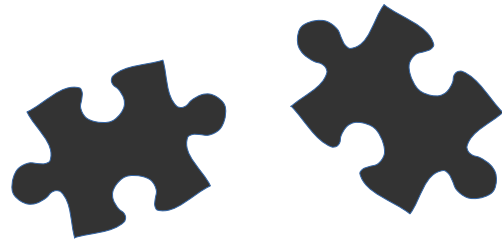


Data Analysis: Introduction

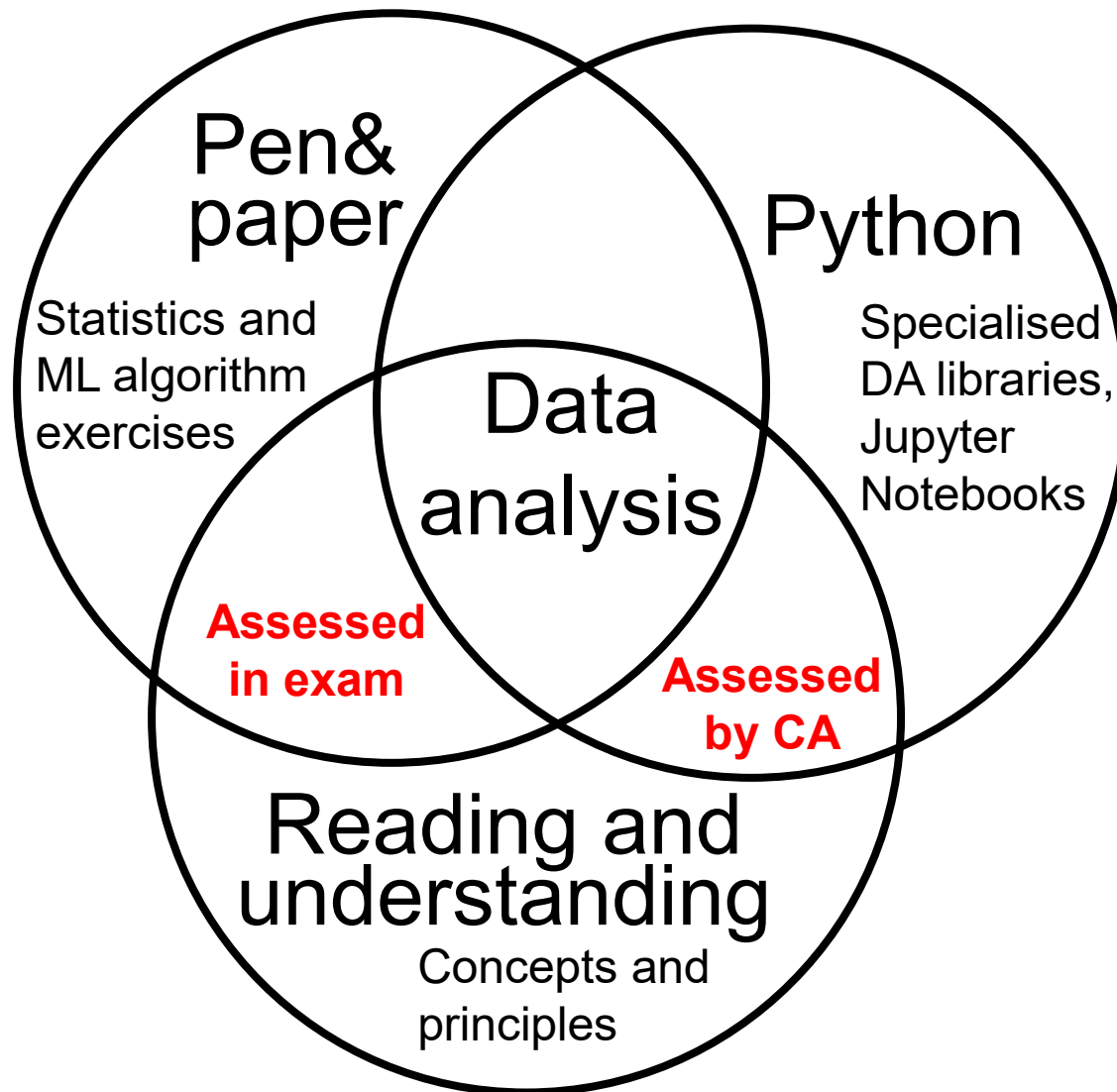
TU Dublin, Tallaght Campus

School of Enterprise Computing and Digital Transformation

In this module



you will find

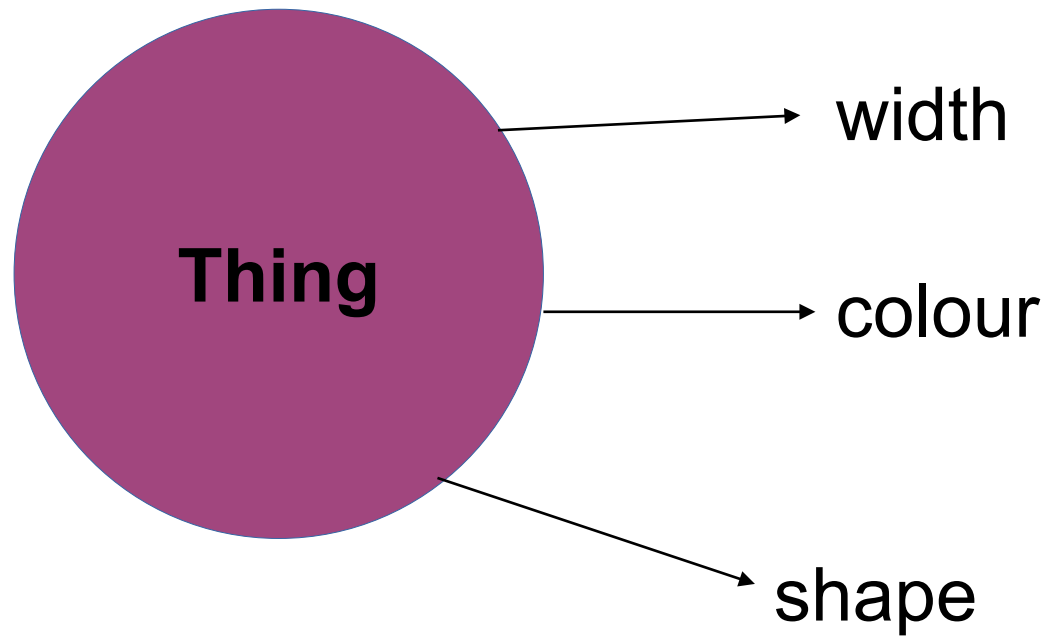


What are we
dealing with

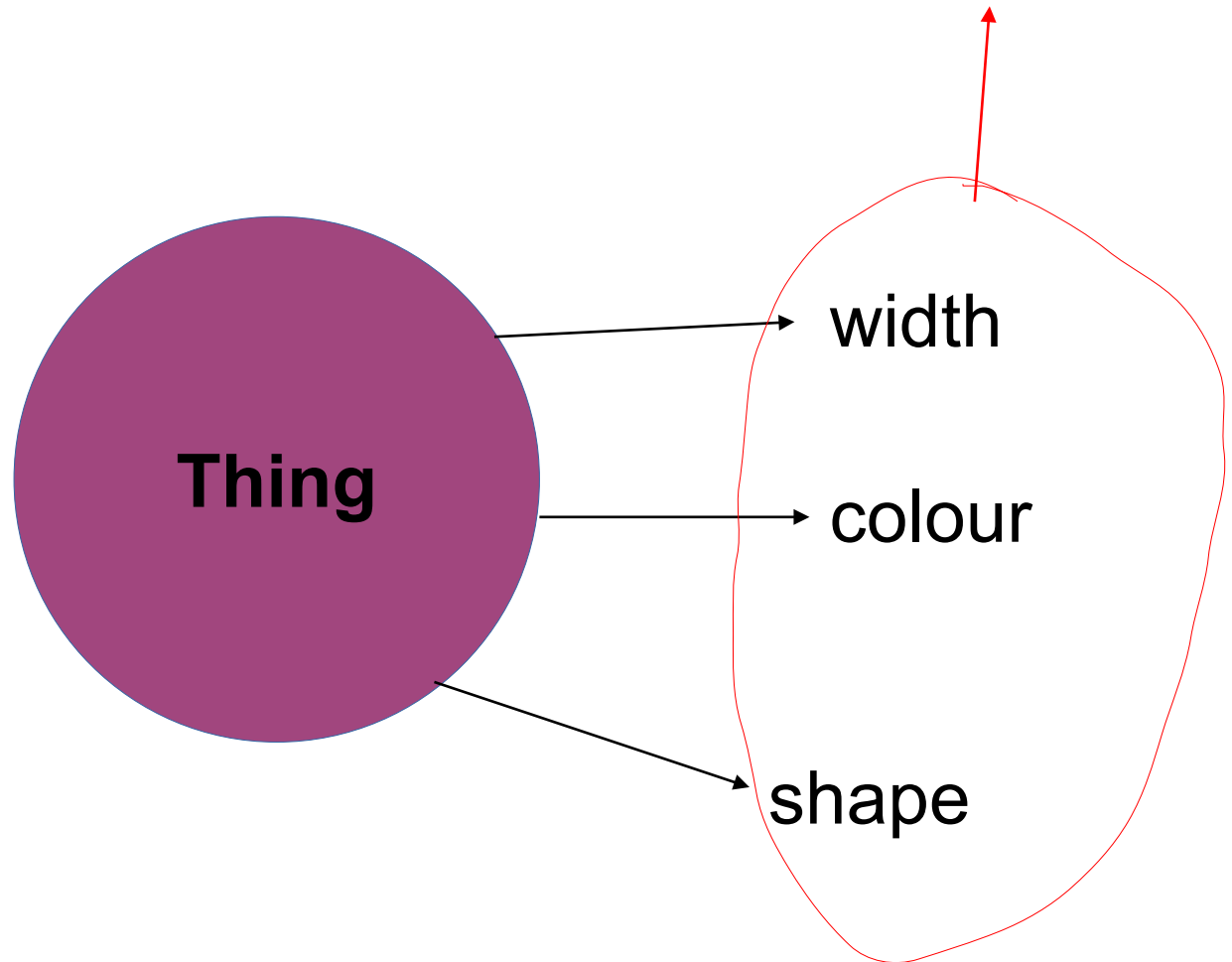




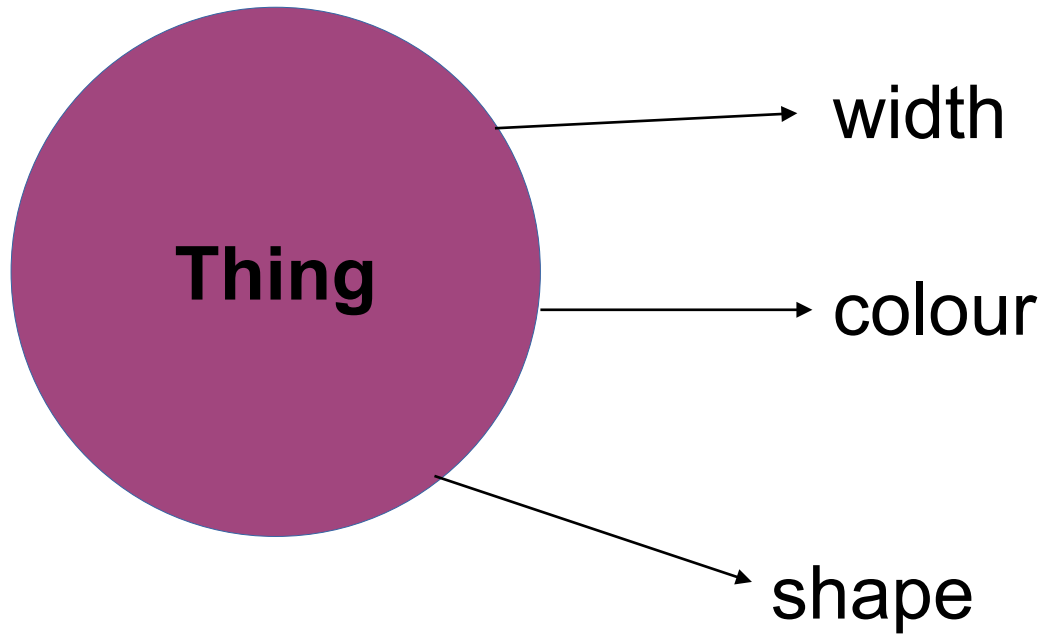
Thing



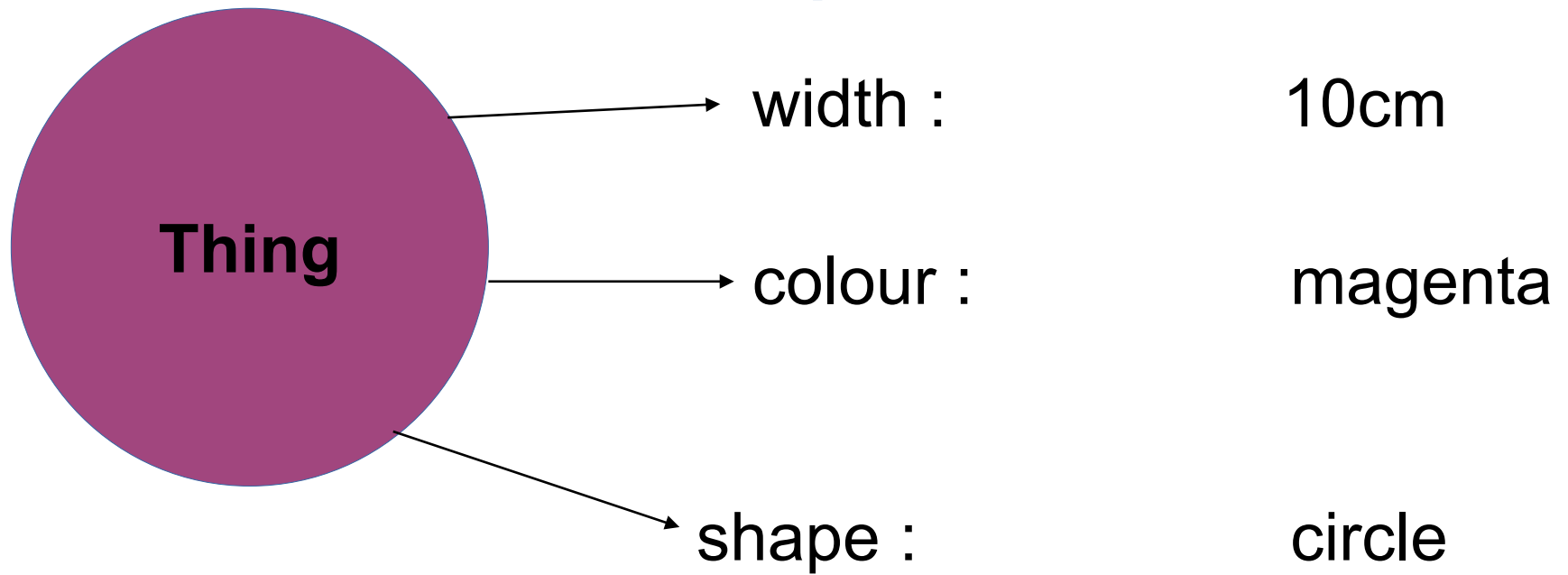
What are these?



