



Samsung Innovation Campus

| Artificial Intelligence Course

Together for Tomorrow!
Enabling People

Education for Future Generations

Chapter 7.

Natural Language Processing and Language Models for Text Mining

AI Course

Chapter Description

◆ Chapter objectives

- ✓ Process input text from sources of various text formats to extract high-quality information.
- ✓ Structure language and derive patterns by natural language processing. And evaluate and analyze these results to utilize them for real-world applications.

◆ Chapter contents

- ✓ Unit 1. Text Mining
- ✓ Unit 2. Text Preprocessing
- ✓ Unit 3. Language Model
- ✓ Unit 4. Natural Language Processing with Keras

Unit 1.

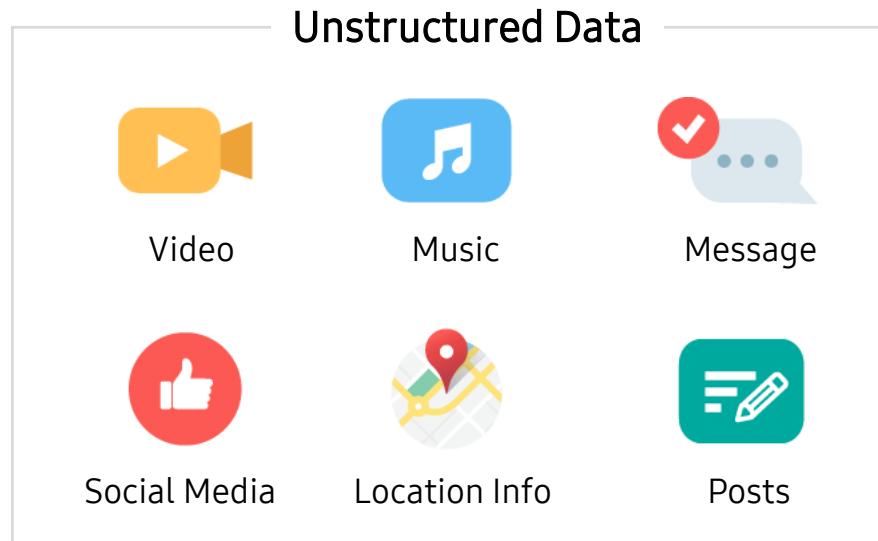
Text Mining

- | 1.1. What is Text Mining?
- | 1.2. Data Collection
- | 1.3. String Manipulation
- | 1.4. Natural Language Processing (NLP)
- | 1.5. Sequential Data
- | 1.6. Corpus

Text Mining

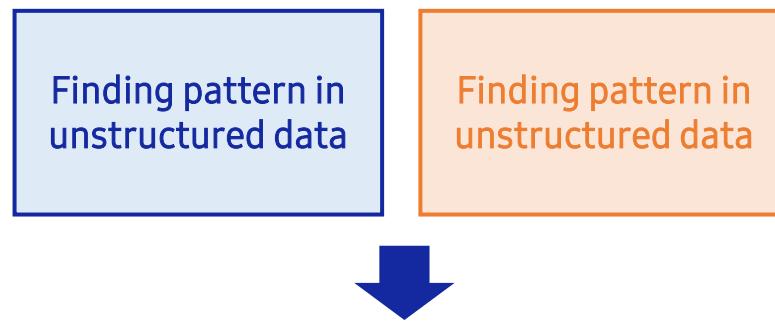
| What is unstructured data?

- ▶ Data that is not yet structured. It does not have a defined data model (structure).
- ▶ Documents, videos, or audio with a large amount of data but with varying structures and forms.
- ▶ Books, journals, documents, metadata, health records, audio, video, analog data, images, files, e-mail messages, webpages, and word-processor documents are all composed of unstructured texts.



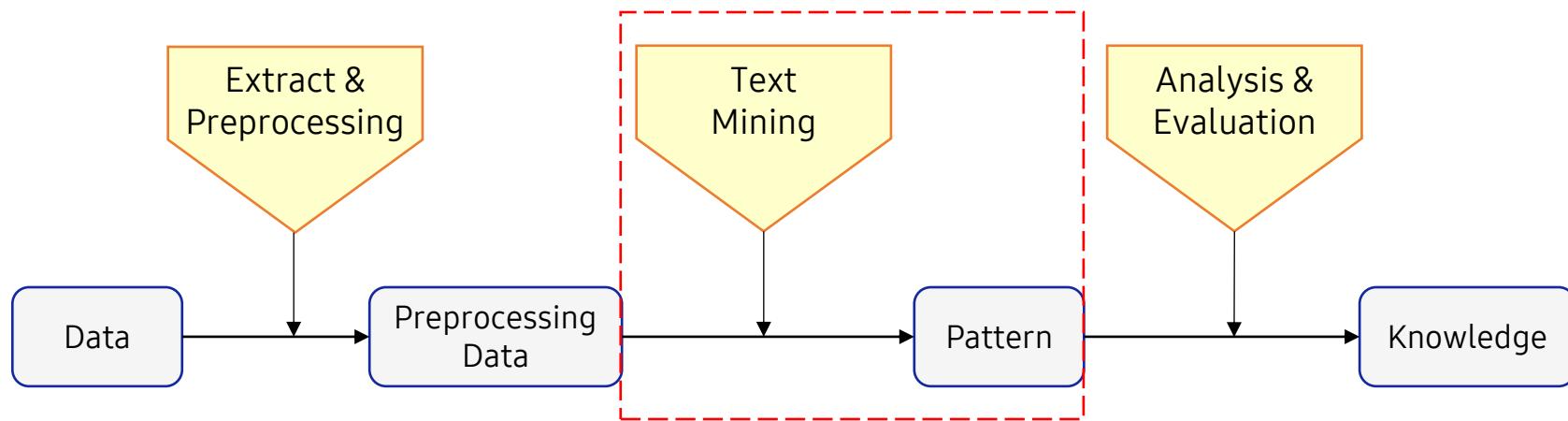
Analyzing unstructured data

- ▶ Manually enter tags into metadata to structure texts.
- ▶ Skilled data structuring based on text-mining uses a method that creates a tag so that a word in the text and a part in the speech correspond.
- ▶ Uses software that builds structures processable by machines.
- ▶ Analysis that enables semantic deduction from texts, syntax, and other small or large patterns. This analysis uses algorithms that inspect all internal structures of human communication in word units by forming them into a linguistic, auditory, and visual structure.



What is Text Mining?

- ▶ Text mining or text analytics is a technology that extracts useful information from unstructured text data.
- ▶ To be more specific, it is finding practical patterns from a large amount of document data by applying mechanical algorithms and statistical techniques.



Text mining vs. data mining

- ▶ Data mining extracts useful and valuable patterns from structured data.
- ▶ In contrast, text mining extracts named entities, patterns, or information on word-sentence relationships from unstructured data composed of natural language.

Text Mining Application Example in Real World

- ▶ **Used for marketing:** Corporations collect and analyze posts on Twitter that mention their brand names or messages with the customers.
 - They examine if users mention certain topics at a certain time; if users have positive or negative connotations; how keywords change within time; if customers' keywords changed before and after a campaign; if a promotion generated word of mouth; if a certain group of customers reacts; and etc.
 - With this information, corporations can closely monitor marketing activities and control reputation management and build competitive strategies against competing brands based on feedback monitoring.
- ▶ **Supporting data for various industries:** It can support data from fields such as politics, environment, medicine, and business areas such as manufacturing, facilities, and marketing.
- ▶ **Factories predicting breakdown:** Factories can predict facilities breakdown from documents recorded by facilities maintenance workers.
- ▶ **Checking product reviews:** Customers can access information about a product's performance or issues from product analysis or reviews.
- ▶ **News analysis:** It can show popular issues within time flow and discover experts of a certain field by analyzing topics of news reports or speech statements.

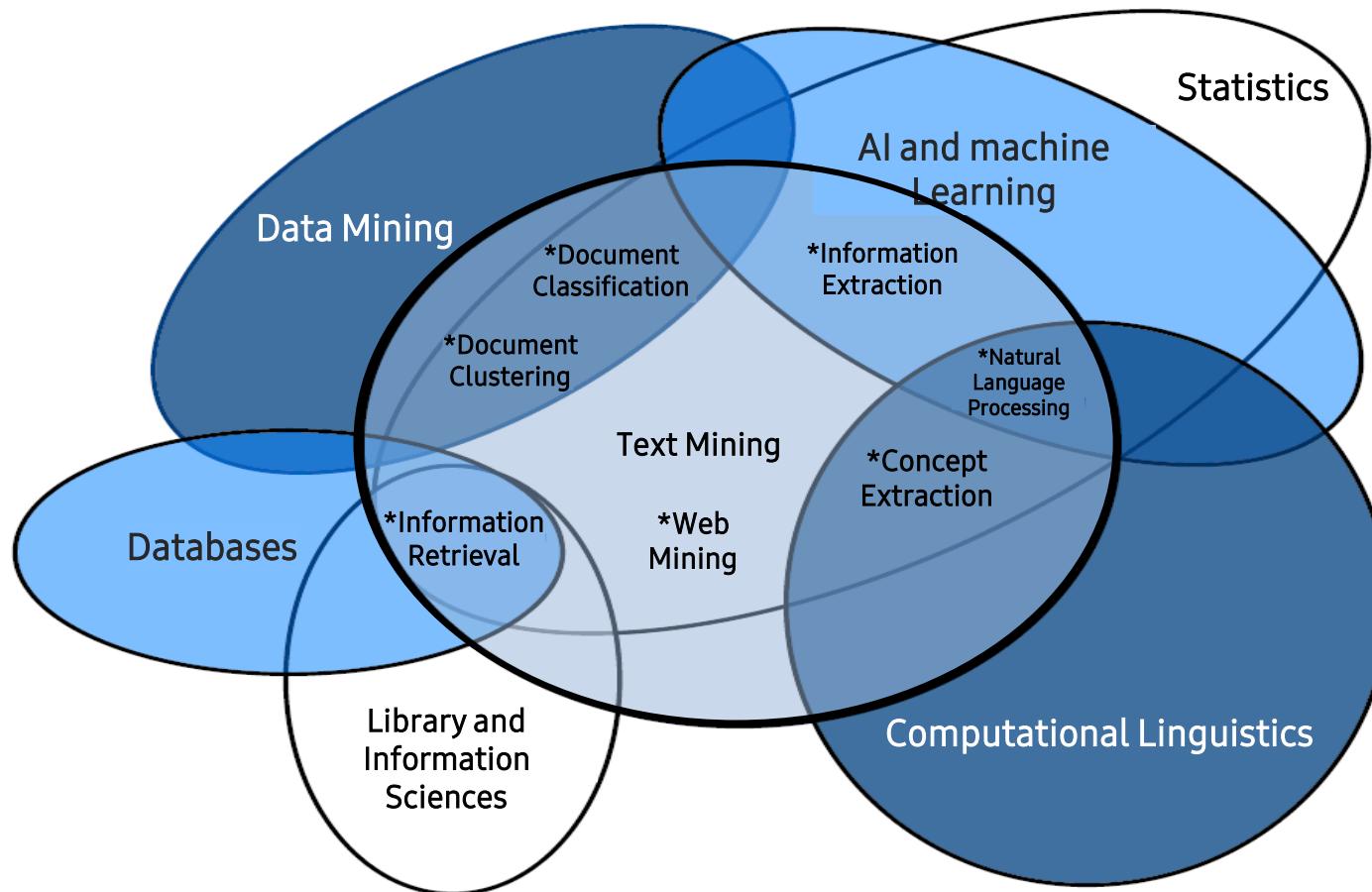
Text Mining vs. Natural Language Processing (NLP) (1/3)

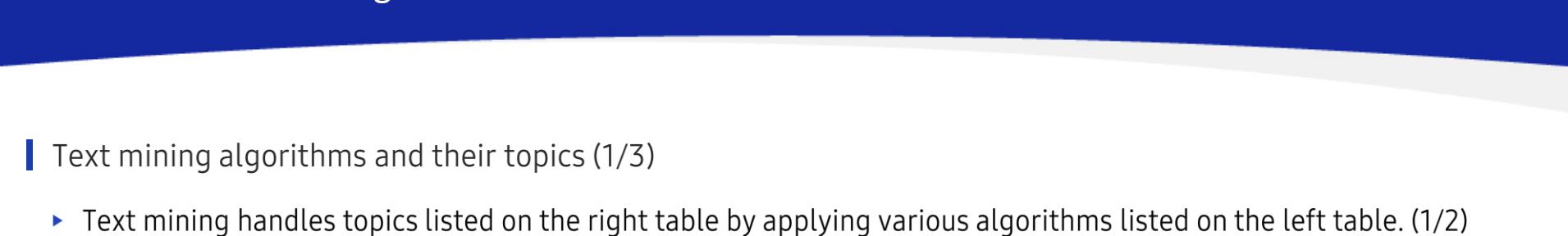
- ▶ **NLP is a field of AI that handles communication.** It enables machines to generate and analyze human language (generate and understand natural language)
 - NLP can process many voice types, including slang, dialects, and grammatical errors. Machine learning constructs the foundation for this methodology. Applications of NLP are search engines, AI chatbots, grammar correction apps, etc.
- ▶ On the other hand, **text mining is a sub-category of data mining science.** Data mining science includes data search, data mining, and machine learning methods.
- ▶ More than 80% of organizations worldwide utilize textual information. **NLP recognizes text and voice, while text mining evaluates the quality of a text.**

Text Mining vs. Natural Language Processing (NLP) (2/3)

- ▶ Different tools are used for text mining and NLP.
 - To construct a high-quality NLP system, you must be proficient in neural networks, deep learning, and NLTK.
 - Text mining system is a technique like Levenshtein Distance, Cosine similarity, or Feature Hashing.
 - You must be familiar with text-processing programming languages and statistical models like Perl or Python.
- ▶ **NLP offers the understanding** of explained emotions and grammatical structure and detects the intent behind a text.
 - This assists fluent translation of a text to another language.
- ▶ Meanwhile, **text mining discovers the relationship** among words in the text.
 - It analyzes the frequency of word use and patterns.

| Text Mining vs. Natural Language Processing (NLP) (3/3)

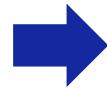




| Text mining algorithms and their topics (2/3)

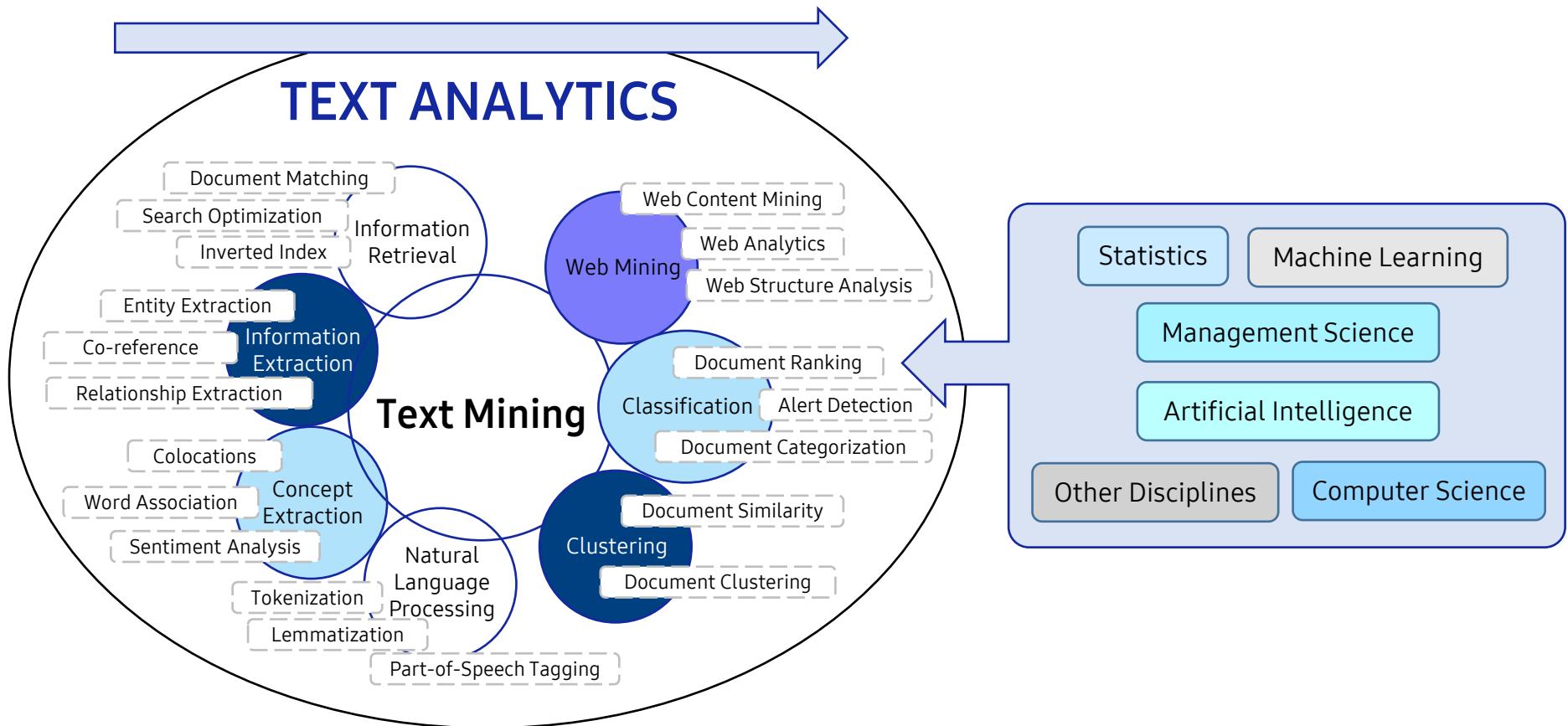
- ▶ Text mining handles topics listed on the right table by applying various algorithms listed on the left table. (2/2)

