

Topics

- ① word2vec
- ② RNN (Simple RNN)
 - ✓ Sequential data
 - Time series data
 - Sentences
 - Theory
 - Time series data
 - Char level prediction
 - Word level prediction
 - Implementation (Keras/Tensorflow)
- ③ LSTM
 - Implementation
 - Theory

Concept

- Theory using example
- Math intuition
- + derivation
- Implementation

① word2vec

- ② RNN
 - ↳ Simple RNN
 - ↳ LSTM

③ Attention models / Transformers

- ④ Real time case studies
 - ↳ Speech recognition
 - ↳ chat bot (CRASA)
 - ↳ Sentiment analysis

Word 2 vec

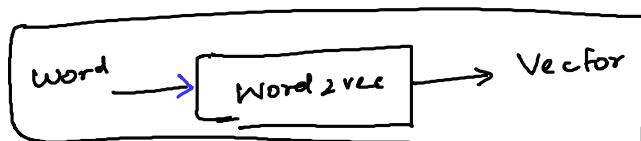
- Document term matrix
- Sparse (98%)
 - High dimension
 - Content



→ LSA

→ Word embeddings

- Word2vec ✓
- Glove -
- fast text



$[d_1, d_2, \dots, d_m]$
 $M \approx 300$

Objective

- Should be dense
- Should be in lower dimension (300)
- Should represent meaning of the word

$\text{Sim(benz, Audi)} \sim \text{high}$

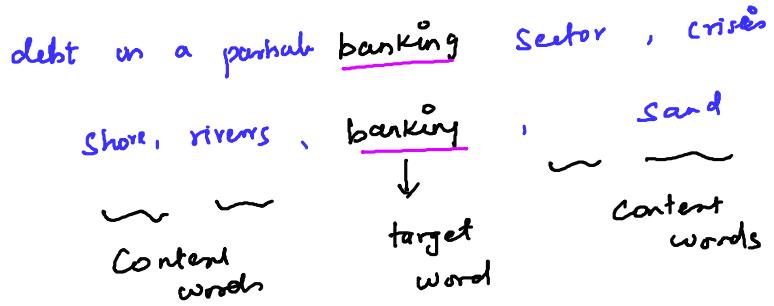
$\text{Sim(King, man)} \sim \text{high}$

Semantic field - hospital(Doctors, nurses, surgeon, ..)

- Should be comparable with each other

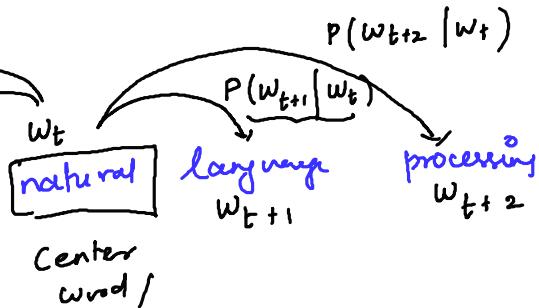
Distributional semantics

"You shall know a word by the company it keeps"
(J.R. Firth 1957)



