# Lecture 12: Single variable data analysis

Instructor: Michael Szell
Oct 6, 2023
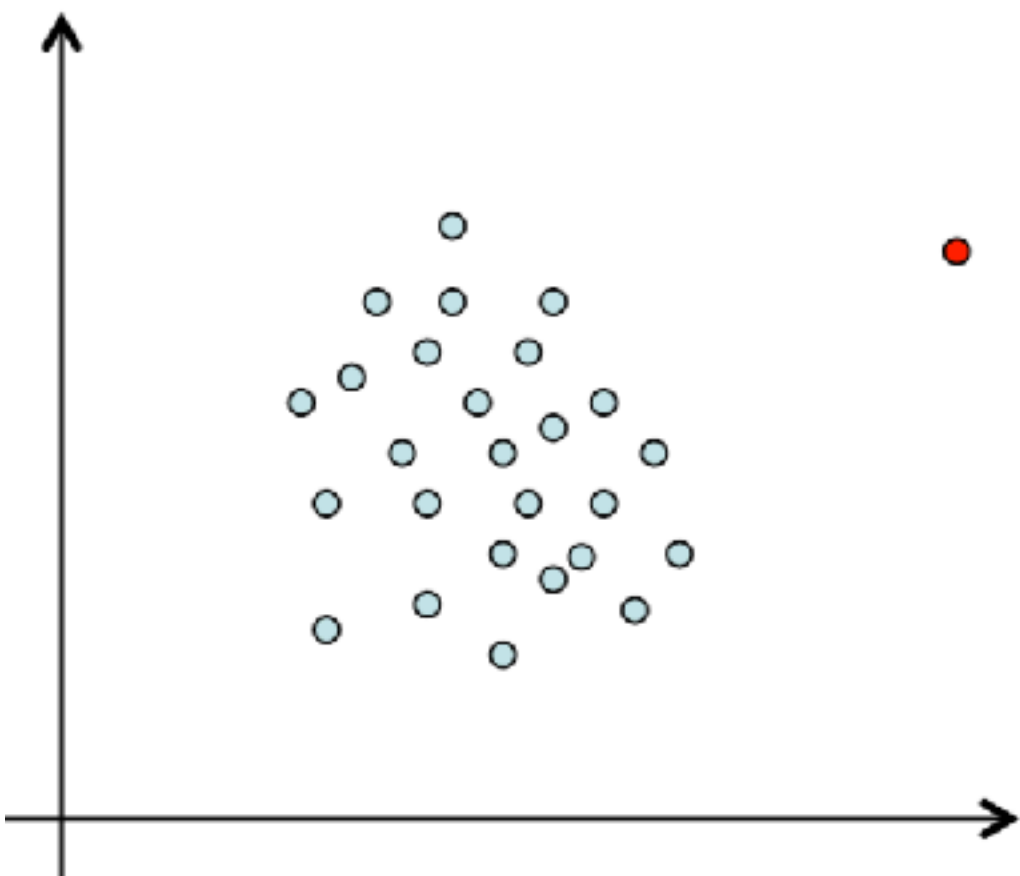


IT UNIVERSITY OF COPENHAGEN

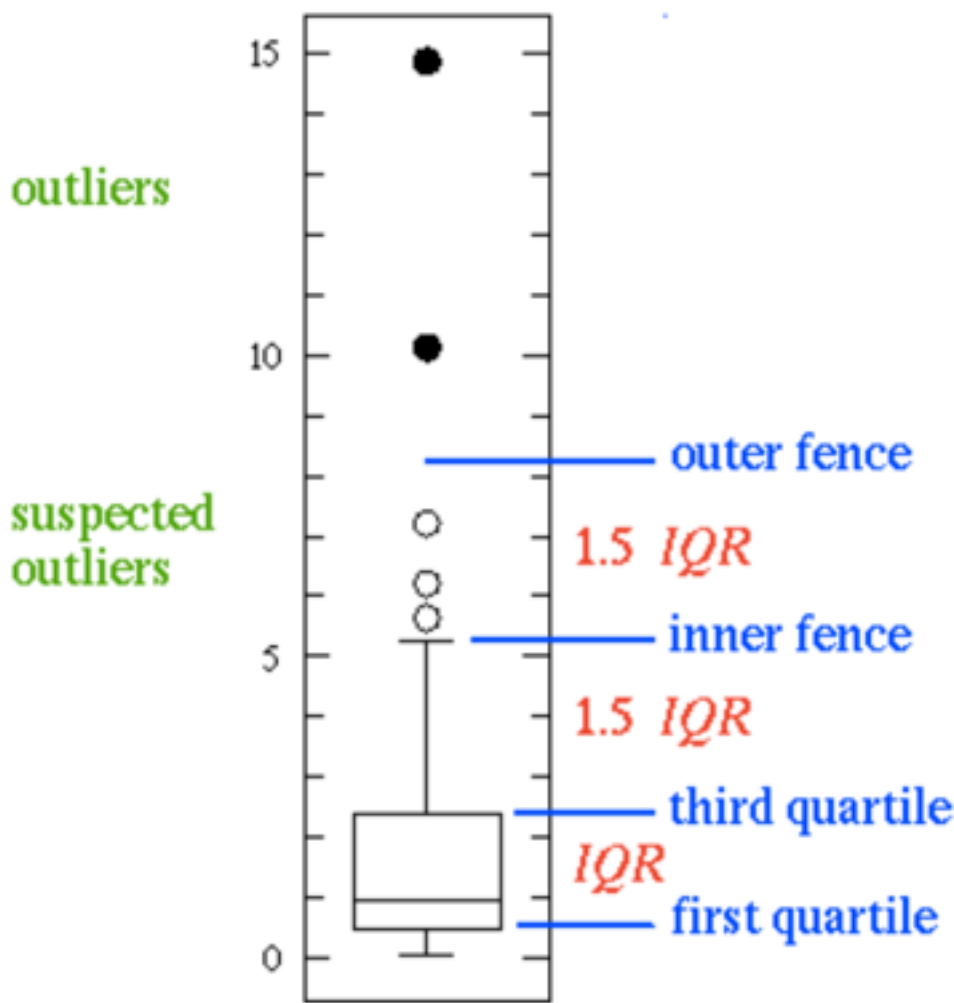# Today you will learn first steps in analyzing data