Properties

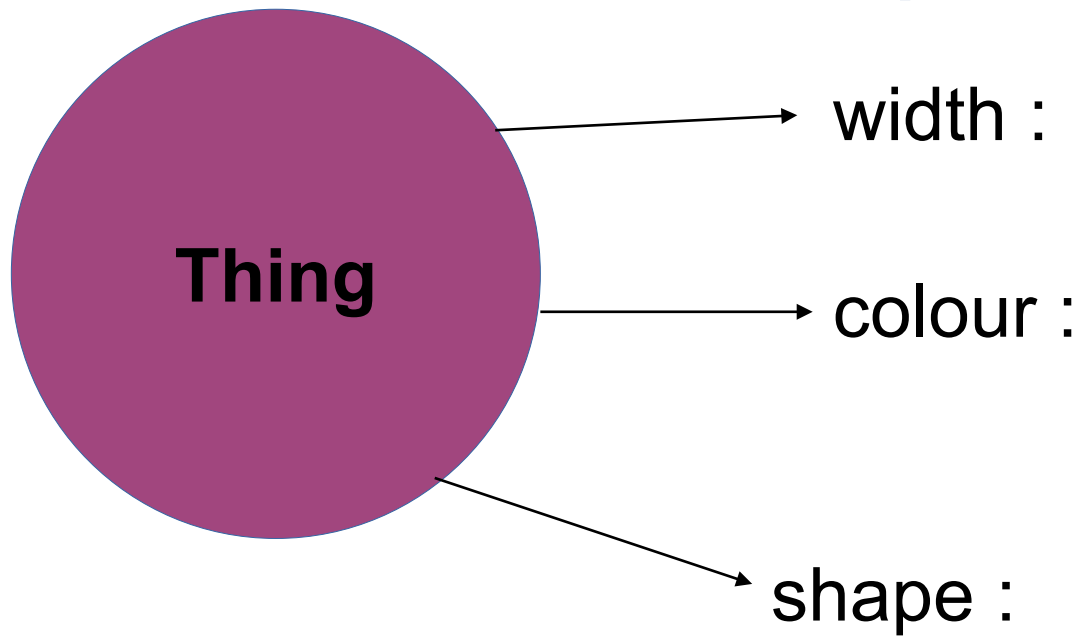


Properties



What are these?

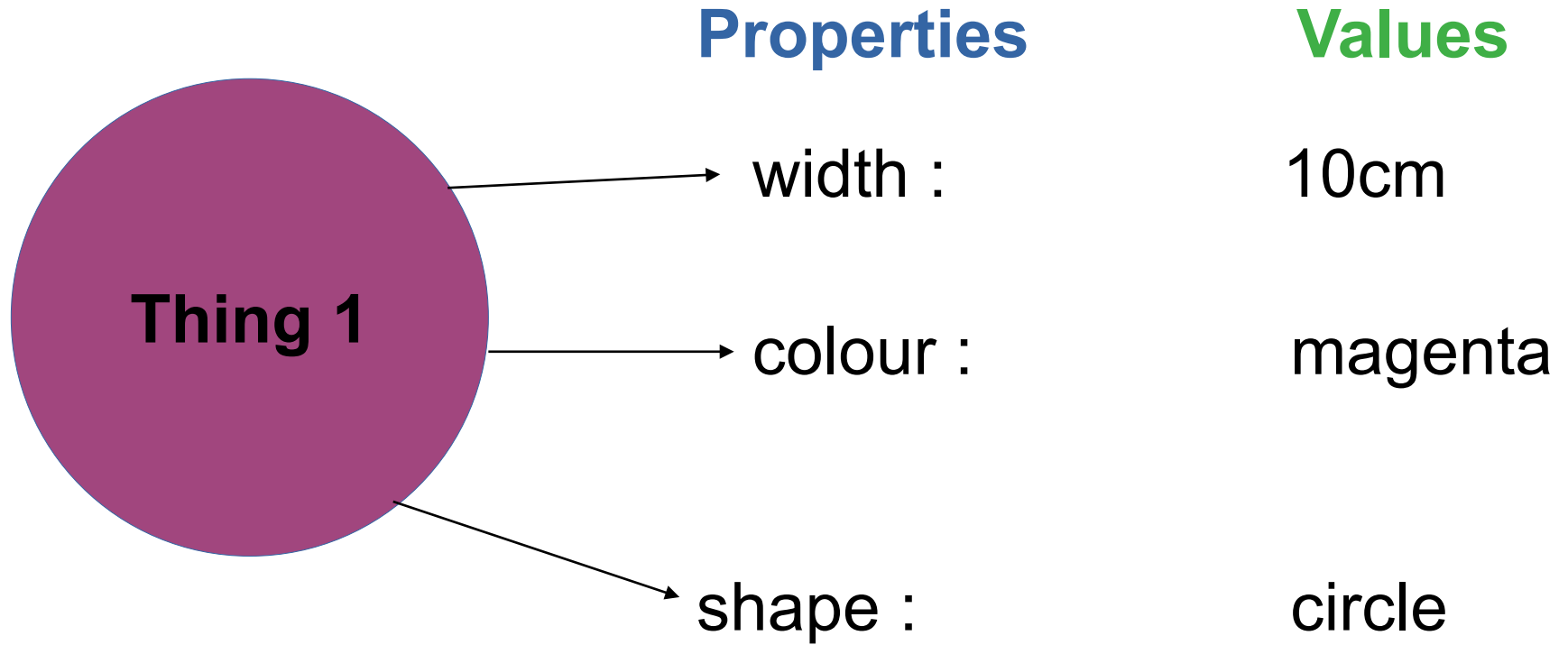
Properties

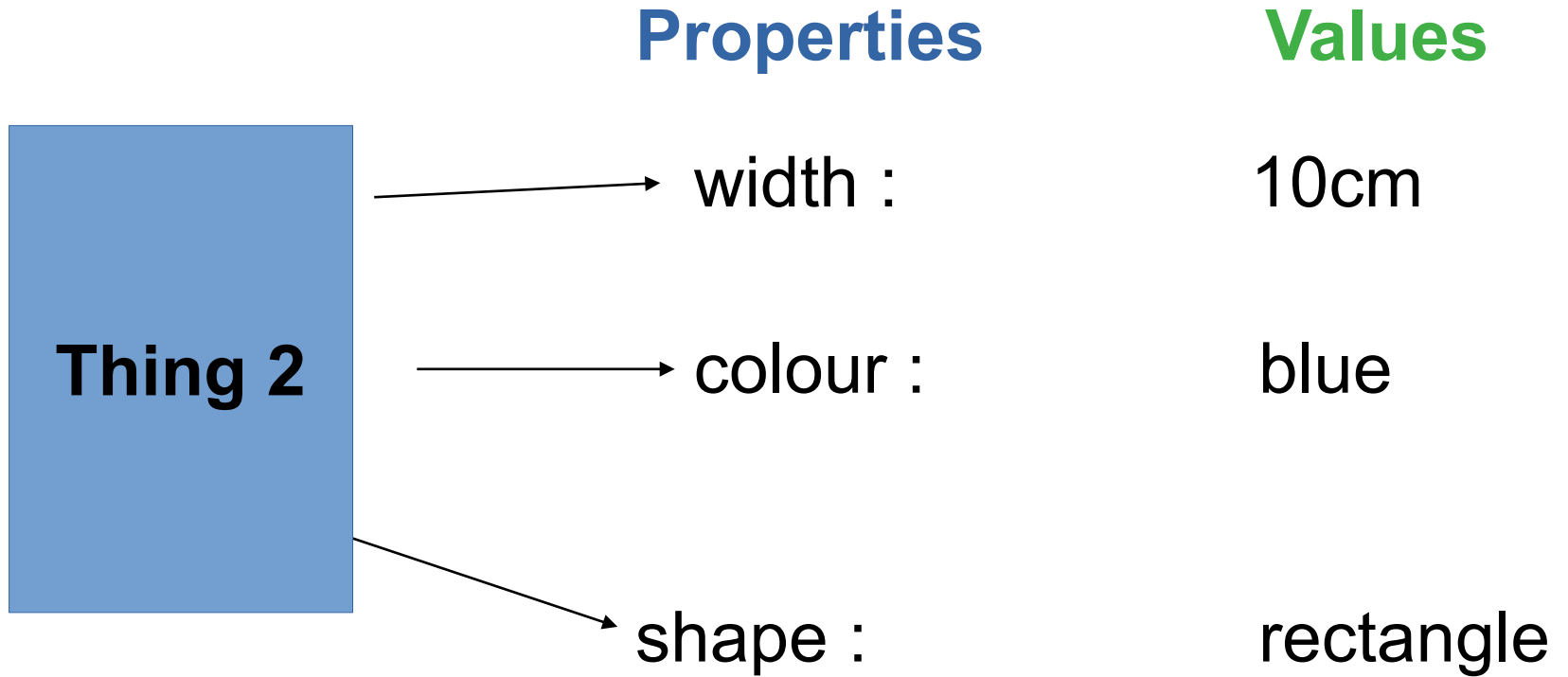


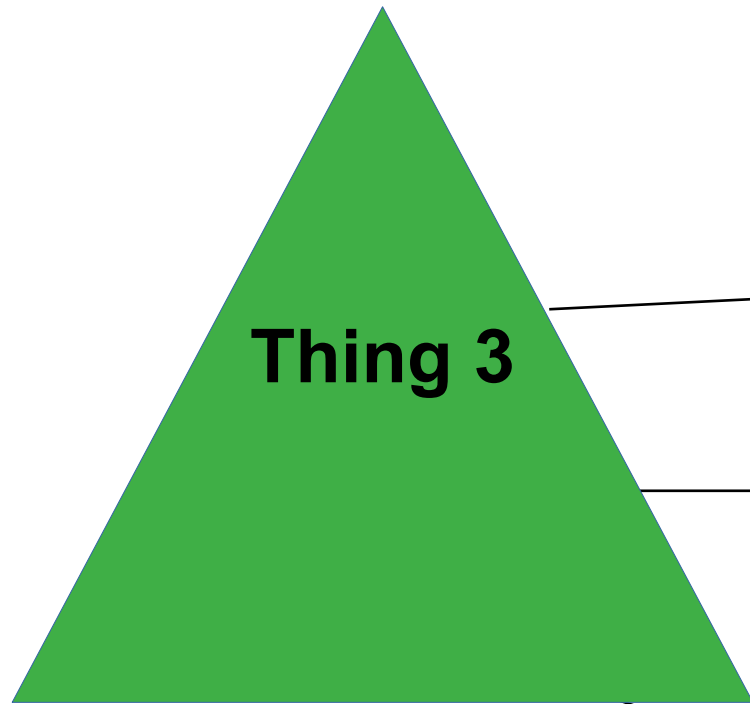
10cm

magenta

circle







Properties

Values

width :

15cm

colour :

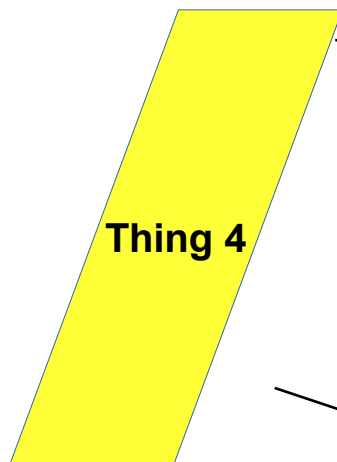
green

shape :

triangle

Properties

Values



width :

