

# Natural Language Processing

Mohammad Amin Samadi

Mohammad Sadegh Akhondzadeh

# What is Natural Language Processing

- A subfield of linguistics, Artificial Intelligence and Cognitive Sciences
- An interdisciplinary subject
- Aim: To build intelligent computers that can interact with human being like human beings

# Why NLP

- Huge amounts of data Internet = at least 20 billions pages.
- Text data: web sites, blogs, tweets.
- Audio data speech.

# Applications

## Information Retrieval

Doc A



Doc 1

Doc 2

Doc 3

## Sentiment Analysis



## Information Extraction



## Machine Translation



# Text Processing

## Question Answering



Human: When was Apollo sent to space?



Machine: First flight - AS-201, February 26, 1966

# Why NLP is HARD

## Cross Language Problems

- Ambiguity
- Compression
- Example sentiment analysis

## Problem of other languages

- German: Donaudampfschiffahrtsgesellschaftskapitän (5 “words”)
- Chinese: 50,000 different characters.
- Japanese: 3 writing systems

# Examples of sequence data

Speech recognition



→ "The quick brown fox jumped over the lazy dog."

Music generation



Sentiment classification

→ "There is nothing to like in this movie."



DNA sequence analysis → AGCCCCTGTGAGGAACTAG

→ AGCCCCTGTGAGGAACTAG

Machine translation

Voulez-vous chanter avec moi?

→ Do you want to sing with me?

Video activity recognition



→ Running

Name entity recognition → Yesterday, Harry Potter met Hermione Granger.

→ Yesterday, **Harry Potter** met **Hermione Granger**.

# Sequence in Language



# Word Representation

How to express words to a computer?



# One-hot Representation

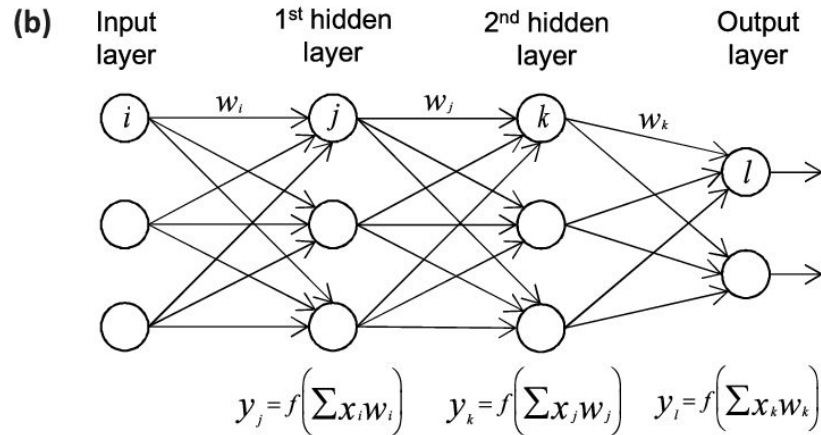
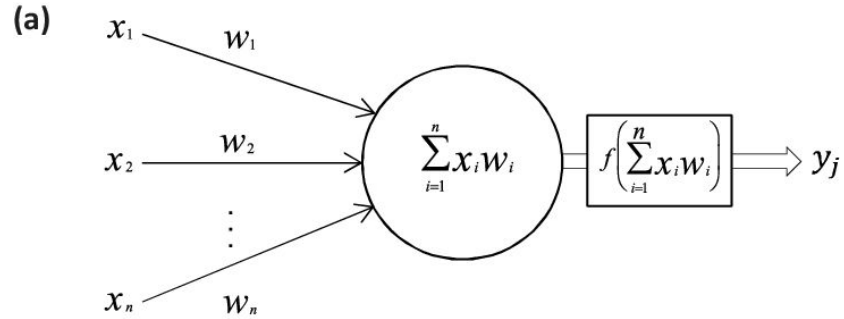
Problems :

- Meaning relations and similarities ignored
- Inefficient Memory consumption
- Vector size dependant on vocabulary size

"a"	"abbreviations"	"zoology"
1	0	0
0	1	0
0	0	0
.	.	.
.	.	.
.	.	.
0	0	0
0	0	1
0	0	0

# Neural Network

- Based on human's brain and nervous system
- Approximates a function from input to output



# Word embedding

- Featurized representation of words
- Each words is embedded to a vector in a 100 or 200 or ... dimension space
- Learned from large text corpus ( 1-100B words)
  - We can train it
  - Download pre-trained
- Helps us learn a context with fewer examples



