
Prepared and Written By
RAZEEN AHMED

CSE-440
[NATURAL LANGUAGE PROCESSING]

HANDWRITTEN NOTE

Computer Science and Engineering
BRAC University

Github: github.com/razeen
LinkedIn: linkedin.com/in/razeenahmed

We will only focus on Natural languages (the language which grew with time)

→ It grew organically.

Ambiguity :-

Reasons:-

- (i) Phonetics. (How we use sound)
- (ii) Morphology.
- (iii) Syntax → A single sentence may have multiple form interpretation.
- (iv) Semantics → A single word may have multiple meaning

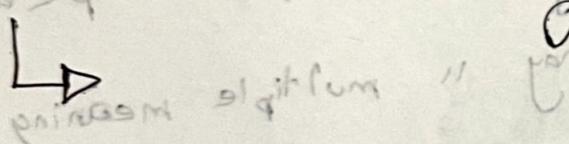
VARIABILITY :-

- (i) Same meaning sentence written in different way.

Supervised Learning Agents: - The agent who learns using examples.

Data annotation: - The inputs which is given to the system so that the system can learn from it. (The first step for any learning)

Data annotation must go through some processes.



Sentence segmentation:

Learning based Model:

Challenges:

(i) Non deterministic soln

Rule based Model

* There are some information we need to encode, we can't generate that.

Challenges:

(i) hard to cover all the rules.

(ii) maintenance

Tokenization:

Each language requires its own tokenization.

Lemmatization: - It is possible to solve to understand the PC that these are same words.

(Never use stemmer) → +, then and loses the meaning

Vectorization: - A word becomes multi-dimensional in space.

challenges:- contexts of similar words

NLP common libraries such as spaCy, NLTK, CoreNLP

(i) ~~spaCy~~ *

(ii) NLTK *

(iv) CoreNLP

(v) Processors *

AI first coined in 1950s.

* Since advent of AI there was a tendency to manipulate language

→ propositional logics is the base of AI

Parts of Speech (POS) Tagging:

→ 36 tags → golden ratio.

Recurrent
Neural
Network

NLP's major challenge is Ambiguity.

(i) Closed class b> (ii) Open class

↳ categories have fixed set of words

↳ limit to

b> (ii) Open class

↳ categories have a growing set of words

a → cont challenges:

→ One word can have different POS tag based on its context.

To maximum overall probability we will use global maximum (greedy algorithm)

Identifying names:-

Why important? →

→ Using the technique's of POS tagging we will find out names.

Name identifying

task →

Rifa

is multi word

Tasnim

split

"Rifa"

split

→ (209)

Tasnim

↑

multiple person

in same university
with this name.

Ambiguity

Washington DC

George Washington

→ 1 word can have multiple names.

words are multidimensional.

10/10/2021

so point A

• multiword challenges → where to start & where to end.

* Co Reference Resolution:

when we face a pronoun which noun we are referring

→ Allows machine to understand the grammatical structure of language.

* Parsing : → Analyze g. structure of a sentence to understand its components and the relationships between them

Challenge → Some sentences can be parsed in multiple ways → Attachment Ambiguity

→ Coordination "

Session-2 :

Supervised L. → entire learning is supervised by examples.

→ Identify the feature

→ map the features into a output.

Classification 1st step is → "which feature we need"

	f_1	f_2	f_3	\dots	f_m	
①	Y	Y	N		N	Y and WNE
②	N	Y	Y		Y	WE
③						

A technique

~~Bar graphs or words~~ of
creating a vocabulary
and creating column

Example: brown

- (1) pretty good, not bad
- (2) pretty bad, not good

f_p f_g f_b f_{not} f_{comma}

(1) 1 1 1 1 1

(2) 1 1 1 1 1

can or call have f₁ f_{ne} f_i f_{on} f_v f_{insub}

(1) ones 1 0 1 0 1 0 0 0 0

(2) 10 1 0 0 1 0 0 0 0

(3) 1 0 0 1 0 0 0 0 0

(4) 1 0 0 1 0 0 0 0 0

technique of contrast set up

tricks in contrast set up

task → contrast " " → contrast set up

multiple choice questions

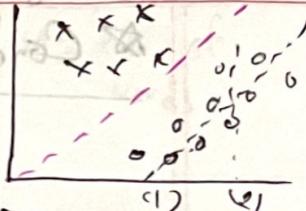
multiple choice questions

multiple choice questions

Machine classifiers
Find out the probability

$$P(Y|O)$$

Discriminative:
Creates an optimal boundary between two different classes



Generative:

e.g. logistic r,
SVM,

Advantage of Discriminative:

- (i) It's very quick and inference type.

When its useful (Discriminative)?

How fast busy
and 1st

Generative classifiers:

example: Naive Bayes

Working idea

- : falls

+ : not falls

* Conditional probability

* Bayes theorem

* Prior probability

$$P(+|O), P(-|O)$$

$$P(+ | f_1(O), f_2(O), \dots, f_m(O))$$

$$= \frac{P(f_1(O) \cap f_2(O) \cap \dots \cap f_m(O) \cap +)}{P(f_1(O) \cap f_2(O) \cap \dots \cap f_m(O))}$$

$$= P(f_1(O), f_2(O), \dots, f_m(O) | +) P(+)$$

* Condi

$$= P(f_1(O) | +) P(f_2(O) | +) \dots P(f_m(O) | +) P(+)$$

for not O

additive rule

* Joint probability

The event which may occur together.

can calculate when they are independent of each other.

Use more training data
than the probability
do increase decrease the
probable probability of
getting unknown words

Test data must be
completely unknown

* Conditional Independence

$$P(A, B | c) = P(A|c) P(B|c)$$

$$= P(A, A \text{ given } B | c)$$

$$= \prod_i P(f_i | B)$$

Classifiers

Step → make vocab table. (Convert each e.g. to feature vector)

+ S = 10	training data	f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9	f_{10}	f_{11}	f_{12}	f_{13}	f_{14}	f_{15}	f_{16}	f_{17}	f_{18}	f_{19}	f_{20}
Sorry! I'll call later		0	1	1	1	0	1	0	0	0	1	1	0	0	0	0	0	1	0	0	
You can call me now		0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	
You have call not		2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
Test																					
Sorry! You can not subscribe		1	0	(0 -)9	0	(0 +)9	0	(0 +)9	0	0	0	0	0	0	0	0	0	0	0	0	

$$\underline{P(+ | 0)}$$

$$P(+ | 1+) = \frac{0}{10} = 0$$

$$P(+ | 1+) * P(sorry + | 0) . 7 / 9 = P(call +)$$

problem

$$(+)9$$

Total words $\rightarrow 12$

Add 1 normalization

$$\text{Solve: } P(!|+) = \frac{O+1}{10+M} \xrightarrow{\text{Since } M=12} \frac{1}{10+12}$$

| M = no. of words

$$\rightarrow P(f_1(0) \cdot f_2(0) \cdots f_R(0) | +) \quad P(+)$$

$$= P(f_1(0)|+) P(f_2(0)|+) \cdots P(f_R(0)|+) \quad | P(+) = \frac{2}{3}$$

no. of exclamation present in Q is 0 (check the file)

$$= \frac{O+1}{10+12} \times \frac{12}{22} \times \frac{2}{22} \times \frac{2}{22} \times \frac{2}{22} \times \frac{2}{3}$$

! sorry v 1 can

$$= .0.000023$$

$$P(-|0) = P(!|-) P(V|-) \quad | P(=) = P(somg|-) P(mL|)$$

$$= \frac{2+1}{7+12} \times \frac{1+1}{7+12} \times \frac{1}{3} \times \frac{1}{7+12} \quad |$$

$$= \frac{1}{7+12} \times \frac{1}{7+12} \times \frac{1}{3} \times \frac{1}{7+12} \quad |$$

\Rightarrow lexicon \rightarrow set of relevant words which is word related

- [with probability \rightarrow if it is present in words]
- [else if it is not present in words]

Morphemes

NOTES

TOKENIZATION:- It is the process of breaking down text into smaller meaningful units.

Challenges :-

(i) Ambiguity

(ii) Handling punctuation

(iii) Compound words

(iv) Contractions

→ "Similar words look different"

Solution:-

(i) Stemming : chopping of suffixes & prefixes from words based on set of rules without understanding context or meaning.

→ Cut out common substrings.

(ii) Lemmatization : it reduces the word to their base root form and considers parts of speech.

→ Strip words to fit a base form

(iii) embeddings : vector representation of words, phrases.

→ similar meaning word have similar vector representation.

→ Faster, inaccurate. Resulting base word may not be a real word in the language, causing loss of meaning

→ Lemmatizers rely on a lexicon to ensure the base form is a real word, which will be meaningful and grammatically correct

→ Slower, accurate. [words with multiple meaning]

[More complexity & processing time as it relies on dictionary or linguistic rules]

↑ Errors in POS ...

POS TAGGING:-

The process of assigning a part of speech to each word in a sentence based on its context and definition. [Assigning grammatical categories to words]

i) → Context of similar words

i) → Requirement of high computational resources & memory.

She reads book

SHE → PRPN
reads → verb
book → noun

PROBLEMS:-

One word can have diff pos tag, e.g. painted the room

The " room
adj
of
verb

NAME ENTITY RECOGNITION (NER)

Identify and classifies entities into predefined categories

Challenges:-

(i) Ambiguity → apple (company & fruit)
Jordan (country & person) need context to be correctly classified.

Solution:-

(i) sequence tagging.

PARSING:-

understanding → sentences grammatical structure
Challenges: → determine → syntactic

→ many sentences have multiple possible interpretation or structure

Bag of words:-

Simplifies text into collection of individual words but keeping track of their frequency.

Problems:-

(i) One feature function per word (large, sparse matrices).

(ii) Completely ignores word order and context.

13/11/29

"Sentence" level representation does not work.

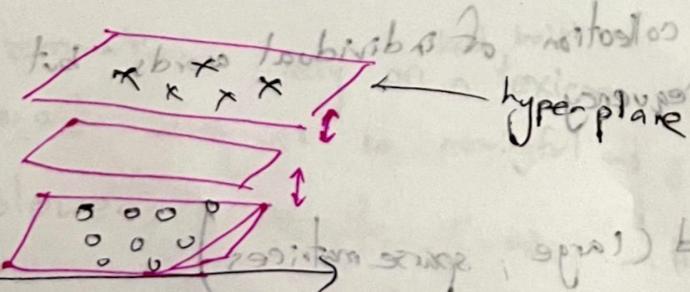
→ How to make ML Algorithm better?

We will use discriminative classifiers.

$x \times \times \times$ ← we have to identify this point
 $\times \times \times$ (this) → (compound words)
 $\times \times \times$ (birds) → (birds)
 $\times \times \times$ (birds & birds) → (birds & birds)
 $y = mx + c$

$x \times \times \times$ (meaning) → have to implement
 $\times \circ \circ \circ$ (sentences) → have to implement quadratic terms
 $x \times$ (sentences) → implement quadratic terms

↓ what if upward projection.



Target → to create a hyperplane in a dimension to create an optimal solution

Loss → To identify if the data is wrong. | why needed?

→ The amount of error.

We will use "Cross Entropy Loss".

Sigmoid Function:

$$\frac{1}{1 + e^{-y}} = \frac{1}{1 + e^{-(m+x+c)}} \quad | \text{Sigmoid}$$

Grain Function $\rightarrow \hat{y}^y (1 - \hat{y})^{(1-y)}$

$$\rightarrow -\log(\hat{y}^y (1 - \hat{y})^{(1-y)})$$

$$\rightarrow -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

(Binary Cross Entropy Loss) of function

$$x_m = 1$$

$$\hat{y}_m = 0.2$$

$$x_n = 0$$

$$\hat{y}_n = 0.2$$

(Good)

$$x_n = 0$$

$$\hat{y}_n = 0.9$$

Goal: minimize

<0 1 > <0 1 >

backward propagation

Every 10 example we will update! or \rightarrow ~~10~~ epochs
↑
This 10 is batch.