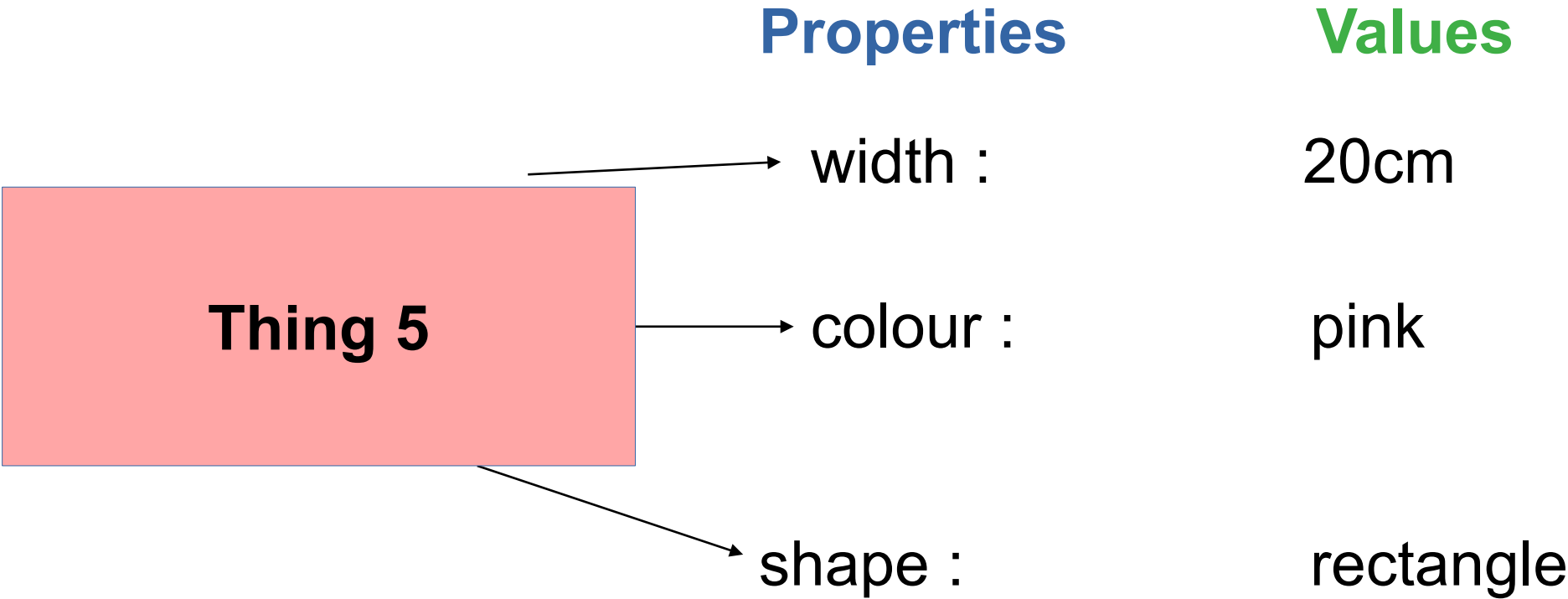
7cm

colour :

yellow

shape :

rhomboid

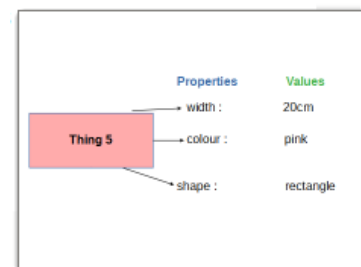
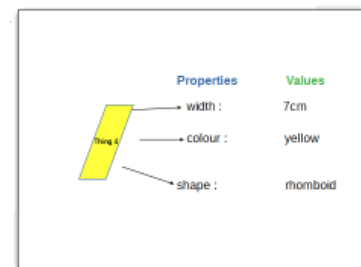
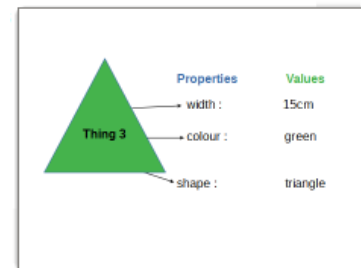
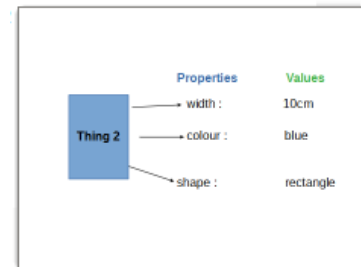
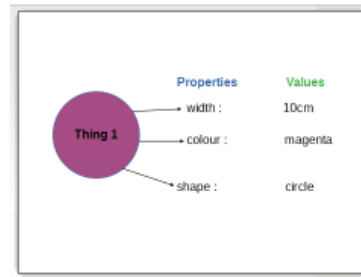


WE HAVE
4 properties (variables,
attributes, features)

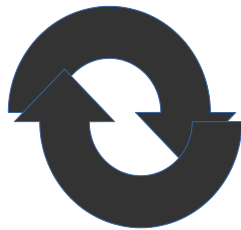
ID	width (cm)	colour	shape
1	10	magenta	circle
2	10	blue	rectangle
3	15	green	triangle
4	7	yellow	rhomboid
5	20	pink	rectangle

This table is what we
analyse

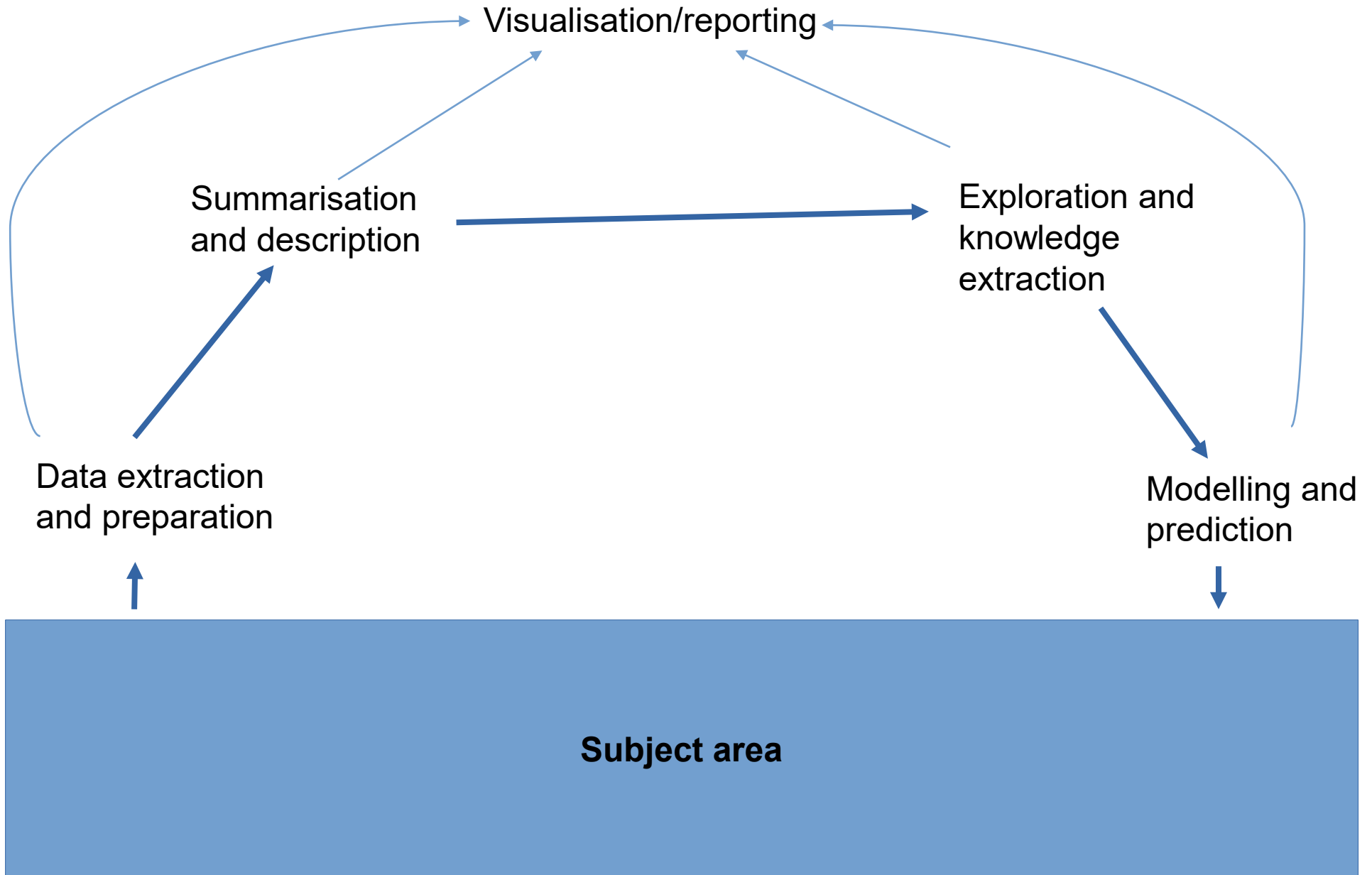
WE HAVE
5 things
(instances,
observations,
examples)



