

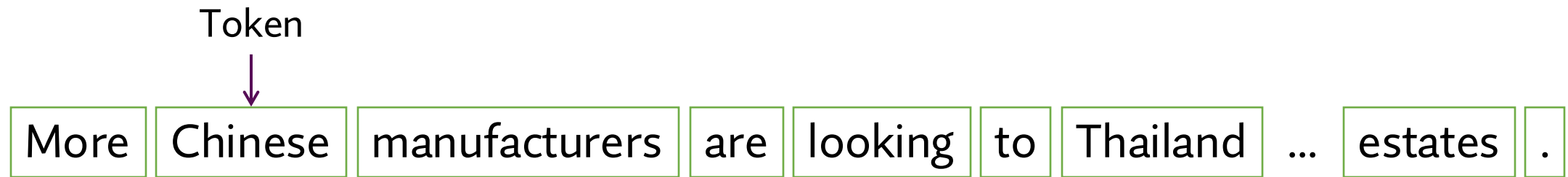
Feature Engineering

Papangkorn Inkeaw, Ph.D.

Feature Extraction

Chapter 5 (Part II) - Feature Extraction for Text Data

Text – Basic Knowledge



Text is a human-readable sequence of word(s).

Text – Basic Knowledge

n-gram A contiguous sequence of n tokens from a given sample of text.

More Chinese manufacturers are looking to Thailand ... estates .

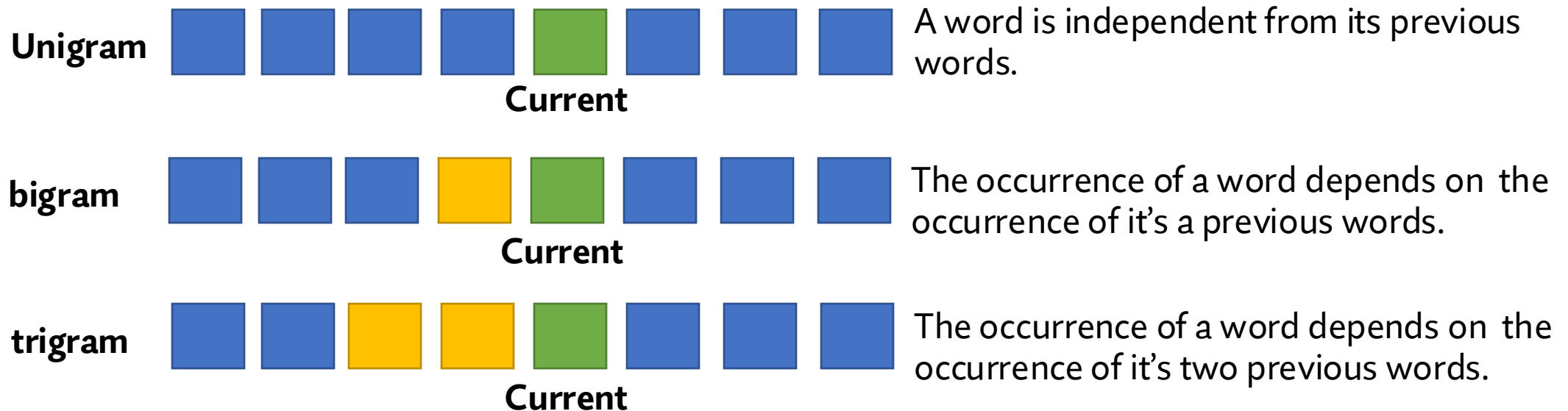
Unigram More Chinese manufacturers are looking to Thailand ...

Bigram More Chinese Chinese manufacturers manufacturers are
are looking looking to to Thailand ...

Trigram More Chinese manufacturers Chinese manufacturers are
manufacturers are looking are looking to looking to Thailand ...

