f and g

current pts W → update viewed as noise around W

E ↑ as var is ~~state~~ ↑ of dist. net

→ goal penaltie such pts