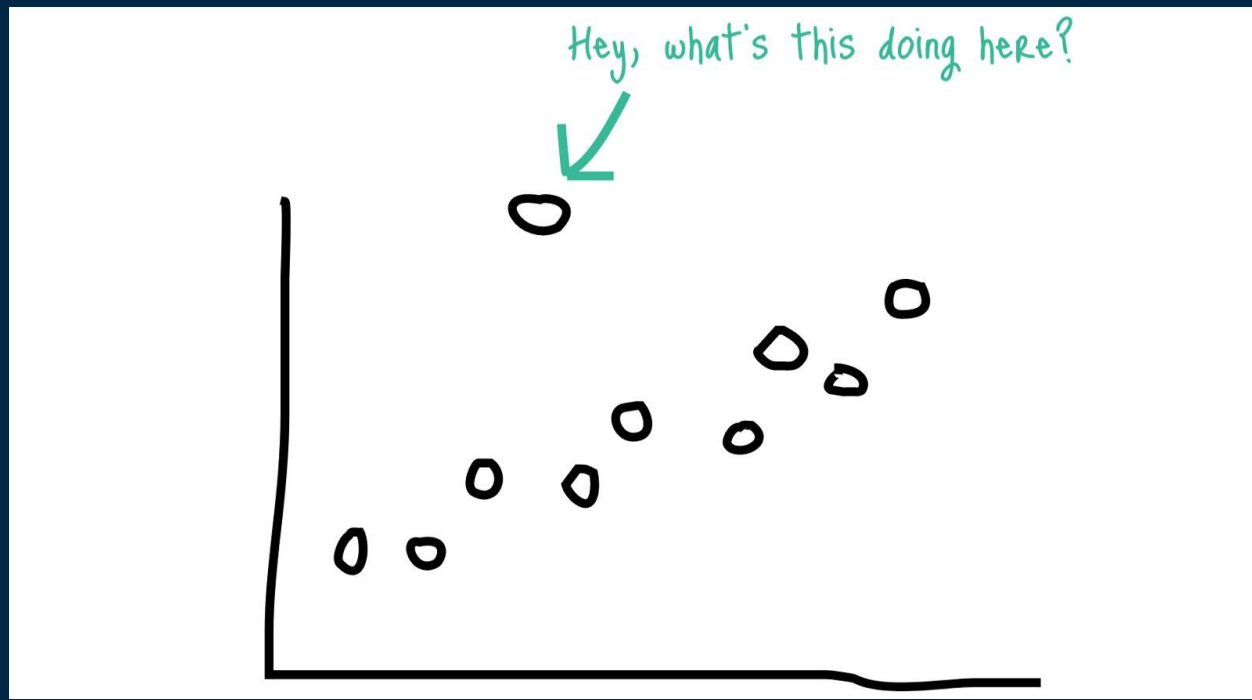# Outlier Detection Analysis

Jiyan Aytek

# If your data is bad, your machine learning tools are useless

# What is an outlier?

# Most Common Causes of Outliers

## 01
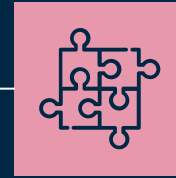### HUMAN ERRORS
DATA ENTRY ERRORS

