

Data Analysis: Introduction

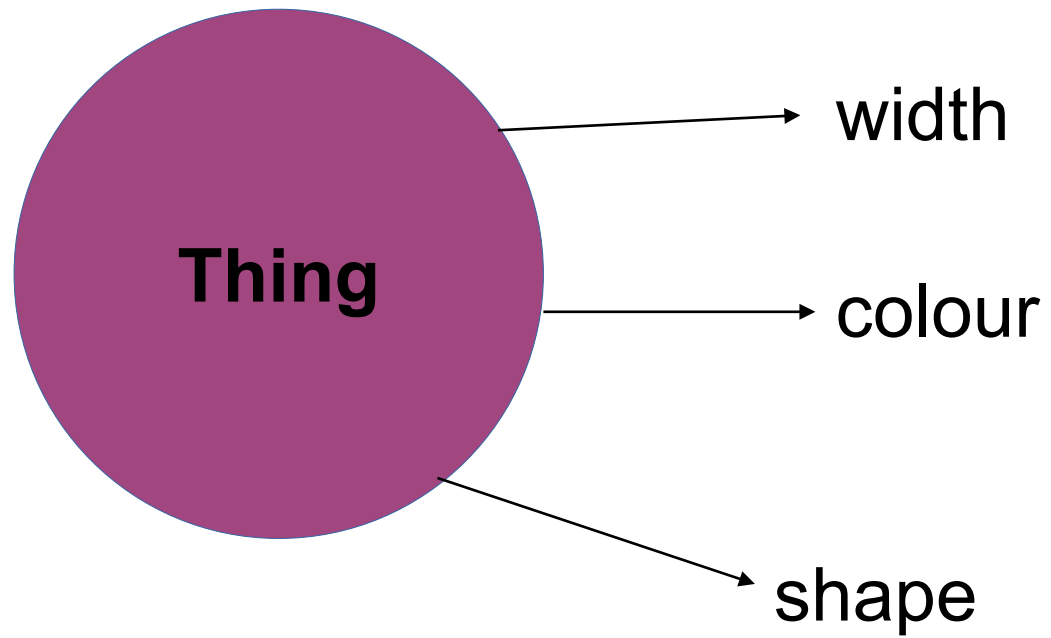
TU Dublin Tallaght Campus,
Department of Computing

What are we
dealing with

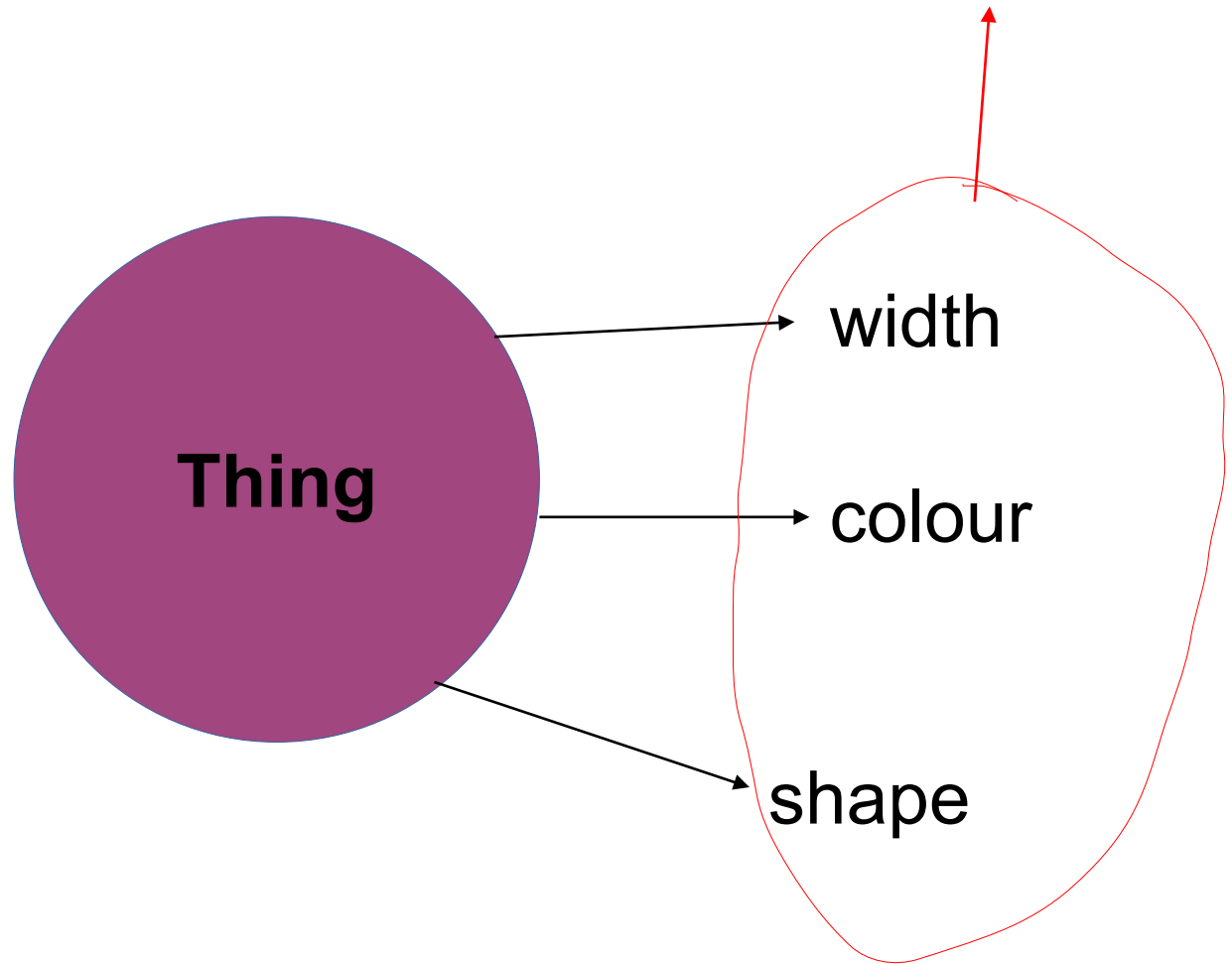




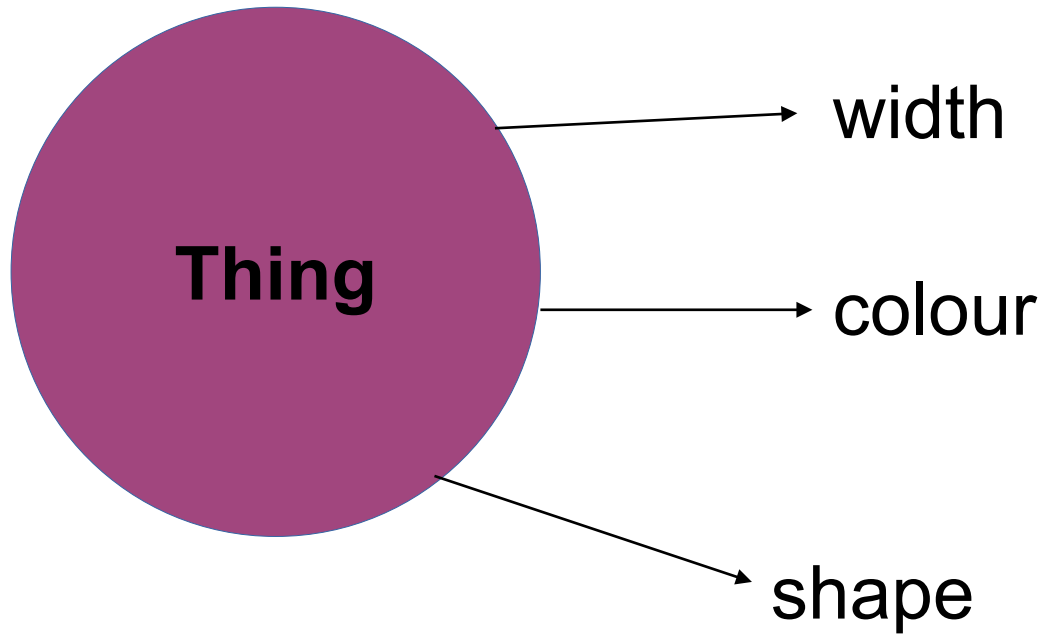
Thing



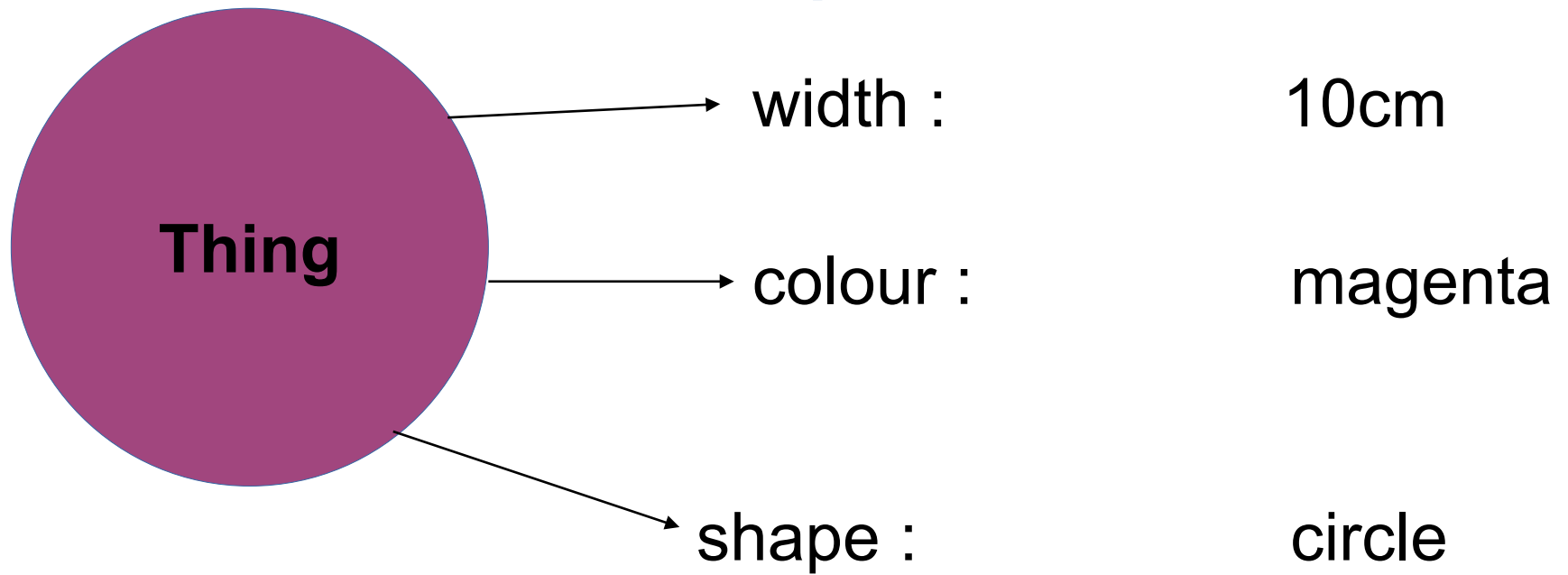
What are these?