Algorithm	Area
Neural network	Document classification
Support vector machines	Document classification
MARSplines	Document classification
Link analysis	Concept extraction
k-nearest neighbors	Document classification
Word clustering	Concept extraction
Regression	Classification

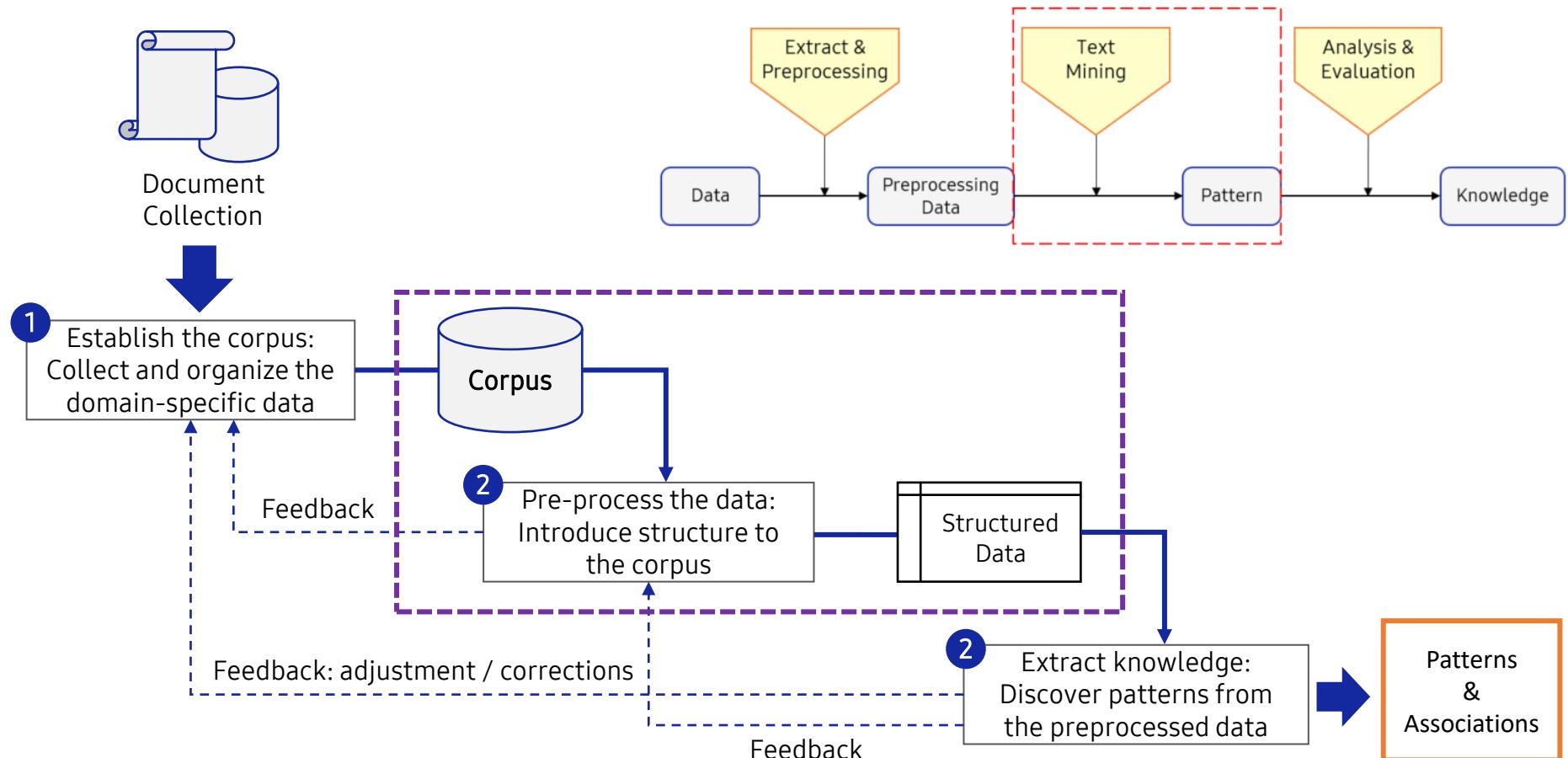


Topic	Practice Area
Web crawling	Web mining
Link analytics	Web mining
Entity extraction	Information extraction
Link extraction	Information extraction
Part of speech tagging	Natural language processing
Tokenization	Natural language processing
Question answering	Natural language processing Search and information retrieval
Topic modeling	Concept extraction
Synonym identification	Concept extraction

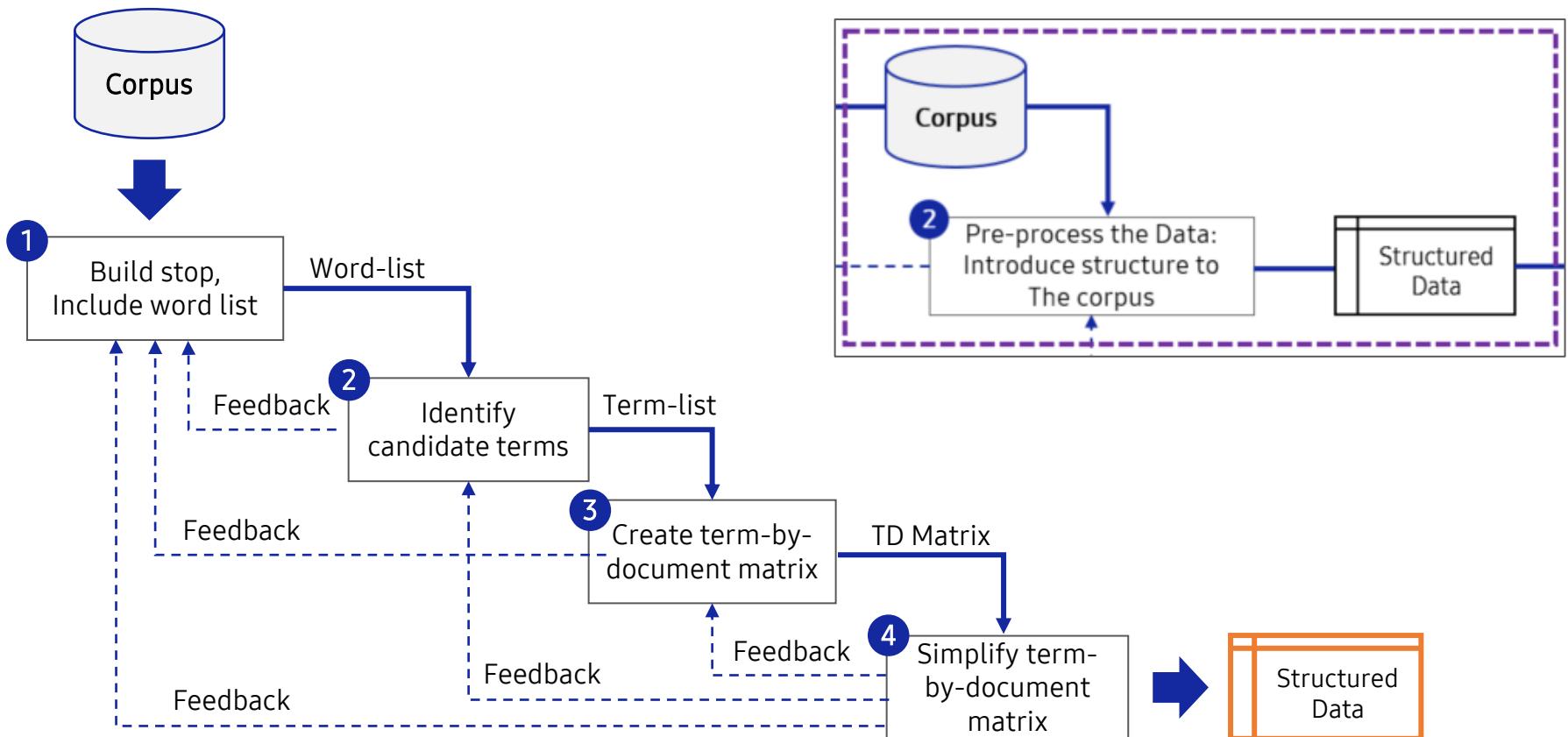
| Text mining algorithms and their topics (3/3)



I Basic procedure of text mining



Basic procedure of text mining



Unit 1.

Text Mining

- | 1.1. What is Text Mining?
- | **1.2. Data Collection**
- | 1.3. String Manipulation
- | 1.4. Natural Language Processing (NLP)
- | 1.5. Sequential Data
- | 1.6. Corpus

Collecting Data from Different Resources

- Collecting text is a process of establishing the plan for collecting and gathering data suitable for the task's objective and characteristics.
 - ▶ Collecting text is an important process that determines the quality of text analysis services.
 - ▶ You make a detailed plan after reviewing the time period, cost, the possibility of personal information infringement, and inclusion of data categories meeting the objective. According to the plan, you perform a pre-test and proceed with data collection from different resources.

Collecting Data from Different Resources

- Collecting text is a process of establishing the plan for collecting and gathering data suitable for the task's objective and characteristics.
- ▶ Various collecting techniques are used according to data type and features. Main techniques are listed below.

Technique	Features	Data Type
Crawling	<ul style="list-style-type: none">- Collects web documents and information on the web, such as social media, news, and web information.- Follows URL link and repetitively collects.	Web document
Scraping	<ul style="list-style-type: none">- Collects information from a single website (or document).	Web document
FTP	<ul style="list-style-type: none">- Transmits and receives files from internet servers using TCP/IP protocol.- Considers using SFTP for reinforced security.- Considers constructing an exclusive network for linked servers.	File
Open API	<ul style="list-style-type: none">- Offers data collecting method with an open API that allows easy access to service, information, and data.	Real-time data
RSS	<ul style="list-style-type: none">- XML-based content distribution protocol that allows sharing of up-to-date web-based information.	Content

Text Data from Websites

| Download a webpage with the Requests library

- ▶ This is okay when the exact URL is known.
- ▶ If log-in is required, use the Selenium library instead.
- ▶ HTML content without parsing.

Ex

```
In [1]: import requests as rq
res = rq.get("https://en.wikipedia.org/wiki/Machine_learning")
print(res.status_code)
print(res.text)
```



Line 1-3

- If the status code is 200, then OK.

Parsing HTML with the BeautifulSoup4 library

- ▶ The downloaded HTML should be parsed to access the desired content.

```
Ex In [1]: import requests, bs4
res = requests.get("https://en.wikipedia.org/wiki/Machine_learning")
soup = bs4.BeautifulSoup(res.text, 'html.parser')
x = soup.find_all('p')
text = ''
for i in range(len(x)):
    text += x[i].text.strip() + '\n'
print(text)
```

Line 1-3

- Return a BeautifulSoup object.

Parsing HTML with the BeautifulSoup4 library

- ▶ The downloaded HTML should be parsed to access the desired content.

```
Ex In [1]: import requests, bs4
res = requests.get("https://en.wikipedia.org/wiki/Machine_learning")
soup = bs4.BeautifulSoup(res.text, 'html.parser')
x = soup.find_all('p')
text = ''
for i in range(len(x)):
    text += x[i].text.strip() + '\n'
print(text)
```

Line 1-4

- Get all the paragraphs.

Parsing HTML with the BeautifulSoup4 library

- ▶ The downloaded HTML should be parsed to access the desired content.

```
Ex In [1]: import requests, bs4
res = requests.get("https://en.wikipedia.org/wiki/Machine_learning")
soup = bs4.BeautifulSoup(res.text, 'html.parser')
x = soup.find_all('p')
text = ''
for i in range(len(x)):
    text += x[i].text.strip() + '\n'
print(text)
```

Line 1-7

- Join all the text contents.

Parsing HTML with the BeautifulSoup4 library

- ▶ The downloaded HTML should be parsed to access the desired content.

Ex The example from the previous slide produces an output as shown below.

Machine learning (ML) is the scientific study of algorithms and statistical models that computer systems use to perform a specific task without using explicit instructions, relying on patterns and inference instead. It is seen as a subset of artificial intelligence. Machine learning algorithms build a mathematical model based on sample data, known as "training data," in order to make predictions or decisions without being explicitly programmed to perform the task.[1][2]:2 Machine learning algorithms are used in a wide variety of applications, such as email filtering and computer vision, where it is difficult or infeasible to develop a conventional algorithm for effectively performing the task.

Machine learning is closely related to computational statistics, which focuses on making predictions using computers. The study of mathematical optimization delivers methods, theory, and application domains to the field of machine learning. Data mining is a field of study within machine learning and focuses on exploratory data analysis through unsupervised learning.[3][4] In its application across business problems, machine learning is also referred to as predictive analytics.

Unit 1.

Text Mining

- | 1.1. What is Text Mining?
- | 1.2. Data Acquisition
- | **1.3. String Manipulation**
- | 1.4. Natural Language Processing (NLP)
- | 1.5. Sequential Data
- | 1.6. Corpus

String Manipulation

Useful functions and methods for string manipulation

Function	Explanation
<code>x.lstrip()</code>	Remove space on the left.
<code>x.rstrip()</code>	Remove space on the right.
<code>x.strip()</code>	Remove space on both sides.
<code>x.replace(str1, str2)</code>	Replace the substring <code>str1</code> by <code>str2</code> .
<code>x.count(str)</code>	Number of occurrences of <code>str</code> in <code>x</code> .
<code>x.find(str)</code>	Find the sub-string <code>str</code> . Returns -1 if not found.
<code>x.index(str)</code>	Find the sub-string <code>str</code> . Throws an error if not found.
<code>y.join(str_list)</code>	Concatenate the elements of <code>str_list</code> using <code>y</code> as separator.
<code>x.split(y)</code>	Break up a string using <code>y</code> as separator.
<code>x.upper()</code>	Convert <code>x</code> into uppercase.
<code>x.lower()</code>	Convert <code>x</code> into lowercase.
<code>len(x)</code>	Returns the string length.

| Useful functions and methods for string manipulation

- ▶ Used for making string patterns.
- ▶ Useful for recognizing and processing complex string patterns ⇒ pre-processing of text data.
- ▶ More powerful than the combinations of the usual string functions or methods.
- ▶ Supported by many languages, including Python.

Regular Expression

| Metacharacters

- ▶ Characters with special meanings in regular expressions.
 . ^ \$ * + ? { } [] \ | ()
- ▶ Can be used to construct patterns.
- ▶ More information can be found at <https://docs.python.org/3/howto/regex.html>.

Metacharacters: []

- ▶ Can enclose a set of characters as a match pattern.
- ▶ Any character can be enclosed by the [].
- ▶ For example, “[abc]” matches with any pattern containing “a” or “b” or “c” character.

RegEx	String	Match?	Explanation
“[abc]”	“a”	Yes	“a” is in the string.
“[abc]”	“before”	Yes	“b” is in the string.
“[abc]”	“dude”	No	There is neither “a” nor “b” nor “c” in the string.

Metacharacters: []

- ▶ We can use a hyphen “-” to indicate a range of characters.

Ex “[0-5]” is the same as “[012345].”

Ex “[0-9]” means the entire set of number digits.

Ex “[a-d]” is the same as “[abcd].”

Ex “[a-zA-Z]” means the entire set of alphabet letters, both uppercase and lowercase.

Metacharacters: [^]

- ▶ Characters that are not in the enclosed set will be matched.
- ▶ “^” has to be the first character within the square brackets.

RegEx	String	Match?	Explanation
“[^abc]”	“a”	No	In the string, there is no other character than “a” or “b” or “c.”
“[^abc]”	“before”	Yes	There are characters other than “a” or “b” or “c” in the string.
“[^abc]”	“dude”	Yes	There are characters other than “a” or “b” or “c” in the string.

Metacharacters: [] and [^]

- ▶ There are shorthand expressions as following:

- “\w” is the same as “[a-zA-Z0-9_].”
- “\W” is the same as “[^a-zA-Z0-9_].”
- “\d” is the same as “[0-9].”
- “\D” is the same as “[^0-9].”
- “\s” means whitespace character.
- “\S” means non-whitespace character.

| Metacharacters: Dot .

- ▶ Dot matches with any character.
- ▶ “\.” is the dot as a character (not a metacharacter).

RegEx	String	Match?	Explanation
“a.b”	“aab”	Yes	“a” in the middle of the string matches with the dot.
“a.b”	“a0b”	Yes	“0” in the middle of the string matches with the dot.
“a.b”	“abc”	No	There is no character in between “a” and “b.”

| Metacharacters: *

- ▶ Pattern that repeats the preceding character any number of times (including 0).

RegEx	String	Match?	Explanation
"ca*t"	"ct"	Yes	"a" does not appear.
"ca*t"	"cat"	Yes	"a" appears once.
"ca*t"	"caaat"	Yes	"a" is repeated three times.

| Metacharacters: **+**

- ▶ Pattern that repeats the preceding character at least once or more times.

RegEx	String	Match?	Explanation
"ca +t "	"ct"	No	"a" does not appear.
"ca +t "	"cat"	Yes	"a" appears once.
"ca +t "	"caaat"	Yes	"a" is repeated three times.

Metacharacters: ?

- Pattern where the preceding character does not appear or appears just once.

RegEx	String	Match?	Explanation
"ca?t"	"ct"	Yes	"a" does not appear.
"ca?t"	"cat"	Yes	"a" appears once.
"ca?t"	"caat"	No	"a" is repeated twice (more than once).

Metacharacters: {m}

- Pattern where the preceding character is repeated m times.

RegEx	String	Match?	Explanation
"ca{2}t"	"ct"	No	"a" does not repeat twice.
"ca{2}t"	"cat"	No	"a" does not repeat twice.
"ca{2}t"	"caat"	Yes	"a" is repeated exactly twice.

| Metacharacters: {m, n}

- ▶ Pattern where the preceding character is repeated from m to n times.