Language models

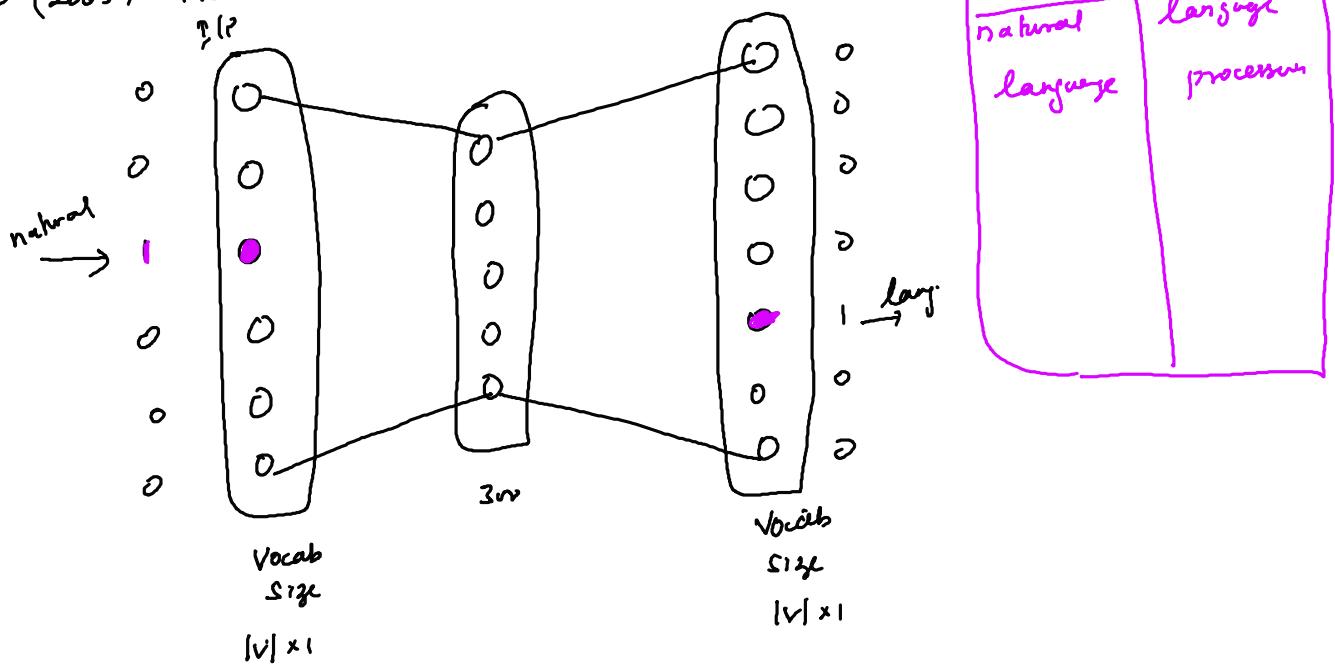
This course is about w_{t-1}



→ Probability model

→ Markov model

→ (2003) Neural network language model (NNLM) - Bengio et al 2003



→ NNLM → Output layer → $|V| \times 1$ → high

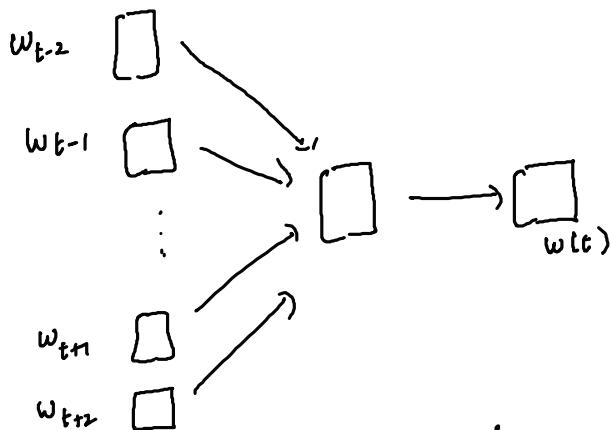
→ 2013 → word2vec

Word2Vec

→ Skip gram

→ Continuous bag of words (CBOW)

