

NLP \rightarrow Introduction

DATE

\Rightarrow NLP \Rightarrow Natural Language Processing

NLP : All the preprocessing done on raw text before NLU is NLP

Natural language \Rightarrow Formal language $\xrightarrow{\text{FT is language of Linear Algebra}}$

FL \Rightarrow First Order Predicate Logic

Compositional Semantics: $NL \xrightarrow{T(.)} FL$ This process is goal in case of
(Worried abt meaning) $\quad \quad \quad \xrightarrow{T(.)}$ (compositional Semantics)

Semantics should have well formed formula.

or

Distributional Semantics : $NL_1 \xrightarrow{T(.)} FL_1$, if $NL_1 \sqsubset NL_2 \Rightarrow f_{sim}(NL_1, NL_2) = Y$
(Worried about similarity scores) $NL_2 \xrightarrow{T(.)} FL_2$ then $FL_1 \sqsubset FL_2 \Rightarrow f_{sim}(FL_1, FL_2) = Y$

FL is something that can be easily interpreted by the computer / BOT

How come we can do $f_{sim}(NL_1, NL_2) = f_{sim}(FL_1, FL_2) = Y$

NL can't be understood by computer - So it is done by human touch but it is highly subjective. To over come this we will keep large no of people as subject for computing $f_{sim}(NL_1, NL_2)$.

One way is take avg of scores of subjects.

In case of Distributional Semantics the $T(.)$ can be loose transformation or fuzzy transform of NL. Here meaning capturing is not important everything should be consistent whether it is mistake or correct take.

\Rightarrow Language \Rightarrow Syntax + Semantics + pragmatic

(Grammar) (Literal Meaning) (Intended meaning)

Homonymy \rightarrow Some word different meaning (Tiger, Ban)

Polysemy \rightarrow Some word slight difference in meaning (root meaning is same)

classmate (Bank)

PAGE

Phrase :- Sequence of words that have collectively a meaning in dictionary.

DATE [] [] []

=> D

1m

For pragmatics, we need more context to understand the intended meaning.

Context Window : The size of window to understand context of the sentence.

C=20 mean t=20 from current sentences

Form : S V O
Subject Verb Object

Syntax

↓ ↓ ↓ ↓
Form Part of speech "Sequence"
(SVO) (Pos) Parse tree
Dependency
(context free) Parsing
(CP)

is it valid or not? \Rightarrow There are some sentences which are grammatically correct but has no meaning.

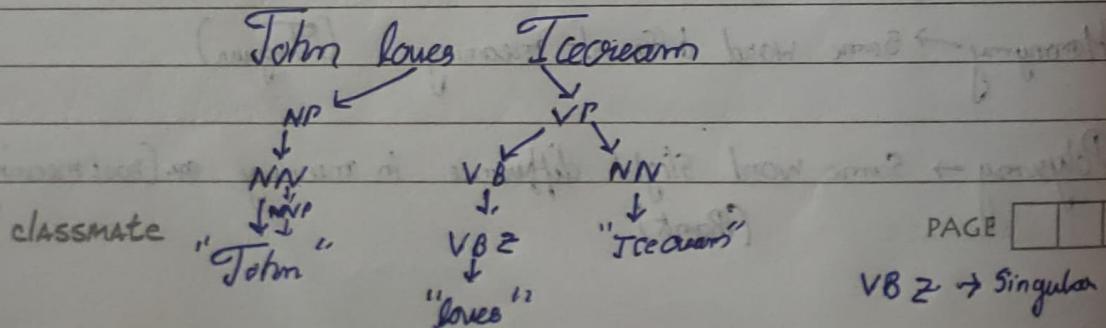
e.g. colorless dreams sleep furiously.

Here the Compositional syntacticians transform the sentences & get back to original sentence after considering proper grammar parsing but it is still of no use.

Any sentence is composed of : Noun phrase & Verb phrase

$|V| = \sum \text{lexicons}$

[Cardinality of Vocabulary]



PAGE [] []

VBZ \rightarrow Singular Verb

If parse tree goes wrong then everything goes wrong

→ When there is usage of any kind of slang & it maintains the same structure as in training data then it will consider similar tag for the unknown phrase.

⇒ In dependency parsing the verb is made the pivot & the rest of the things are figured out

If we can assign truth value to a sentence then we can understand that it is meaningful or not even though it is grammatically correct.

What is most important is A word with no semantic role is useless
A verb with no semantic role is also useless
A verb with a semantic role is useful

semantic role with verb form no auxiliary & verb &

Verb with full semantic role & A descriptive of anything else
or verb related to action form (whose) situation (condition etc.)
verb present or verb past form no auxiliary under all
verb like see eat hit play etc.

verb related to situation description

verb related to time temporal

verb as "time temporal" or temporal

verb related to time temporal
already, happening, happened

initial, final

⇒ If two things are similar they should share some similar properties/characteristics.