RegEx	String	Match?	Explanation
"ca{2,5}t"	"cat"	No	"a" is repeated less than two times.
"ca{2,5}t"	"caat"	Yes	"a" is repeated three times.
"ca{2,5}t"	"caaaaaat"	No	"a" is repeated more than five times.

Metacharacters: ^

- Pattern after the ^ matches with the beginning of a string or text.
- Not the same meaning as the first hat character within the square brackets "[^]."

RegEx	String	Match?	Explanation
"^Life"	"Life is boring"	Yes	"Life" pattern is found at the beginning of the string.
"^Life"	"My Life is boring"	No	"Life" pattern is not found at the beginning of the string.

| Metacharacters: \$

- ▶ Pattern before the \$ matches with the end of a string or text.

RegEx	String	Match?	Explanation
"Python\$"	"Python is easy"	No	"Python" pattern is not found at the end of the string.
"Python\$"	"You need Python"	Yes	"Python" pattern is found at the end of the string.

| Metacharacters: |

- ▶ Used to join patterns by the logical **or**.
- ▶ More than two patterns can be concatenated by the logical **or**.

RegEx	String	Match?	Explanation
“love hate”	“I love you”	Yes	“love” pattern found in the string.
“love hate”	“I hate him”	Yes	“hate” pattern found in the string.
“love hate”	“I like you”	No	Neither “love” nor “hate” pattern found in the string.

Matching group patterns

- ▶ We can group patterns by enclosing them with `()`.

Ex

```
In [1]: import re  
my_regex = re.compile("[0-9]+[0-9]+([0-9]+)")  
m = my_regex.search("Anna is 15 years old and John is 12 years old.")  
print(m.group(0))  
print(m.group(1))  
print(m.group(2))
```

```
15 years old and John is 12  
15  
12
```

- ▶ In the example, an equivalent regular expression is: `my_regex = re.compile("(\\d+)\\D+(\\d+)")`.

Matching group patterns

- ▶ We can group patterns by enclosing them with `()`.

Ex “Hide the telephone number”

```
In [1]: import requests as rq
res = rq.get("https://en.wikipedia.org/wiki/Machine_learning")
print(res.status_code)
print(res.text)
```

 Line 1-3

- If the status code is 200, then OK.

Matching group patterns

- ▶ We can group patterns by enclosing them with `()`.

Ex “Hide the telephone number”

```
In [1]: import requests as rq
res = rq.get("https://en.wikipedia.org/wiki/Machine_learning")
print(res.status_code)
print(res.text)
```

 Line 1-4

- Not easily intelligible yet.

Matching group patterns

- ▶ We can group patterns by enclosing them with `()`.

Ex “Extract only the phone number”

```
In [2]: my_regex = re.compile("(\\D+)(\\d+)\\D+(\\d+)\\D+(\\d+)")
m = my_regex.search("John 010-1234-5678")
print("Phone number: " + m.group(2))
```

Unit 1.

Text Mining

- | 1.1. What is Text Mining?
- | 1.2. Data Acquisition
- | 1.3. String Manipulation
- | **1.4. Natural Language Processing (NLP)**
- | 1.5. Sequential Data
- | 1.6. Corpus

Natural Language Processing

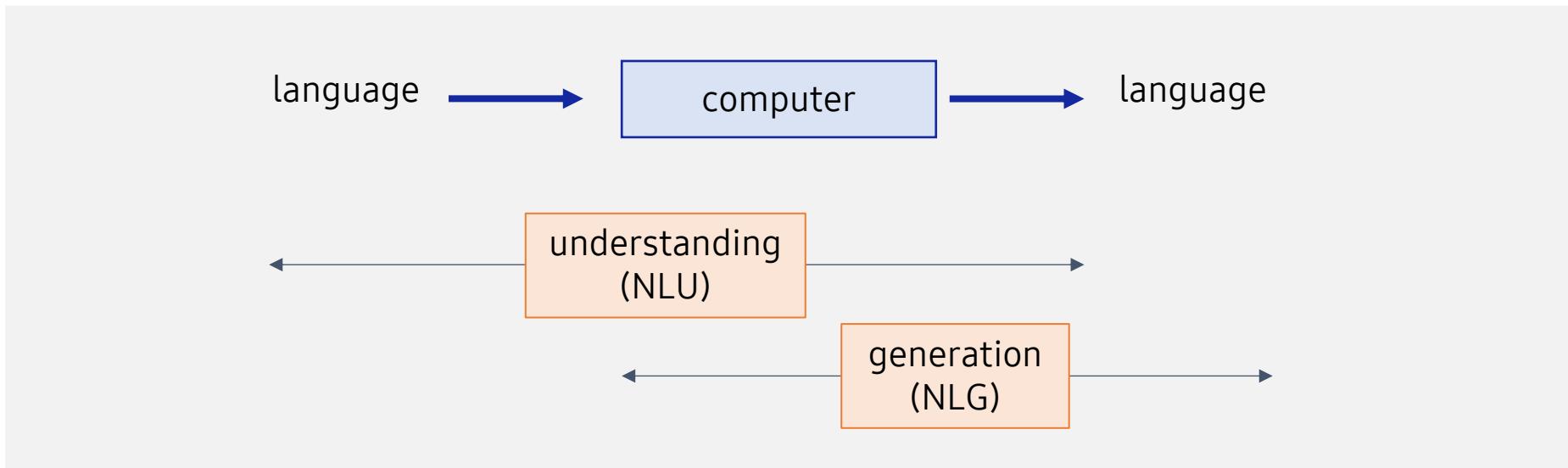
| What is the Natural Language Processing (NLP)?

Natural language refers to language people naturally use in their daily lives.

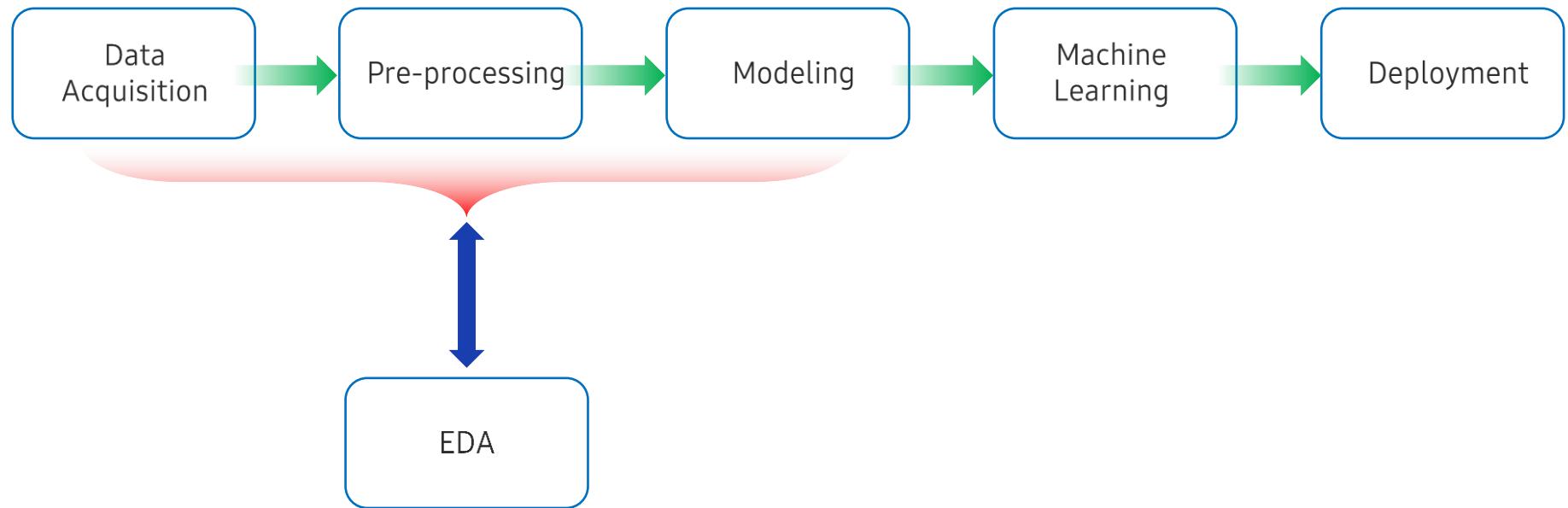
Natural Language Processing (NLP) is an academic field that enables the computer to understand and generate natural language.

- ▶ Extracts features from text data to classify, summarize, cluster, and do sentiment analysis
- ▶ Intersection of linguistics and AI
- ▶ Based on statistical models
- ▶ Different from the way humans understand the language
- ▶ Requires transformation into a structured model

- A narrow meaning of Natural Language Processing is the processing mechanism used by programs using natural language as input and output.



NLLP Workflow



Difficulties of Natural Language Processing

- ▶ NLP is receiving much attention and is widely used, but the procedure is complex.
- ▶ The performance of machine translation has improved, but machine translation still often produces awkward translation results.
- ▶ NLP is complex because the input data is not a numerical value but human language. The input of human language makes data processing extremely complex and uncertain.
 - Even the same words can have various interpretations depending on context. This is called linguistic ambiguity.
 - Like idioms that take a different meaning once various words assemble, there is always an exception to how phrases, words, or morphemes construct a sentence.
 - Since language is flexible and open to expansion, modeling language always entails uncertainty. Also, as time goes by, new words are created, and some words become unused.

I Research paradigm of NLP: (1) Rule-based approach

- ▶ The rule-based approach defines beforehand the grammatical rules of the language and processes natural language based on those rules.
 - It decides the part of speech (POS) for a given word based on linguistic phenomena rather than statistical methods.
 - In the sentence below, it grasps the meaning of the sentence with the first verb and figures out the object of the instruction and its subject matter with 'to' or 'that.'

"Send a message to Michael that I will be late for the meeting."

- ▶ The problem with the rule-based system is that it is impossible to establish the rules prior.
- ▶ Currently, it is only used within the combination of other methods because grammatical rules cannot be entirely neglected in language processing.

I Research paradigm of NLP: (2) Statistics-based approach

- ▶ The statistics-based approach decides the part of speech (POS) by computing lexical probabilities and contextual probabilities within reference to a huge volume of dictionaries to eliminate the ambiguity of POS.
 - Lexical probability is a probability that a certain POS applies to a word. This can be expressed as $P(\text{POS} | \text{word})$ in mathematical form.
 - Contextual possibility is a possibility that a certain POS of a word will show with the POS of the next word. Mathematically, this is $P(\text{POS} | \text{POS})$.
 - Then, it labels the POS that produces the highest result of the multiplication of linguistic probability and contextual probability as the most appropriate for words of semantic ambiguity.
- ▶ There has been much progress as computers analyze sentences much faster with performance improvement, but human intervention is still necessary.
- ▶ Such issues are being tackled by deep learning techniques nowadays.

I Research paradigm of NLP: (3) Deep Learning-based approach

- ▶ While statistical analysis is based on peripheral analysis, such as frequency of word appearance, the Deep Learning approach enables in-depth analysis based on the composite connection among data.

- NLP models that understand a sentence or the overall context of sentences could have been created after deep learning-based NLP was initiated.
- If you create an artificial neural network component that connects with all parts of a sentence, this variable, after learning, contains information about the whole sentence. The accuracy is constantly increasing.

Coding Exercise #0501~0508



Follow practice steps on 'ex_0501~8.ipynb' file

Unit 1.

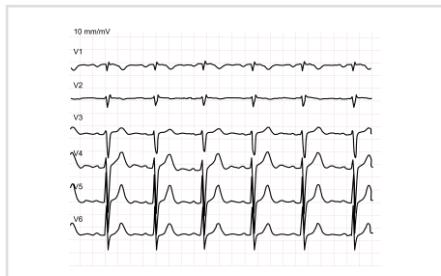
Text Mining

- | 1.1. What is Text Mining?
- | 1.2. Data Acquisition
- | 1.3. String Manipulation
- | 1.4. Natural Language Processing (NLP)
- | **1.5. Sequential Data**
- | 1.6. Corpus

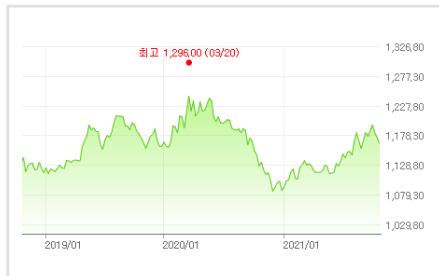
Sequential Data

| Data of temporal property is widely used around the world. (1/2)

- ▶ For example, there are stock prices, human voice, or ECG signals.
- ▶ These data have sequences. You must utilize this feature and temporal information to obtain high performance from sequential data.



(a) ECG signal



(b) Stock price



(c) Vocal signal

Climbing down did I see, the flower I
couldn't see on the way up

(d) sentence

590 600 610
ACTTCAAAATYACGGATTACGCCGAAGAGCTGC
ACTTCAAAATCACGGATTACGCCGAAGAGCTGC
ACTTCAAAATTACGGATTACGCCGAAGAGCTGC

(e) DNA sequence

- | Data of temporal property is widely used around the world. (2/2)
- ▶ Temporal data is dynamic because they change over time; it generally has a variable length.
 - Recurrent neural network, which you will study later, is a learning model that effectively processes such temporal data.
 - ▶ Lately, it can process extremely long patterns that occur in daily life. For example, while only 10 words were the maximum in the past, machine translators can now translate sentences of more than 30 words.
 - ▶ To translate a long sentence, it is necessary to understand the context between two words that are far apart. This is called long-term dependency.
 - Since standard RNN cannot fully process long-term dependency, LSTM, which supplemented selective memory function to standard RNN, is widely used. Selective memory is the capacity to discern memory for the long term and short term.

Unit 1.

Text Mining

- | 1.1. What is Text Mining?
- | 1.2. Data Acquisition
- | 1.3. String Manipulation
- | 1.4. Natural Language Processing (NLP)
- | 1.5. Sequential Data
- | 1.6. Corpus**

Corpus

| Corpus

- ▶ Refers to a set of text data subject to analysis.
 - a) Raw corpus: text data stored in a database
 - b) Tagged corpus: text data where words and phrases have been labeled according to a model

Unit 2.

Text Preprocessing

| 2.1. Tokenization

| 2.2. Stop Words

| 2.3. Lemmatization and Stemming

| 2.4. POS Tagging

| 2.5. Integer Encoding

| 2.6. Padding

| 2.7. One-Hot Encoding

Pre-processing

| Pre-processing

- ▶ Tokenization: break down the raw text into words or sentences.
- ▶ Cleaning:
 - Remove punctuation marks, excess spaces, special characters, etc. (*)
 - Also remove excessively short or infrequent words, etc.
- ▶ Normalization:
 - Conversion to lowercase.
 - Remove the stop words.
 - Stemming and Lemmatization.
 - Expansion of abbreviations. (*)
- ▶ (*) Regular expressions can be useful.

Tokenization

| Often the first step in the data pre-processing

▶ a) Tokenization into sentences:

- Usually, the sentence boundary is marked by a period, exclamation mark, question mark, etc.
- But sometimes, a period does not mark the end of a sentence: abbreviations, for example.

Ex “He got his *M.D.* from the University of Wisconsin.”

- A good sentence tokenizer should recognize these exceptions.

▶ b) Tokenization into words:

- Usually splitting by white space or punctuation mark works.
- In some cases, there is ambiguity: apostrophe, for example.

Ex “*Don't* lose hope. *Everything's* possible.”

⇒ How do we tokenize “*Don't*” and “*Everything's*”? ⇒ It depends on the tokenizer.

Tokenization

```
In [ ]: # NLTK offers tools necessary for tokenization of English corpus.
```

```
In [1]: from nltk.tokenize import word_tokenize  
print(word_tokenize("Don't be fooled by the dark sounding name, Mr. Jone's Orphanage is as cheery as cheery goes for a past  
['Do', "n't", 'be', 'fooled', 'by', 'the', 'dark', 'sounding', 'name', ',', 'Mr.', 'Jone', "'s", 'Orphanage', 'is', 'as',  
'cheery', 'as', 'cheery', 'goes', 'for', 'a', 'pastry', 'shop', '.']")
```

```
In [ ]: # word_tokenize separated Don't into Do and n't, but Jone's into Jone and '
```

```
In [2]: from nltk.tokenize import WordPunctTokenizer  
print(WordPunctTokenizer().tokenize("Don't be fooled by the dark sounding name, Mr. Jone's Orphanage is as cheery as cheery  
['Don', "", 't', 'be', 'fooled', 'by', 'the', 'dark', 'sounding', 'name', ',', 'Mr', '.', 'Jone', "", 's', 'Orphanage',  
'is', 'as', 'cheery', 'as', 'cheery', 'goes', 'for', 'a', 'pastry', 'shop', '.']")
```

```
In [ ]: # Keras, also as a tokenization tool, supports text_to_word_sequence
```

```
In [3]: from tensorflow.keras.preprocessing.text import text_to_word_sequence  
print(text_to_word_sequence("Don't be fooled by the dark sounding name, Mr. Jone's Orphanage is as cheery as cheery goes fo  
['don't', 'be', 'fooled', 'by', 'the', 'dark', 'sounding', 'name', 'mr', "jone's", 'orphanage', 'is', 'as', 'cheery', 'a  
's', 'cheery', 'goes', 'for', 'a', 'pastry', 'shop']")
```

Unit 2.

Text Preprocessing

| 2.1. Tokenization

| **2.2. Stop Words**

| 2.3. Lemmatization and Stemming

| 2.4. POS Tagging

| 2.5. Integer Encoding

| 2.6. Padding

| 2.7. One-Hot Encoding

Stop Words

Stop words and keywords

- ▶ Stop words are commonly used words that do not contain semantic information:

Ex Definite and indefinite articles: “the,” “a,” “an”

Ex Prepositions: “on,” “with,” “into,” “upon,” etc.