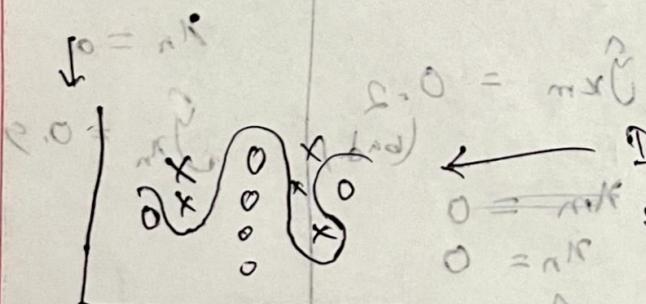
Full iteration in one entire dataset is 1 epoch.
biomips

when we stop, $\frac{\Delta}{\Delta} = \frac{1}{1}$

→ At one point the graph will hit a plateau
then we will stop.

we will reduce the loss function (gradient descent & BGD)

until we hit plateau.



It will be a very complex solution. So, its not recommended.

we want to have a simple solution to reach optimal soln.

[editors]

Goal:
Minimum loss
as simple as
possible decision boundary

DATA SPLITTING:

→ Training data \rightarrow train model

→ validation / Evaluation data \leftarrow evaluate

→ Test case

(i) type : hyper parameters

(ii) : initial values

good
Not good

to get the validation data :-

→ cross validation fit test at N fold validation.

Change parameters

change of learning

environment

updates

hyper parameters

Try to split the data as early as possible.

→ blonde mailing set

→ notulor no

→ all tip angles to flexors

→ notAno

the amount of splitters are no parameters

→ easier to train and it is not down out probability

High recall :-

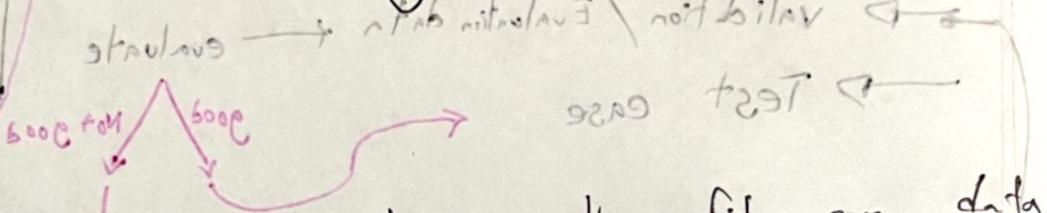
$$\text{TP} = \frac{\text{actual terms to on}}{\text{confusing terms}} - \text{positive}$$

$$\text{TP} = \frac{\text{actual terms to on}}{\text{confusing terms} + \text{FP}}$$

What is my fit?

We want to fit our model to the training data, which contains structures.

Underfitting :-



Overfitting :- When we create something fit our data 100%.

* Idea is to have to best fit not able to perfect fit, with bias in.

development / evaluation / validation.

Solutions:-

→ Non-linear equation

→ Increase dimensions.

→ The problem should fit our solution or else it will be overfitting.

PROBLEM :-

* Providing too much data with high amount of noise & outliers.

$$\text{Accuracy} = \frac{\text{no. of correct predictions}}{\text{All examples}} = \frac{TP + TN}{TP + TN + FP + FN}$$

Accuracy is not a good metric when the goal is to discover something rare.

When accuracy score is not a good choice?

- if the classes are imbalanced.
- if we want to focus on one single class.

Undersampling:- matching both intent & reality

problem - data are reducing how many times it is

Oversampling:- added in I think too much data are added

problem :- we are losing the real life information.
commit to overfitting

Precision:- we have to reduce the FP in case of $P_+ = \frac{TP}{TP + FP}$

2 cars \rightarrow 2 cars are flawless.
are produced

$$P_- = \frac{TN}{TN + FN}$$

The amount of operations we are performing should be more precise.

High recall:

$$R_+ = \frac{TP}{TP + FN}$$

$$R_- = \frac{TN}{TN + FP}$$

noise

TP

FP

FN

TN

Better Solution :-

$$\text{Harmonic mean P-F1 score} = \frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{2PR}{P+R}$$

which metric should be presented?

* Can we claim our model is better?

→ Model performance has to be significantly better, to significant no. of times.

P value :- sensitivity value.

Hypothesis testing :-

If H_0 \rightarrow reject the null hypothesis

End of session 3

Dimension

→ If the root words are same.

$\frac{TP}{TP+FP}$
→ In words

→ It's good idea.
Sentence to convert to number is not good idea.

→ Capture synonyms.
→ " antonyms

→ Looks same, different meaning

Performance metrics

Distributional semantics :-

Term Term Co occurrence semantics:-

 $|V| \times |V|$ matrix

vocab (unique words)

	a	an	and	the	is	of	to	in	on	at	for	by	with	from	as	near	near	near	near
a	125	10	25	95	200	125	200	125	200	125	200	125	200	125	200	125	200	125	
an	10	125	25	95	200	125	200	125	200	125	200	125	200	125	200	125	200	125	
and	25	95	125	10	200	125	200	125	200	125	200	125	200	125	200	125	200	125	
the	95	200	10	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
is	200	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
of	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
to	200	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
in	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
on	200	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
at	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
for	200	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
by	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
with	200	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
from	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
as	200	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	
near	125	200	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	125	

↳ how many times both the words came together in a sentence or paragraph

↳ true positive or TP

↳ proportion of correctly predicted connections with bag of words

Problems:-

initially it was lab a word 210000
 deal for all instances problem how out of lab
 # no. transitions has given distribution
 + mislabeling + high false negative rate

Overall

Rather better between tree (A) or C
 sub print is and actual 6/10
 something wrong with that was measures the ability of

first answer to correctly predict the identity of a word

→ P.A

(A) correctly predicted (1)
 (B) correctly denied (2)

Underfitting:- When a model has not learned the patterns in the training data well and is unable to generalize well on the new (given) data.

* Poor generalization

→ When a model is too simple to capture data complexities.

→ It is a scenario where a data model is unable to capture the relationship between input and output accurately.

PROBLEM :-

(i) Poor performance in both training and test data.

Overfitting:- Overfitting occurs when a model learns the training data too well, ^{perfectly} including noise and irrelevant patterns rather than the underlying patterns [/] generalization trends.

* Excellent performance on training data, but poor performance on unseen test data.

Problem :-

Poor performance on unseen/new data.

e.g. →
(1) Healthcare Diagnosis (Protecting life)
(2) Search and rescue operation

lexicon :- list of words that have been annotated for a task

Performance metrics:-

Accuracy :- It measures the proportion of the correctly predicted predicted instances out of the total no. of instances.

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN}$$

Problem :-

- (i) The distribution of classes or categories are imbalanced.
- (ii) If want to focus on one single classes.

Precision :-

Ratio between the true positive or (+) and predicted (+) or (-)

→ Proportion of correctly predicted instances out of all instances predicted as positive.

Problem :-

- Not ideal for all scenario's
- e.g. → cancer detection.
- ignores false negatives.

Recall :-

Ratio between true (+) or (-) and actual (+) or (-)

→ That measures the ability of a model to correctly identify all relevant instances.

Problem :-

Unbalanced focus

→ does not ensure that the predicted positives are correct

→ not ideal for all scenario's

→ does not account how many incorrect (+) prediction it makes.

and that about 20% to test :- no false/ sent a lot of information

- contains something

F1 score: Harmonic mean of precision & recall.
put out to two combines p. and r. score
written to .01

Eg - if weapon detection in airport.

$\frac{2 \times PR}{P+R}$

Statistical Significance:

It is a measure used in hypothesis testing to determine whether the observed effect or relationship is likely to have occurred by chance or represents a true underlying effect.

→ why need? consequences to mitigate

P-value:

probability of obtaining the observed results if the null hypothesis is true.

⇒ Avoid random errors
⇒ Decision making
⇒ Programmatic

test value → p value
(reject null)

to put into test conclusion fall
another transfer H₀ fitting is a fit of model

not based on given to 2006
from 2006 to 2006
from 2006 to 2006

any big cat dog

	any	big	cat	dog
any	1	1	1	0
big	1	2	2	0
cat	1	2	3	1
dog	0	0	1	1

Individual Sentence → any big cat
Cat dog cat.

any big cat dog

	any	big	cat	dog
S1	1	-1	1	0
S2	0	1	1	0
S3	0	0	2	1

W.V S₁ S₂ S₃

	any	big	cat	dog
S ₁	1	0	0	0
S ₂	0	1	1	0
S ₃	0	0	1	1

$$\begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

elimination of grad

$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 2 & 2 & 0 \\ 1 & 2 & 3 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 2 & 2 & 0 \\ 1 & 2 & 3 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}$$

They are not connected at all.
 $\|v\|$
 $\|w\|$

Cosine similarity

It is better function to

calculate the separation between

two matrix.

To check whether our vector

are cosine are normalized.

In NLP length does not matter.

In NLP, because its

scalar multiplication



$$u = [0, 1, 0, 1] \xrightarrow{\text{normalize}} u = [0, 0, 1, -1]$$

$$v = [1, 0, 1, 0, 0]$$

$$w = [3, 0, 3, 0]$$

$$\begin{bmatrix} 0 & 1 & 1 \\ -5 & 2 & 3-8 \\ 1 & 1 & 0 & 0 \end{bmatrix}$$

have to normalize
with the 4th vector

We want something more semantically make sense?

TF-IDF vectorizing word

($\frac{\text{weight}}{(\alpha, w)} \cdot \text{pol}$)

= Assigning important values of a word in a certain context.

Prototype: $\{a, b, c, d, e\} = D$, sentence

→ \vec{v}_1

messi → football

raw count is not smart
idea.

petra → polter

→ gives more informative vector of a certain sentence.

(Term frequency): $b_{nn} \rightarrow$
calc.

$\therefore \vec{v}_{nn}$ not to get a small val

TF → $(\epsilon_{88}^b, b_{nn}) \neq$

$tf(w, d) = \log(\frac{1 + f(w, d)}{f(w, d)})$

any | big

ant | dog

S1

$(c+e) \text{pol} = (\epsilon_{88}^b, p_{nn}) \neq$

$\frac{\partial}{\partial} \text{pol}$

$\sum \epsilon_{nn} \cdot 0 =$

↑ raw frequency

$\frac{\partial}{\partial} \text{pol} \times \frac{\partial}{\partial} \text{pol} = \text{primitiv}$

$(\frac{\partial}{\partial} \text{pol}) \text{pol} = (0, p_{nn}) \neq$

$\sum \epsilon_{nn} \cdot 0 =$

$\Delta =$

$\Delta =$

$\Delta = \Delta \times \Delta = \text{primitiv}$

IDF for multiplier

$$\log \left(\frac{N}{f(w, D)} \right)$$

ekta goto kom context e present
hoy star names or particular
context e ei word tk important

Example:-

$$N = 1000 \text{ sentence}, D = \{ d_0, d_1, \dots, d_{999} \}$$

$$f(\text{sorry}, d_{783}) = 2 \quad | \quad f(\text{and}, d_{783}) = 10$$

$$f(\text{sorry}, D) = 100 \quad | \quad f(\text{and}, D) = 900$$

For sorry :

$$tf(\text{sorry}, d_{783}) = \log(2+2) \\ = \log 4 \\ = 0.6989$$

$$IDF(\text{sorry}, D) = \log\left(\frac{1000}{100}\right) \\ = 1.3979$$

$$\log 3 = 0.4778$$

$$\text{Multiplying} = 1 \times 0.4778 = 0.4778$$

For and : (important)

$$tf(\text{and}, d_{783}) = \log 10$$

$$IDF(\text{and}, D) = \log\left(\frac{1000}{900}\right) \\ = \log \frac{10}{9}$$

$$\text{Multiplying} = \log 10 \times \log \frac{10}{9} \\ = 0.0479$$

	w_0	w_1	w_2	\dots	w_{100000}
d_0				\dots	
d_1				\dots	
d_2	gibberg	get	How	work	fitting
\vdots					
diggs					

PMG \rightarrow Book
 M.R. \rightarrow MID

weighted form of term-term co-occurrence.

