

---

# MODULE 3: DATA ANALYSIS AND VISUAL STORYTELLING

CT ACADEMY | DATA ACTION LAB

---

# 8. DATA ANALYSIS

DATA ANALYSIS AND VISUAL STORYTELLING

# CONTINGENCY/PIVOT TABLES

**Contingency table:** examines the relationship between two categorical variables via their relative (cross-tabulation).

**Pivot table:** a table generated by applying operations (sum, count, mean, etc.) to variables, possibly based on another (categorical) variable.

Contingency tables are special cases of pivot tables.

	Large	Medium	Small
Window	1	32	31
Door	14	11	0

Type	Count	Signal avg	Signal stdev
Blue	4	4.04	0.98
Green	1	4.93	N.A.
Orange	4	5.37	1.60

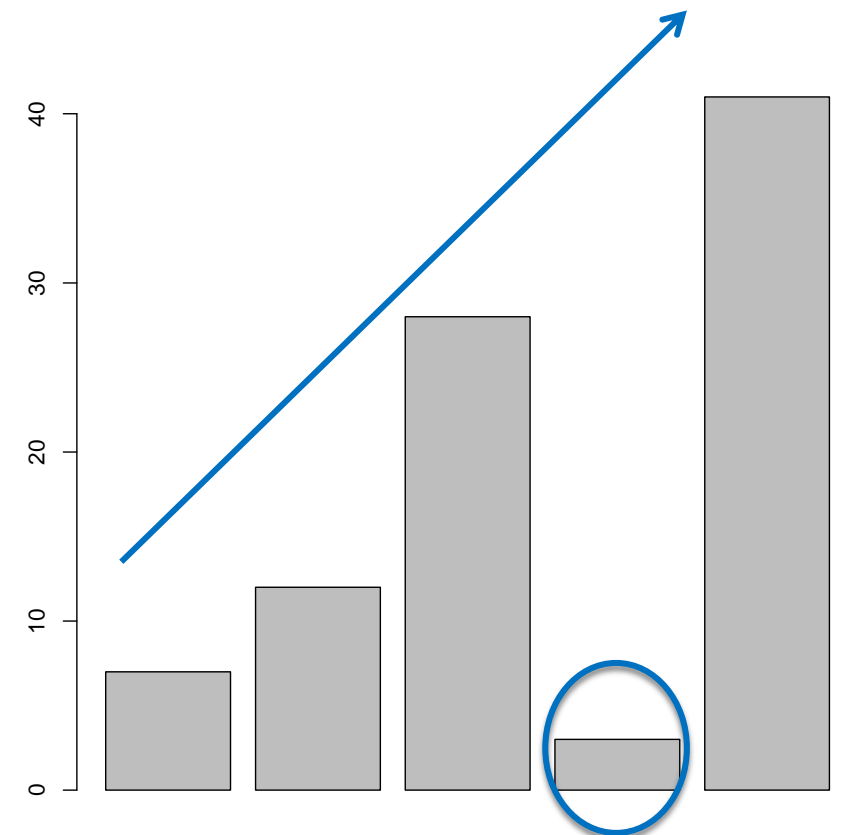
# ANALYSIS THROUGH VISUALIZATION

## Analysis (broad definition):

- identifying patterns or structure
- adding meaning to these patterns or structure by interpreting them in the context of the system.

**Option 1:** use analytical methods to achieve this.

**Option 2:** visualize the data and use the brain's analytic power (perceptual) to reach meaningful conclusions about these patterns.



# NUMERICAL SUMMARIES

In a first pass, a variable can be described along 2 dimensions: **centrality** & **spread** (skew and kurtosis are also used).

**Centrality measures** include:

- median, mean, mode

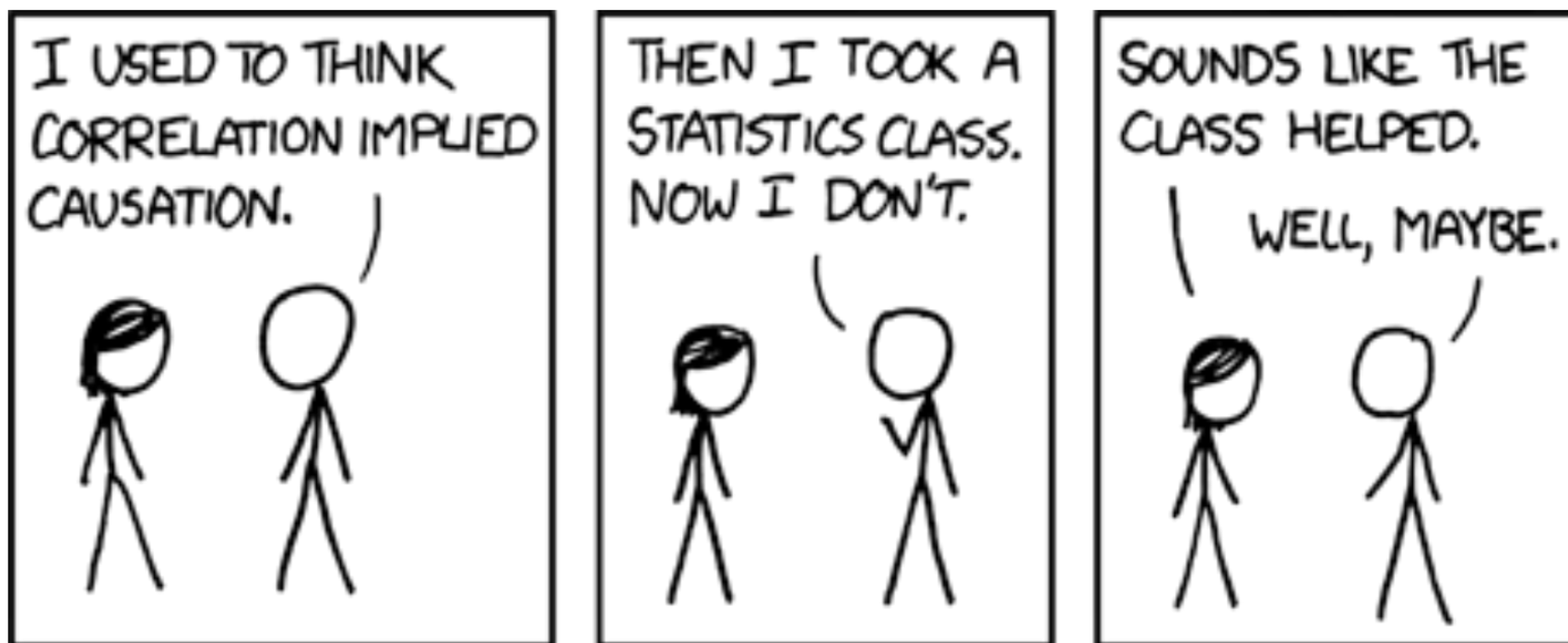
**Spread (or dispersion) measures** include:

- standard deviation (sd), variance, quartiles, range, etc.

The median, range and the quartiles are easily calculated from **ordered lists**.



# CORRELATION



Correlation doesn't imply causation, but it does waggle its eyebrows suggestively and gesture furtively while mouthing 'look over there'.

# LINEAR REGRESSION

The basic assumption of **linear regression** is that the dependent variable  $y$  can be approximated by a linear combination of the independent variables:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

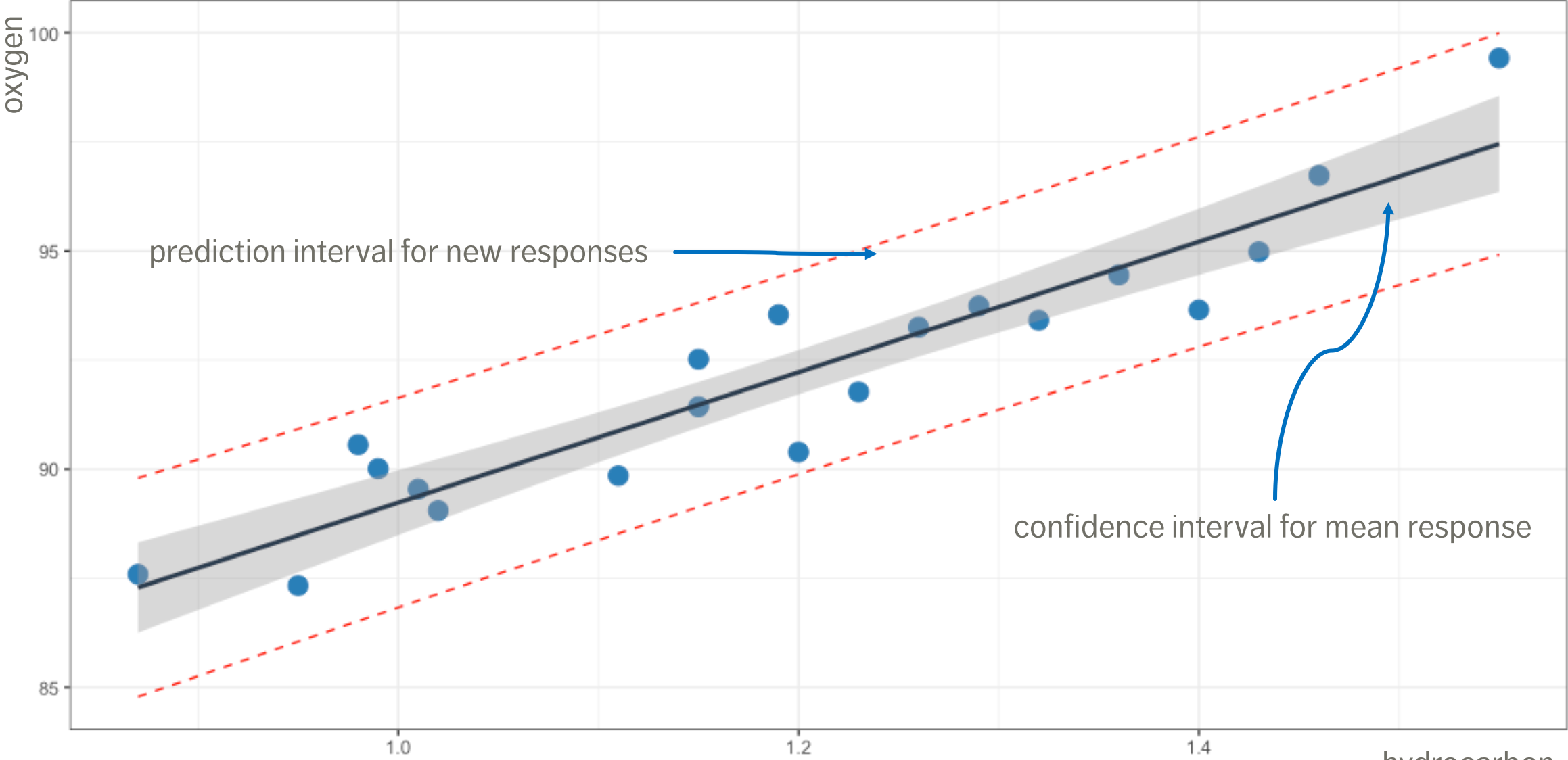
where  $\boldsymbol{\beta} \in \mathbb{R}^p$  is to be determined based on the **training set**, and for which

$$E(\boldsymbol{\varepsilon}|\mathbf{X}) = \mathbf{0}, \quad E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T|\mathbf{X}) = \sigma^2\mathbf{I}.$$

Typically, the errors are also assumed to be **normally distributed**:

$$\boldsymbol{\varepsilon}|\mathbf{X} \sim N(\mathbf{0}, \sigma^2\mathbf{I}).$$

$$\text{oxygen} = 14.95 \times \text{hydrocarbon} + 74.28$$





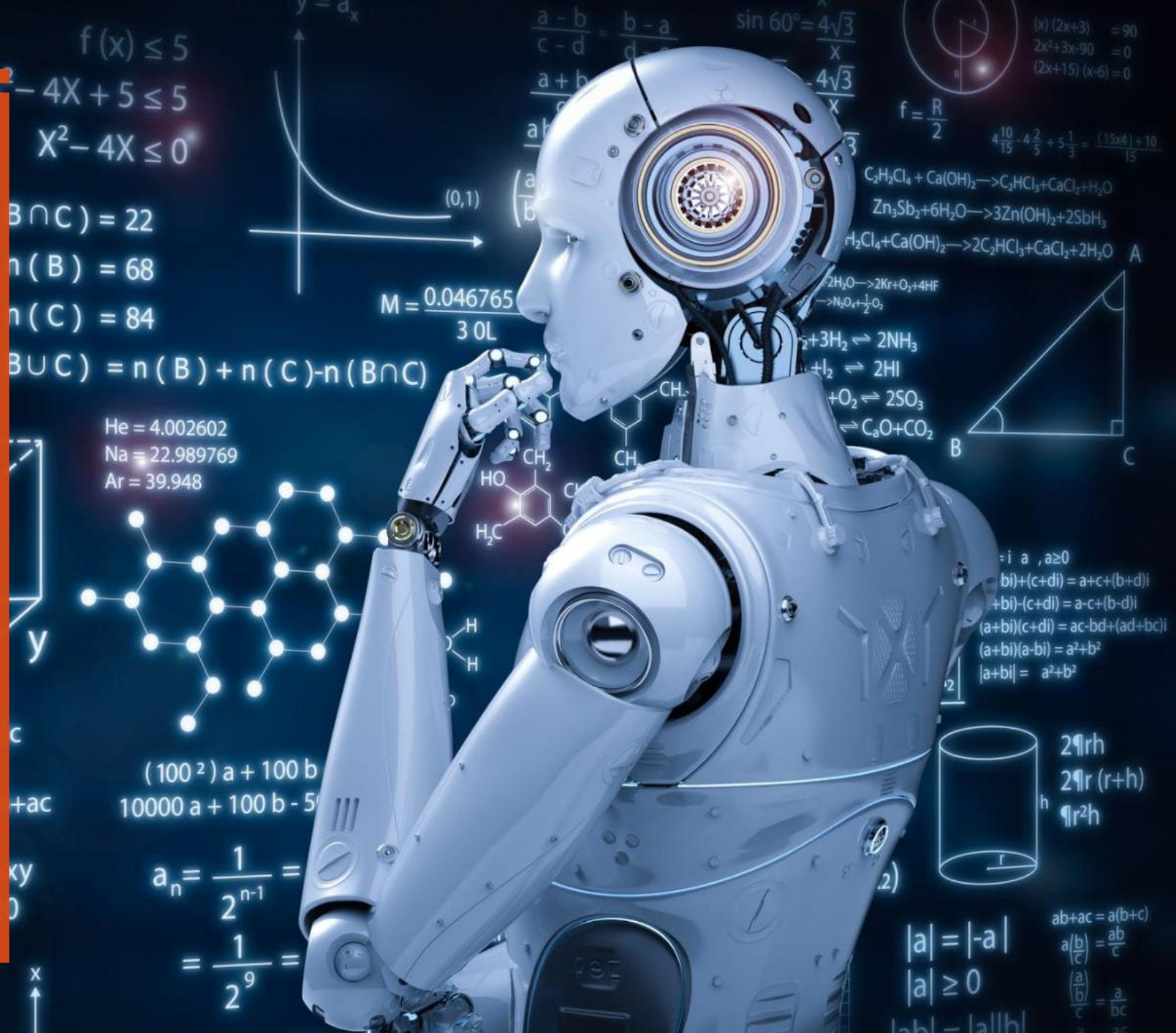
# MACHINE LEARNING TASKS

**Classification, class probability estimation:** which clients are likely to be repeat customers?

**Clustering:** do customers form natural groups?

**Association rule discovery:** what books are commonly purchased together?

Others: **value estimation** (how much is a client likely to spend in a restaurant); **profiling and behaviour description**; **link prediction**; **data reduction**; **influence/causal modeling**; **similarity matching** (which prospective clients are similar to a company's best clients?), etc.