## 02
### INSTRUMENT ERRORS
MEASUREMENT ERRORS

## 03
### EXPERIMENTAL ERRORS
DATA EXTRACTION OR EXECUTING ERRORS

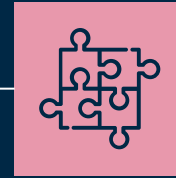# Most Common Causes of Outliers

**04**

DATA PROCESSING ERRORS

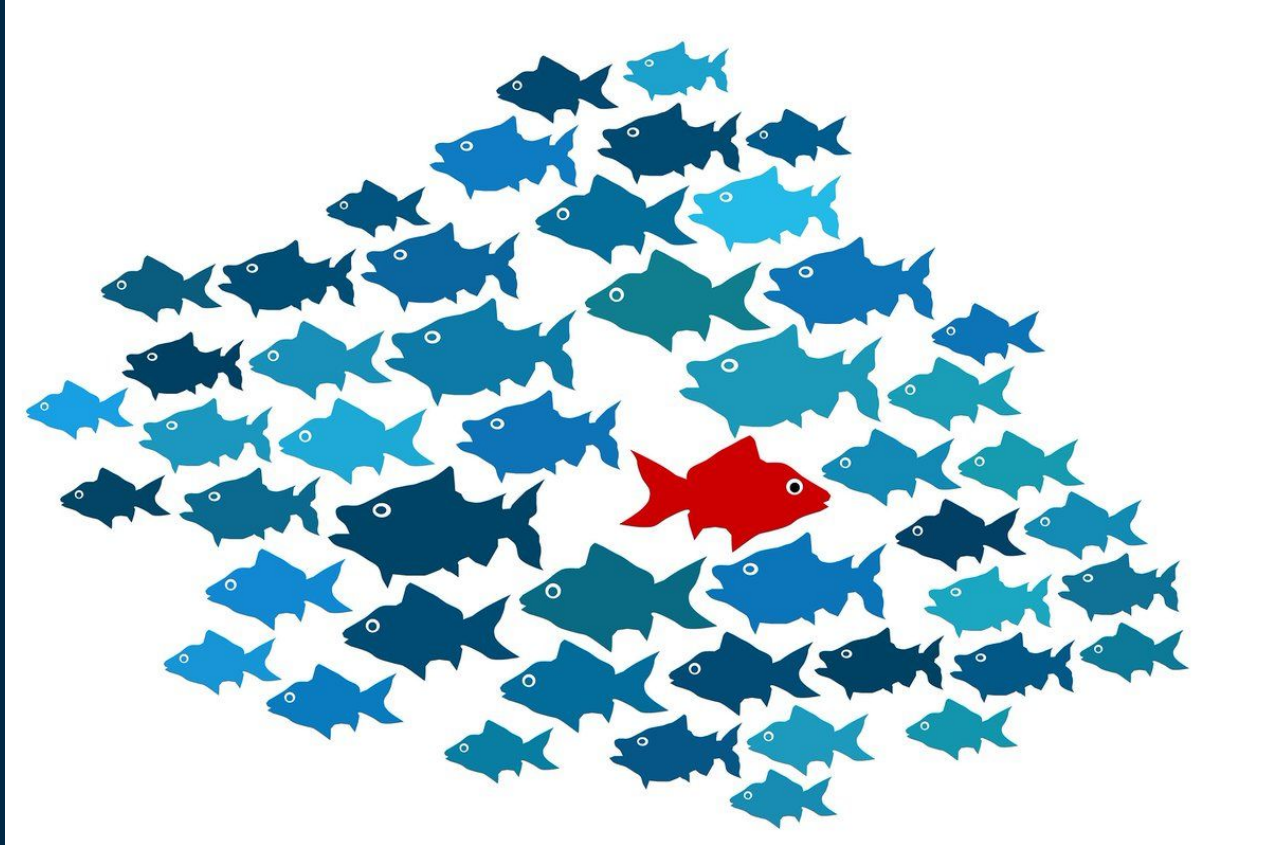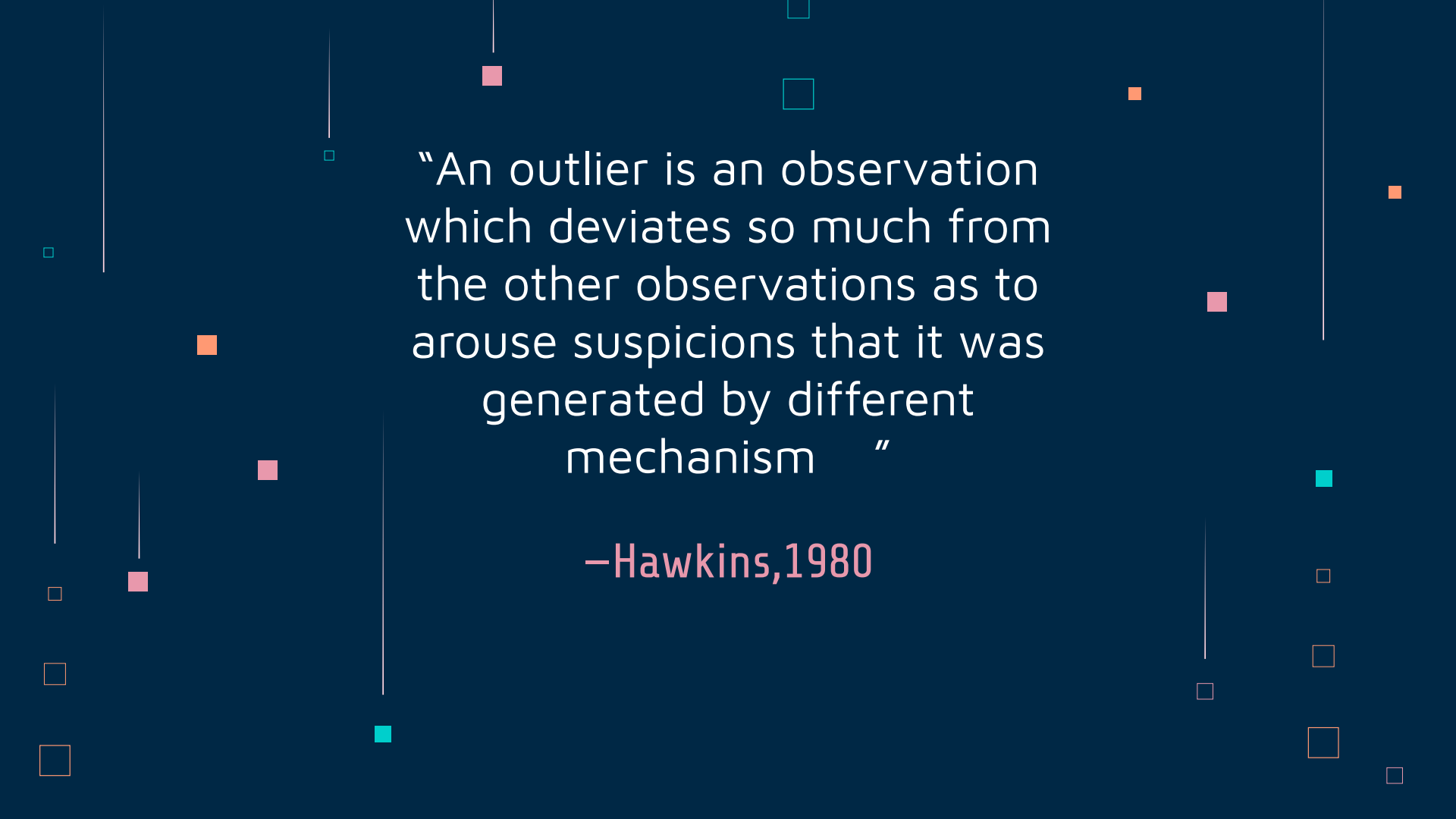DATA MANIPULATION