$[w_{j-2} \rightarrow w_{j-1} \rightarrow \text{house} \rightarrow w_{j+1} \rightarrow w_{j+2}]$ } might happen

$[w_{i-2} \rightarrow w_{i-1} \rightarrow \text{rent} \rightarrow w_{i+1} \rightarrow w_{i+2}]$

Similar frequency distribution

- If there is overlapping of terms to great extent then it can be claimed that two terms are similar otherwise we will call it related.

- Some places ↳ Whenever the term A is occurring term B should occur
- Some friends ↳ Words in context window A & B are same as of B
- Some frequency ↳ The histogram of frequency of A & B are same

If these 3 conditions are met then they are similar.

- The problem is even though A & B are similar but the model (Distributional Semantic Model) won't consider it similar due to the above 3 conditions are not fulfilled in training data. it will tell A & B are related.

→ Now onwards similarity is relatedness.

⇒ Vector Space Modelling

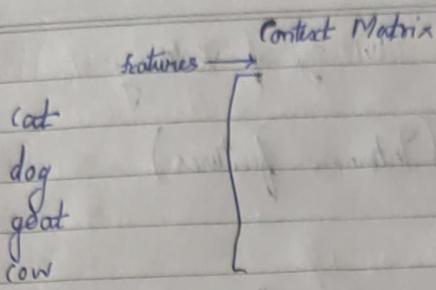
- Represent a "linguistic unit" as a vector.

Linguistic unit : character, word, phrase, clause, sentences, short paragraph, long paragraph, document.

Context Matrix :

name of the matrix is set on basis of columns.

DATE



The features / contact comes from contact window.

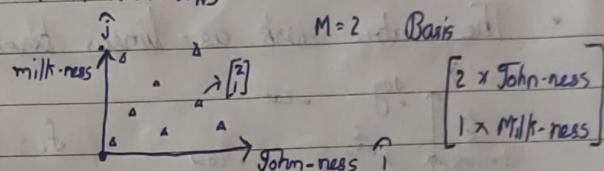
- Features from basis.

Given a set of features & set of datapoints
(f_1, f_2, \dots, f_M) $\{w_1, w_2, \dots, w_N\}$

f_1 : John

f_2 : milk

↳ denotes the words associated with each one of datapoints.



[Initially]: John = Johness

Contact - Matrix

	John	milk
cat	2	1
dog	3	0
cow	1	4

John is in 3D

Cat is in 2D

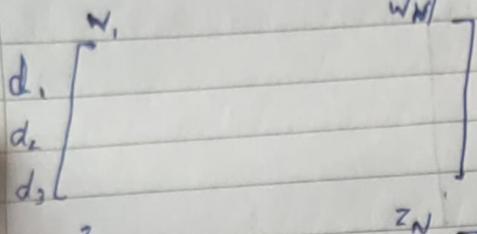
$$\text{cat} = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \times 2 = 2 \begin{bmatrix} 2 \\ 3 \\ 1 \end{bmatrix} \times 1 \begin{bmatrix} 1 \\ 0 \\ 4 \end{bmatrix} = \begin{bmatrix} 5 \\ 6 \\ 6 \end{bmatrix} \quad \text{Now cat in 3D}$$

The features you take due to co-occurrence then it is an extrinsic features

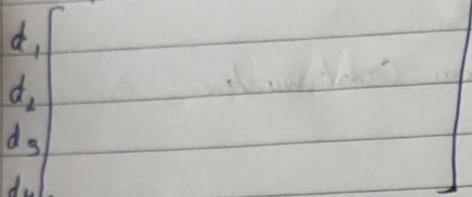
When we use representation of contents inside them are ~~are~~ internal features.

W size of Vocabulary

(Internal features)



z_N

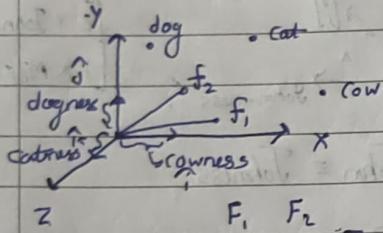


$z_i = i^{\text{th}}$ zip file

(External features)

When modelling of documents is done based on external features most of times it is not useful.