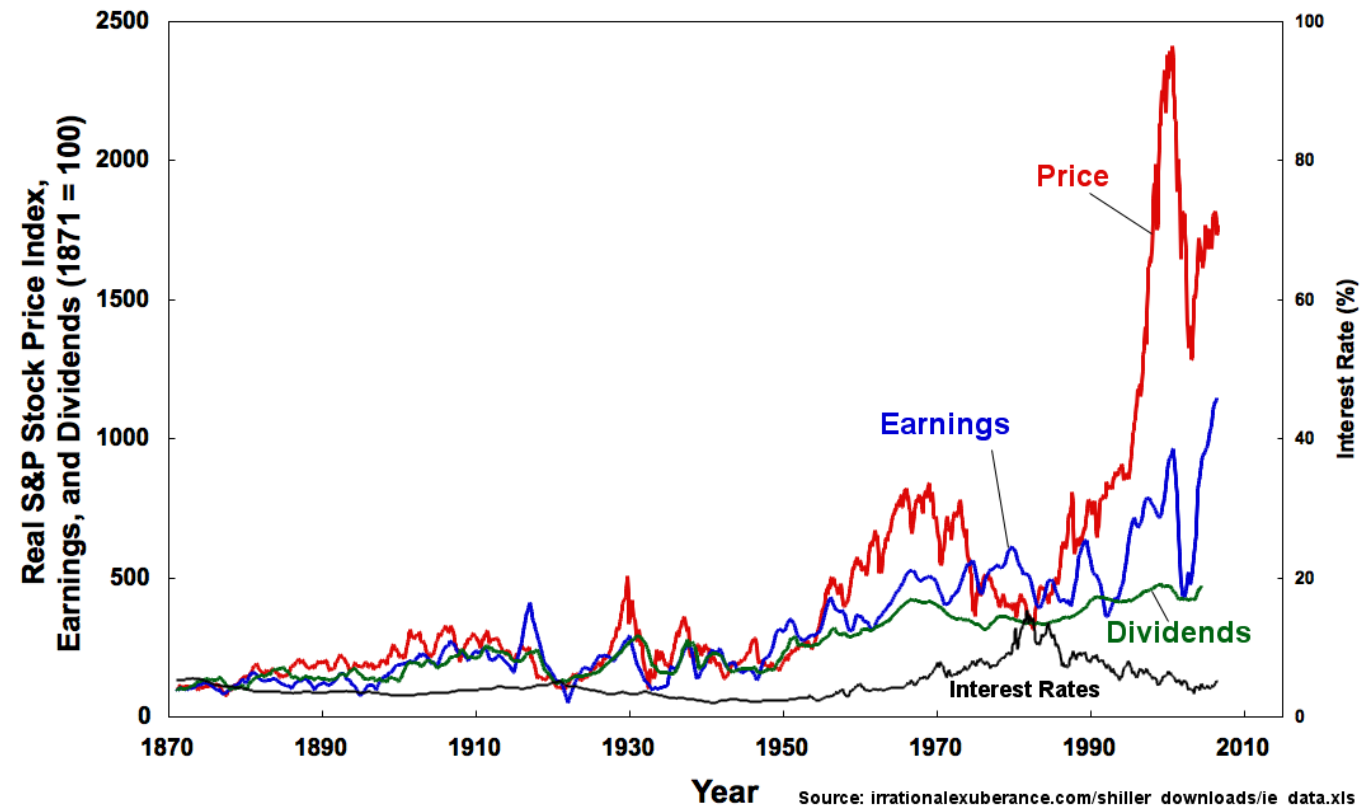
# TIME SERIES ANALYSIS

A simple **time series**:

- has two variables: time + 2<sup>nd</sup> variable
- the second variable is *sequential*

What is the **pattern of behaviour** of this second variable over time?  
Relative to other variables?

Can we use this to **forecast the future behaviour** of the variable?



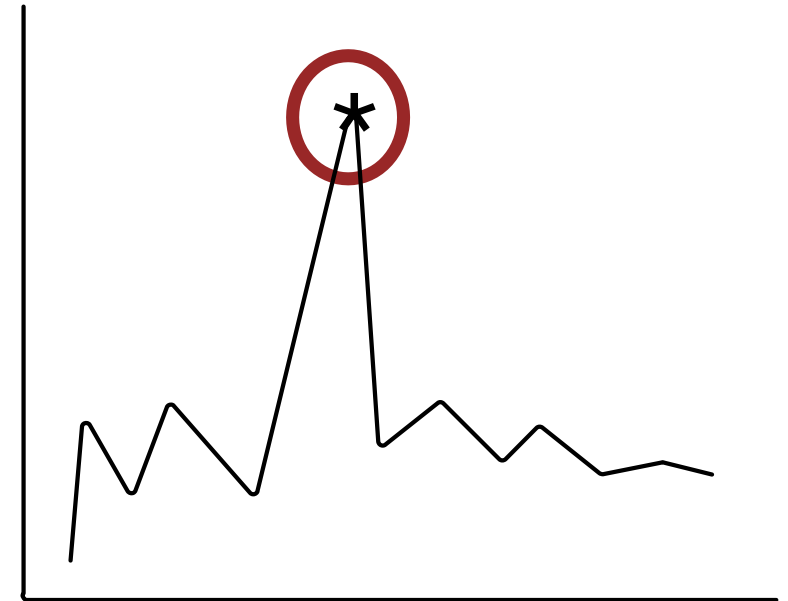
# ANOMALY DETECTION

**Anomaly:** an unexpected, unusual, atypical or statistically unlikely event

Wouldn't it be nice to have a data analysis pipeline that alerted you when things were out of the ordinary?

Many different analytic approaches to take!

- clustering
- classification
- ensemble techniques, etc.



# EXERCISES

1. Are there opportunities for computations like correlation, linear regression, and times series analysis in your workplace datasets?
2. Are there opportunities for machine learning tasks or anomaly detection in your workplace datasets?

# THE HOT MESS

“Data is messy, you know.”  
“Even after it’s been cleaned?”  
“*Especially* after it’s been cleaned.”

Data **cleaning, processing, wrangling** are essential aspects of data science projects; analysts may spend **up to 80%** of their time on **data preparation**.

# DATA WRANGLING AND TIDY DATA

**Tidy data** has a specific structure:

- each variable is in a single column
- each observation is in a single row
- each type of observational unit is in a single table

Country	2011	2012	2013
FR	7000	6900	7000
DE	5800	6000	6200
US	15000	14000	13000

VS.

Country	Year	n
FR	2011	7000
DE	2011	5800
US	2011	15000
FR	2012	6900
DE	2012	6000
US	2012	14000
FR	2013	7000
DE	2013	6200
US	2013	13000

# DATA WRANGLING FUNCTIONALITY

Data wrangling functions should allow the analyst to:

- extract a subset of variables from the data frame
- extract a subset of observations from the data frame
- sort the data frame along any combination of variables in increasing or decreasing order
- to create new variables from existing variables
- to create (so-called) pivot tables, by observation groups
- database functionality (joins, etc.)
- etc.

# EXERCISE

Turn the dataset found in the file [cities.txt](#) into a tidy dataset.



# APPROACHES TO DATA CLEANING

There are 2 **philosophical** approaches to data cleaning and validation:

- methodical
- narrative

The **methodical** approach consists of running through a **check list** of potential issues and flagging those that apply to the data.

The **narrative** approach consists of **exploring** the dataset and trying to spot unlikely and irregular patterns.



# APPROACHES TO DATA CLEANING

## Methodical (syntax)

- Pros: checklist is **context-independent**; pipelines **easy to implement**; common errors and invalid observations **easily identified**
- Cons: may prove **time-consuming**; cannot identify new types of errors

## Narrative (semantics)

- Pros: process may simultaneously yield **data understanding**; false starts are (at most) as costly as switching to mechanical approach
- Cons: may miss important sources of errors and invalid observations for datasets with **high number of features**; domain knowledge may bias the process by neglecting uninteresting areas of the dataset

# DATA SOUNDNESS

The ideal dataset will have as few issues as possible with:

- **validity:** data type, range, mandatory response, uniqueness, value, regular expressions
- **completeness:** missing observations
- **accuracy and precision:** related to measurement and data entry errors; target diagrams (accuracy as bias, precision as standard error)
- **consistency:** conflicting observations
- **uniformity:** are units used uniformly throughout?

Checking for data quality issues at an early stage can save headaches later in the analysis.

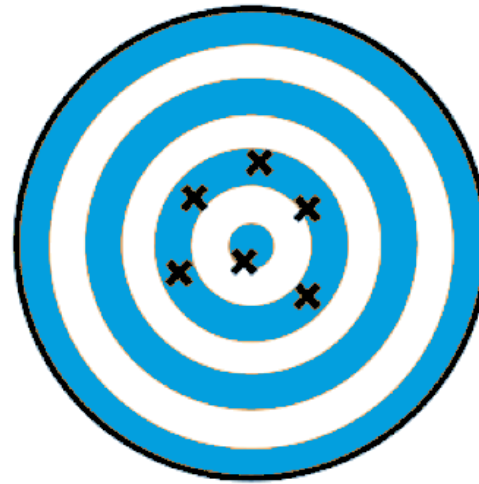
# DATA SOUNDNESS



accurate and  
precise



precise but  
not accurate



accurate but  
not precise

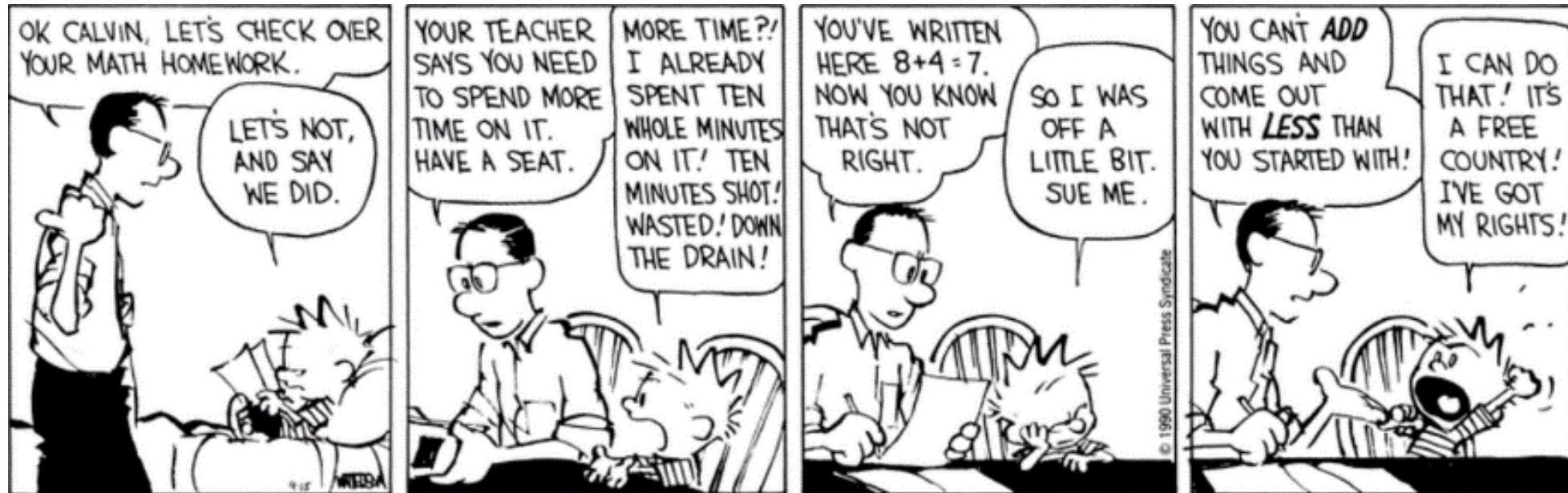


neither accurate  
nor very precise

# COMMON ERROR SOURCES

When dealing with **legacy**, **inherited** or **combined** datasets (i.e., datasets over which there is no collection and initial processing control):

- missing data given a code
- 'NA'/'blank' given a code
- data entry error
- coding error
- measurement error
- duplicate entries
- heaping



# DETECTING INVALID ENTRIES

Potentially invalid entries can be detected with the help of:

- **univariate descriptive statistics**  
count, range, z-score, mean, median, standard deviation, logic check
- **multivariate descriptive statistics**  
*n*-way table, logic check
- **data visualization**  
scatterplot, scatterplot matrix, histogram, joint histogram, etc.

## DETECTING INVALID ENTRIES

Univariate tests do not always tell the **whole** story.

This step might allow for the identification of potential outliers.