Q.P. Pattern without straight pattern and J.E.S. \leftarrow

Three q. \rightarrow theory questions \rightarrow symmetry

\rightarrow analytical question \rightarrow mathematical

utilizing better, complex q. \rightarrow real life situation

find out which is better in this particular scenario.

Cross VALIDATION

It is a technique in ML used to evaluate the performance of a model and prevent overfitting.

It splits the available datasets into multiple subsets to train and test the model iteratively.

Goal:- Not overly reliant on specific portion of training set.

CROSS ENTROPY LOSS:

measures the performance of classification model

* loss function

→ It is used to quantify how well the predicted probabilities match the true labels. **Penalizes** incorrect predictions by comparing predicted probabilities to true class labels.

→ C.E.L. is a mathematical loss function used to evaluate the performance of classification models.

It measures how well the predicted probability distribution aligns with the true labels by quantifying the uncertainty of the predictions.

Abide two bnd
Formula:

$$L = -[y \log(\hat{y}) + (1-y) \cdot \log(1-\hat{y})]$$

We will reduce the loss function until we hit a plateau.

BATCH:

Subset of the training data used during one iteration of training a model.

EPOCH:

One full iteration through the entire training dataset.

there is

Bag of words problems \rightarrow Alternate method. mat mat #

* No features overlap if we try to calculate similarity.

Perfect word representation :-

would ideally encode a word in such a way that all the semantic, syntactic and contextual information relevant to its meaning and usage are captured.

How we can do that?

- (i) Share lemmas ; e.g. \rightarrow mouse/mice,
- (ii) different word senses ; e.g. \rightarrow pe mouse vs pet mouse,
- (iii) Synonyms ; e.g. \rightarrow couch/sofa, car/automobile
- (iv) antonyms ; e.g. \rightarrow long/short, light/dark
- (v) word similarity ; e.g. \rightarrow cup/coffee, ear/road
- (vi) word relatedness. ; e.g. \rightarrow cat/dog, doctor/nurse.
- (vii) word arousal ; e.g. \rightarrow excited & angry \rightarrow high arousal
 \rightarrow Does not occur in every context relaxed & depressed \rightarrow low arousal

* Words can be represented as vectors, each entry in the vector represents a dimension.

Term Term Co Occurrence matrix:

It is a representation used in NLP to capture the frequency with which a pair of words occur together within a defined context such as within a single document, a sentence, + n words of each other.

It is:

$$X = |V| \times |W|$$

↑
no. of words
in vocabulary

Build steps:-

(i) build binary BOW for the sentences

(ii) Transpose it

(iii) Multiply it with

Cosine similarity:

$$\cos(V, W) = \frac{V \cdot W}{\sqrt{V^2} \sqrt{W^2}}$$

PROBLEMS:-

- (i) very sparse
- (ii) does not carry any contextual info
- (iii) " " represent how important a word is in a sentence.

$$TF-IDF_{wt} = tf(w,d) \times idf(w,d)$$

TF-IDF: Assign important values of the word MAX

It is a statistical measure used in NLP to evaluate how important a word w is to a document d within a collection of documents.

TF:

* TF → occurs many times in ~~in~~ ^{freq} ~~in~~ ^{count} specific context.

★ High TF vector \rightarrow informative

★ High IDF $\downarrow \rightarrow$ uninformative

* IDF → Does appears in ~~in~~ ^{freq} ~~in~~ ^{count} context.

Formula:-

$$tf(w,d) = \log \left(\frac{1 + f(w,d)}{\text{raw frequency of the count}} \right)$$

If a term appears frequently in a document, it's likely to be important in that document.

IDF:

→ Does not occur in every context.

→ Measures how important or unique a term is across the entire document.

Formula:-

$$idf(w,d) = \log \left(\frac{\text{no. of sentence}}{\frac{N}{f(w,D)}} \right)$$

↑
word no. of this
document

→ Rare words carry more information
→ Common words are less informative
e.g. → the, and, is.

$$(b,w)T_{bi} \times (b,w)T_f = T_{bi} - T_f$$

EXAMPLES:- \rightarrow regular trigrams apical : TGB - TF

stemming of equivalence \rightarrow equivalence facilitate a si TF
 of trigrams organization \rightarrow organization \rightarrow organ below \rightarrow trigrams
 it follows a pattern stomach

lemmatization :-

leaves \rightarrow leaf \rightarrow FT4
 falling \rightarrow fall \rightarrow FT*

embedding:- met \rightarrow dog \rightarrow word embedding \rightarrow TGB

$$\text{dogs} = [0.5 \quad 0.4 \quad 0.4]$$

$$((b,w)T + \Delta)_{\text{pool}} = (b,w)T$$

: TDE #

types for each in every context

FINAL \rightarrow word trigrams document

$$((b,w)T)_{\text{pool}} = (b,w)T_{bi}$$

Word Representation: Word Vectors

1) Skip gram (word2vec) | 2) Glove

(i) Sparse vectors

(ii) Dense vectors

WORD EMBEDDINGS:-

It transforms words or phrases from their raw textual form into continuous vector representation in a high dimensional space.

These representations capture semantic relationships between words.

two types:-

(i) CBOW :- (Continuous Bag of Words)

Predict the target word given the neighbouring words.

(ii) Skip-Gram :-

Predict the neighbouring words given the target word.

→ Skip-Gram Classifier:-

- it's not hardware intensive.
- it's a simple technique.
-

MS/CS/80

Glove is designed to generate dense vector representation of words.

PROBLEMS:

- (i) They are symmetrically dense. | (iv) Word2Vec cannot deal with unknown words.
- (ii) How can we combine the word vectors so that we can classify the sentence. — S_1
- (iii) The word vectors cannot capture the contextual information in a particular situation. — S_2

Solution:-

- (S1) → we can calculate the centroid vector
knowing If we use antonyms both of them will banish each other in to a particular dimension.

(S2) [better solution] → RNN can convert a word vector to a sentence level representation.

(S3) → using language models.

Language Modelling:-

• Artificial neural network
• generative model for π_{ti}
• discriminative model for π_{ti}

LANGUAGE MODELLING:-

- To assign probabilities to sequence of words or generate plausible word sequence.
- It is a statistical model that learns the likelihood of a sequence of words occurring in a particular language.

Why do we need it?

- (i) Understanding context & meaning.
- (ii) Text generation.
- (iii) sentiment & emotion analysis.
- (iv) Speech recognition.

Semantic properties of embeddings:

(i) Different types of similarity or words with similar meaning.

- Based on the context window words similarity changes.
⇒ Hogwarts — Harry, Hermione, Potter

(ii) Analogy / relational similarity.

- encode meaningful relationships between words.

ex — man is to King
woman is to Queen.

(iii) Historical context.

- The meaning of a certain word may vary based on what time period we are using it.

IMP FOR EXAM for exam :-

- relational similarity vector given at utilid.org
- prove that here relational similarity vector isn't working

There is no engineering technique which can remove the bias from a model.

SEQUENCE LEARNING

It is very important to capture the sequence

How? "Using POS tagging"

Fully furnished condo in the beautiful catalina foothills
Adv Adj Noun Prep Det Adj Prop

There to identify that this farming an connection
identify that ADJ comes before nountag with ADJ

Challenge ambiguity

Solution dictionary

WSJ we will assign tags for each word
but there is a problem 1 tag 15%. we might get only 70%
2 tag 55% accuracy at best

Solution 2 → We will create a big table for the ambiguous words and store the most commonly used tag.

This simple solution will provide a good accuracy.

MAXIMUM ENTROPY

HIDDEN MARKOV MODEL :-

\$ used for sequence modeling predict sequences of hidden state

For a certain event occurring at the present depends only on what happened before

* For word prediction we will look at the immediate words

If we want to predict the present words tags & look at past words → tag, lemma, casing info

STEPS OF ME :-

(i) If we want to predict the tag of a word we will make feature vector of the past words. We will name it $x_1, x_2, x_3, \dots, x_n$

(ii) then use $y = m_1 x_1 + m_2 x_2 + \dots + m_n x_n$ → linear regression classifier.

Revised MEMM set steps:

- (i) Make feature of every word of the sentence.
 Vector
 Using casing info, suffix, prefix
 - (ii) Then train the model.
 - (iii) Input a new sentence now predict the tags of the next words.
- One we will predictors of input
 bidirectionally. Far from left \rightarrow right, right \rightarrow left
 then compare both prediction and create a final prediction.

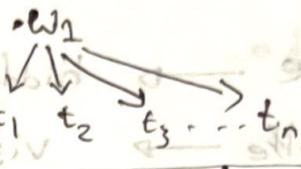
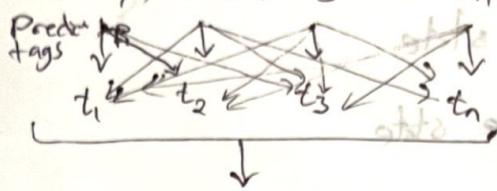
[Because, there is a less chance to make the same error from both direction]

If we incorrectly predict the previous word then it will nearly effect the prediction of the next words

MM = 29378

HIDDEN MARKOV MODEL:-

words $w_1, w_2, w_3, \dots, w_n$



PROBLEM :-

It will create exponential time complexity

we will calculate the probability of every assignments

we will predict the tags, pos probability of every tag assignments

the probability which will provide us make out that will be our final tag assignment

Select the maximum tag probability
Final tag selected

2nd approximation:-

★ Words Depend only on their POS tags.

Not their surrounding POS tags

Calculate the probability of $w_1, w_2, w_3, \dots, w_n$ given $t_1, t_2, t_3, \dots, t_n$

$$P(w_1, w_2, w_3, \dots, w_n | t_1, t_2, t_3, \dots, t_n) = P(t_1, t_2, t_3, \dots, t_n | w_1, w_2, w_3, \dots, w_n)$$