CBOW

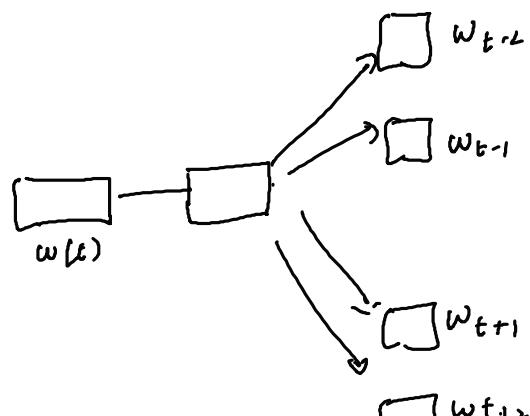


→ Predict the center words given the context words

→ extremely faster compared to skip gram

→ Good representation for frequent words

Skip Gram



→ Predict the context words given the target word

→ slower compared to CBOW

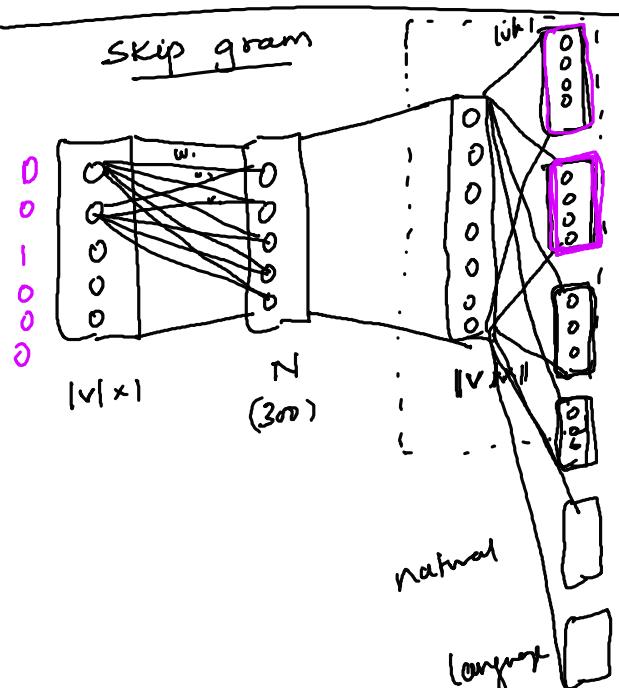
→ Good representation for rare words

CBOW & skip gram → predict words ✓

Vector representation of words

window size = 1

Skip gram



thin

Sent1 → Thin
Sent2 → Next

is

course

about

natural

language

next

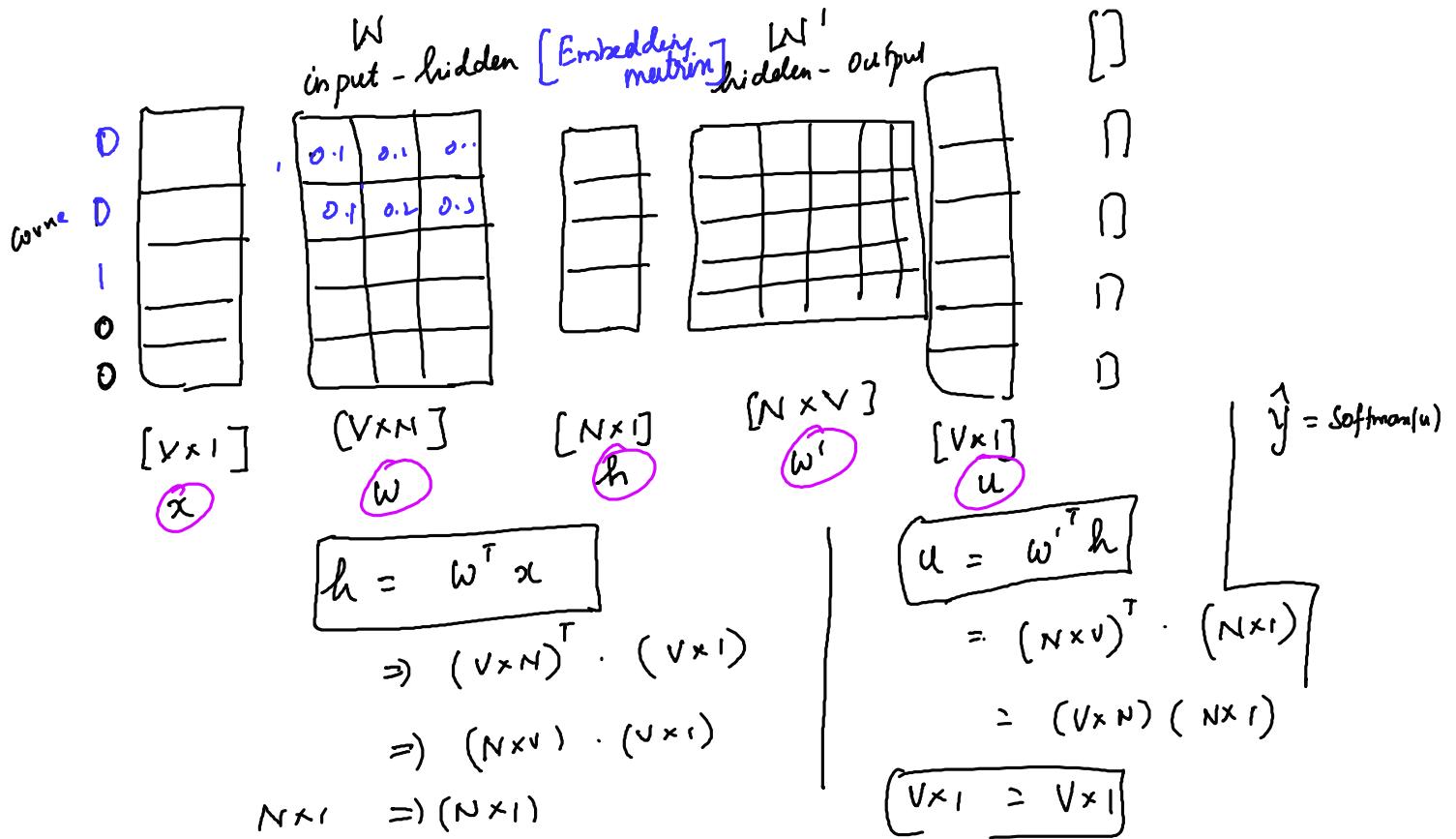
big

data

course in about natural language
course is big data

Center word	Content words
This ✓	course ✓
course	(thin, is)