**05**

SAMPLING ERRORS

VARIOUS SOURCES

**06**

NATURAL

NOT AN ERROR

# Outlier Example

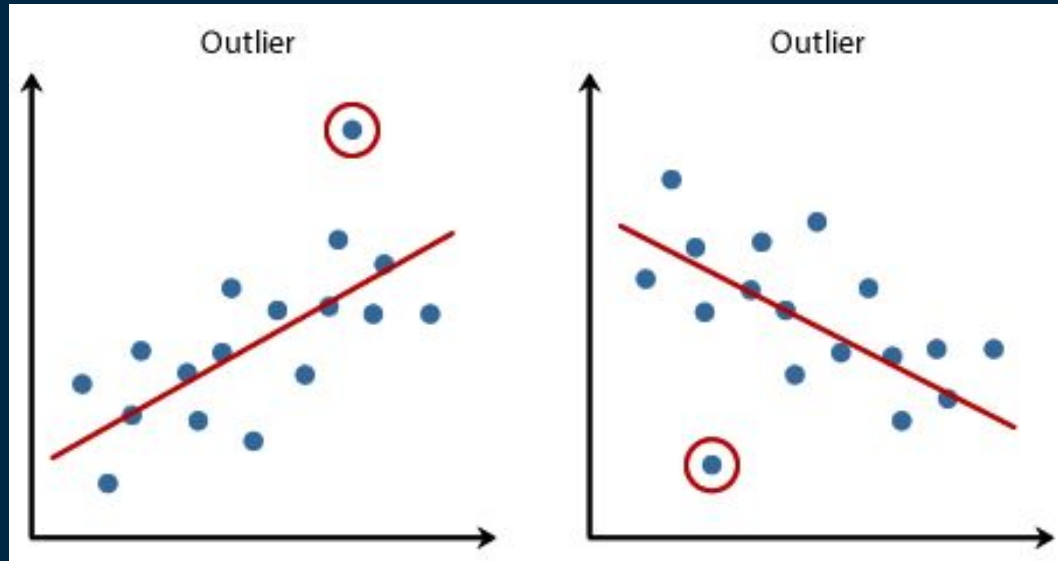"An outlier is an observation which deviates so much from the other observations as to arouse suspicions that it was generated by different mechanism "