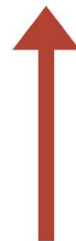
[drove] [my] [high] [speed] vehicle [down] [the] [road] [today]



Context



Centre Word



Context

# Skip-gram

- Fake task explanation
- Tries to predict the context based on the center word

Source Text

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

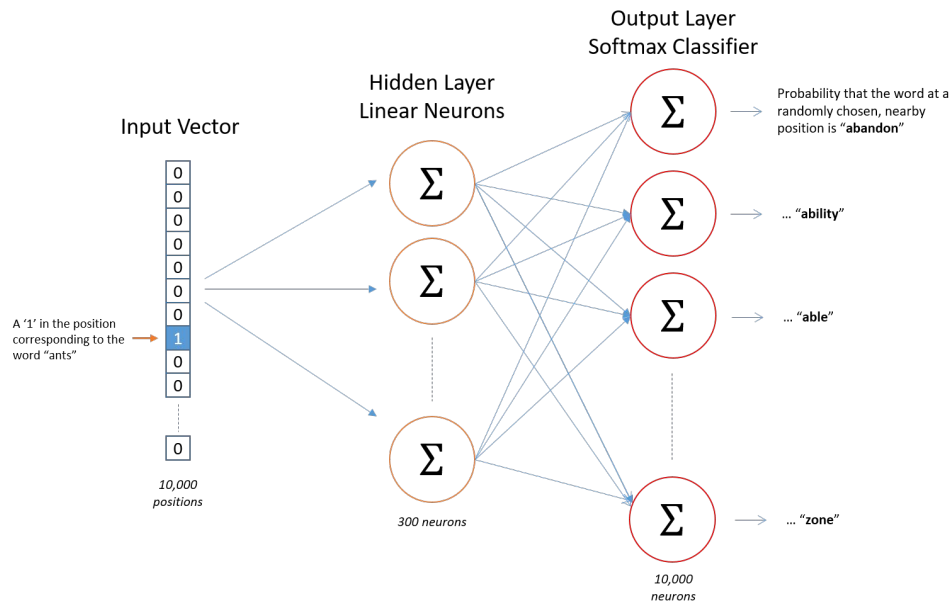
Training Samples

(the, quick)  
(the, brown)

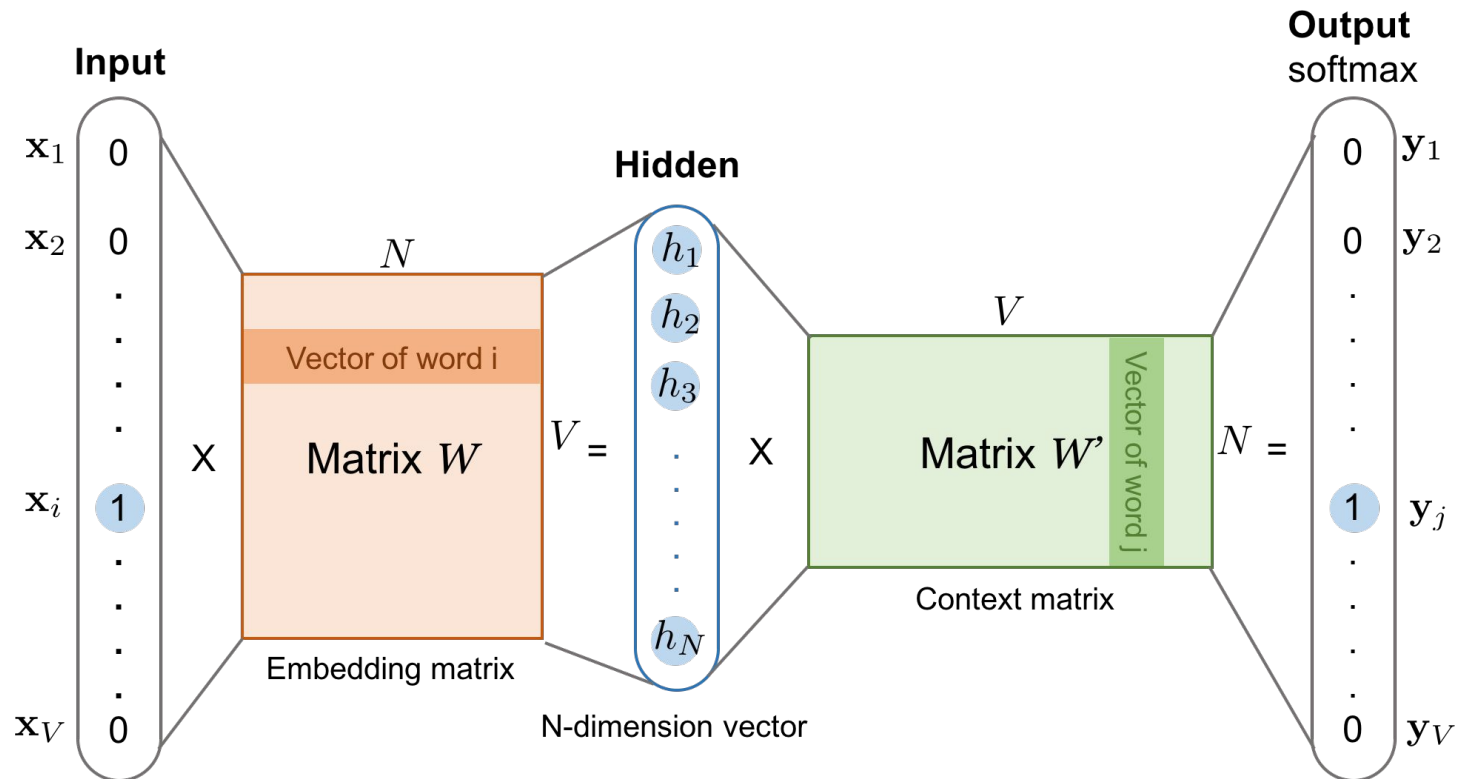
(quick, the)  
(quick, brown)  
(quick, fox)