is	(thin, course, about natural)
course	(next, is)

Thin	1	0	0	0	0	0
is	0	1	0	0	0	0
course	0	0	1	0	0	0
about	0	0	0	1	0	0
natural	0	0	0	0	1	0
language	0	0	0	0	0	1
next						
big						
data						



$$\begin{bmatrix} w^T x \\ \vdots \\ w^T x \end{bmatrix}_{N \times V} \cdot \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}_{V \times 1} \Rightarrow \begin{bmatrix} 0.1 & 0 & 0 & 0 \\ 0.2 & 0 & 0 & 0 \\ -0.1 & 0 & 0 & 0 \\ -0.1 & 0 & 0 & 0 \end{bmatrix}$$

Word Embedding matrix

V_1	$\begin{bmatrix} \cdot & \cdot \end{bmatrix}_{V \times N}$
V_2	
V_3	
:	
V_n	

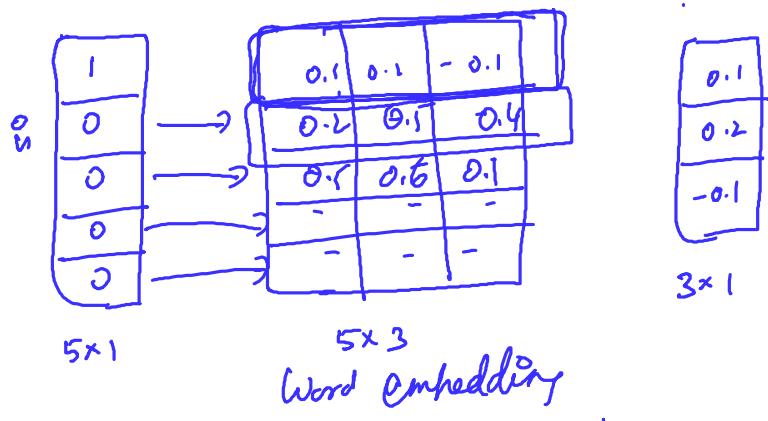
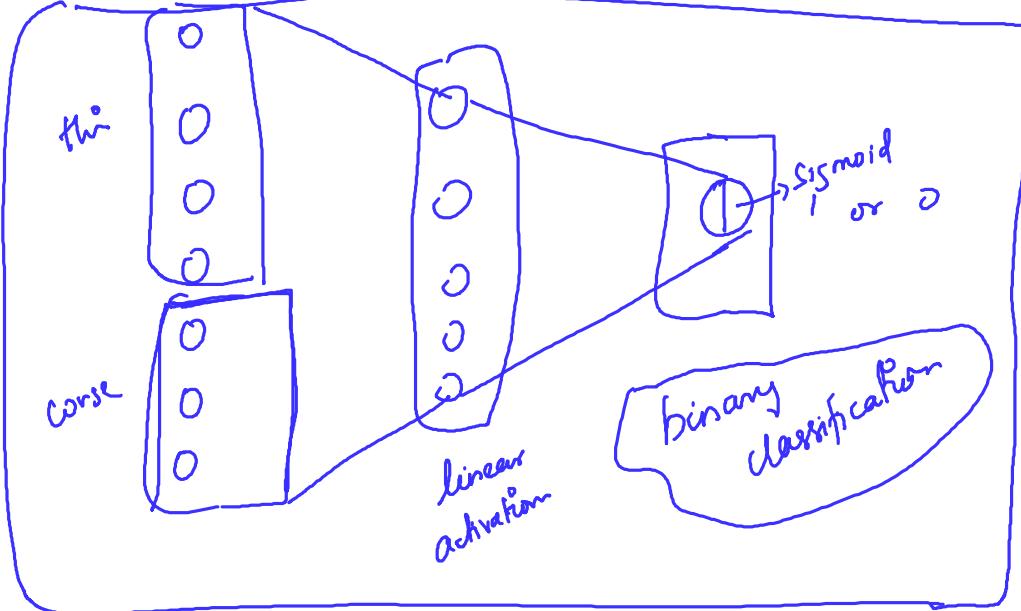
Center word	Context word
this	course
course	(this, b)
b	this, math