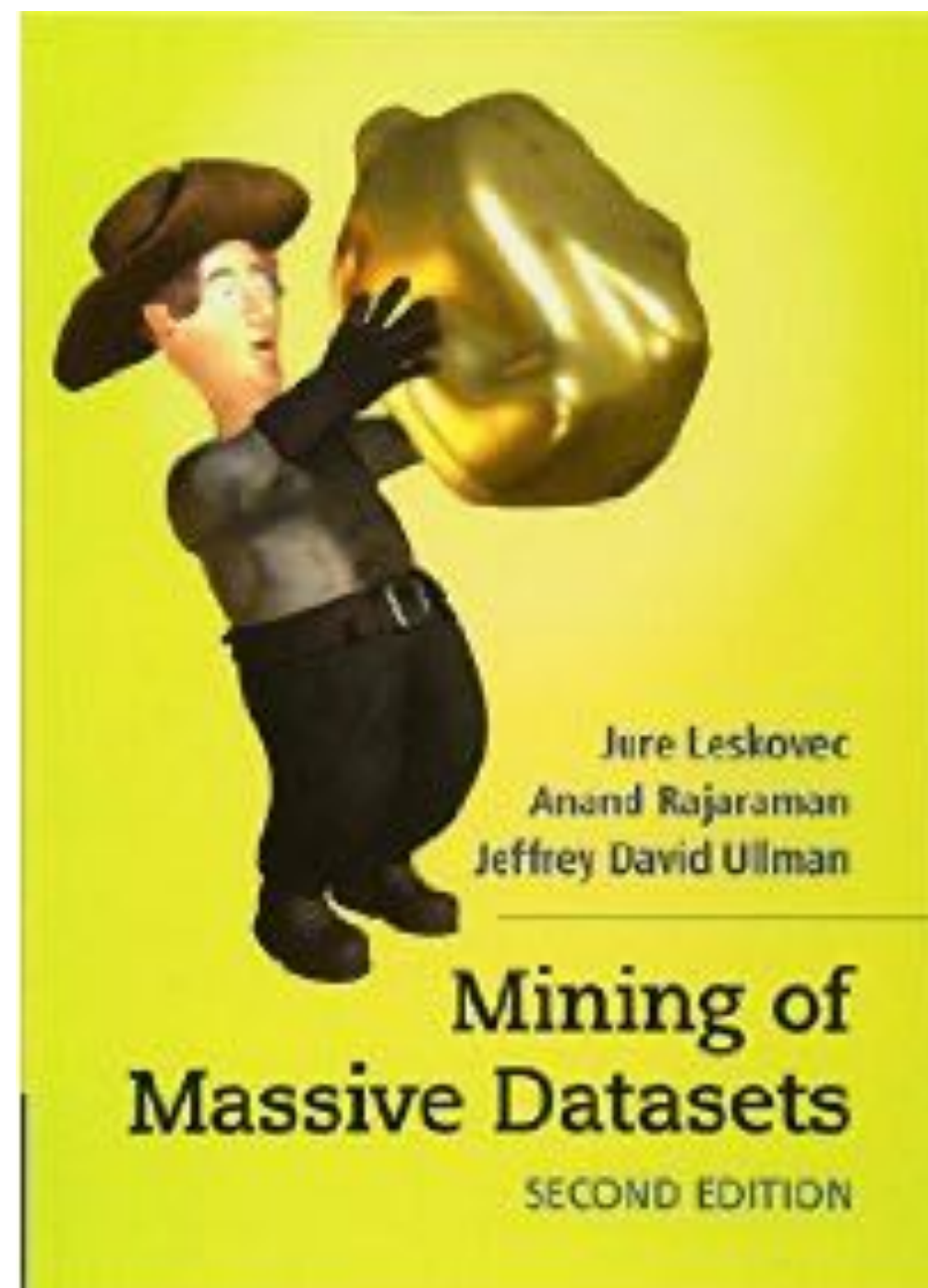## Variable types



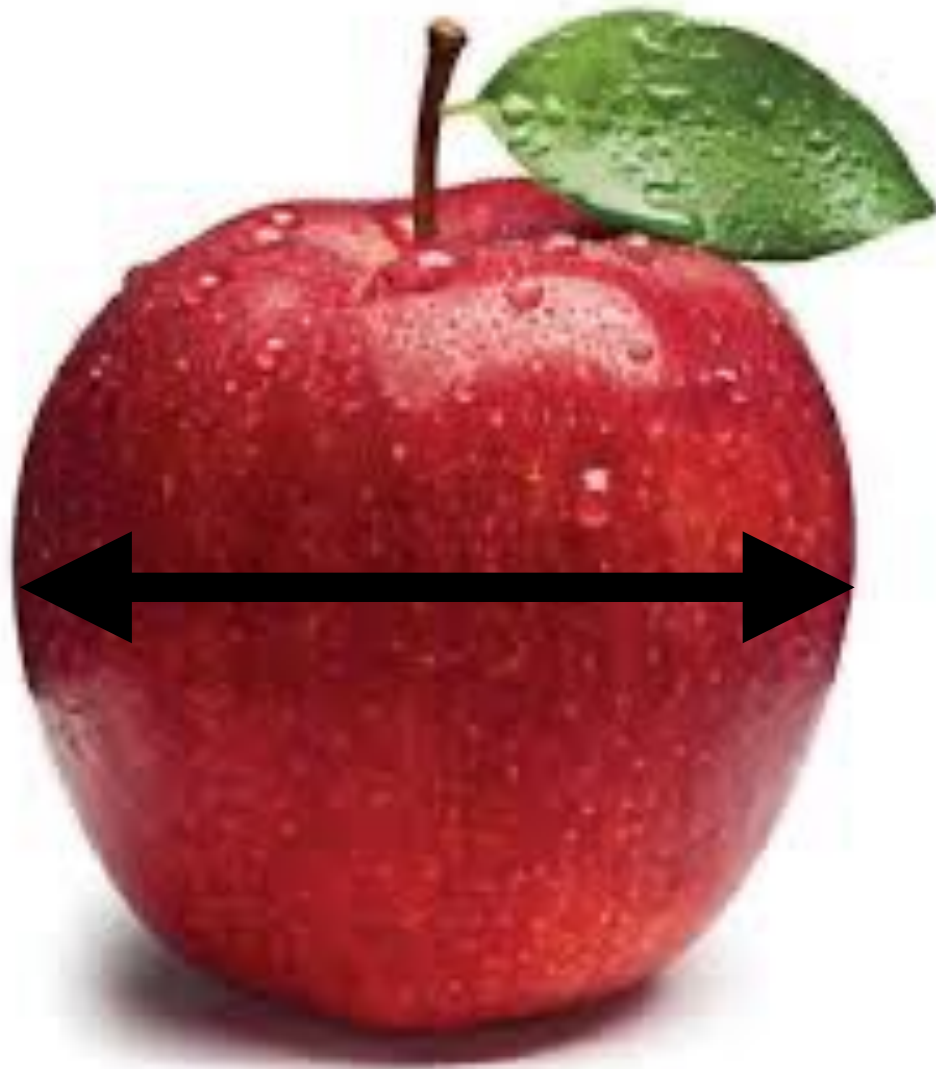## Exploratory data analysis



## Describing data

# Data analysis is the process of:

Cleaning, transforming, exploring and/or modeling data

with the goal of discovering useful information,
informing conclusions or decision-making
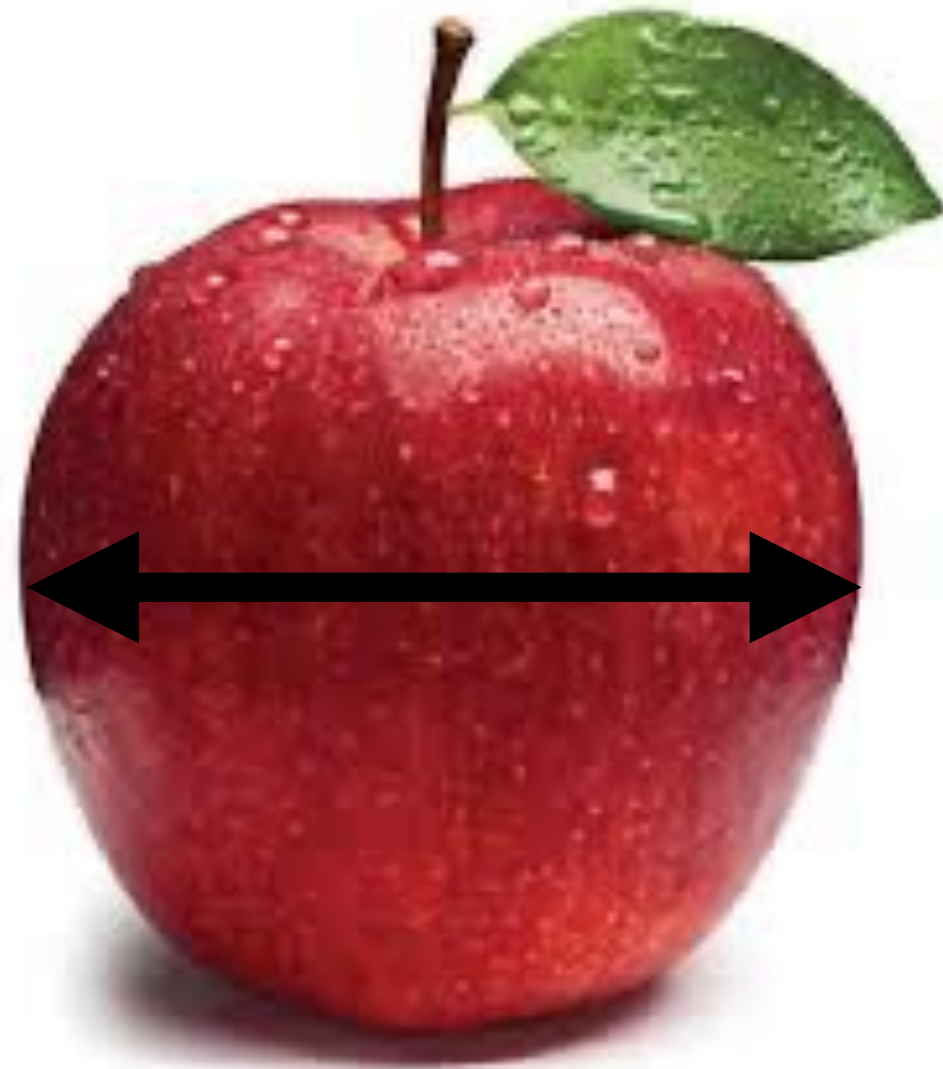
1) Descriptive statistics

