

LECTURER: TAI LE QUY

INTRODUCTION TO DATA SCIENCE

TOPIC OUTLINE

Introduction to Data Science

1

Data

2

Data Science in Business

3

Statistics

4

Machine Learning

5

Summary session

6

UNIT 2

DATA



On completion of this unit, you will have learned ...

- what is meant by data and information.
- the different types and shapes of data.
- the typical sources of data.
- the 5Vs of big data.
- the issues concerning data quality.
- the challenges associated with the data engineering process.



- What do “data” and “information” mean?
- What data types and data shapes are there?
- Where can data be collected?
- What are the important criteria for data quality?
- What are the issues concerning data quality? How to fix them?
- What are the challenges and benefits of data processing?

DEFINITION

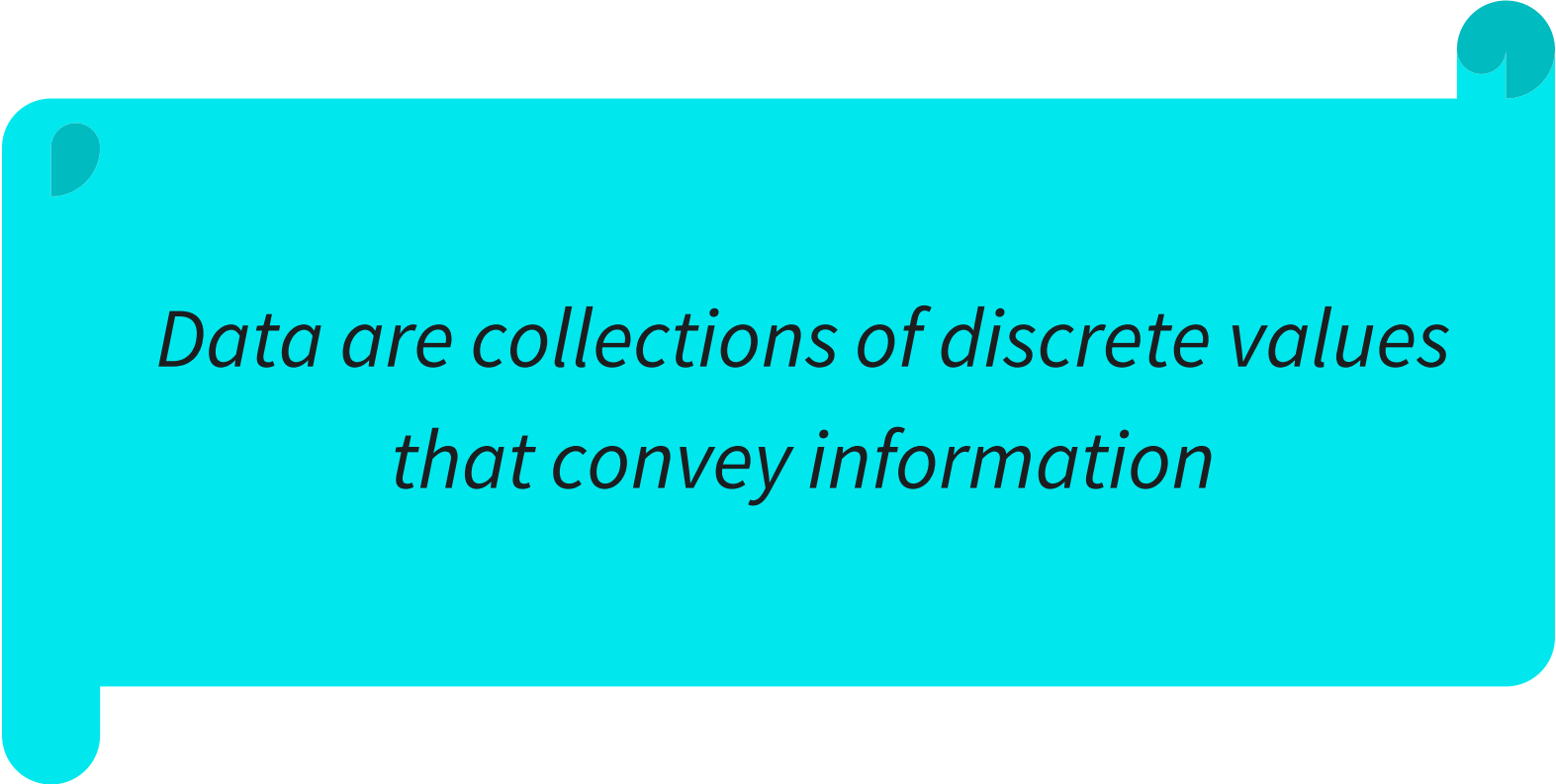
DATA

Facts, observations,
assumptions, or
incidences

Data describe
quantity, quality,
statistics, symbols, or
other units of meaning

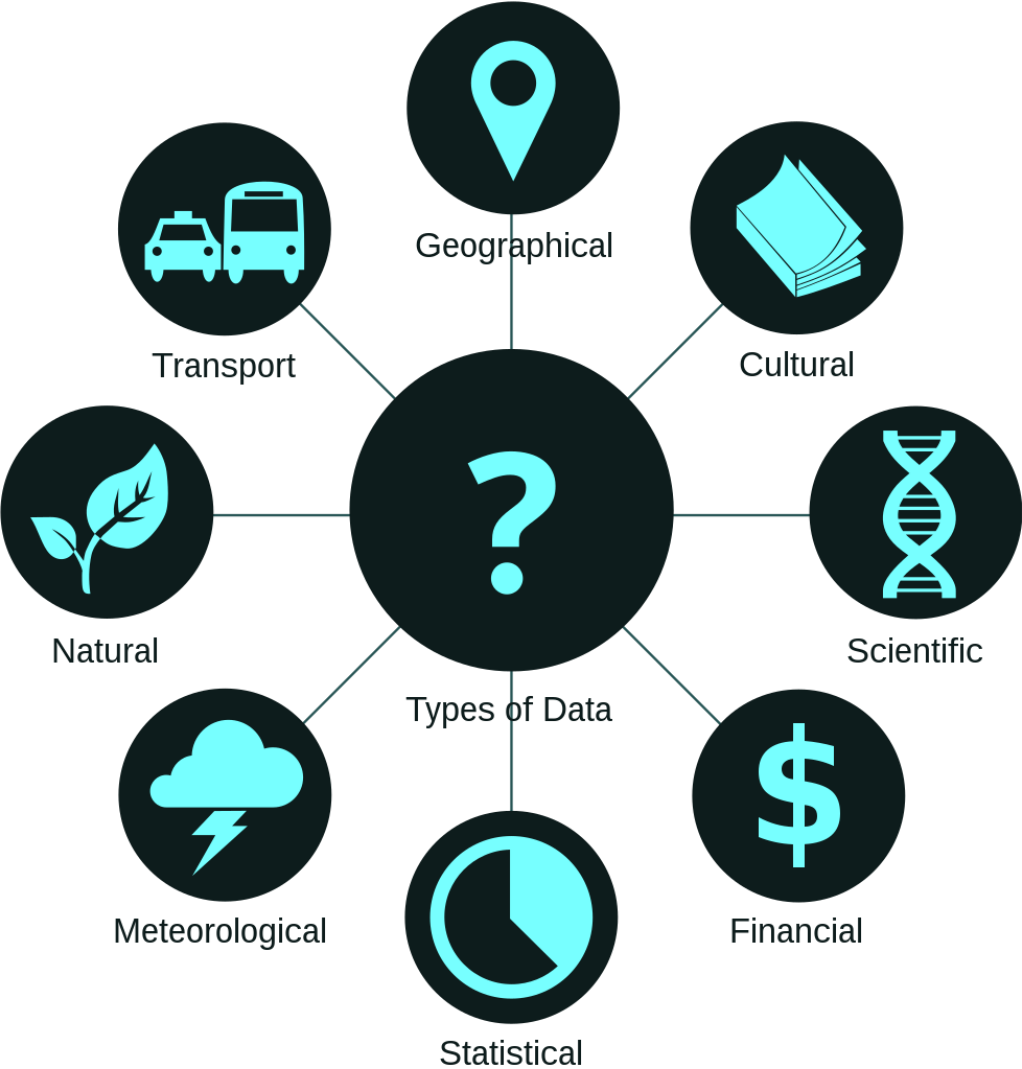
Data are processed to
return information

DEFINITION



*Data are collections of discrete values
that convey information*

TYPES OF DATA



Qualitative Data

Describe qualities or characteristics

Cannot be counted

Words, objects, pictures,
observations, and symbols

Answer to questions:
What characteristics? What property?