negative sampling

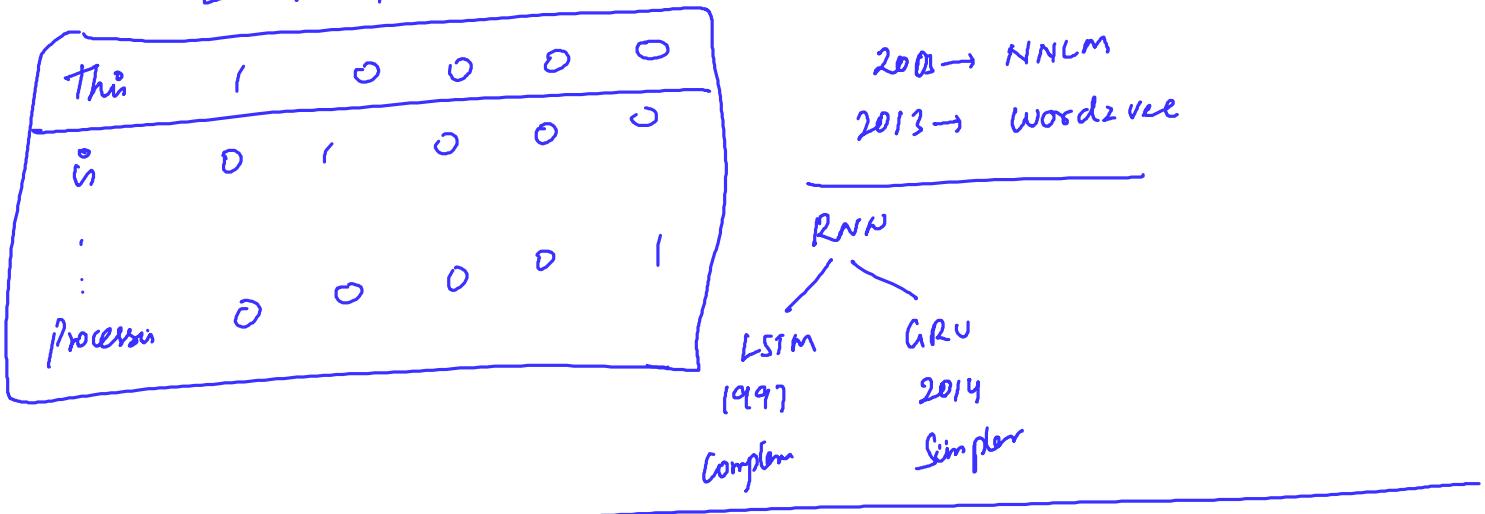
Center word	Context word	
this	course	positive pair
this	math	→ negative pair
this	natural	→ negative pair
course	this	→ positive pair
course	is	→ positive pair
math	natural	→ negative pair

Word2Vec

Skipgram with negative sampling



Vocabs $[w_1, w_2, v, \text{natural}, \{c\}y\text{a}g\text{e}, w_5, \text{process}]$



\rightarrow [This, course, is, on, nlp]
 \rightarrow [1, 2, 100, 5, 300]

Vocab
 this - 1
 course - 2

① Tokenization

② Pad sequences to make size of all document equal

Doc 1 [This is about nlp]

\leftarrow [1, 100, 300, 5 0 0 0 0 0 0]

Doc 2 [This course teacher nlp basic and con...]