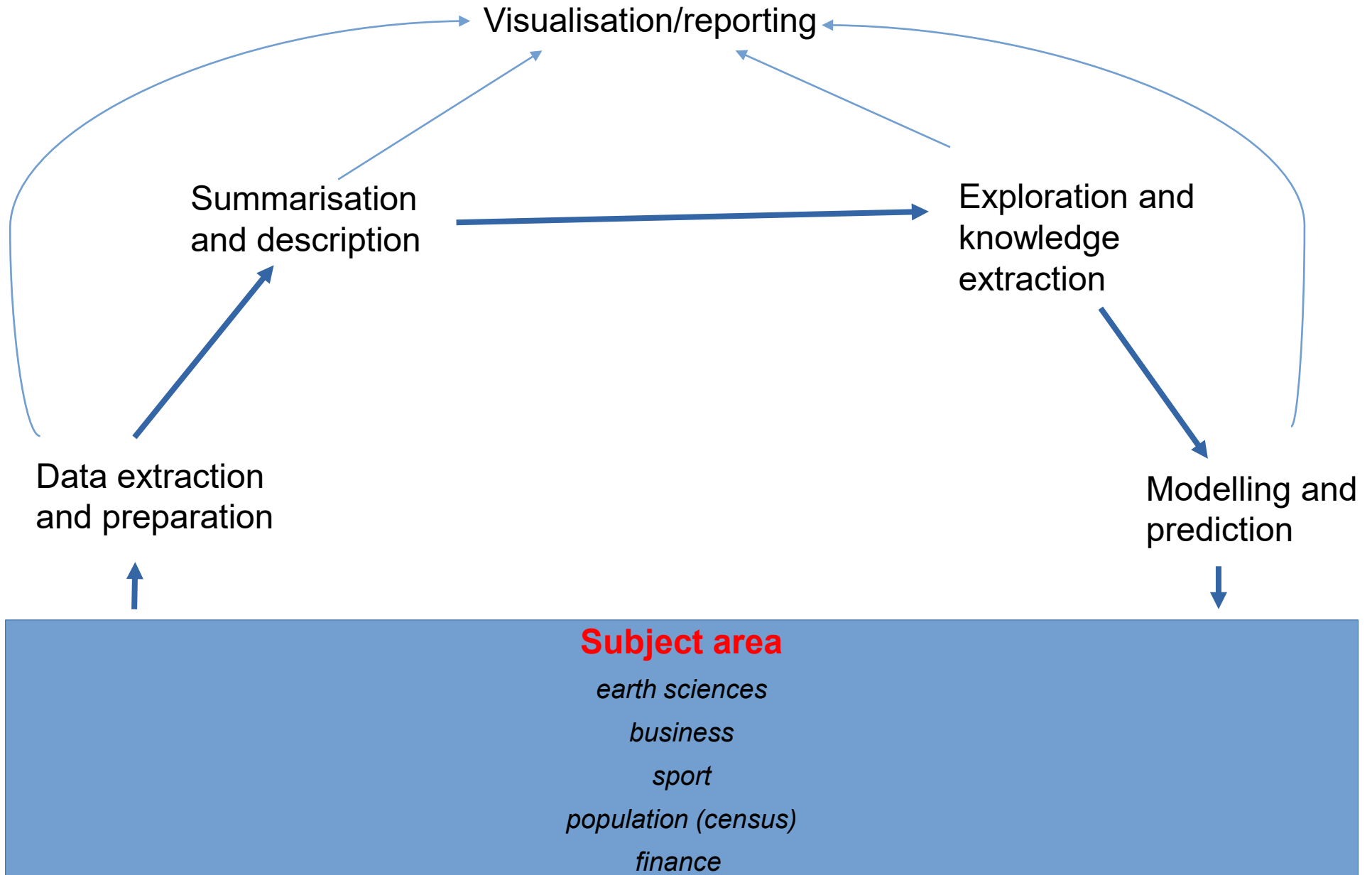
The analysis cycle



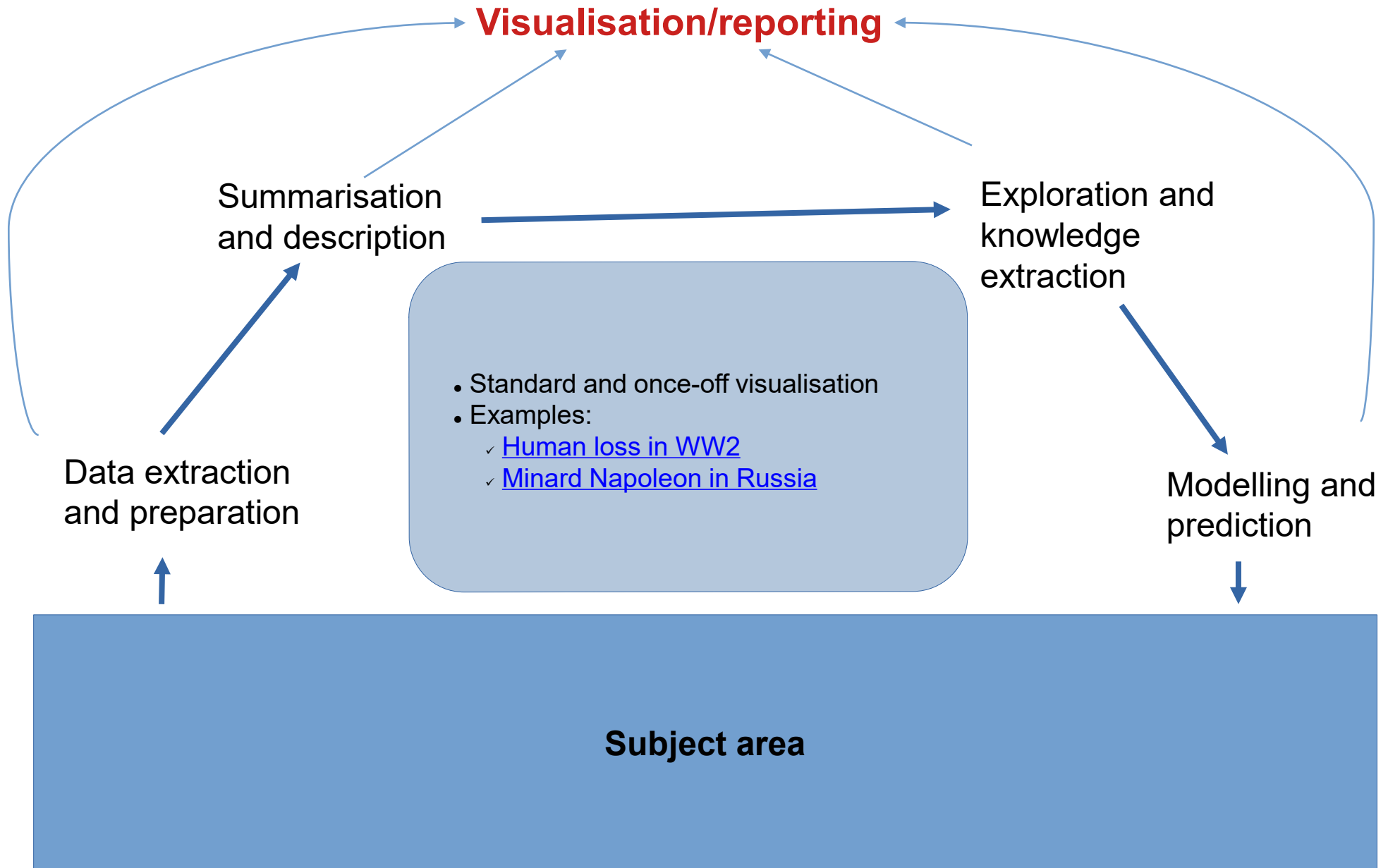
The data analysis cycle



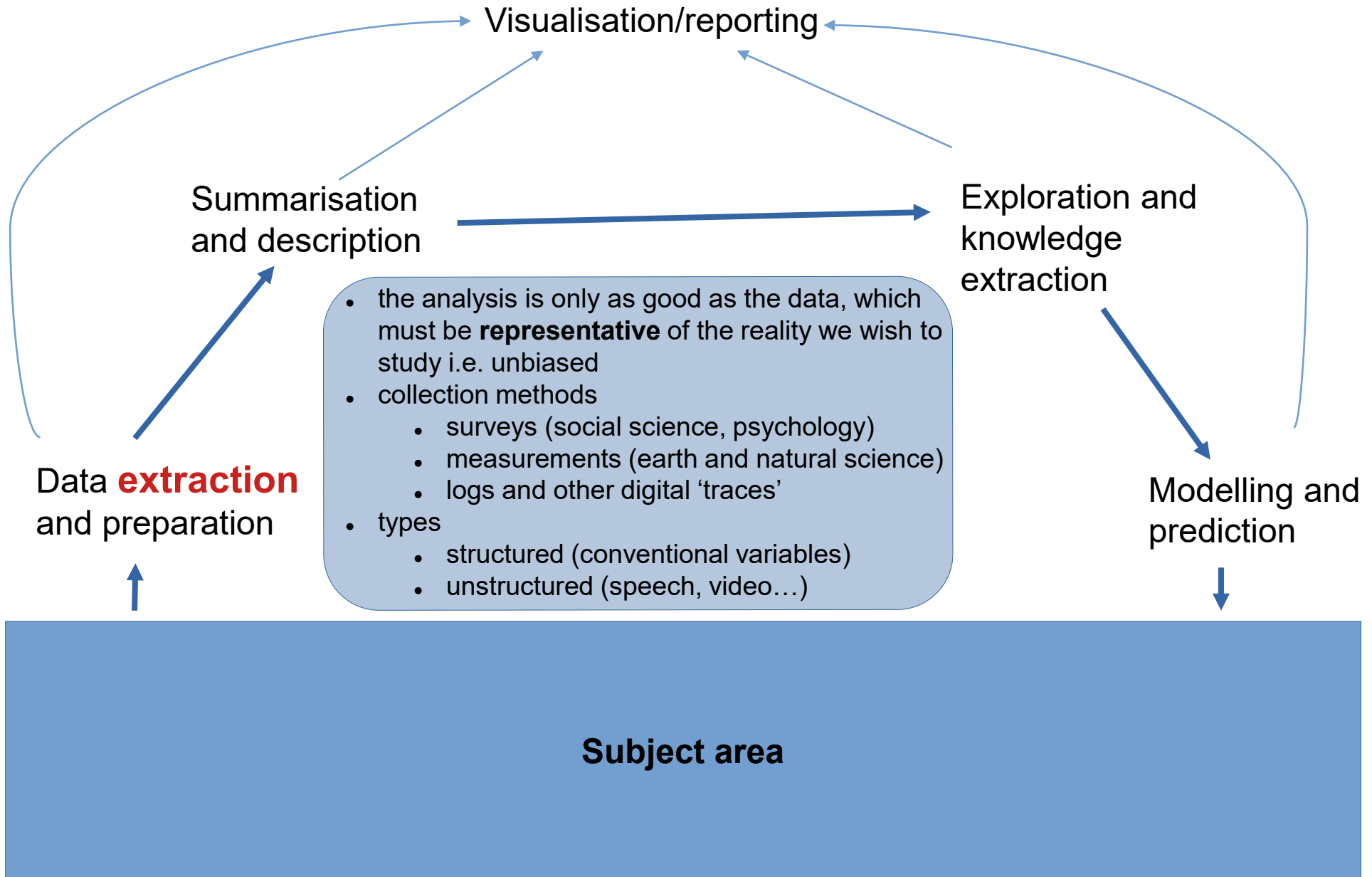
The data analysis cycle



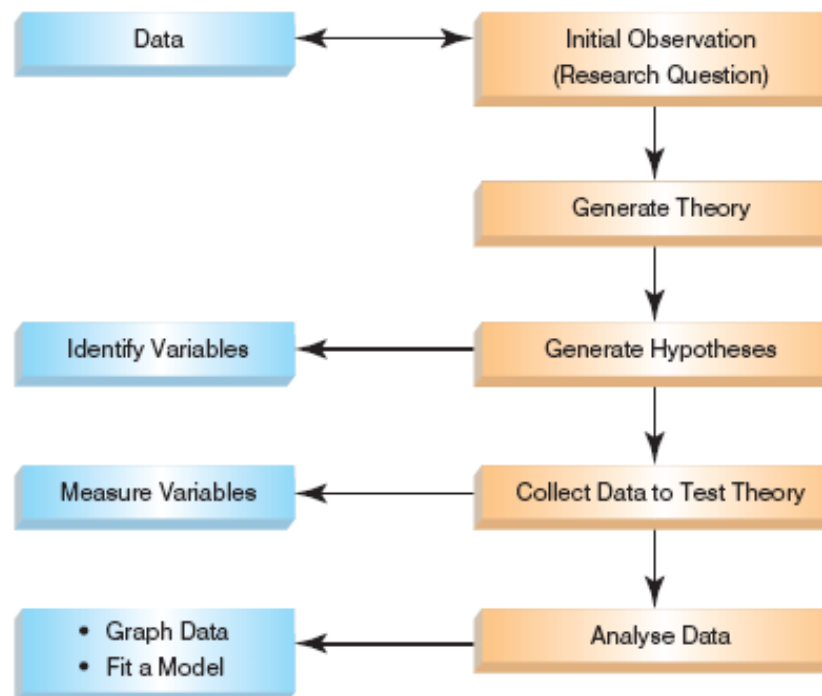
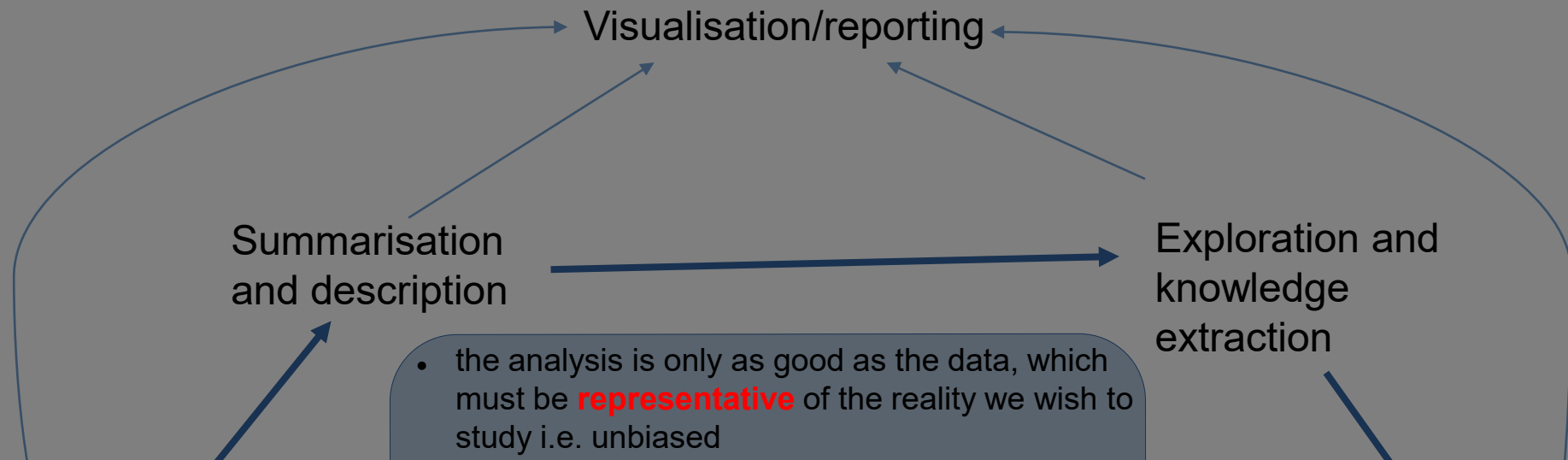
The data analysis cycle



The data analysis cycle



The data analysis cycle

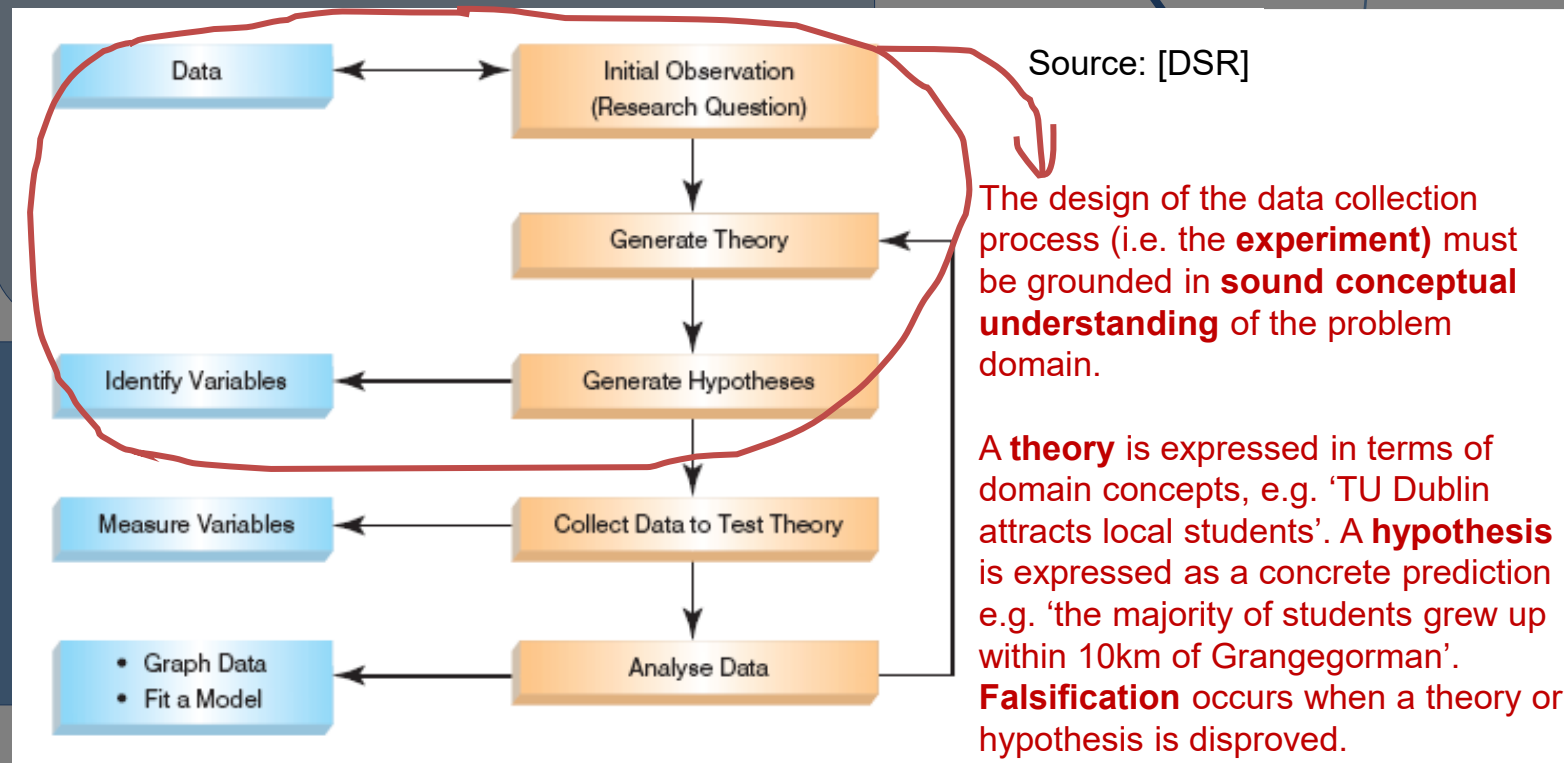
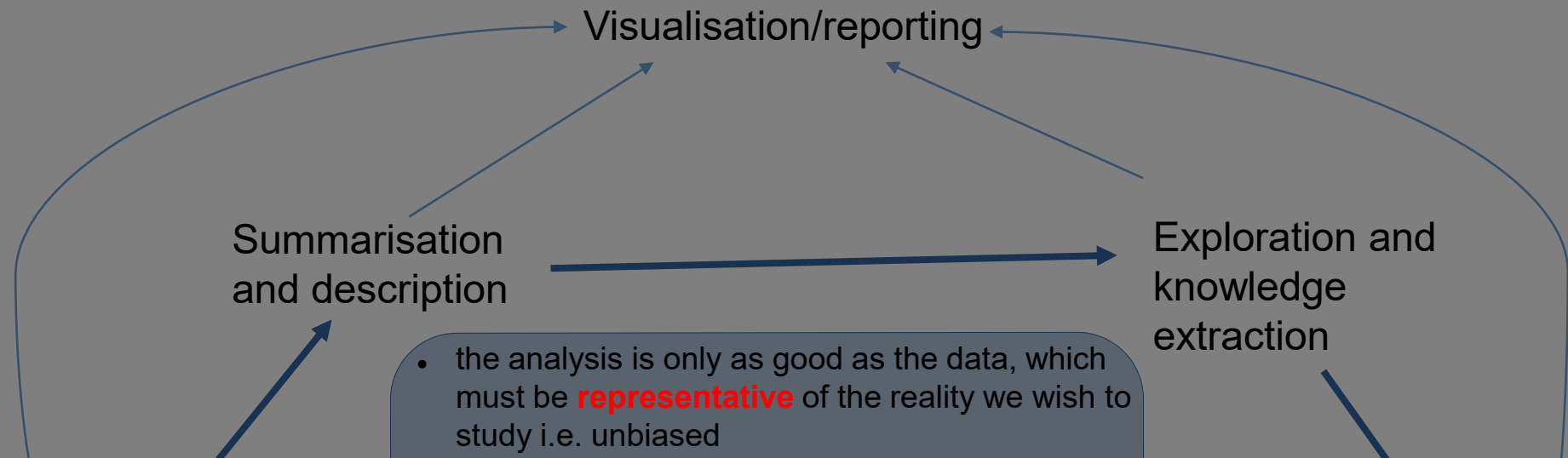