# There are 3 types of data analysis
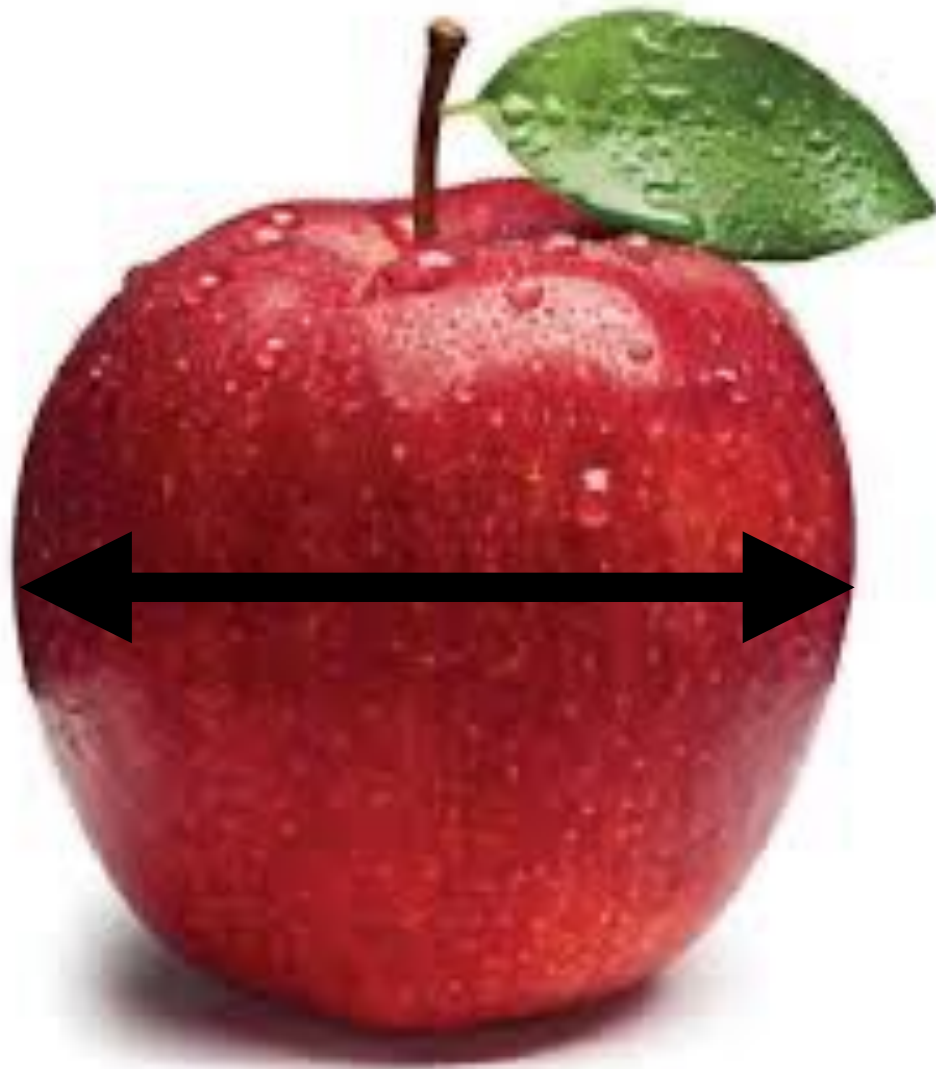
1) Descriptive statistics
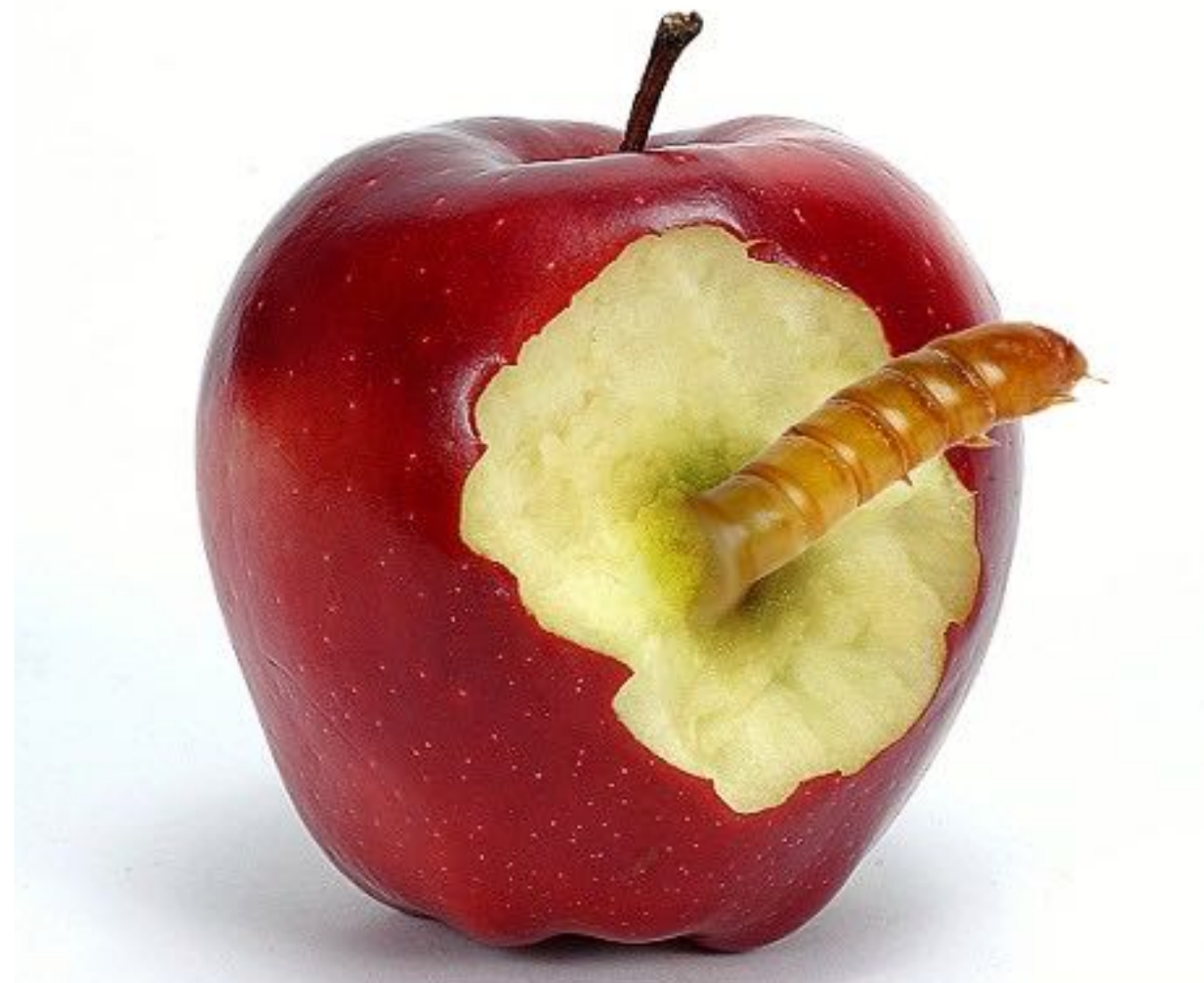
2) Exploratory