Text – Basic Knowledge

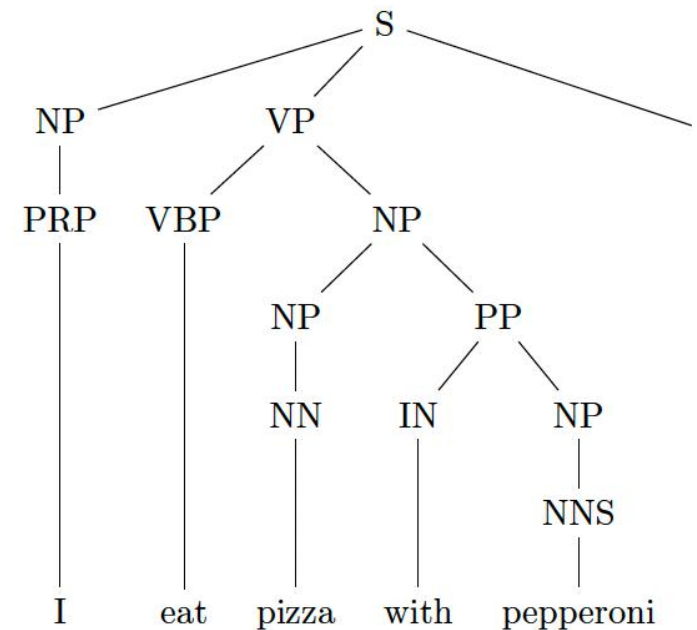
n-gram model predicts the occurrence of a word based on the occurrence of its $N - 1$ previous words.



Text – Basic Knowledge

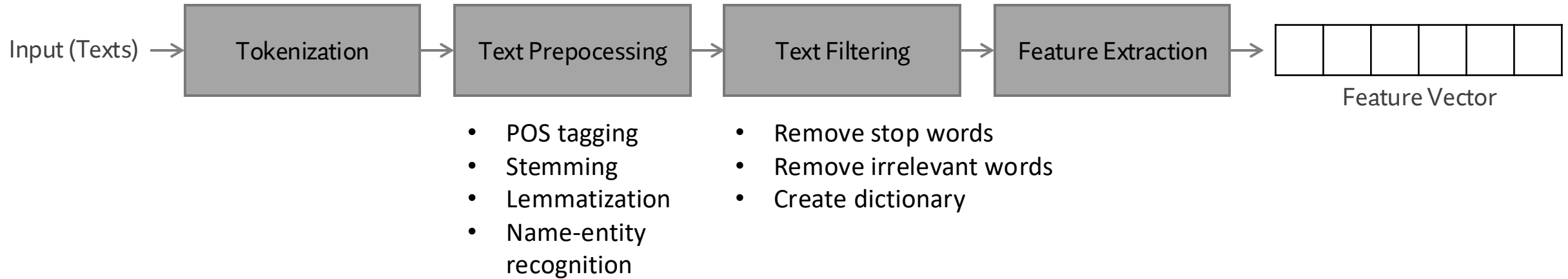
Part-of-speech (POS) The grammatical function of each word in a sentence

We can use NLP tool to categorizing words in a text in correspondence with a particular part of speech, depending on the definition of the word and its context.



From Text to Feature Vector

General Feature Extraction Process



One-hot Encoding

- A representation of categorical variables as binary vectors.
- Each word is represented as a binary vector that is:
 - All zero values
 - Except the index of the word, which is marked with a 1.

	All possible words
	. It a cat is
Word	It = [0., 1., 0., 0., 0.],
	is = [0., 0., 0., 0., 1.],
	a = [0., 0., 1., 0., 0.],
	cat = [0., 0., 0., 1., 0.],
	. = [1., 0., 0., 0., 0.]

Bag of Words

- What about full texts instead of single words?
- The vector representation of a text is simply the vector sum of all the words it contains:

	All possible words
	. It a cat is
Word	It = [0., 1., 0., 0., 0.],
	is = [0., 0., 0., 0., 1.],
	a = [0., 0., 1., 0., 0.],
	cat = [0., 0., 0., 1., 0.],
	. = [1., 0., 0., 0., 0.]
	[1., 1., 1., 1., 1.] ←

Vector sum of all words represents the text "It is a cat."

Bag of Words

- In practice, it's much more convenient to use a dictionary instead of an actual vector.
- This is known as a **bag-of-words**, and word order is discarded.

		Dictionary			
		cat	dog	bird	panda
Word	They =	[0.,	0.,	0.,	0.],
	are =	[0.,	0.,	0.,	0.],
	cat =	[1.,	0.,	0.,	0.],
	and =	[0.,	0.,	0.,	0.],
	dog =	[0.,	1.,	0.,	0.]

They are cat and dog = [1., 1., 0., 0.]

