



Session 2

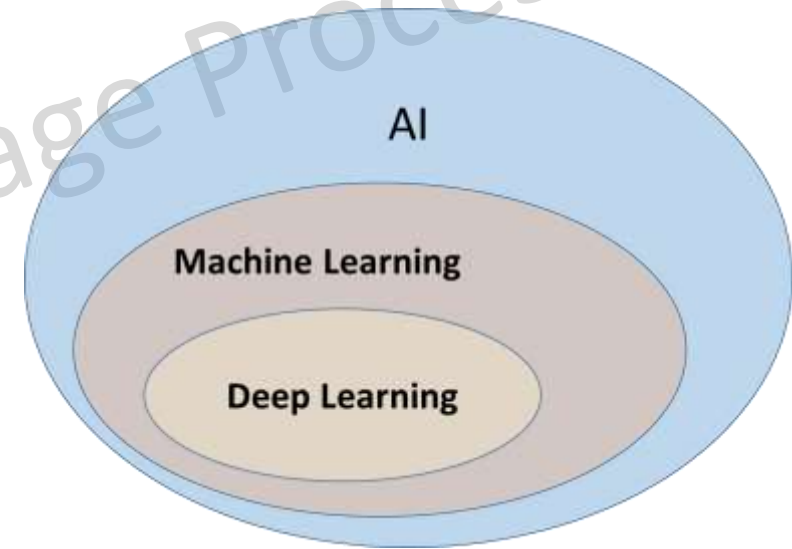
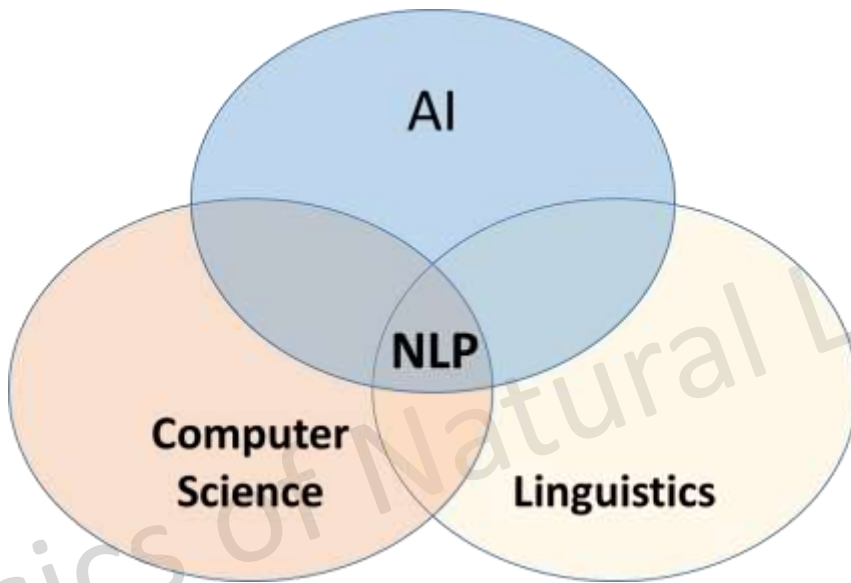
# Basics of Linguistics

Azam Rabiee, PhD

August 23, 2020

# Roadmap

We will review ML/DL methods for NLP



# Digikala Academy NLP Events

## Basic

- **Session 1.** Introduction
- **Session 2.** Basics of Linguistics
- **Session 3.** Basics of ML
- **Session 4 (Lab).** Effective Word Representation by python

## Intermediate

- TBA

## Advanced

- TBA

# Outline

## Session 1: Introduction

- Applications
- Tasks
- Approaches

## Session 2. Basics of Linguistics

## Session 3. Basics of ML

## Session 4 (Lab). Effective Word Representation by python

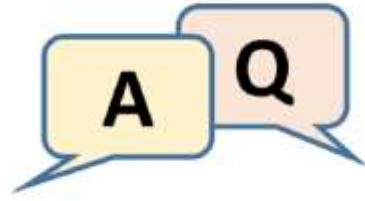
# Review: Applications



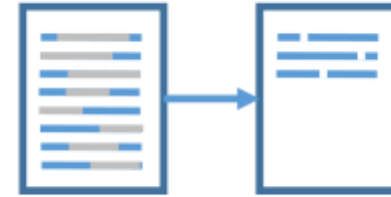
Machine Translation



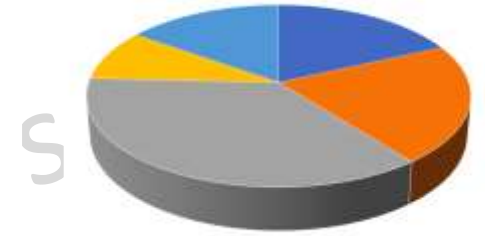
Sentiment Analysis



Question Answering



Automatic Summarization



AI Marketing



Text / Document Classification



Speech Recognition



↓  
Give me

Handwritten Character Recognition

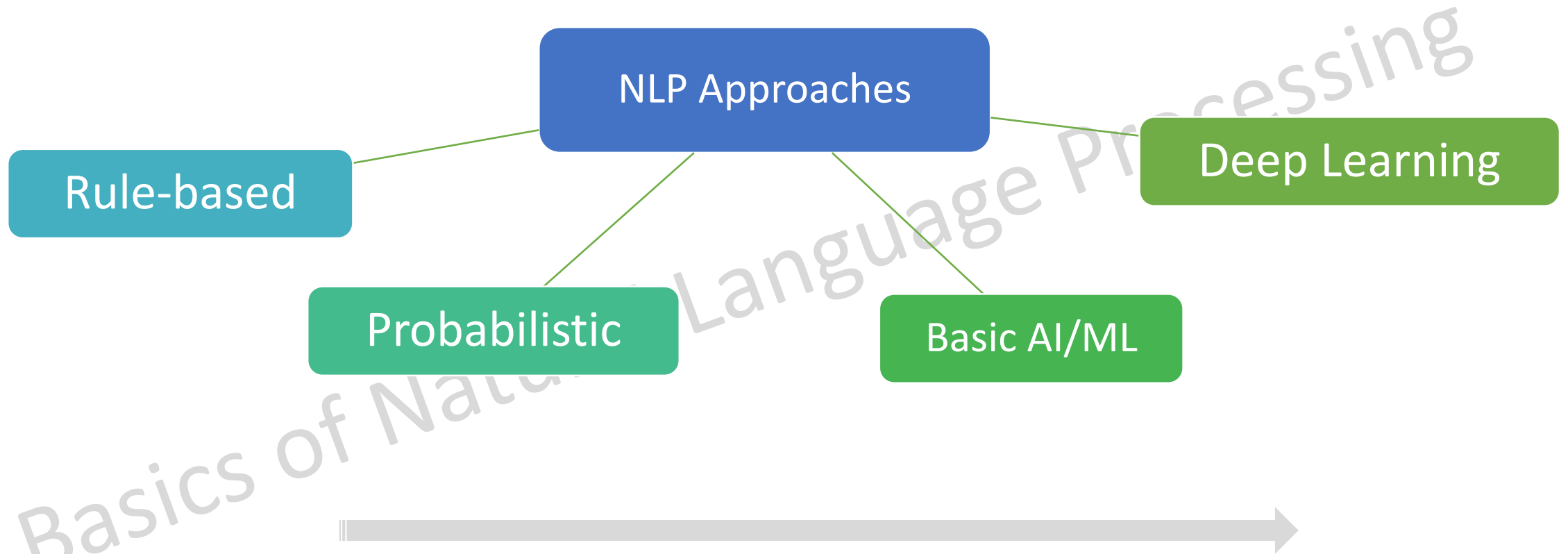


Spell Checking

# Review: Tasks



# Approaches



# Outline

Session 1: Introduction

**Session 2. Basics of Linguistics**

- Components
- Challenges
- Vectorization

**Session 3. Basics of ML**

**Session 4 (Lab). Effective Word Representation by python**



# Outline

Session 1: Introduction

Session 2. Basics of Linguistics

- **Components**
- Challenges
- Vectorization

Session 3. Basics of ML

Session 4 (Lab). Effective Word Representation by python

# Components

## Tokenization

Breaking up the sequence of characters into sentences and words

A dog is chasing a boy on the playground.

# Components

## Tokenization

Tokenizer also extracts token features, such as:

- Capitalization

- Inclusion of digits

- Punctuation

- Special characters

The MRI, I took on 2019, priced at \$500!!