Properties

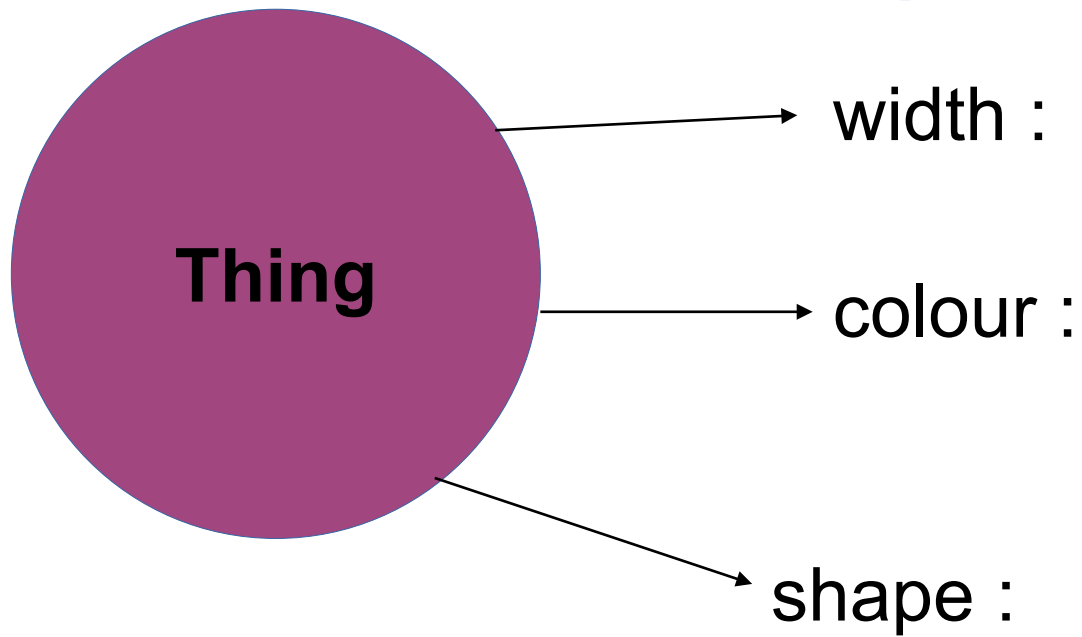


Properties



What are these?

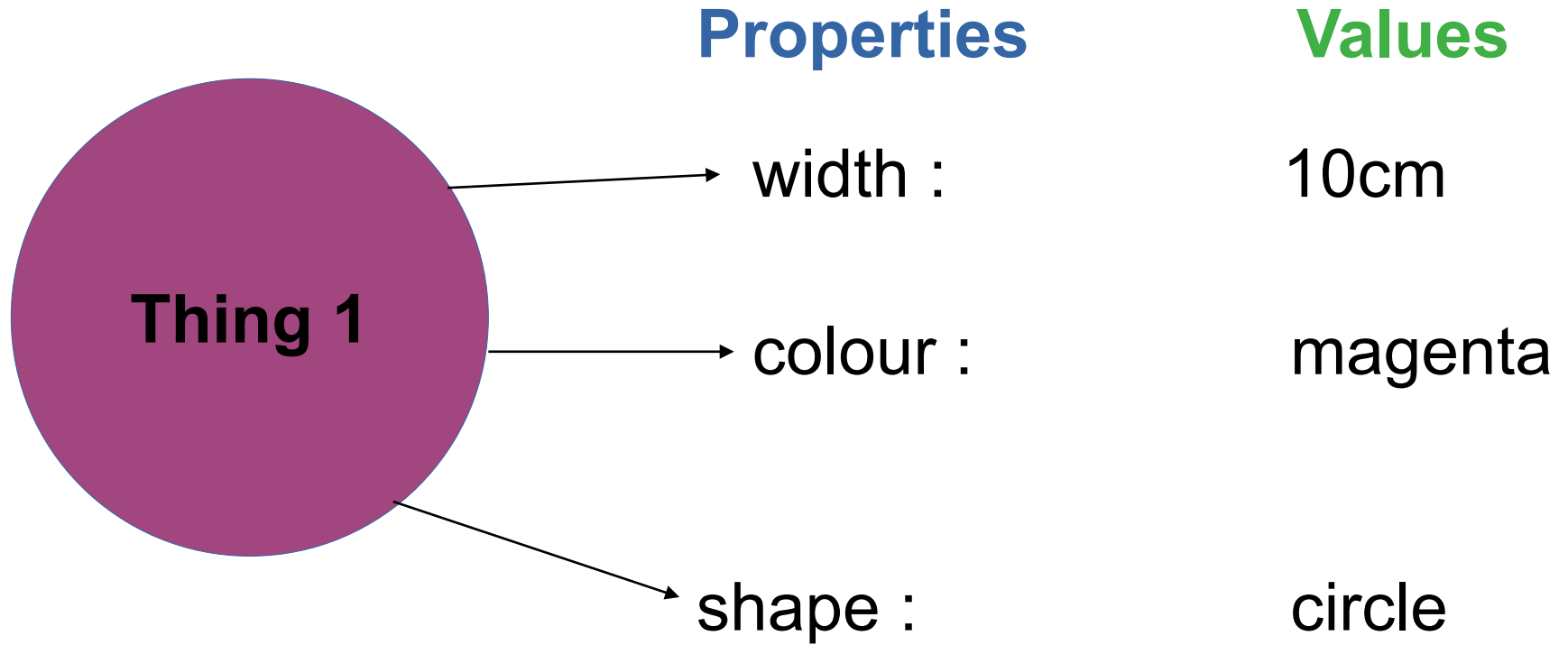
Properties

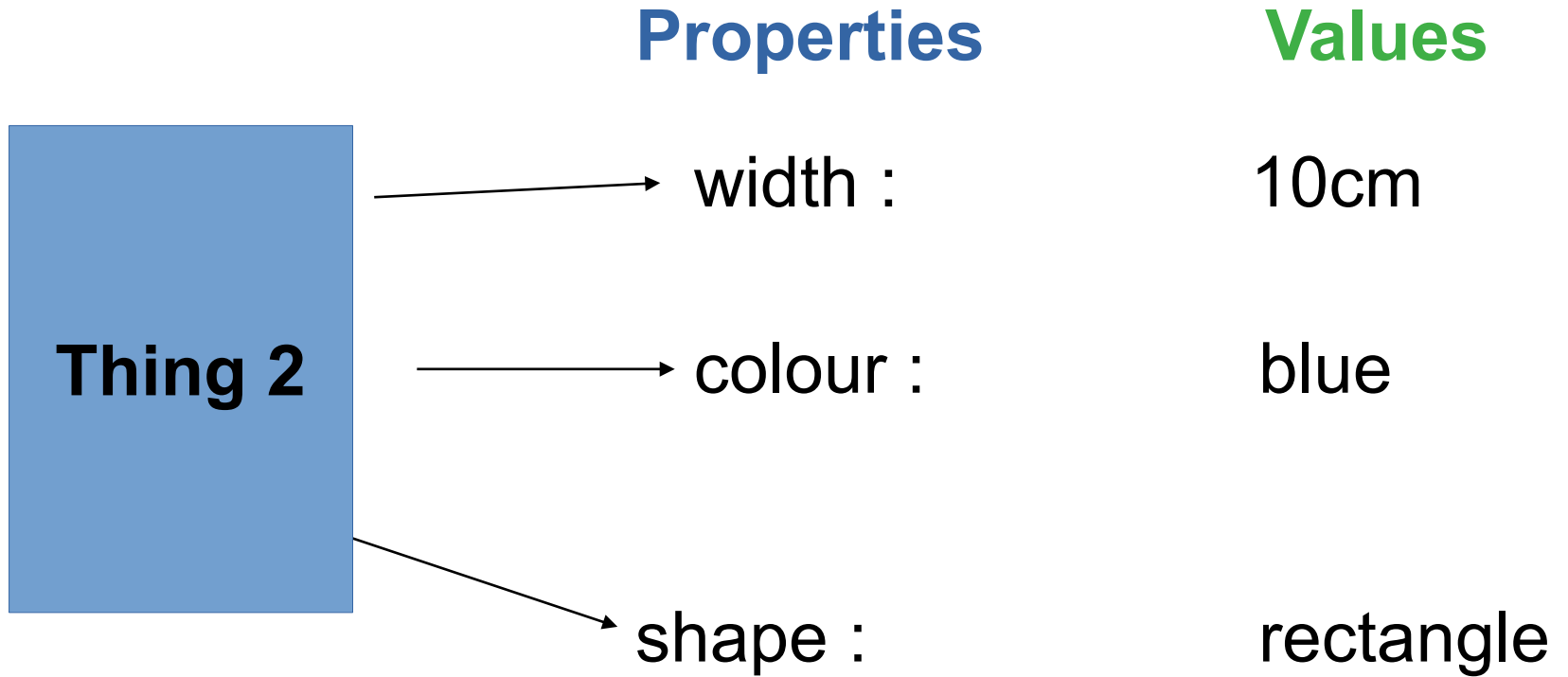


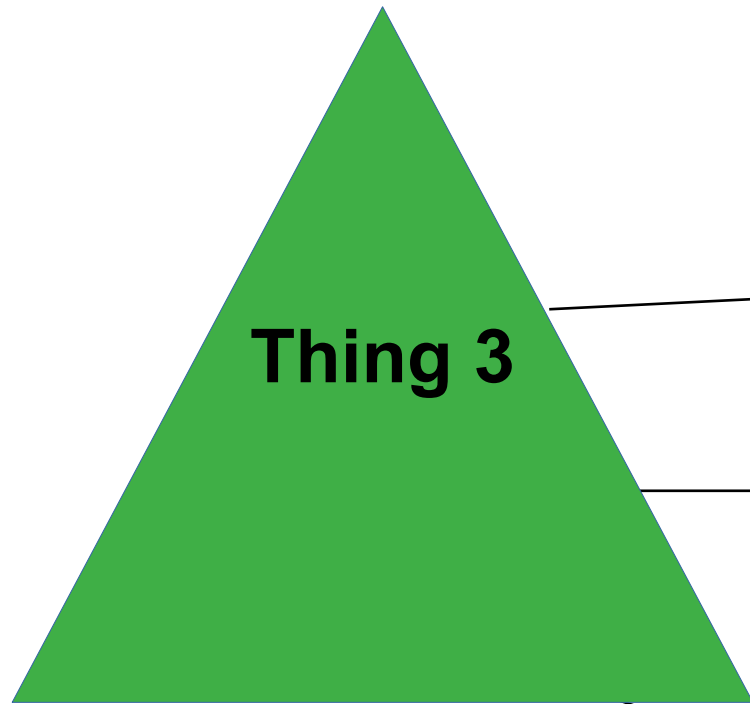
10cm

magenta

circle







Properties

Values

width :

15cm

colour :

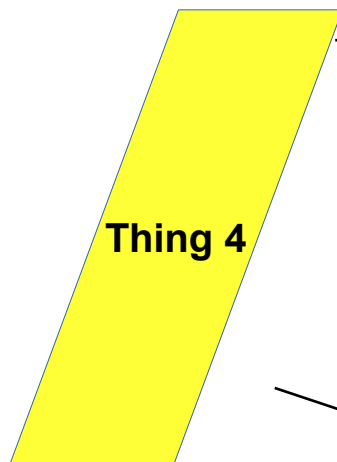
green

shape :

triangle

Properties

Values



width :