← **Bag of words**
represents the text
"They are cat and dog"

TF-IDF

- The **term frequency** $tf(t,d)$ of term t in document d is defined as the relative frequency of times that t occurs in d .

$$tf(t, d) = \frac{f_{t,d}}{\sum_{t' \in d} f_{t',d}}$$

- Logarithmically scaled frequency is defined as $tf(t, d) = \log(1 + f_{t,d})$
- Document frequency: Rare terms are more informative than frequent terms
- The **inverse document frequency** is a measure of how much information the word provides.
- The $idf(t,D)$ is the ratio of the total number of documents D to the number of documents containing the term t .

$$idf(t, D) = \frac{N}{n_t}$$

TF-IDF

- Then, the Term frequency–inverse document frequency $\text{tfidf}(t,d,D)$ is calculated as

$$\text{tfidf}(t, d, D) = \text{tf}(t, d) \times \text{idf}(t, D)$$

- A high weight in tf–idf is reached by a high term frequency (in the given document).
- A low document frequency of the term in the whole collection of documents; the weights hence tend to filter out common terms.

Word Embeddings

“ Words with similar meanings should occur in similar contexts. ”

- The distributional hypothesis in linguistics

- From a word we can get some idea about the context where it might appear.

__ bird __ __

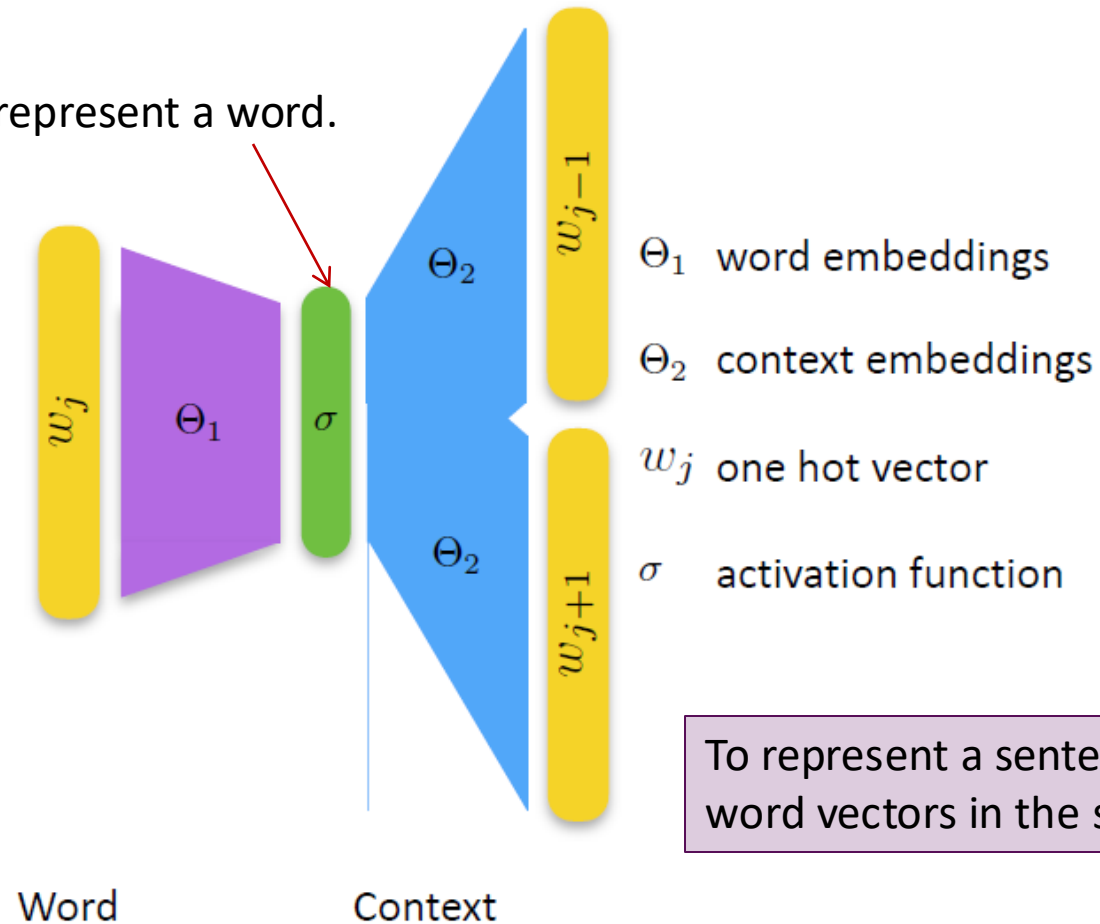
- From the context, we have some idea about possible words.