# Components

## Part-of-Speech (PoS) Tagging

A   dog   is   chasing   a   boy   on   the   playground  
article   noun   aux   verb   article   noun   preposition   article   noun

## Stemming

Normalizing the word to its stem word  
(not necessarily the dictionary word)

`stemmer(`studies`) → `studi``

`stemmer(`studied`) → `studi``

`stemmer(`studying`) → `studi``

`stemmer(`study`) → `studi``

# Components

## Lemmatization

Normalizing the word to its stem dictionary word

`lemmatizer(`studies`)` ➔ ``study``

`lemmatizer(`studied`)` ➔ ``study``

`lemmatizer(`studying`)` ➔ ``study``

# Stemming vs. Lemmatization

stemmer(`am`) → `am`  
stemmer(`was`) → `was`  
stemmer(`were`) → `were`  
stemmer(`is`) → `is`  
stemmer(`are`) → `are`

lemmatizer(`am`, pos=`v`) → `be`  
lemmatizer(`was`, pos=`v`) → `be`  
lemmatizer(`were`, pos=`v`) → `be`  
lemmatizer(`is`, pos=`v`) → `be`  
lemmatizer(`are`, pos=`v`) → `be`

stemmer(`studying`) → `studi`

lemmatizer(`studying`, pos=`v`) → `study`  
lemmatizer(`studying`, pos=`n`) → `studying`

# Stemming vs. Lemmatization

Stemming	Lemmatization
Does not consider the context	Considers the context & PoS
Fast	Slow
Less accurate	Accurate



# Components

## Chunking / Parsing

Extract meaningful phrases, such as:

noun phrases

verb phrases

prepositional phrases

adjective phrases

adverbial phrases

to make sure if a sentence **syntactically** is correct

# Components

## Chunking / Parsing

A   dog   is   chasing   a   boy   on   the   playground  
article   noun   aux   verb   article   noun   preposition   article   noun

noun phrase

complex verb

noun phrase

noun phrase

verb phrase

prepositional phrase

verb phrase

# Components

## Chunking / Parsing

I am feeling hungry

Uni-Gram	I	am	feeling	hungry
Bi-Gram	I am	am feeling	feeling hungry	
Tri-Gram	I am feeling	am feeling hungry		

Parsing one word at a time

Parsing two words at a time

**N-gram: Parsing N word(s) at a time**

# Components

## Chinking

Removing stop words

<u>A</u>	<u>dog</u>	<u>is</u>	<u>chasing</u>	<u>a</u>	<u>boy</u>	<u>on</u>	<u>the</u>	<u>playground</u>
article	noun	aux	verb	article	noun	preposition	article	noun

**dog      chasing      boy      playground**

# Components

---

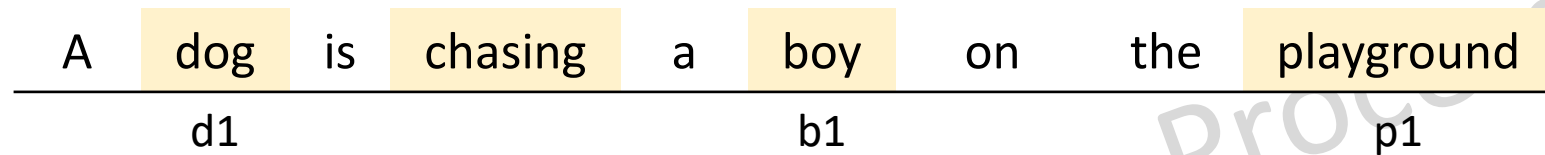
## Semantic Analysis

Understanding the meaning, relationships, and interpretation of words

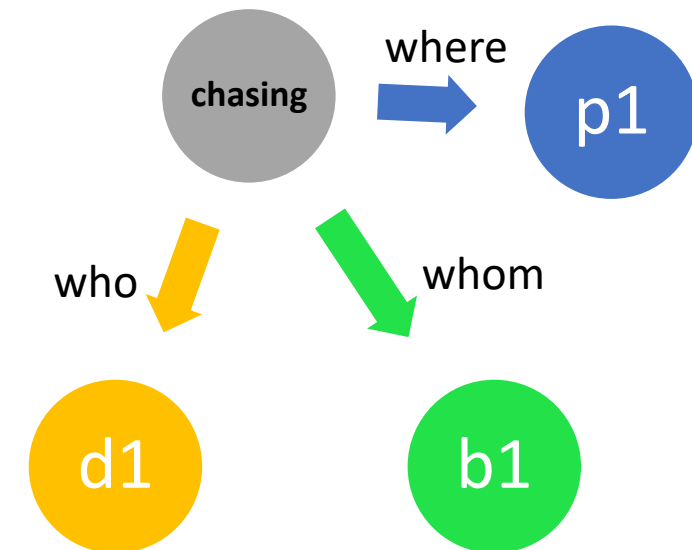
Basics of Natural Language Processing

# Components

## Semantic Analysis



Dog(d1)  
Boy(b1)  
Playground(p1)  
Chasing(d1, b1, p1)



# Components

## Disclosure Integration

Integrating sentences to take into account the context

I know Martin Cook. He works at Google.

A sense of the context: Martin Cook works at Google.

## Pragmatic Analysis

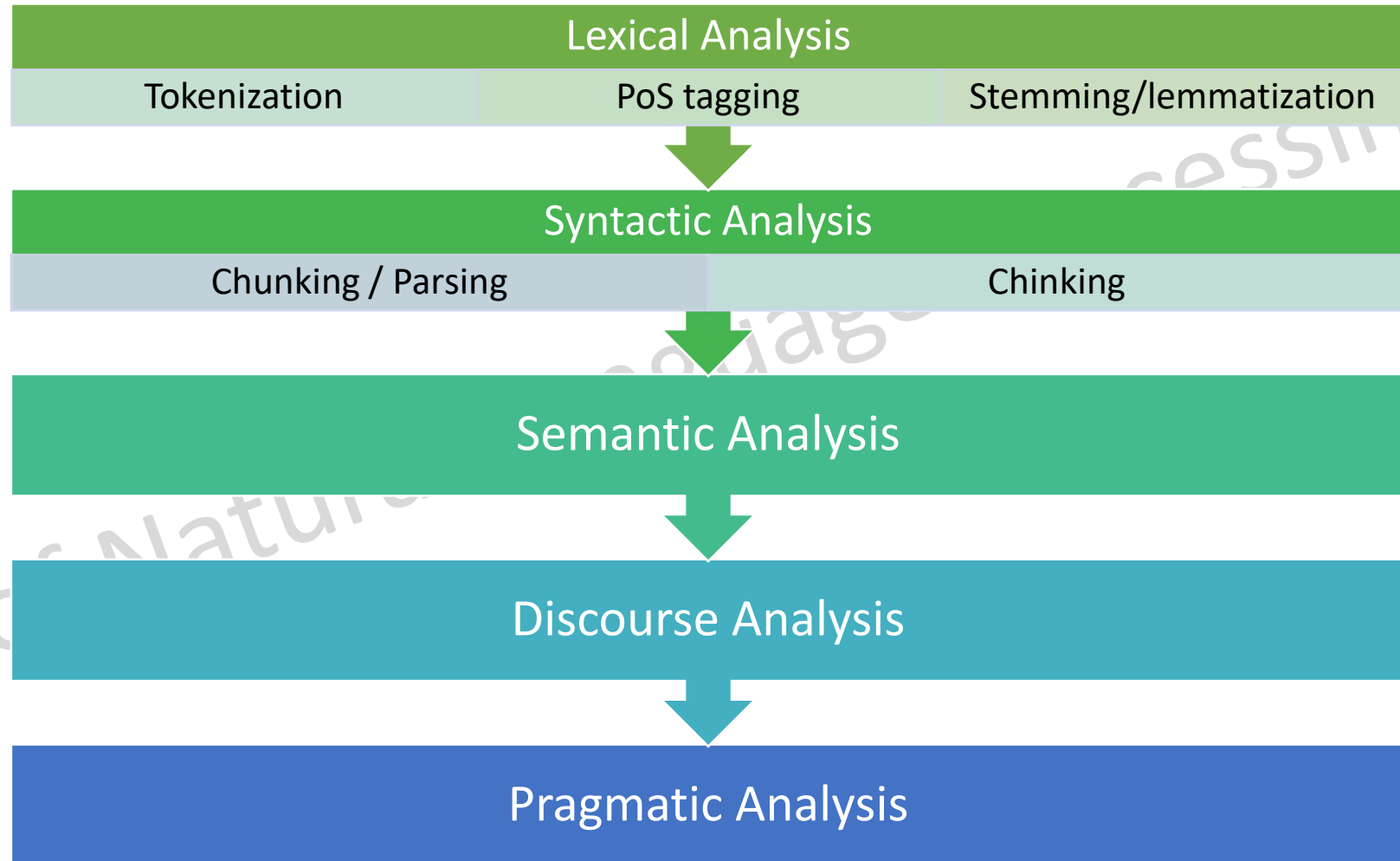
**A:** - Do you know what time it is?  
- Yes. it's 6:00 pm.