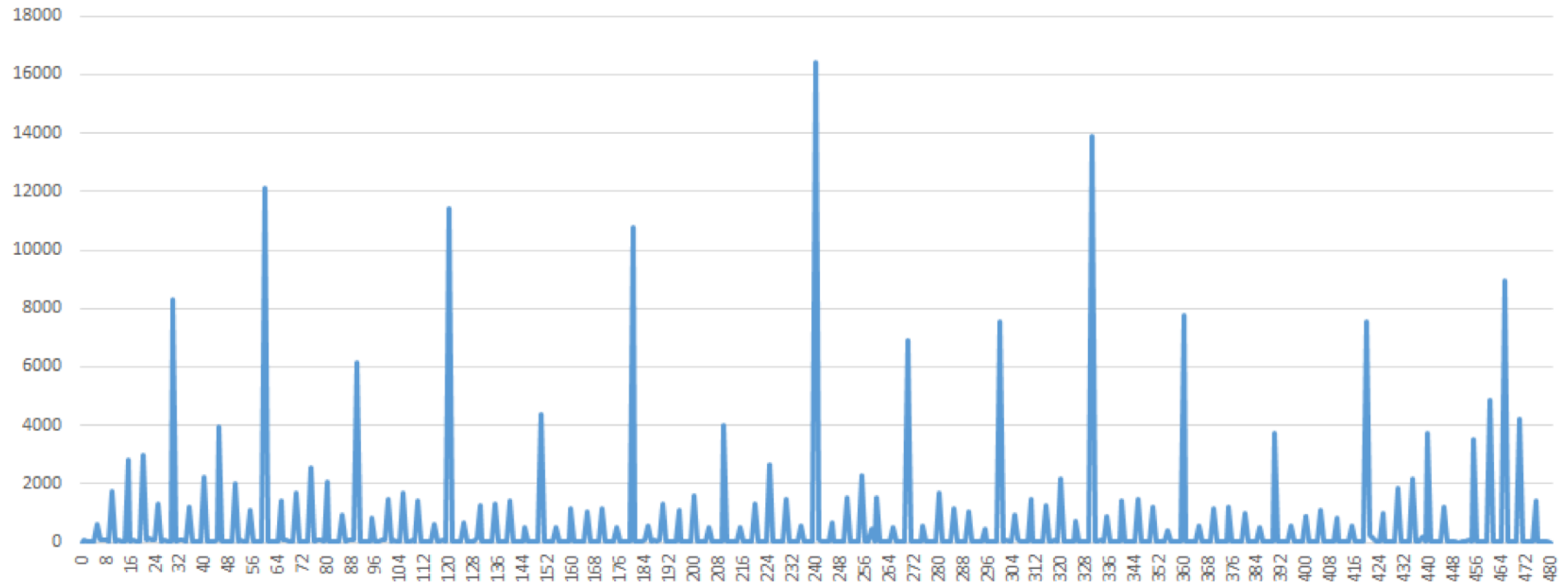
Failure to detect invalid entries  $\neq$  all entries are valid.

Small numbers of invalid entries recoded as “missing.”

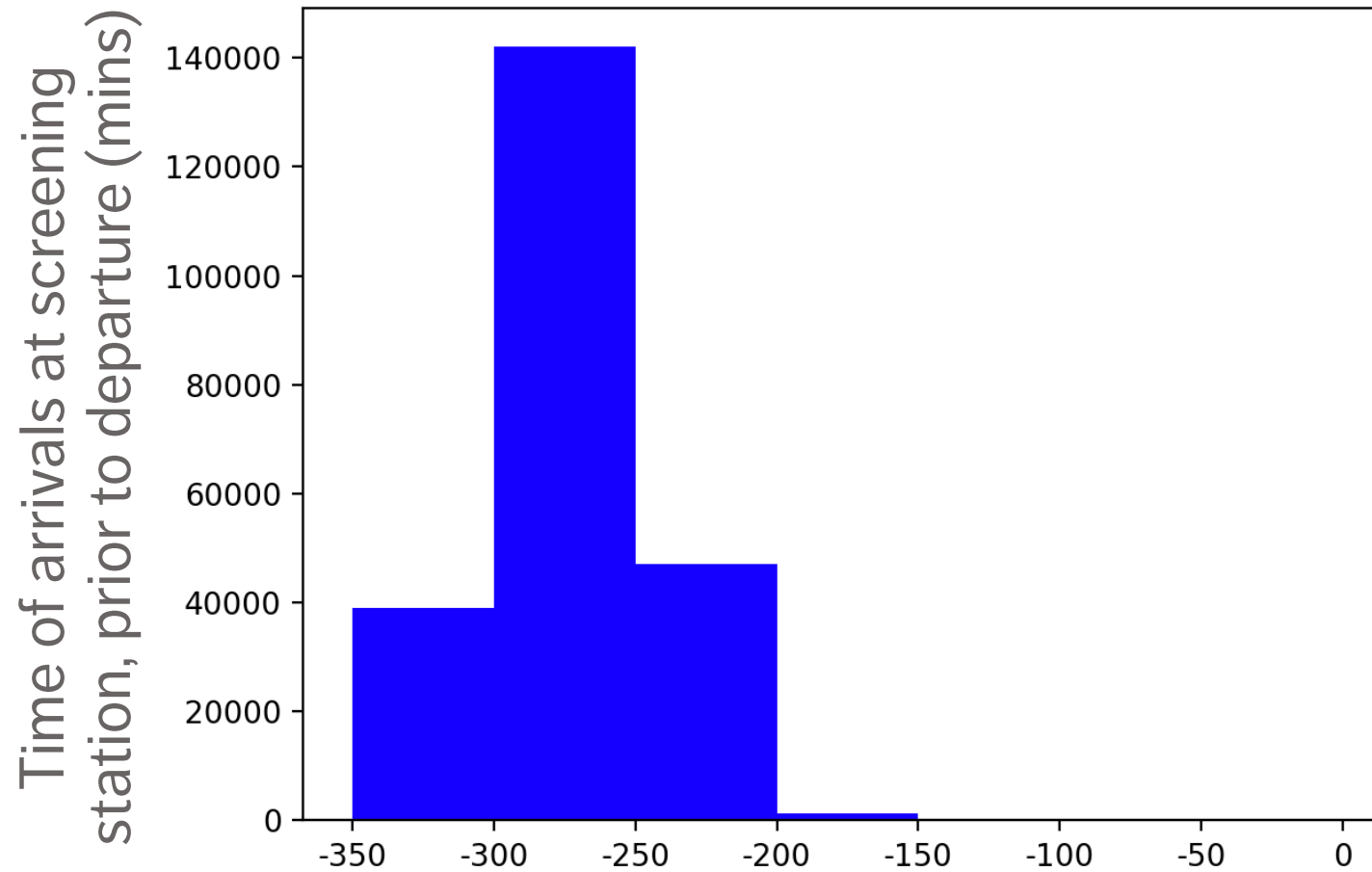




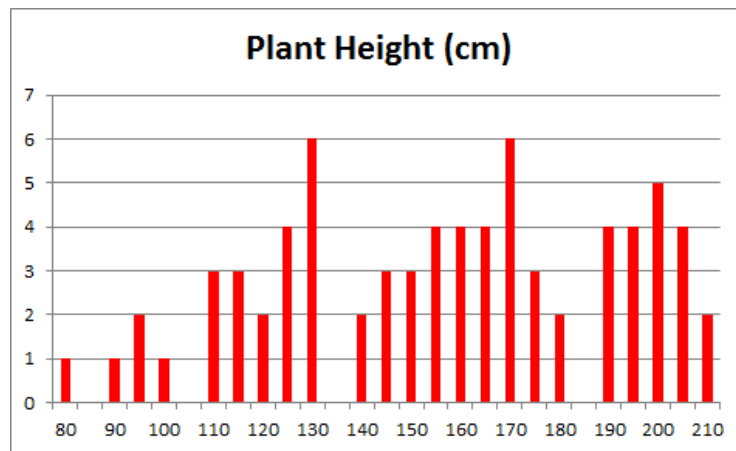
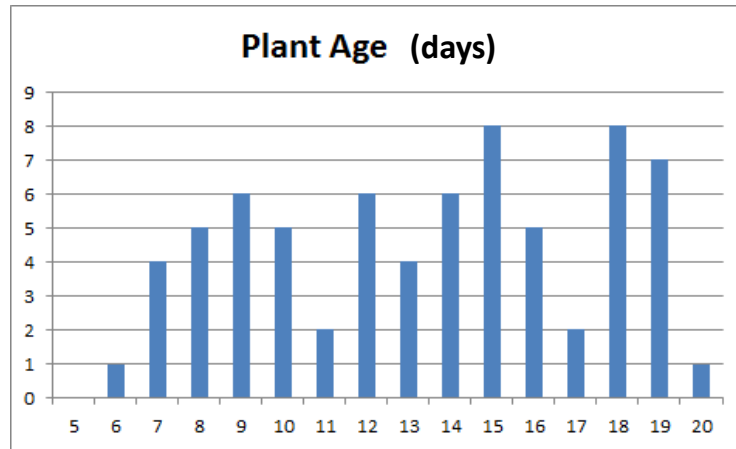
# DETECTING INVALID ENTRIES



