

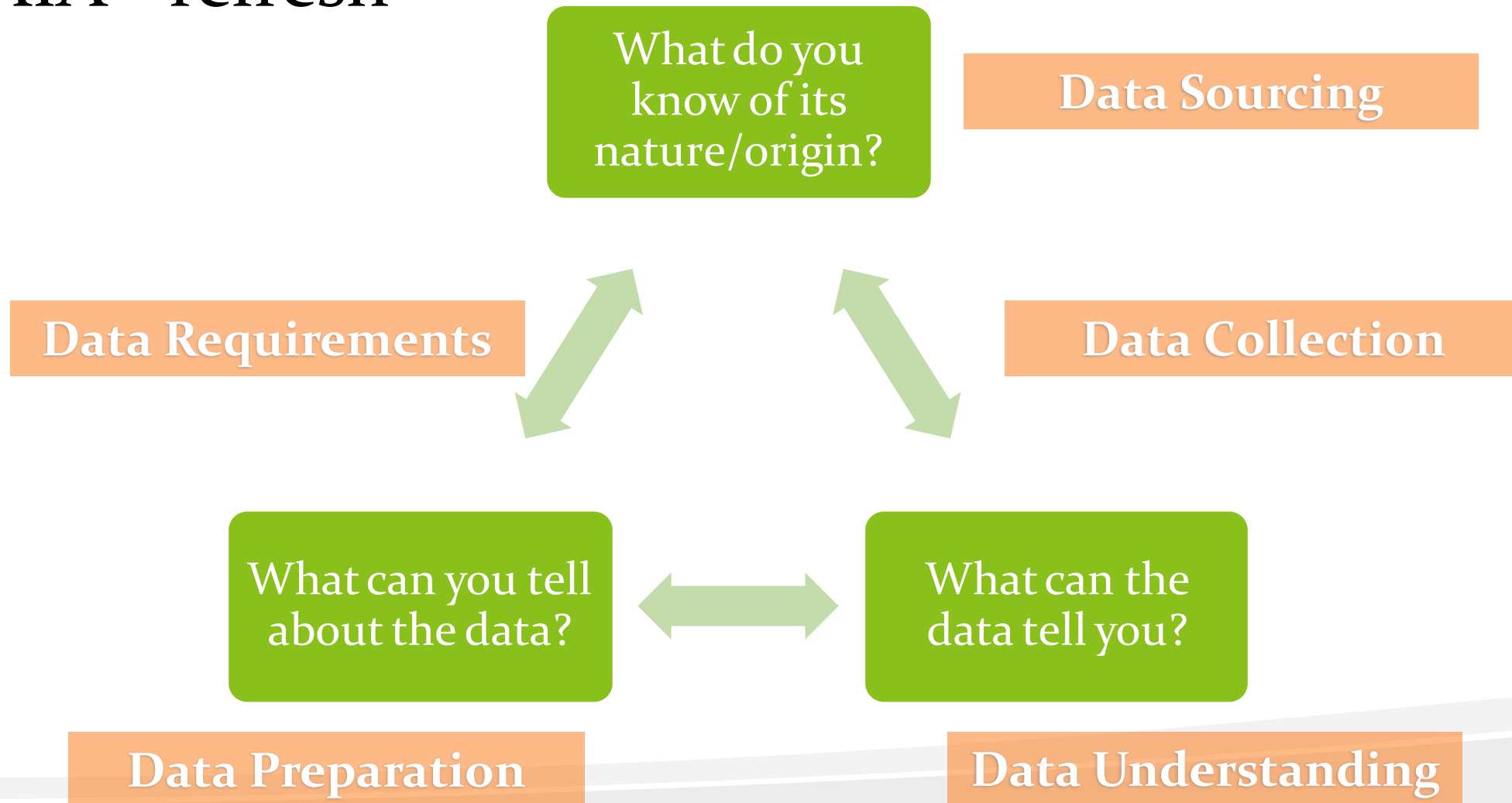


# DAIA

Visualizations (basics) – week 2



# DAIA - refresh



# AI project methodology: your roadmap



Target variable, model and domain requirements, data source(s), ...  
**What do you want to achieve (predict)?**

What data (quality) is required?  
Define a data dictionary

How do you get (generate) and combine your data? Capture the process.