Identify conceptual
framework in an area of study

Quantitative Data

Can be quantified or
expressed as a number

Can be counted

Numbers and statistics

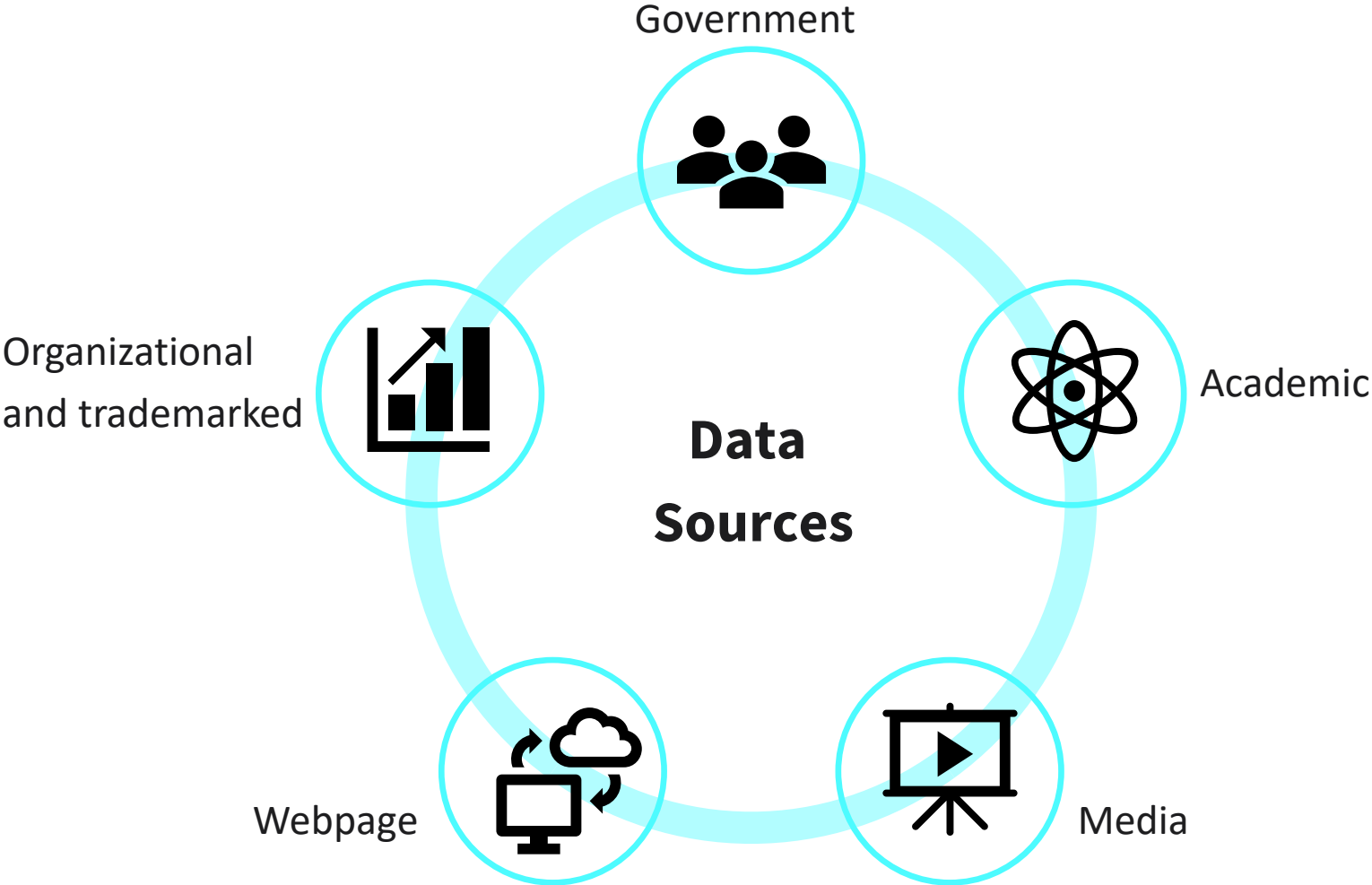
Answer to questions:
How much? How often?