- ▶ Stop words are removed during the “data normalization” process.
- ▶ Keywords are selected among the available words after removing the stop words.
 - Often, the most frequent words can be selected as keywords.
 - Sometimes, we should also take into account the purpose and context.

| Stop words and keywords

- ▶ In the NLTK library, there are 179 English stop words:

i, me, my, myself, we, our, ours, ourselves, you, you're, you've, you'll, you'd, your, yours, yourself, yourselves, he, him, his, himself, she, she's, her, hers, herself, it, it's, its, itself, they, them, their, theirs, themselves, what, which, who, whom, this, that, that'll, these, those, am, is, are, was, were, be, been, being, have, has, had, having, do, does, did, doing, a, an, the, and, but, if, or, because, as, until, while, of, at, by, for, with, about, against, between, into, through, during, before, after, above, below, to, from, up, down, in, out, on, off, over, under, again, further, then, once, here, there, when, where, why, how, all, any, both, each, few, more, most, other, some, such, no, nor, not, only, own, same, so, than, too, very, s, t, can, will, just, don, don't, should, should've, now, d, ll, m, o, re, ve, y, ain, aren, aren't, couldn, couldn't, didn, didn't, doesn, doesn't, hadn, hadn't, hasn, hasn't, haven, haven't, isn, isn't, ma, mightn, mightn't, mustn, mustn't, needn, needn't, shan, shan't, shouldn, shouldn't, wasn, wasn't, weren, weren't, won, won't, wouldn, wouldn't

Stop words and keywords

```
In [9]: from nltk.corpus import stopwords  
stopwords.words('english')[:10]
```

```
Out[9]: ['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're"]
```

```
In [11]: # NLTK defines words such as 'i', 'me', 'my' as stop words
```

```
In [10]: from nltk.corpus import stopwords  
from nltk.tokenize import word_tokenize
```

```
example = "Family is not an important thing. It's everything."  
stop_words = set(stopwords.words('english'))
```

```
word_tokens = word_tokenize(example)
```

```
result = []  
for w in word_tokens:  
    if w not in stop_words:  
        result.append(w)
```

```
print(word_tokens)  
print(result)
```

```
['Family', 'is', 'not', 'an', 'important', 'thing', '.', 'It', "'s", 'everything', '.']  
['Family', 'important', 'thing', '.', 'It', "'s", 'everything', '.']
```

Unit 2.

Text Preprocessing

- | 2.1. Tokenization
- | 2.2. Stop Words
- | **2.3. Lemmatization and Stemming**
- | 2.4. POS Tagging
- | 2.5. Integer Encoding
- | 2.6. Padding
- | 2.7. One-Hot Encoding

Lemmatization and Stemming

| Lemmatization

```
In [12]: from nltk.stem import WordNetLemmatizer  
n=WordNetLemmatizer()  
words=['policy', 'doing', 'organization', 'have', 'going', 'love', 'lives', 'fly', 'dies', 'watched', 'has', 'starting']  
print([n.lemmatize(w) for w in words])  
  
['policy', 'doing', 'organization', 'have', 'going', 'love', 'life', 'fly', 'die', 'watched', 'ha', 'starting']  
  
In [13]: n.lemmatize('dies', 'v')  
Out[13]: 'die'  
  
In [14]: n.lemmatize('watched', 'v')  
Out[14]: 'watch'  
  
In [15]: n.lemmatize('has', 'v')  
Out[15]: 'have'
```

Stemming

```
In [16]: from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize
s = PorterStemmer()
text="This was not the map we found in Billy Bones's chest, but an accurate copy, complete in all things--names and heights
words=word_tokenize(text)
print(words)
```

```
[This', 'was', 'not', 'the', 'map', 'we', 'found', 'in', 'Billy', 'Bones', "'s", 'chest', ',', 'but', 'an', 'accurate',
'copy', ',', 'complete', 'in', 'all', 'things', '--', 'names', 'and', 'heights', 'and', 'soundings', '--', 'with', 'the',
'single', 'exception', 'of', 'the', 'red', 'crosses', 'and', 'the', 'written', 'notes', '.']
```

```
In [17]: print([s.stem(w) for w in words])
```

```
['thi', 'wa', 'not', 'the', 'map', 'we', 'found', 'in', 'billi', 'bone', "'s", 'chest', ',', 'but', 'an', 'accur', 'copi',
',', 'complet', 'in', 'all', 'thing', '--', 'name', 'and', 'height', 'and', 'sound', '--', 'with', 'the', 'singl', 'excep
t', 'of', 'the', 'red', 'cross', 'and', 'the', 'written', 'note', '.']
```

Stemming

In [18]: *#Stemmization of porter algorithm has rules as below.*

```
#ALIZE → AL  
#ANCE → Delete  
#ICAL → IC
```

In [19]: words=['formalize', 'allowance', 'electrical']
print([s.stem(w) for w in words])

```
['formal', 'allow', 'electric']
```

In [20]: from nltk.stem import PorterStemmer
s=PorterStemmer()
words=['policy', 'doing', 'organization', 'have', 'going', 'love', 'lives', 'fly', 'dies', 'watched', 'has', 'starting']
print([s.stem(w) for w in words])

```
['polici', 'do', 'organ', 'have', 'go', 'love', 'live', 'fli', 'die', 'watch', 'ha', 'start']
```

In [21]: from nltk.stem import LancasterStemmer
l=LancasterStemmer()
words=['policy', 'doing', 'organization', 'have', 'going', 'love', 'lives', 'fly', 'dies', 'watched', 'has', 'starting']
print([l.stem(w) for w in words])

```
['policy', 'doing', 'org', 'hav', 'going', 'lov', 'liv', 'fly', 'die', 'watch', 'has', 'start']
```

Unit 2.

Text Preprocessing

| 2.1. Tokenization

| 2.2. Stop Words

| 2.3. Lemmatization and Stemming

| 2.4. POS Tagging

| 2.5. Integer Encoding

| 2.6. Padding

| 2.7. One-Hot Encoding

POS Tagging

| Part of Speech (POS) tag set in English (1/2)

Tag Set	Description	Example
CC	Coordinating conjunction	
CD	Cardinal digit	
DT	Determiner	
EX	Existential ‘there’	there is
FW	Foreign word	
IN	Preposition/subordinating conjunction	
JJ	Adjective	Big
JJR	Adjective, comparative	Bigger
JJS	Adjective, superlative	Biggest
LS	List marker	1)
MD	Modal	could, will
NN	Noun, singular	Desk
NNS	Noun, plural	Desks
NNP	Proper noun, singular	Harrison
NNPS	Proper noun, plural	Americans
PDT	Predeterminer	‘all’ the kids
POS	Possessive ending	parent’s

I Part of Speech (POS) tag set in English (2/2)

Tag Set	Description	Example
PRP	Personal pronoun	I, he, she
PRP\$	Possessive pronoun	my, his, hers
RB	Adverb	very, silently
RBR	Adverb, comparative	better
RBS	Adverb, superlative	best
RP	Particle	give up
TO	To	go 'to' the store
UH	Interjection	errrrrrrm
VB	Verb, base form	take
VBD	Verb, past tense	took
VBG	Verb, gerund/present participle	taking
VBN	Verb, past participle	taken
VBP	Verb, sing. Present, non-3d	take
VBZ	Verb, 3rd person sing. Present	takes
WDT	Wh-determiner	which
WP	Wh-pronoun	who, what
WP\$	Possessive wh-pronoun	whose
WRB	Wh-abverb	where, when

Part of Speech (POS) tag set in English (Example)

```
In [22]: # Test sentence.  
my_sentence = "The Colosseum was built by the emperor Vespassian"  
  
In [24]: import nltk  
  
In [25]: # Simple pre-processing.  
my_words = nltk.word_tokenize(my_sentence)  
for i in range(len(my_words)):  
    my_words[i] = my_words[i].lower()  
my_words  
  
Out[25]: ['the', 'colosseum', 'was', 'built', 'by', 'the', 'emperor', 'vespassian']  
  
In [26]: # POS tagging.  
# OUTPUT: A list of tuples.  
my_words_tagged = nltk.pos_tag(my_words)  
my_words_tagged  
  
Out[26]: [('the', 'DT'),  
          ('colosseum', 'NN'),  
          ('was', 'VBD'),  
          ('built', 'VBN'),  
          ('by', 'IN'),  
          ('the', 'DT'),  
          ('emperor', 'NN'),  
          ('vespassian', 'NN')]
```

NLTK Library

| NLTK (Natural Language Toolkit)

- ▶ One of the most used Python libraries for the NLP:
 - a) Tokenization: sent_tokenize, word_tokenize, etc.
 - b) Stop words: 179 in total. The list of stop words needs to be downloaded.
`nltk.download('stopwords')`
 - c) Stemming: PorterStemmer, LancasterStemmer, SnowballStemmer, etc.
 - d) Lemmatization: WordNetLemmatizer.
 - e) POS tagging: Uses the “Penn Treebank tag set.”
 - f) Lexical resource: WordNet, SentiWordNet, etc. Useful for sentiment analysis.

Visualization

| Word Cloud

- ▶ Visualization of the keywords where the size is related to the frequency.



Word Cloud

- ▶ We can create a word cloud following the steps below:
 - 1) Tokenize by words.
 - 2) Apply cleaning and normalization.
 - 3) Make a frequency table of the words.
 - 4) Sort by frequency and keep only the top keywords. The number of keywords is adjustable.
 - 5) Customize the arguments of the WordCloud() function.

Unit 2.

Text Preprocessing

| 2.1. Tokenization

| 2.2. Stop Words

| 2.3. Lemmatization and Stemming

| 2.4. POS Tagging

| **2.5. Integer Encoding**

| 2.6. Padding

| 2.7. One-Hot Encoding

Integer Encoding

| Exercise

```
In [1]: import nltk  
nltk.download('punkt')
```

```
[nltk_data] Downloading package punkt to  
[nltk_data]   C:\Users\emcast\AppData\Roaming\nltk_data...  
[nltk_data]   Package punkt is already up-to-date!
```

```
Out[1]: True
```

```
In [2]: from nltk.tokenize import sent_tokenize  
from nltk.tokenize import word_tokenize  
from nltk.corpus import stopwords
```

```
In [3]: text = "A barber is a person. a barber is a good person. a barber is a huge person. he Knew A Secret!. The Secret He Kept i
```

```
In [4]: text = sent_tokenize(text)  
print(text)
```

```
['A barber is a person.', 'a barber is a good person.', 'a barber is a huge person.', 'he Knew A Secret!.', 'The Secret He  
Kept is huge secret.', 'Huge secret.', 'His barber kept his word.', 'a barber kept his word.', 'His barber kept his secre  
t.', 'But keeping and keeping such a huge secret to himself was driving the barber crazy.', 'the barber went up a huge mou  
ntain.']}
```

Integer Encoding

```
In [4]: # Cleaning and word tokenization
vocab = {} # Python's dictionary datatype
sentences = []
stop_words = set(stopwords.words('english'))

for i in text:
    sentence = word_tokenize(i) # executes word tokenization
    result = []

    for word in sentence:
        word = word.lower() # Reduce the number of words by lowercasing all words
        if word not in stop_words: # Remove stop words in case of word tokenization.
            if len(word) > 2: # Remove additional words if the length of the word is lower or equal to 2
                result.append(word)
            if word not in vocab:
                vocab[word] = 0
                vocab[word] += 1
    sentences.append(result)
print(sentences)

[['barber', 'person'], ['barber', 'good', 'person'], ['barber', 'huge', 'person'], ['knew', 'secret'], ['secret', 'kept', 'huge', 'secret'], ['huge', 'secret'], ['barber', 'kept', 'word'], ['barber', 'kept', 'word'], ['barber', 'kept', 'secret'], ['keeping', 'keeping', 'huge', 'secret', 'driving', 'barber', 'crazy'], ['barber', 'went', 'huge', 'mountain']]
```

```
In [5]: print(vocab)

{'barber': 8, 'person': 3, 'good': 1, 'huge': 5, 'knew': 1, 'secret': 6, 'kept': 4, 'word': 2, 'keeping': 2, 'driving': 1, 'crazy': 1, 'went': 1, 'mountain': 1}
```

```
In [6]: print(vocab["barber"])
```

Aligning in the order of frequency

```
In [7]: vocab_sorted = sorted(vocab.items(), key = lambda x:x[1], reverse = True)
print(vocab_sorted)

[('barber', 8), ('secret', 6), ('huge', 5), ('kept', 4), ('person', 3), ('word', 2), ('keeping', 2), ('good', 1), ('knew', 1), ('driving', 1), ('crazy', 1), ('went', 1), ('mountain', 1)]
```

Assigning low integer index for words of high frequency

```
In [8]: word_to_index = {}
i=0
for (word, frequency) in vocab_sorted :
    if frequency > 1 : # Low frequency words are deleted. (It was learned in the Cleaning section.)
        i=i+1
        word_to_index[word] = i
print(word_to_index)

{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5, 'word': 6, 'keeping': 7}
```

Using top 5 words

```
In [9]: vocab_size = 5
words_frequency = [w for w,c in word_to_index.items() if c >= vocab_size + 1] # Remove words of index higher than 5
for w in words_frequency:
    del word_to_index[w] # Remove index information from the corresponding word
print(word_to_index)

{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5}
```

Changing each word from sentences stored after tokenization to an integer using word_to_index

```
In [11]: word_to_index['OOV'] = len(word_to_index) + 1
```

```
In [12]: encoded = []
for s in sentences:
    temp = []
    for w in s:
        try:
            temp.append(word_to_index[w])
        except KeyError:
            temp.append(word_to_index['OOV'])
    encoded.append(temp)
print(encoded)
```

```
[[1, 5], [1, 6, 5], [1, 3, 5], [6, 2], [2, 4, 3, 2], [3, 2], [1, 4, 6], [1, 4, 6], [1, 4, 2], [6, 6, 3, 2, 6, 1, 6], [1, 6, 3, 6]]
```

Using Counter

```
In [13]: from collections import Counter
```

```
In [14]: print(sentences)
```

```
[['barber', 'person'], ['barber', 'good', 'person'], ['barber', 'huge', 'person'], ['knew', 'secret'], ['secret', 'kept', 'huge', 'secret'], ['huge', 'secret'], ['barber', 'kept', 'word'], ['barber', 'kept', 'word'], ['barber', 'kept', 'secret'], ['keeping', 'huge', 'secret', 'driving', 'barber', 'crazy'], ['barber', 'went', 'huge', 'mountain']]
```

```
In [15]: words = sum(sentences, [])
# Above task can be executed by words=np.hstack(sentences)
print(words)
```

```
['barber', 'person', 'barber', 'good', 'person', 'barber', 'huge', 'person', 'knew', 'secret', 'secret', 'kept', 'huge', 'secret', 'huge', 'secret', 'barber', 'kept', 'word', 'barber', 'kept', 'word', 'barber', 'kept', 'secret', 'keeping', 'huge', 'secret', 'driving', 'barber', 'crazy', 'barber', 'went', 'huge', 'mountain']
```

```
In [16]: vocab = Counter(words) # You can easily count all frequency of the words by using Python's Counter module.
print(vocab)
```

```
Counter({'barber': 8, 'secret': 6, 'huge': 5, 'kept': 4, 'person': 3, 'word': 2, 'keeping': 2, 'good': 1, 'knew': 1, 'driving': 1, 'crazy': 1, 'went': 1, 'mountain': 1})
```

```
In [17]: print(vocab["barber"]) # Print the frequency of the word 'barber'
```

```
8
```

Assigning low integer index for words of high frequency

```
In [18]: vocab_size = 5
vocab = vocab.most_common(vocab_size) # Save top 5 words with the highest frequency
vocab
```

```
Out[18]: [('barber', 8), ('secret', 6), ('huge', 5), ('kept', 4), ('person', 3)]
```

```
In [19]: word_to_index = {}
i = 0
for (word, frequency) in vocab :
    i = i+1
    word_to_index[word] = i
print(word_to_index)
```

```
{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5}
```

Using NLTK's FreqDist

```
In [20]: from nltk import FreqDist
import numpy as np

In [21]: # Remove sentence section with np.hstack and use it as input. Ex) 'barber', 'person', 'barber', 'good' ...
vocab = FreqDist(np.hstack(sentences))

In [22]: print(vocab["barber"])

8

In [23]: vocab_size = 5
vocab = vocab.most_common(vocab_size) # Save top 5 words with the highest frequency
vocab

Out[23]: [('barber', 8), ('secret', 6), ('huge', 5), ('kept', 4), ('person', 3)]

In [24]: word_to_index = {word[0] : index + 1 for index, word in enumerate(vocab)}
print(word_to_index)

{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5}
```

Using Enumerate

```
In [25]: test=['a', 'b', 'c', 'd', 'e']
for index, value in enumerate(test): # Assign index starting from 0 in the order of input
    print("value : {}, index: {}".format(value, index))

value : a, index: 0
value : b, index: 1
value : c, index: 2
value : d, index: 3
value : e, index: 4
```

Using Keras

```
In [27]: from tensorflow.keras.preprocessing.text import Tokenizer
```

```
In [28]: sentences=[['barber', 'person'], ['barber', 'good', 'person'], ['barber', 'huge', 'person'], ['knew', 'secret'],  
    < ... >
```

```
In [29]: tokenizer = Tokenizer()  
tokenizer.fit_on_texts(sentences)
```

```
In [30]: print(tokenizer.word_index)
```

```
{'barber': 1, 'secret': 2, 'huge': 3, 'kept': 4, 'person': 5, 'word': 6, 'keeping': 7, 'good': 8, 'knew': 9,  
'driving': 10, 'crazy': 11, 'went': 12, 'mountain': 13}
```

```
In [31]: print(tokenizer.word_counts)
```

```
OrderedDict([('barber', 8), ('person', 3), ('good', 1), ('huge', 5), ('knew', 1), ('secret', 6), ('kept', 4),  
('word', 2), ('keeping', 2), ('driving', 1), ('crazy', 1), ('went', 1), ('mountain', 1)])
```

```
In [32]: print(tokenizer.texts_to_sequences(sentences))
```

```
[[1, 5], [1, 8, 5], [1, 3, 5], [9, 2], [2, 4, 3, 2], [3, 2], [1, 4, 6], [1, 4, 6], [1, 4, 2], [7, 7, 3, 2, 10,  
1, 11], [1, 12, 3, 13]]
```

Coding Exercise



Follow practice steps on 'Integer Encoding.ipynb' file

Unit 2.