Explore your data (EDA and EDV)

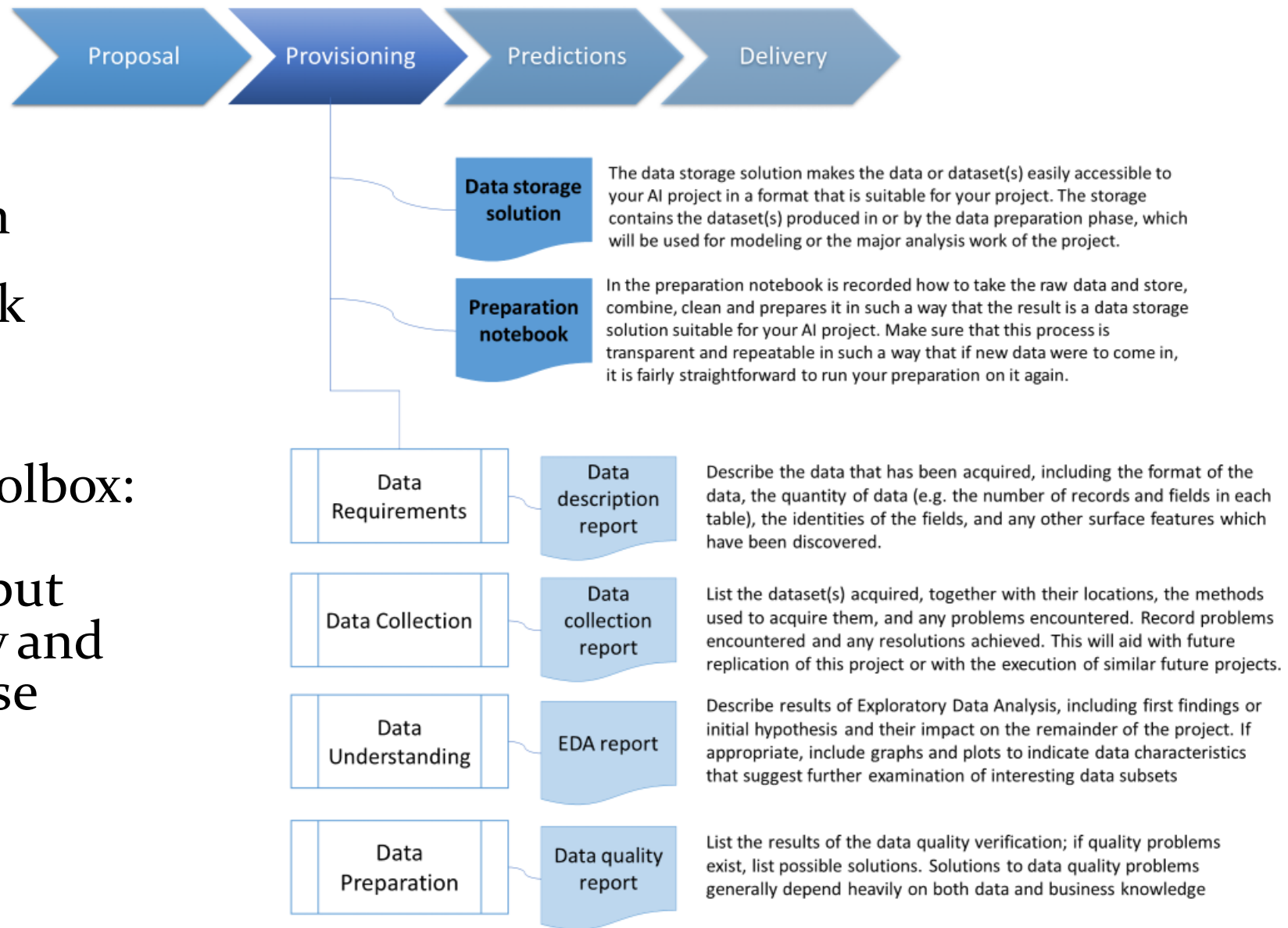


# Deliverables (beroepsproducten)

Prepared dataset:

- Data storage solution
- Preparation notebook

Provisioning is a big toolbox:  
you don't need to use/  
demonstrate all tools, but  
select tools consciously and  
show how to apply those  
properly.





# The exploration process

- Start with **data**
- Decide what your target audience is and what **questions** are interesting to investigate
- Explore the data to get **insights**, guided by the questions
- Start with looking at **data distribution, relationships and timelines**
- Enrich with **new datasets**, ask **new questions**
- Use descriptive statistics and advanced charts