\leftarrow [1, 50 350 1 0 0 0]

doc \Rightarrow SD

Vocab = ["a", "the", "movie", "go", "ba", "..."] \rightarrow 44,000 words

	d_1	d_2	\dots	d_{300}
v_1	0	0	0	0
v_2	0	0	0	0
\vdots				
v_{44000}	0	0	0	0

embedding ['a'] = []
 embedding ['the'] = []

Word embedding matrix

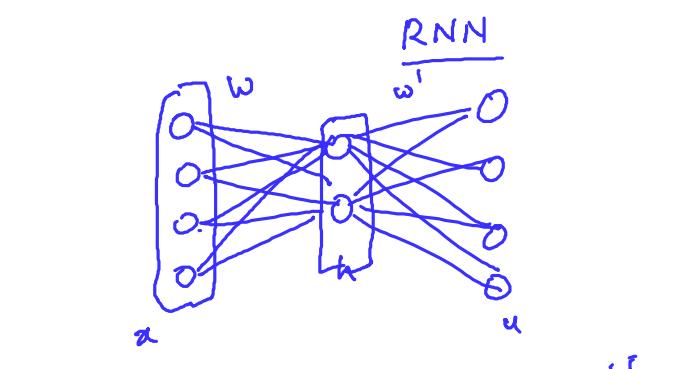
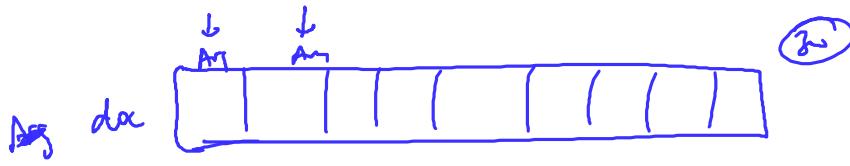
44000×300

$2^{44000} \times 300$

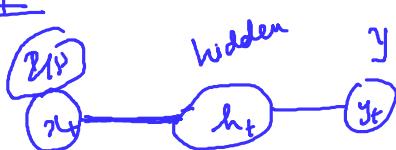
doc1 $\left[\begin{array}{|c|c|c|c|c|} \hline r_1 & r_2 & & \dots & r_v \\ \hline \end{array} \right]$

doc \Rightarrow ("This a nlp course")

{This							
is							
nlp							
course							

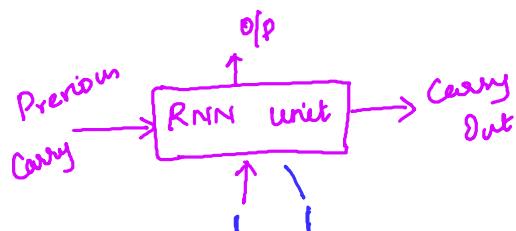


Forward pass $h = W^T x$ $u = W^T h$.

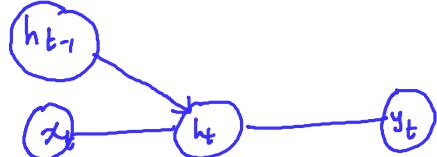
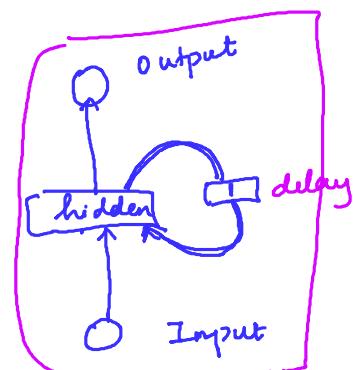
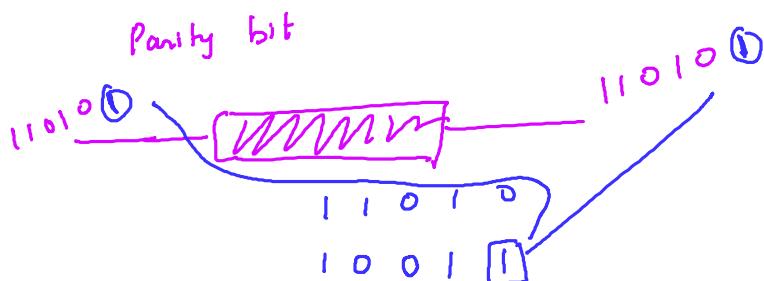