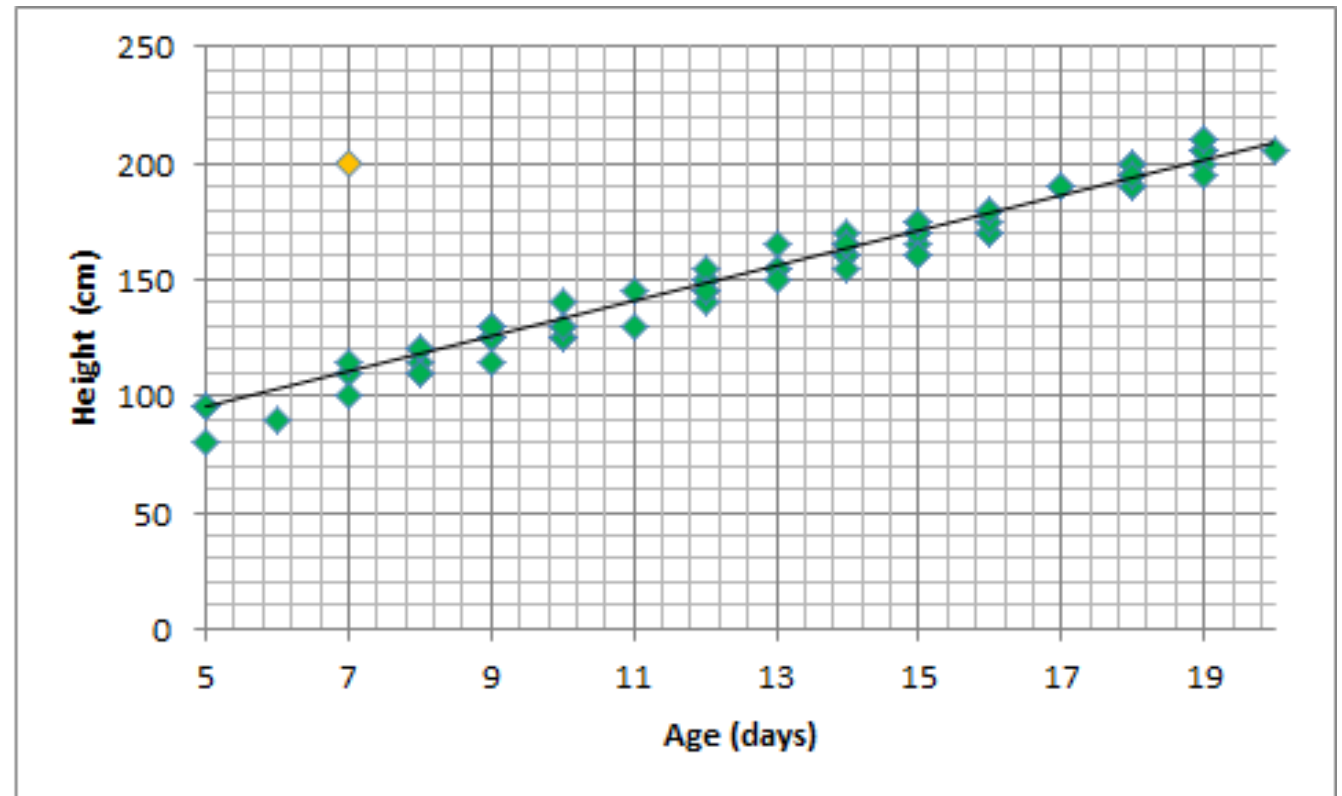
# DETECTING INVALID ENTRIES



# DETECTING INVALID ENTRIES

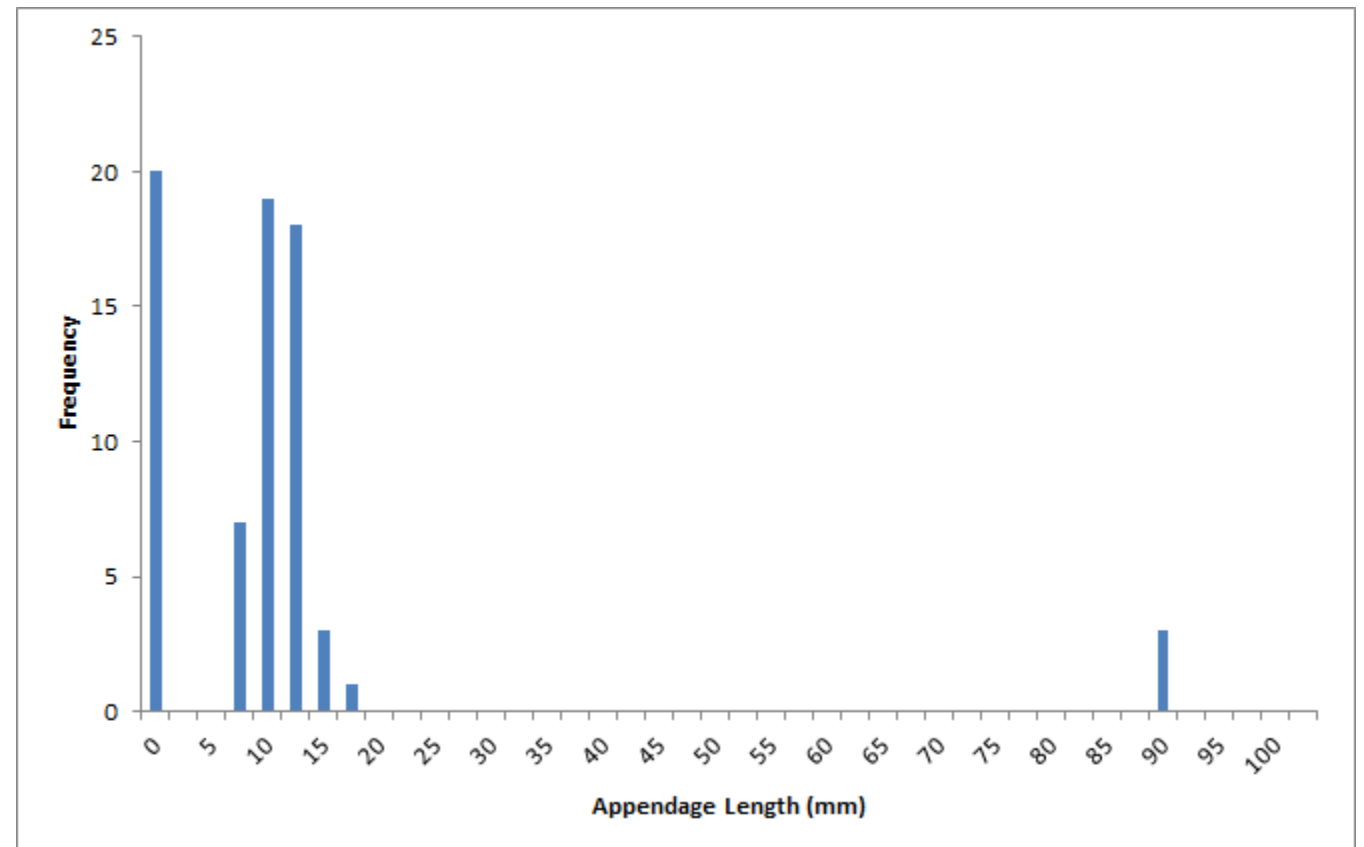


VS.



# DETECTING INVALID ENTRIES

<i>Appendage length (mm)</i>	
Mean	10.35
Standard Deviation	16.98
Kurtosis	16.78
Skewness	4.07
Minimum	0
First Quartile	0
Median	8.77
Third Quartile	10.58
Maximum	88
Range	88
Interquartile Range	10.58
Mode	0
Count	71



# EXERCISES

1. Does the dataset found in the file [cities.txt](#) appear to be of good quality (is it sound? does it have invalid entries?)
2. Create a list of items that could be used in a methodical data cleaning checklist. Use data that you have encountered in the past as inspiration (numerical, categorical, text data).

# TYPES OF MISSING OBSERVATIONS

Blank fields come in 4 flavours:

- **nonresponse**  
an observation was expected but none had been entered
- **data entry issue**  
an observation was recorded but was not entered in the dataset
- **invalid entry**  
an observation was recorded but was considered invalid and has been removed
- **expected blank**  
a field has been left blank, but expectedly so

# TYPES OF MISSING OBSERVATIONS

Too many missing values (of the first three type) can be indicative of **issues with the data collection process** (more on this later).

Too many missing values (of the fourth type) can be indicative of **poor questionnaire design**.

Finding missing values can help you deal with other data science problems.

---

# THE CASE FOR IMPUTATION

---

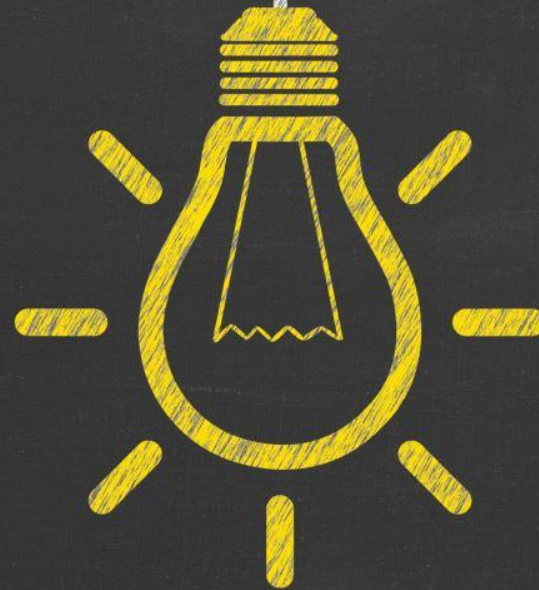
Not all analytical methods can easily accommodate missing observations:

- **discard** the missing observation
  - not recommended, unless the data is MCAR in the dataset as a whole
  - acceptable in certain situations (e.g., small number of missing values in a large dataset)
- come up with a **replacement (imputation) value**
  - main drawback: we never know what the true value would have been
  - often the best available option

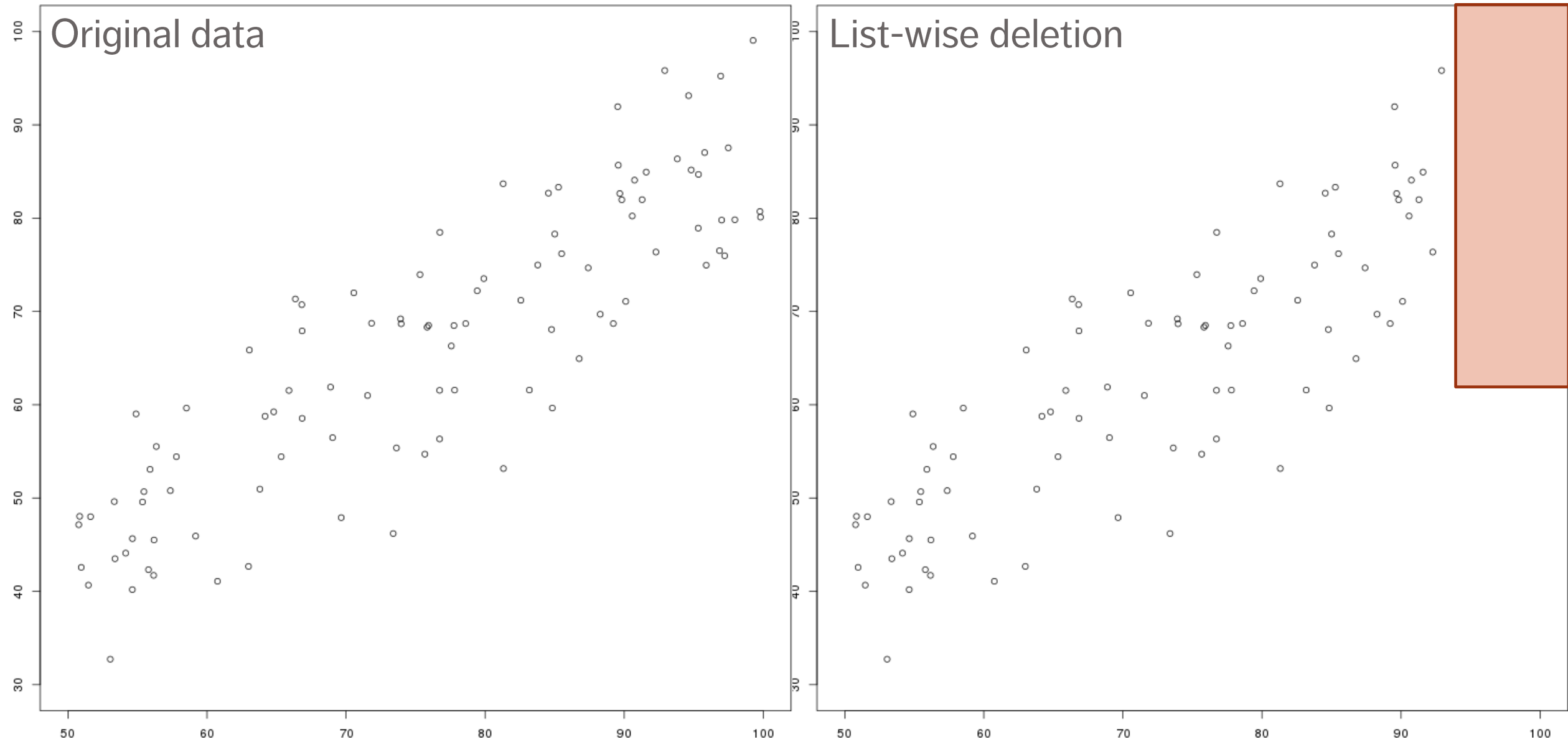


# IMPUTATION METHODS

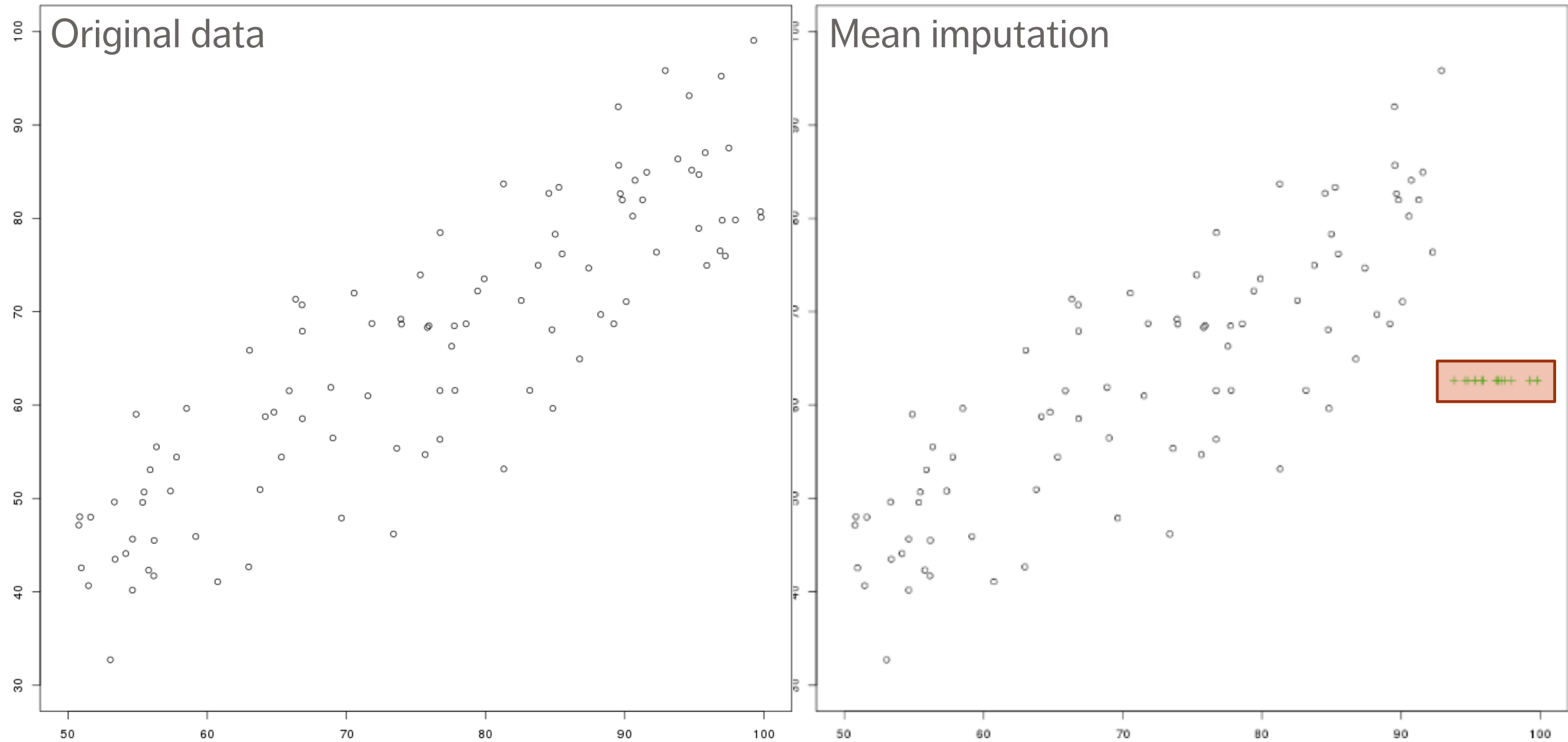
- list-wise deletion
- mean or most frequent imputation
- regression or correlation imputation
- stochastic regression imputation
- last observation carried forward
- next observation carried backward
- $k$ -nearest neighbours imputation
- multiple imputation
- etc.