3rd approximation:-

★ A tag is dependent only on the prev tag

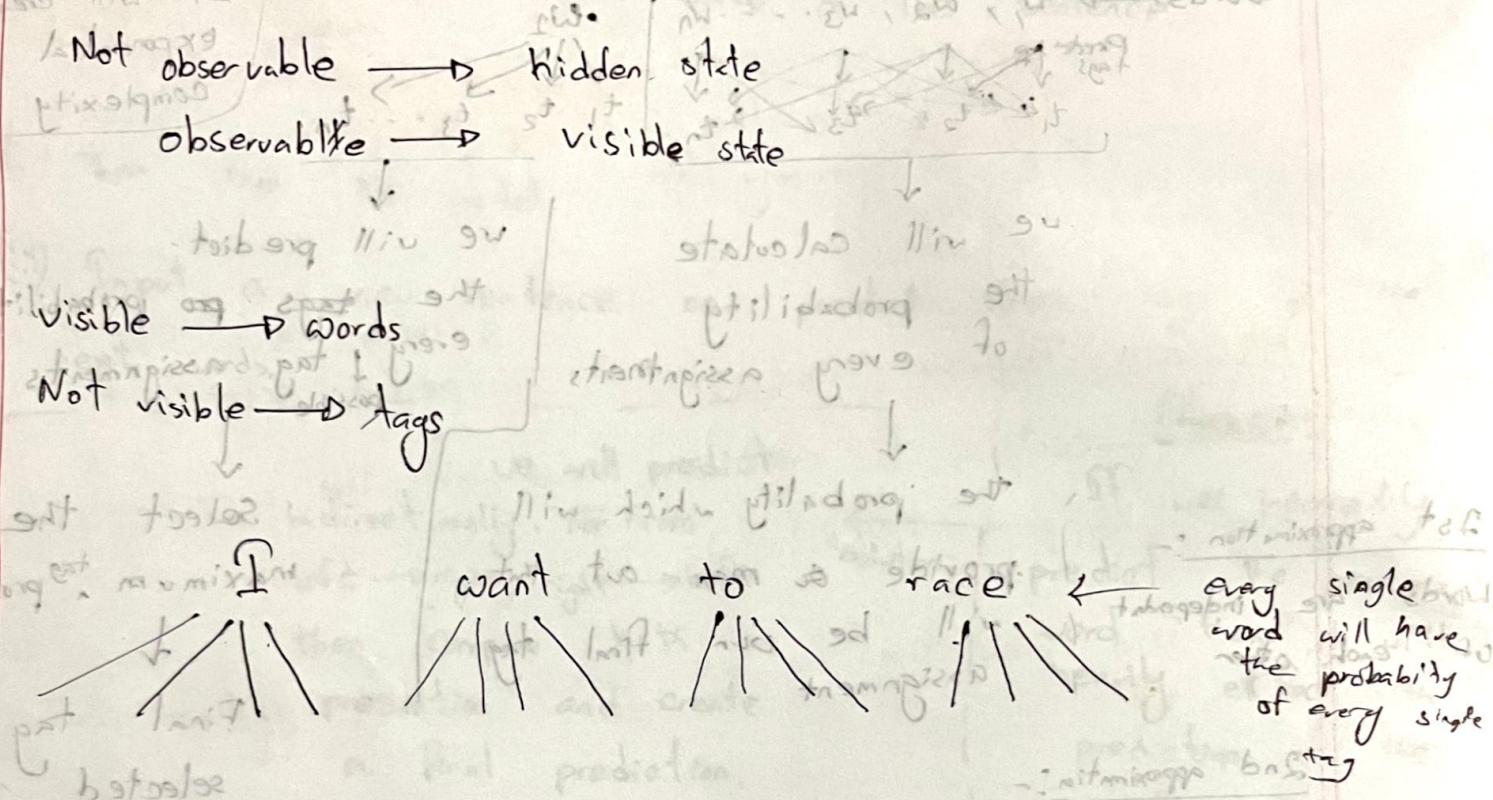
$w_1, w_2, w_3, \dots, w_n$

$t_1, t_2, t_3, \dots, t_n$

$$P(t_1, t_2, t_3, \dots, t_n | w_1, w_2, w_3, \dots, w_n) = P(t_1) \prod_{i=2}^n P(t_i | t_{i-1}, w_i)$$

HIDDEN MARKOV MODEL

* To infer hidden state of things from visible things.



↳ Goal → Find the globally best possible assignment.
 (we won't see local maxima)

$$\hat{t}_1 = \text{AM}(P(w_1^n | e_1^n))$$

$$= \text{AM}(P(w_1^n | e_1^n) P(e_1^n))$$

$$\approx \text{AM} \prod_{i=1}^N \left(P(w_i | t_i) \prod_{j=2}^i P(t_j | e_{j-1}) \right) P$$

The words are independent.

$$\hat{t}_i^n = \operatorname{argmax} \left(\underbrace{P(w_i^n | t_i^n)}_{\text{Words likelihood P.}} \cdot \underbrace{P(t_i | t_{i-1})}_{\text{Tag transition P.}} \right)$$

1st

For all possible (w_i, t_i) we create a table.

$$- P(w_i | t_i)$$

word likelihood table

Start

2nd

$$(w_1 | t_{now})^q$$

$$(av | task)^q$$

$$(uu | t_{now})^q$$

$$(avv | t_{now})^q$$

Trans Create a transition probability table.

$$(av | PT)^q$$

$$(or | av)^q$$

$$(ot | or)^q$$

$$(uu | ov)^q$$

INFERENCE PART [VITERBI ALGORITHM]

1

$$w_{ant}(or)^q$$

to

Initialization:-

$$P(t | \Omega) = P(\Omega | t) P(t | \langle s \rangle)$$

$$\Rightarrow P(\Omega | NN) \quad P(NN | s)$$

$$P(\Omega | VB) \quad P(VB | s)$$

$$P(\Omega | PPSS) \quad P(PPSS | s)$$

$$P(\Omega | TO) \quad P(TO | s)$$

1 want

PPSS

NN

TO

VB

values
for
loop
x
the
maximum
arrow
for
arrow
there
will
be
optimal
value
for
race

NN/PPSS

Optimal value for
race

race

race

race

race

race

race

prev x trans x words likely

$0.055 \times 0.0012 \times 0.025 = 3.3 \times 10^{-5}$

$0.055 \times 0.0012 \times 0.025 = 3.3 \times 10^{-5}$

$0.055 \times 0.0012 \times 0.025 = 3.3 \times 10^{-5}$

$0.055 \times 0.0012 \times 0.025 = 3.3 \times 10^{-5}$

$0.055 \times 0.0012 \times 0.025 = 3.3 \times 10^{-5}$

$0.055 \times 0.0012 \times 0.025 = 3.3 \times 10^{-5}$

$$\left(\frac{P(\text{I} | \text{want})}{\text{Want}} + \frac{P(\text{I} | \text{not want})}{\text{Not want}} \right) \times \text{sample} = \frac{\text{I}}{\text{Total}}$$

$$P(\text{I} | \text{want}) = P(\text{I} | \text{Want}) + P(\text{Want} | \text{I}) \cdot P(\text{I} | \text{not want})$$

$$= P(\text{Want} | \text{I}) \quad P(\text{Want} | \text{I}) \cancel{P(\text{Want} | \text{I})}$$

$$P(\text{Want} | \text{PPS})$$

$$P(\text{Want} | \text{NN})$$

$$P(\text{VB} | \text{VB})$$

$$P(\text{TO} | \text{VB})$$

$$P(\text{VB} | \text{TO})$$

$$P(\text{TO} | \text{TO})$$

$$P(\text{VB} | \text{NN})$$

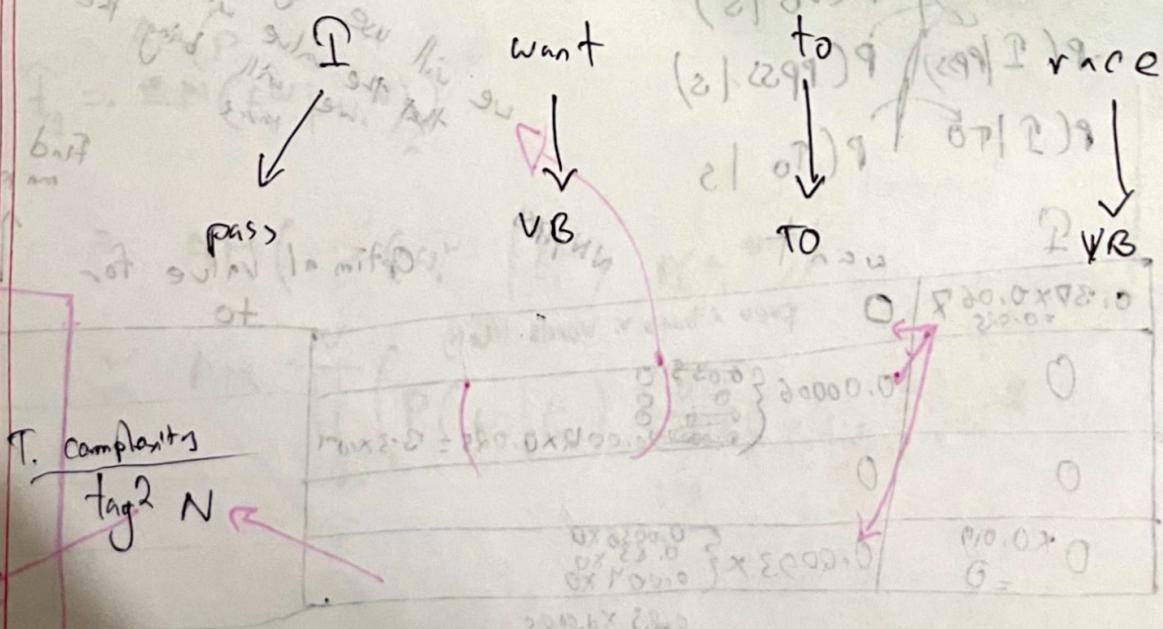
$$P(\text{TO} | \text{NN})$$

$$P(\text{VB} | \text{PPS})$$

$$P(\text{TO} | \text{PPS})$$

Ques
Speech language processing \rightarrow Durafsky, 8.2, 8.4

$$(I | \text{Want}) = \frac{P(\text{Want} | \text{I})}{P(\text{Want} | \text{I}) + P(\text{not want} | \text{I})}$$



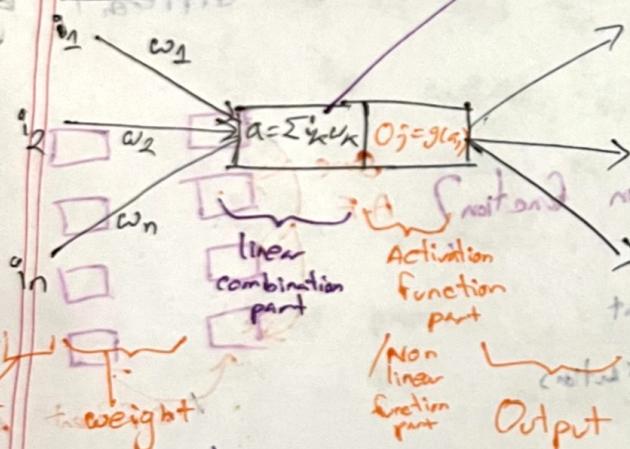
Disadvantage:-

- The number gets smaller as it moves forward.
- Without unknown word zero will not propagate forward.
- Hard to add features in the model. \therefore Slow to learn.

NEURAL NETWORKS:

NEURON:

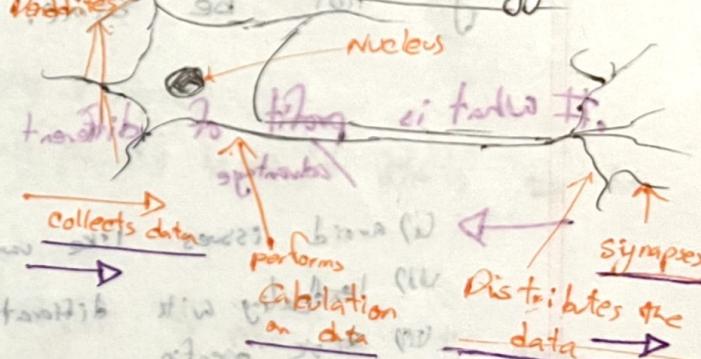
In terms of CS:



$$a = i_1 w_1 + i_2 w_2 + \dots + i_n w_n$$

12/12/24

In terms of Biology:-



Activation
function

$$O_j = g(a)$$

$$\Rightarrow O_j = \max(O_i, O_j)$$

TENSOR:

It is an N -dimensional / Multidimensional array of data.

The Neural Network :-

Fully connected neural network :- When every neuron stays connected with every neuron of the next layer.

★ Every \rightarrow The activation function of every neuron of the same layer will be same.

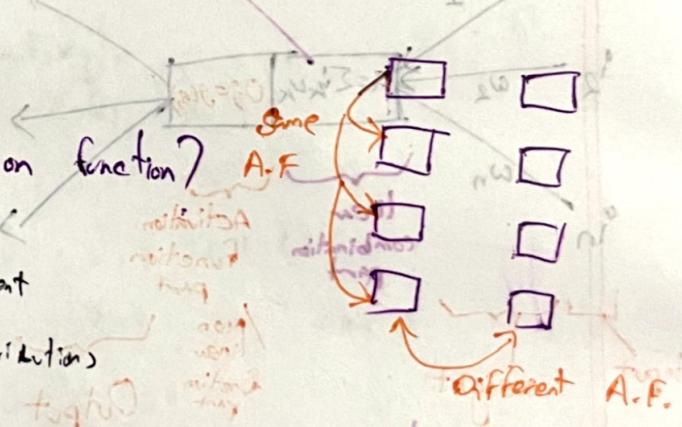
★ The activation function of neurons of different layers might not be same.