- The matrix work as linear transformation



$$f_1 = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$f_2 = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

$$CM = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \quad \text{if } \text{not } CM \neq \text{Cat} \Rightarrow \text{Cat}'$$

$$\begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \times \begin{bmatrix} x \\ y \end{bmatrix} = x \begin{bmatrix} F_1 \\ F_2 \end{bmatrix} + y \begin{bmatrix} F_1 \\ F_2 \end{bmatrix}$$

$$3 \times 2 : s - 2 \times 1 = (3 \times 1)$$

- If you are not able to get better performance then the culprit is the Content Matrix as it can't transform things efficiently.

- Therefore, the content matrix is needed to be learnt.

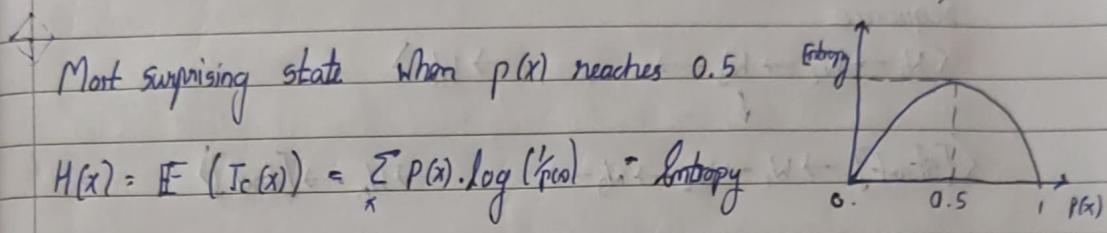
- You can randomly initialize CM (May be values from gaussian dist) then start training & check formation.

• $TF\text{-}idf \rightarrow$ Term Frequency - inverse document frequency.

If we know that a particular feature will surprise you / or help you to classify similar & non-similar content.

Q How to compute information content?

$$\Rightarrow I_c(f_i) = \log\left(\frac{1}{p(f_i)}\right) = I_c(x) = \log\left(\frac{1}{p(x)}\right)$$



$$H(x) = \mathbb{E}(I_c(x)) = \sum_x p(x) \cdot \log\left(\frac{1}{p(x)}\right) = \text{Entropy}$$

• $idf = \frac{1}{\frac{\text{num of docs}}{N}} = \frac{N}{d_i \cdot \ln N_{\text{docs}}} \propto \frac{1}{p(f_i)}$ now if we take \log it is similar.

• TF : We get I_c of (f_i) but we want expected value for f_i .
 (Surprise value of F_i , knowing we have cat) $= \frac{\sum_i w_i \cdot \text{cat } E_i}{\sum_i w_i E_i}$

⇒ Initializing a context matrix

$$cm = \begin{bmatrix} f_1 \\ \vdots \\ f_m \end{bmatrix} \quad \begin{bmatrix} \# \\ \vdots \\ \# \end{bmatrix}$$

Calculation of each cell is independent of row & column
 (i.e. assuming all the features are independent of each other)
 [Orthogonal] \downarrow (BOW = Bag of Words)

But it is false as the language models depends on the assumption that words are dependent so are the features which will help in predicting that the missing words.

John loves \leftarrow (correlation) \rightarrow Prediction.

There is a flow in the language otherwise those modelling techniques won't exist.

We can initizlize CM with help of TF/IDF it is better than random. If you are doing so then you are adding inductive bias.

The main contributor behind this is idf.

⇒ Language Model

Sequence

$$\langle w_1 \rightarrow w_2 \rightarrow w_3 \rightarrow \dots \rightarrow w_n \rangle$$

↓

$$\langle [\text{start}] \rightarrow w_1 \rightarrow w_2 \rightarrow w_3 \rightarrow \dots \rightarrow w_n \rightarrow [\text{End}] \rangle$$

$\hat{p}(w_i \rightarrow w_n) = ?? \rightarrow$ This estimation is what language model does

$\hat{p}(w_n | w_1 \rightarrow w_{n-1}) = ??$ (Yes can be computed)

Next word prediction

$\hat{p}(w_i \rightarrow w_2 | w_1 \rightarrow w_{i-1}) = ??$ (Yes)

Next Sequence prediction

Note : $p(w_1 \rightarrow w_2 \rightarrow w_3) \neq p(w_1, w_2, w_3)$

Joint distribution ⇒ Chain Rule

$$p(w_1) p(w_2 | w_1) p(w_3 | w_2) = p(w_1) p(w_2 | w_1) p(w_3 | w_2) = p(w_3) p(w_2 | w_3) p(w_1 | w_2)$$

Therefore, in RHS Order doesn't matter

$$P(w_n | \leftarrow w_{n-1}, \leftarrow w_{n-2} \leftarrow \dots, \leftarrow w_1, \leftarrow \text{Start})$$

N-gram

Modelling But if we consider Markov's chain Rule & use $P(w_n | \leftarrow w_{n-1})$ then this is bi-gram model.