**B:** - Do you know what time it is?  
- Sorry. I won't be late anymore.





## Components of NLP





Which type of analysis can find mistakes in the following sentences?

A: Avoid playing when you are **in** tired

B: **Hot** ice-cream

1. A: Syntactic Analysis, B: Semantic Analysis

2. A: Semantic Analysis, B: Syntactic Analysis



**Coreference resolution** is the task of finding all expressions that refer to the same entity in a text. The task is mostly related to:

1. Tokenization
2. Lemmatization
3. Discourse Analysis

*"I voted for Nader because he was most aligned with my values," she said.*

A diagram illustrating coreference resolution. Three curved arrows connect the pronouns to their referents: one from "I" to "she", one from "he" to "Nader", and one from "my" to "she".

AI-hard problem

# Outline

## Session 1: Introduction

## Session 2. Basics of Linguistics

- Components
- **Challenges**
- Vectorization

## Session 3. Basics of ML

## Session 4 (Lab). Effective Word Representation by python

# Challenges

## Tokenization

The main challenge for the sentence tokenizer:

Is it one sentence or two?

A dog is chasing Mr. Cook on the playground.

# Challenges

## PoS Tagging

PoS is crucial for syntactic and semantic analysis

Can you help me with the can?

# Challenges

## Name Entity Recognition (NER)

**Cook** said: “he is coming.”

# Challenges

## Context Understanding / Pragmatic Analysis

**A:** - Do you know what time it is?

- Yes. it's 6:00 pm.

**B:** - Do you know what time it is?

- Sorry. I won't be late anymore.



# Challenges

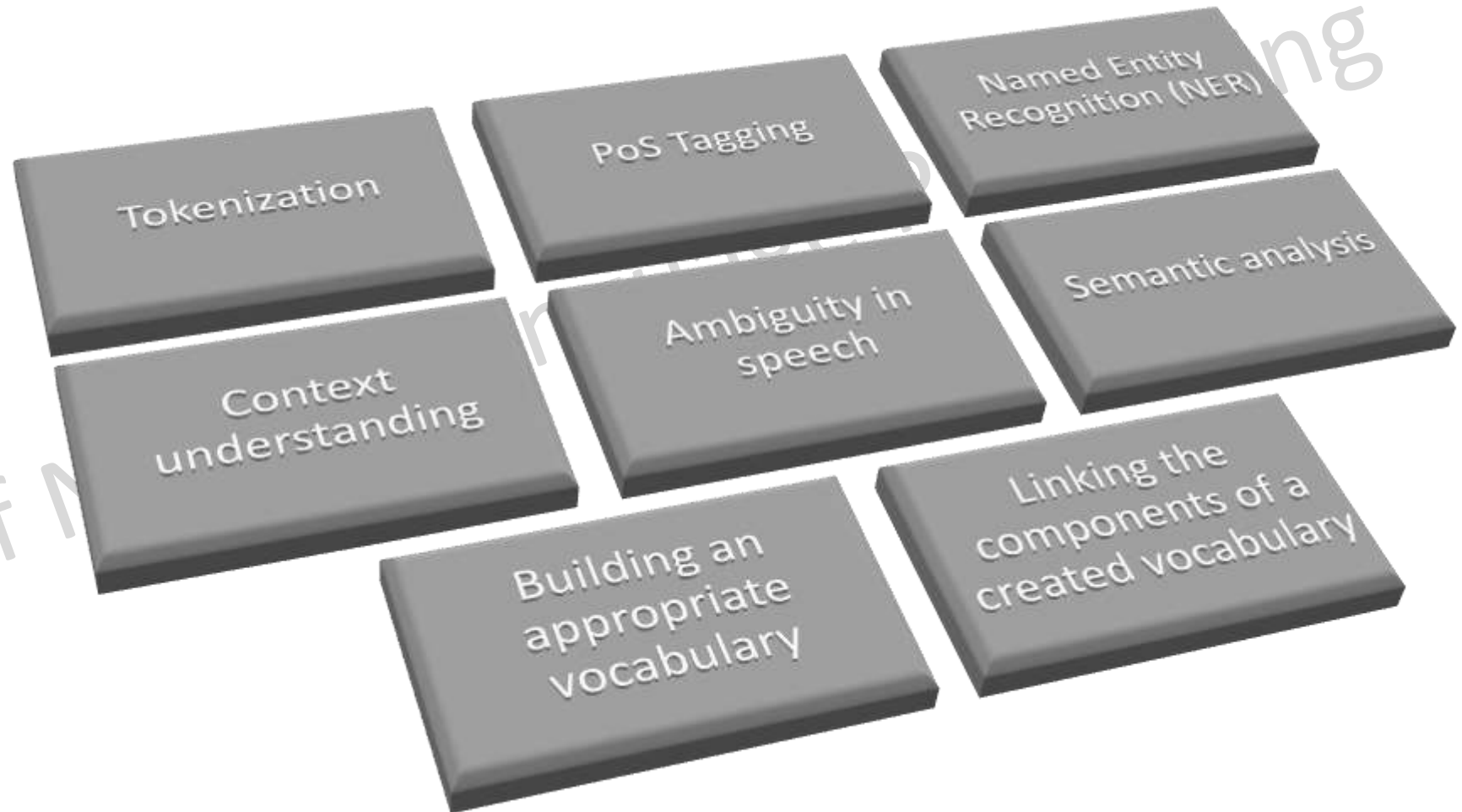
## Ambiguity

I saw a man on a hill with a telescope.

- There is a man on the hill, and I watched him with my telescope.
- There is a man on the hill, and he has a telescope.
- I'm on a hill, and I saw a man using my telescope.
- I'm on a hill, and I saw a man who has a telescope.
- There is a man on a hill, and I saw him something with my telescope.



## Current NLP Challenges

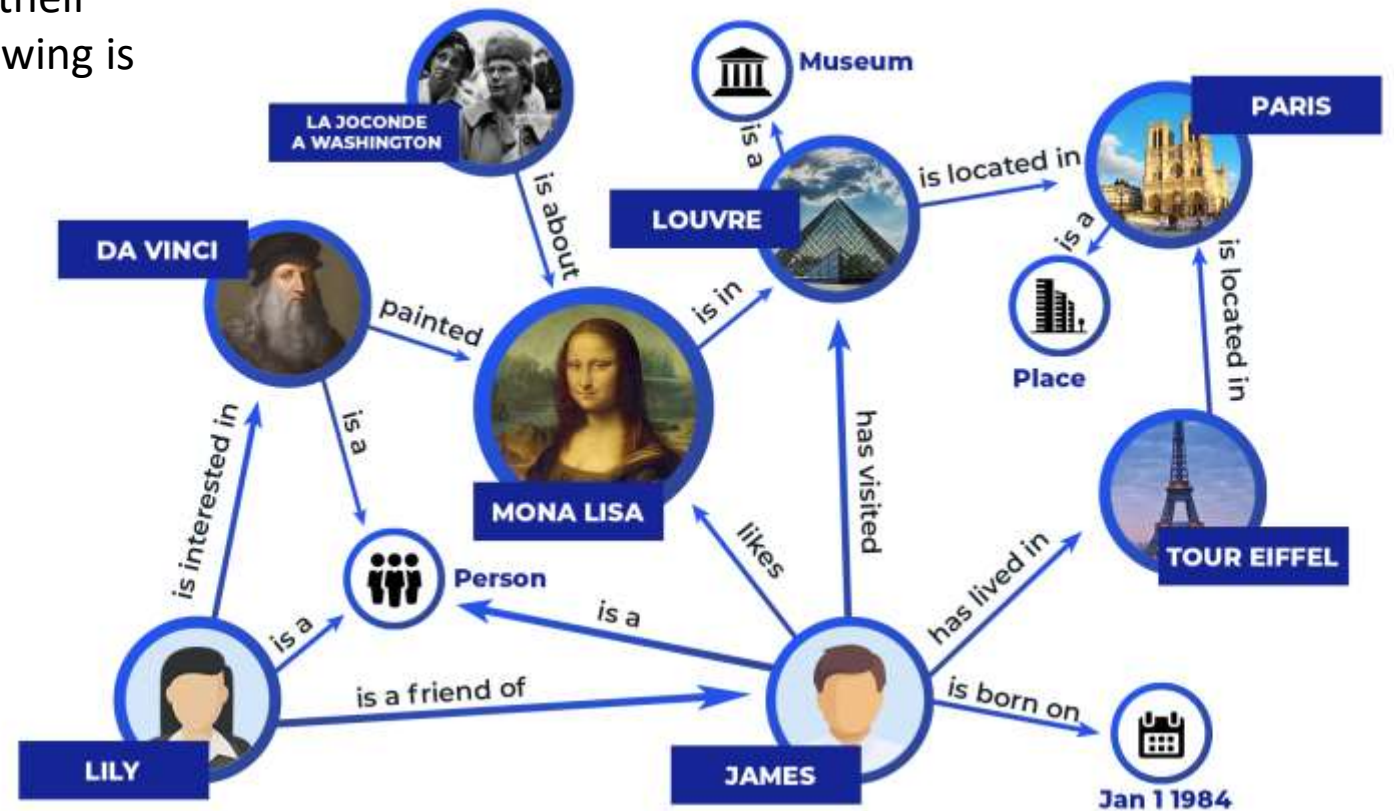


# QUIZ



**Knowledge Graph (KG)** is a network of entities, their properties, and relationships. Which of the following is highly related to KG?