Test hypotheses, analyses the
connection between cause and effect.

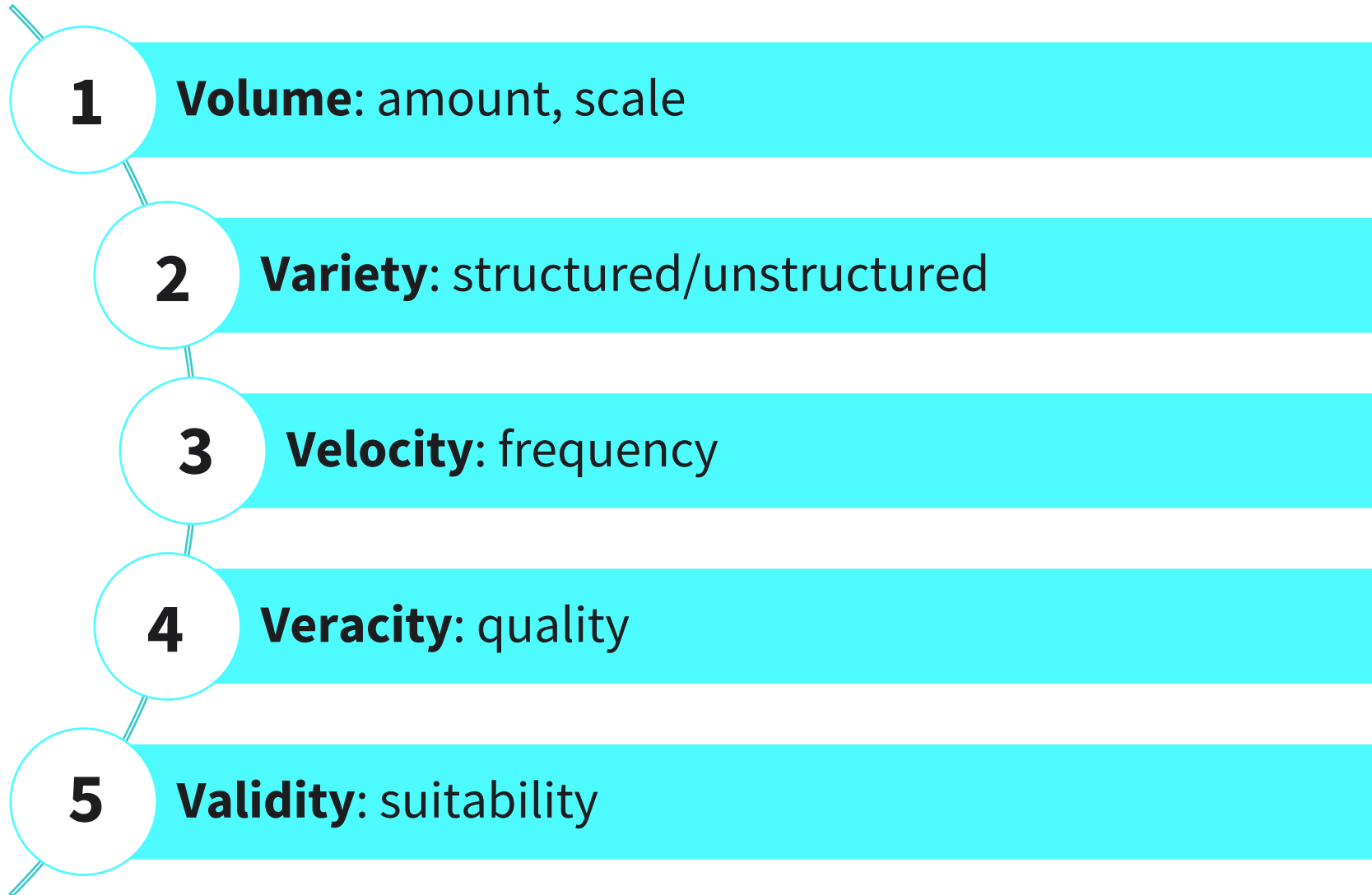
SHAPES OF DATA

	Structured Data	Unstructured Data	Streaming Data
Characteristics	<ul style="list-style-type: none">• predefined data models• usually, text or numerical• easy to search	<ul style="list-style-type: none">• no predefined data models• can be text, images, audio, or other formats• difficult to search	<ul style="list-style-type: none">• continuously generated• by sensors• large amount• processed incrementally
Applications	<ul style="list-style-type: none">• inventory control• airline reservation systems	<ul style="list-style-type: none">• word processing• tools for editing media	<ul style="list-style-type: none">• state monitoring• process control
Examples	<ul style="list-style-type: none">• phone numbers• customer names• transaction information	<ul style="list-style-type: none">• reports• imagery• email/messages	<ul style="list-style-type: none">• real-time sensor values• stock updates