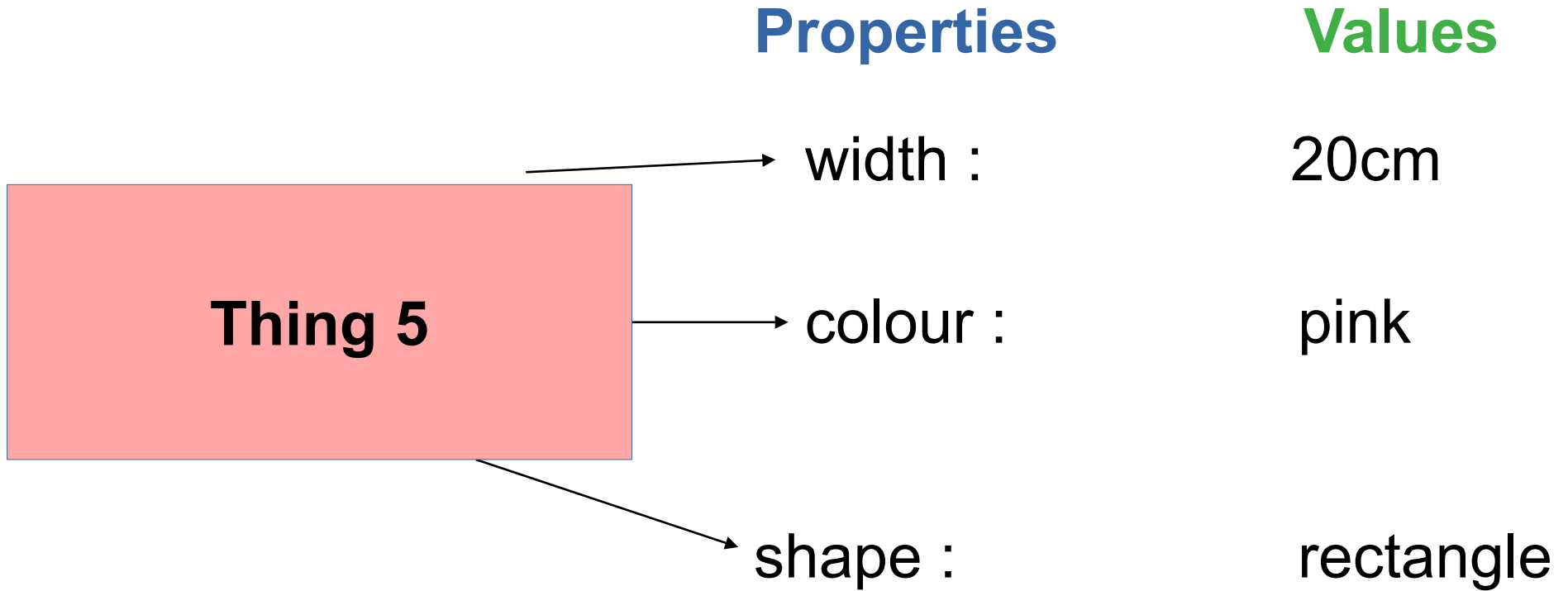
7cm

colour :

yellow

shape :

rhomboid

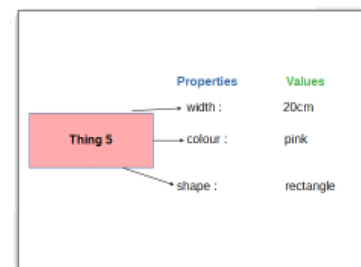
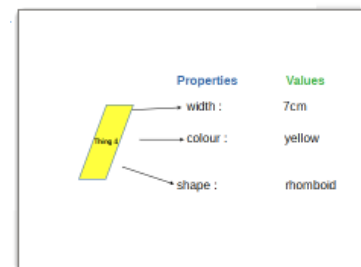
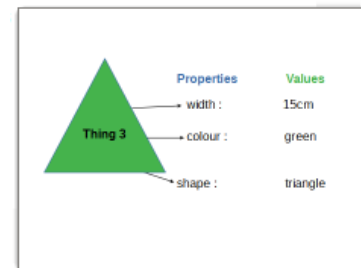
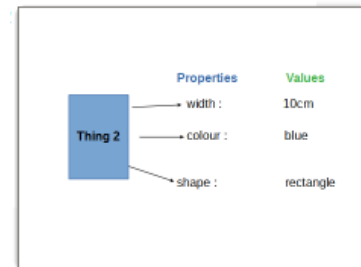
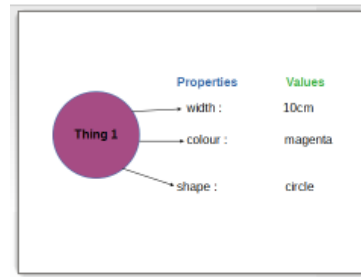


WE HAVE
4 properties (variables,
attributes, features)

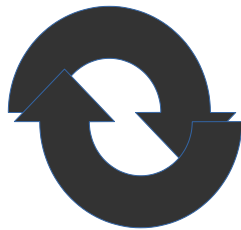
ID	width (cm)	colour	shape
1	10	magenta	circle
2	10	blue	rectangle
3	15	green	triangle
4	7	yellow	rhomboid
5	20	pink	rectangle

This table is what we
analyse

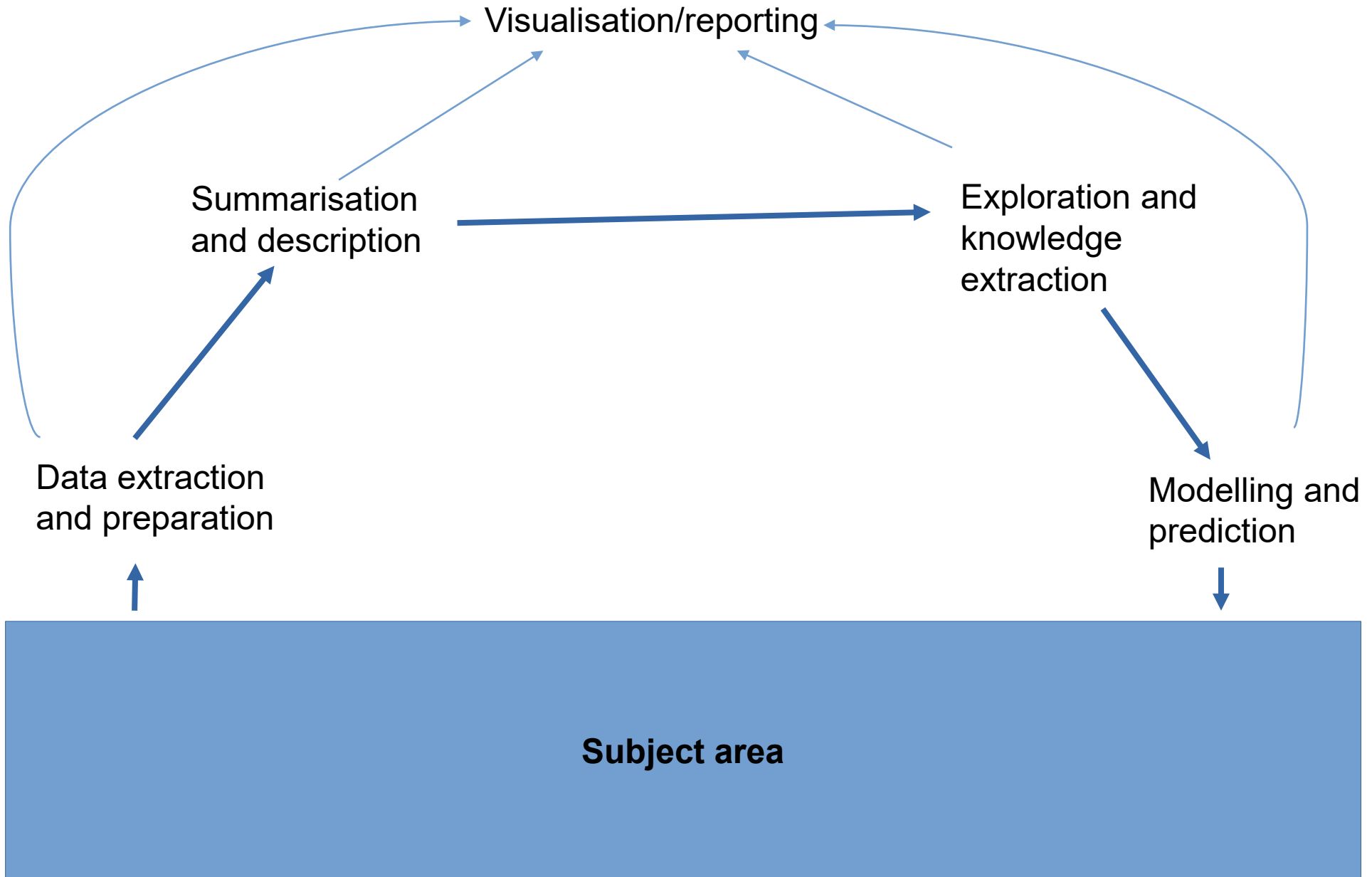
WE HAVE
5 things
(instances,
observations,
examples)