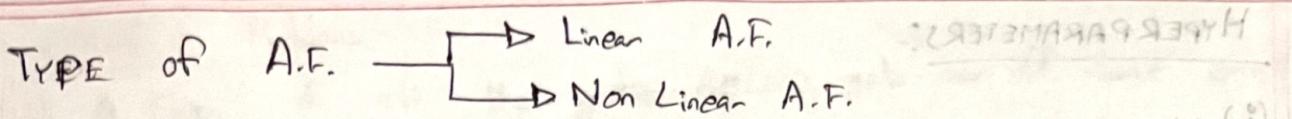
What is profit of different activation function?

- (i) avoid issues like vanishing gradient
- (ii) dealing with different data distributions
- (iii) task specific.

Why do we need Activation Function?

→ A.F. is a type of gate. It works as a gate in the neural network which decides if the inputs/values ~~or~~ should pass forward or not.





- * Why do we need non linear A.F.? → enable the model to learn and represent complex relationships in data.
- Our real world is too complex to model using linear functions. It is often impossible to classify real world data using linear function. [Real world problem often have complex non-linear relationships between features and labels]
 - The combination of two linear function is another linear function. We can't introduce any non linearity in the linear function.

- * Where would we find the weights?
- ⇒ Initially the weights are randomly assigned
 - ⇒ As the time passes, the weights are updated and it starts to make better predictions

The training process updates the weight to minimize the error between predicted outputs and the actual targets.

set the ML engineers

HYPERPARAMETERS:-



(i) Set HP \rightarrow how many layers would there be in the NN

(ii) In every layer, how many neurons do we need.

depends on the number of features \Rightarrow (iii) In every layer, which activation function should we use.

\Rightarrow $N \uparrow F \uparrow$ As it goes forward to next layers, size of hidden layer should be decreased

- (i) Number of neurons in each layer
- (ii) Number of layers
- (iii) Activation function.
- (iv) The Optimizer we need to use.

RELU:

Advantage :-

(i) It perfectly acts as a gate.

(ii) It is really easy to implement

PROBLEM

(i) Bad luck! (There might be a scenario where capacity,

that every neuron of the

neurons gets off, then the neuron

cannot possibly pass the data

forward to the next layer.

As a result, the network will

stop learning.

\rightarrow Non-differentiable at $x=0$ / breaking point

why differentiation needed?

\rightarrow For optimization.

Solution of Problems:-

- (i) Even if all the gates ReLUs gets set off. Some small amount of value will still pass through.
- (ii) We use ReLU to make the entire function differentiable by making the elbow like a small curve.

SIGMOID :-

→ switches between 0 to 1.

→ we use this as a output activation function

provides a probability value

TANH :-

→ Works very good in the intermediate layer as a activation function.

→ switches between -1 and 1.

provides a probability value

→ even though it has a negative slope, it is still useful.

→ it has a negative slope and is still useful.

[but it's not a good idea to use it]

Q.1

Q.2

Q.3 → ^{Q3} model to mitigate

OPTIMIZER :-

It is such type of algorithm whose task is to make the learnable parameters to learn using the optimum weight. It tries to reduce the loss.

Q.3

$P(\text{water} | VB) \cdot P(VB | NB) \cdot \text{prev.}$

If the situation given is we don't need to use prev. Because prev is included in the denominator.

REGULARIZER :- set of strategies to reduce the generalization error

↳ Engineering Hacks
↳ Mathematical regularization

(i) [DROP OUT]
Suppose, we have 100 neurons in a single situated layer. In every epoch we randomly turn off some neurons.

(ii) Early stopping

⇒ When validation reaches a flat surface we stop there.

[stops when performance gets saturated]

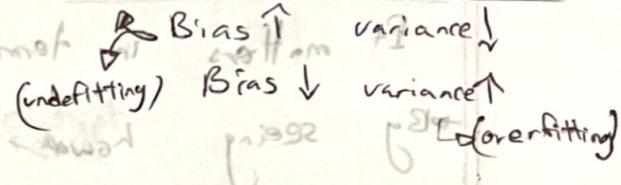
ADAM Optimizer (Advantage)

- (i) Easy to implement
- (ii) Needs no manual tuning
- (iii) Little memory requirements
- (iv) suited for large data.

HOPAD

(v) Data augmentation. (dotted note) → Bias vs Variance trade

⇒ Introduce new training data with more variation.



Mathematical Regularization :-

→ L1 Regularization:-

weight → modulus

add all weights

$$L_1 = |w_1| + |w_2| + \dots + |w_n|$$

[selects the best features]

Bias →

partiality without any observation

Ridge

Bias ↑ observation

what less observation

Bias is the higher

Variance

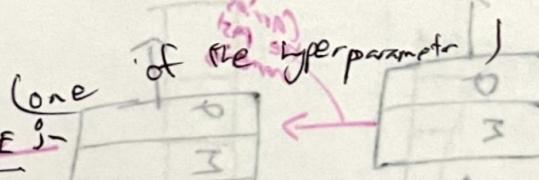
To reduce overfit
introduce a little bit of bias.

→ L2 Regularization:-

[Priorizes the best features.]

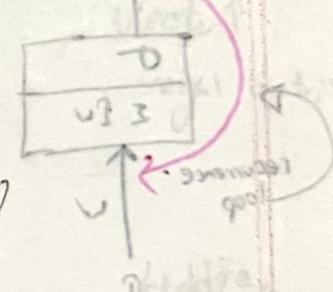
$$L_2 = w_1^2 + w_2^2 + \dots + w_n^2$$

LEARNING RATE :-



It is the lambda value of the regularization term?

→ how quickly a network updates its parameters.



Epochs :-

One complete pass through the entire training dataset during the training of a model.

point known on about (ii) training of gen (i)
about equal not better (iii) changing program diff (iv)

BATCH

~~BATCH~~ SIZE :- (Mini batch), not training set (i)

It matters in terms of time, subset of training data.

By seeing however

→ not developed position (AM)

or can

RNN ← sequence learning classification AS A

we have to handle sequential data.

18/01/2018

chap 11.0 bbs

last test set 200/200

student

A

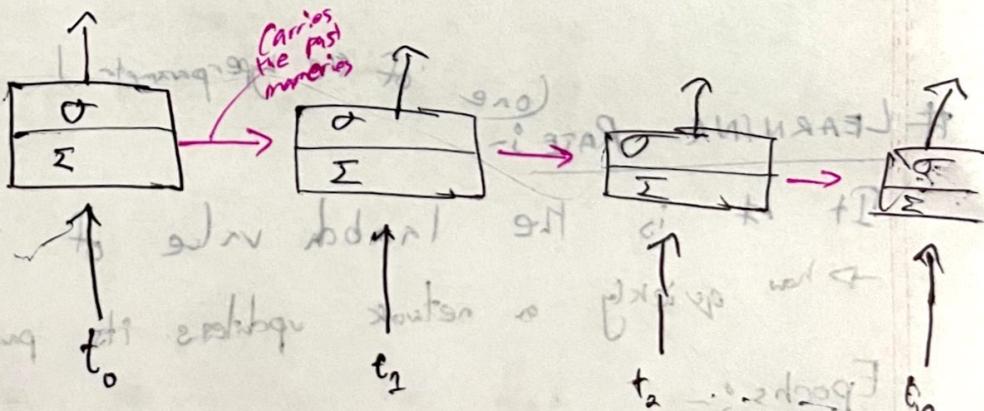
AM

A STUDENT

SJ

[student + AM = position]

$$\sigma_w + \dots + \sigma_2 + \sigma_1 = s$$



last all points taken against entire set avoid any delay in one lesson to 70

Every element
Own neuron.

Encoder architecture

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \Delta X$$

Decoder architecture: One to many \rightarrow output

Encoder - decoder architecture: e.g. summarizer
(many to many)

EXAMPLE Many RNN form

HYPERTPARAMETERS complex t_{age}

$$h_t = \tanh(wx_t + Uh_{t-1})$$

(i) what will be the activation function.

$$= \text{act}(Uh_{t-1} + wx_t)$$

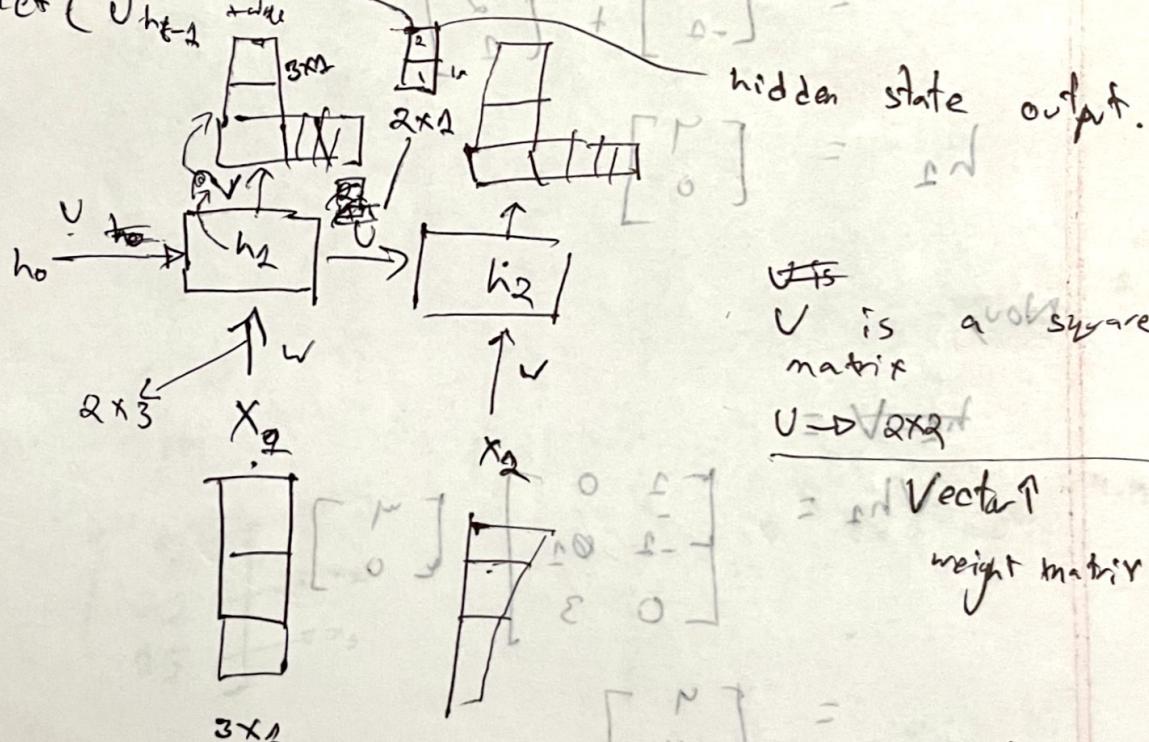
Example :-

$$t_{\text{gs}} = 3$$

A0)