**Artificial data:** the  $y$  values of all points for which  $x > 92$  have been erased by mistake.

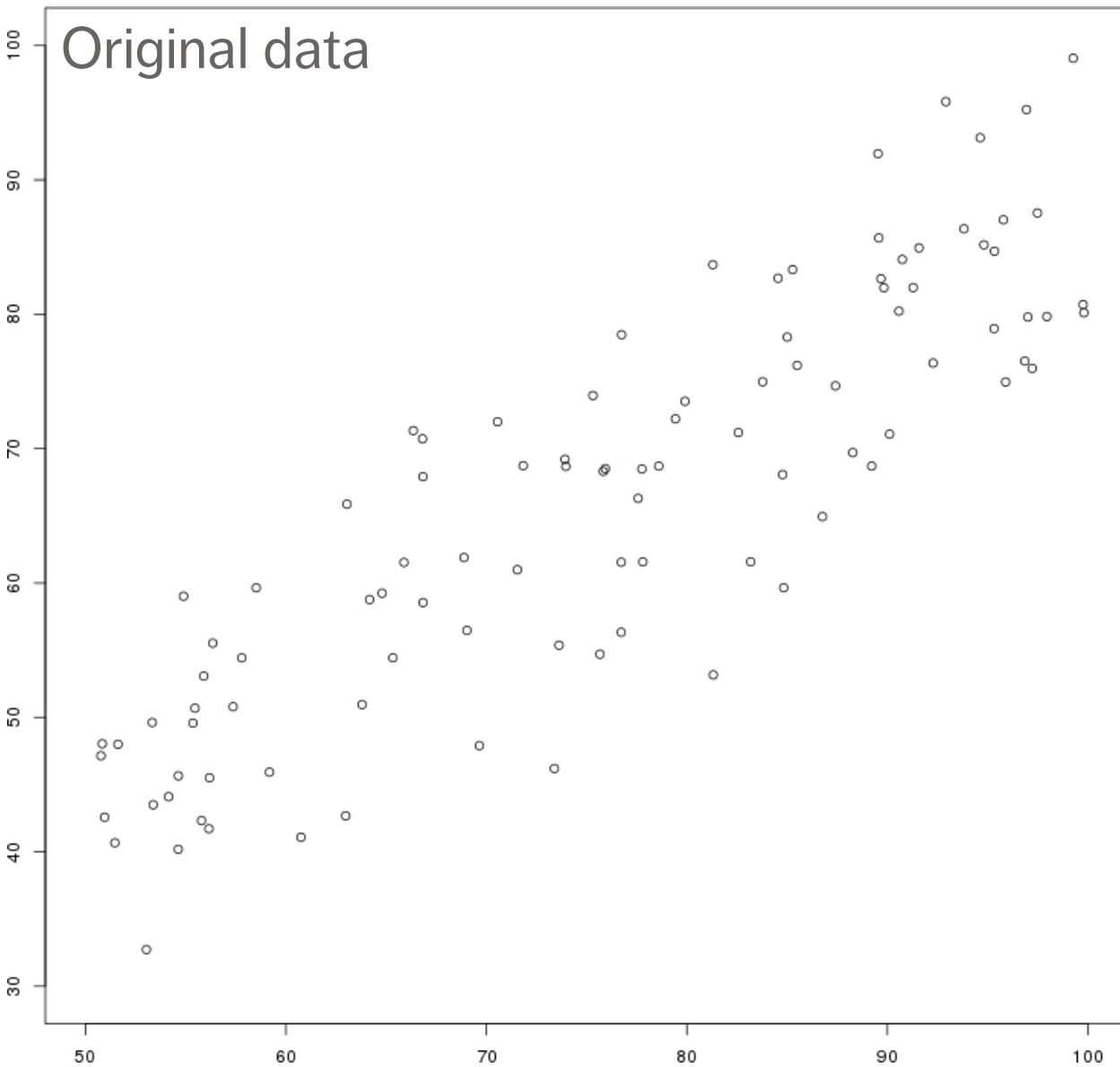


**Artificial data:** the  $y$  values of all points for which  $x > 92$  have been erased by mistake.

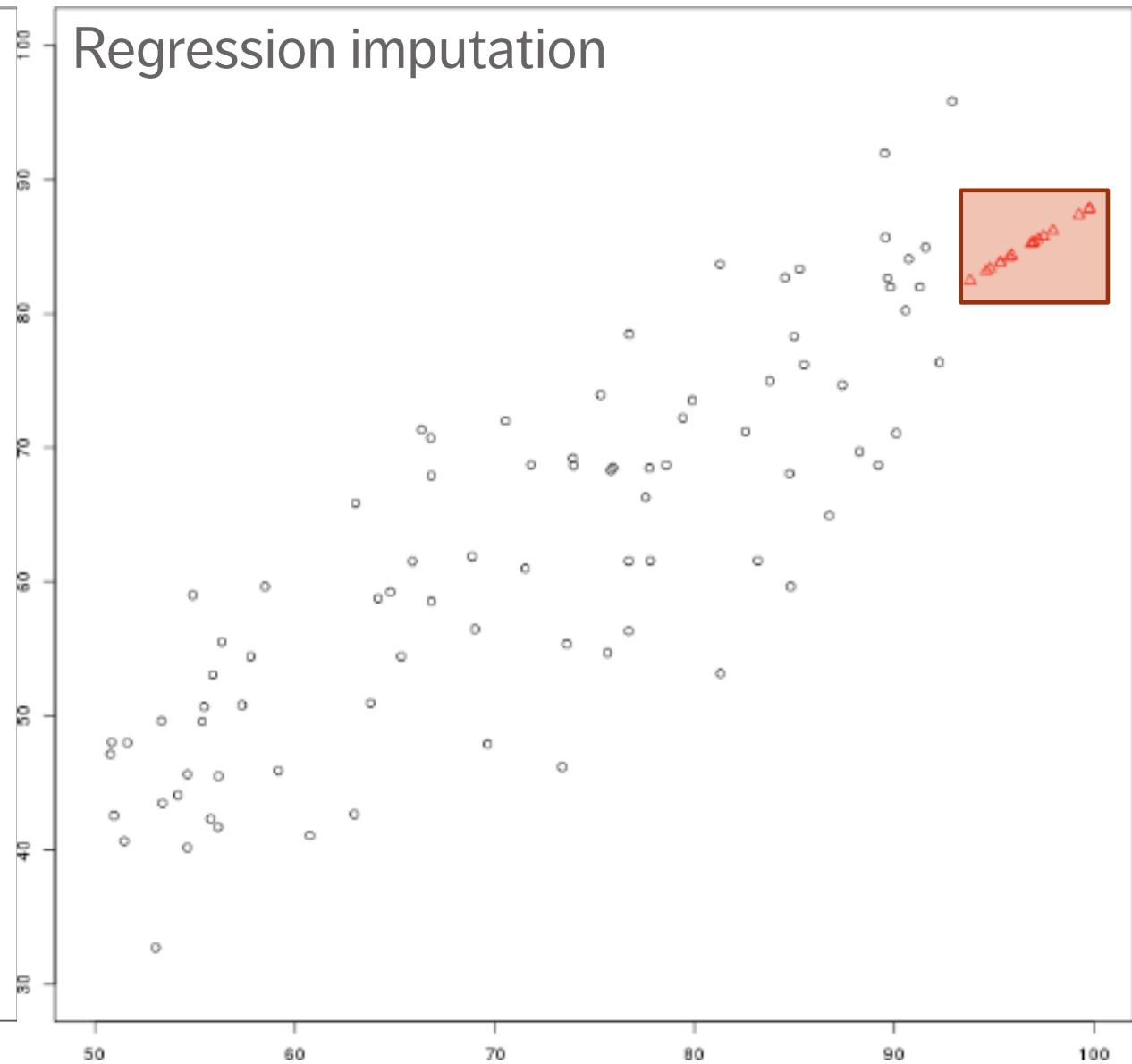


**Artificial data:** the  $y$  values of all points for which  $x > 92$  have been erased by mistake.

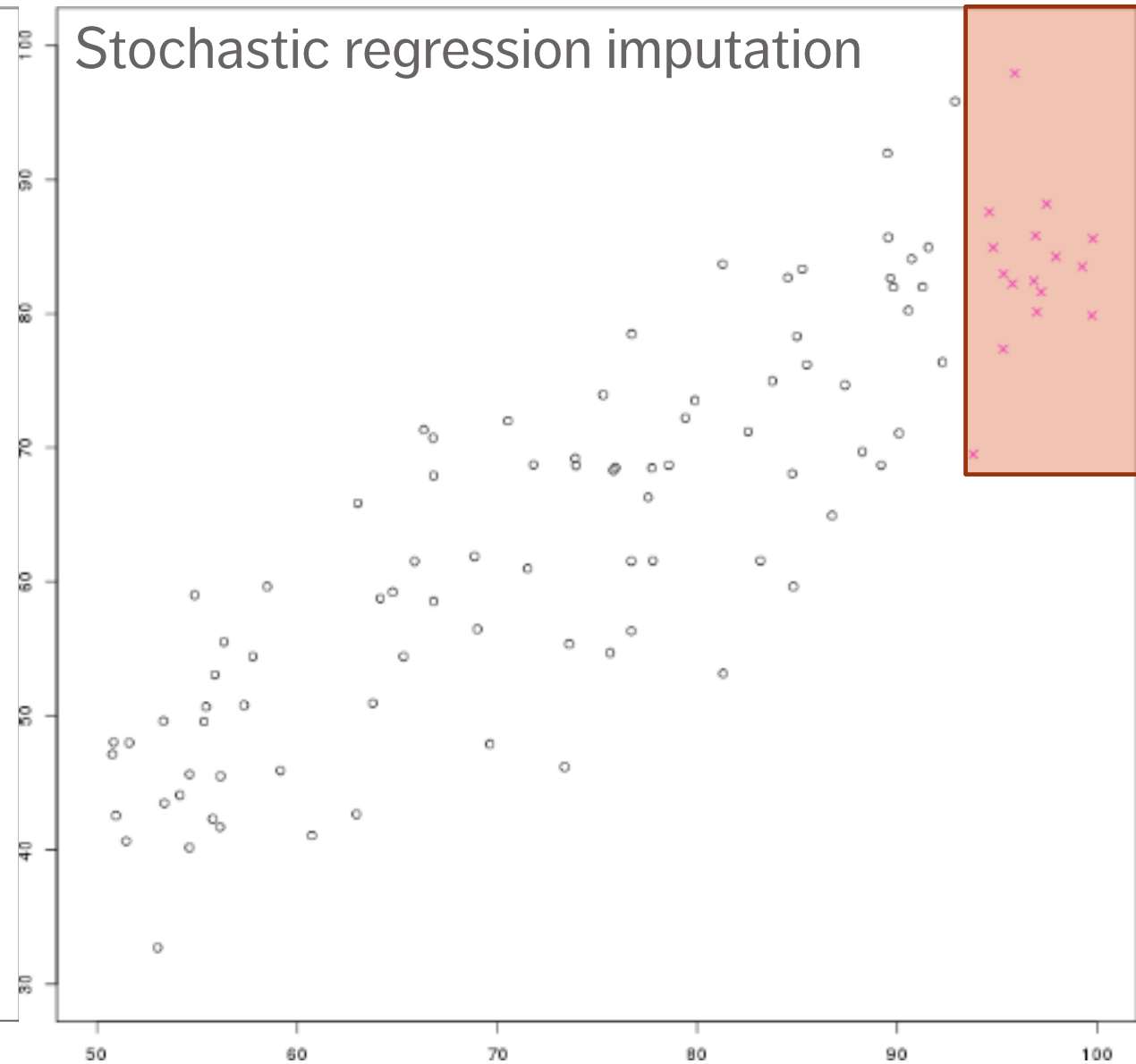
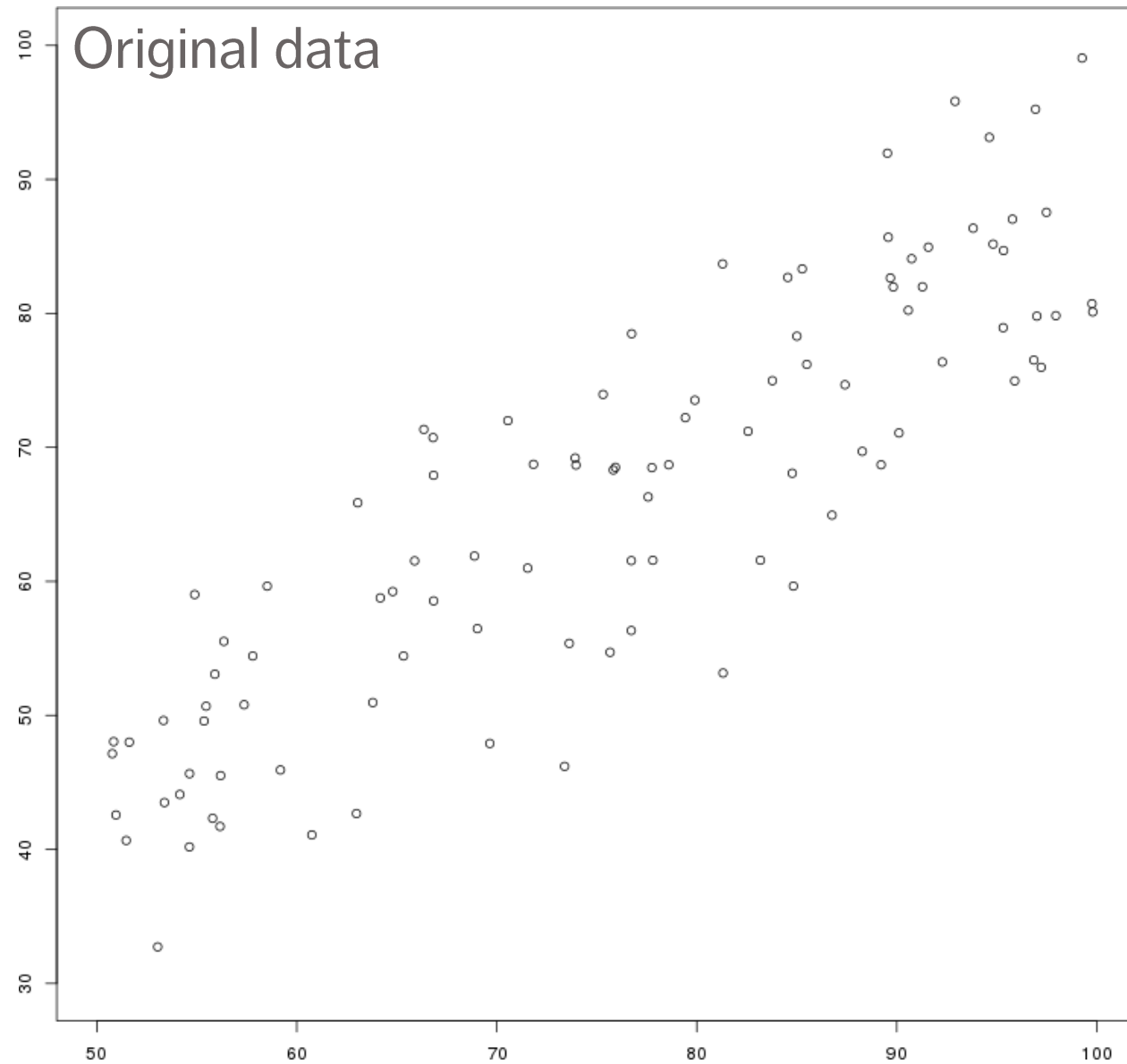
Original data



Regression imputation



**Artificial data:** the  $y$  values of all points for which  $x > 92$  have been erased by mistake.



## TAKE-AWAYS

Missing values **cannot simply be ignored**.