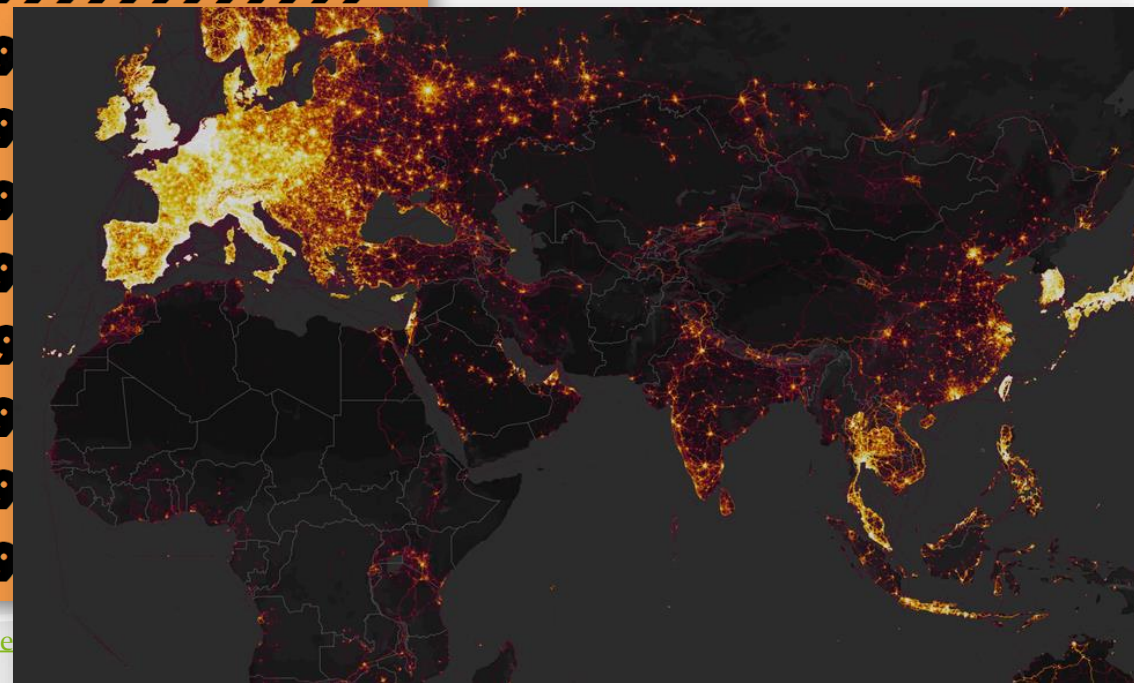
## CAN YOU FIND 8 IN BELOW TEXT?