Source: [DSR]

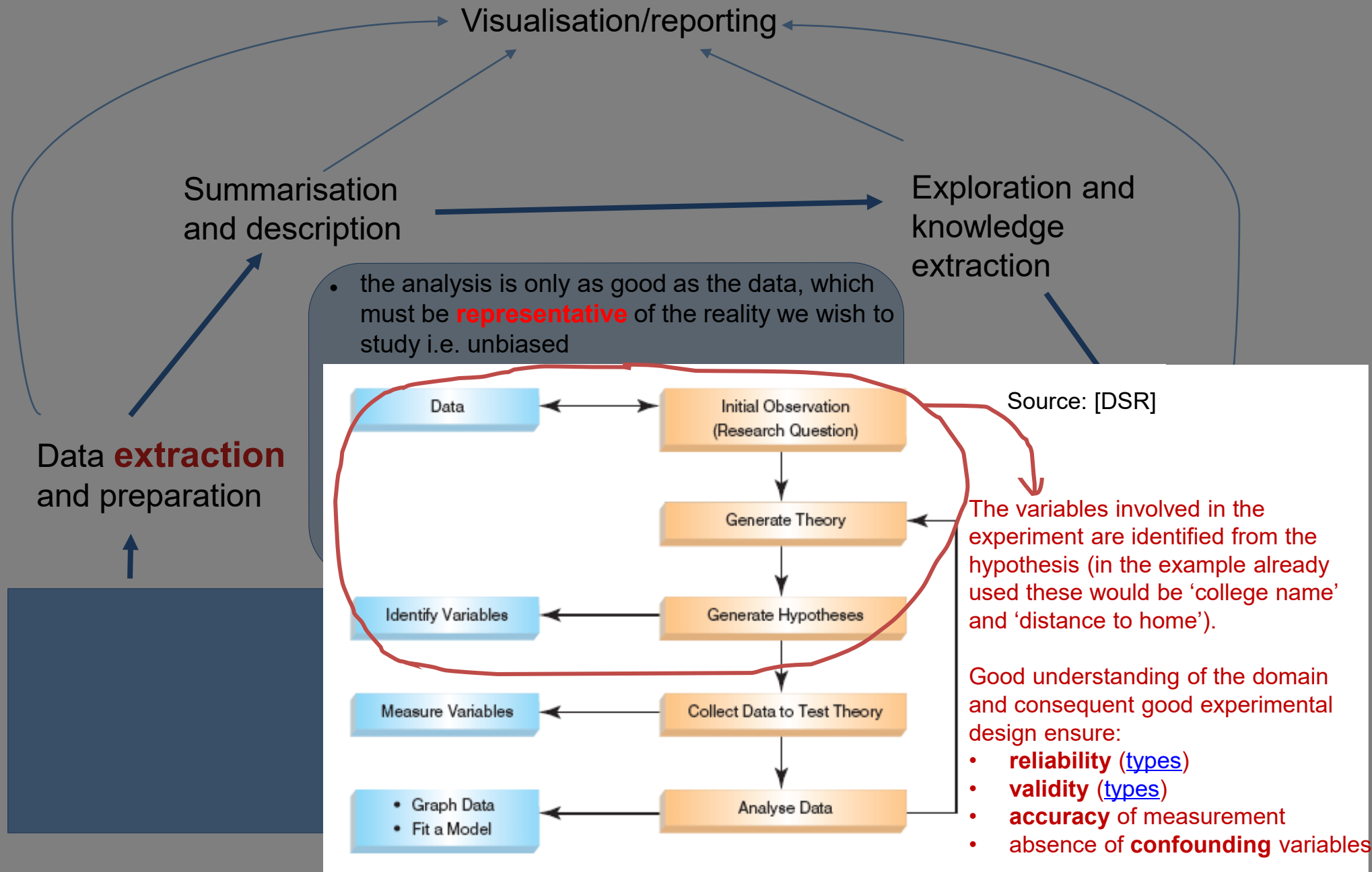
Often, data analysis is performed simply because **data are available**. This is the case with enterprise data warehouses, which are explored for knowledge hidden in data already available as a by-product of a company's commercial activity.

The standard scientific approach is to start with a theory that needs to be tested and collect data with that aim. The picture shows this process, which needs to be paralleled in **sampling assessment** for analysis, even if the data is pre-collected.

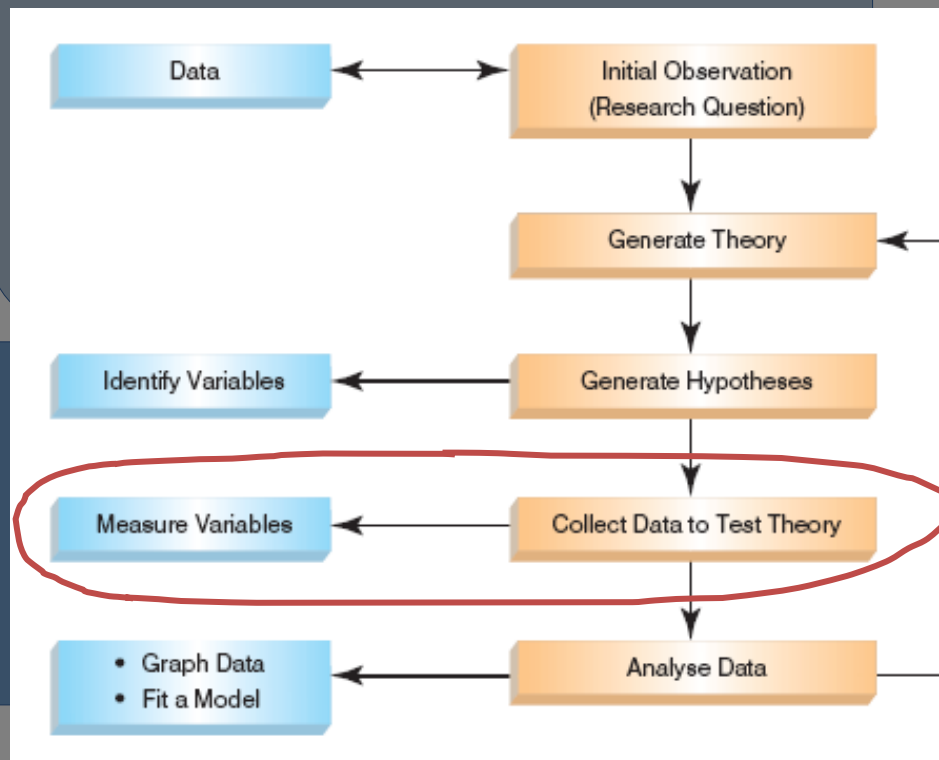
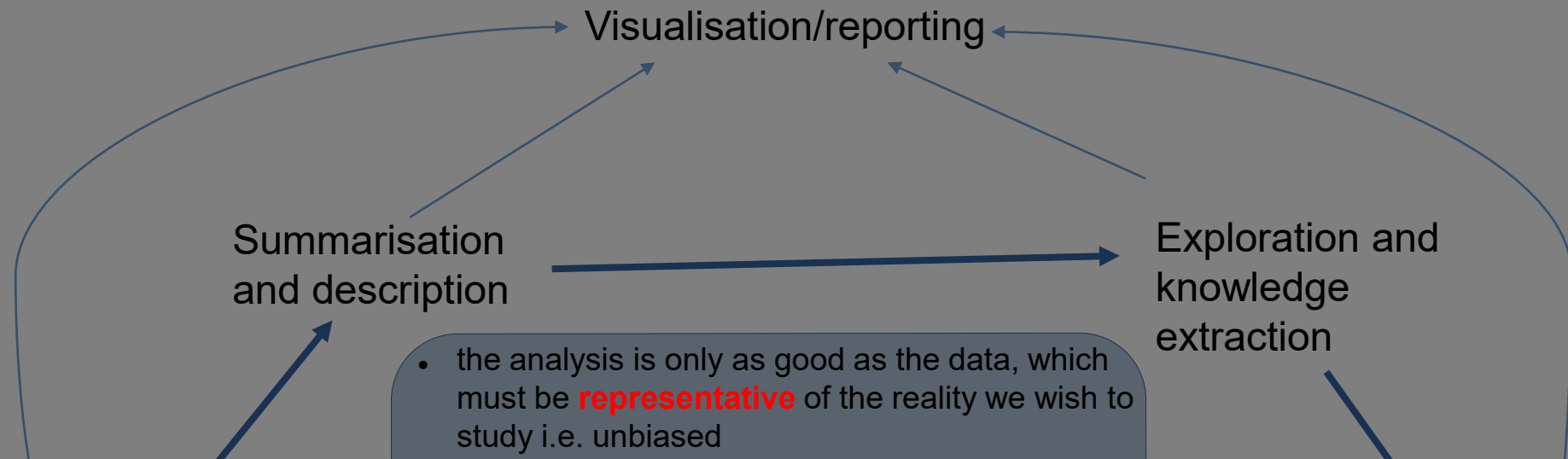
The data analysis cycle