DET

NN



h

hidden layer dimensionality

$$\text{Sigmoid Function} = \frac{1}{1 + e^{-x}}$$

$$x_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

\rightarrow activation vector

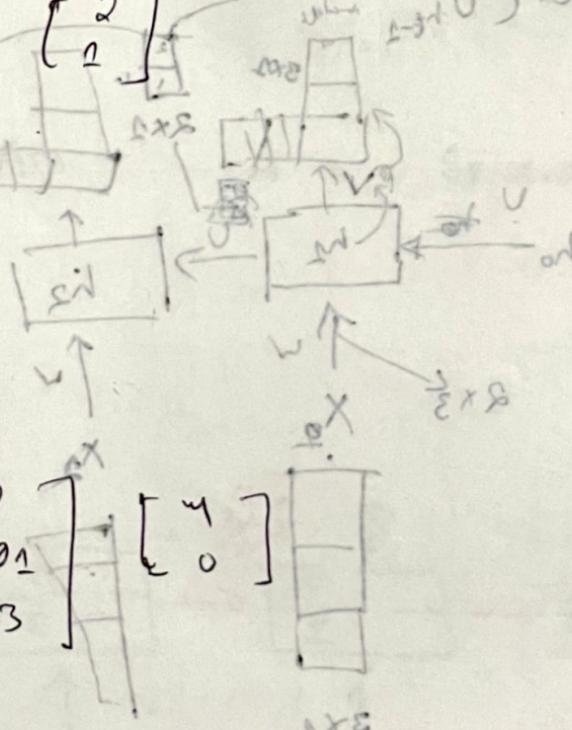
t_{10} bias \leftarrow plus at end; output vector

$$h_2 = U \times h_{t-1} + w x_1$$

$$= \begin{bmatrix} 2 & 0 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 & 0 & -1 \\ 2 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 2 \\ 2 \end{bmatrix} + \begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$$

$$h_2 = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$$



$$\cancel{h_2} = U$$

$$\cancel{h_2} = V$$

$$\cancel{h_2} = h_2$$

$$\text{bp directional RNN} \rightarrow o(t) = \frac{\sigma(x^t)}{\sum e^{x_i}} \begin{bmatrix} v_h \\ h_t \end{bmatrix}$$

$$h_2 = h_{t-1} + w \times 2$$

$$= \begin{bmatrix} 2 & 0 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 4 \\ -2 \\ 0 \\ 2 \end{bmatrix} + \begin{bmatrix} 2 & 0 & -1 \\ 2 & 2 & -2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 2 \end{bmatrix}$$

$$= \begin{bmatrix} 8 \\ 9 \end{bmatrix} + \cancel{\begin{bmatrix} 2 & -2 \\ -1 & 1 \end{bmatrix}} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 2 \end{bmatrix}$$

$$h_2 = \begin{bmatrix} 8 \\ 5 \end{bmatrix}$$

$$Vh_2 = \begin{bmatrix} 1 & 0 \\ -1 & 1 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} 7 \\ 5 \end{bmatrix}$$

$$= \begin{bmatrix} 7 \\ -7+5 \\ 5 \end{bmatrix}$$

$$Vh_2 = \begin{bmatrix} 7 \\ -2 \\ 5 \end{bmatrix}$$

$$\text{Sigmoid} = \frac{e^{-7}}{e^{-7} + e^{-2} + e^{25}}$$

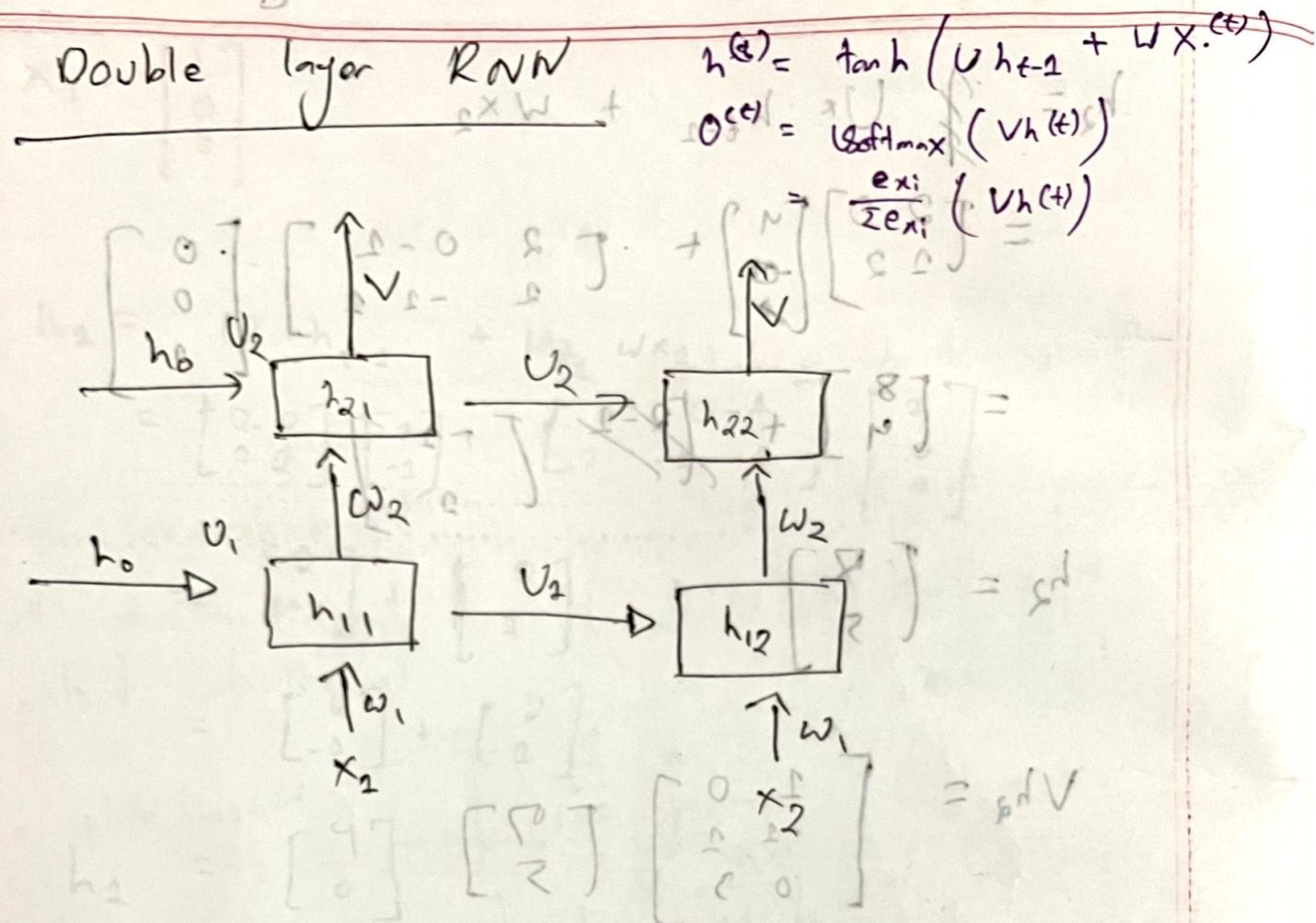
$$= 3.35 \times 10^{-1} \quad o^{(2)} \Rightarrow$$

$$\frac{e^{25}}{e^{-7} + e^{-2} + e^{25}} = 4.14 \times 10^{-3}$$

Input gate:

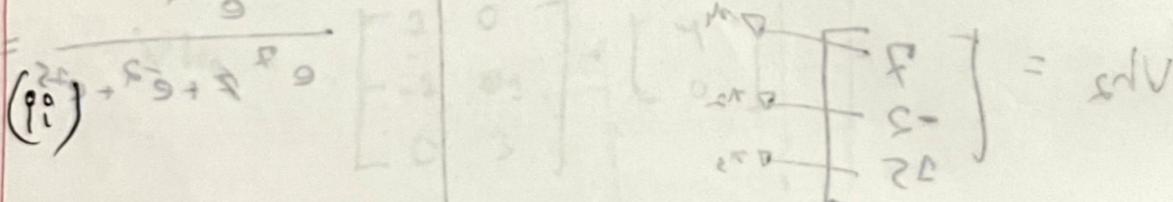
ind=0	3.35×10^{-1}
ind=1	$0.996 \leftarrow$ max.
ind=2	4.14×10^{-3}

No. of hidden layer depends on the size of input.



Challenges:-

(i) It is unidirectional $\xrightarrow{\text{solve}} \text{make it bi directional}$



(ii) Problem of vanishing gradient $\xrightarrow{\text{value}} \text{LSOM;ompic}$

$$r = 0.1 \times 28.8 \quad \text{obit} = 2 \quad \text{bit} = 1 \quad b = 4$$

Drawbacks:-

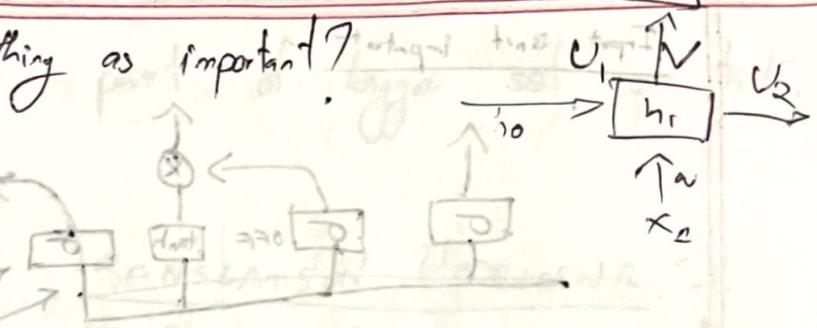
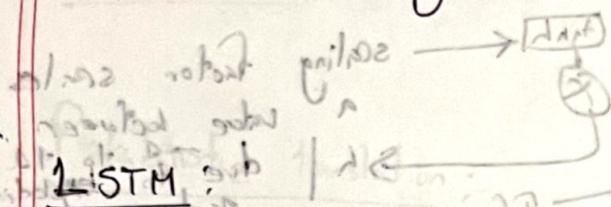
\rightarrow Only see the past, not the future

\rightarrow Must forget the same amount of history at each time step

↳ Two inputs are step input \rightarrow step input

23/01/2024

How to identify something as important?



LSTM:

- New three gates are introduced
- element wise operations

⇒ If the information is important enough we will keep it in the persistent memory.

The shape has to be exactly the same to perform these operations

Forget gate:

If ^{any} info is important, the forget gate free up some memory. To clear up the memory.

If it is very important info if it is needed it will fully empty up the memory.

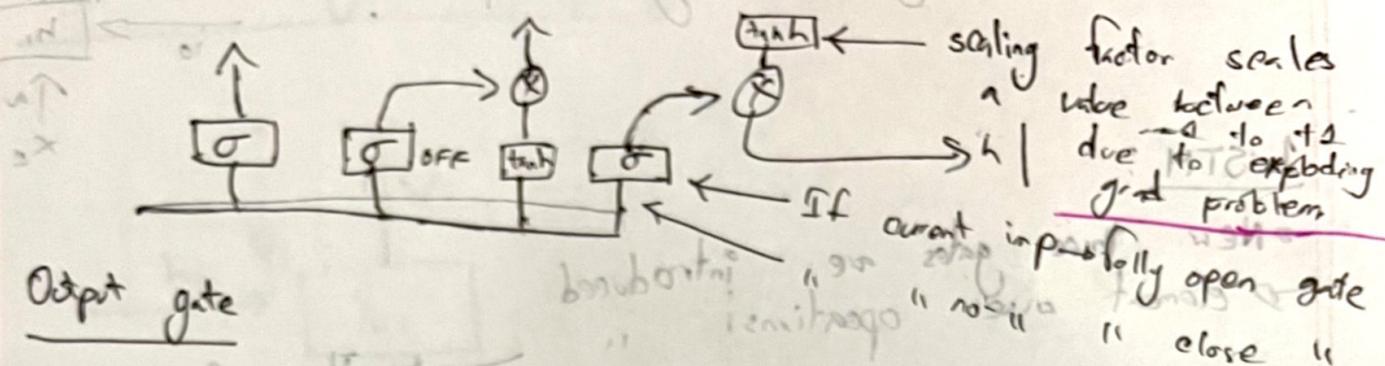
Mathematical function that works like a gate.

Input gate:-

- how much we ^{will} open or close the gate
- " " " " pass the data.

target gate & input gate are reciprocal.