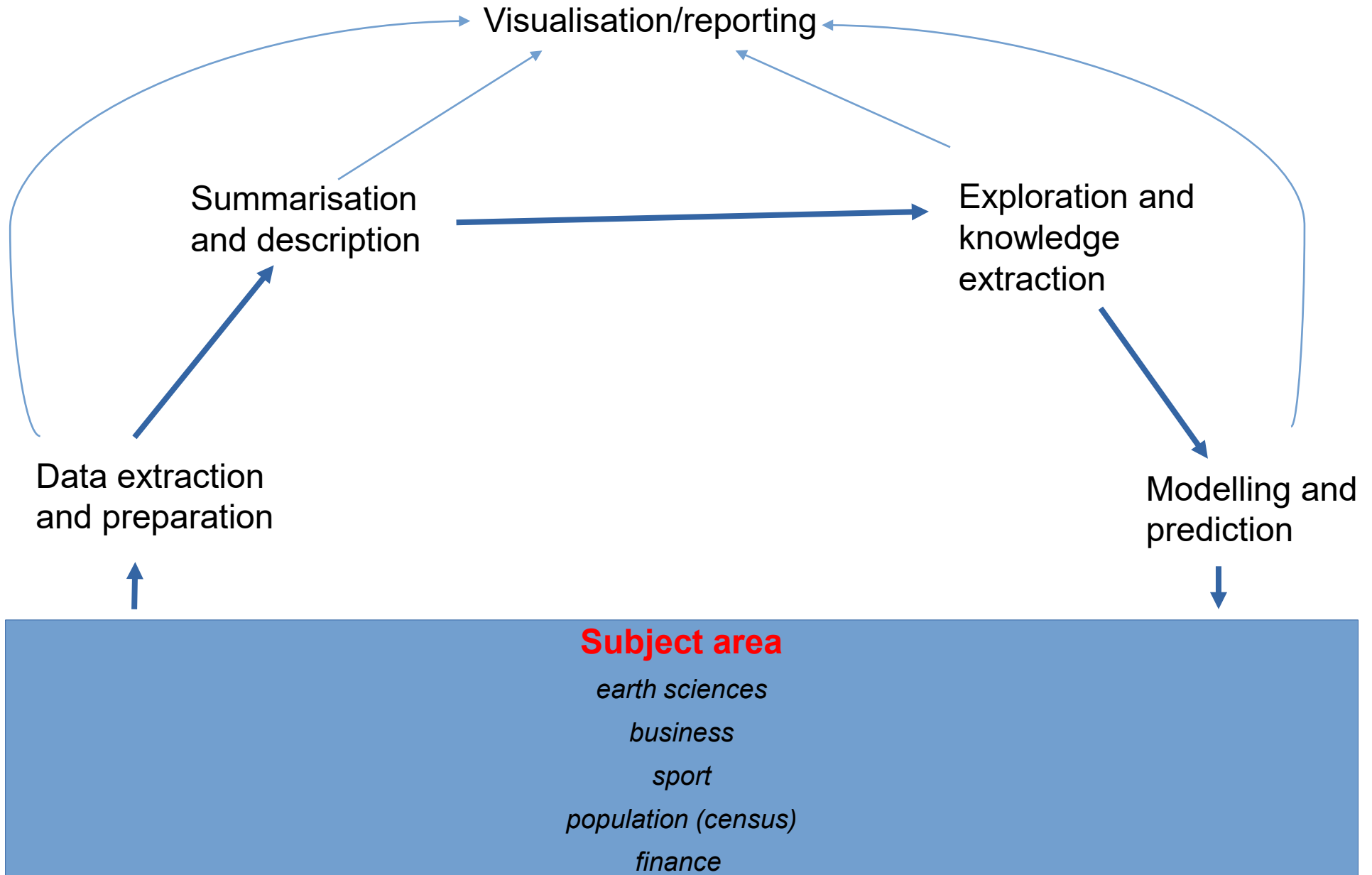
The analysis cycle



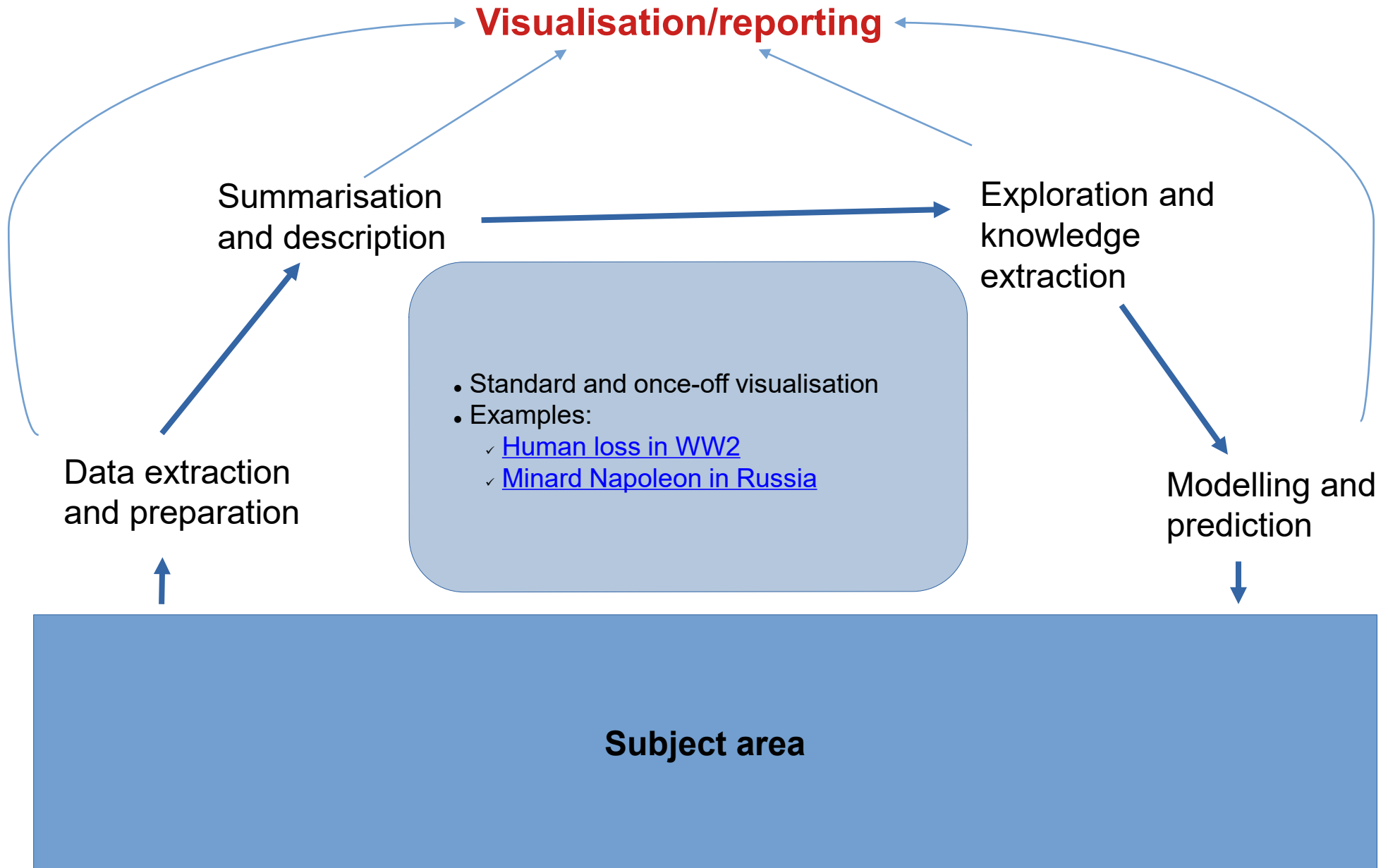
The data analysis cycle



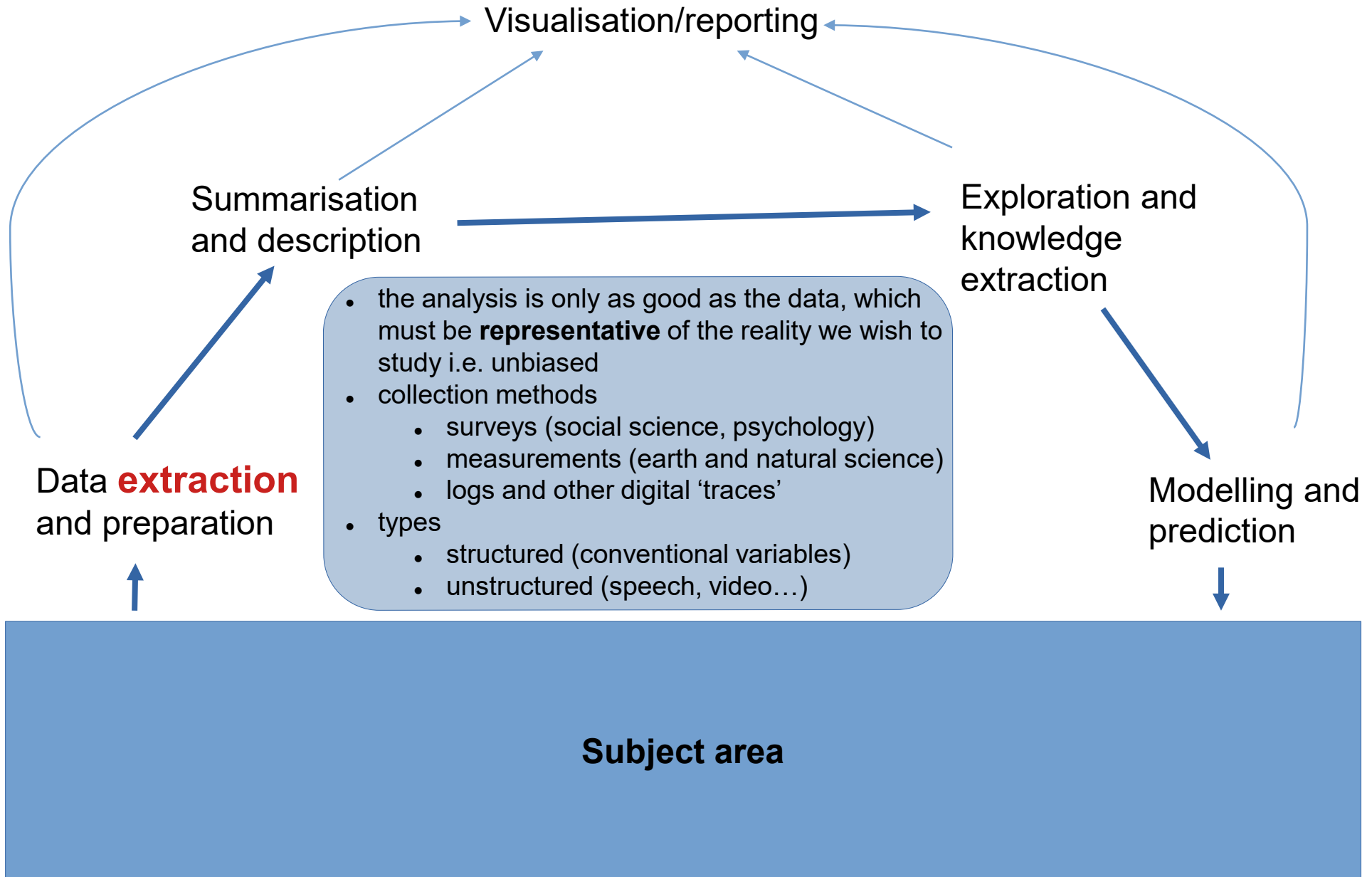
The data analysis cycle



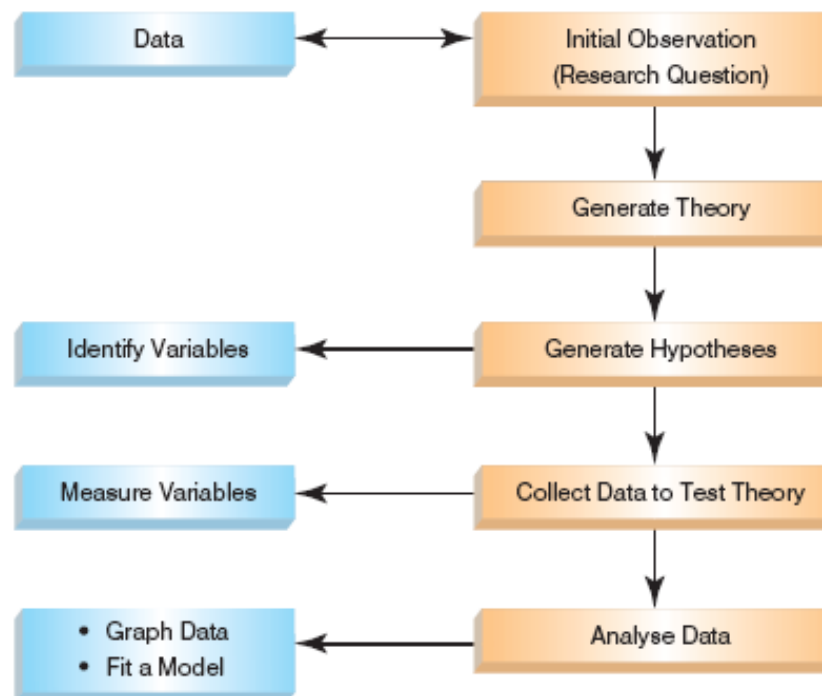
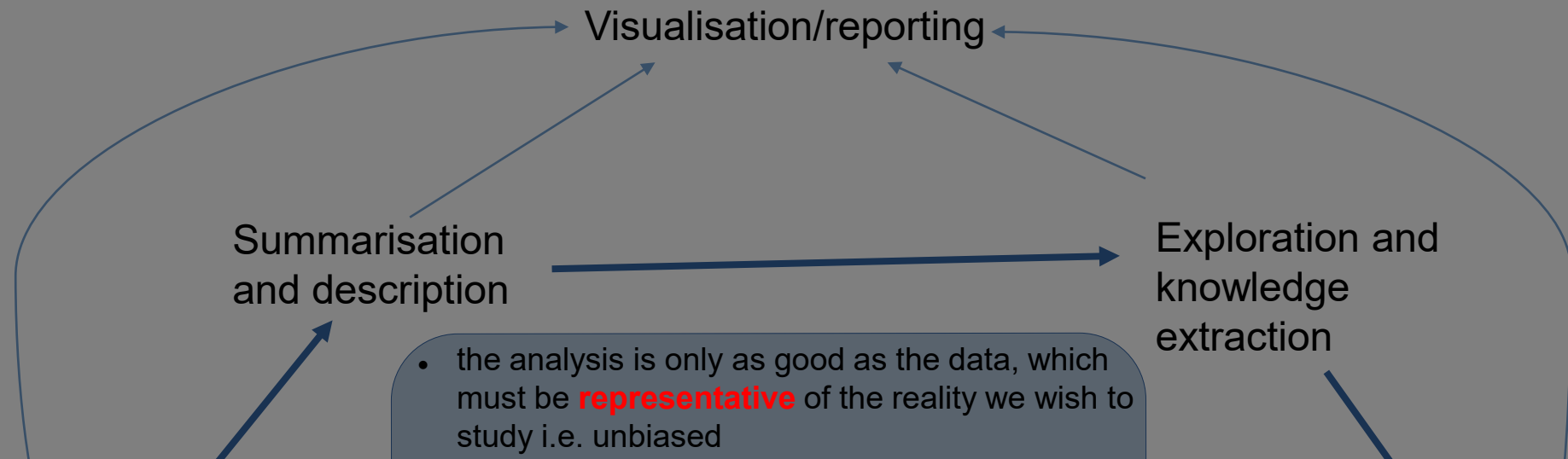
The data analysis cycle



The data analysis cycle



The data analysis cycle