—Hawkins,1980

# Clustering

# Regression

# Invalid or Outlier Data

| Name       | Gender | Year of Birth |
|------------|--------|---------------|
| Andrew NG  | M      | 1976          |
| Sarah Tan  | F      | 2976          |

# Detection Methods :
# Uni-variate Methods

# Boxplot



$(x > (Q3 + \textbf{1.5*IQR})) \lor (x < (Q1 - \textbf{1.5*IQR})) \rightarrow (x \text{ is an } \textbf{outlier})$

$(x > (Q3 + \textbf{3*IQR})) \lor (x < (Q1 - \textbf{3*IQR})) \rightarrow (x \text{ is an } \textbf{extreme-value})$
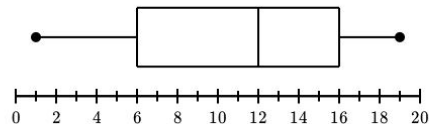
# Boxplot

14, 6, 3, 2, 4, 15, 11, 8, 1, 7, 2, 1, 3, 4, 10, 22, 20

He wants to create a graph that helps him understand the spread of distances (and the median distance) that people travel. What kind of a graph should he create?

2.5        median   12.5

1, 1, 2, 2, 3, 3, 4, 4, 6, 7, 8, 10, 11, 14, 15, 20, 22

0  2.5  56   10  12.5 15   20  22   25   30   35

---

Which data set could be represented by the box plot shown below?



0  2  4  6  8  10  12  14  16  18  20

Choose 1 answer:

(A)  1, 3, 6, 8, 10, 12, 13, 13, 16, 18, 20

(B)  1, 3, 6, 8, 10, 12, 13, 13, 16, 18, 19

(C)  1, 3, 6, 8, 10, 11, 13, 13, 18, 18, 19

(D)  1, 3, 6, 8, 10, 11, 13, 13, 16, 18, 19

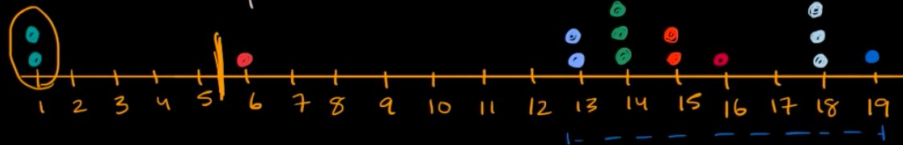# Identifying outliers in a dataset (with Boxplot)