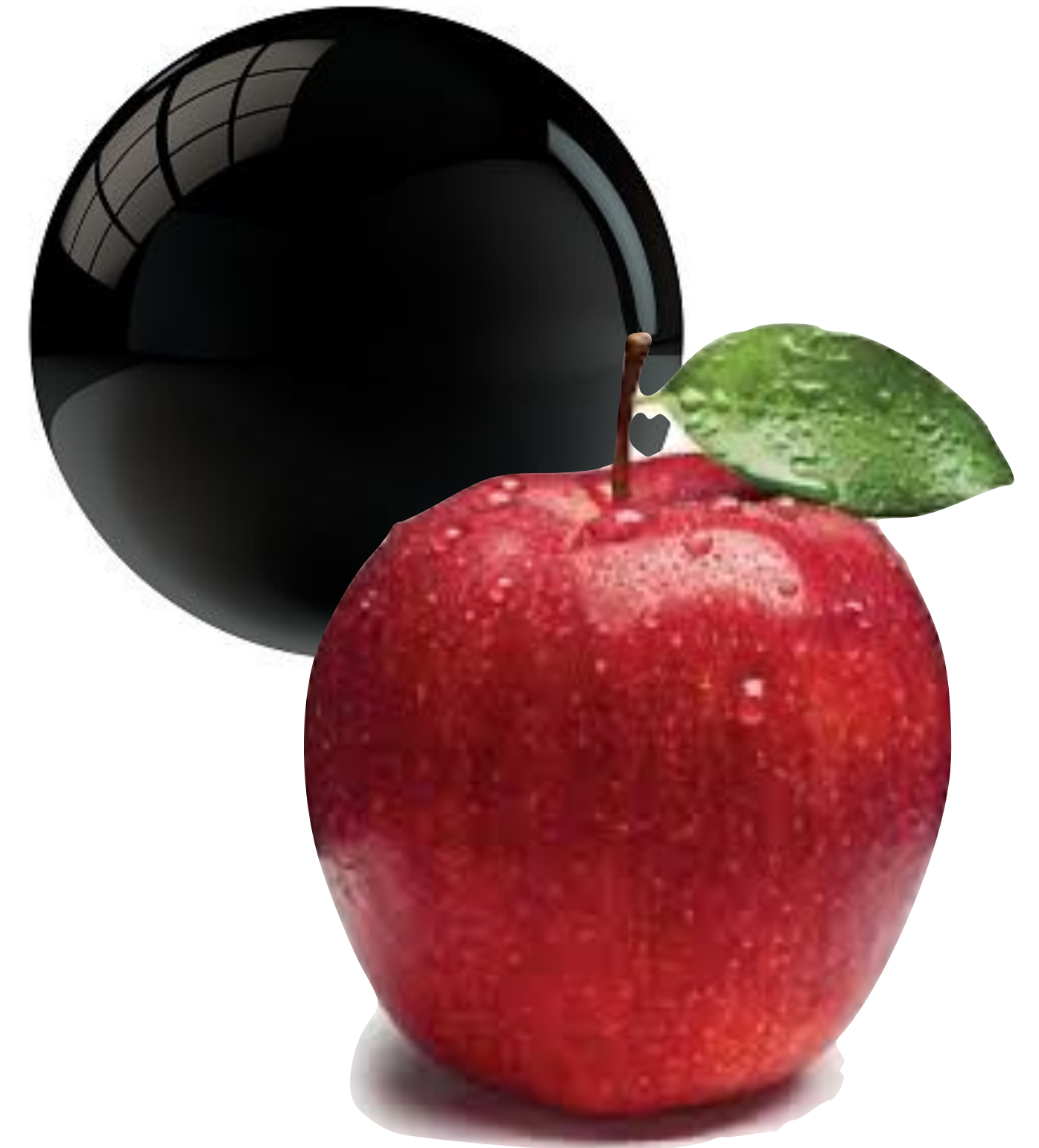# There are 3 types of data analysis

1) Descriptive statistics

2) Exploratory

3) Inferential statistics
(Hypothesis testing)

# Today we focus on descriptive statistics and exploration

1) Descriptive statistics

2) Exploratory

3) Inferential statistics
(Hypothesis testing)

# There are many steps in data mining/analysis

Data Mining is short for Knowledge Discovery from Data (KDD):

Input data

→

Preprocessing
1. Data cleaning
2. Data integration
3. Data selection
4. Data transformation

Other names:
Data wrangling,
Data munging,
Data preparation

Data mining
5. Data mining

Postprocessing
6. Pattern evaluation
7. Knowledge presentation

→ Information

# Data sets have objects and attributes

Data set

| Student ID | Year | Grade Point Average (GPA) | ... |
|---|---|---|---|
| ⋮ | | | |
| ▸ 1034262 | Senior | 3.24 | ... |
| 1052663 | Sophomore | 3.51 | ... |
| 1082246 | Freshman | 3.62 | ... |
| ⋮ | | | |

# Data sets have objects and attributes

Data set

Attributes

| Student ID | Year | Grade Point Average (GPA) | ... |
|---|---|---|---|
| ⋮ | | | |
| ▸ 1034262 | Senior | 3.24 | ... |
| 1052663 | Sophomore | 3.51 | ... |
| 1082246 | Freshman | 3.62 | ... |
| ⋮ | | | |

Data object

Data object = record, individual, point, event, observation, vector, entity

Attribute = field, feature, variable, dimension, characteristic

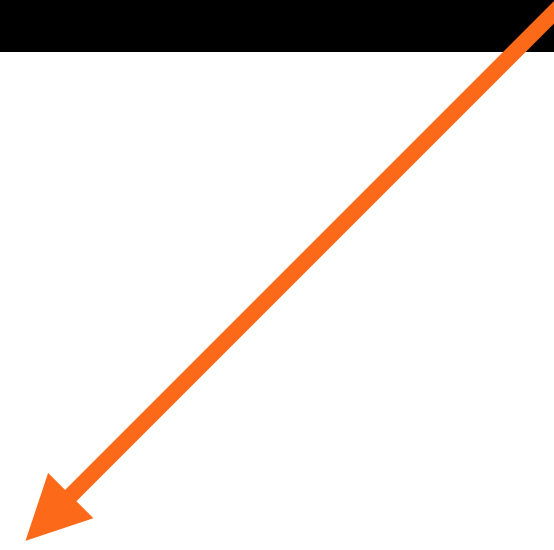# In this lecture we will deal with single-variable analysis

Data set

| Student ID | Year | Grade Point Average (GPA) | ... |
|---|---|---|---|
| ⋮ | ⋮ | | |
| ▸ 1034262 | Senior | 3.24 | ... |
| 1052663 | Sophomore | 3.51 | ... |
| 1082246 | Freshman | 3.62 | ... |
| ⋮ | ⋮ | | |

Data object = record, individual, point, event, observation, vector, entity

Attribute = field, feature, variable, dimension, characteristic

# There are two types of variables: categorical and quantitative

Places an individual into one of several categories

# Categorical variables can be nominal or ordinal

Places an individual into one of several categories



No order                    Order

# There are two types of variables: categorical and quantitative

Places an individual into one of several categories

Takes values for which arithmetic operations make sense

# Quantitative variables can be interval or ratio

Places an individual into one of several categories

Takes values for which arithmetic operations make sense



Differences meaningful

Ratios also meaningful

| Zip code | Street number | C° | Age |
| Student ID | | | K° |

| Nominal | Ordinal | Interval | Ratio |

# Jupyter

# Outliers can be a sign for low data quality

Outliers (anomalous objects or values):
1) Data objects that have characteristics different from most others, or
2) Values of an attribute that are unusual

# Outliers can be a sign for low data quality

Outliers (anomalous objects or values):
1) Data objects that have characteristics different from most others, or
2) Values of an attribute that are unusual

This is not just noise! An outlier is an event that is suspected of not being generated by the same mechanisms as the rest of the data.