1. Stemming
2. Semantic Analysis
3. Tokenization



# Outline

Session 1: Introduction

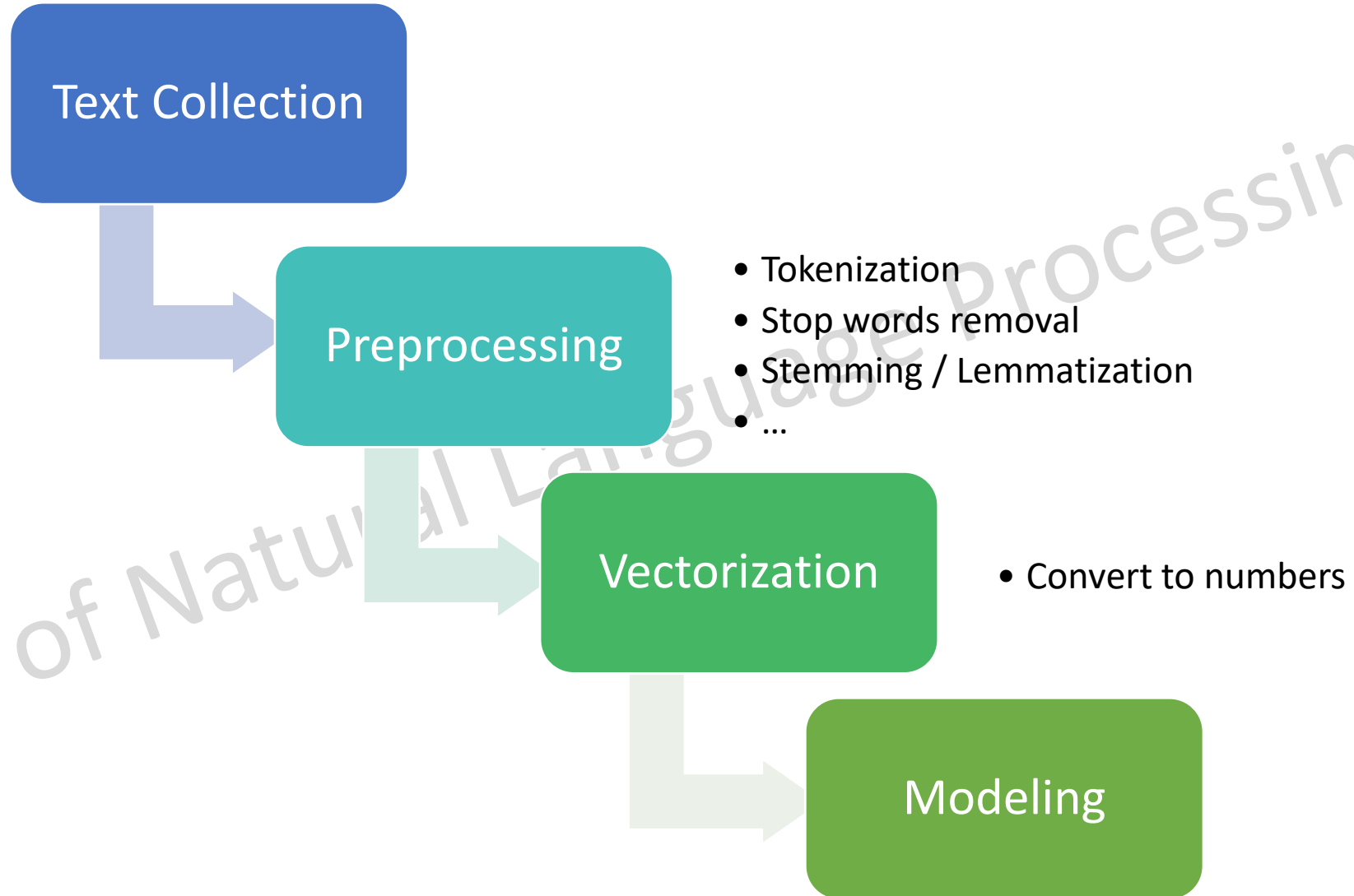
Session 2. Basics of Linguistics

- Components
- Challenges
- **Vectorization**

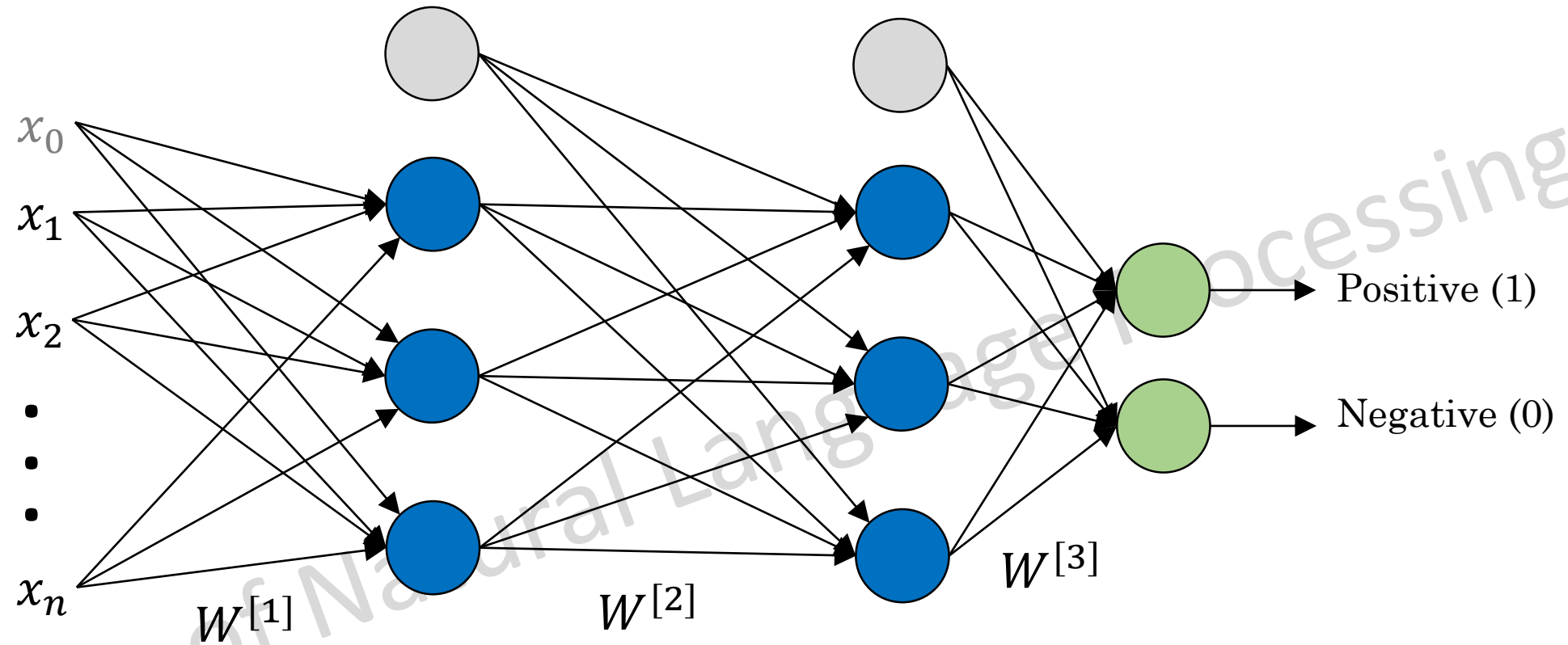
Session 3. Basics of ML

Session 4 (Lab). Effective Word Representation by python

# Steps of ML/DL Projects with Text Data



# ANN Example for Sentiment Analysis



$\mathbf{x} = [x_1, x_2, \dots, x_n]$ : tweets

Example: "This movie was almost good"