# Boxplot

```python
import numpy as np

def outliers_iqr(ys):
    quartile_1, quartile_3 = np.percentile(ys, [25, 75])
    iqr = quartile_3 - quartile_1
    lower_bound = quartile_1 - (iqr * 1.5)
    upper_bound = quartile_3 + (iqr * 1.5)
    return np.where((ys > upper_bound) | (ys < lower_bound))
```
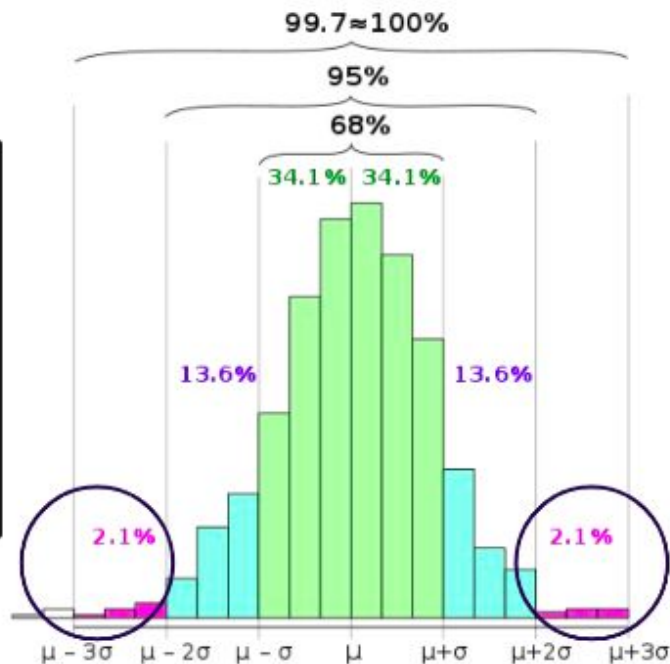
# Standard Deviation

# Z Score



A Normal Distribution

The Standard Normal Distribution

950  970  990  1010  1030  1050  1070

−3  −2  −1  0  +1  +2  +3

$$z = \frac{X - \mu}{\sigma}$$

### Z-scores in the Normal Distribution

0.001350

0.001350

Z-scores

Density

# Hard Edges Method

Data yielding outside of the ($1th$ - $99th$) quantile/percentile interval will be evaluated as *outlier*.

# Why use Hard Edges Method?

- No calculate std, mean, median
- Basic and quick
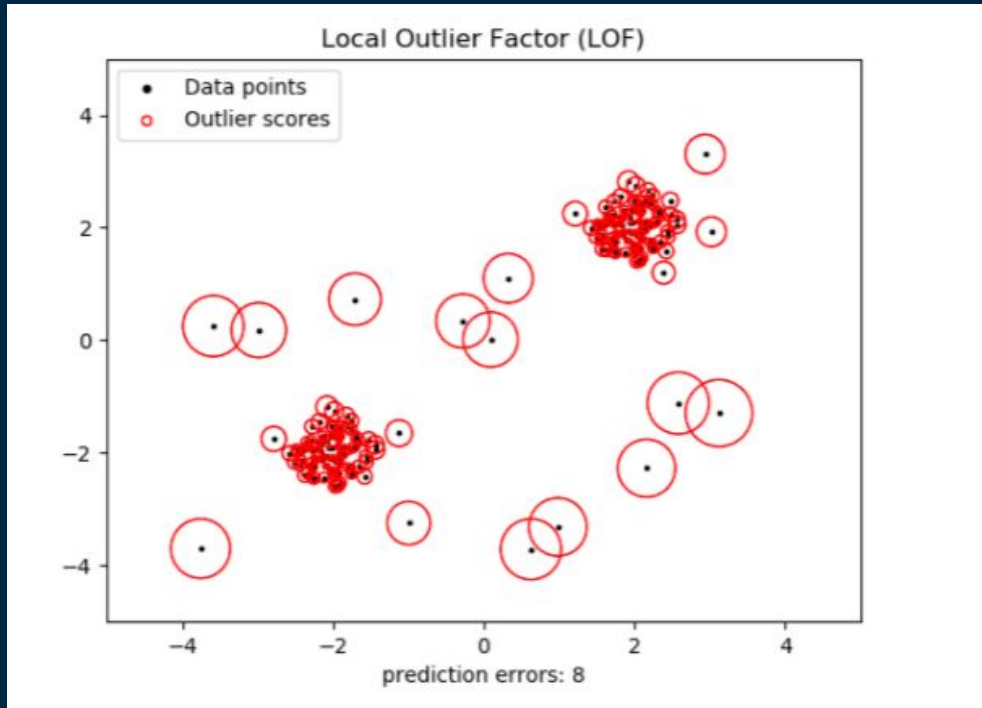- Appropriate for big dataset ( for example : 300.000 rows)

# Detection Methods : Multi-variate Methods

# LOCAL OUTLIER FACTOR

THANK YOU