Input isn't input



Output gate

- (i) There are no two channels
- (ii) There is a scaling factor

exists

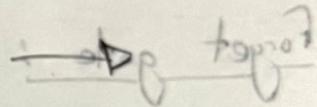
Gated RNN problem \rightarrow $\frac{\partial L}{\partial h_t}$ \rightarrow $\frac{\partial L}{\partial h_{t-1}}$ \rightarrow \dots

can't solve

- (i) Back propagation through time
- (ii) RNN's are not parallelizable.

Redundancy of LSTM

LSTM \rightarrow forget g. & input g.



stop fugue

stop \rightarrow $\frac{\partial L}{\partial h_t}$ \rightarrow $\frac{\partial L}{\partial h_{t-1}}$ \rightarrow \dots

" "

" "

" "

TRANSLATION :- Is a part of bigger set of task.

Probabilistic translation :-

Language Modelling :-

IBM MODEL

Solutions :-

(i) Factor regression

IBM Model

Expectation maximization algol algor. (EM)

Example

(i) ~~b c~~ b c

(ii)

Step-1

$$P_1(x|b) = \frac{1}{2}$$

$$P_2(y|b) = \frac{1}{2}$$

$$P_3(z|c) = \frac{1}{2}$$

TRANSLATION MODELLING :-

checks if correctly identifying word level translation

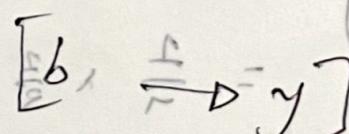
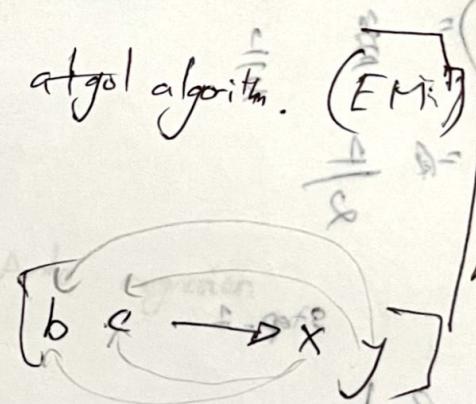
Disadvantage :-

(i) factorial complexity

(ii) extreme inefficient process

(iii) Ambiguity word handling

Dictionary set



$$P(z|x,y) = \frac{1}{2}$$

Because specific words meaning may change

$$\begin{array}{c} \overbrace{a} \\ \overbrace{b} \\ \overbrace{c} \end{array}$$

6 alignments
permutation.

Step-2

~~Step-2~~
Normalizing

$$P\left(\begin{array}{c} b \\ \text{x} \\ \text{y} \end{array} \mid \text{odd}\right) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} = \frac{1}{2} \times 2 = \frac{1}{2}$$

$$P\left(\begin{array}{c} b \\ \text{x} \\ \text{y} \end{array} \mid \text{even}\right) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} = \frac{1}{2} \times 2 = \frac{1}{2}$$

$$P(b) = P\left(\begin{array}{c} b \\ \text{x} \\ \text{y} \end{array}\right) = \frac{1}{2} \times 2 = 1$$

Step-3

- spanning

$$P(x|b) = \frac{1}{2} \xrightarrow{\text{Normalizing}} \frac{1}{3}$$

(because weights are normalized)

$$P(y|b) = \frac{1}{2} + \frac{1}{2} = \frac{1}{2}$$

(probability exceeds 0.5)

$$P(x|c) = \frac{1}{2}$$

(each word)

Iteration \Rightarrow

$$P\left(\begin{array}{c} b \\ \text{x} \\ \text{y} \end{array} \mid \text{odd}\right) = \frac{1}{2} \times \frac{1}{2} \xrightarrow{\text{divide by 3}} \frac{1}{8} \xrightarrow{\text{normalize}} \frac{1-2/3}{9} = \frac{1}{9}$$

$$P\left(\begin{array}{c} b \\ \text{x} \\ \text{y} \end{array} \mid \text{even}\right) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} = \frac{1}{2} \times 2 = \frac{1}{2}$$

$$P\left(\begin{array}{c} b \\ \text{x} \\ \text{y} \end{array} \mid \text{odd}\right) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} = \frac{1}{2} \times 2 = \frac{1}{2}$$

Step-2

$$P(X|b) = \frac{2}{7}$$

$$P(Y|b) = \frac{2}{8} = \frac{1}{4}$$

$$P(x|e) = \frac{3}{7}$$

$$P(y|e) = \frac{1}{5}$$

Challenge:- leads

- (i) If translates to multiple words, so it won't be able to able to translate. Can't translate more than equal words
- (ii)

→ will work in low resource language.

Solution to work better :-

→ Religious

books [no dictionary present for low resource language]

LM models idea is Auto regression

Probabilistic translation model → Neural machine translation

LSTM

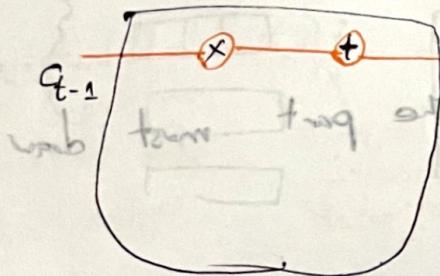
[Long & Short Term Memory Networks]

A technique to know
remember which we should
remember which we should
forget and
forget

→ Designed to avoid Long term dependency problem.

Characteristics :- Remembering information for long period of time

LSTM STRUCTURE :-

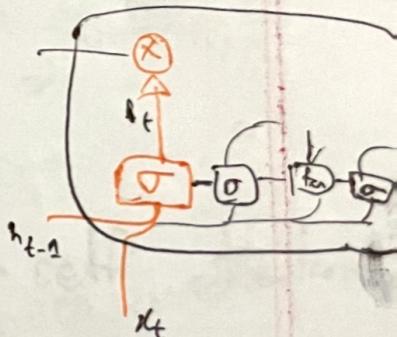


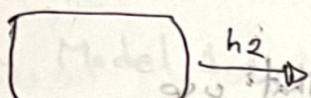
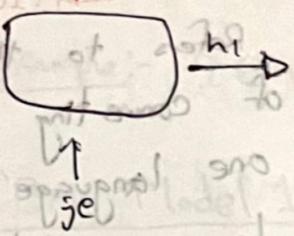
$C_t \Rightarrow$ cell state → kind of conveyor belt.
have the ability to remove or add info to them cell state.

Forget gate layer [sigmoid layer] :-

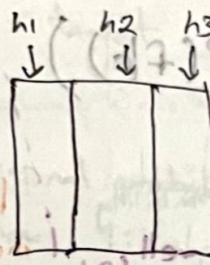
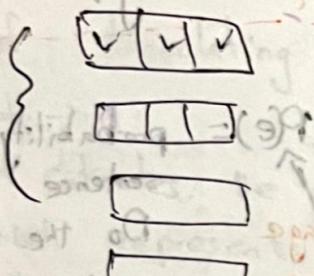
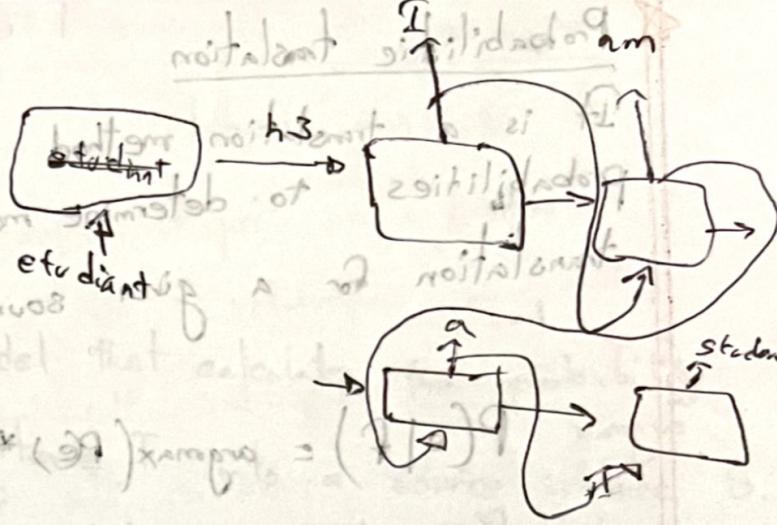
→ what info. we are going to throw away from the cell state.

→ output a number between 0 & 1





suiv



Attention matrix M

TRANSFORMER

① will look at all possible combinations. \rightarrow self attention.

(r) no (T/2)9 utilisation limitations of estimators + P
no sometimes forget something new & how we print things
no sometimes forget something new & how we print things
no sometimes forget something new & how we print things

TRANSLATION

380
350
250

395 → 750

★ Probabilistic translation

It is a translation method that uses probabilities to determine the most likely translation for a given source text.

$$\text{argmax } P(e|f) = \text{argmax}(P(e) * P(f|e))$$

$P(e)$ = language model

$P(f|e)$ = Translation model probability

LM measures how well the source sentence f aligns with the target sentence e .

language modelling estimates it

★ Language Modelling:

It involves creating models that assign probabilities to sequences of words in a language.

A language model estimates $P(e)$.

If e is the probability that a sentence e is an English sentence.

TRANSLATION MODELLING:

It is a probabilistic framework that defines the relationship between a source language S and a target language T .

It estimates the conditional probability $P(S|T)$ or $P(T|S)$, representing how well a given source sentence corresponds to a target sentence.

TRANSLATION:

Refers to the process of converting text from one language to another language while preserving meaning.

NMT

$P(e)$ = probability of English sentence

Do the words follow English order?

In short:-

LM estimates probability of triplets sequence of words.

techniques are:-

(i) n-gram LM

(ii) RNN's

To note a sentence is a sequence of words.

$$P(T|S) = P(S|T) P(T)$$

T = target sentence
S = source "

$P(S|T)$ → measures how well the source sentence aligns with the target sentence

$P(T)$ → language model, measures

example:- IBM Model 1 input is I bora paa A
output is X gaa kaa Kaa no n baa

IBM Model 1: probabilistic

It is a machine translation model that calculates the probability of translating a target sentence T into a source sentence S.

Goals →

(i) It estimates the conditional probability $P(S|T)$, which represents how likely a source sentence S is as a translation of a target sentence T.

→ consider all possible word alignments

Combine the word probabilities

(ii) Align each source word to a target word

$$P(f|e) = \sum_{(e_i, f_i, e_a)} \prod P(f_i | e_i)$$

$P(f_i | e_i)$ = for each English word in the probability of it being translated as French

prob not when target word is not found

of bad no of at part starte with

and target words in output of last year

Attention allowed decoder to access all hidden states of encoder to teach it? $T/29 = (8/1)^9$
Info loss \rightarrow addressed by attention mechanism $(T/29)$

Sequence to Sequence Model

A seq2seq model is a type of neural architecture map? designed to map an input sequence X to an output sequence Y .

two main components:-

Encoder:
→ Encodes the input sequence into a fixed-length representation.

→ Encoder processes the input sequence word by word and encodes it into a fixed dimensional context vector.
→ A context vector that summarizes the entire input sequence.

→ At the end of the input sentence, the encoder produces a context vector that captures the sentence's meaning.

Decoder: Decodes the encoded representation into the output sequence.

→ The decoder generates the output sequence one word at a time using the summarized context vector.

Problem:

The encoder compresses the entire input sequence into a fixed length context vector. For long input sequences, this single may lead to information loss / may fail to fully capture all the important info.

Solve:

At Attention mechanism.