$$\mathbf{z}^{[i]} = \mathbf{W}^{[i]} \mathbf{a}^{[i-1]}$$

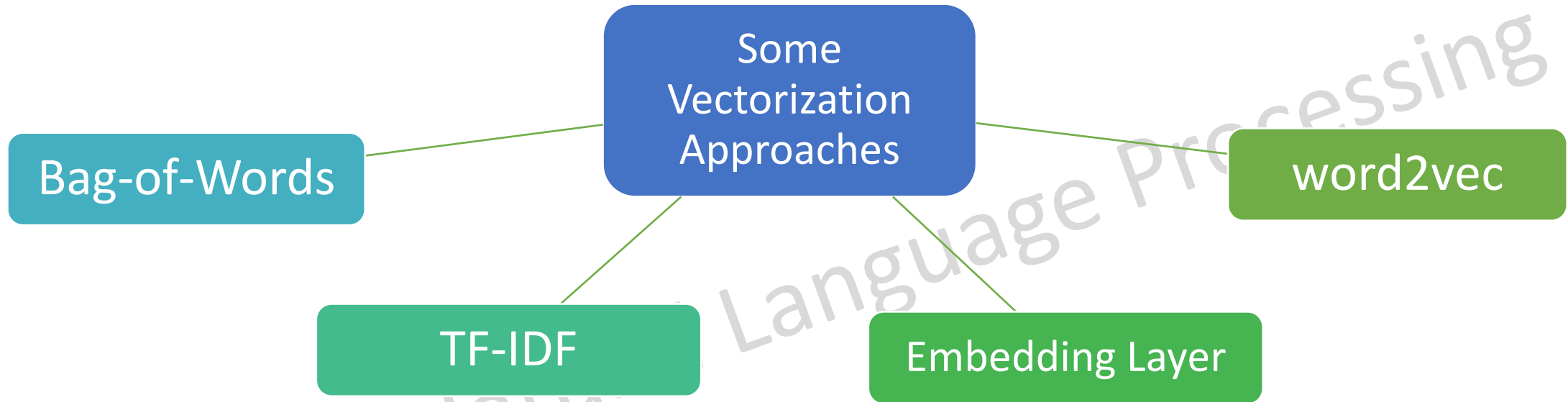
$$\mathbf{a}^{[i]} = f(\mathbf{z}^{[i]})$$



**Vectorization** (converting words to numbers) is unavoidable for ML

Basics of Natural Language Processing

# Vectorization





# Bag-of-Words (BoW)

**S1:** Jim and Pam traveled by bus.

**S2:** The train was late.

**S3:** The flight was full; Traveling by flight is expensive.

	and	bus	by	expensive	flight	full	is	jim	late	pam	the	train	travel	was
<b>S1</b>	1	1	1	0	0	0	0	1	0	1	0	0	1	0
<b>S2</b>	0	0	0	0	0	0	0	0	1	0	1	1	0	1
<b>S3</b>	0	0	1	1	2	1	1	0	0	0	1	0	1	1

# Bag-of-Words (BoW)

## Pros

- Simple
- Fast

## Cons

- Sparse
- High dimension
- No distinction between rare and common words
- No word relation (context)



# TF-IDF

Term Frequency — Inverse Document Frequency

Frequency of the word  $\times$  Importance of the word

Basics of Natural Language Processing

# TF-IDF

Term Frequency — Inverse Document Frequency

Frequency of the word × Importance of the word

$$\frac{\text{no. of times it appears}}{\text{total no. of words}} \times \log_{10} \frac{\text{total no. of docs}}{\text{no. of docs in which word appears}}$$

# TF-IDF

Term Frequency — Inverse Document Frequency

Frequency of the word × Importance of the word

$$\frac{\text{no. of times it appears}}{\text{total no. of words}} \times \log_{10} \frac{\text{total no. of docs}}{\text{no. of docs in which word appears}}$$

Examples:

$$\text{TF}(\text{"the"}) = 0.1$$

**Q1: What does it mean?**

**Q2: If "the" appears in all 10 documents, what is IDF("the")?**

# TF-IDF

Term Frequency — Inverse Document Frequency

Frequency of the word × Importance of the word

$$\frac{\text{no. of times it appears}}{\text{total no. of words}} \times \log_{10} \frac{\text{total no. of docs}}{\text{no. of docs in which word appears}}$$

Examples:

$$\text{TF}(\text{"the"}) = 0.1$$

$$\text{IDF}(\text{"the"}) = 0$$

$$\rightarrow \text{TF-IDF}(\text{"the"}) = 0$$

# TF-IDF

Term Frequency — Inverse Document Frequency

Frequency of the word × Importance of the word

$$\frac{\text{no. of times it appears}}{\text{total no. of words}} \times \log_{10} \frac{\text{total no. of docs}}{\text{no. of docs in which word appears}}$$

Examples:

$$\text{TF}(\text{"NLP"}) = 0.1$$

**Q: If "NLP" appears in only one document, what is IDF("NLP")?**

# TF-IDF

Term Frequency — Inverse Document Frequency

Frequency of the word × Importance of the word

$$\frac{\text{no. of times it appears}}{\text{total no. of words}} \times \log_{10} \frac{\text{total no. of docs}}{\text{no. of docs in which word appears}}$$

Examples:

$$\text{TF}(\text{"NLP"}) = 0.1$$

$$\text{IDF}(\text{"NLP"}) = 1$$

$$\rightarrow \text{TF-IDF}(\text{"NLP"}) = 0.1$$



# TF-IDF

**D1:** Jim and Pam traveled by bus.

**D2:** The train was late.

**D3:** The flight was full. Traveling by flight is expensive.

	and	bus	by	expensive	flight	full	is	jim	late	pam	the	train	travel	was
<b>D1</b>	0.4	0.4	0.3	0	0	0	0	0.4	0	0.4	0	0	0.3	0
<b>D2</b>	0	0	0	0	0	0	0	0	0.5	0	0.4	0.5	0	0.4
<b>D3</b>	0	0	0.2	0.3	0.6	0.3	0.3	0	0	0	0.2	0	0.2	0.2

# TF-IDF vs. BoW

BoW

	and	bus	by	expensive	flight	full	is	jim	late	pam	the	train	travel	was
<b>S1</b>	1	1	1	0	0	0	0	1	0	1	0	0	1	0
<b>S2</b>	0	0	0	0	0	0	0	0	1	0	1	1	0	1
<b>S3</b>	0	0	1	1	2	1	1	0	0	0	1	0	1	1

TF-IDF

	and	bus	by	expensive	flight	full	is	jim	late	pam	the	train	travel	was
<b>D1</b>	0.4	0.4	0.3	0	0	0	0	0.4	0	0.4	0	0	0.3	0
<b>D2</b>	0	0	0	0	0	0	0	0	0.5	0	0.4	0.5	0	0.4
<b>D3</b>	0	0	0.2	0.3	0.6	0.3	0.3	0	0	0	0.2	0	0.2	0.2