The data analysis cycle



The data analysis cycle

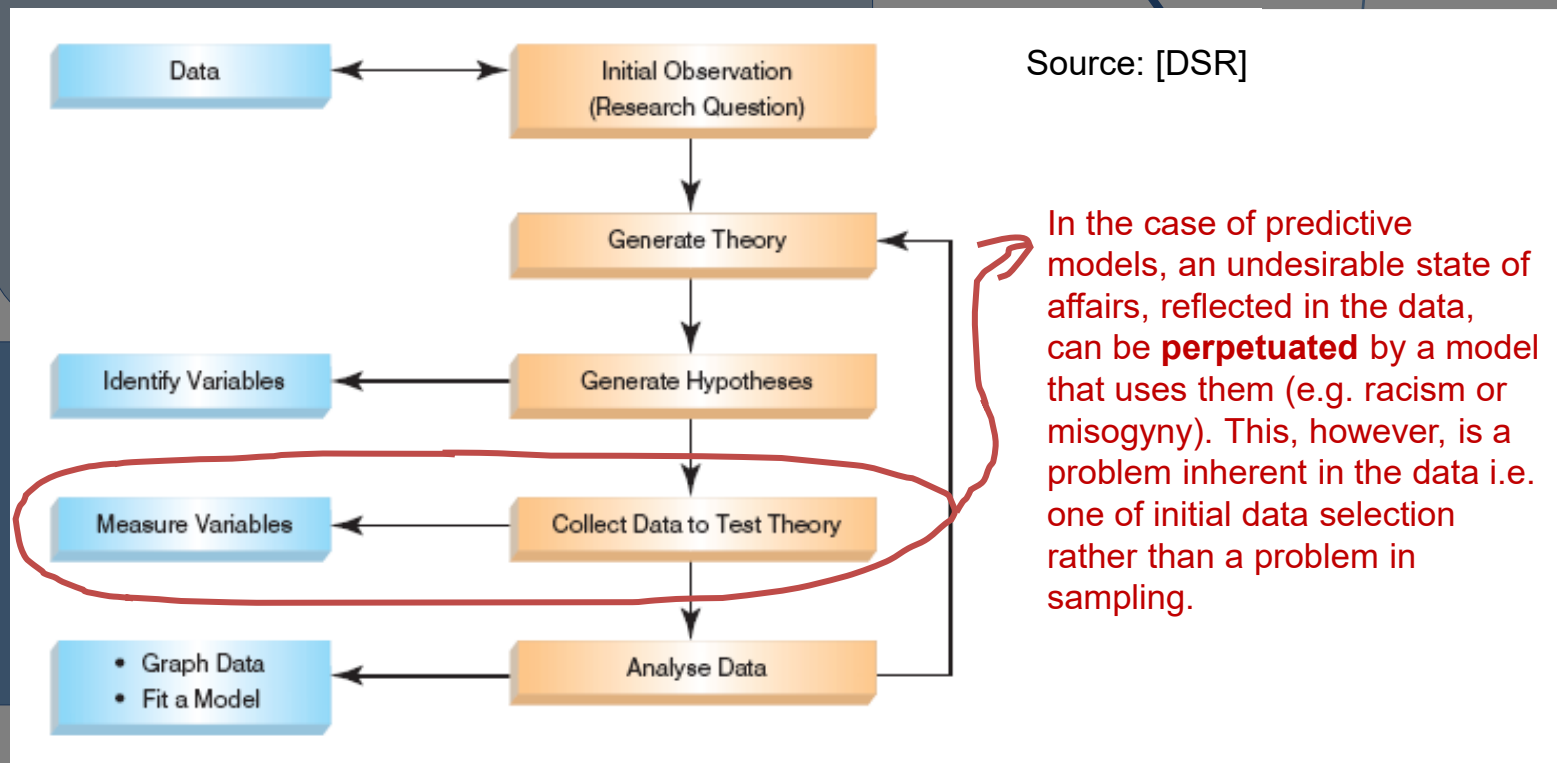
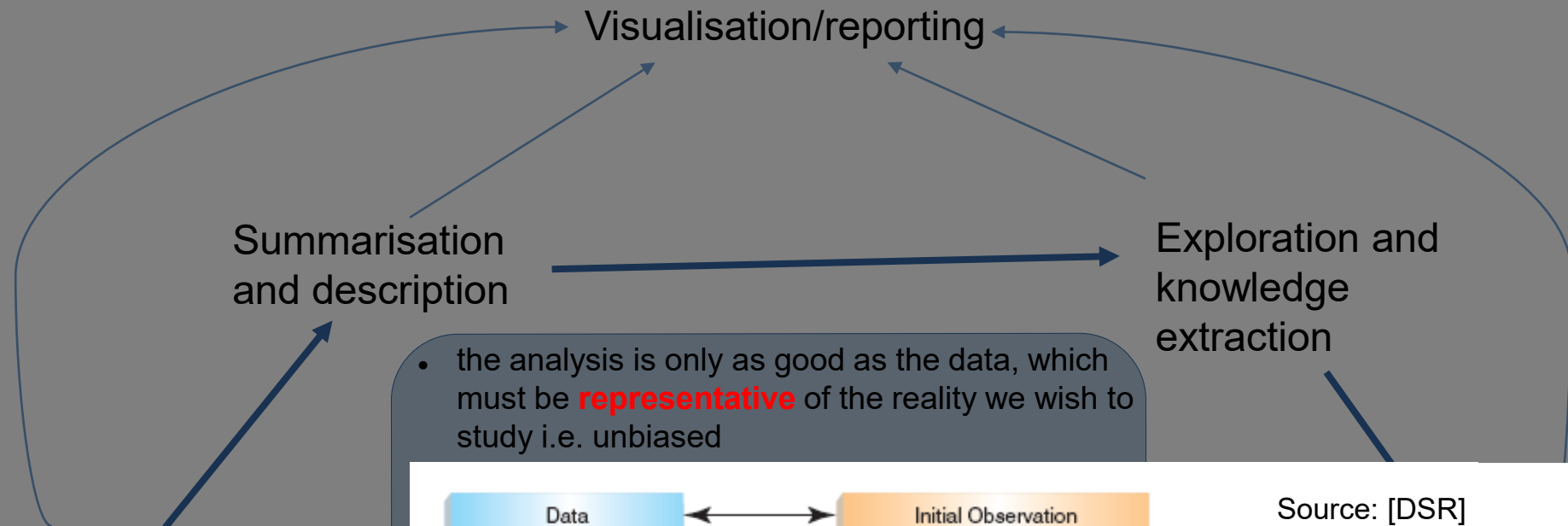


Source: [DSR]

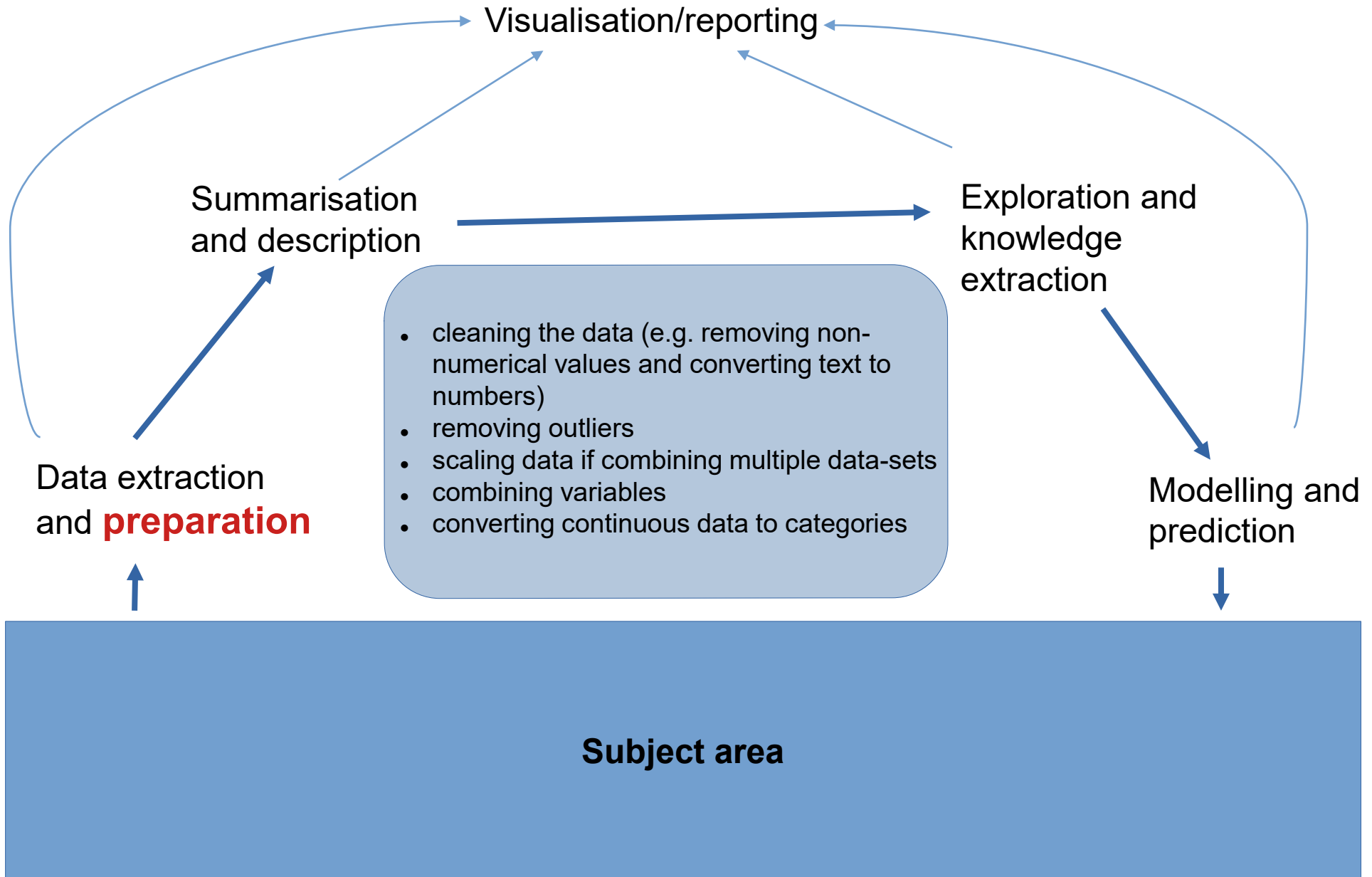
The instances for which data is collected represent a **sample** of the **population** of all instances. Two important (but not only) factors in how unbiased a sample will be are its **size** and **sampling method**. Some type of random sampling is desirable but other methods can be acceptable ([types](#)).

Small samples may lead to wrong conclusions. Large samples can be wasteful.