Text Preprocessing

- | 2.1. Tokenization
- | 2.2. Stop Words
- | 2.3. Lemmatization and Stemming
- | 2.4. POS Tagging
- | 2.5. Integer Encoding
- | **2.6. Padding**
- | 2.7. One-Hot Encoding

Padding

Padding using NumPy

```
In [1]: import numpy as np
from tensorflow.keras.preprocessing.text import Tokenizer

In [2]: sentences = [['barber', 'person'], ['barber', 'good', 'person'], ['barber', 'huge', 'person'],
                 ['knew', 'secret'], ['secret', 'kept', 'huge', 'secret'], ['huge', 'secret'],
                 ['barber', 'kept', 'word'], ['barber', 'kept', 'word'], ['barber', 'kept', 'secret'],
                 ['keeping', 'keeping', 'huge', 'secret', 'driving', 'barber', 'crazy'],
                 ['barber', 'went', 'huge', 'mountain']]

In [3]: tokenizer = Tokenizer()
tokenizer.fit_on_texts(sentences) # fit_on_texts(), with corpus as input, generates vocabulary based on word frequency

In [4]: encoded = tokenizer.texts_to_sequences(sentences)
print(encoded)

[[[1, 5], [1, 8, 5], [1, 3, 5], [9, 2], [2, 4, 3, 2], [3, 2], [1, 4, 6], [1, 4, 6], [1, 4, 2], [7, 7, 3, 2, 10, 1, 11], [1,
12, 3, 13]]]

In [5]: max_len = max(len(item) for item in encoded)
print(max_len)
```

7

Padding using NumPy

```
In [6]: for item in encoded: # For each sentence
    while len(item) < max_len: # If smaller than max_len
        item.append(0)
```

```
padded_np = np.array(encoded)
padded_np
```

```
Out[6]: array([[ 1,  5,  0,  0,  0,  0,  0],
   [ 1,  8,  5,  0,  0,  0,  0],
   [ 1,  3,  5,  0,  0,  0,  0],
   [ 9,  2,  0,  0,  0,  0,  0],
   [ 2,  4,  3,  2,  0,  0,  0],
   [ 3,  2,  0,  0,  0,  0,  0],
   [ 1,  4,  6,  0,  0,  0,  0],
   [ 1,  4,  6,  0,  0,  0,  0],
   [ 1,  4,  2,  0,  0,  0,  0],
   [ 7,  7,  3,  2, 10,  1, 11],
   [ 1, 12,  3, 13,  0,  0,  0]])
```

Padding using Keras preprocessing tool

```
In [8]: from tensorflow.keras.preprocessing.sequence import pad_sequences
```

```
In [10]: encoded = tokenizer.texts_to_sequences(sentences)
print(encoded)
```

```
[[1, 5], [1, 8, 5], [1, 3, 5], [9, 2], [2, 4, 3, 2], [3, 2], [1, 4, 6], [1, 4, 6], [1, 4, 2], [7, 7, 3, 2, 10, 1, 11], [1, 12, 3, 13]]
```

```
In [11]: padded = pad_sequences(encoded)
padded
```

```
Out[11]: array([[ 0,  0,  0,  0,  0,  1,  5],
 [ 0,  0,  0,  0,  1,  8,  5],
 [ 0,  0,  0,  0,  1,  3,  5],
 [ 0,  0,  0,  0,  0,  9,  2],
 [ 0,  0,  0,  2,  4,  3,  2],
 [ 0,  0,  0,  0,  0,  3,  2],
 [ 0,  0,  0,  0,  1,  4,  6],
 [ 0,  0,  0,  0,  1,  4,  6],
 [ 0,  0,  0,  0,  1,  4,  2],
 [ 7,  7,  3,  2,  10,  1,  11],
 [ 0,  0,  0,  1,  12,  3,  13]])
```

Coding Exercise



Follow practice steps on 'Padding.ipynb' file

Unit 2.

Text Preprocessing

- | 2.1. Tokenization
- | 2.2. Stop Words
- | 2.3. Lemmatization and Stemming
- | 2.4. POS Tagging
- | 2.5. Integer Encoding
- | 2.6. Padding
- | **2.7. One-Hot Encoding**

One-Hot Encoding

One-hot encoding using Keras preprocessing tool

```
In [1]: text='I want to go to lunch with me. The lunch menu is hamburgers. Hamburgers are the best'
```

```
In [2]: from tensorflow.keras.preprocessing.text import Tokenizer  
from tensorflow.keras.utils import to_categorical
```

```
In [3]: t = Tokenizer()  
t.fit_on_texts([text])  
print(t.word_index) # Print encoding result for each word
```

```
{'to': 1, 'lunch': 2, 'the': 3, 'hamburgers': 4, 'i': 5, 'want': 6, 'go': 7, 'with': 8, 'me': 9, 'menu': 10, 'is': 11, 'are': 12, 'best': 13}
```

```
In [4]: encoded=t.texts_to_sequences([text])  
print(encoded)
```

```
[[5, 6, 1, 7, 1, 2, 8, 9, 3, 2, 10, 11, 4, 4, 12, 3, 13]]
```

One-hot encoding using Keras preprocessing tool

```
In [5]: one_hot = to_categorical(encoded)
print(one_hot)

[[[0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0.]
 [0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0.]
 [0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 0.]
 [0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 1. 0.]
 [0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]
 [0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 1.]
 [0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0.]]]
```

Coding Exercise



Follow practice steps on 'One-hot Encoding.ipynb' file

Unit 3.

Language Model

| 3.1. Language Model

| 3.2. Representation Model

| 3.3. Classification Analysis

| 3.4. Topic Modeling

Language Model

Language model refers to a model that predicts or generates the next component by assigning the probability to elements of language (letter, word, morpheme, string (sentence), paragraph, etc.).

- ▶ Language model is divided into statistical language model (SLM) and deep learning language model based on artificial neural network. Essentially language models, based on a given word, predicts the next word or combination of words. Such functions can solve numerous natural language processing problems like document generation, machine translation, document summarization, etc.

About language model

- ▶ Predicts the probability of a sequence: $P(w_1, w_2, w_3, \dots, w_i)$

Caution: The sub-index of w means the actual time order that cannot be changed.

- ▶ Given a sequence of words $\{w_1, w_2, w_3, \dots, w_{(i-1)}\}$ what is the probability of w_i ?

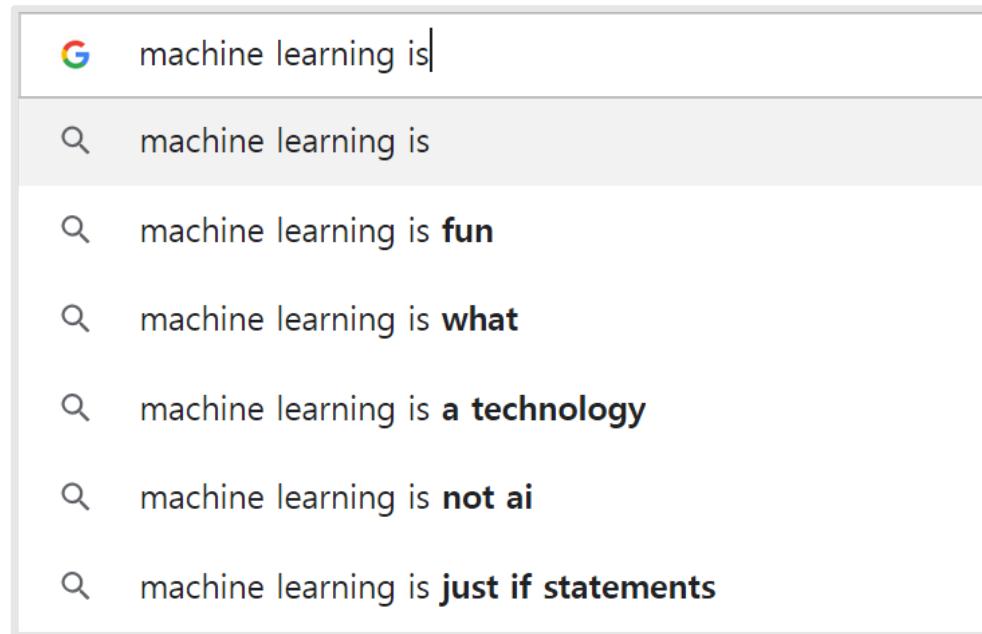
$P(w_i | w_1, w_2, w_3, \dots, w_{(i-1)})$?

- ▶ Data sparsity is a major problem because most (long) sequences appear very infrequently.

- ▶ Practical applications: machine translation, speech recognition, spell correction, autofill, etc.

| About language model

Ex In the search engine:



Regular expression

- ▶ A joint probability can be expanded as following:

$$\begin{aligned} P(w_1, w_2, w_3, \dots, w_m) &= P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)P(w_4|w_1, w_2, w_3) \cdots P(w_m|w_1, w_2, \dots, w_{m-1}) \\ &= \prod_{i=1}^m P(w_i|w_1, \dots, w_{i-1}) \end{aligned}$$

Ex $P(\text{three little pigs lived happily})?$

⇒

w_1	w_2	w_3	w_4	w_5
“three”	“little”	“pigs”	“lived”	“happily”

$P(\text{three little pigs lived happily})$

$$= P(\text{three})P(\text{little}|\text{three})P(\text{pigs}|\text{three little})P(\text{lived}|\text{three little pigs})P(\text{happily}|\text{three little pigs lived})$$

n-Grams

- Given a text sequence, n-Grams can be constructed by sliding a “moving window” of length = n .

Ex “three little pigs lived happily”

→ $n = 1$, Unigrams = [“three,” “little,” “pigs’, “lived,” “happily”]

→ $n = 2$, Bigrams = [“three little,” “little pigs,” “pigs lived,” “lived happily”]

→ $n = 3$, Trigrams = [“**three little pigs**,” “**little pigs lived**,” “**pigs lived happily**”]

“three **little pigs** lived happily”

“three **little pigs lived** happily”

“three little **pigs lived happily**”

n-Gram approximations

- As the sequence grows, the probabilities become harder to estimate due to the data sparsity:

$$P(w_i | w_1, w_2, w_3, \dots, w_{i-1}) = \frac{\text{Count}(w_1, w_2, w_3, \dots, w_i)}{\text{Count}(w_1, w_2, w_3, \dots, w_{i-1})}$$

- Instead of an exact estimation of probabilities, we can do the so-called n -Gram approximation:

$$P(w_1, w_2, w_3, \dots, w_m) \approx \prod_{i=1}^m P(w_i | w_{i-(n-1)}, \dots, w_{i-1})$$

which can be compared with the following exact relation.

$$P(w_1, w_2, w_3, \dots, w_m) = \prod_{i=1}^m P(w_i | w_1, \dots, w_{i-1})$$

- => Usually the n above is a small positive number $\cong 1, 2, 3, \dots$

| n-Gram approximations

- ▶ When $n=1$, it is the Unigram approximation:

$$P(w_1, w_2, w_3, \dots, w_m) \approx P(w_1)P(w_2)P(w_3) \cdots P(w_m)$$

- ▶ When $n=2$, it is the Bigram approximation:

$$P(w_1, w_2, w_3, \dots, w_m) \approx P(w_1)P(w_2|w_1)P(w_3|w_2) \cdots P(w_m|w_{m-1})$$

- ▶ When $n=3$, it is the Trigram approximation:

$$P(w_1, w_2, w_3, \dots, w_m) \approx P(w_1)P(w_2|w_1)P(w_3|w_2, w_1)P(w_4|w_3, w_2) \cdots P(w_m|w_{m-1}, w_{m-2})$$

Ex Bigram approximation for **Sequence** = “*three little pigs lived happily*”

$$P(\text{Sequence}) \approx P(\text{three})P(\text{little}|\text{three})P(\text{pigs}|\text{little})P(\text{lived}|\text{pigs})P(\text{happily}|\text{lived})$$

Coding Exercise #0509



Follow practice steps on 'ex_0509.ipynb' file

Unit 3.

Language Model

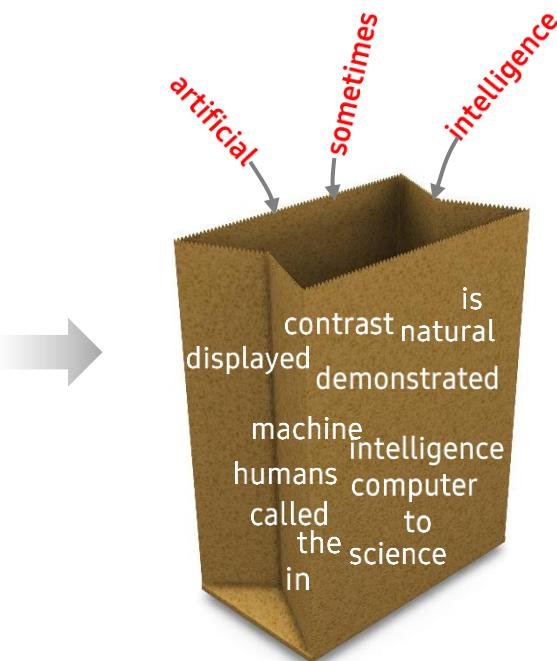
- | 3.1. Language Model
- | **3.2. Representation Model**
- | 3.3. Classification Analysis
- | 3.4. Topic Modeling

Representation Model

| Bag-of-Words (BOW) model

- ▶ A document is represented by a collection of its words.
- ▶ Word ordering and grammar are ignored.
- ▶ Only the word frequencies matter.

Ex “In computer science, artificial intelligence, sometimes called machine intelligence, is intelligence demonstrated by machines, in contrast to the natural intelligence displayed by humans.”



| Bag-of-Words (BOW) model

Ex “It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness.”

⇒ after removing the stop words such as “it,” “the,” and “of,” the BOW can be expressed as an array.

age	best	foolishness	times	was	wisdom	worst
2	1	1	2	1	1	1

Document-Term Matrix (DTM) and Term-Document Matrix (TDM)

- ▶ Documents expressed as BOW are the rows of DTM.
- ▶ Documents expressed as BOW are the columns of TDM.

	Feature #1	Feature #2	Feature #3	Feature #4	⋮
Document #1	1	0	1	0	0
Document #2	0	0	2	0	0
Document #3	0	1	0	0	1
Document #4	0	0	0	1	0
⋮	0	0	1	0	0

DTM

	Document #1	Document #2	Document #3	Document #4	⋮
Feature #1	1	0	0	0	0
Feature #2	0	0	1	0	0
Feature #3	1	2	0	0	1
Feature #4	0	0	0	1	0
⋮	0	0	1	0	0

TDM

| Document-Term Matrix (DTM) and Term-Document Matrix (TDM)

Ex Let's suppose the following pre-processed documents:

Document #1: "learning intelligence machine learning statistics"

Document #2: "machine classification learning performance"

Document #3: "machine classification, machine learning, machine performance"

DTM =

	learning	intelligence	machine	statistics	classification	performance
Doc. #1	2	1	1	1	0	0
Doc. #2	1	0	1	0	1	1
Doc. #3	1	0	3	0	1	1

| Document-Term Matrix (DTM) and Term-Document Matrix (TDM)

Ex Let's suppose the following pre-processed documents:

Document #1: "learning intelligence machine learning statistics"

Document #2: "machine classification learning performance"

Document #3: "machine classification, machine learning, machine performance"

TDM =

	Doc. #1	Doc. #2	Doc. #3
learning	2	1	1
intelligence	1	0	0
machine	1	1	3
statistics	1	0	0
classification	0	1	1
performance	0	1	1

| Term Frequency (TF)

- ▶ Indicates the relative importance of each word (term) within a document.
- ▶ A frequently occurring word within a short document would have a large TF value.

$$TF(\text{word}, \text{document}) = \frac{\text{Frequency of the } \text{word} \text{ within the } \text{document}}{\text{The } \text{document} \text{ length}}$$

- ▶ TF has to be calculated per **word** and per **document**.

| Team Frequency (TF)

Ex Let's suppose the following pre-processed documents:

Document #1: "learning intelligence machine learning statistics"

⇒ length = 5

Document #2: "machine classification learning performance"

⇒ length = 4

Document #3: "machine classification, machine learning, machine performance"

⇒ length = 6

	Doc. #1	Doc. #2	Doc. #3
learning	$2/5 = 0.4$	$1/4 = 0.25$	$1/6 = 0.17$
intelligence	$1/5 = 0.2$	0	0
machine	$1/5 = 0.2$	$1/4 = 0.25$	$3/6 = 0.5$
statistics	$1/5 = 0.2$	0	0
classification	0	$1/4 = 0.25$	$1/6 = 0.17$
performance	0	$1/4 = 0.25$	$1/6 = 0.17$

TF =

| Document Frequency (DF) and Inverse Document Frequency (IDF)

- ▶ DF: the number of documents where a particular word appears
- ▶ IDF: a measure of rarity and information carried by a particular word

$$IDF(\text{word}) = \log \left(\frac{\text{Total number of documents}}{\text{Number of documents that include the word}} \right)$$

- ▶ IDF is a property of the corpus and has to be calculate per word only.

| Document Frequency (DF) and Inverse Document Frequency (IDF)

Ex Let's suppose the following pre-processed documents:

Document #1: "learning intelligence machine learning statistics"

Document #2: "machine classification learning performance"

Document #3: "machine classification, machine learning, machine performance"

IDF =

	DF	IDF
learning	3	$\text{Log}(3/3) = 0$
intelligence	1	$\text{Log}(3/1) = 0.48$
machine	3	$\text{Log}(3/3) = 0$
statistics	1	$\text{Log}(3/1) = 0.48$
classification	2	$\text{Log}(3/2) = 0.18$
performance	2	$\text{Log}(3/2) = 0.18$

TF IDF representation

Ex Let's suppose the following pre-processed documents:

Document #1: "learning intelligence machine learning statistics"

Document #2: "machine classification learning performance"

Document #3: "machine classification, machine learning, machine performance"

TF IDF =

	Doc. #1	Doc. #2	Doc. #3
learning	0.4	0.25	0.17
intelligence	0.2	0	0
machine	0.2	0.25	0.5
statistics	0.2	0	0
classification	0	0.25	0.17
performance	0	0.25	0.17

✗

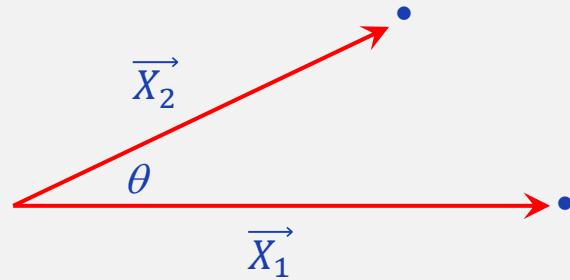
	IDF
learning	0
intelligence	0.48
machine	0
statistics	0.48
classification	0.18
performance	0.18

=

	Doc. #1	Doc. #2	Doc. #3
learning	0	0	0
intelligence	0.095	0	0
machine	0	0	0
statistics	0.095	0	0
classification	0	0.044	0.03
performance	0	0.044	0.03

| Cosine similarity

- ▶ Documents are vectors. The similarity between two documents can be quantified.



$$\text{Cosine similarity is } \text{Cos}(\theta) = \frac{\vec{X}_1 \cdot \vec{X}_2}{|\vec{X}_1| |\vec{X}_2|}$$

Coding Exercise #0510



Follow practice steps on 'ex_0510.ipynb' file

Unit 3.