# TF-IDF

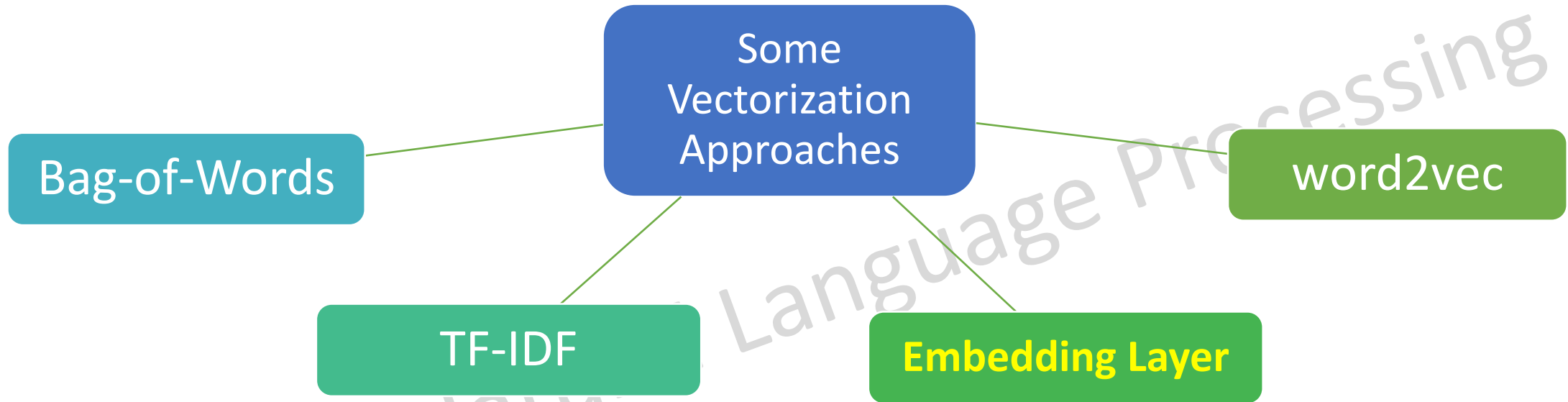
## Pros

- Simple
- Fast
- Considers the importance of words

## Cons

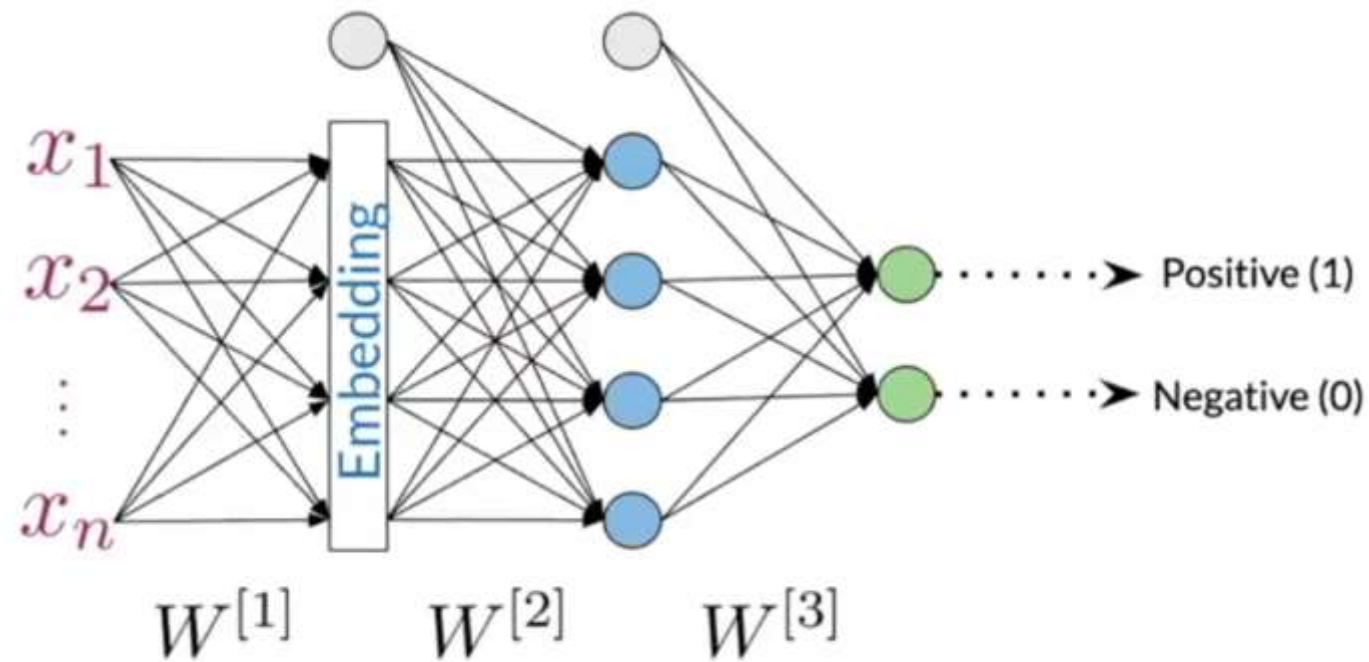
- Sparse
- High dimension
- No word relation (context)

# Vectorization



# Embedding Layer

## Example: ANN for Sentiment Analysis



$x = [x_1, x_2, \dots, x_n]$ : tweets

Example: "This movie was almost good"

# Embedding Layer

## Input Representation

Word	Number
a	1
able	2
about	3
...	...
movie	680
...	...
zebra	1000

# Embedding Layer

## Input Representation

Word	Number
a	1
able	2
about	3
...	...
movie	680
...	...
zebra	1000

$\mathbf{x} = [x_1, x_2, \dots, x_n]$ : tweets

Example: “This movie was almost good”

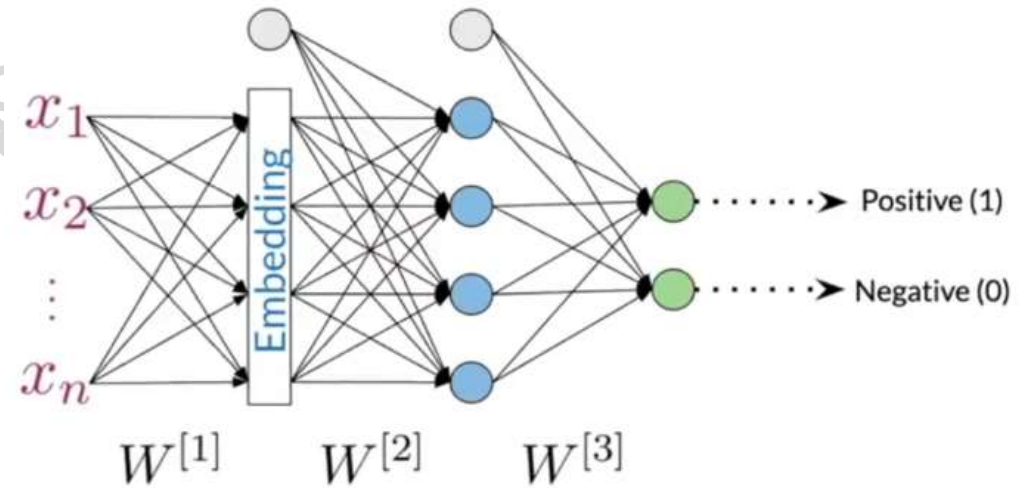
$\mathbf{x} = [700, 680, 720, 20, 55]$



In the ANN for sentiment analysis, how can we handle variable-length tweets?

Tweet 1: “this movie was almost good”

Tweet 2: “it was terrible”





# Embedding Layer

## Input Representation

Word	Number
a	1
able	2
about	3
...	...
movie	680
...	...
zebra	1000

$x = [x_1, x_2, \dots, x_n]$ : tweets

Example: "This movie was almost good"

$x = [700, 680, 720, 20, 55]$



Zero-padding to match size of largest tweet

$x = [700, 680, 720, 20, 55, 0, 0, 0, 0, 0, 0, 0, 0, 0]$

# Embedding Layer

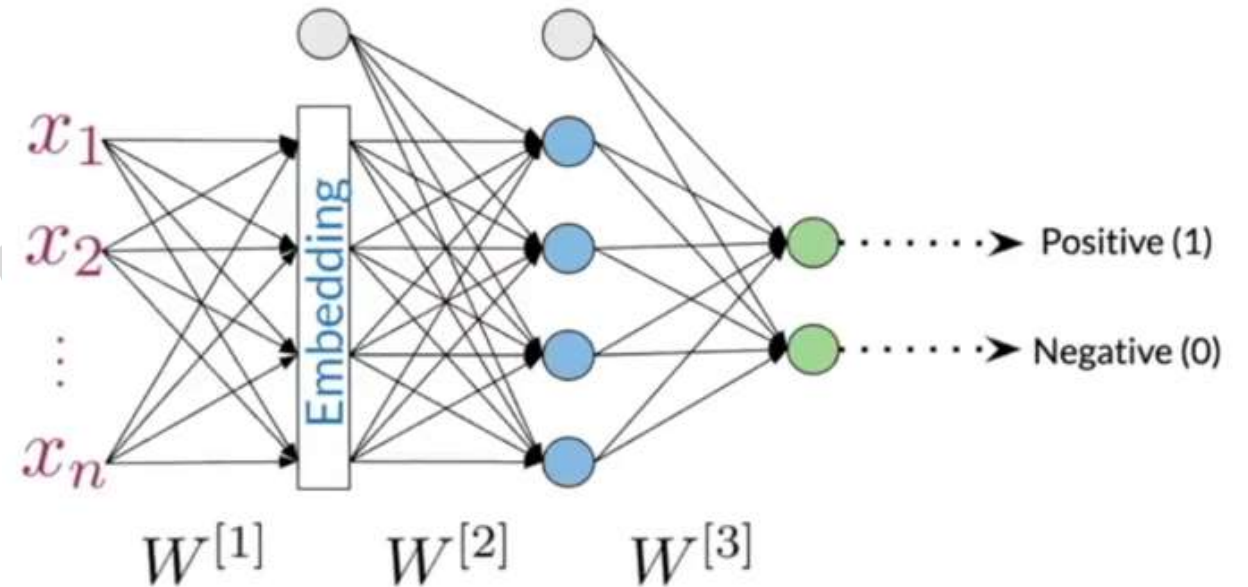
**ANN learns the 4-D embedding representation of sentences**