(brown, the)  
(brown, quick)  
(brown, fox)  
(brown, jumps)

(fox, quick)  
(fox, brown)  
(fox, jumps)  
(fox, over)

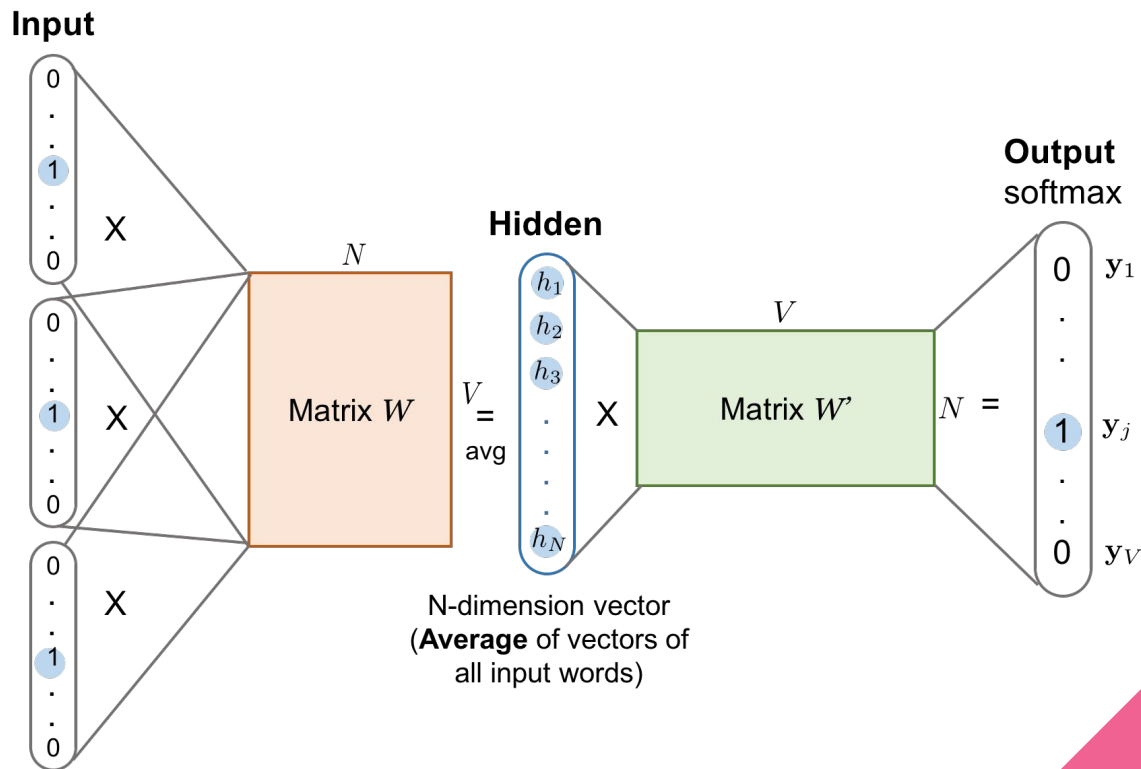


# Representation Learning (Skip-gram)



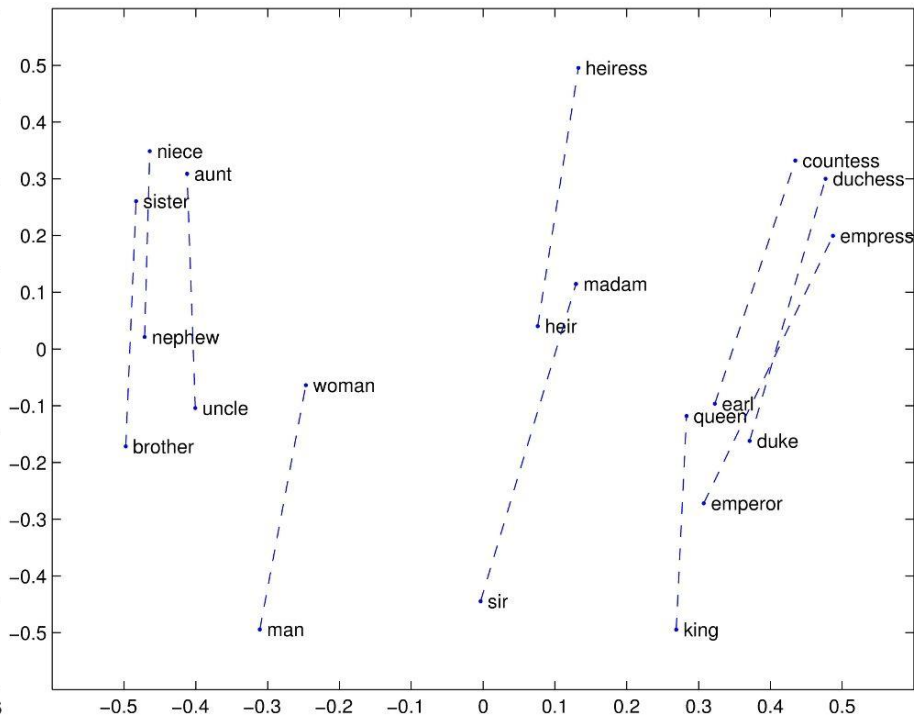