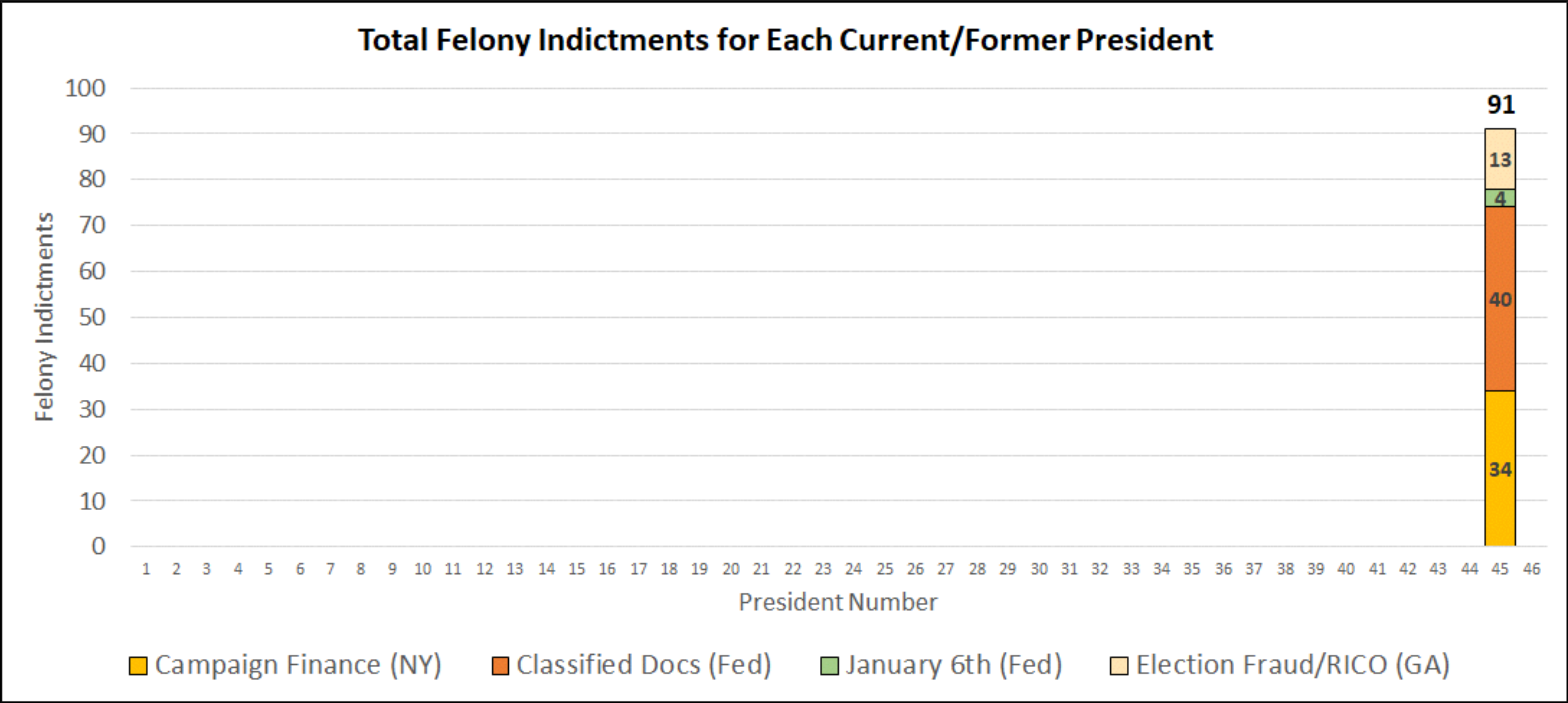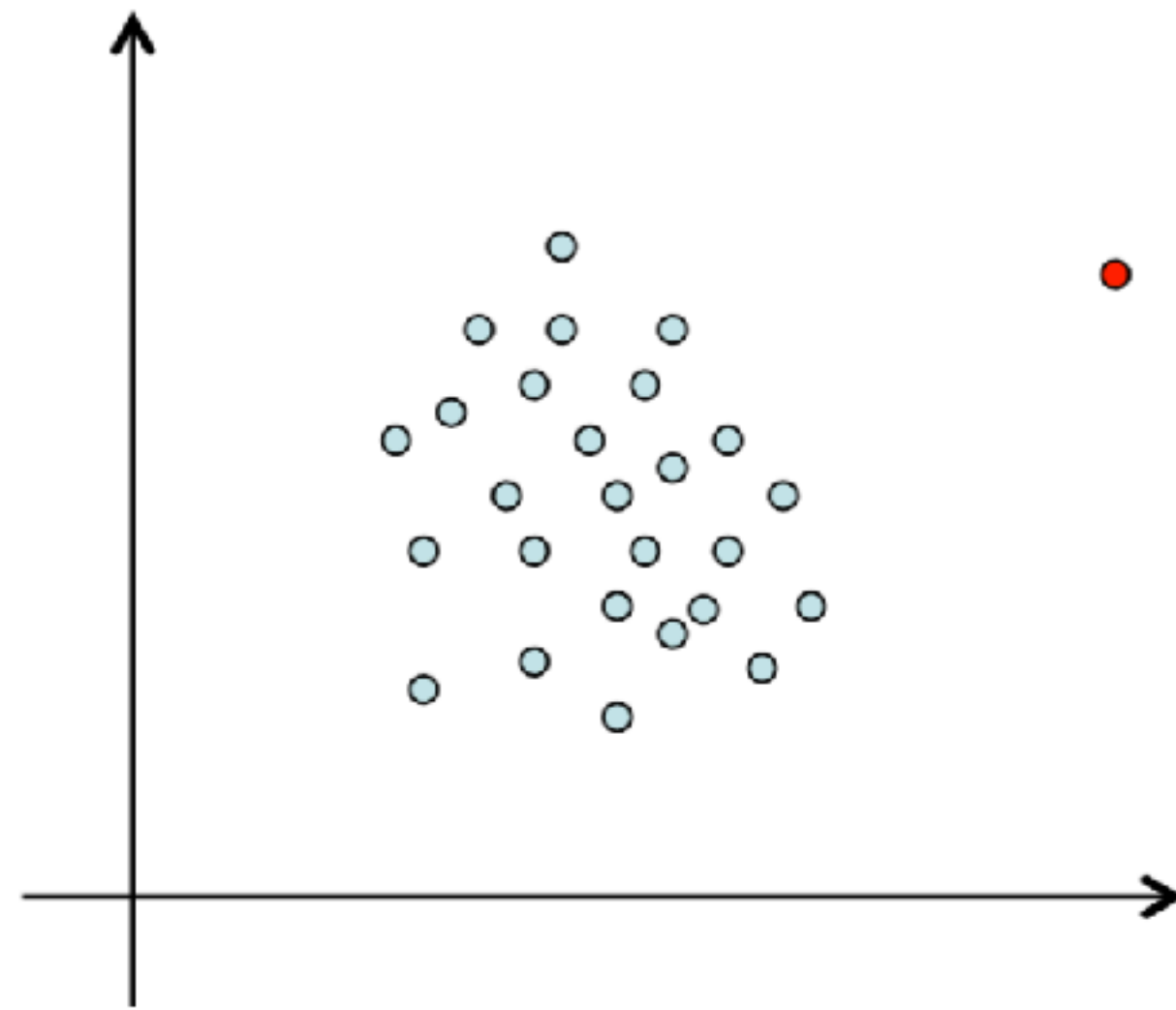# Outliers can be legitimate, interesting objects