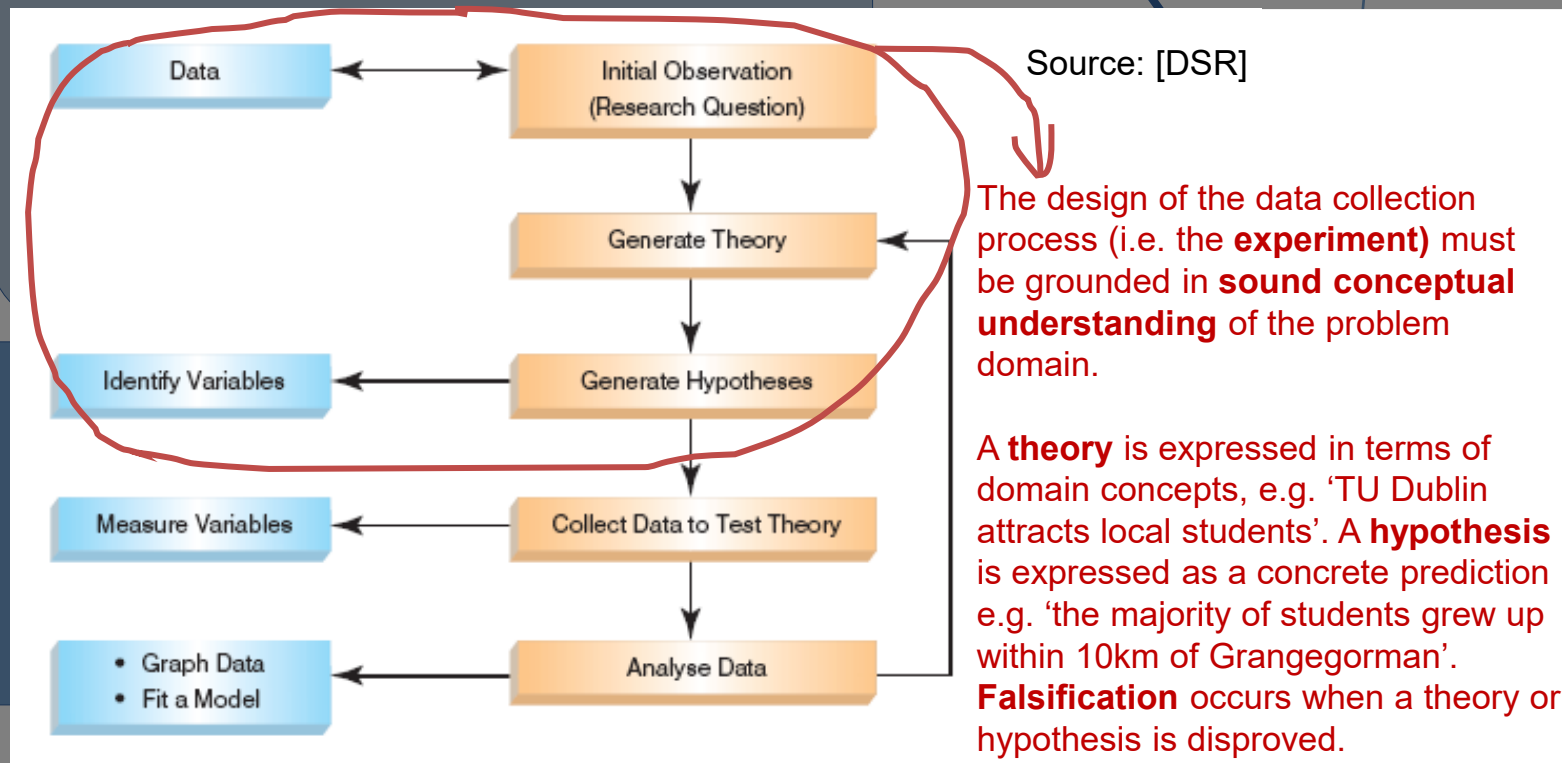
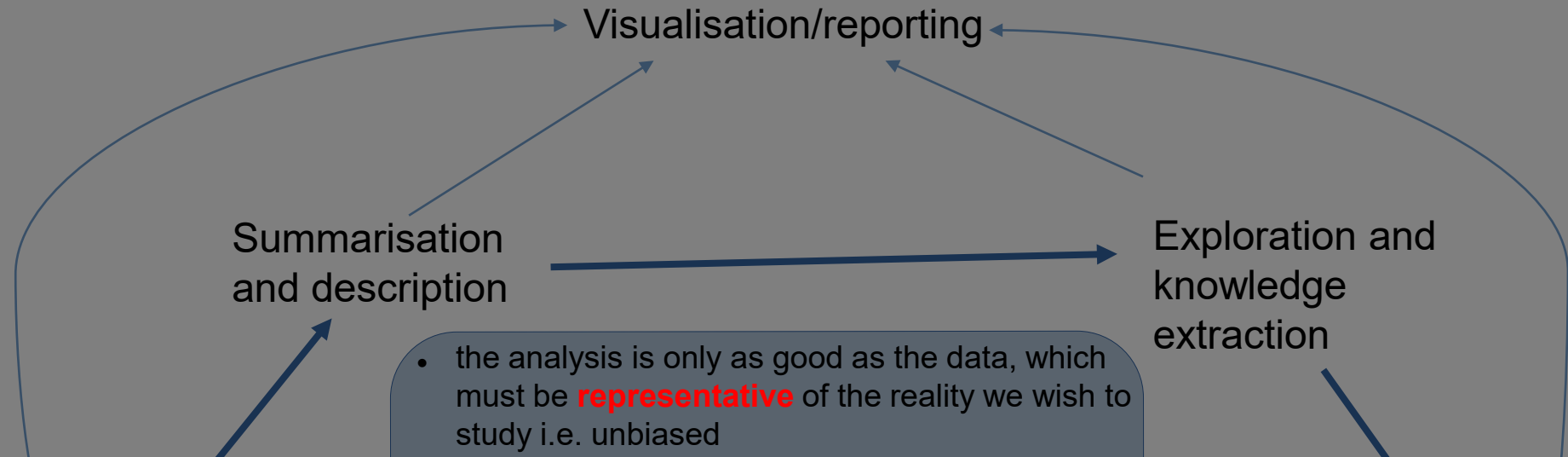


Source: [DSR]

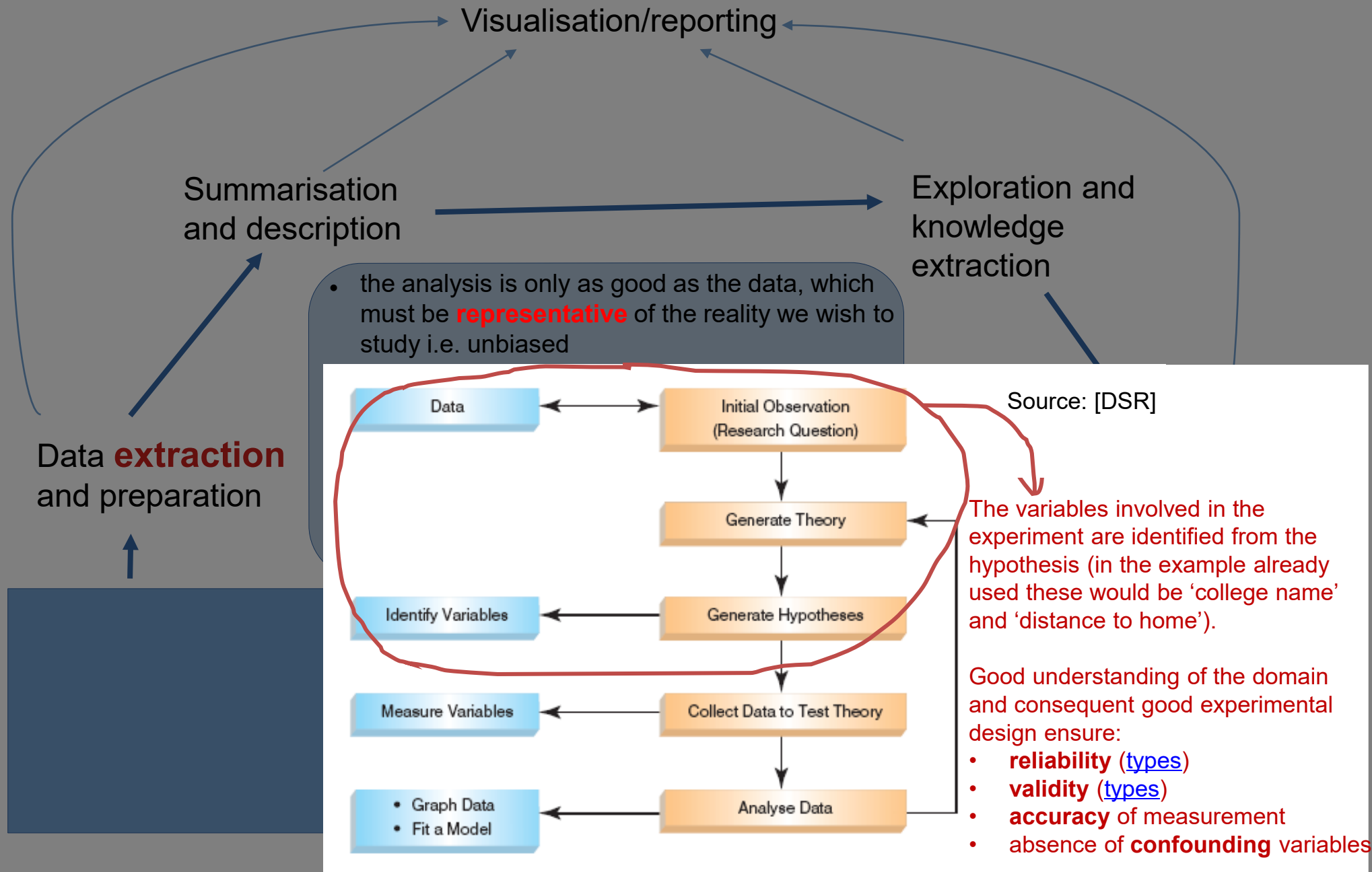
Often, data analysis is performed simply because **data are available**. This is the case with enterprise data warehouses, which are explored for knowledge hidden in data already available as a by-product of a company's commercial activity.