$$h_t = f(x_t, \theta)$$

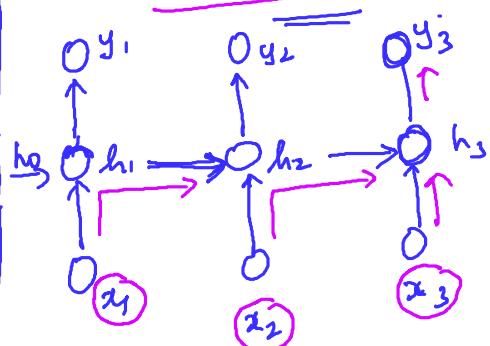
$$y_t = \theta h_{t-1} + \text{error}$$



$$\begin{array}{r} 1101 \\ 1010 \\ \hline 11010 \end{array}$$



Unfolded version



$$h_1 = f(x_1, h_0)$$

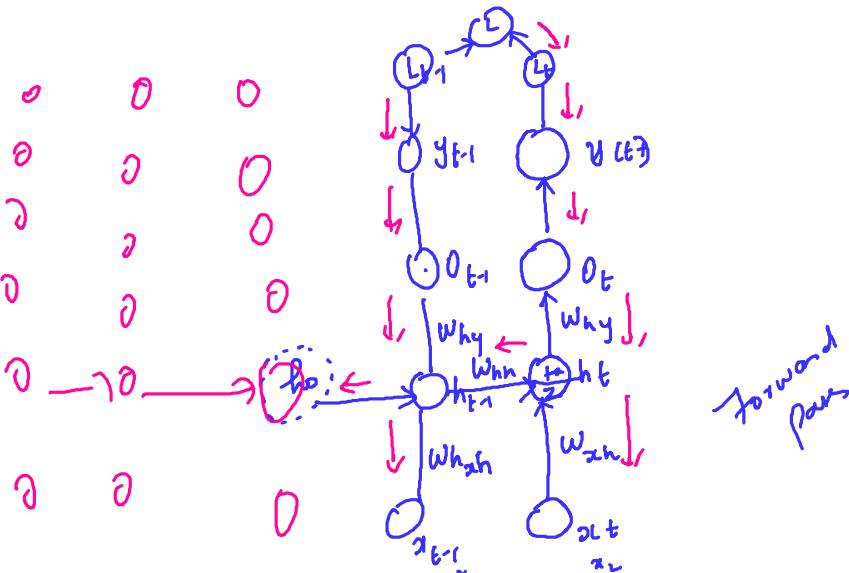
$$h_2 = f(x_2, h_1)$$

$$h_3 = f(x_3, h_2)$$

$$h_3 = f(x_3, f(x_2, f(x_1, h_0)))$$

$$h_n = f(x_n, f(x_{n-1}, f(x_{n-2}, \dots, f(x_1, h_0))))$$