Fraud



Innovation

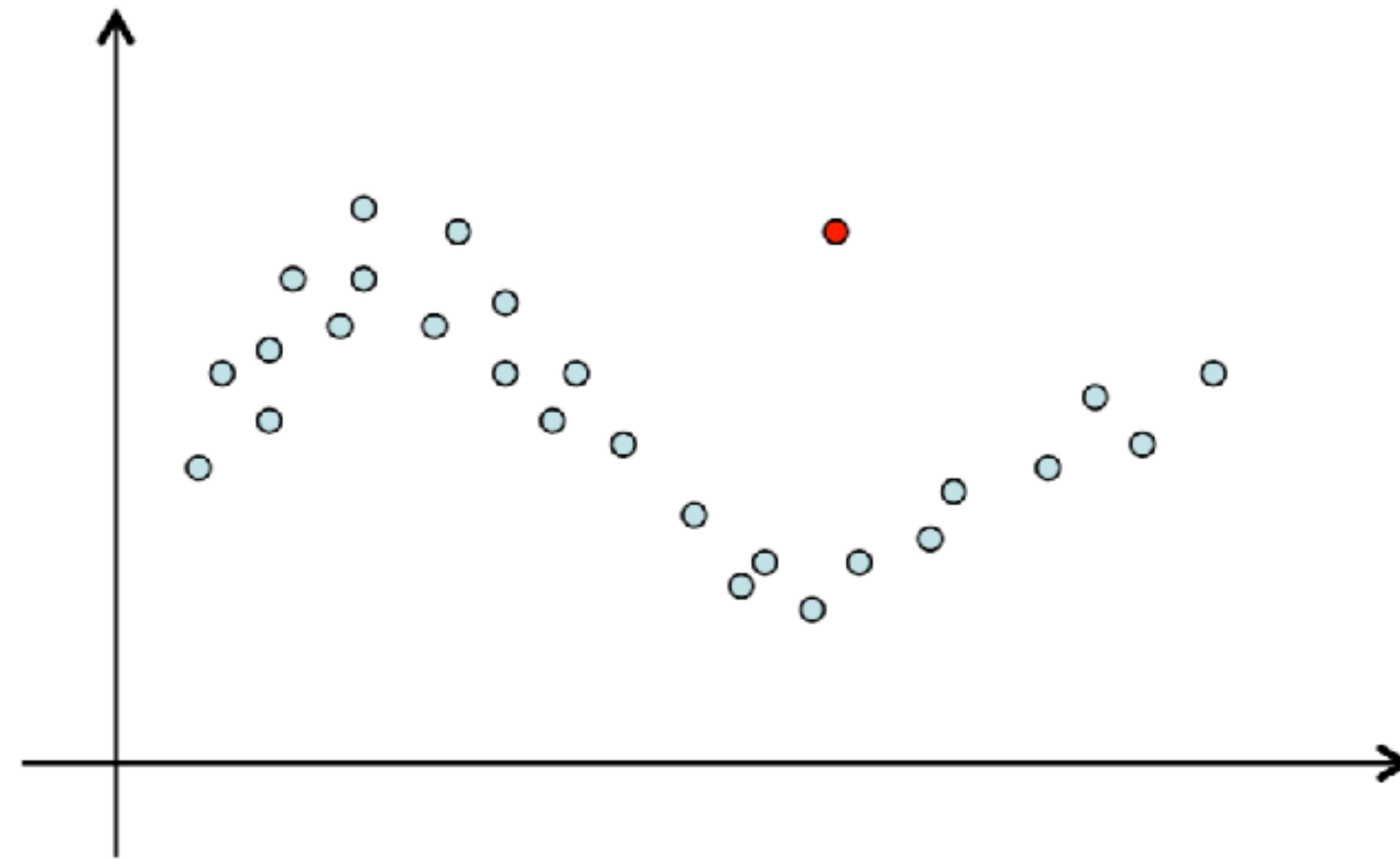# The average US president has been charged with 2 felonies

# The average US president has been charged with 2 felonies

## Total Felony Indictments for Each Current/Former President



**91**

13
4
40
34

Felony Indictments: 0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100

President Number: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46

☐ Campaign Finance (NY)  ☐ Classified Docs (Fed)  ☐ January 6th (Fed)  ☐ Election Fraud/RICO (GA)
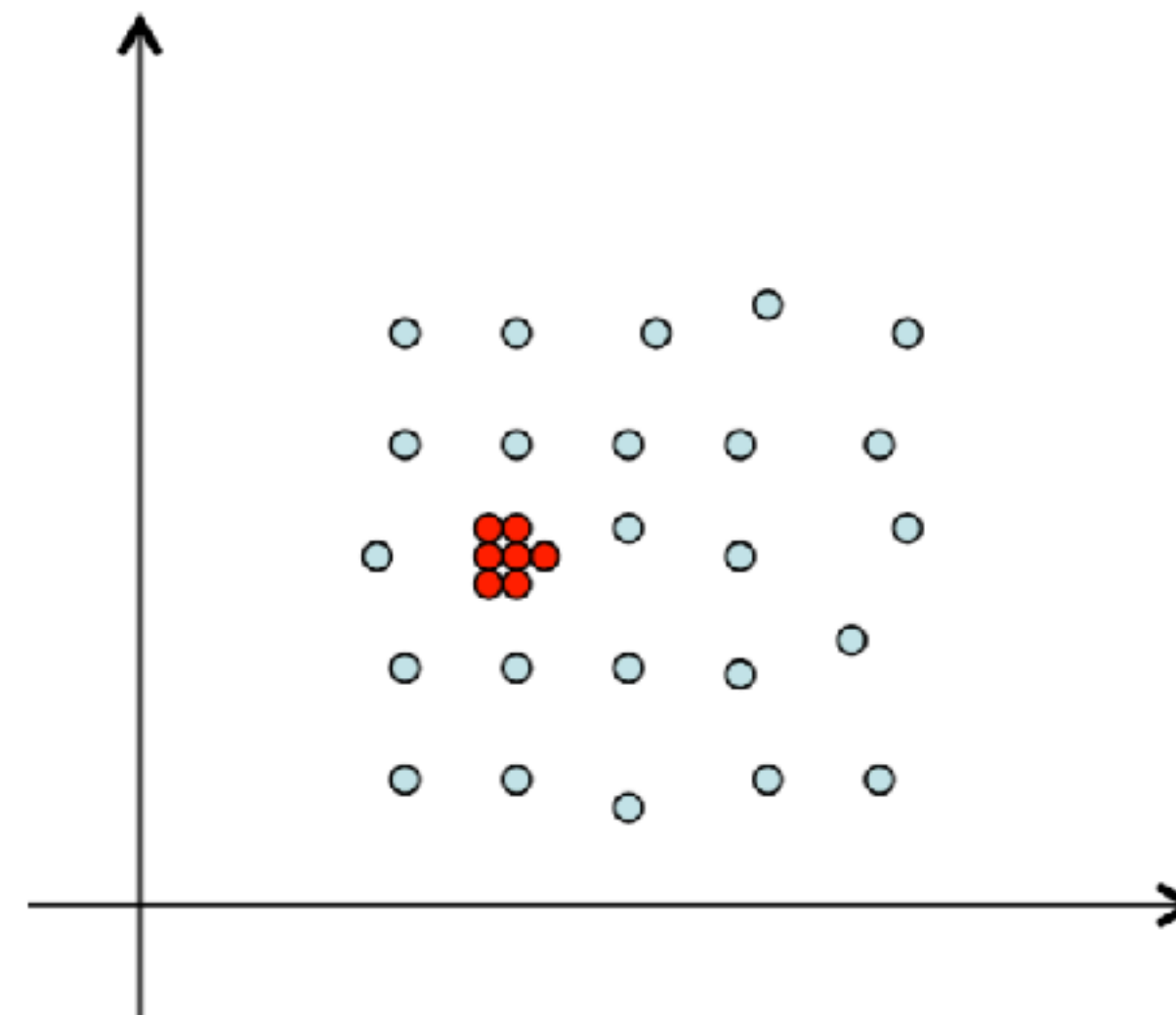
# There are different kinds of outliers