The older 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.

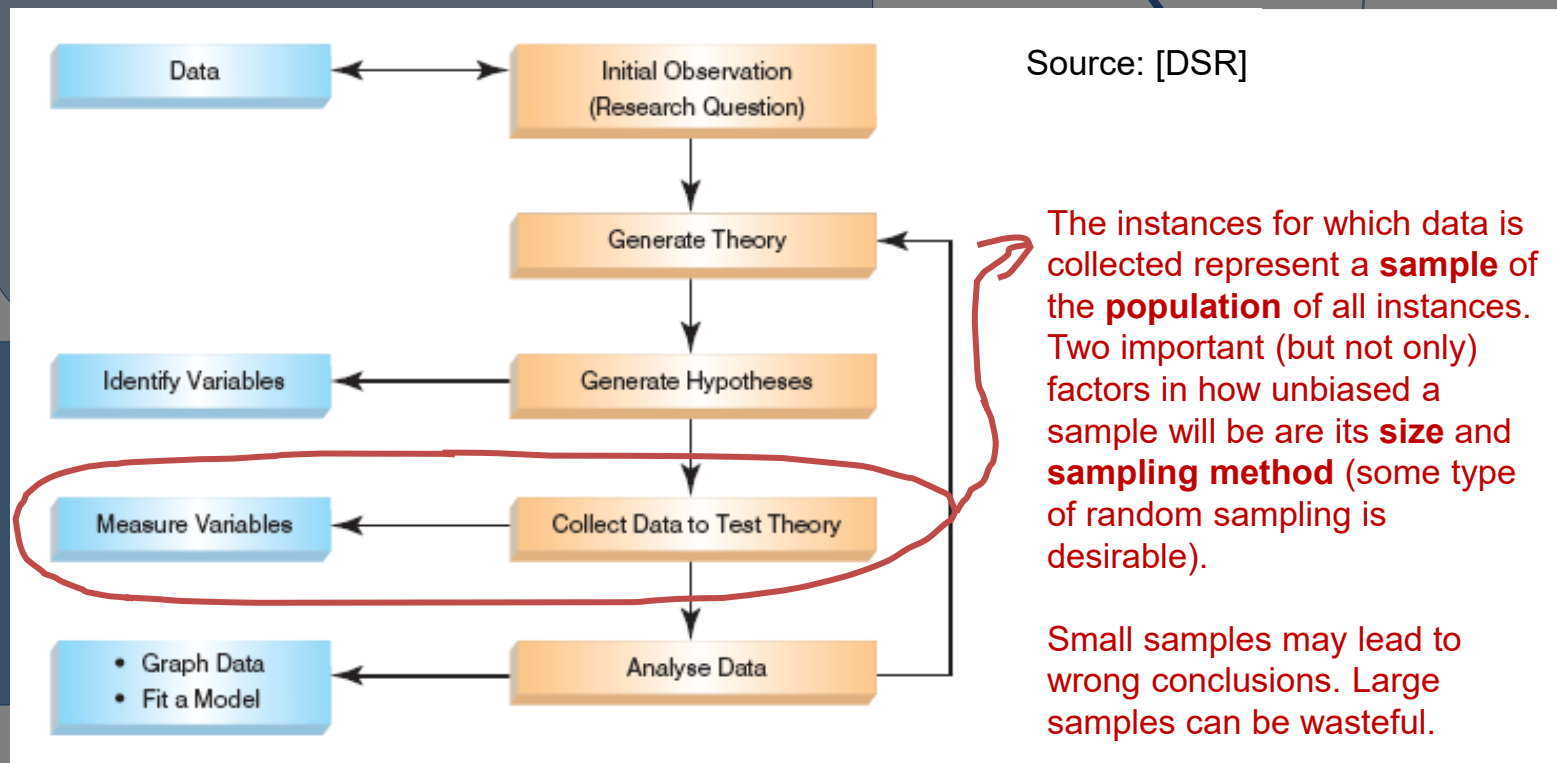
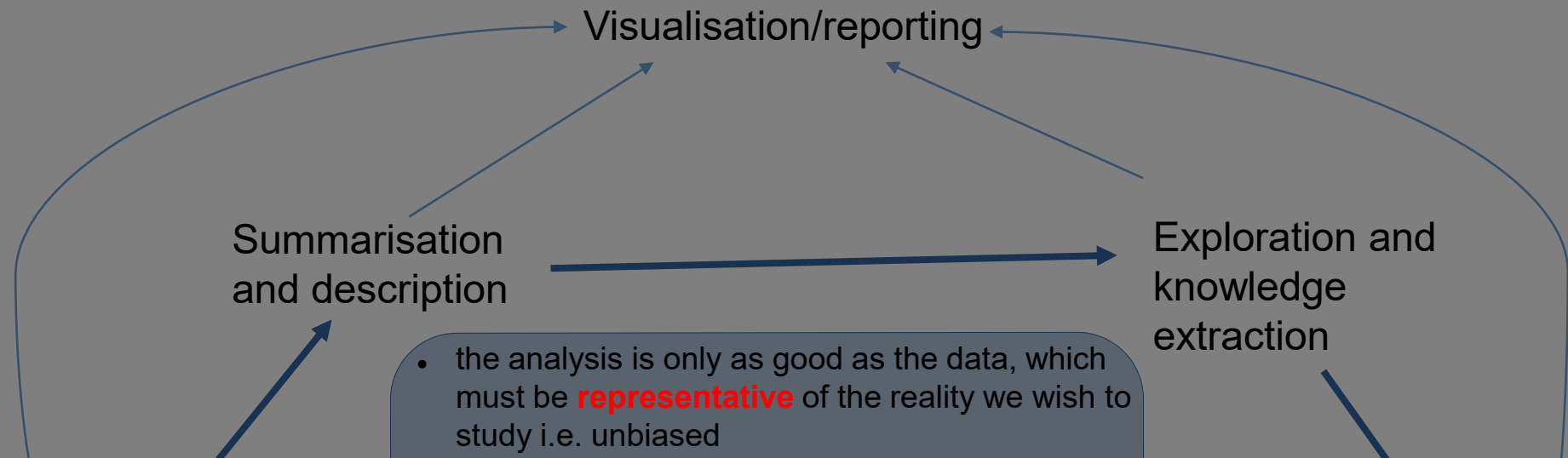
The data analysis cycle



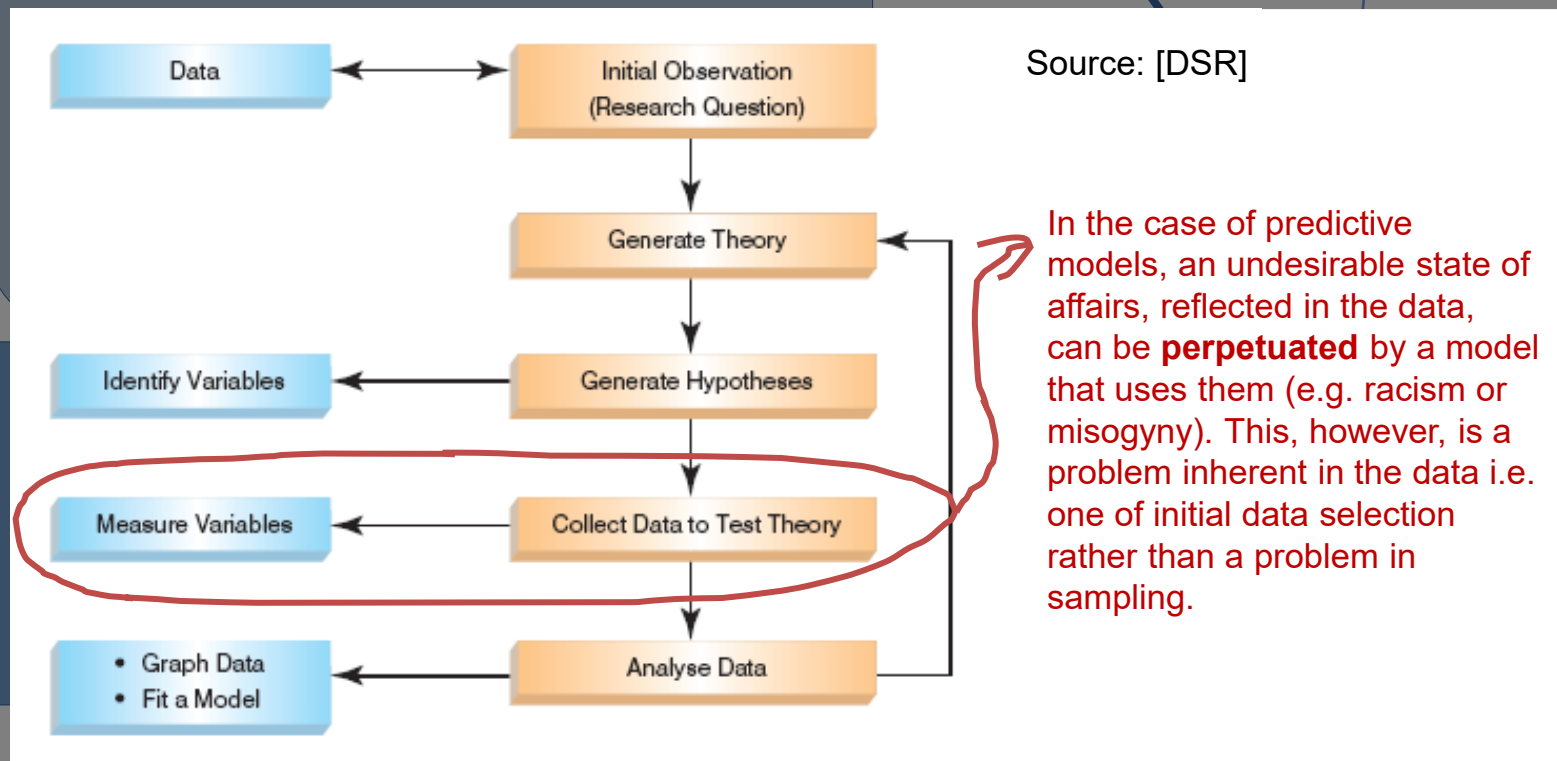
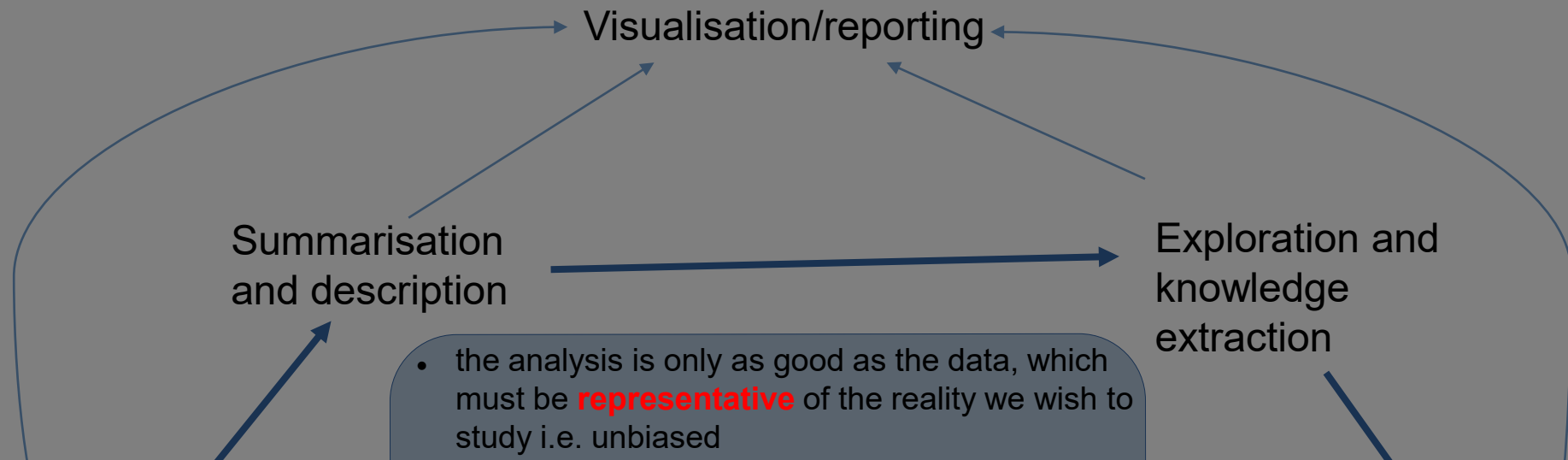
The data analysis cycle