Example: "This movie was almost good"

$x = [700, 680, 720, 20, 55, 0, 0, 0, 0, 0, 0, 0, 0, 0]$



$a^{[1]} = [0.4 \ 0.01 \ 0.3 \ 0.21]$



# Embedding Layer

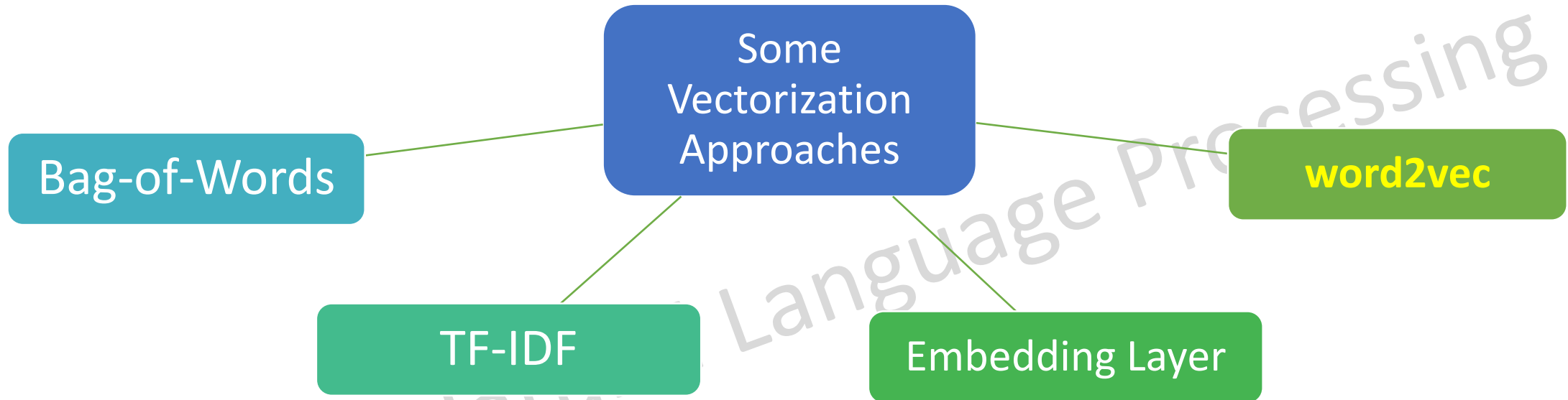
## Pros

- Simple
- Fast
- Considers the importance of words
- Condense and low dimension

## Cons

- No word relation
- Meaningless math operation between words

# Vectorization



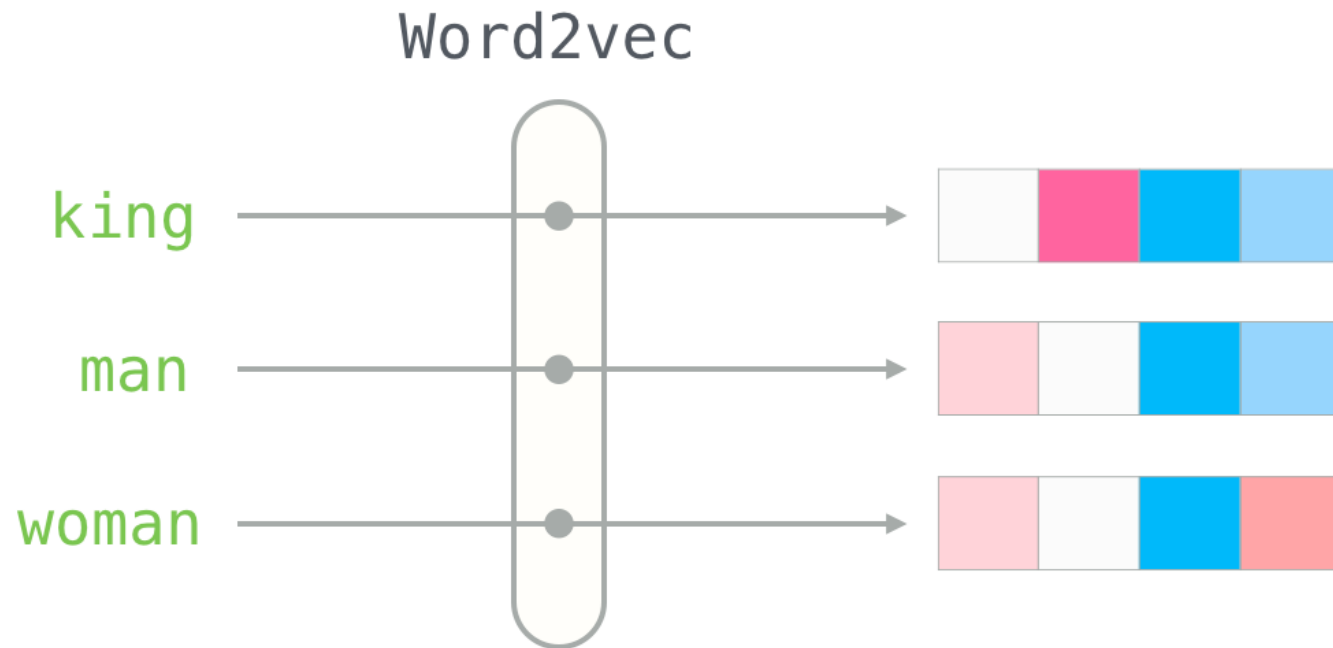
**Word2vec** is a model provided by Google in 2013  
for the **effective word embedding**.

Basics of Natural Language Processing

[Mikolov, Tomas, et al. "Distributed representations of words and phrases and their compositionality." *Advances in neural information processing systems*. 2013]

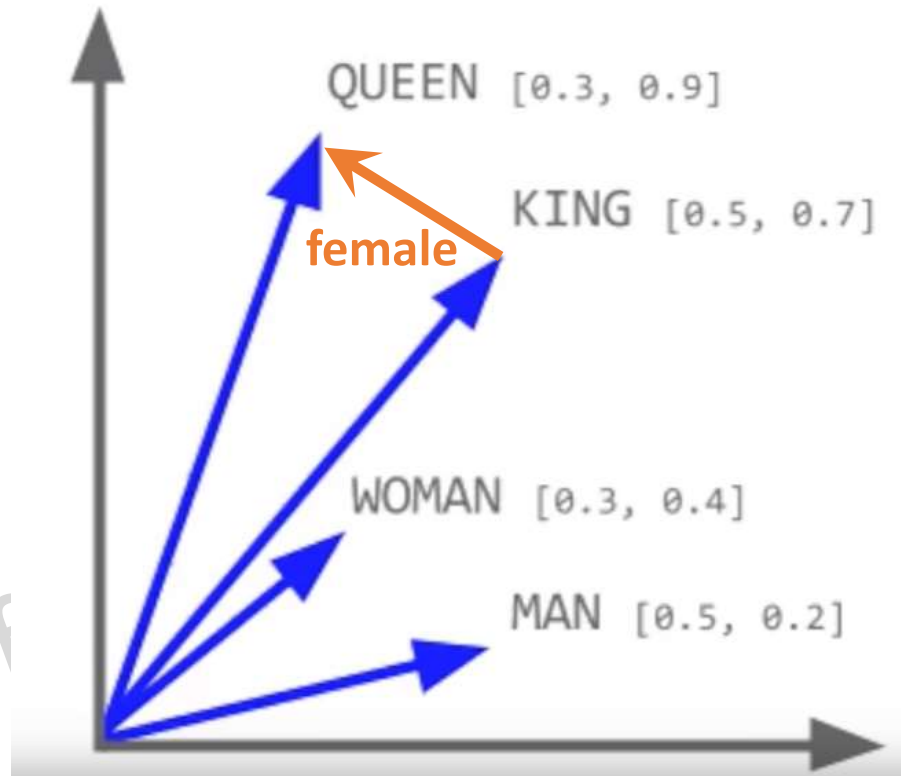
# word2vec

**Word2vec** is a model provided by Google in 2013  
for the **effective word embedding**.



[Mikolov, Tomas, et al. "Distributed representations of words and phrases and their compositionality." *Advances in neural information processing systems*. 2013]