We also have to compute $p(\text{End} | w_n)$. It is needed to be calculated because we need to ensure that the sentence is ending there and there is no extension to it.

⇒ Why not bigram? It can't sustain longer dependencies.

Eg John, who lives in NY and works at Apple, loves ice-creams

(Bi-gram) \Rightarrow

(Tri-gram) \Rightarrow

still context is not understood.

But if we consider large degree for consideration then there is no meaning to consider Markov's rule.

$$P(JC | \leftarrow \text{loves}) = \frac{c(\text{loves} \rightarrow JC)}{c(\text{loves})}$$

This is something we are interested in finding

⇒ Bag of Words model can be used when we don't have to consider sequence

⇒ Skip-gram Modelling (Word2Vec)

If we have to estimate $\hat{P}(w_1 \rightarrow w_2 \rightarrow \dots \rightarrow w_n)$.

In skip gram modelling it says if we have to figure out $w_i < w_i < w_n$ then we have to consider sequence $w_i \rightarrow \dots \rightarrow w_{i-1}$ as well as $w_{i+1} \rightarrow \dots \rightarrow w_n$.

Problem with n-gram: John → loves → Ice-cream

$$P(JC | \leftarrow \text{loves} \leftarrow \text{John}) = \frac{c(\text{John} \rightarrow \text{loves} \rightarrow \text{Ice-cream})}{c(\text{John} \rightarrow \text{loves})}$$

- There might be few John → loves in the data so denominator problem
- There are some John → loves → Ice-cream counts so the prob decrease.
- In some cases we can say directly from John that ice-cream should come but due to n-gram we are penalizing it.

$$P(w^c / w^t)$$

$w^t \Rightarrow$ target word

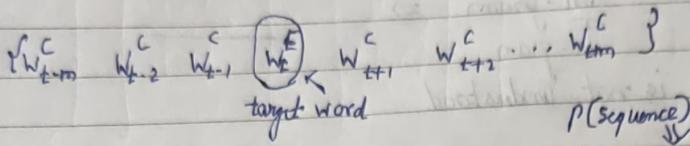
w^c : content word

$$w^c = w^c_{j+1}$$

the content word is at j offset from target word.

In case i want to take bigger context window then.

$$\prod_{j=-m}^T P(w_{t+j}^c | w^t) \Rightarrow \text{Context window} = 2m$$



So we are trying to find probability of $P(w_{t-m}^c \rightarrow w_t^c \rightarrow w_{t+m}^c)$

Here we are taking assumption that they are conditionally independent otherwise if they are totally independent then it would be Bag of words

Conditional independence of all w^c 's

If you have multiple target words then

$$\prod_{t=1}^T \left[\prod_{j=-m}^{+m} P(w_{t+j}^c | w^t) \right] \rightarrow \text{Maximize this}$$

\Rightarrow Word2Vec takes target word as input & gives set of context words as output. While learning this it learns the weight values for target words [embeddings]

\rightarrow Word2Vec doesn't attempt LM

\rightarrow We can extend Word2Vec to do Language Model

→ Similarity function function (Word2Vec) : Dot product

$$w_i^t \cdot w_j^t = s_{ij} \quad \text{Similarity score}$$

Input hot-vector $[0, 0, 0, 0, 1, 0, 0]$ w_i^t

Desired output: $\{w_j^c\}$ → this will be found by calculating similarity with w_j^c

Now we have similarity matrix $\begin{bmatrix} \end{bmatrix}$ so we will convert it into $\begin{bmatrix} \frac{s_1}{\sum s_i} \\ \frac{s_2}{\sum s_i} \\ \vdots \\ \frac{s_n}{\sum s_i} \end{bmatrix}$ probability matrix

but there is a chance that one similarity can dominate rest we may get skewed matrix, so instead of this we can use softmax $\begin{bmatrix} \exp(s_1) \\ \exp(s_2) \\ \vdots \\ \exp(s_n) \end{bmatrix}$

$$w_i^{t'} = w_i^t * w_i^t = L(w_i^t) - ① \quad w_i^t \rightarrow \text{are embeddings in word2vec}$$

$$s_i^{t'} = (w_i^c)^T * w_i^{t'} = f_d(w_i^{t'}) - ②$$

$$\hat{y}^c = \hat{p}(w_j^c / w_i^{t'}) = \text{softmax}(s_i^{t'}) - ③ \quad w_i^c \rightarrow \text{Context matrix}$$

$$\begin{array}{c} \text{decision} \\ \text{Vector} \\ \hline \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_{i-1} \\ p_i \\ \vdots \\ p_n \end{bmatrix} \xrightarrow{\text{winess}} \begin{bmatrix} & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \end{bmatrix} \\ \hat{y}^c \quad y^c \end{array}$$