Language Model

- | 3.1. Language Model
- | 3.2. Representation Model
- | **3.3. Classification Analysis**
- | 3.4. Topic Modeling

Naïve Bayes Classifier

| Naïve Bayes classifier using the BOW model

- ▶ For convenience, suppose there are two document types, A and B. The document types can be, for example, A= “spam” and B = “no spam.”
- ▶ Let’s apply the BOW model: bags A and B contain the tokenized words.
- ▶ Applying Bayes’ theorem, we have:

$$P(\textcolor{red}{A}|w_1, w_2, w_3, \dots) = \frac{P(w_1, w_2, w_3, \dots | \textcolor{red}{A})P(\textcolor{red}{A})}{P(w_1, w_2, w_3, \dots)}$$

$$P(\textcolor{blue}{B}|w_1, w_2, w_3, \dots) = \frac{P(w_1, w_2, w_3, \dots | \textcolor{blue}{B})P(\textcolor{blue}{B})}{P(w_1, w_2, w_3, \dots)}$$

Caution: The sub-index of w serves only a labeling purpose. Here, the words are not ordered.

Naïve Bayes classifier using the BOW model

- ▶ Prediction based on the comparison between $P(A|w_1, w_2, w_3, \dots)$ and $P(B|w_1, w_2, w_3, \dots)$.
- ▶ For the comparison, only the relative difference matters. Question: Which probability is higher?
- ▶ For the comparison, we do not need the common denominator $P(w_1, w_2, w_3, \dots)$.

$$P(\textcolor{red}{A}|w_1, w_2, w_3, \dots) \sim P(w_1, w_2, w_3, \dots | \textcolor{red}{A})P(\textcolor{red}{A})$$

$$P(\textcolor{blue}{B}|w_1, w_2, w_3, \dots) \sim P(w_1, w_2, w_3, \dots | \textcolor{blue}{B})P(\textcolor{blue}{B})$$

- ▶ In the BOW model, the words occur independently from each other. Thus, we can expand in the following way:

$$P(\textcolor{red}{A}|w_1, w_2, w_3, \dots) \sim P(w_1|\textcolor{red}{A})P(w_2|\textcolor{red}{A})P(w_3|\textcolor{red}{A}) \cdots P(\textcolor{red}{A})$$

$$P(\textcolor{blue}{B}|w_1, w_2, w_3, \dots) \sim P(w_1|\textcolor{blue}{B})P(w_2|\textcolor{blue}{B})P(w_3|\textcolor{blue}{B}) \cdots P(\textcolor{blue}{B})$$

Naïve Bayes classifier using the BOW model

- ▶ Instead of comparing the probabilities, we can compare the logarithms of probabilities.
- ▶ Applying the *Log()* on both sides of the equal sign, we have:

$$\text{Log}(P(\textcolor{red}{A}|w_1, w_2, w_3, \dots)) \sim \text{Log}(P(w_1|\textcolor{red}{A})) + \text{Log}(P(w_2|\textcolor{red}{A})) + \text{Log}(P(w_3|\textcolor{red}{A})) + \dots + \text{Log}(P(\textcolor{red}{A}))$$

$$\text{Log}(P(\textcolor{blue}{B}|w_1, w_2, w_3, \dots)) \sim \text{Log}(P(w_1|\textcolor{blue}{B})) + \text{Log}(P(w_2|\textcolor{blue}{B})) + \text{Log}(P(w_3|\textcolor{blue}{B})) + \dots + \text{Log}(P(\textcolor{blue}{B}))$$

- ▶ If we balance the training set such that the number of type **A** = number of type **B**, then $\text{Log}(P(\textcolor{red}{A})) = \text{Log}(P(\textcolor{blue}{B}))$. So we can also **drop** these terms in the comparison.

Naïve Bayes classifier using the BOW model

- ▶ Training steps:

- 1) For each word in bag A, calculate the probabilities $P(w_i|A)$ and their logarithms $\text{Log}(P(w_i|A))$.
- 2) For each word in bag B, calculate the probabilities $P(w_i|B)$ and their logarithms $\text{Log}(P(w_i|B))$.
- 3) Save for the later use of the logarithmic probabilities calculated in steps 1) and 2).



Naïve Bayes classifier using the BOW model

- Prediction steps:

- Given a test document made up of words w'_1, w'_2, w'_3, \dots add their logarithmic probabilities:

$$\text{LogProbA} = \text{Log}(P(w'_1|A)) + \text{Log}(P(w'_2|A)) + \text{Log}(P(w'_3|A)) + \dots$$

$$\text{LogProbB} = \text{Log}(P(w'_1|B)) + \text{Log}(P(w'_2|B)) + \text{Log}(P(w'_3|B)) + \dots$$

- If $\text{LogProbA} > \text{LogProbB}$: then the test document is predicted as type A.

If $\text{LogProbA} < \text{LogProbB}$: then the test document is predicted as type B.

Classification Analysis

| Classification analysis using the TF IDF model

- ▶ In the TF IDF model:
 - document \cong observation
 - word (W_i) \cong explanatory variable (X_i)
- ▶ If the data is labeled (response Y), we can do predictive analysis using classification algorithms like logistic regression, KNN, decision tree, Random Forest, etc.

	X_1	X_2	X_3	...	Y
Obs. #1
Obs. #2
Obs. #3
:					



	W_1	W_2	W_3	...	Y
Doc. #1
Doc. #2
Doc. #3
:					

Coding Exercise #0511



Follow practice steps on 'ex_0511.ipynb' file

Unit 3.

Language Model

- | 3.1. Language Model
- | 3.2. Representation Model
- | 3.3. Classification Analysis
- | 3.4. Topic Modeling**

Latent Semantic Analysis (LSA)

| About the LSA

- ▶ Extracts common topics from a set of documents.
- ⇒ With the TF IDF matrix, carry out the SVD and extract the principal components.
- ▶ Let's suppose that a TF IDF matrix \mathbf{M} has the following size:

$$\text{Size}(\mathbf{M})=m \times n$$

m = Number of documents

n = Number of features

- ▶ Each “topic vector” is a principal component of \mathbf{M} .
- ▶ By sorting the singular values, we can extract the most salient topics.

LSA can be used for:

- ▶ Clustering of documents
- ▶ Studying the relationship between the documents
- ▶ Labeling of documents for a search engine

| Singular value decomposition (SVD)

- ▶ A matrix \mathbf{M} is decomposed as $\mathbf{M} = \mathbf{U}\Sigma\mathbf{V}^T$.

where,

$$\text{Size}(\mathbf{M})=m \times n$$

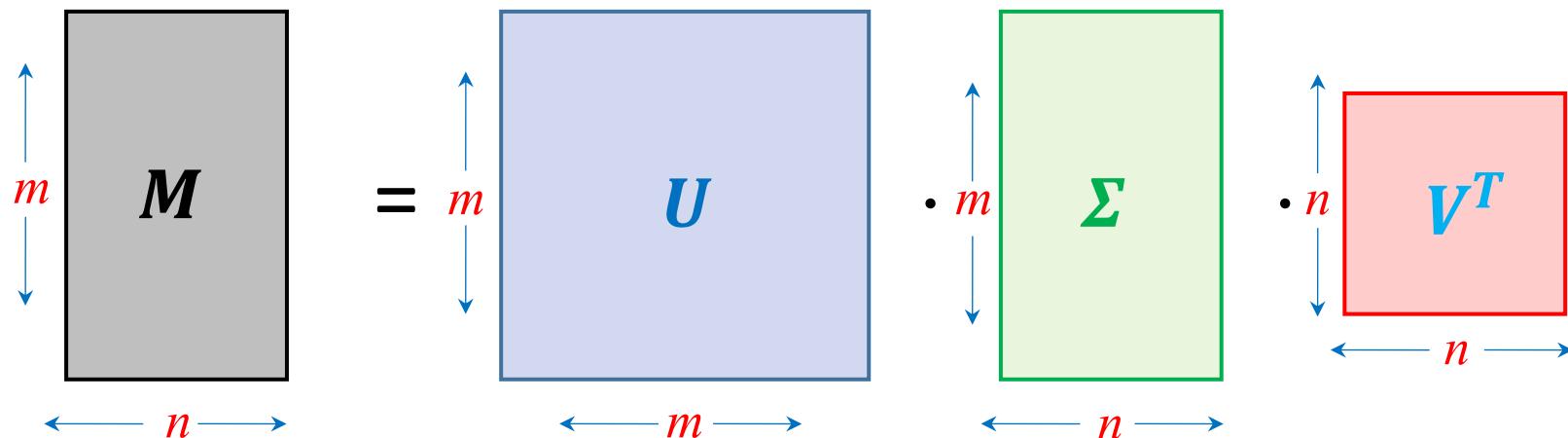
$$\text{Size}(\mathbf{U})=m \times m$$

$$\text{Size}(\Sigma)=m \times n$$

$$\text{Size}(\mathbf{V})=n \times n$$

| Singular Value Decomposition (SVD)

- ▶ A matrix M is decomposed as $M = U\Sigma V^T$.



| Singular value decomposition (SVD)

- ▶ A matrix \mathbf{M} is decomposed as $\mathbf{M} = \mathbf{U}\Sigma\mathbf{V}^T$.
- ▶ Here, Σ contains the singular values as diagonal elements.
- ▶ These singular values are ordered from the largest to the smallest: $\sigma_1 > \sigma_2 > \dots > \sigma_m$.

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 & 0 \\ 0 & \sigma_2 & \cdots & 0 & 0 \\ \vdots & \ddots & \ddots & & \vdots \\ 0 & 0 & \cdots & \sigma_m & 0 \end{bmatrix}$$

| Singular value decomposition (SVD)

- ▶ A matrix \mathbf{M} is decomposed as $\mathbf{M} = \mathbf{U}\Sigma\mathbf{V}^T$.
- ▶ The columns of \mathbf{U} are the “left singular vectors.”
- ▶ The columns of \mathbf{V} are the “right singular vectors.”

$$= \begin{bmatrix} \uparrow & \cdots & \uparrow \\ \mathbf{u}_1 & \cdots & \mathbf{u}_m \\ \downarrow & \cdots & \downarrow \end{bmatrix}$$

$$\mathbf{V} = \begin{bmatrix} \uparrow & \cdots & \uparrow \\ \mathbf{v}_1 & \cdots & \mathbf{v}_n \\ \downarrow & \cdots & \downarrow \end{bmatrix}$$

- ▶ => Between a set of the left and right singular vectors and their singular value we have:

$$\mathbf{M} \mathbf{v}_i = \sigma_i \mathbf{u}_i$$

- ▶ => Between any two singular vectors, we have the following orthogonality condition:

$$\begin{aligned} \mathbf{v}_i \cdot \mathbf{v}_j &= \delta_{ij} & \Leftrightarrow \quad \mathbf{V}\mathbf{V}^T &= \mathbf{V}^T\mathbf{V} = \mathbf{I} \\ \mathbf{u}_i \cdot \mathbf{u}_j &= \delta_{ij} & \Leftrightarrow \quad \mathbf{U}\mathbf{U}^T &= \mathbf{U}^T\mathbf{U} = \mathbf{I} \end{aligned}$$

| Truncated SVD and topic vector

- ▶ If r = number of topics, a dimension is reduced as following:

$$\text{Size}(\mathbf{M})=m \times n$$

$$\text{Size}(\mathbf{U})=m \times m \quad \rightarrow \text{reduced to } m \times r$$

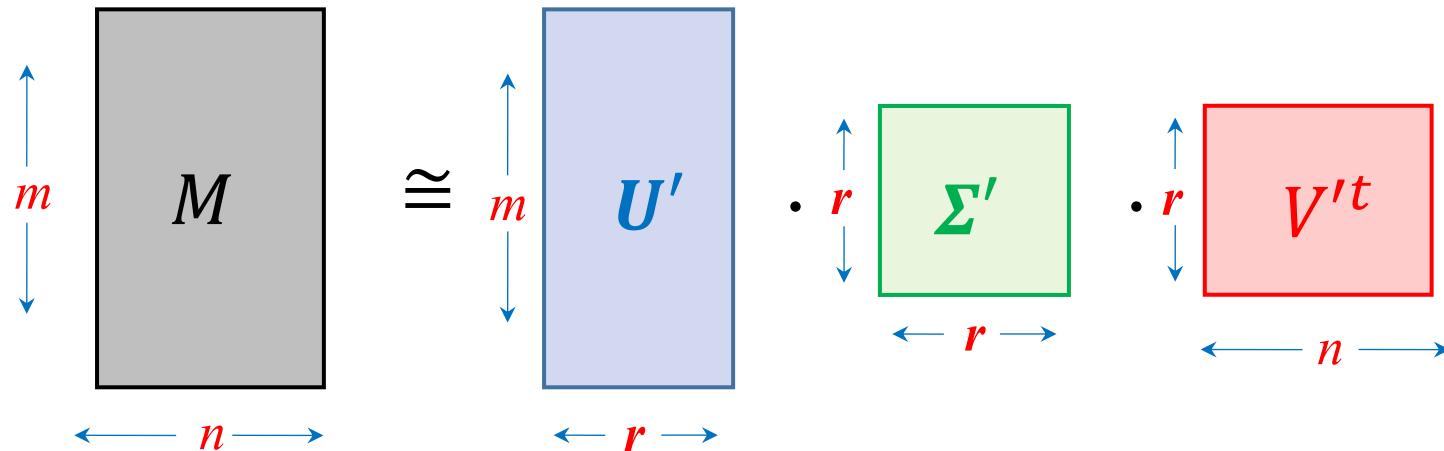
$$\text{Size}(\boldsymbol{\Sigma})=m \times n \quad \rightarrow \text{reduced to } r \times r$$

$$\text{Size}(\mathbf{V}) r=n \times n \quad \rightarrow \text{reduced to } n \times r$$

- ▶ We do not need all the possible topics. We only need a few “important” topics.

| Truncated SVD and topic vector

- ▶ If r = number of topics, the dimensional reduction as following is carried out:



- ▶ The columns of V' are the topic vectors.

Coding Exercise #0512



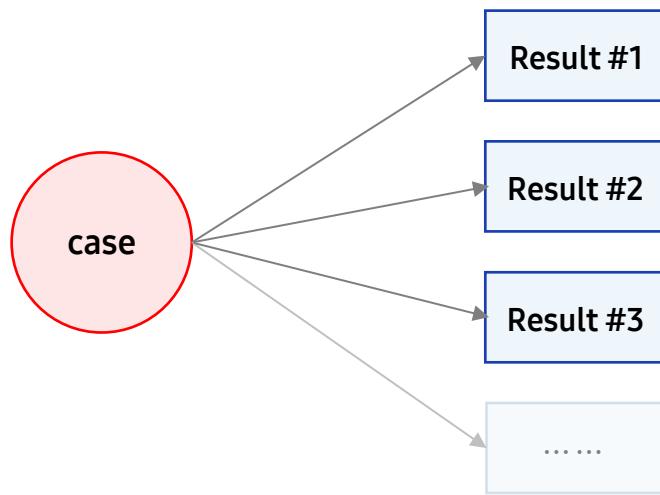
Follow practice steps on 'ex_0512.ipynb' file

Latent Semantic Analysis (LDA)

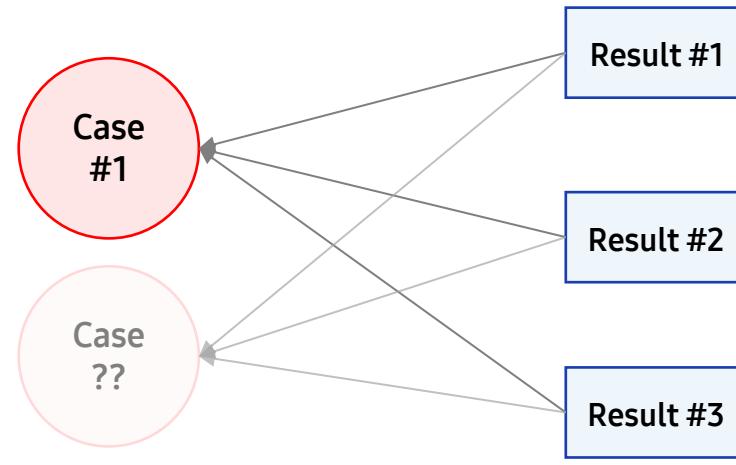
About the LDA

- ▶ Developed by D. Blei, A. Ng, and M. I. Jordan.
- ▶ One of the most representative topic modeling algorithms.
- ▶ Achieves clustering by calculating the topic distributions.
- ▶ Based on a Bayesian model.
- ▶ More information can be found at:
<http://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf>
<https://ai.stanford.edu/~ang/papers/nips01-lda.pdf>

Deductive reasoning vs. Inductive reasoning



Deductive reasoning



Inductive reasoning

- ▶ In inductive reasoning, the results are observed, and the causes are hidden.