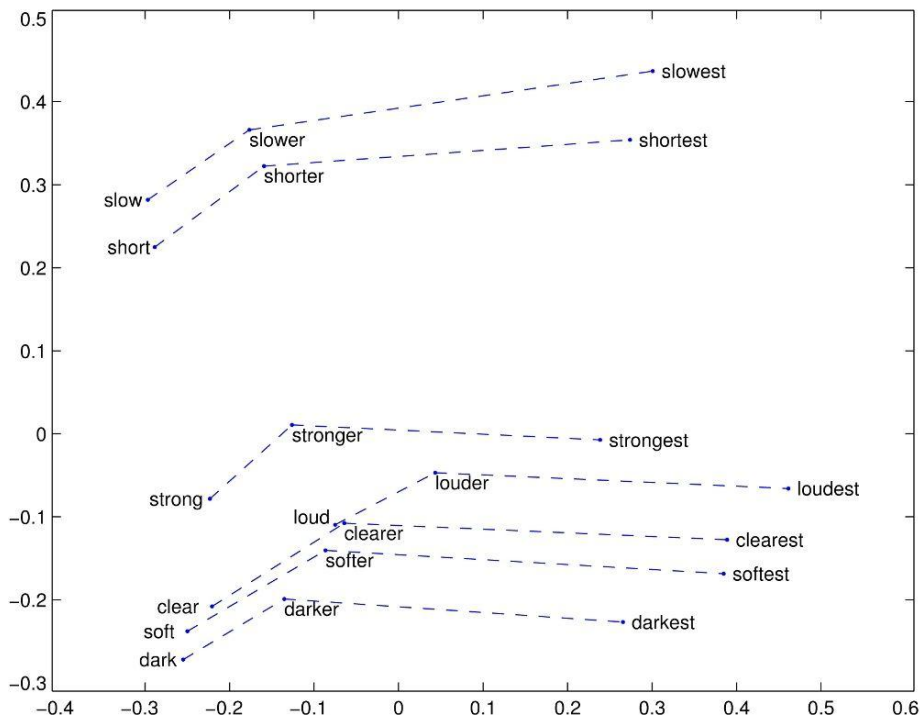
# Representation Learning (CBow)

Tries to predict the center word based on the context



# Vector offset

can encode semantic and syntactic relations between words





# Dimension Reduction