Global outliers

Contextual outliers

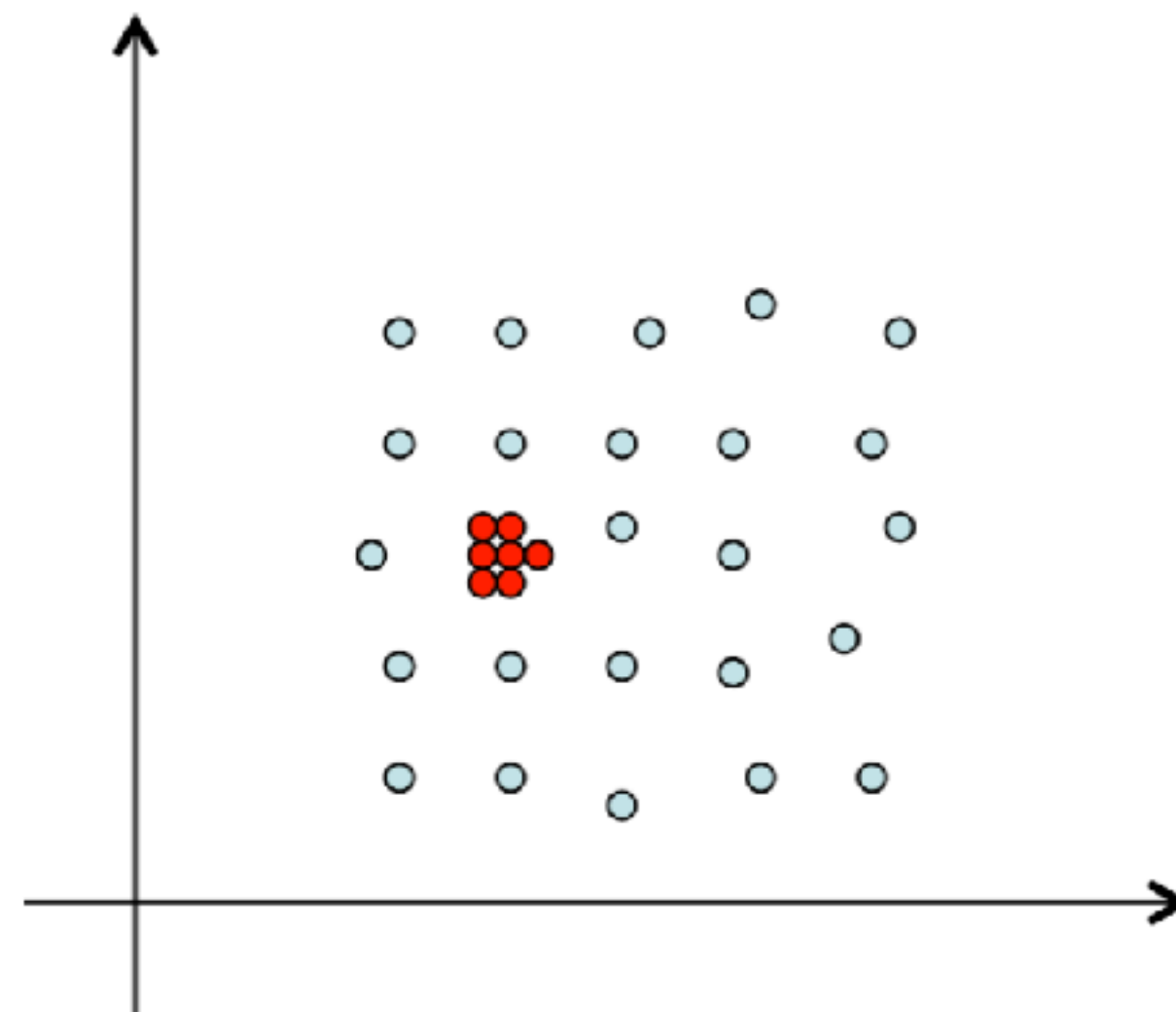Collective outliers

# There are different kinds of outliers


Global outliers


Contextual outliers

deviates significantly with respect to a given context of the object

deviates significantly from the rest of the data. Also called: point anomaly
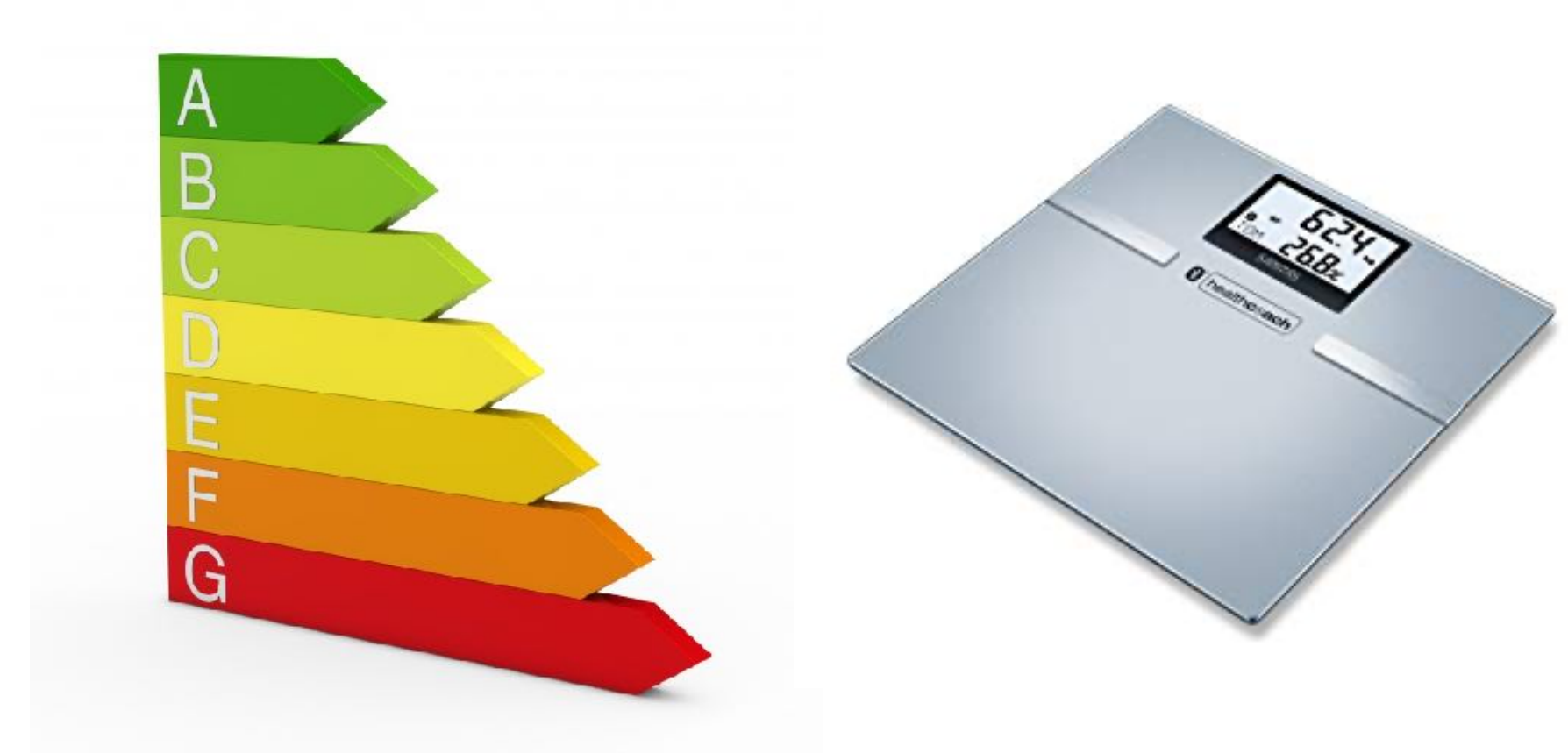

Collective outliers

a subset of data objects that as a group deviate significantly from the typical behavior of the entire data set.