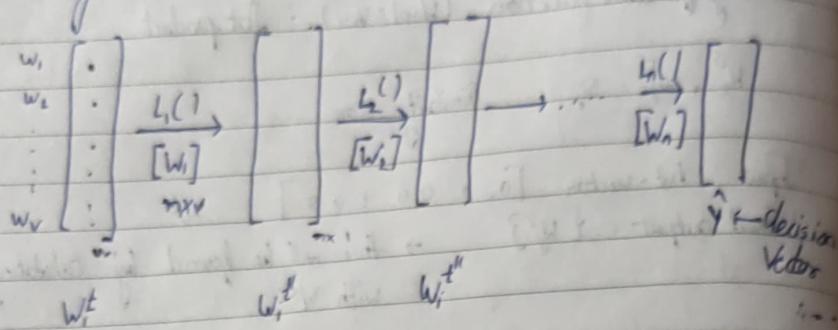
You know how many 1's should be there at which position in y^c .

Now we have two vectors so we will decide about loss. $\mathcal{L}(\hat{y}^c, y^c)$

$$[N_{\text{final}}^t] * w_i^t = w_i^t \quad \text{loss is cross entropy}$$

⇒ discuss
hyperbolic
functions
consider

→ Some goal with FFN



$$w_i^{t''} = w_i \star w_i^{t'} + \vec{b}_i$$

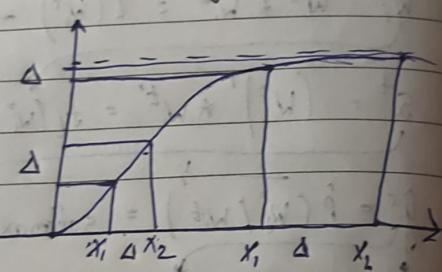
bias: safety parameter

⇒ What is happening at hidden layers

- By applying activation function we are squashing the vectors

$$w_i^{t''} = \text{fact}(w_i \star w_i^{t'} + \vec{b}_i)$$

$$\text{fact} \left(\begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \right) = \begin{bmatrix} f(1) \\ f(2) \\ f(3) \end{bmatrix}$$

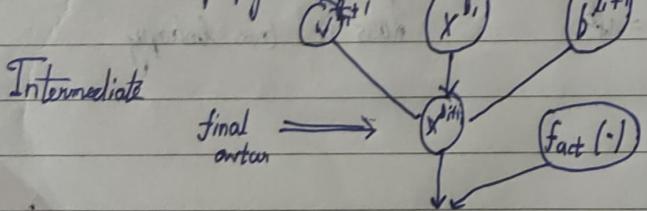


- it will squash more the unimportant features & squash less the important features
(this is why activations are important)

all 5 type of activation can take care of this exaggeration

Sigmoid, tanh & softmax

⇒ • Loss & Backpropagation



final output -
normalized or squared

final output

decision function

classmate

PAGE

- Q How to know about the role of each component in giving wrong decision?
 A So we will partially derivate the loss function w.r.t all components.

$$\frac{\partial \hat{L}}{\partial W^{l+1}} = \frac{\partial x^{l+1}}{\partial W^{l+1}} * \frac{\partial x_{norm}^{l+1}}{\partial x^{l+1}} * \frac{\partial \hat{x}}{\partial x_{norm}^{l+1}} * \frac{\partial \hat{L}}{\partial \hat{x}}$$

This is backpropagation of error

So now error correction or updation rule

$$(New) \quad W^{l+1} = W^{l+1} - \alpha \frac{\partial \hat{L}}{\partial W^{l+1}} \quad - (1) \Rightarrow \text{Stochastic Gradient Descent}$$

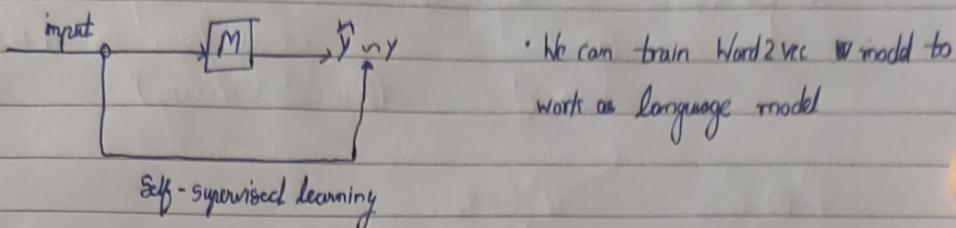
$$x_{norm}^{l_i} = x_{norm}^{l_i} - \alpha \frac{\partial \hat{L}}{\partial x_{norm}^{l_i}} \quad - (2)$$

$$b^{l+1} = b^{l+1} - \alpha \frac{\partial \hat{L}}{\partial b^{l+1}} \quad - (3)$$