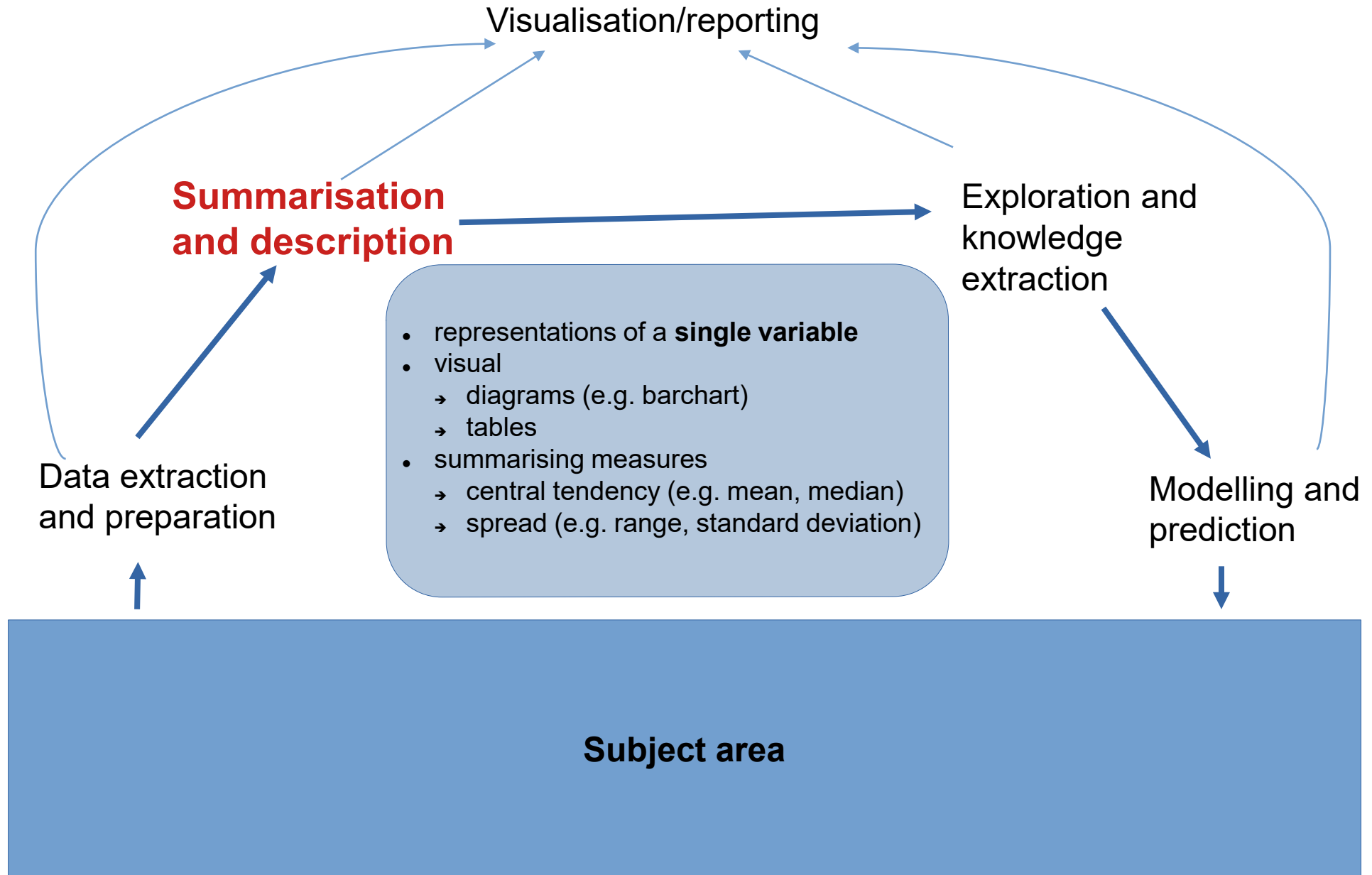
The data analysis cycle



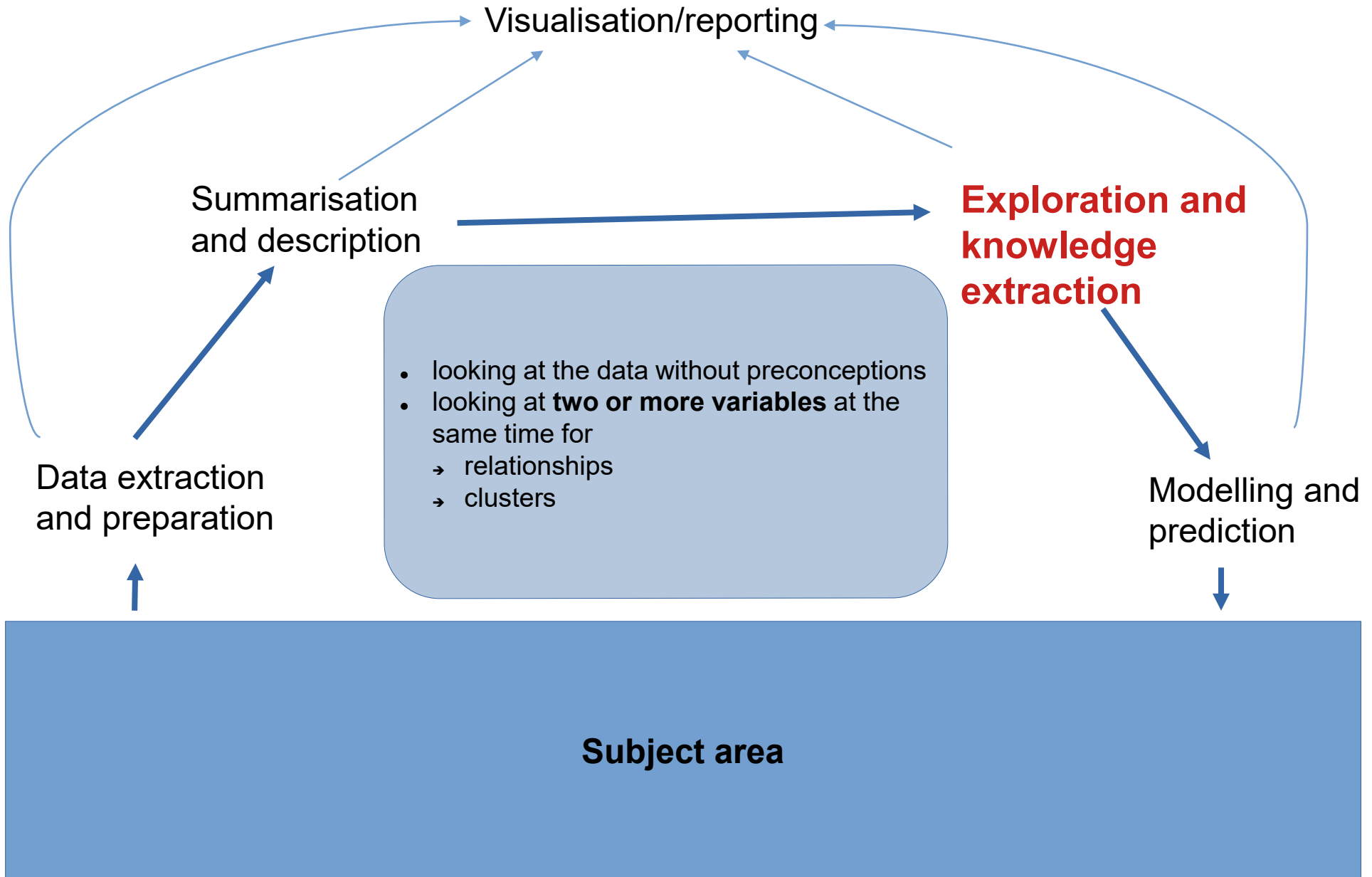
The data analysis cycle



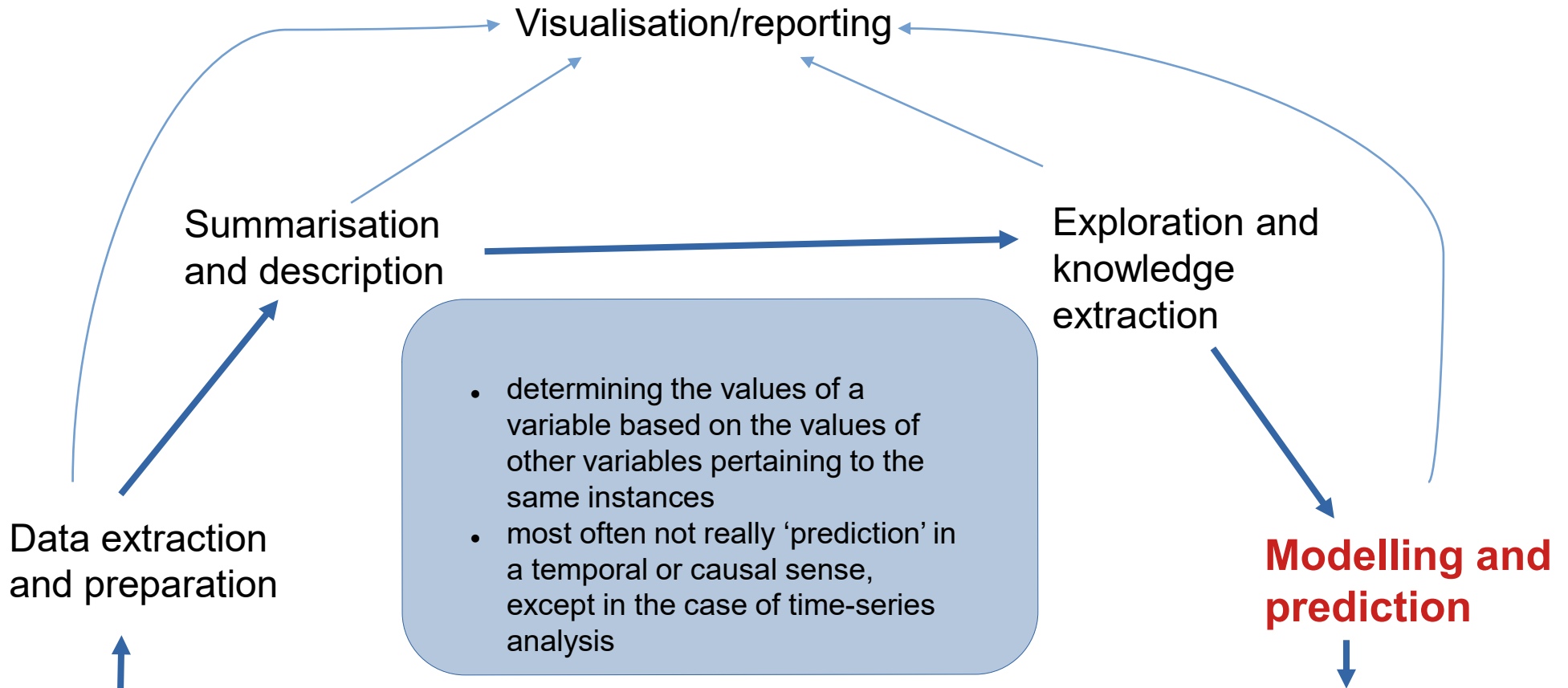
The data analysis cycle



The data analysis cycle



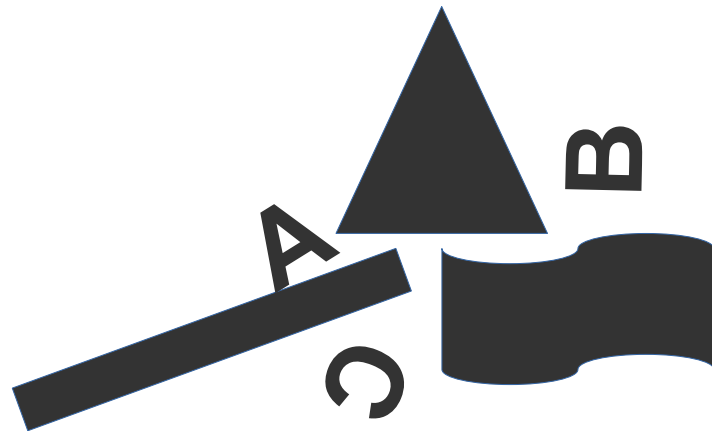
Subject area



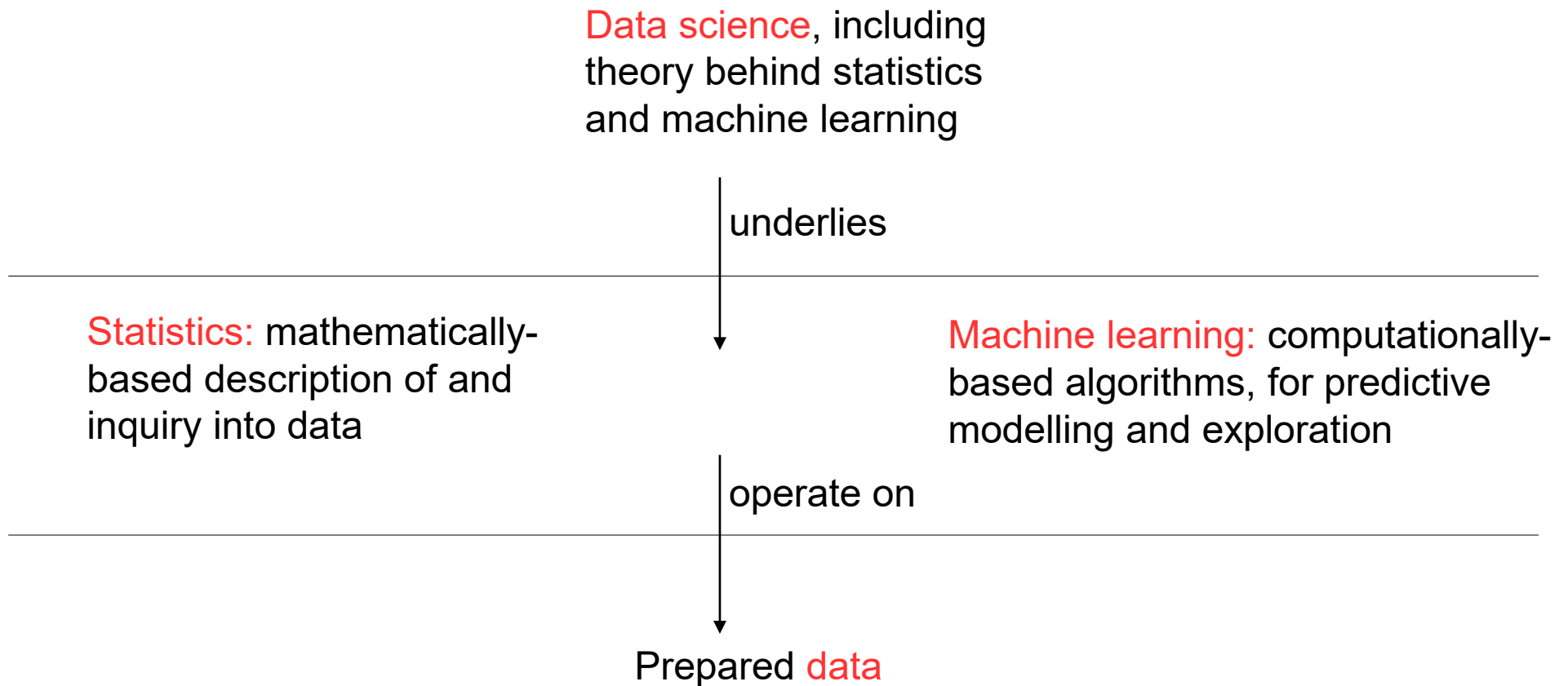
[illegible]

- **statistical inference:**
 - using a sample of instances to understand an entire population and its parameters
 - understanding how much we can trust such conclusions about the population
- relevant in in some form at **all stages** of the data cycle

The landscape and terminology of data analysis



The data analysis landscape



Landscape and terminology

- Sets out the principles and theory for understanding and using data
- Studies how these principles and techniques should be applied in each individual case

Data science, including theory behind statistics and machine learning

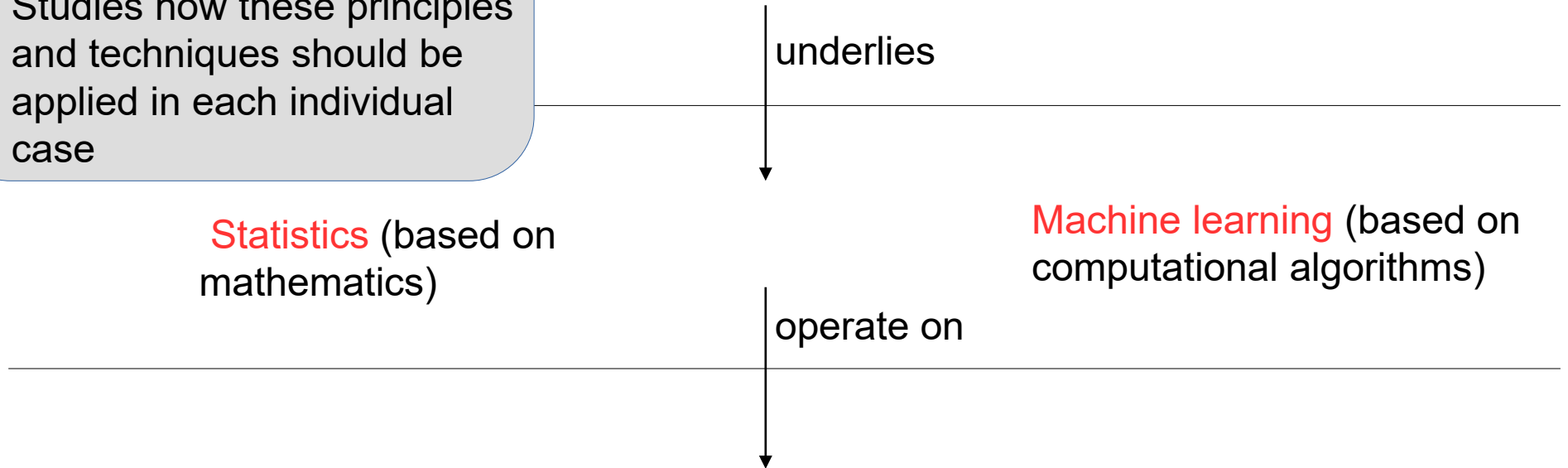
underlies

Statistics (based on mathematics)

Machine learning (based on computational algorithms)

operate on

Prepared **data**



Landscape and terminology

- The science and practice of analysing numerical data, particularly with the purpose of understanding the properties of a large population by analysing a representative sample.