$$y_n \sim f(x_n, x_{n-1}, x_{n-2}, \dots, x_1)$$



$$\theta = (W_{xh}, W_{hy}, W_{hh}, y)$$

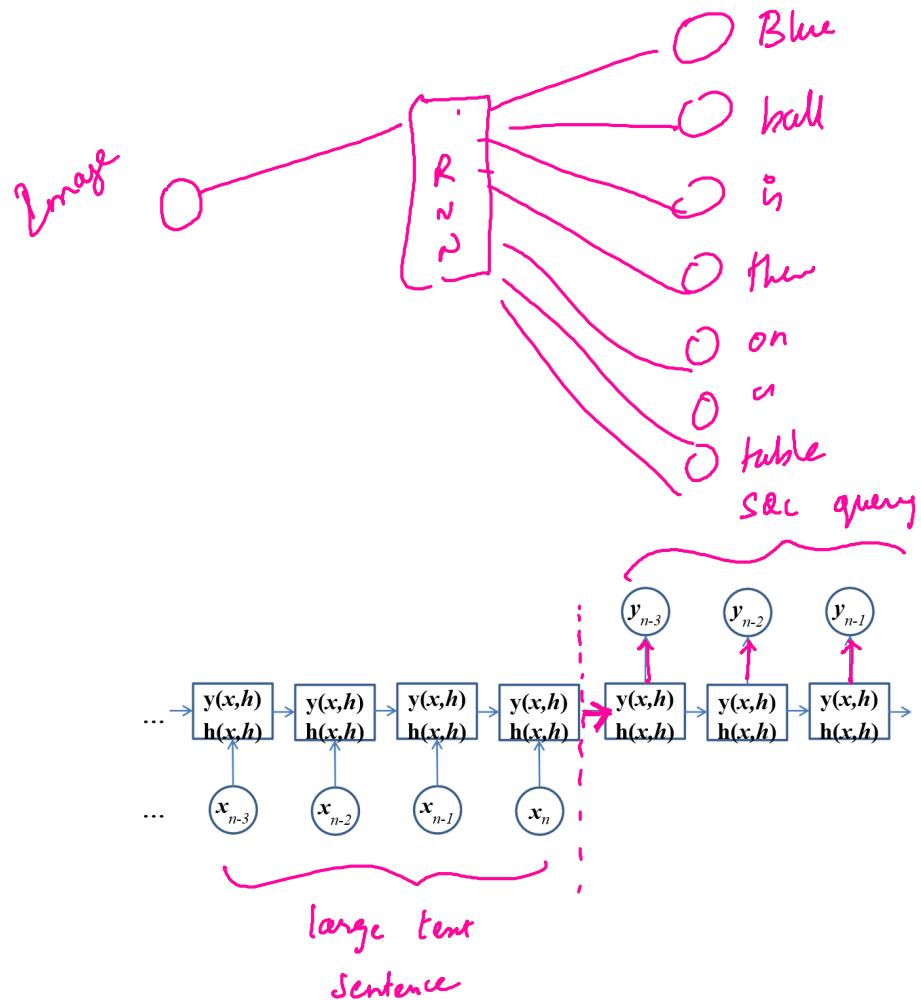
$$z_t = W_{xh}x + W_{hh}h_{t-1}$$

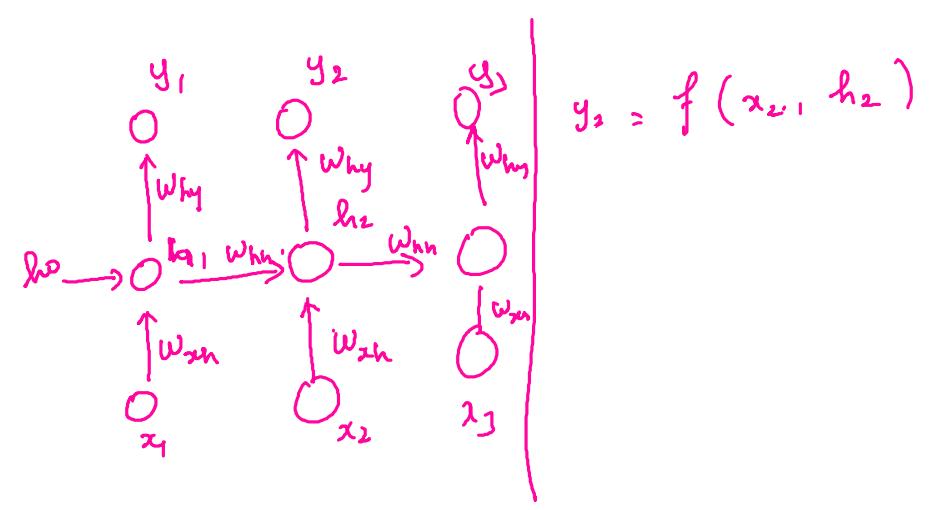
$$h_t = \tanh(z_t)$$

$$o_t = W_{hy}h_t$$

$$y_t = \text{softmax}(o_t)$$

Backward propagation through time (BPTT)





$$y_2 = f(x_{2,1}, h_2)$$