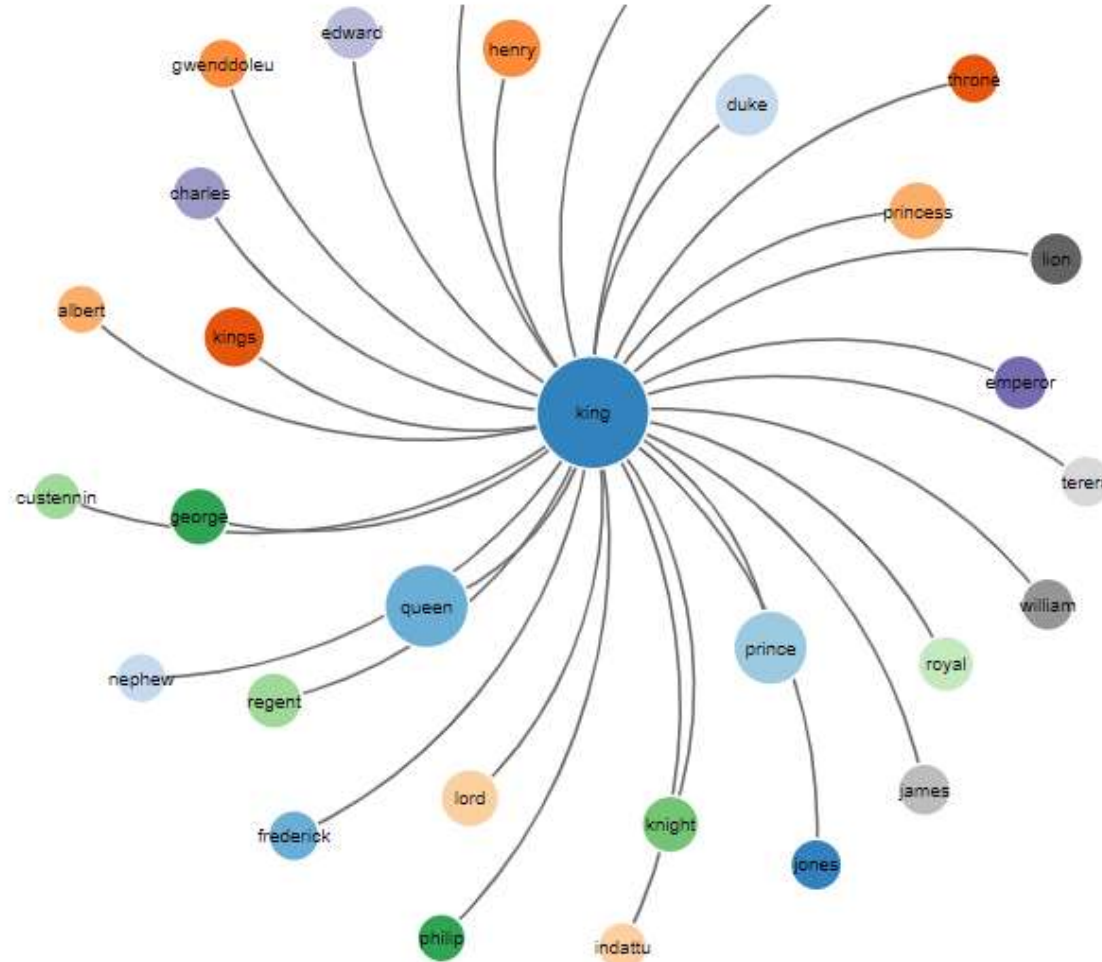
# Word Embedding



King + female = Queen  
Man + female = Woman  
Queen – royal = Woman

# Word Similarity

Top 30 analogous words or synonyms for king



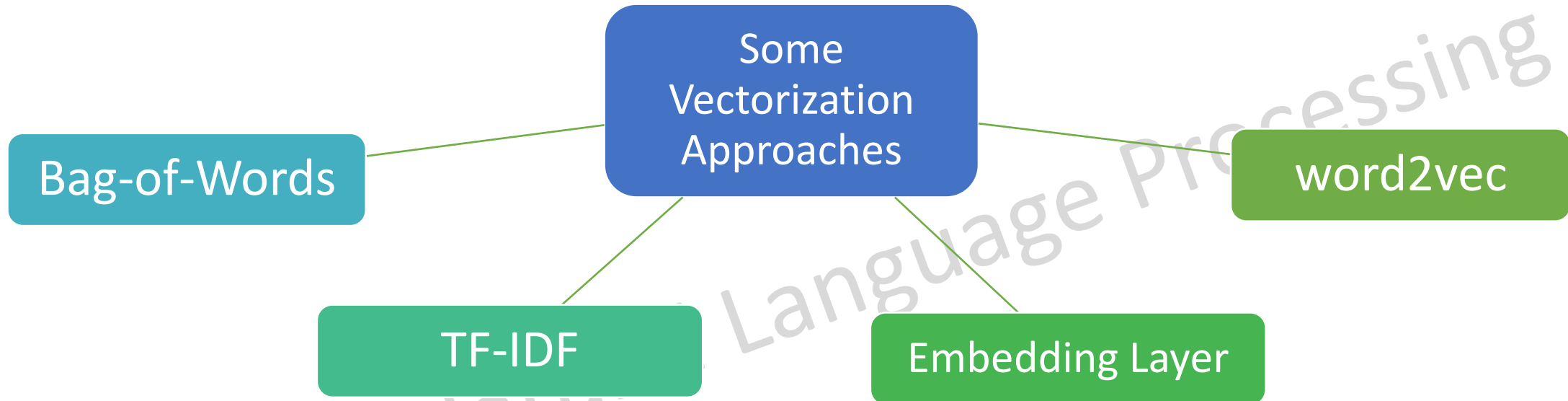
Top 10 similar words or synonyms for king

queen	0.745687
prince	0.659387
duke	0.574773
kings	0.544352
henry	0.532862
princess	0.524547
lord	0.518185
george	0.512252
knight	0.505063
regent	0.499793

<https://wordsimilarity.com/>

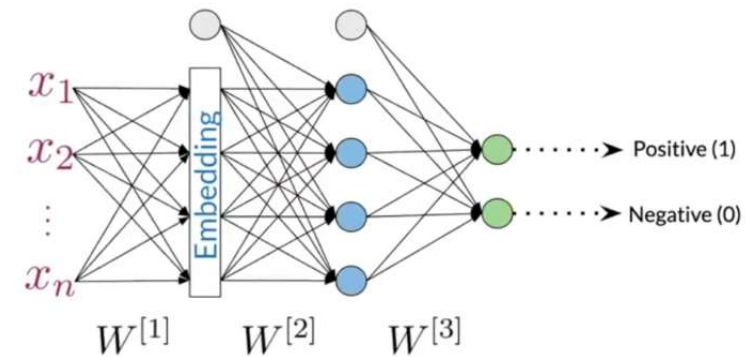


# Vectorization: Wrap up



	and	bus	by	expensive	flight	full	is	jim	late	pam	the	train	travel	was
S1	1	1	1	0	0	0	0	1	0	1	0	0	1	0
S2	0	0	0	0	0	0	0	0	1	0	1	1	0	1
S3	0	0	1	1	2	1	1	0	0	0	1	0	1	1

	and	bus	by	expensive	flight	full	is	jim	late	pam	the	train	travel	was
D1	0.4	0.4	0.3	0	0	0	0	0.4	0	0.4	0	0	0.3	0
D2	0	0	0	0	0	0	0	0	0.5	0	0.4	0.5	0	0.4
D3	0	0	0.2	0.3	0.6	0.3	0.3	0	0	0	0.2	0	0.2	0.2



# Wrap up

Session 1: Introduction

**Session 2. Basics of Linguistics**

- Components
- Challenges
- Vectorization

**Session 3. Basics of ML**

**Session 4 (Lab). Effective Word Representation by python**

# Digikala Academy NLP Events

## Basic

- **Session 1.** Introduction
- **Session 2.** Basics of Linguistics
- **Session 3.** Basics of ML
- **Session 4 (Lab).** Effective Word Representation by python

## Intermediate

- TBA

## Advanced

- TBA

# References of this session



## Natural Language Processing (NLP) with Python — Tutorial

<https://medium.com/towards-artificial-intelligence/natural-language-processing-nlp-with-python-tutorial-for-beginners-1f54e610a1a0#2847>