SOURCES OF DATA

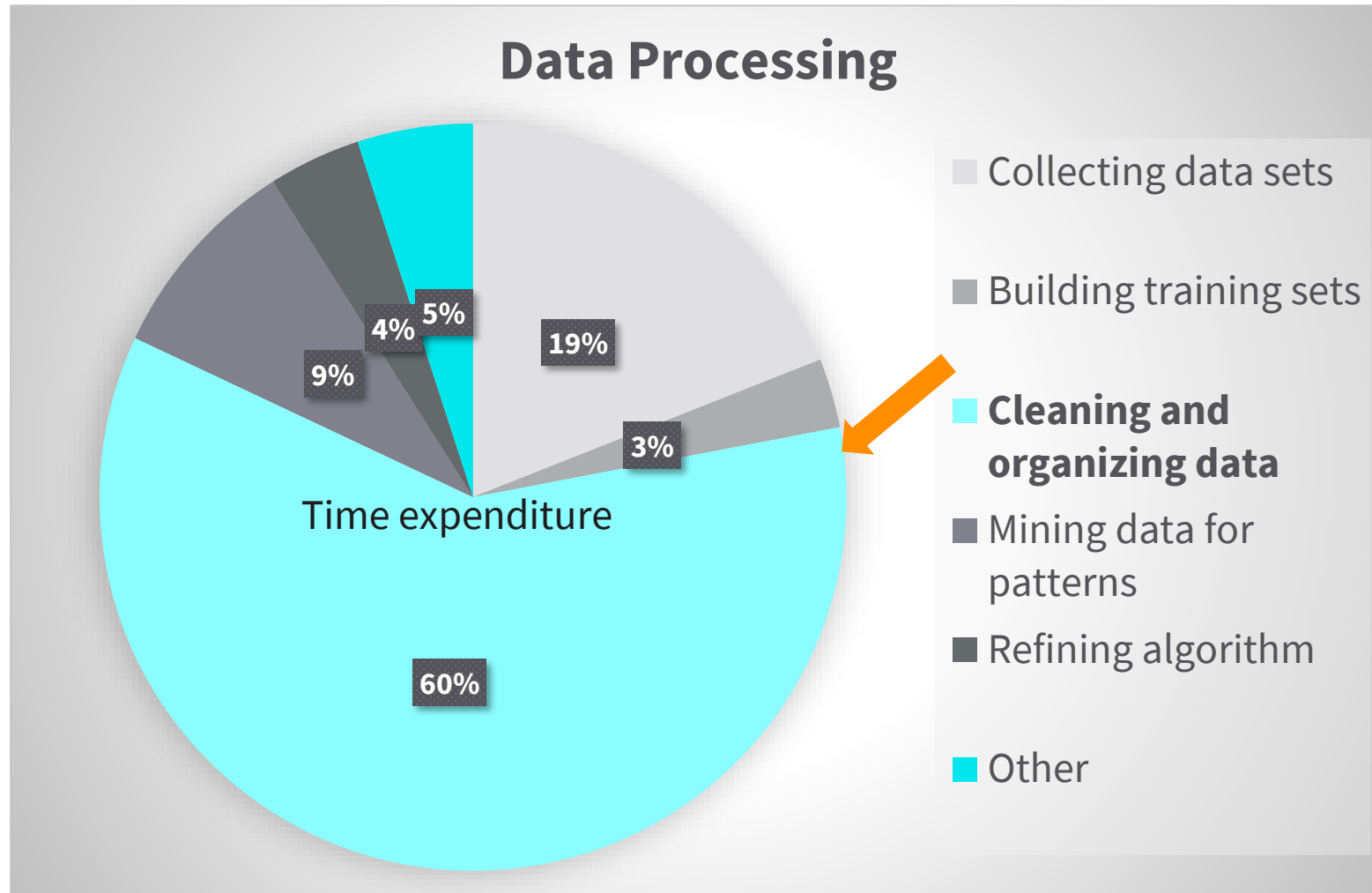


THE 5VS OF BIG DATA



DATA QUALITY

Issues	Causes	Solutions
Missing values and outliers	<ul style="list-style-type: none">- not observed- incorrectly observed	<ul style="list-style-type: none">- removal- replacement
Duplicate records	<ul style="list-style-type: none">- duplicate observation	<ul style="list-style-type: none">- removal
Redundancy	<ul style="list-style-type: none">- collected more than needed	<ul style="list-style-type: none">- correlation analysis- removal