Similarly

$$W^{l_i} \rightarrow \frac{\partial \hat{L}}{\partial W^{l_i}} = \frac{\partial x^{l_i}}{\partial W^{l_i}} * \frac{\partial x_{norm}^{l_i}}{\partial x^{l_i}} + \frac{\partial \hat{L}}{\partial x_{norm}^{l_i}}$$

⇒ How Word2Vec uses skipgram modelling?



* n-gram Language Model using RNN

RNN → 1990 Word2vec → 2013

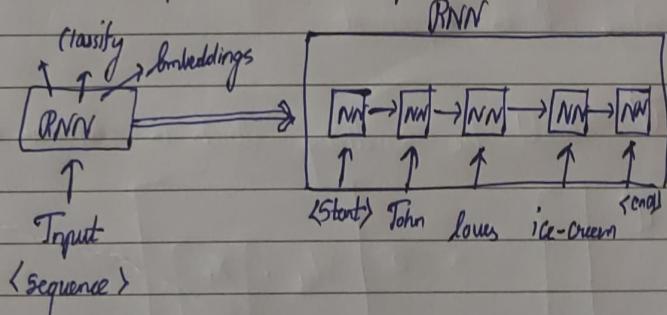
RNN is not a fixed or particular type of model. It considers sequential data.

We can create our version's of RNN. Mostly used RNN is LSTM

^(Recurrent)
R + Recurrent Neural Networks

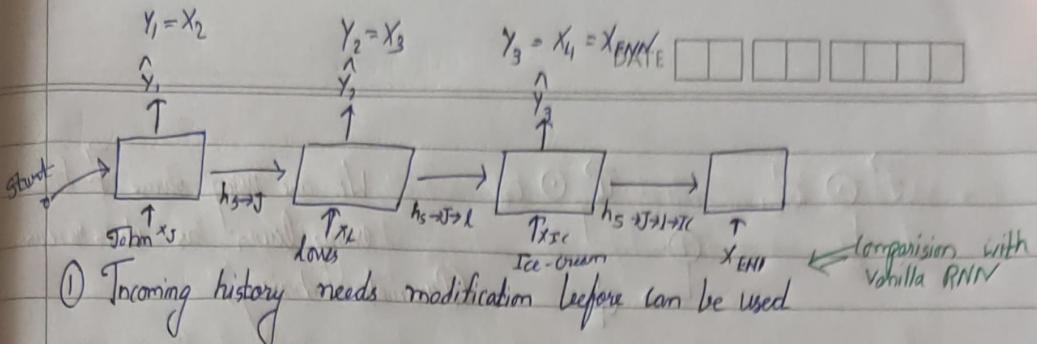
Word2Vec considers part of the sequence at a time as input. So next loss update is independent of previous loss update.

It assumes that the words follow IID (Independent & Identical Distribution)



There will be different embeddings for many in many loves John & John loves many in RNN while as in Word2Vec it will be same.

Structure architecture will be same throughout the sequence.



- ① Incoming history needs modification before can be used
- ② Basic structure remains same, knock off structures which can confuse the model. The past can be carried to certain point continuing further will confuse the model. This the job of forget gate.

Forget gate: It is a kind of vector, with which we are doing element wise multiplication of Incoming history.

$$\vec{F} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \odot \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

c = Content inside previous cell
 c' = " " modified cell

Forget gate Incoming history $\vec{h} \xrightarrow{h_t} \{c'\}$

(Human shadow Reduce has make h much neess. ti)

↳ We are forgetting certain neess & we are keeping the rest as it is

Therefore, \vec{F} should be like a reducing agent i.e. squashing function.