The missing mechanism **cannot typically be determined** with any certainty.

Imputation methods work best when values are **MCAR** or **MAR**, but imputation methods tend to produce biased estimates.

In single imputation, imputed data is treated as the actual data; multiple imputation can help reduce the noise.

Is stochastic imputation best? In our example, yes – but ... **No-Free Lunch theorem!**

# ANOMALOUS OBSERVATIONS

In practice, an **anomalous observation** may arise as

- a **“bad” object/measurement**: data artifacts, spelling mistakes, poorly imputed values, etc.
- a **misclassified observation**: according to the existing data patterns, the observation should have been labeled differently;
- an observation whose measurements are found in the **distribution tails** of a large enough number of features;
- an **unknown unknown**: a completely new type of observations whose existence was heretofore unsuspected.

---

# ANOMALOUS OBSERVATIONS

---

Observations could be anomalous in one context, but not in another:

- a 6-foot tall adult male is in the 86th percentile for **Canadian males** (tall, but not unusual);
- in **Bolivia**, the same man would be in the 99.9th percentile (very tall and unusual).

Anomaly detection points towards interesting questions for analysts and SMEs: in this case, **why is there such a large discrepancy** in the two populations?



# OUTLIERS

**Outlying observations** are data points which are **atypical** in comparison to

- the unit's remaining features (*within-unit*),
- the field measurements for other units (*between-units*)

Outliers are observations which are **dissimilar to other cases** or which **contradict known dependencies** or rules.

Careful study is needed to determine whether outliers should be retained or removed from the dataset.

# DETECTING ANOMALIES

Outliers may be anomalous along any of the unit's variables, or in combination.

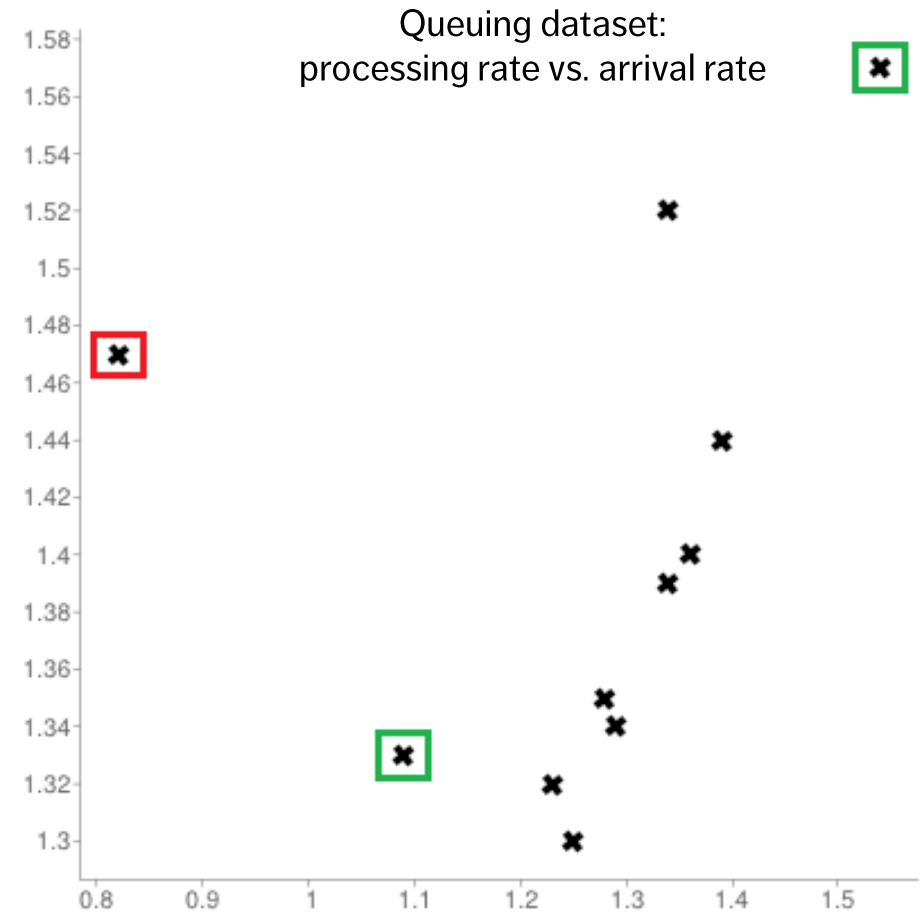
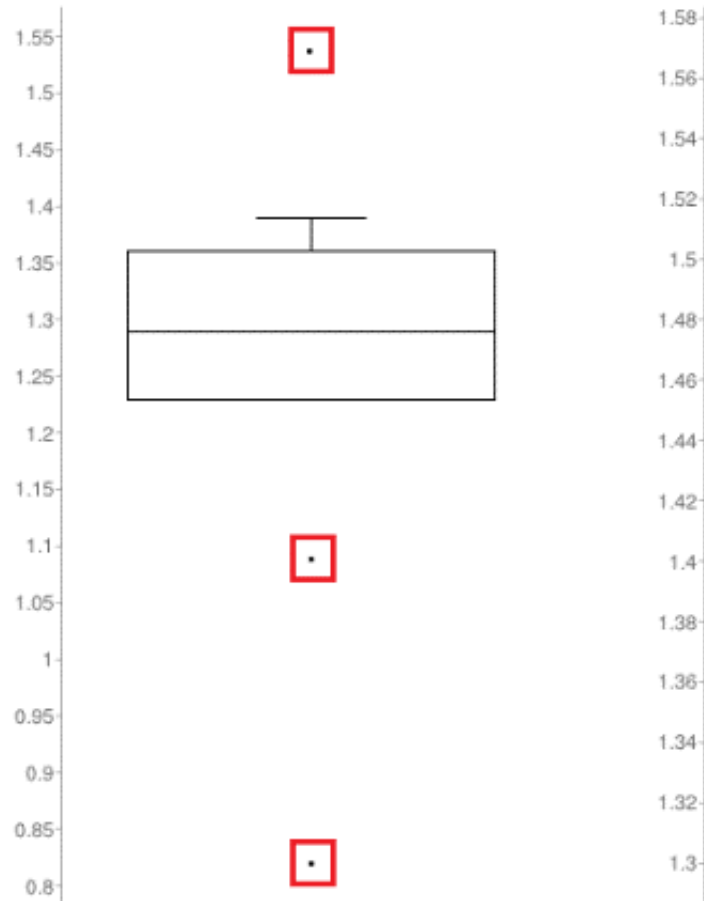
Anomalies are by definition **infrequent**, and typically shrouded in **uncertainty** due to small sample sizes.

Differentiating anomalies from noise or data entry errors is **hard**.

Boundaries between normal and deviating units may be **fuzzy**.

Anomalies associated with malicious activities are typically **disguised**.

# VISUAL OUTLIER DETECTION



# DETECTING ANOMALIES

Numerous methods exist to identify anomalous observations; **none of them are foolproof** and judgement must be used.

Graphical methods are easy to implement and interpret:

- **outlying observations**

box-plots, scatterplots, scatterplot matrices, 2D tour, Cooke's distance, normal qq plots

- **influential data**

some level of analysis must be performed (leverage)

**Careful:** once anomalous observations have been removed from the dataset, previously “regular” units may become anomalous.

# ANOMALY DETECTION ALGORITHMS

**Supervised methods** use a historical record of labeled anomalous observations:

- domain expertise is required to tag the data
- classification or regression task
- rare occurrence problem

		Predicted Class	
		Normal	Anomaly
Actual Class	Normal	<i>TN</i>	<i>FP</i>
	Anomaly	<i>FN</i>	<i>TP</i>

**Unsupervised methods** don't use external information:

- traditional methods and tests
- can also be seen as a clustering or association rules problem

# ANOMALY DETECTION ALGORITHMS

The mis-classification cost is often assumed to be symmetrical, which can lead to **technically correct but useless** outputs.

For instance, most (99.999+%) air passengers do not bring weapons with them on flights; a model that predicts that no passenger is smuggling a weapon would be 99.999+% accurate, but it would miss the point completely.

For the **security agency**, the cost of wrongly thinking that a passenger is:

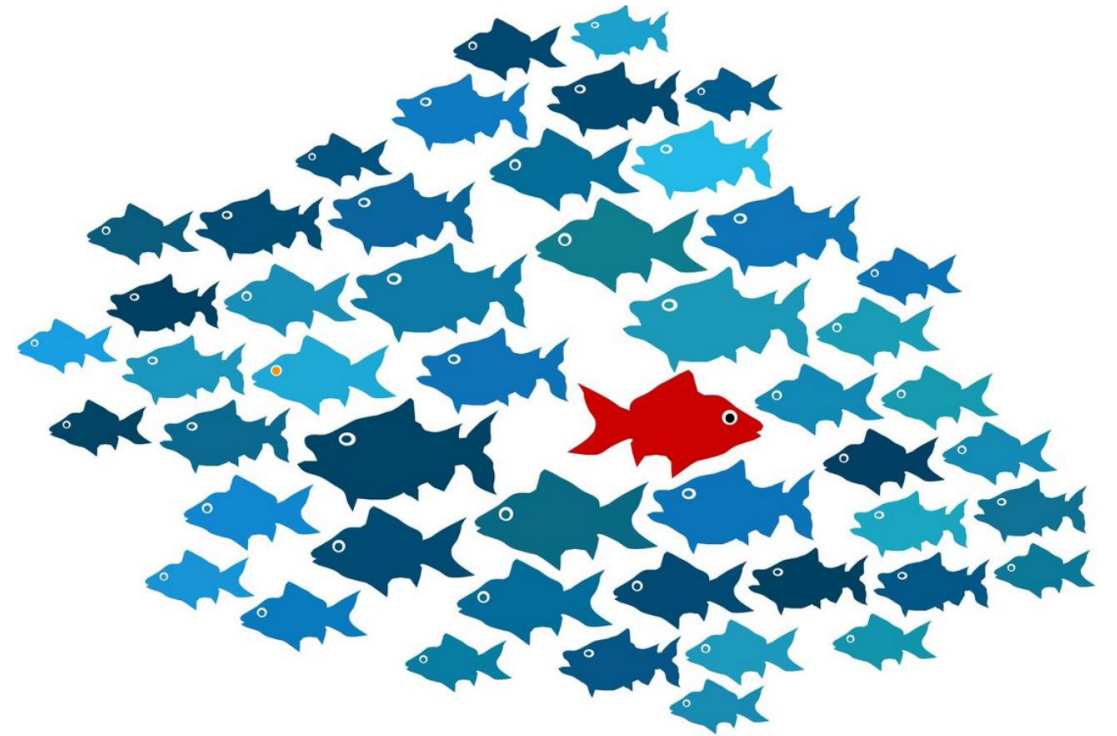
- smuggling a weapon  $\Rightarrow$  cost of a single search
- NOT smuggling a weapon  $\Rightarrow$  catastrophe (potentially)

But **wrongly targeted individuals** may have a different take on this!

# ANOMALY DETECTION ALGORITHMS

If all participants in a workshop except for one can view the video conference lectures, then the one individual/internet connection/computer is **anomalous** – it behaves in a manner which is different from the others.

But this **DOES NOT MEAN** that the different behaviour is necessarily the one we are interested in...



# INFLUENTIAL OBSERVATIONS



**Influential data points** are observations whose absence leads to **markedly different** analysis results.



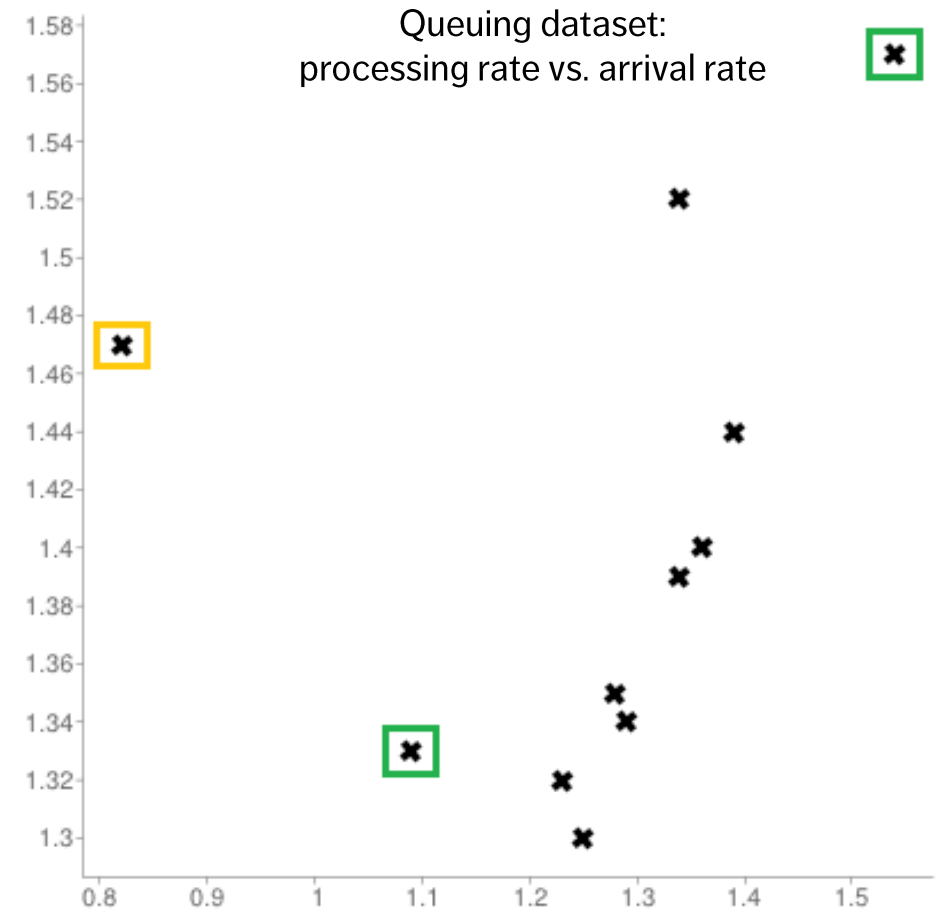
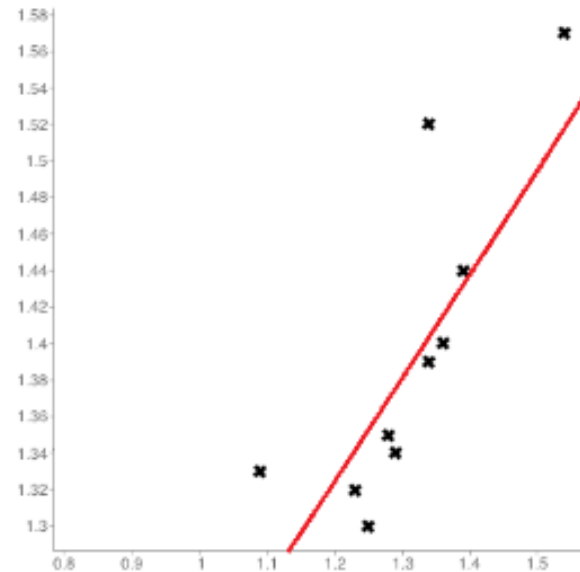
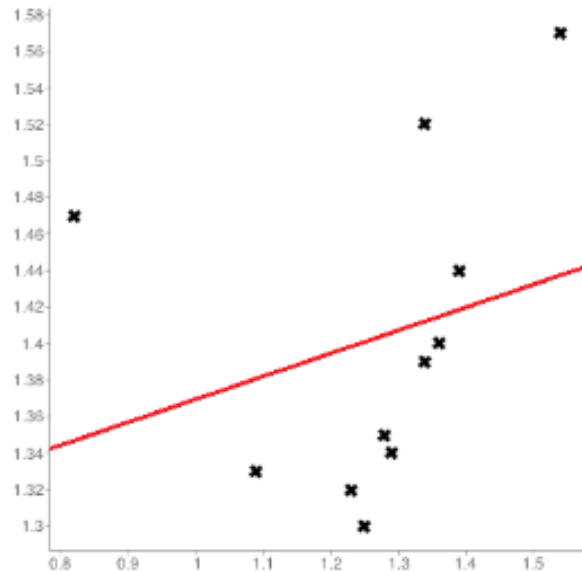
When influential observations are identified, **remedial measures** (such as data transformations) may be required to minimize their undue effects.