Benefits of data processing

Improvement of **analysis** & **demonstration** of data

Reduction of meaningless data

Easier **storage** and **distribution** of data

Simplified **report** creation

Enhanced **productivity** and increased **profits**

Accurate **decision**-making

Data Transformation Methods

Variable scaling

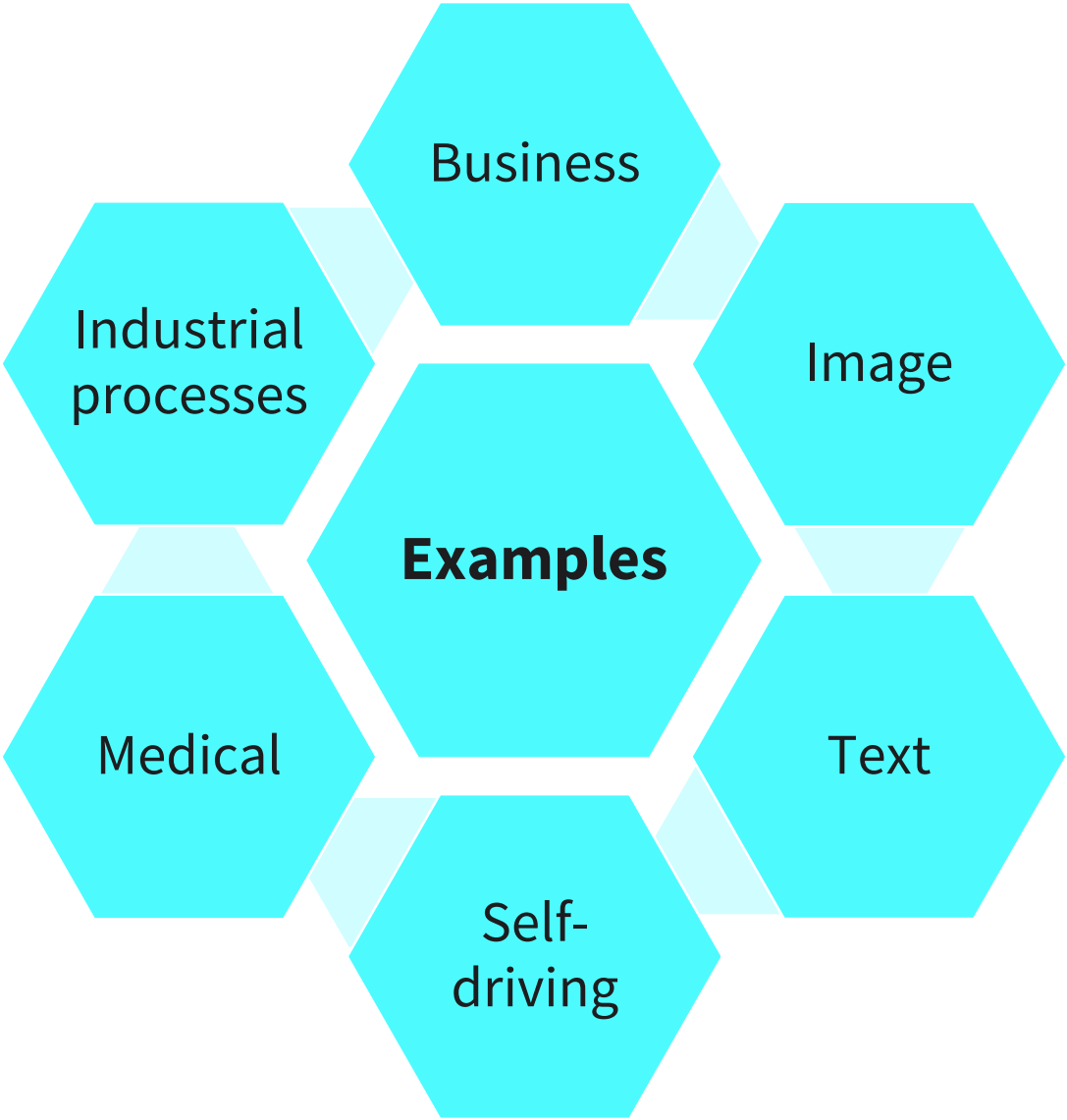
- Scaled values: $\{0,1\}$ or $\{-1,1\}$
- To ensure equal weighting