LDA algorithm

- ▶ The accumulated counts of random experiments with more than two possible discrete outcomes can be described by a Multinomial distribution.
- Ex** Rolling a dice n times, the count probability of the sides follows a Multinomial distribution.


$$\Leftarrow P(x_1, x_2, x_3, x_4, x_5, x_6; n, p_1, p_2, p_3, p_4, p_5, p_6)$$

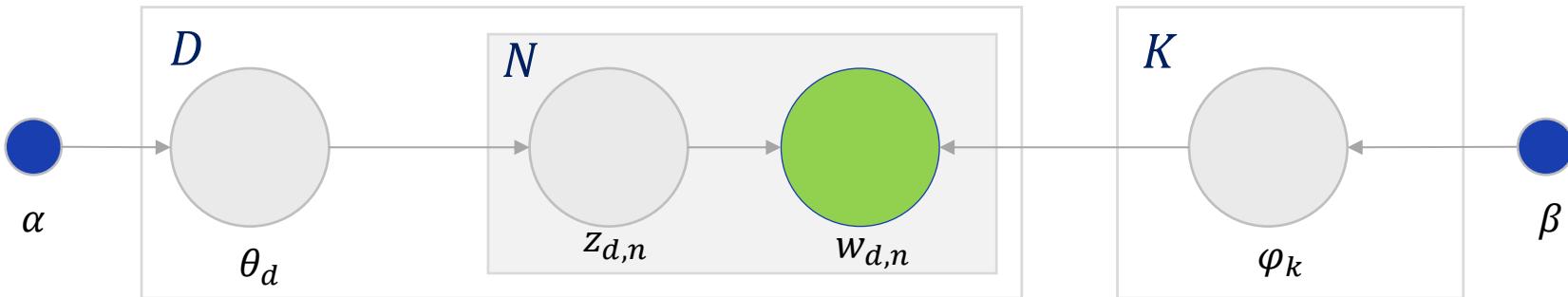
- ▶ In Bayes' theorem, the Dirichlet is “prior conjugate” of the Multinomial distribution.

Multinomial posterior \longrightarrow $P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$ \longleftarrow Dirichlet prior

| LDA algorithm

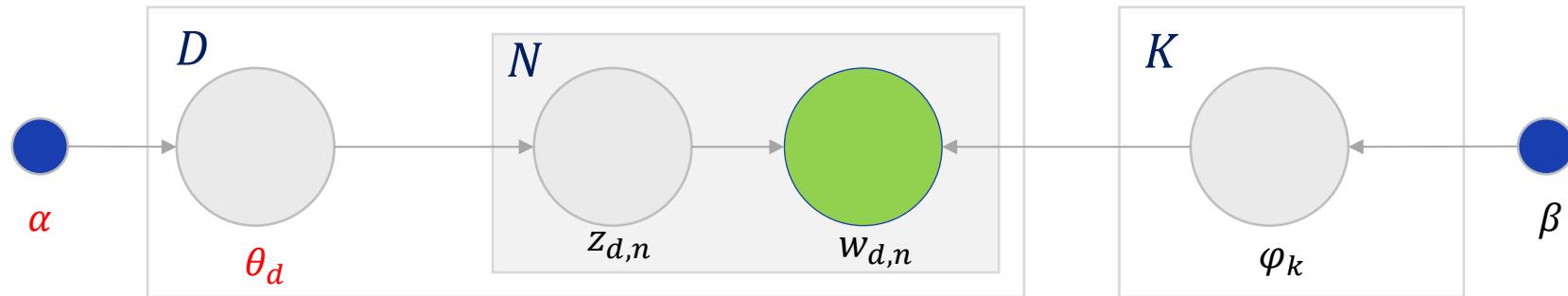
- ▶ In LDA, we assume that:
 - Document distribution follows a Multinomial with the Dirichlet prior conjugate.
 - Word distribution follows a Multinomial with the Dirichlet prior conjugate.
- ▶ What we can actually observe are the documents that contain words. The topics remain hidden.
- ▶ We also assume that the documents are “bags-of-words” where the ordering does not matter.
- ▶ By inductive reasoning (Bayes’ theorem), LDA extracts the distribution of the hidden topics.

| LDA algorithm

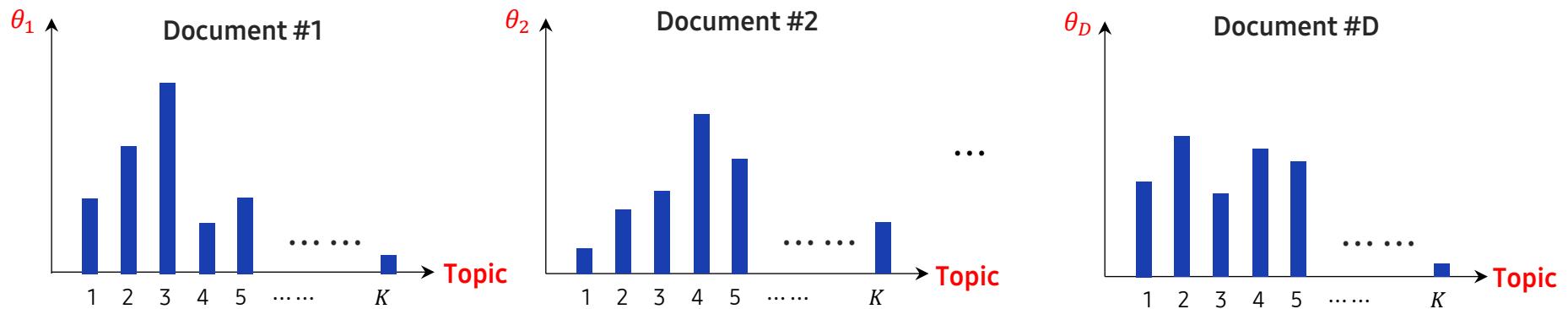


- ▶ Illustrated in a “plate notation.”
- ▶ Rectangles mean **repetitive steps**.
- ▶ We assume that there are in total:
 - D documents such that $d=1,\dots,D$
 - N words such that $n=1,\dots,N$
 - K topics such that $k=1,\dots,K$
- ▶ Also, we assume that the number of different words (features) is F .

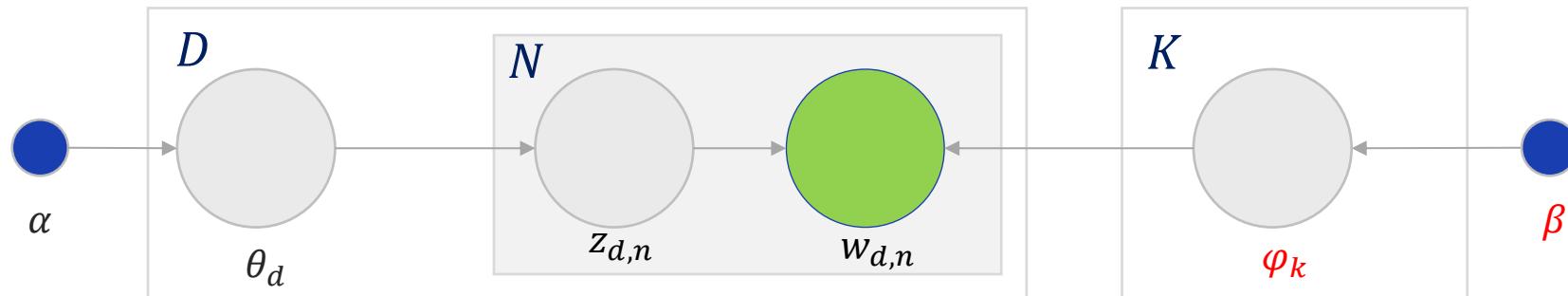
| LDA algorithm



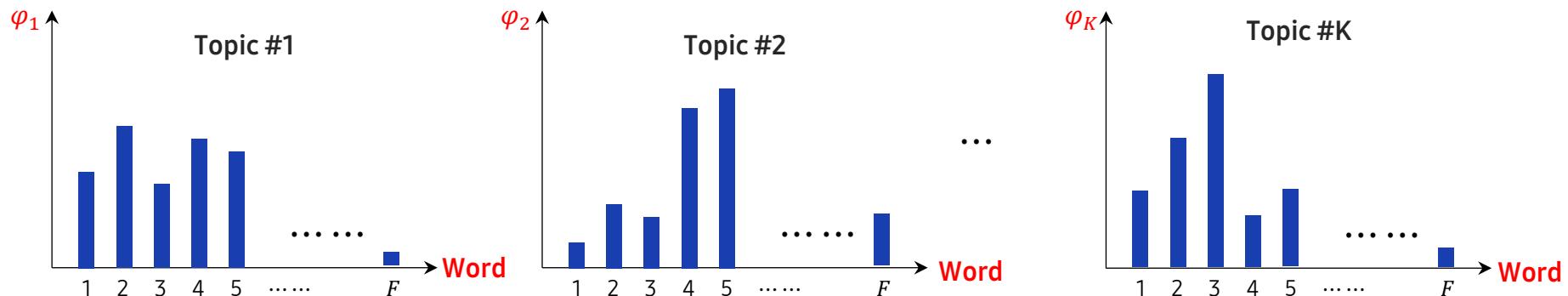
- ▶ α is the parameter of Dirichlet prior for the per-document topic distribution θ_d (with $d=1,\dots,D$).



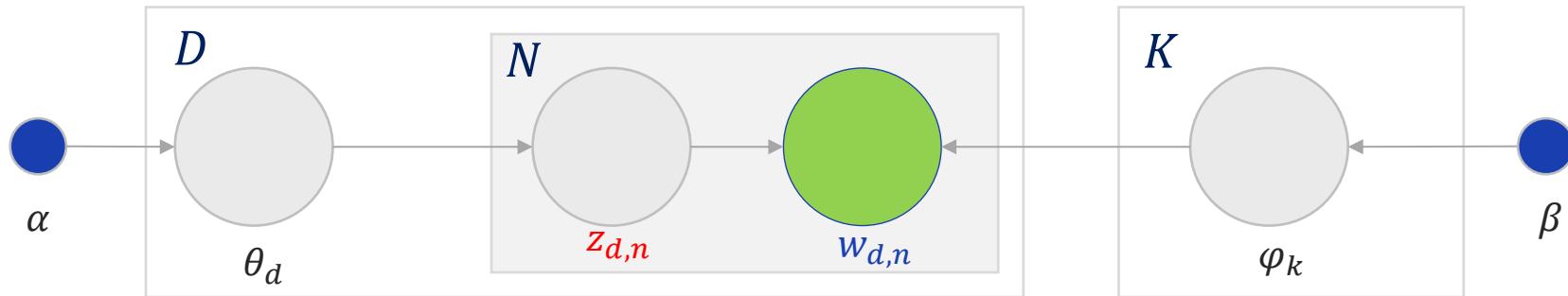
| LDA algorithm



- ▶ β is the parameter of Dirichlet prior for the per-topic word distribution φ_k (with $k=1,\dots,K$).



| LDA algorithm

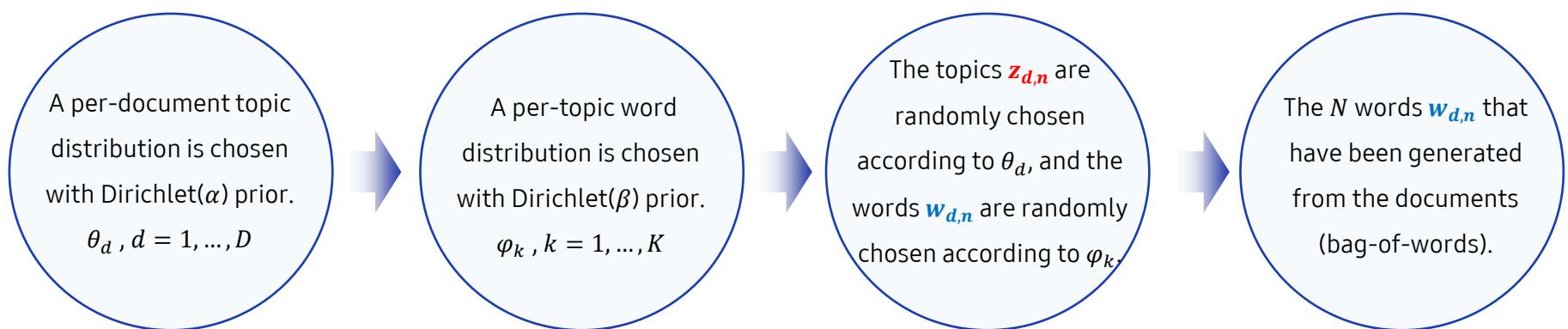


- ▶ $z_{d,n}$ is the topic for the n-th word in the d-th document. It is an integer from the range $[1,2,\dots,K]$.
- ▶ $w_{d,n}$ is the n-th word in the d-th document. It is an integer from the range $[1,2,\dots,F]$.

Generative process

- ▶ Corresponds to **deductive** reasoning.
- ▶ Suppose that we know the model.
- ▶ Documents that are word collections (bag-of-words) can be generated as following:

K = number of topics, D = number of documents, N = number of words.



Inference of the topics

- ▶ Corresponds to **inductive** reasoning.
- ▶ We know neither the model nor the topics.
- ▶ We would like to infer the topics from the dataset (documents, that is, word collections).
- ▶ The generative model was based on the joint probability $P(W, Z, \theta, \varphi; \alpha, \beta)$.
- ▶ Steps for the inference:
 - 1) Marginalize out the θ and φ . Get $P(W, Z; \alpha, \beta) = \int \int (W, Z, \theta, \varphi; \alpha, \beta) d\theta d\varphi$
 - 2) Using the dataset (W), get the probability $P(Z | W; \alpha, \beta)$.
 - 3) Infer about the α and β .

Coding Exercise #0513



Follow practice steps on 'ex_0513.ipynb' file

Unit 4.

Natural Language Processing with Keras

4.1. Natural Language Processing with Keras

Natural Language Processing with Keras

One-hot-encoding vs. Word embedding

- ▶ We notice obvious problems with the one-hot-encoding representation.
- ▶ To improve, “word embedding” is introduced, which is a distributed representation method.

One-Hot-Encoding	Embedding
The dimension of the vector space is large. The dimension is as large as the vocabulary size.	The dimension of the vector space is limited.
Vectors are sparse and mostly filled with 0s that carry no information.	Vectors are dense. Every vector element carries some information.
No semantic relationship among the vectors. The vectors are orthogonal to each other.	Semantic relationship among the vectors.

- ▶ There are also “paragraph embedding” and “document embedding” representations.
- ▶ We will call “dense vector” or “embedding vector” interchangeably.

One-hot-encoding vs. Word embedding

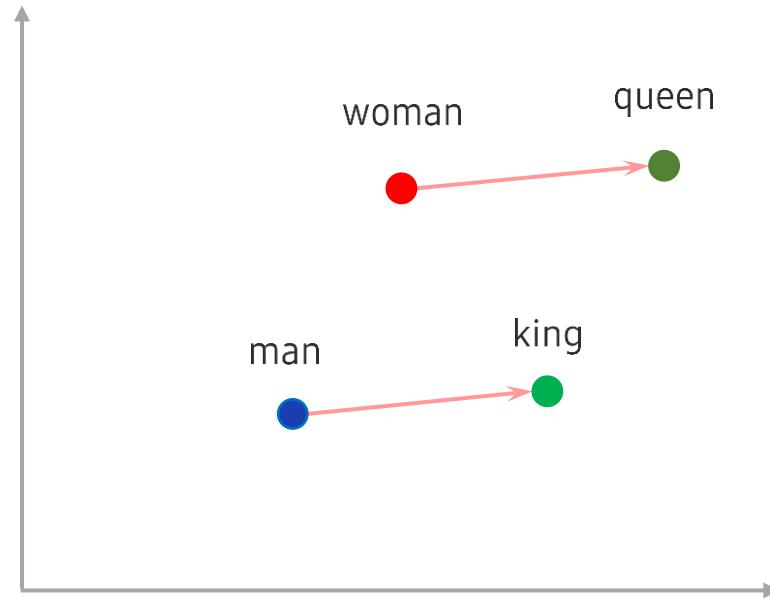
Ex Given a sentence, “I eat an apple every morning,” suppose that the words are indexed as:

I	:	3
Eat	:	0
An	:	2
Apple	:	1
Every	:	4
Morning	:	5

The words would have the following one-hot-encoding representations:

I	:	[0 0 0 1 0 0]
Eat	:	[1 0 0 0 0 0]
An	:	[0 0 1 0 0 0]
Apple	:	[0 1 0 0 0 0]
Every	:	[0 0 0 0 1 0]
Morning	:	[0 0 0 0 0 1]

Word embedding (Word2Vec)



- Among the dense vectors, relationships such as following are established:

$$\text{queen} - \text{woman} = \text{king} - \text{man}$$

Word embedding (Word2Vec)

- ▶ Use CBOW (Continuous Bag of Words) and/or Skip-Gram to build the embedding vectors:
 - 1) Build a predictive model based on the Softmax regression (multi-class logistic regression).
 - 2) We assume one-hot-encoded input and output vectors.
 - 3) Extract the embedding vectors from the trained weight matrices.