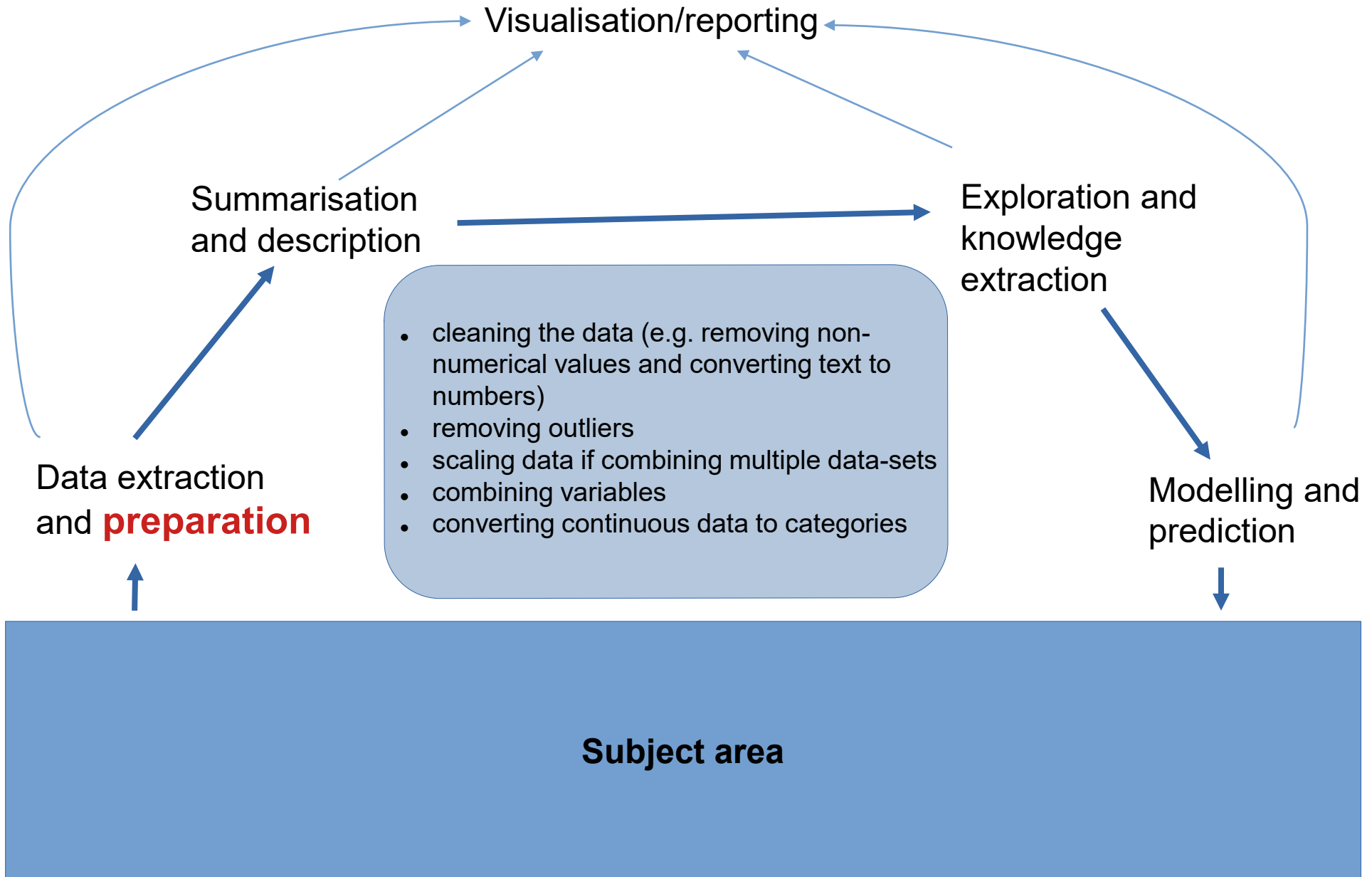
The data analysis cycle



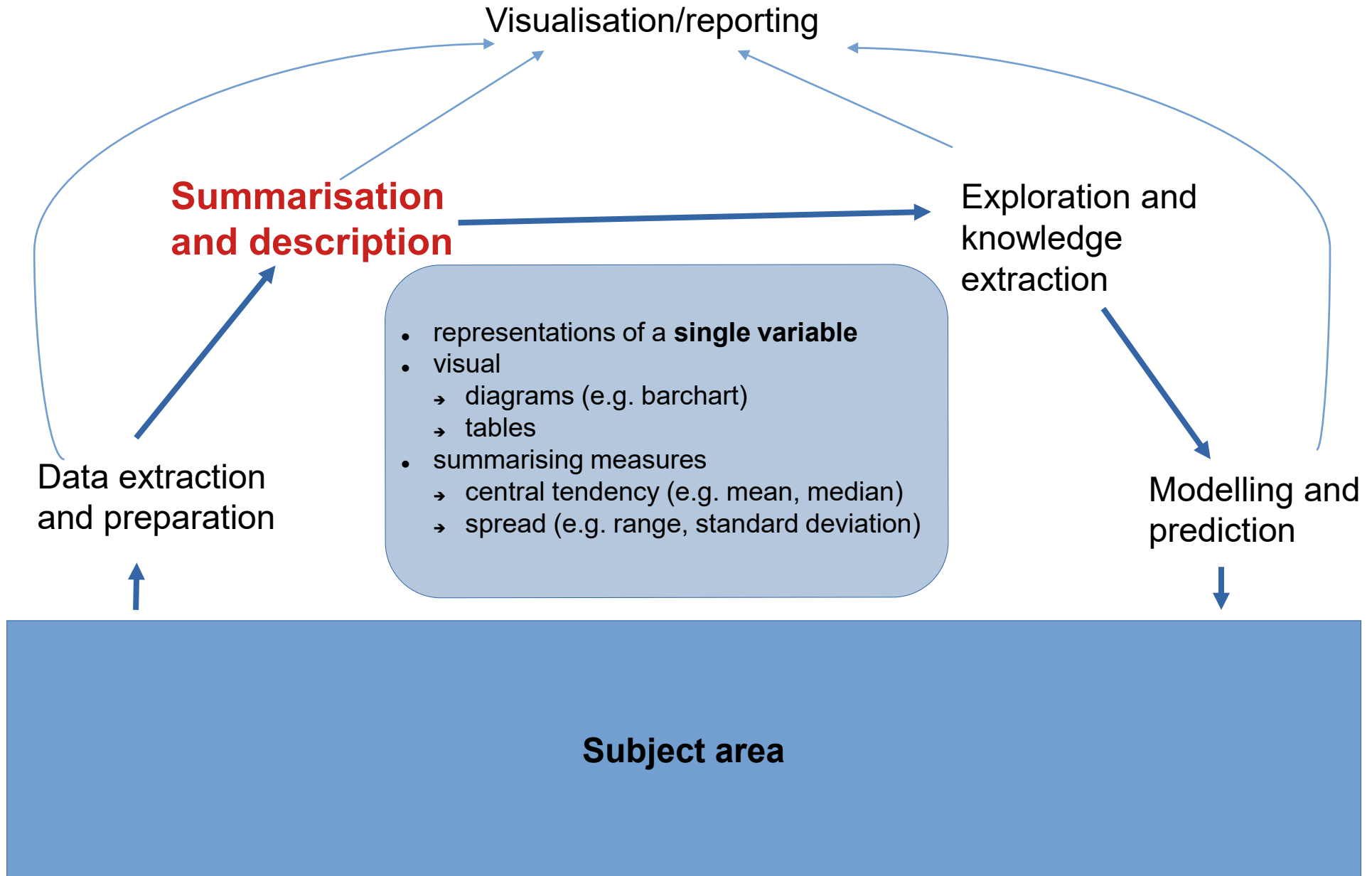
The data analysis cycle



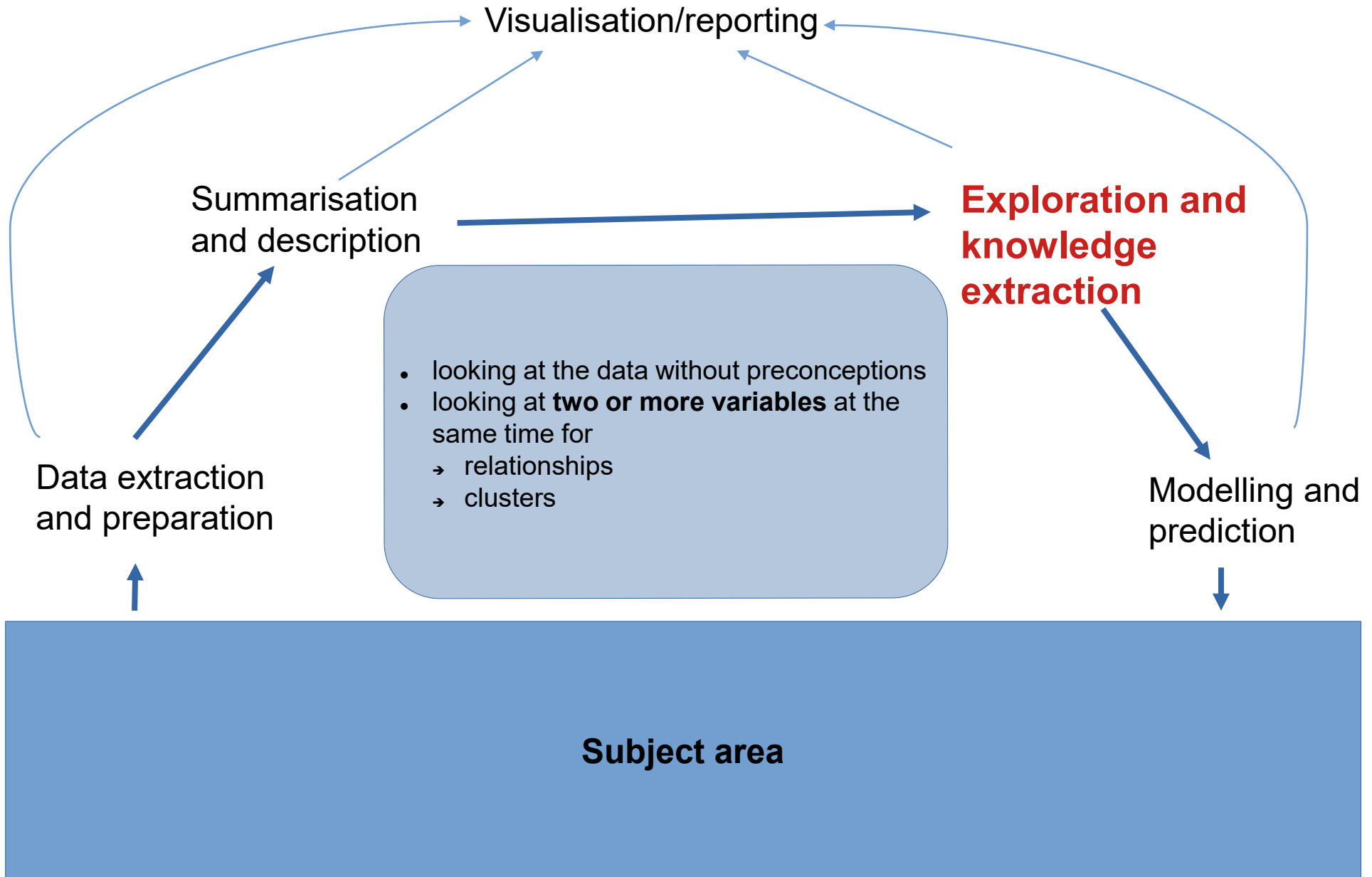
The data analysis cycle



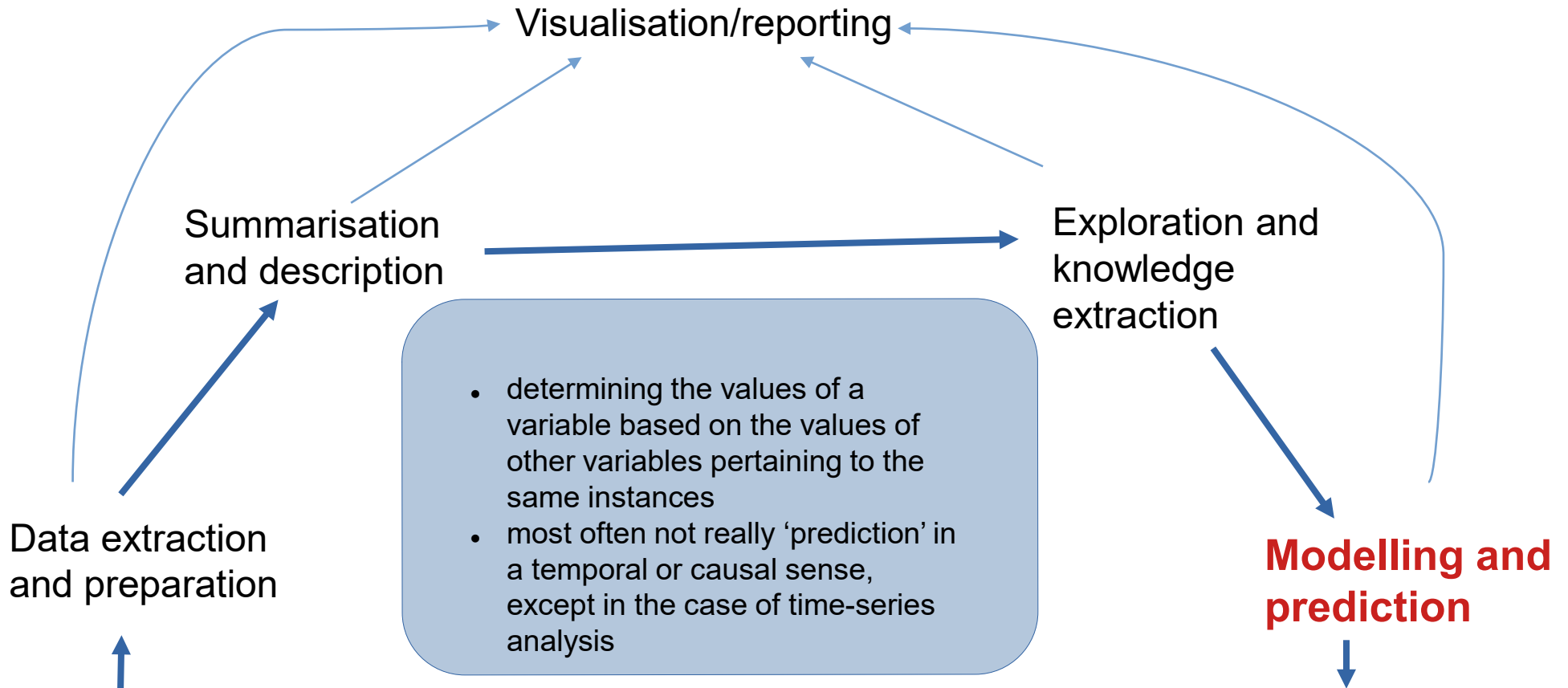
The data analysis cycle



The data analysis cycle



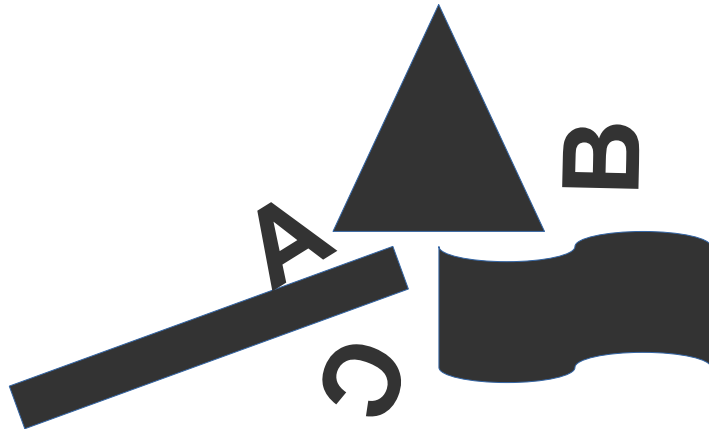
Subject area



[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

Statistics in the data analysis landscape



The data analysis landscape

Data science: concepts and theory behind the analysis of data

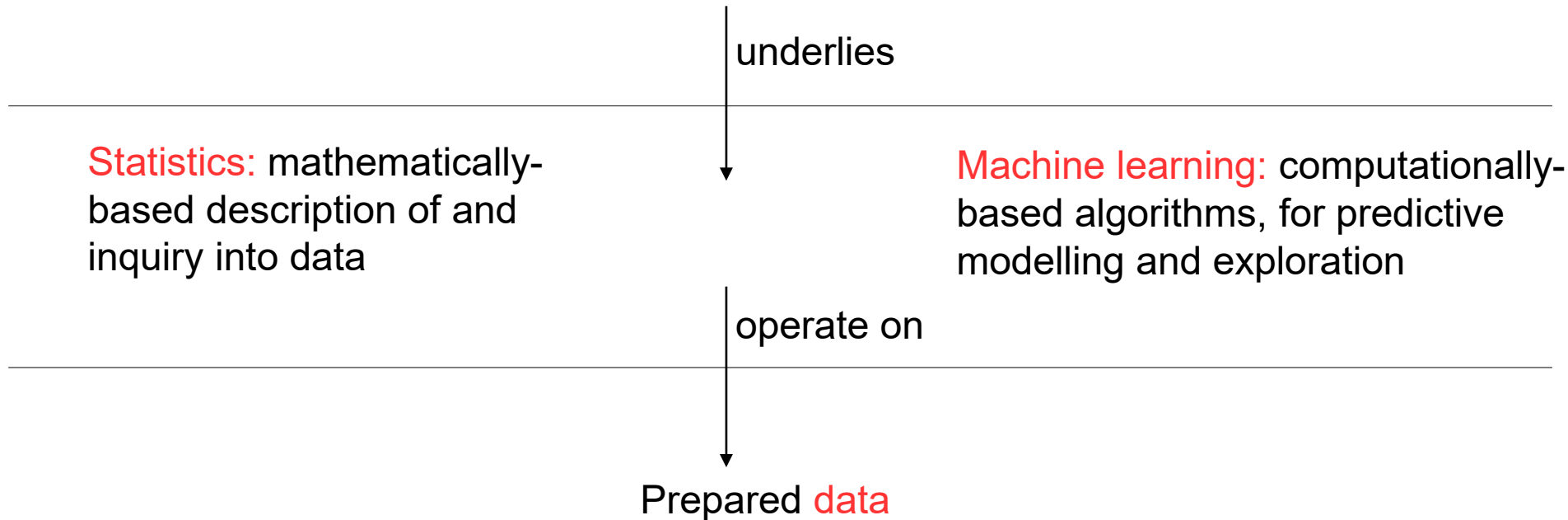
underlies

Statistics: mathematically-based description of and inquiry into data

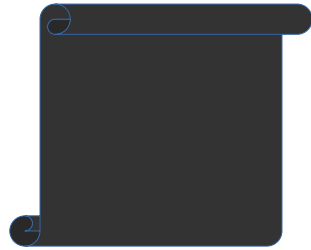
Machine learning: computationally-based algorithms, for predictive modelling and exploration

operate on

Prepared **data**



In this module



you will learn

to...

- Present the **question** you are interested in
in a way that makes sense to conduct a statistical analysis
- Inspect and **prepare** the data you have
to support a statistical analysis
- **Describe** the data you have
in a way that demonstrates the constraints placed on the analysis by the data
- Conduct a **statistical analysis**
using appropriate statistical tests
- **Interpret the outcomes** of your statistical analysis
drawing appropriate conclusions
- **Report** on the findings of your statistical analysis
in a way that makes sense for your consumer

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.