Statistics (based on mathematics)

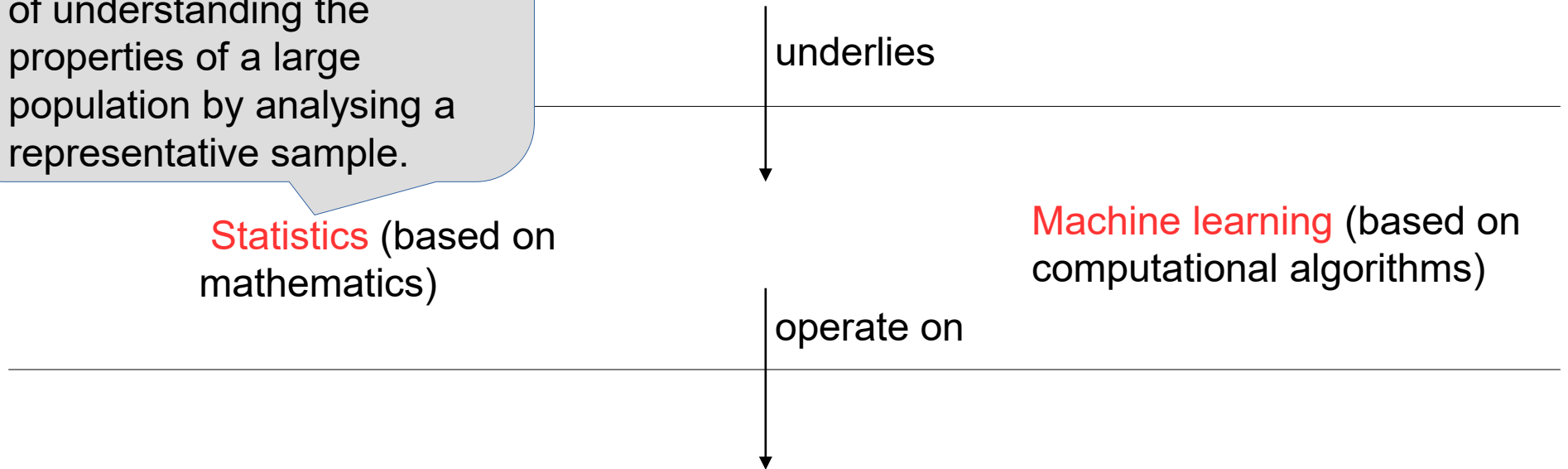
Data science, including theory behind statistics and machine learning

underlies

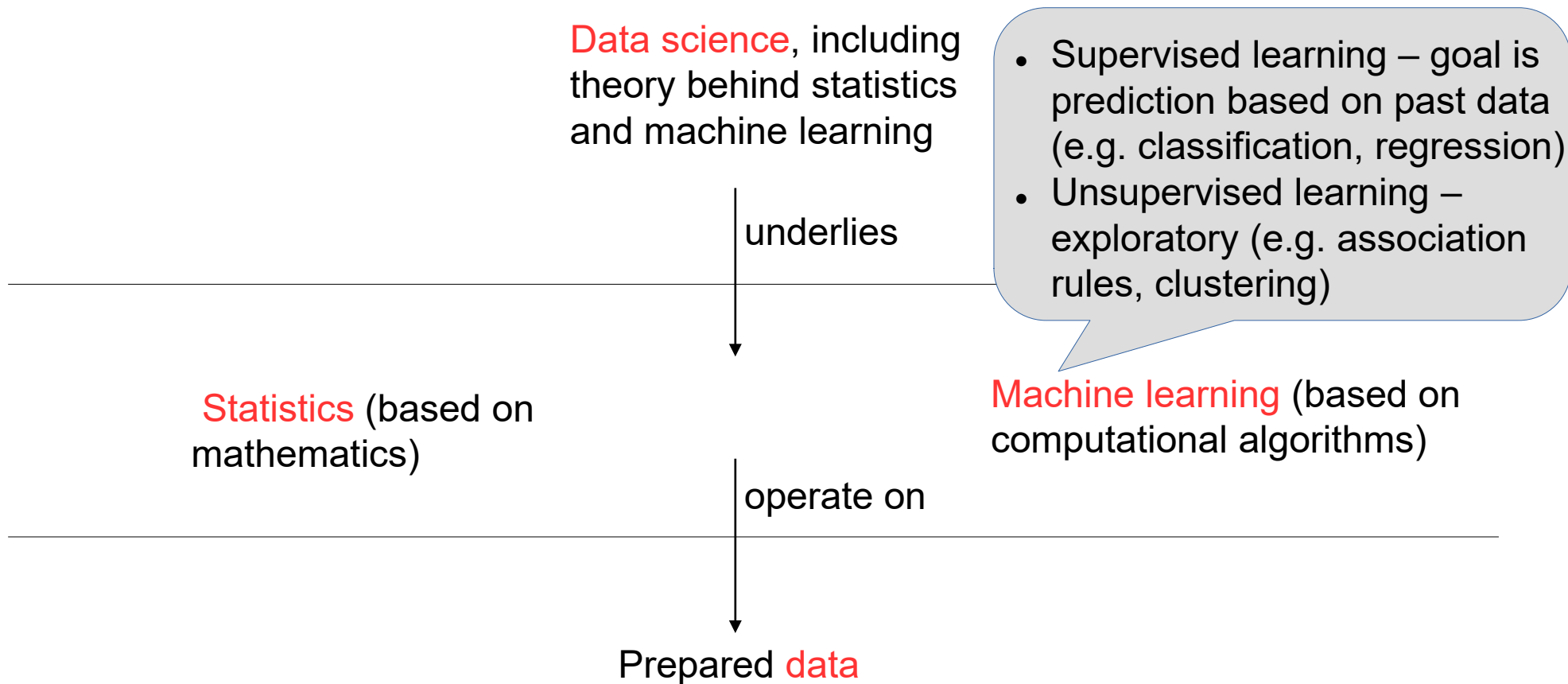
Machine learning (based on computational algorithms)

operate on

Prepared **data**

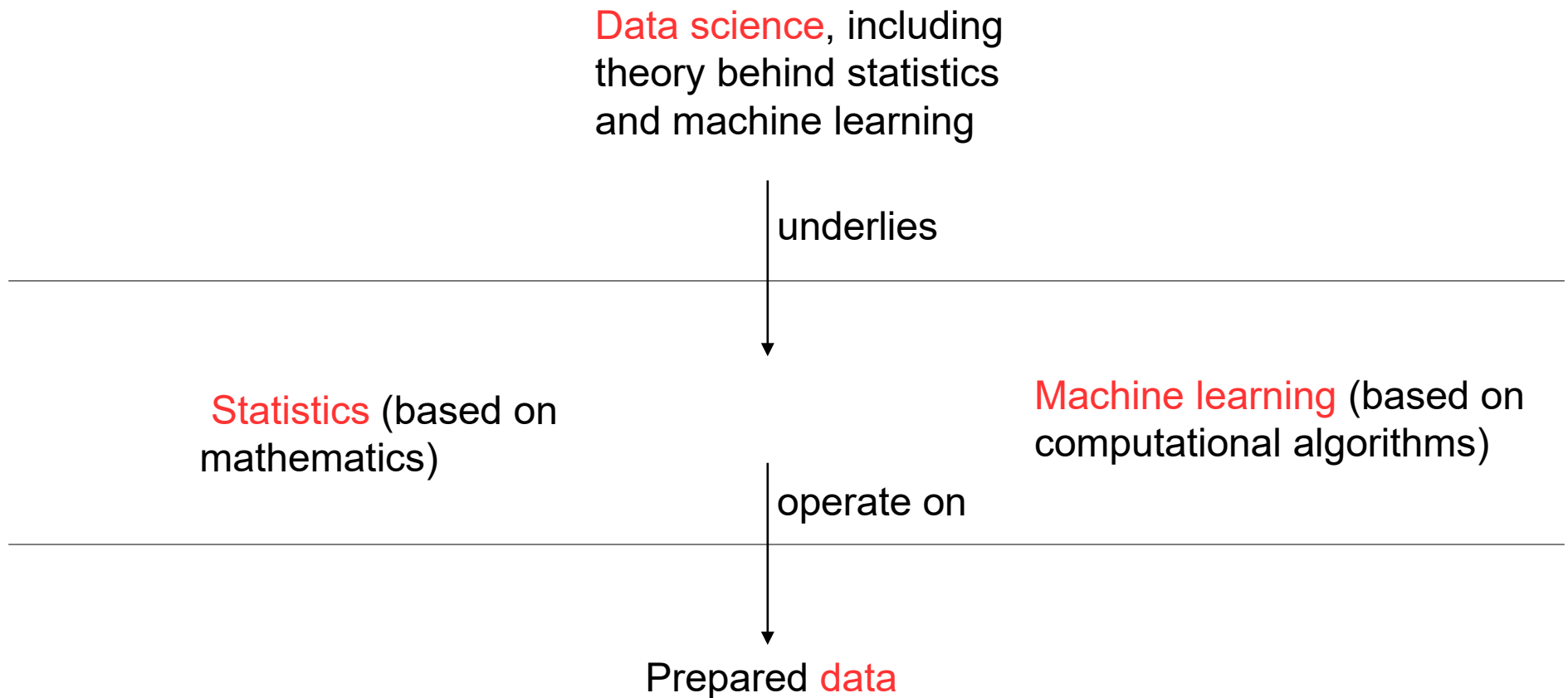


Landscape and terminology



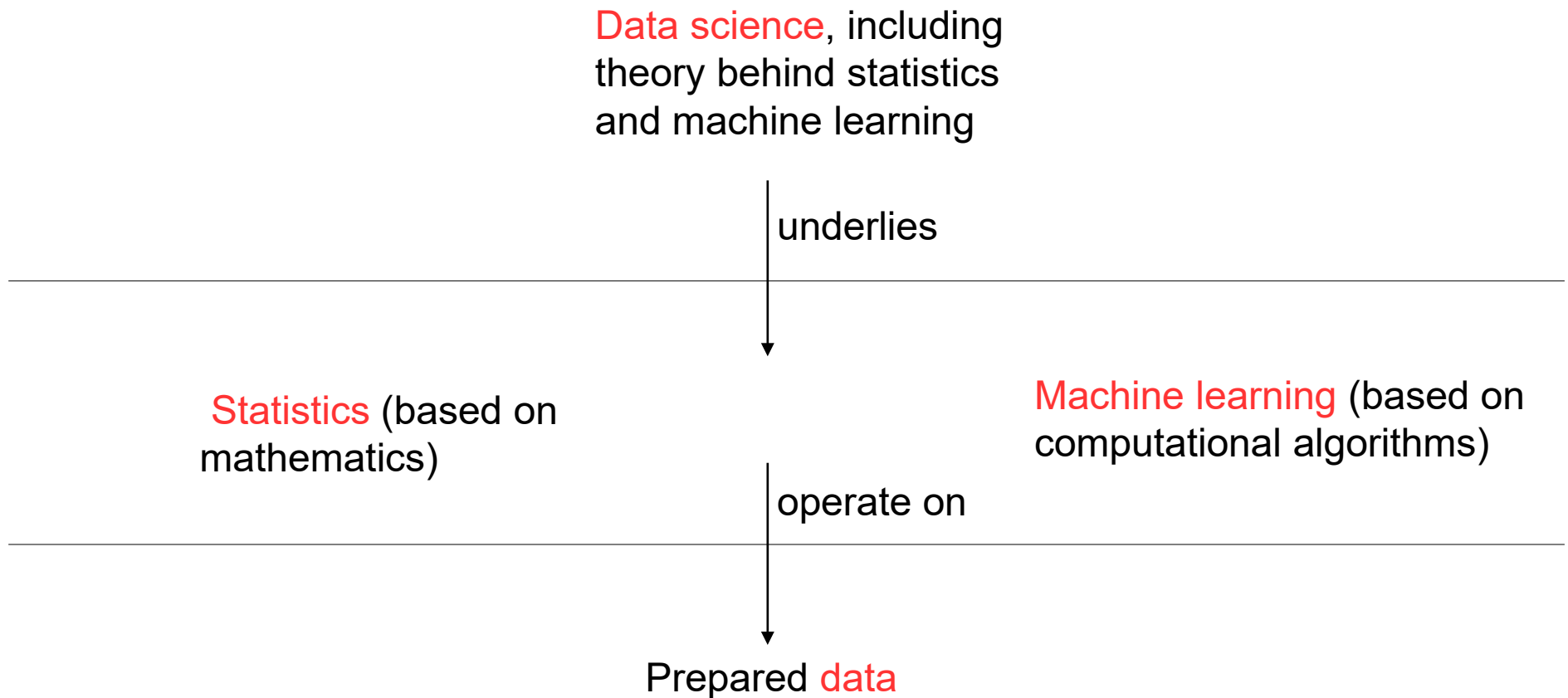
Landscape and terminology

- **Analytics** – techniques and activities based around using data in a problem domain e.g. business analytics, financial analytics (the middle of the ‘sandwich’ below).



Landscape and terminology

- **Analytics** – techniques and activities based around using data in a problem domain e.g. business analytics, financial analytics (the middle of the ‘sandwich’ below).
- **Data analysis** – the activity of analysing data in any way.



Landscape and terminology

- **Analytics** – techniques and activities based around using data in a problem domain e.g. business analytics, financial analytics (the middle of the ‘sandwich’ below).
- **Data analysis** – the activity of analysing data in any way.
- **Data mining** – exploration of (predominantly) enterprise data warehouses (1990s)

Data science, including
theory behind statistics
and machine learning

underlies

Statistics (based on
mathematics)

Machine learning (based on
computational algorithms)

operate on

Prepared **data**

Landscape and terminology

- **Analytics** – techniques and activities based around using data in a problem domain e.g. business analytics, financial analytics (the middle of the ‘sandwich’ below).
- **Data analysis** – the activity of analysing data in any way.
- **Data mining** – exploration of (predominantly) enterprise data warehouses (1990s)
- **Big data** (on next page)

Data science, including
theory behind statistics
and machine learning

underlies

Statistics (based on
mathematics)

Machine learning (based on
computational algorithms)

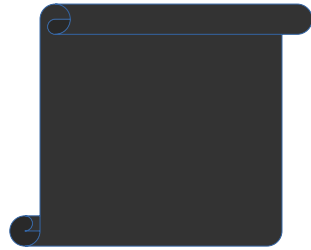
operate on

Prepared **data**

Big Data

- In the last 20 years the data cycle is 'intensifying'
 - Growing processing power
 - Almost limitless storage capacity
 - Connectivity with large bandwidths
 - Techniques have developed on this new wave of possibilities
-
- Big data is at a scale that cannot be processed with conventional technologies.
 - 4 Vs
 - [IBM 4Vs of Big Data](#)
 - New technologies:
 - Hadoop (Apache)
 - MapReduce (Google)
 - MongoDB etc.
 - The data science principles are the same as 'normal sized' data

In this module



you will learn

to...

Ask **questions** that can be answered by data analysis

Design and assess the **sampling and data collection** processes

Inspect and **prepare data** for analysis

Statistically **describe and summarise data** to quantify properties of things

Understand and conduct **exploratory analysis** to investigate property interrelations

Understand and conduct **predictive analysis** by making use of the interrelations

Interpret the outcomes of data analysis and draw appropriate conclusions

Report on the findings of data analysis

References

Some pictures in this presentation were taken from the following books. They are cited using the keys shown in square brackets.

[DSR] *Discovering Statistics Using R*, by Andy Field, Jeremy Miles and Zöe Field, Sage, 2012.