Outliers may be influential data points; influential data points need not be outliers (and *vice-versa*).



# INFLUENTIAL OBSERVATIONS



# EXERCISE

Find anomalous observations in the [cities.txt](#) and [HR\\_2016\\_Census\\_simple.xlsx](#) datasets (if applicable).

# DIMENSIONALITY OF DATA

In data analysis, the **dimension** of the data is the number of attributes that are collected in a dataset, represented by the **number of columns**.

We can think of the number of variables used to describe each object (row) as a vector describing that object: the dimension is simply the **size** of that vector.

**(Note:** “dimension” is used differently in business intelligence contexts)

# HIGH DIMENSIONALITY AND BIG DATA

Datasets can be “big” in a variety of ways:

- too large for the **hardware** to handle (cannot be stored, accessed, manipulated properly due to # of observations, # of features, the overall size)
- dimensions can go against **modeling assumptions** (# of features  $\gg$  # observations)

## Examples:

- multiple sensors recording 100+ observations per second in a large geographical area over a long time period = **very big dataset**
- in a corpus' *Term Document Matrix* (cols = terms, rows = documents), the number of terms is usually substantially higher than the number of documents, leading to **sparse data**

# SAMPLING OBSERVATIONS

**Question:** does every row of the dataset need to be used?

If rows are selected randomly (with or without replacement), the resulting sample might be **representative** of the entire dataset.

## **Drawbacks:**

- if the signal of interest is rare, sampling might drown it altogether
- if aggregation is happening down the road, sampling will necessarily affect the numbers (passengers vs. flights)
- even simple operations on a large file (finding the # of lines, say) can be taxing on the memory – **prior information on the dataset structure can help**

# FEATURE SELECTION

Removing **irrelevant/redundant** variables is a common data processing task.

## Motivations:

- modeling tools do not handle these well (variance inflation due to multicollinearity, etc.)
- dimension reduction ( $\#$  variables  $\gg$   $\#$  observations)

## Approaches:

- filter vs. wrapper
- unsupervised vs. supervised

# COMMON TRANSFORMATIONS

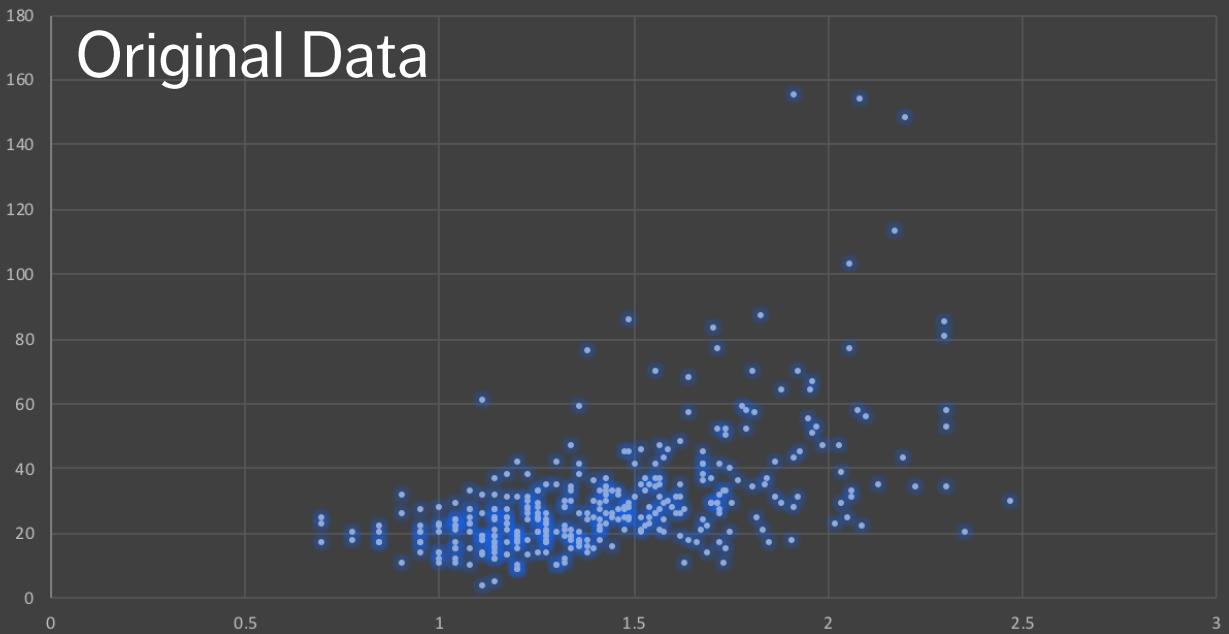
Models sometimes require that certain data assumptions be met (normality of residuals, linearity, etc.).

If the raw data does not meet the requirements, we can either:

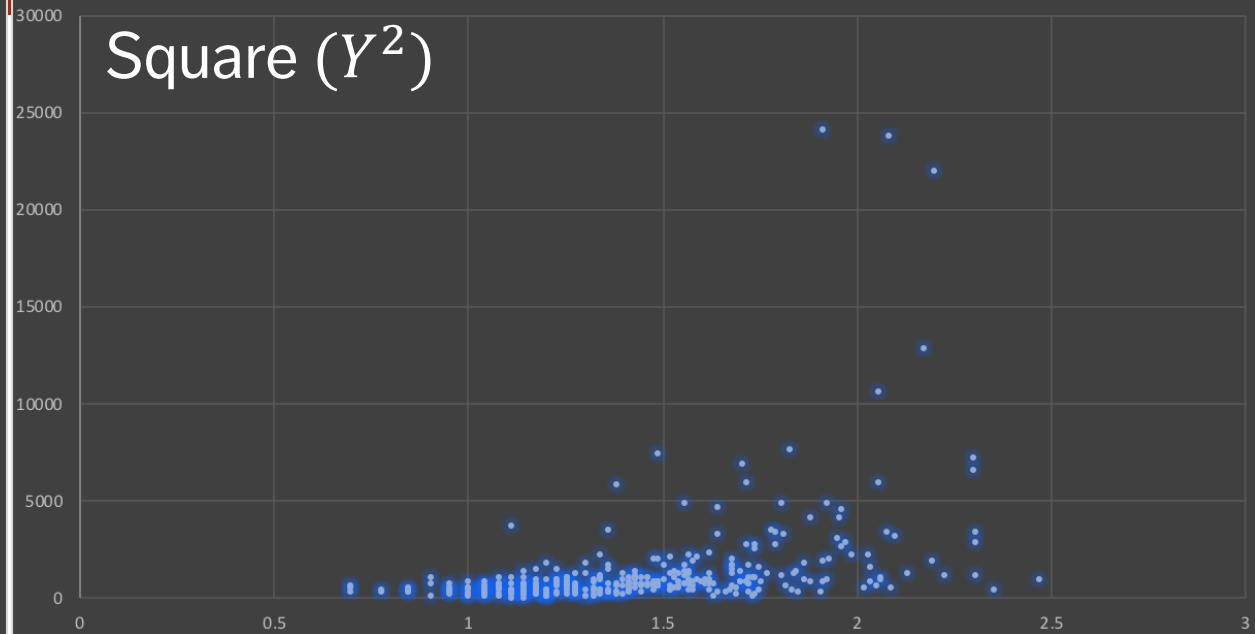
- abandon the model
- attempt to **transform** the data

The second approach requires an **inverse transformation** to be able to draw conclusions about the **original data**.

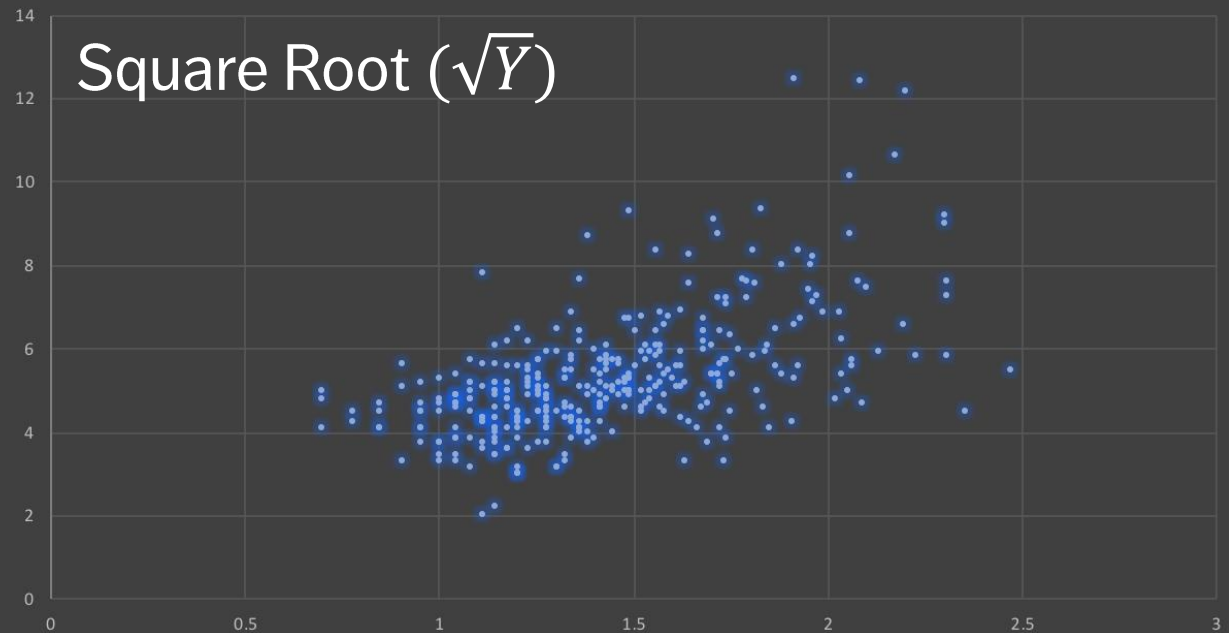
# Original Data



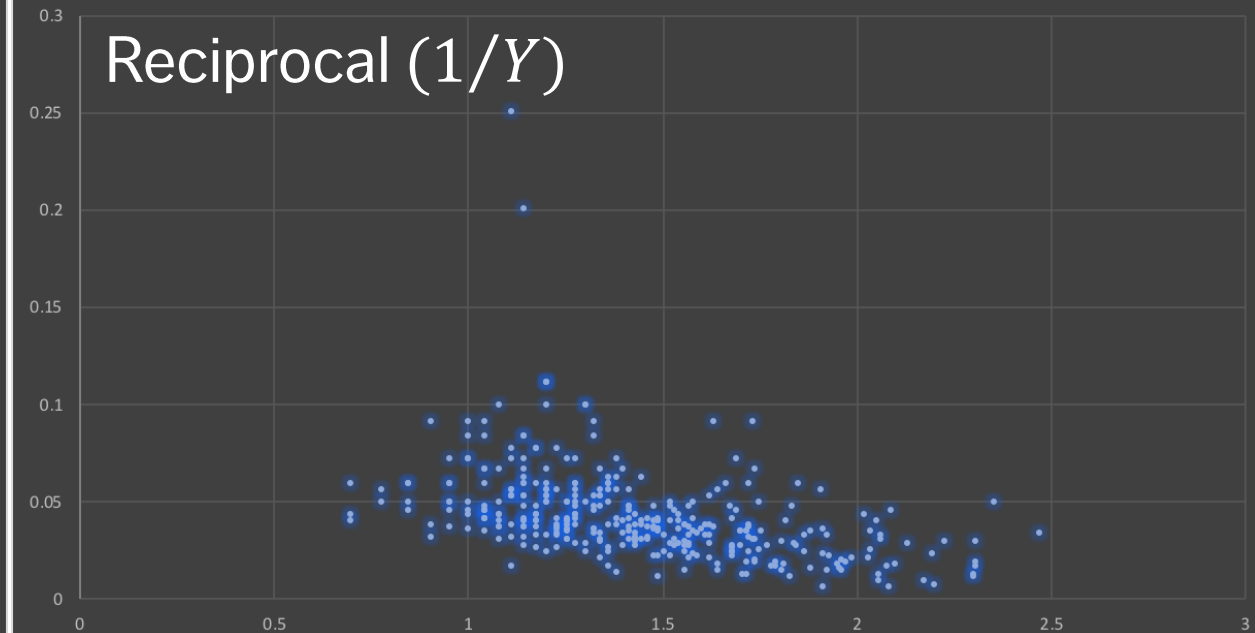
# Square ( $Y^2$ )