Hmm

Review.

$$t_i = \arg \max_{t_i^n} (P(t_i^n | w_i^n)) \rightarrow T_i^n \text{ large} \rightarrow \text{large} / \text{large}$$

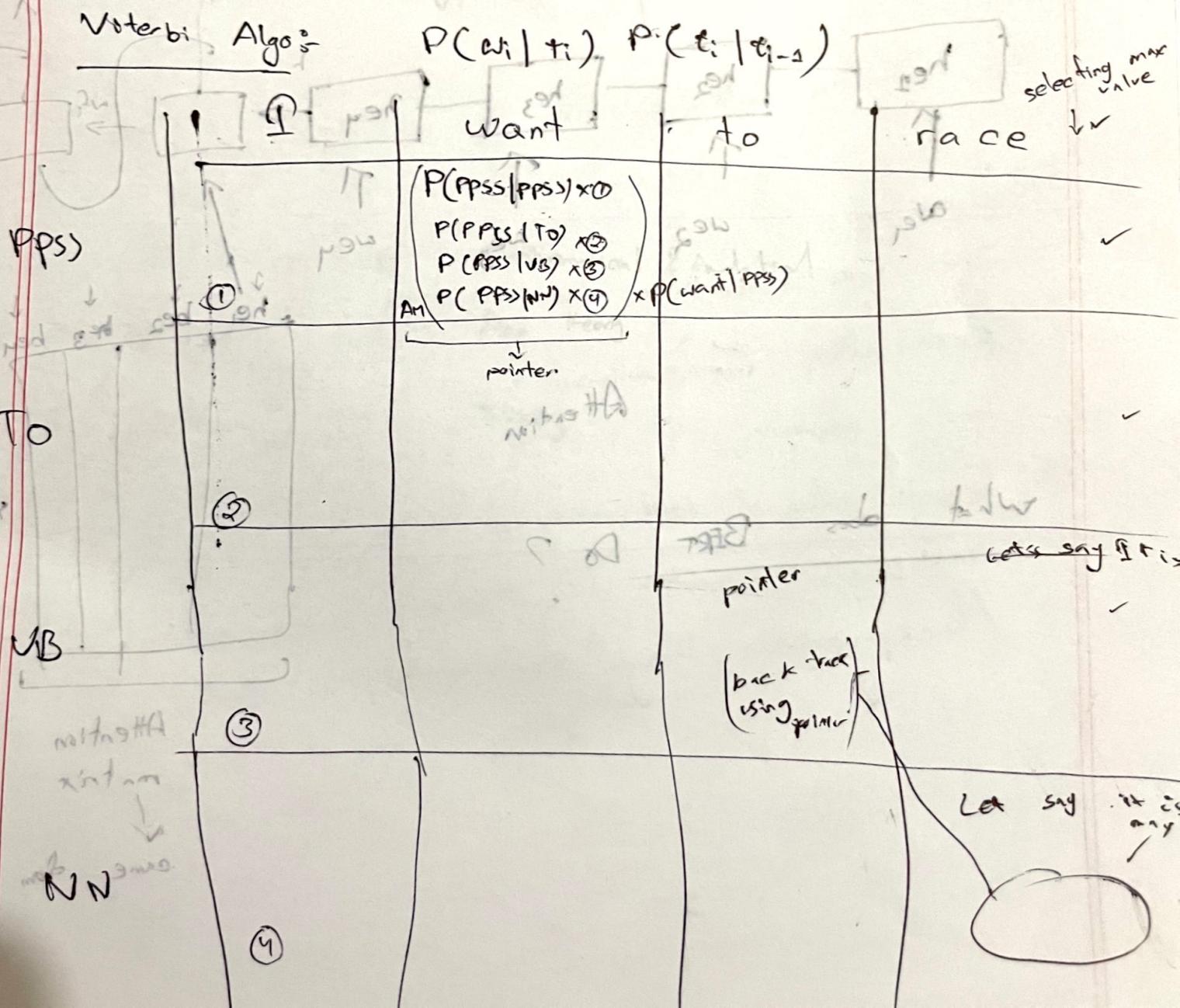
$P(w_i^n | t_i^n) P(t_i^n)$

$P(w_i | t_i) P(t_i | t_{i-1})$

tails down to

viterbi.

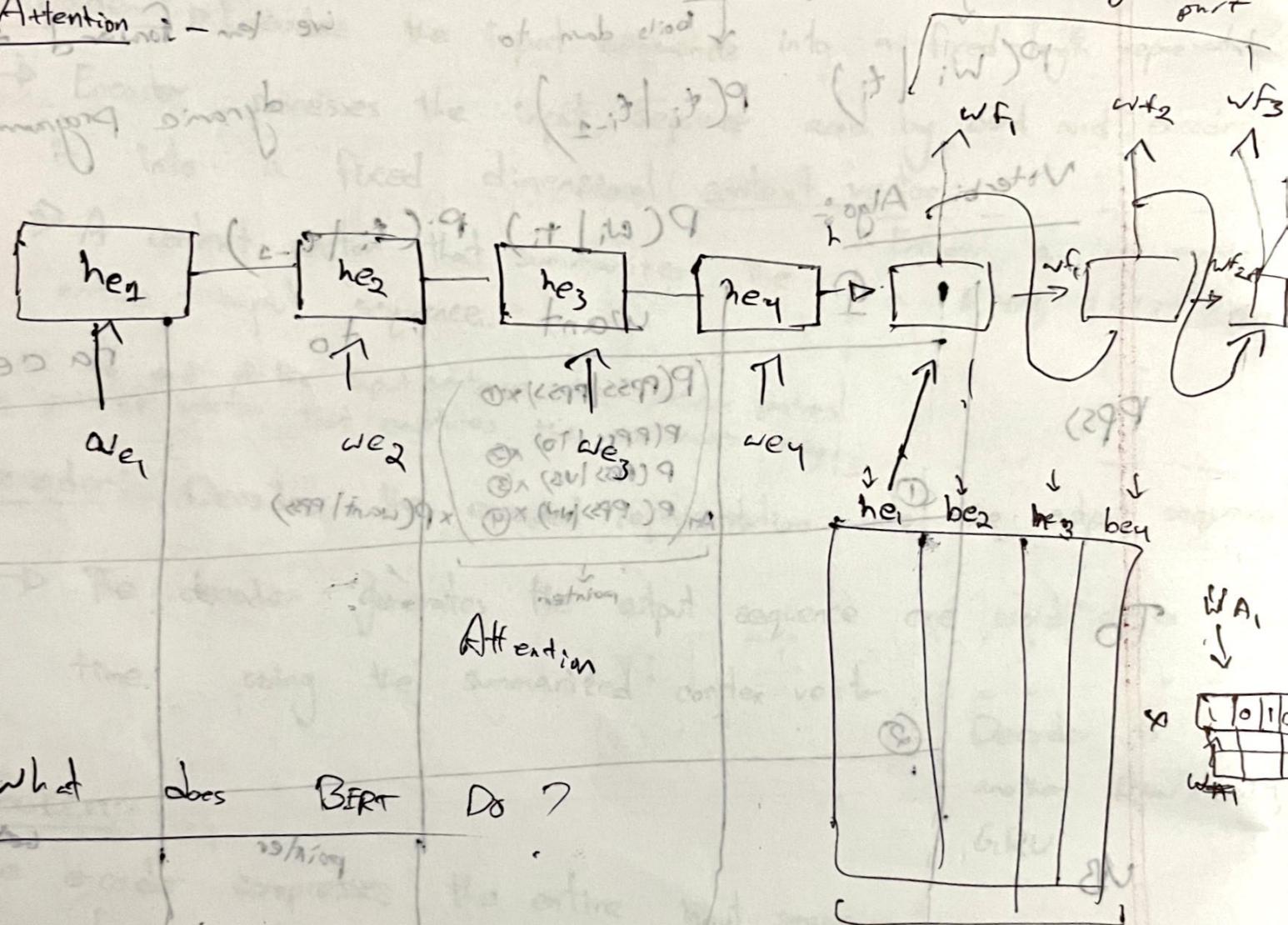
we can forward to dynamic programming



Difference between HMM & MEMM.

LSTM

Attention



What does BERT Do?

② Attention matrix

came from encoder.

Seq2Seq using Attention:- Instead of passing one single context vector.

- Attention allows the decoder to access all hidden states of the encoder, rather than relying on single fixed length context vector.
- At each decoding step, the model computes a weighted sum of the encoder's hidden states, and focuses on most relevant part of all the hidden states is passed in a attention matrix to the decoder.

The information bottleneck for long sentences solved

PROBLEMS:-

- (i) RNN works sequentially, the sequential dependency prevents parallelization. (results slower training)
- (ii) RNN suffers from vanishing gradient issues when processing long sequences.
- (iii) Sequential nature of RNN makes it difficult to leverage modern GPU.
- (iv) RNN struggle to retain info from very distant word in long sequence / Backpropagation through time is still an issue.

Solve:-

- to solve all these questions we have to understand the mechanism of transformer.
- Transformers use self attention to process all input words in parallel.
 - Using attention mechanism and feed forward layers, transformer avoid recurrence to eliminate vanishing gradient problem.
 - By processing all words simultaneously, transformers can take full advantage of GPU.
 - self attention calculates relationship between all words in sentence, regardless of distance.

Why need GPU:-

- It is very good at matrix multiplication.
- More computing units \rightarrow better thread parallelization.
- Matrix multiplication requires more computational operations, GPU's provide more memory for computation operations.
- have faster access to RAM.

TRANSFORMERS

Transformers removes the reliance on Recurrent state structures and uses only attention mechanism to process sequences. Introduced self attention as its core operation.

→ Transformers have \rightarrow no. of encoders \rightarrow no. of decoders
→ no. of layers \rightarrow long size cuz \rightarrow calculates the effect of one word on another word multiple times

Characteristics:

- (i) Size
- (ii) Self attention \rightarrow calculate the effect of every word on all of its surrounding words.
- (iii) Positional encoding \rightarrow doesn't work sequentially \rightarrow so how to capture sequence?
 - \hookrightarrow provides the position of the word
 - \hookrightarrow using this position we find the connection/relationship between every words.
- (iv) Byte pair encoding \rightarrow instead of limit the vocabulary size transformer use byte pair encoding. \rightarrow we use this instead of word tokenization
 - \hookrightarrow reduces the need for large vocabulary
 - \hookrightarrow handles rare and unseen words effectively
- (v) Multi-headed attention. \rightarrow each attention head learns different relationships, parallel. \rightarrow Calculates attention from different angles

* self attention allows each word to understand its relationship to every other word.

TRAINING A TRANSFORMER:-

- Not enough hardware
- Not enough time. / long training times
- Large computational requirements

Training two parts :-
Pretraining :- those who have enough resources, and time train a transformer to learn the language representation

Fine tuning :- to test try to train the pretrained transformer about different tasks.

BERT :-

It is a pre-trained transformer based, developed by Google. Designed for various NLP tasks.

Characteristics:-

(i) Bidirectional context understanding:- BERT processes text in both direction simultaneously.

(ii) Pretraining & Fine tuning:- BERT is pretrained. can be fine tuned on specific NLP tasks.

1. SEMI-SUPERVISED :- training a model to use patterns of language.

- Training on large amounts of text.
- Trained on a certain task enables it to grasp patterns in language.

2. Supervised :-

- Training on a specific task with a labeled dataset.

Pretrained on 2 unsupervised tasks

(i) Masks Masked LM → Replaces a token with a [MASK]

is a task where 80% of the time, with a random token portion of the input tokens is randomly masked 10% of the time and does not replace

(ii) Next sentence prediction → Learns to understand relations between pair of sentences

LSTM

Introduces a memory cell and other gates to selectively retain, forget or update information.

Addresses :-

1. Vanishing gradient problem.
2. RNNs struggle to retain info from long earlier time steps of long sequences. [Difficulty in capturing Long Term Dep]
3. Lack of effective memory mechanism.

PROBLEMS :-

1. Computationally expensive
2. Processes sequentially, can't be parallelized.
3. Vanishing gradient problem.

GRU

Simplified version of LSTM. merges Forget g. & input g. into a single gate. faster to train.

Addresses - [LSTM have multiple gate]

- (i) Simplifies gates by using only update & reset gate.
(ii) Much better at using long distance info
Add (iii) faster training.

PROBLEMS

- (i) Can't solve back propagation through time
- (ii) Not parallelizable.