Word embedding (Word2Vec)

- ▶ CBOW: Using the context words, predict the (missing) center word.

Training Sentence	Center Word	Context Words
I eat an apple every morning	I	eat
I eat an apple every morning	eat	I, an
I eat an apple every morning	an	eat, apple
I eat an apple every morning	apple	an, every
I eat an apple every morning	every	apple, morning
I eat an apple every morning	morning	every

- ▶ We assumed a “sliding window” over the training sentence.

| Word embedding (Word2Vec)

- ▶ CBOW: Using the context words, predict the (missing) center word.

Ex Suppose that the vector dimension of the one-hot-encoded words = 6.

Let's also suppose we would like to find dense vectors of dimension = 3.

We will consider two context words (one from the left and another from the right).

So we have a situation as the following:

I eat an apple every morning



I eat an _____ ? every morning

Word embedding (Word2Vec)

- ▶ CBOW: Using the context words, predict the (missing) center word.

Ex (continues from the previous page)

Then, we build a Softmax regression model.

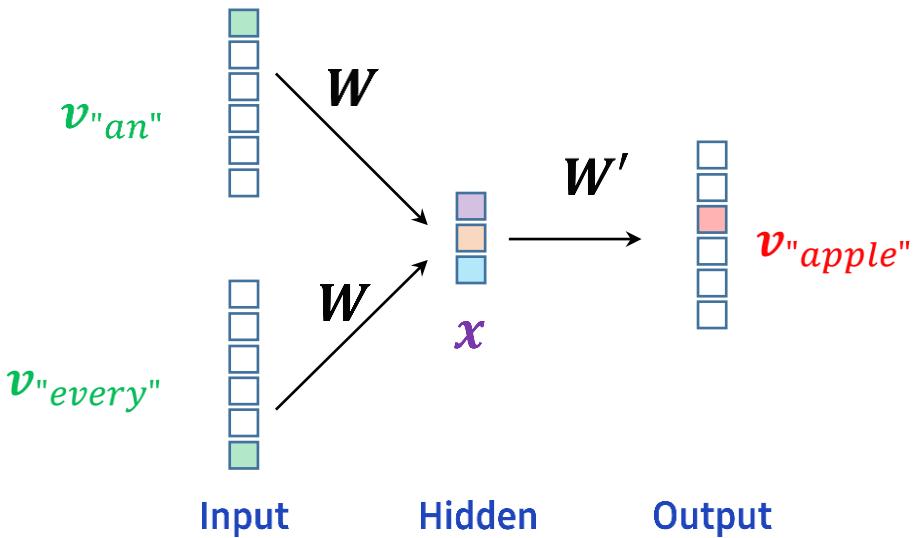
For the vector inputs of “an” and “every,”

we would like to train the weights W

and W' such that the predicted output is

the vector “apple.”

One-hot-encoded words : $v^{“an”} = [1 0 0 0 0 0]$, $v^{“every”} = [0 0 0 0 0 1]$, $v^{“apple”} = [0 0 1 0 0 0]$



Word embedding (Word2Vec)

- ▶ CBOW: Using the context words, predict the (missing) center word.

Ex (continues from the previous page)

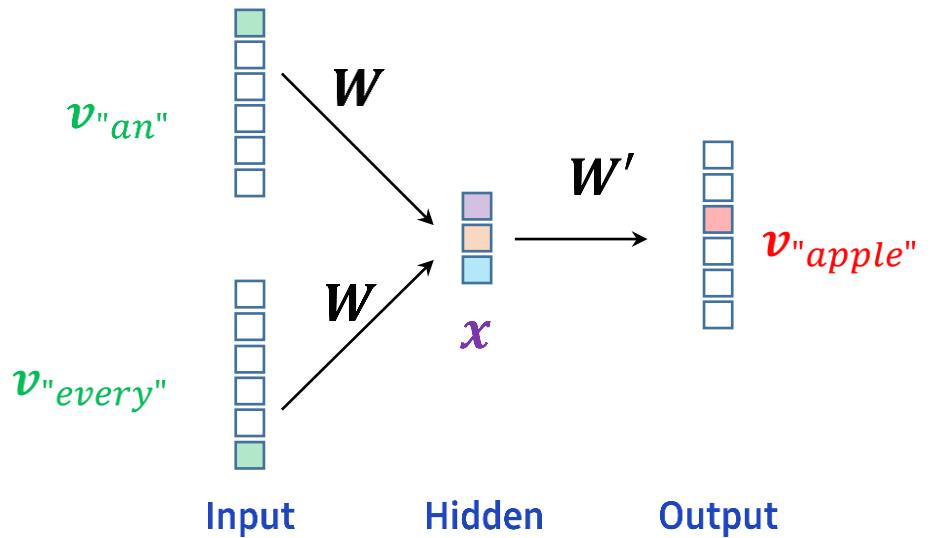
We have the following sizes:

Size of the matrix $W=3 \times 6$

Size of the matrix $W'=6 \times 3$

Dimension of the vector $x= 3$

Dimension of the input and output = 6



Word embedding (Word2Vec)

- ▶ CBOW: Using the context words, predict the (missing) center word.

Ex (continues from the previous page)

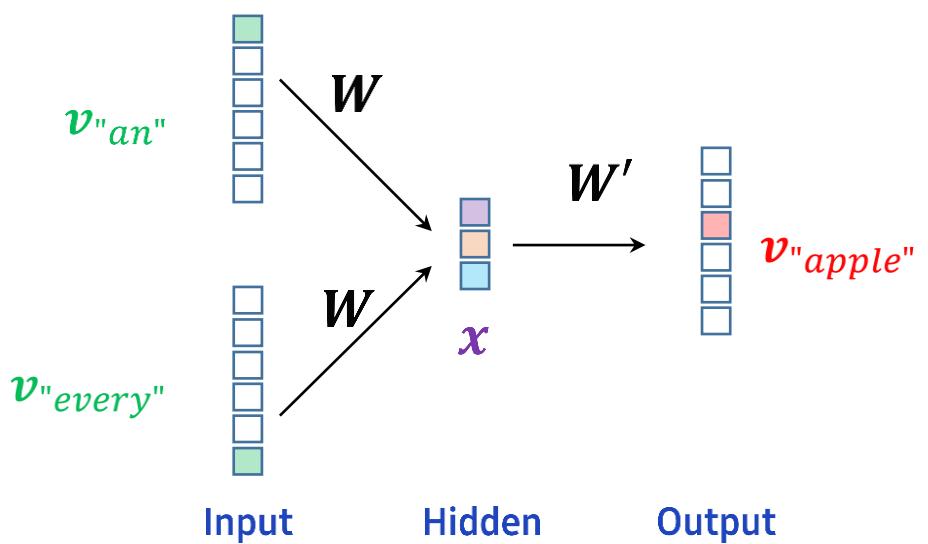
We propagate forward from the input layer

to the hidden layer (a single node):

$$x^{an} = W \cdot v^{an}$$

$$x^{every} = W \cdot v^{every}$$

$$x = \frac{x^{an} + x^{every}}{2} \quad \leftarrow \text{Average for the hidden node}$$



Word embedding (Word2Vec)

- ▶ CBOW: Using the context words, predict the (missing) center word.

Ex (continues from the previous page)

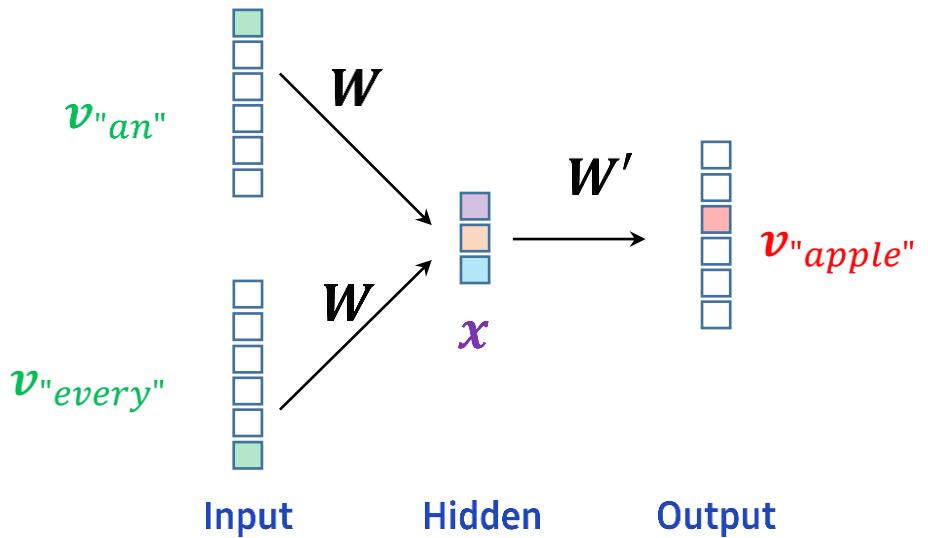
We propagate forward to the output layer:

$$\hat{v} = W' \cdot x$$

We should train the weights W and W' :

$\text{argmax}(\hat{v})$ and $\text{argmax}(v^{\text{"apple"}})$

are minimized.



| Word embedding (Word2Vec)

- ▶ CBOW: Using the context words, predict the (missing) center word.

Ex (continues from the previous page)

Now, let us interpret the result.

a) When we propagate from the input layer to the hidden layer (by matrix multiplication),

the one-hot-encoded input vectors $v_{\text{"an"}}$ and $v_{\text{"every"}}$ are picking the columns 0 and 5

of \mathbf{W} and projecting them to the hidden layer with the target dimension = 3.

b) So, the dense vectors for “an” and “every” are **columns** 0 and 5 of the **trained \mathbf{W}** .

c) Analogously, we can extract dense vectors from the **rows** of the **trained \mathbf{W}'** .

| Word embedding (Word2Vec)

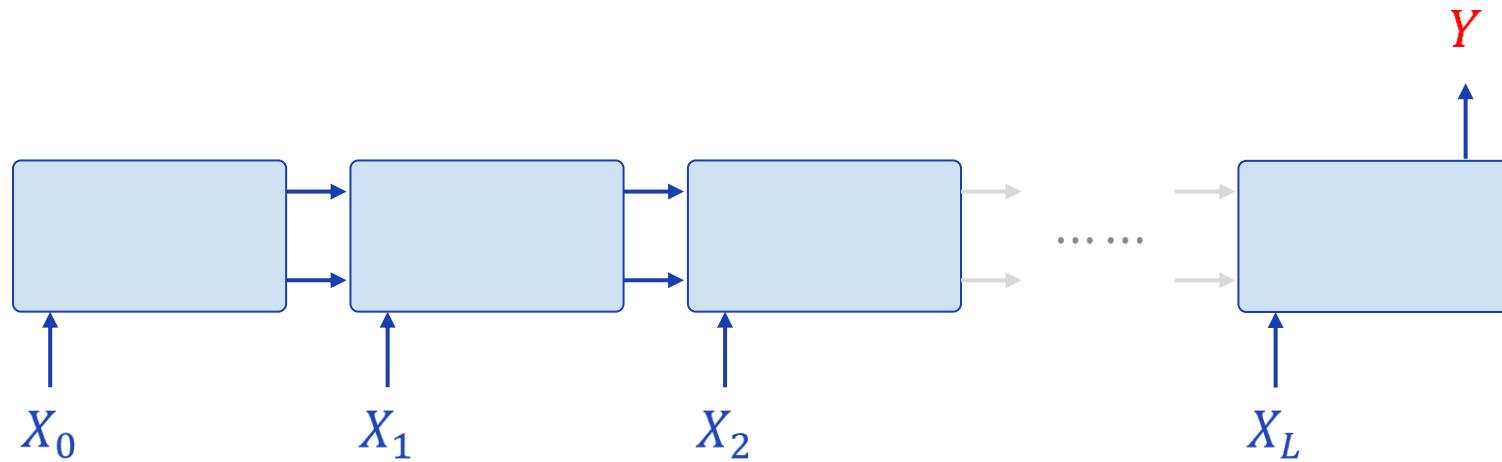
- ▶ Skip-Gram: Using a center word, predict the (missing) context words.

I eat **an apple every** morning
↓
I eat **? apple ?** morning

- Similar to the CBOW, we also train a Softmax regression to predict the missing words.
- We extract the dense vectors (embedding vectors) from the trained weight matrices.

LSTM network for document classification

- ▶ “Sequence in and Vector out” model
- ▶ Embedding representation of the words



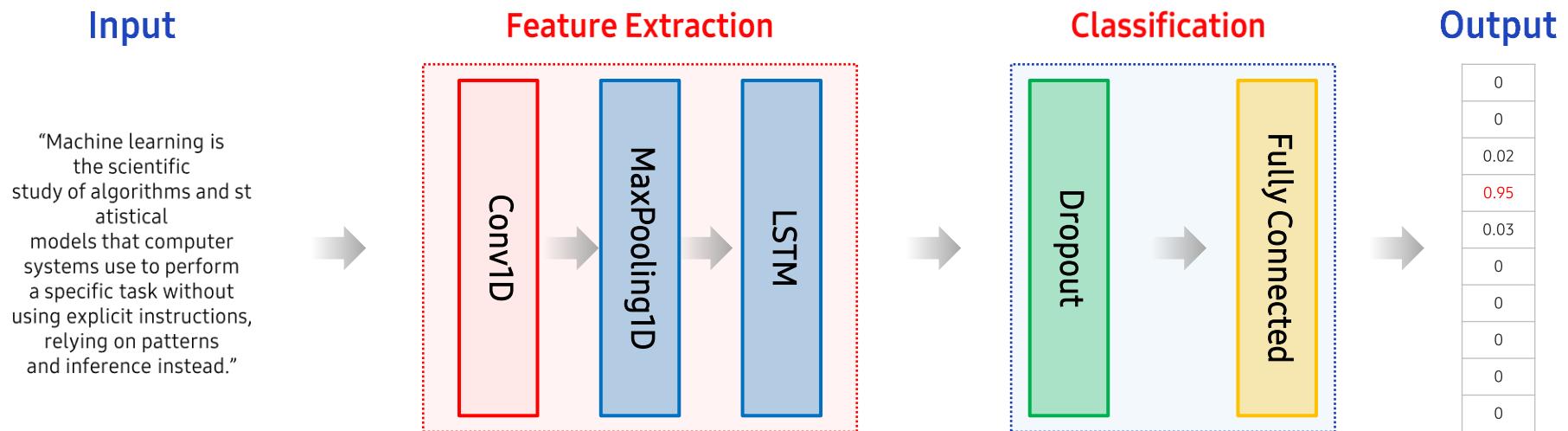
LSTM network for document classification: a code example

```
# Import the necessary classes.  
from tensorflow.keras.models import Sequential  
from tensorflow.keras.layers import Dense, LSTM, Embedding  
  
# We will use the Sequential API.  
  
n_words = 128  
n_input = 32  
n_neurons=64  
  
# Build a model by adding the layers.  
my_model = Sequential()  
my_model.add(Embedding(n_words, n_input))  
my_model.add(LSTM(units=n_neurons, return_sequences=False, input_shape=(None, n_input), activation='tanh'))  
my_model.add(Dense(1, activation='sigmoid'))
```

- ▶ In `LSTM()`, we should set `return_sequences=False` for a “Sequence in and Vector out” model.

Deep learning model for document classification

- ▶ 1D convolution + 1D max pooling + LSTM for the feature extraction
- ▶ Localized sequence patterns picked up by the 1D convolution



Deep learning model for document classification: a code example

```
# Import the necessary classes.  
from tensorflow.keras.models import Sequential  
from tensorflow.keras.layers import Dense, LSTM, Embedding, Conv1D, MaxPool1D, Dropout  
  
n_filters=32  
k_size=3  
stride_size=1  
hold_prob=0.5  
  
# Build a model by adding the layers.  
my_model = Sequential()  
my_model.add(Embedding(n_words, n_input)) # Embedding layer.  
my_model.add(Conv1D(filters=n_filters, kernel_size = k_size, strides=stride_size, padding='valid',activation='relu'))  
my_model.add(MaxPool1D(pool_size = 2))  
my_model.add(LSTM(units=n_neurons, return_sequences=False, input_shape=(None, n_input), activation='tanh'))  
my_model.add(Dropout(rate=hold_prob))  
my_model.add(Dense(1, activation='sigmoid'))
```

Coding Exercise #0514



Follow practice steps on 'ex_0514.ipynb' file

Coding Exercise #0515



Follow practice steps on 'ex_0515.ipynb' file

Coding Exercise #0516



Follow practice steps on 'ex_0516.ipynb' file

Coding Exercise #0517



Follow practice steps on 'ex_0517.ipynb' file

A photograph of a person working at a desk. They are wearing an orange long-sleeved shirt and are holding a brown paper coffee cup with a black lid in their left hand. In their right hand, they are holding a black pen and are pointing it towards a black computer keyboard. On the desk, there are two large computer monitors, one on the left and one above the keyboard. To the right of the keyboard, there is an open notebook with some handwritten notes. In the foreground, there is a stack of papers or books. The background shows a window with vertical blinds.

End of Document



Together for Tomorrow! Enabling People

Education for Future Generations

©2023 SAMSUNG. All rights reserved.

Samsung Electronics Corporate Citizenship Office holds the copyright of book.

This book is a literary property protected by copyright law so reprint and reproduction without permission are prohibited.

To use this book other than the curriculum of Samsung Innovation Campus or to use the entire or part of this book, you must receive written consent from copyright holder.