The red __ is nice

Word Embeddings

Skipgram

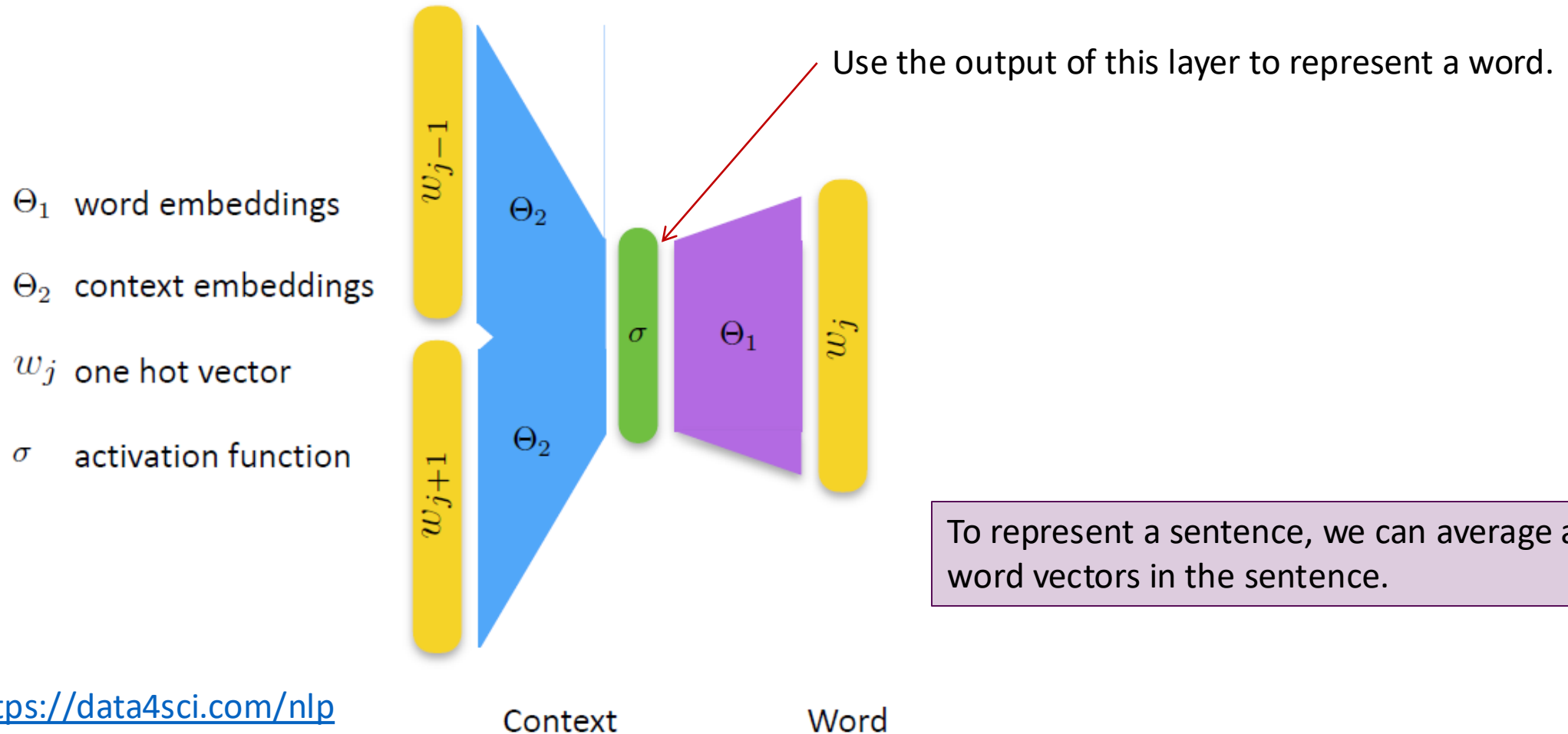
Use the output of this layer to represent a word.



To represent a sentence, we can average all word vectors in the sentence.

Word Embeddings

Continuous Bag of Words



References & Study Resources

- Guozhu Dong and Huan Liu. (2020). *Feature Engineering for Machine Learning and Data Analytics*. CRC Press.
- <https://data4sci.com/nlp>