\vec{F} used in LSTM for squashing is (Sigmoid)

Q Why element wise multiplication? \Rightarrow Bcz we need to work selectively

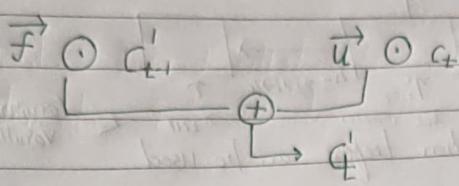
$$\Rightarrow \text{our } \vec{F} \text{ is sigmoid over something} = \sigma \left[w_f^t \cdot \begin{bmatrix} x_t \\ h_{t-1} \end{bmatrix} \right]$$

$$= \sigma \left[w_f^t * x_t + w_f^t * h_{t-1} + b_f \right]$$

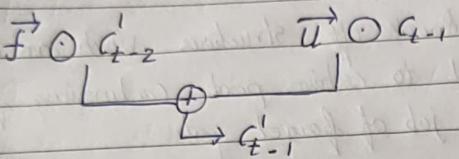
\vec{F} needs to be learnt

If i get \vec{F} , I will element wise multiplication with c'_{t-1}

for $c_t \rightarrow$ modified history at time t



We are using same f' for every layer as it is recommended
(Parameter sharing)



In RNN

$$h_t = \tanh \left((W_h * h_{t-1} + b_h) + (W_x * x_t + b_x) \right)$$

- h_t in RNN $\approx c_t$ in LSTM (Raw New content or new history)

But we are not happy with raw version, we want some part to be removed

the u' is important bcoz it helps in redundancy removal

To choose more $f' \odot c'_t$, Re through data whoi some choose $u' \odot c_t$ due to no redundancy to jayega. isslye u' is important.

u' should be something which is similar to

$$u' \begin{bmatrix} ? \end{bmatrix} = \sigma \begin{bmatrix} W_u \cdot \begin{bmatrix} x_t \\ h_{t-1} \end{bmatrix} \end{bmatrix}$$

- If our f' & u' are identity matrix vectors it will be special case of LSTM.

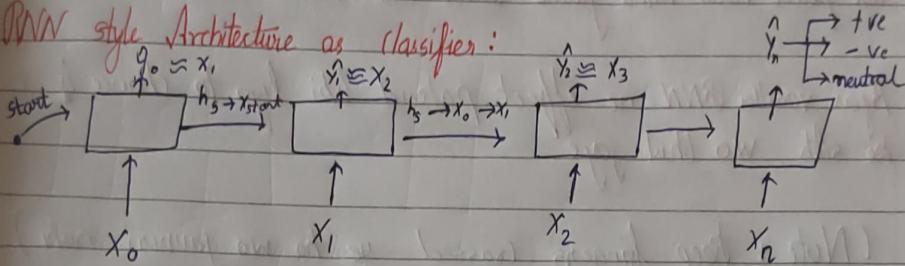
★ If f' is $\vec{0}$ & u' becomes $\vec{1}$ then LSTM will be Vanilla RNN

$\vec{\Theta} : \text{Output filter}$: It decides kitna shadow rakhna h ya kitna shadow chota kar dena h. Wo dekhenge ki noise dobara na aajaye.

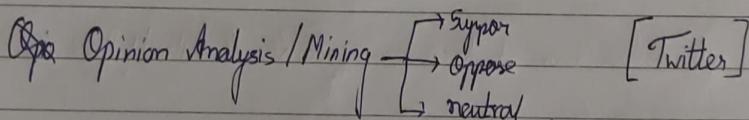
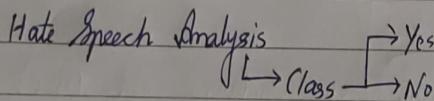
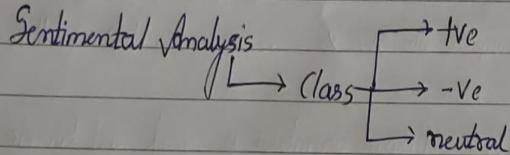
$$\vec{\Theta} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} = \sigma \begin{bmatrix} w_0 \\ x_t \\ h_{t-1} \end{bmatrix}$$

$$h_t'' = \tanh(\vec{\Theta} \odot \vec{c}_t) \quad h_t''' = \vec{\Theta} \odot \tanh(c_t)$$

\Rightarrow RNN style Architecture as classifier:



first example :



The degree of some class is a regression problem.

When doing classification the modifications are:

- Removing y_0, y_1, y_2 part. i.e. We are skipping $y = w_y * h_s \rightarrow x_0$ thik
- He just want the y in the last cell but we want it to predict classes instead of predicting the next word.

- We don't look at all elements but actually there are certain prominent elements which will define the outcome for the classification task

[Certain vectors ka zyada contribution hoga snowball banane me]

Training ke waqt non essential elements ka contribution karni hona hi chahiye
ye hame chahiye

- We prefer other RNN structure rather than Vanilla RNN bcz use patni nahi kitna jaune dena hai kitna nahi
- Agar vanilla RNN ne effect parakhi liya to change h time ke sath wo bhul jaye

(Note : Keep fingers crossed ye ho wo jawari nahi)

=> Conditional LM

$$\hat{P}(w_t | w_{t-1}, \dots, w_1, \text{Condition}) = ? \quad \text{Plain LM}$$

$$\hat{P}(w_t | w_{t-1}, \dots, w_1, \text{Condition}; C) = ?$$

Types of condition

i) should be present in context matrix
Window [Extrinsic]

(Don't be confused with context as condition)

(You can set condition that it should come before or after something)?