# Square Root ( $\sqrt{Y}$ )

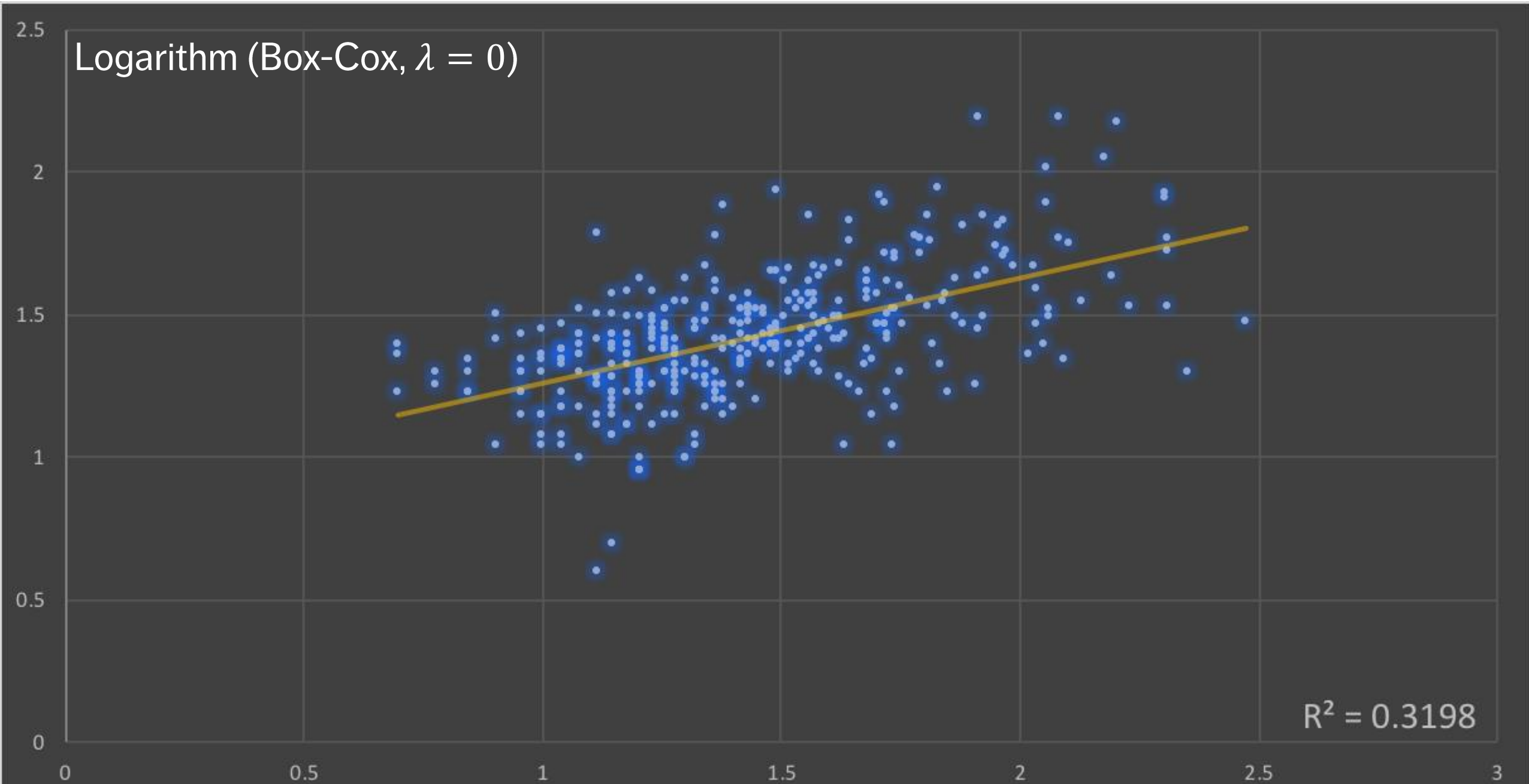


# Reciprocal ( $1/Y$ )





Logarithm (Box-Cox,  $\lambda = 0$ )



# SCALING

Numeric variables may have different **scales** (i.e., weights and heights).

The variance of a large-range variable is typically greater than that of a small-range variable, introducing a bias (for instance).

**Standardization** creates a variable with mean 0 and std. dev. 1:

$$Y_i = \frac{X_i - \bar{X}}{s_X}$$

**Normalization** creates a new variable in the range [0,1]:  $Y_i = \frac{X_i - \min X}{\max X - \min X}$

# DISCRETIZING

To reduce computational complexity, a numeric variable may need to be replaced by an **ordinal** variable (from *height* value to “*short*”, “*average*”, “*tall*”, for instance).

**Domain expertise** can be used to determine the bins’ limits (although that may introduce unconscious bias to the analyses)

In the absence of such expertise, limits can be set so that either

- the bins each contain the same number of observations
- the bins each have the same width
- the performance of some modeling tool is maximized

# CREATING VARIABLES

New variables may need to be introduced:

- as **functional relationships** of some subset of available features
- because modeling tool may require **independence of observations**
- because modeling tool may require **independence of features**
- to simplify the analysis by looking at **aggregated summaries** (often used in text analysis)

Time dependencies → time series analysis (lags?)

Spatial dependencies → spatial analysis (neighbours?)

# EXERCISE

Scale, discretize, and/or create new variables out of the [cities.txt](#) and [HR\\_2016\\_Census\\_simple.xlsx](#) datasets.

---

# SUPPLEMENTAL MATERIAL

## 8. DATA ANALYSIS

# SIMPLE OUTLIER TESTS

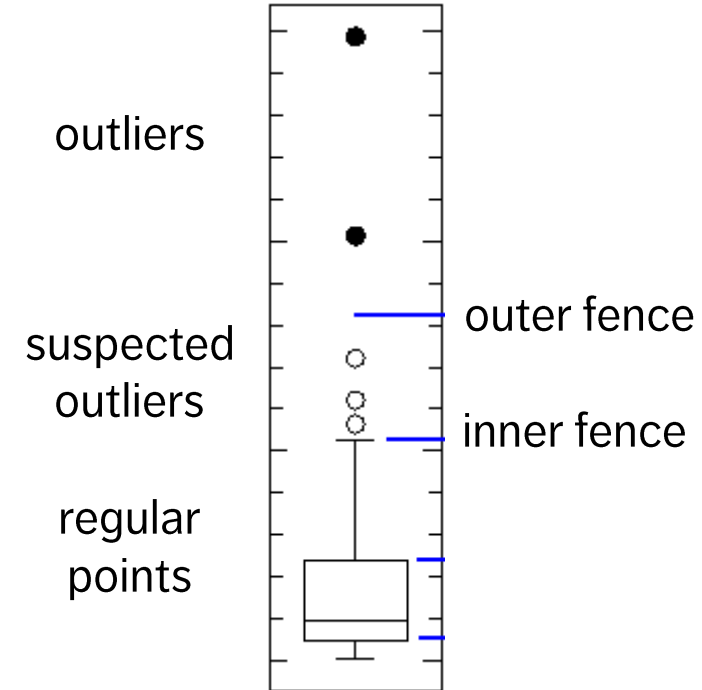
**Tukey's Boxplot test:** for normally distributed data, regular observations typically lie between the **inner fences**

$$Q_1 - 1.5 \times (Q_3 - Q_1) \text{ and } Q_3 + 1.5 \times (Q_3 - Q_1).$$

**Suspected outliers** lie between the **inner fences** and the **outer fences**

$$Q_1 - 3 \times (Q_3 - Q_1) \text{ and } Q_3 + 3 \times (Q_3 - Q_1).$$

**Outliers** lie beyond the **outer fences**.



# SIMPLE OUTLIER TESTS

The **Dixon Q Test** is used in experimental sciences to find outliers in (extremely) small datasets (dubious validity).

The **Mahalanobis Distance** (linked to the leverage) can be used to find multi-dimensional outliers (when relationships are linear).

Other simple tests:

- **Grubbs** (univariate)
- **Tietjen-Moore** (for a specific # of outliers)
- **generalized extreme studentized deviate** (for unknown # of outliers)
- **chi-square** (outliers affecting goodness-of-fit)



# MISSING VALUES MECHANISM

## Missing Completely at Random (MCAR)

- item absence is independent of its value or of auxiliary variables
- **example:** an electrical surge randomly deletes an observation in the dataset

## Missing at Random (MAR)

- item absence is not completely random; can be accounted by auxiliary variables with complete info
- **example:** if women are less likely to tell you their age than men for societal reasons, but not because of the age values themselves)

# MISSING VALUES MECHANISM

## Not Missing at Random (NMAR)

- reason for nonresponse is related to item value (also called **non-ignorable non-response**)
- **example:** if illicit drug users are less likely to admit to drug use than teetotallers

In general, the missing mechanism **cannot be determined** with any certainty; we may need to make assumptions (domain expertise can help).

# IMPUTATION METHODS

**List-wise deletion:** remove units with at least one missing values

- **assumption:** MCAR
- **cons:** can introduce bias (if not MCAR), reduction in sample size, increase in standard error

**Mean/most frequent imputation:** substitute missing values by average/most frequent value

- **assumption:** MCAR
- **cons:** distortions of distribution (spike at mean) and relationships among variables

# IMPUTATION METHODS

**Regression/correlation imputation:** substitute missing values using fitted values based on other variables with complete information

- **assumption:** MAR
- **cons:** artificial reduction in variability, over-estimation of correlation

**Stochastic regression imputation:** regression/correlation imputation with a random error term added

- **assumption:** MAR
- **cons:** increased risk of type I error (false positives) due to small std error

# IMPUTATION METHODS

**Last observation carried forward:** substitute the missing values with latest previous values (in a longitudinal study)

- **assumption:** MCAR, values do not vary greatly over time
- **cons:** may be too “generous”, depending on the nature of study

**$k$  nearest neighbour imputation ( $k$ NN):** substitute the missing entry with the average from the group of the  $k$  most similar complete cases

- **assumption:** MAR
- **cons:** difficult to choose appropriate value for  $k$ ; possible distortion in data structure

# MULTIPLE IMPUTATION

Imputations increase the noise in the data.

In **multiple imputation**, the effect of that noise can be measured by consolidating the analysis outcome from multiple imputed datasets

## Steps:

1. repeated imputation creates  $m$  versions of the dataset
2. each of these datasets is analyzed, yielding  $m$  outcomes
3. the  $m$  outcomes are pooled into a single result for which the mean, variance, and confidence intervals are known

# MULTIPLE IMPUTATION

## Advantages

- **flexible**; can be used in a various situations (MCAR, MAR, even NMAR in certain cases)
- accounts for **uncertainty** in imputed values
- fairly easy to implement

## Disadvantages

- $m$  may need to be fairly **large** when there are many missing values in numerous features, which slows down the analyses
- if the analysis output is not a single value but some complicated mathematical object, this approach is unlikely to be useful

# ANOMALY DETECTION REMARKS

Identifying influential points is an **iterative process** as the various analyses have to be run numerous times.

Fully automated identification and removal of anomalous observations is **NOT recommended**.

Use data transformations if the data is **NOT normally distributed**.

Whether an observation is an outlier or not depends on **various factors**; what observations end up being influential data points depends on the **specific analysis to be performed**.

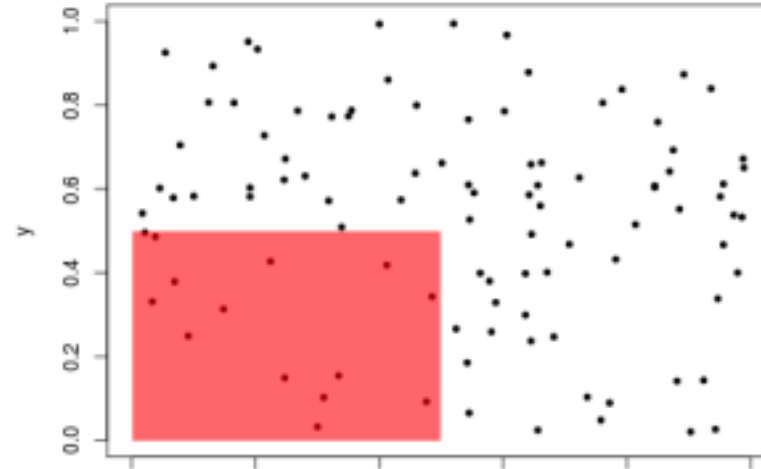


# CURSE OF DIMENSIONALITY

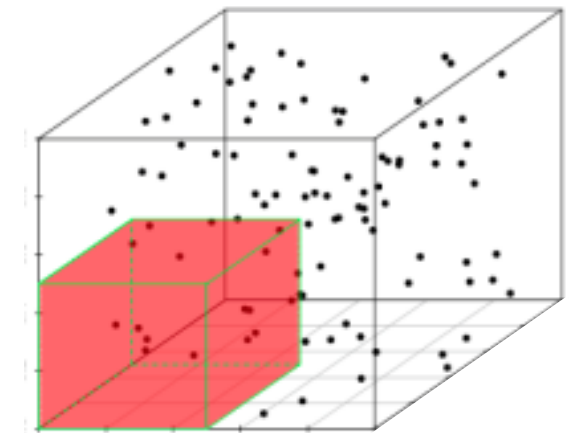
42% of data is captured



14% of data is captured



7% of data is captured



$N = 100$  observations, uniformly distributed on  $[0, 1]^d$ ,  $d = 1, 2, 3$ .  
% of observations captured by  $[0, 1/2]^d$ ,  $d = 1, 2, 3$ .

# COMMON TRANSFORMATIONS

In the data analysis context, transformations are **monotonic**:

- logarithmic
- square root, inverse, power:  $W^k$
- exponential
- Box-Cox, etc.

Transformations on  $X$  may achieve linearity, but usually at some price (correlations are not preserved, for instance). Transformations on  $Y$  can help with non-normality and unequal variance of error terms.