





# Vector data base

AWN, RAG  
maps sentence of  
paris to 504 dim  
vector.

N-gram lag model

push over words  
Bi-gram, Bi-gram  
N-gram  
(content vector)  
→ content sig. overweigh  
by concatenation of 60 vector

NA - large content ↑  
complex ↑  
expensive ↑  
spare data  
activation → information  
loss → sparse activation  
padding → free padding

Retrieval aug generation (RAG)

hallucination,  
Substitute knowledge  
Non-thesaurus  
with external memory  
Indexing  
Retrieval &  
generation

Disambiguation

external challenges  
generation with  
any combination  
knowledge

RNN

time series analysis,  
requester data  
(NLP), autokey (ARIMA)  
b for forecasting &  
w<sub>n</sub> bias are of 1

w<sub>x</sub> w<sub>h</sub> w<sub>y</sub> same  
in all hidden  
layer.

w<sub>n</sub> = tanh(w<sub>n</sub>h<sub>t-1</sub> + w<sub>n</sub>x<sub>t</sub>)

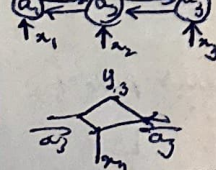
Signal - varying gradient  
ReLU - expressing gradient  
lost for long term  
due to vanishing gradient which push off all  
goes to out put.

Limitation

not from future  
K long term  
only L → t

Bidirectional model

L → R + R → L



update gate

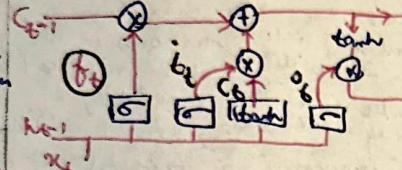
$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

one to one - image classification  
one to many - text generation,  
image captions  
many to one - sentiment  
analysis  
many to many - text summarization  
(x<sub>t</sub>, y<sub>t</sub>) machine language.

many to many - entity recognition  
(x<sub>t</sub>, y<sub>t</sub>) part of speech  
tagging

LSTM (Long term short  
memory)

gate & all memory.



forget gate how much to forget

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

activation = sigmoid  
[0, 1] remember  
forget  
current info  
previous info.

input gate

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

candidate layer

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

$$c_t = f_t \cdot c_{t-1} + i_t \cdot \tilde{c}_t$$

[Cell state]

output gate

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t \cdot \tanh(c_t)$$

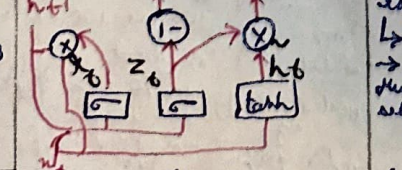
(filtered version)

Signal - varying gradient  
ReLU - expressing gradient  
lost for long term  
due to vanishing gradient which push off all  
goes to out put.

GRU (Gated Recurrent unit)

3 gates & no h<sub>t</sub> a cell state

forget & input gate =  
update gate



update gate (how much info to pass)

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

2> forget gate

can control the info flow  
how much to forget

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

3> candidate hidden state

current memory contents  
h<sub>t</sub> = tanh(W<sub>c</sub> · [h<sub>t-1</sub>, x<sub>t</sub>])

$$h_t = (1 - z_t) \cdot h_{t-1} + z_t \cdot h_t$$

control both previous  
step & current step  
info using update gate

→ updat gate → 0  
(forget everything)

1 → remember  
past all.

no of parameter

$$RNN = U(U + i + 1) \text{ param}$$

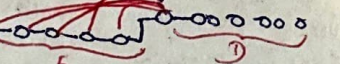
$$LSTM = 4 \times U(U + i + 1) \text{ param}$$

U = hidden unit  
i = input unit  
GRU = 3 U(U + i + 1)

BLEU model is used for NLP  
and other related task

Sequence & sequence model

CNN + LSTM - image captioning  
& vision transformer  
(one to many)



skip connection,  
to preserve previous info  
& to learn early word  
knowledge.

Transformer

- NN arch + attention based

- self attention to compute  
input sequence representation,  
capturing long-term dependency  
& enabling effective parallelization

- multiple machine translations

- Encoder & decoder

[multi head self attention]  
Add & normalized  
feed forward

also skip connection/  
residual connection &  
layer normalization

→ to prevent over fitting  
→ positional encoding (PE)

multi head - allows model to  
self attention from on input  
simultaneously

feed forward - helps to learn  
complex relationship  
coming from MHSA.

info. is processed  
parallelly as to loss order

"Words for words are not sequential"

$$x_1 = [J_{512} + [0, \dots, 0]_{512}]$$

$$x_2 = [J_{512} + [1, \dots, 1]_{512}]$$

embedding will be insignificant as PE ↑  
as E ↓ → D will not be good.

Positional encoding

due to parallel processing → loss order &  
sequential info → so it uses PE.

- uses sinusoidal function - series of  
series of sine & cosine wave with varying  
frequency. pos → position PE = sin(2π \* pos / 10000 \* 2d)

i → dim index PE = cos(2π \* pos / 10000 \* 2d)

d → total dim.

- no. of xops = no. of embedding

- multiple capture to get multiple  
position of a word.

→ word embedding + positional encoding

dim. of PE = dim of E

→ don't lose meaning of word after applying  
word.

# of encoder decoder → hyper parameter.

word → tokens (T) words (J<sub>512</sub>)

5k / 8 (head count) 64

attention (Q & V) correlation along word

→ semantic embedding for each word.  
(H<sub>1</sub>, H<sub>2</sub>, ..., H<sub>n</sub>) → concat → skip connect - FF

BERT (bidirectional encoder rep from transformer)

- drop decoder  
- multi-layer module  
→ pre-trained with task

1> stacked language modeling (LM) - 15% marked

2> strength of bidirectional & content aware

3> short sentence prediction (NSP)

→ build relationship with sentences.

NLP → cross entropy loss

NSP → binary entropy (Isent of word sent)

Addition of new tokens

1> CLS - classification

2> SEP - separation of sentence

IFP → token emb + segment + position  
emb for each word initial token position of each token in the sequence

(Transformer)

Global loss - rank (H<sub>1</sub>, H<sub>2</sub>) + NSP loss

Union Transformer (ViT)

divide the image into patch