Eg: <John → loves → (TC); New York>

There will be cases the condition suppresses the probability?

Eg <John → loves → (TC) if *frankthiragar*>

less chance of occurrence of John with TC or with any other word.

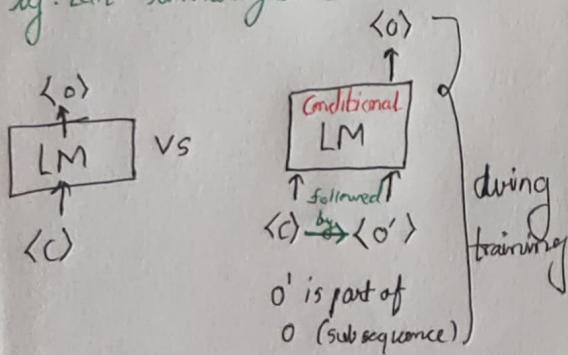
ii) The condition should occur somewhere inside the sequence [Intrinsic] ← *Pere*

• happens in certain generative LMs

iii) The condition is the input itself & the sequence is part of output & we have to complete it further

[Generating sequence of output given c input (prompt)]

Eg: text summarization



$\langle C \rangle \rightarrow \text{Question}$

$\langle 0 \rangle \rightarrow \text{A}$

$\langle d \rangle \rightarrow \text{subsequence of A}$

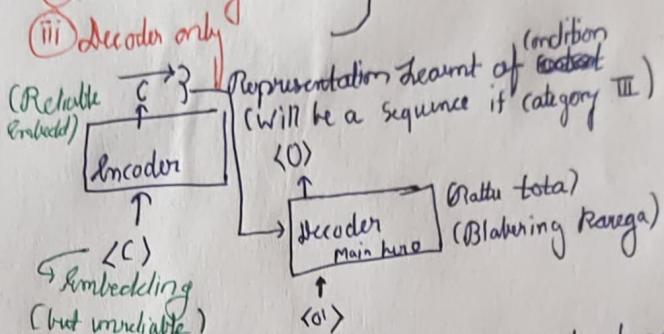
$\langle 0 \rangle \rightarrow \text{can just be start (in case of minimum subsequence)}$

Eg of ① Translating sentence from 1st language to 2nd but you want to translate particular part of sentence

English → Hindi Artificial Intelligence will take over the world - Artificial Intelligence dumya ne kabza kar lega.

Three types of LM

- ① Encoder - Decoder } How to make it
- ② Encoder only } (conditional LM)
- ③ Decoder only }

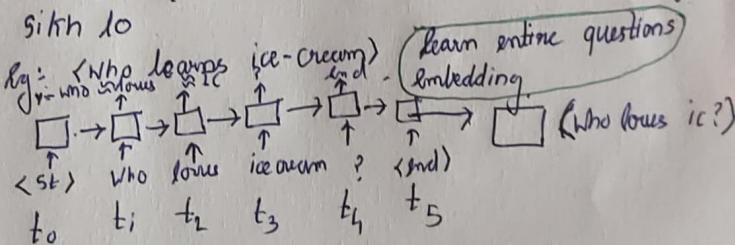


We are injecting question or condition to decoder through encoder.

In toughest case scenario the '*O*' should be start (We want to train it perfect)

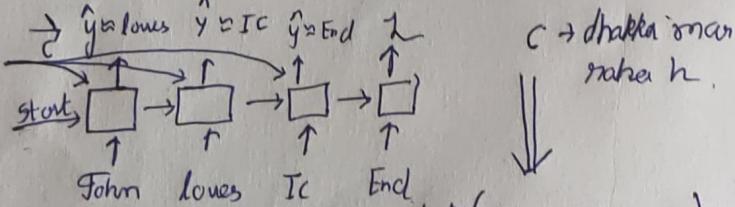
- Condition narrates how the decoder flavours the information

• We can use vanilla RNN as encoder & decoder
→ first train (learn) encoder. Then '*C*' via embedding



You will be passing a lot of question learning W_x^c, W_y^c, W_z^c P encoder

Once trained, you will plug it into decoder.
(It will be similar to encoder but with two inputs)



$$h'_t = W_h + h_{t-1} + b_h \quad h_t = \text{tanh}((h'_{t-1} + \vec{c}) + \vec{x}_t)$$

$$x'_t = W_x * x_t + b_x$$

Prayag dhakka