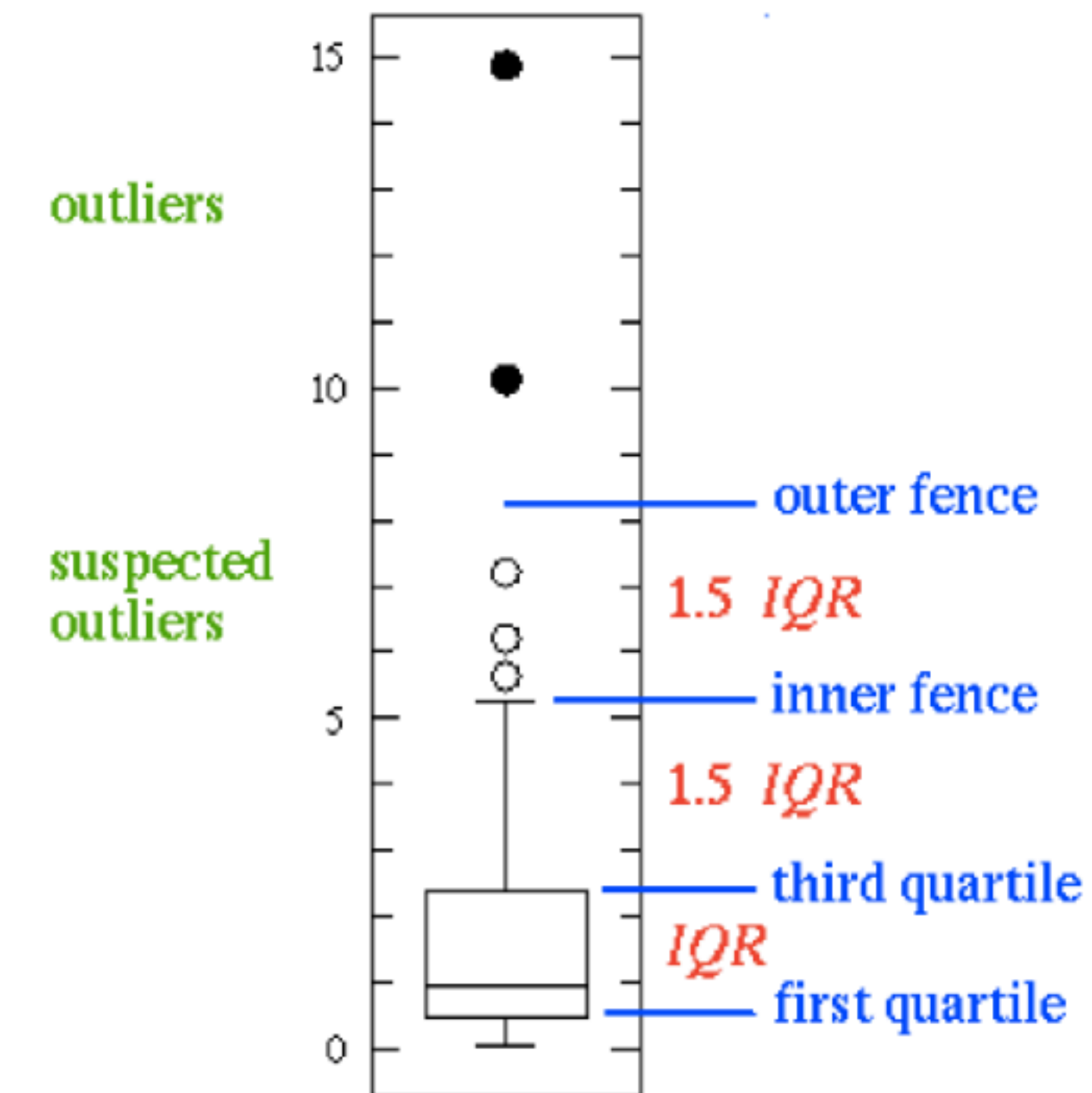An individual object of these collective outliers might not be an outlier itself.

# Jupyter

# Take home message for today
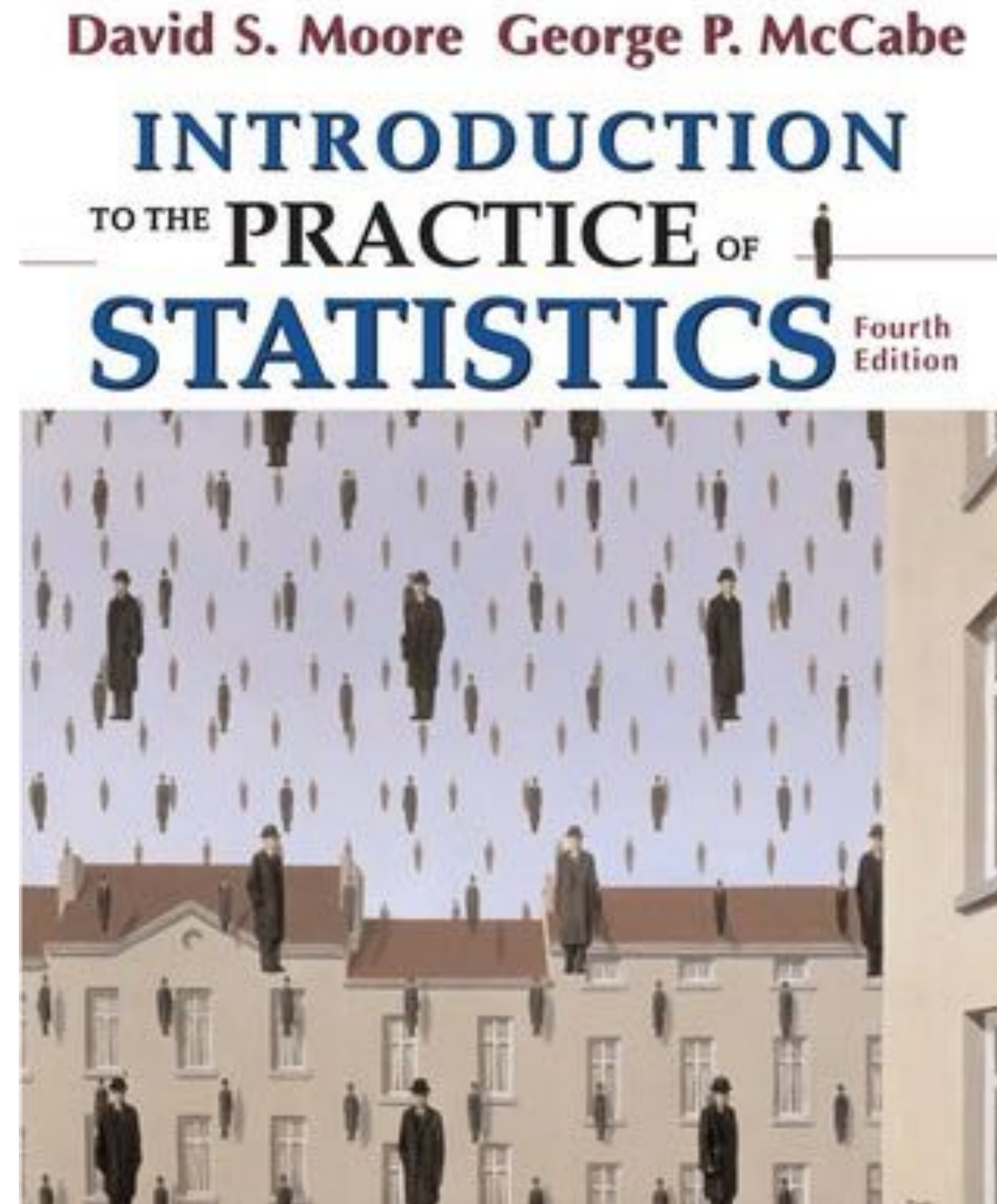
Variables can be categorial or quantitative

The 5 number summary gives a quick quantitative description of a data set

Visualize your data with matplotlib

David S. Moore   George P. McCabe
**INTRODUCTION**
TO THE **PRACTICE** OF
**STATISTICS** Fourth Edition

Chapter 1