Variable decomposition

- One variable \rightarrow multiple variables
- To improve data representation

Variable aggregation

- Some variables \rightarrow one variable
- To reduce data volume





- what is meant by data and information
- the different types and shapes of data
- the typical sources of data
- the 5Vs of big data
- the issues concerning data quality
- the challenges associated with the data engineering process

SESSION 2

TRANSFER TASK

TRANSFER TASK
PRESENTATION OF THE RESULTS

Please present your
results.

The results will be
discussed in plenary.



Scenario

You are going to apply **data science** at an online shop to analyze and optimize the customer experience.

Questions

1. What data should be collected?
2. Which data from your dataset are quantitative, and which are qualitative?
3. Which types of data are there in your dataset?
4. Where can you collect your data?
5. How to guarantee the quality of your data? Give an example of your data quality issues and your solution.
6. How do you schedule your time for the task of data processing?
7. Which data transformation methods do you need for data engineering in your project?
8. Which benefits can you get from the data processing results?



1. Which shape of data has these characteristics: Easy to search, predefined data models, and numerical data?
 - a) structured data
 - b) unstructured data
 - c) streaming
 - d) all of the above



2. Which of the following is incorrect regarding data transformation methods?
- a) variable scaling
 - b) variable decomposition
 - c) variable aggregation
 - d) none of the above



3. Which activity consumes most of the data processing time?

- a) collecting data
- b) cleaning and organizing data
- c) building training sets
- d) mining data for patterns

LIST OF SOURCES

João Batista Neto. (2015). File: Data_types_-_en.svg. *Wikimedia Commons*. <https://en.wikipedia.org/wiki/Data#Meaning>

Wikipedia. (2023, May 14). Data. *Wikipedia Commons*. <https://en.wikipedia.org/wiki/Data#Meaning>

Zöllner, T. (2022). *Introduction to Data Science*. IU International University of Applied Science.

© 2022 IU Internationale Hochschule GmbH

This content is protected by copyright. All rights reserved.

This content may not be reproduced and/or electronically edited, duplicated, or distributed in any kind of form without written permission by the IU Internationale Hochschule GmbH.