

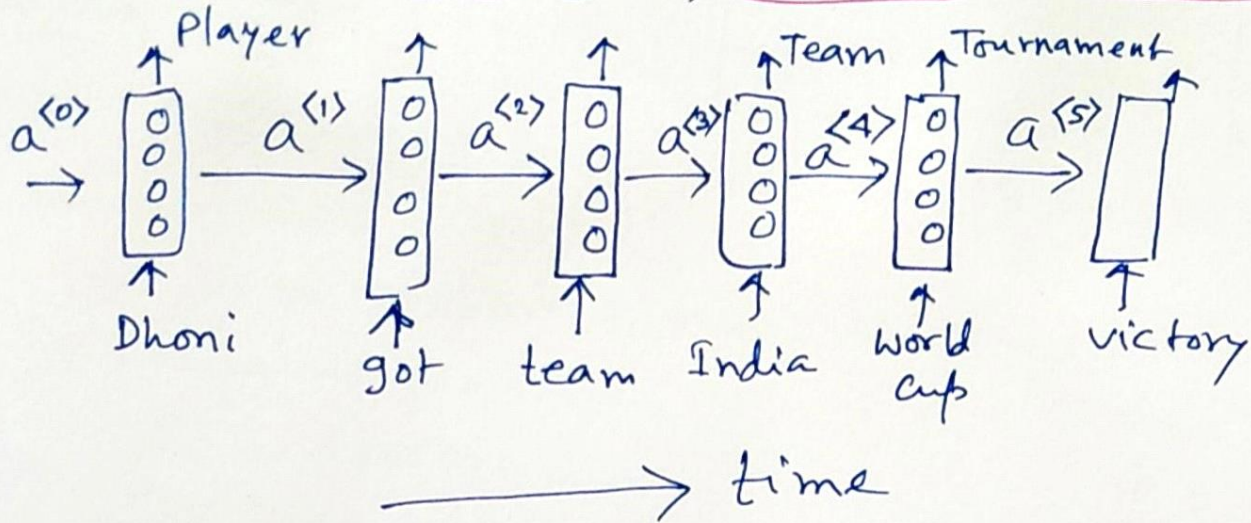
Word Embedding

Page 1

Suppose we want to build a NLP model that can recognise entities in text:

In 2011, ^{Player} Dhoni got team ^{Team name} India a ^{Tournament} world cup victory.

In the last ^{Tournament} Ashes game, ^{Player} Cummins led ^{Team name} Australia to a victory.



PROB: The machine can't understand text

→ We have to convert them to numbers.

For training, we can scrap the internet.

→ Make a vocabulary: list of words.

Ashes	1
Axar	2
bat	3
Dhoni	7
team	1023
:	
Zimbabwe	50,000

Tokenization

Now we can put "7" in place of "Dhoni" to the RNN.

Prob: These numbers are random. They don't capture relationship between words.

Second option: one hot encoding

Probs: - (i) They don't capture relationship between words

(ii) Computationally inefficient.

Word embedding

Pg 2

	Ashes	Axar	Bat	...	Dhoni		Zimbabwe
Ashes	1	0	0	...	0		0
Axar	0	1	0	...	0		0
...				...			
Zimbabwe	0	0	0	...	0		1

Word embedding: How to capture similarities between two words?

How do we compare similarities between two homes?

By comparing features:

- (i) Bedrooms
- (ii) Area
- (iii) # of bathrooms etc.

Your	mine	Palace
4	3	20
1850	1500	10,000
3	2	7

Hand-crafted features:

	Dhoni	Cummins	Australia
Person	1	1	1
Healthy/Fit	0.9	0.87	0.7
Location:	0	0	1
Has two eyes.	1	1	0
Has government	0	0	1

↑
↑
↑
vectors

This is called word-embedding

Automatic features

	Ashes	Australia	Cummins
Person	0	0.02	0.95
Country	0	0.97	0
Healthy & Fit	0	0	0.87
event	1	0.1	0
gear	0	0	0

↑ obtained through training.