999

999

**999999999999999999999999999999999999**

999



nd-puzzle



# The Exploration Process



**Data**

categorical and numerical *variables*, (not) clean, sufficiently diverse

**Descriptive statistics**

measures for *center* and *dispersion* of the dataset

**Basic charts**

Histogram, boxplot, barchart, scatterplot, linechart

**Relevant insights**

Outliers, Correlations, Clusters, Trends

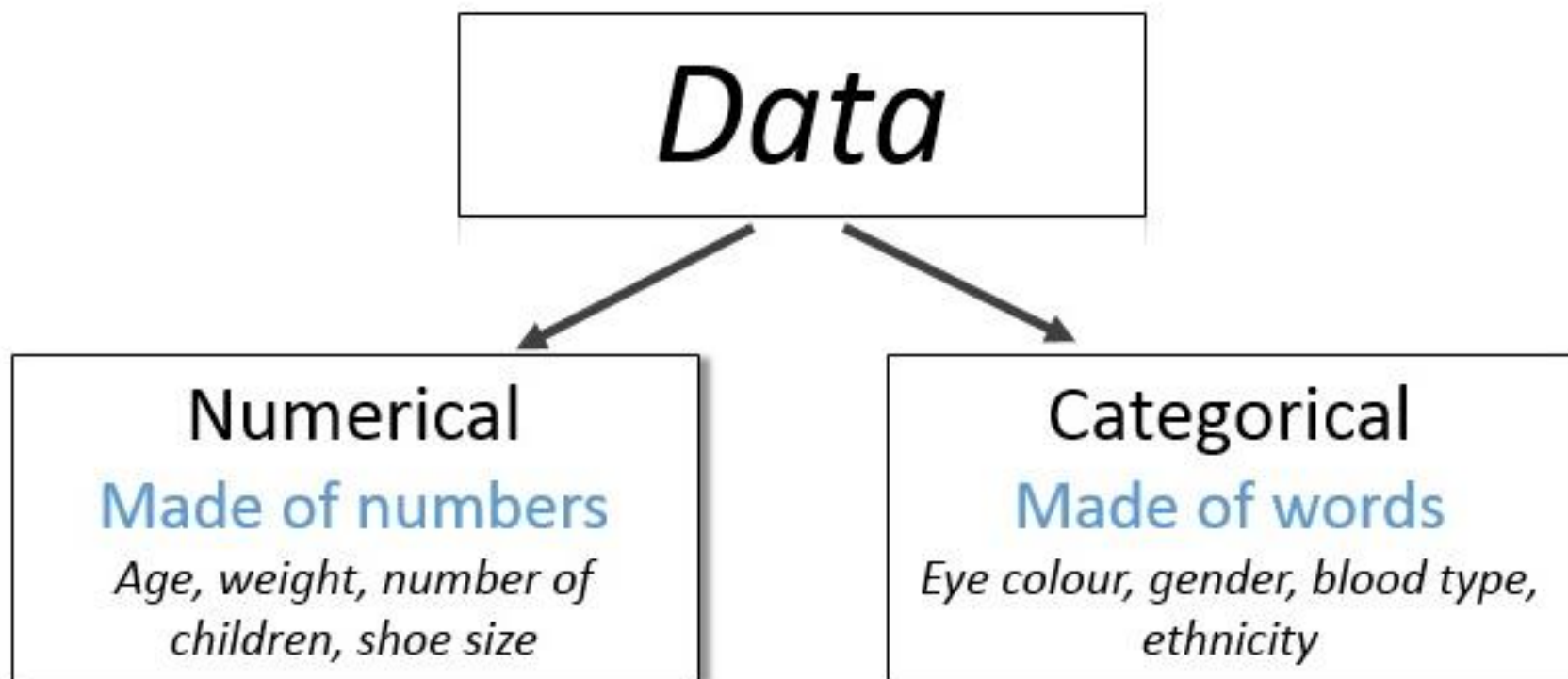
**Reporting**

Dashboards, clear visuals, well-explained answers, advice

explore

focus

# Types of Data



What about **discrete**  
and **continuous**  
numbers?

What about **ordinal** and  
**nominal** data?





# Basic guidelines for data visualization

## Descriptive Statistics

- central tendency
- Spread

## Basic Charts

1. **Histogram** -> distribution, central tendency, spread
2. **Boxplot** -> distribution, central tendency, spread, comparisons
3. **Bar Chart** -> outliers, comparisons
4. **Scatter Plot** -> outliers, correlations
5. **Line Chart** -> trends



# Descriptive statistics: Central Tendency

Measures to describe the **center** of the set

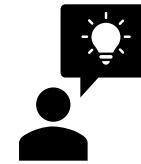
**Mean** = average

**Median** = middle value

**Mode** = most frequent value

Only for  
**numerical**  
data

Typically for  
**categorical**  
data



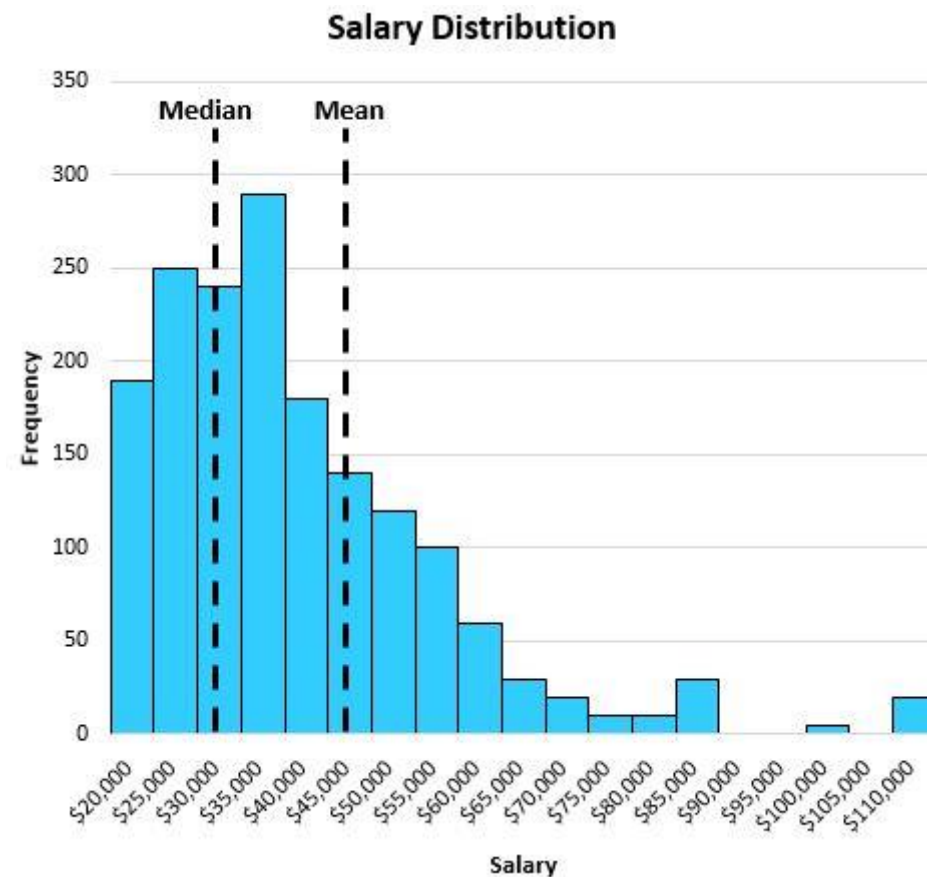
“Elon Musk walks into a bar, and everyone inside becomes a millionaire, on average.”



# Descriptive statistics: Central Tendency



VS

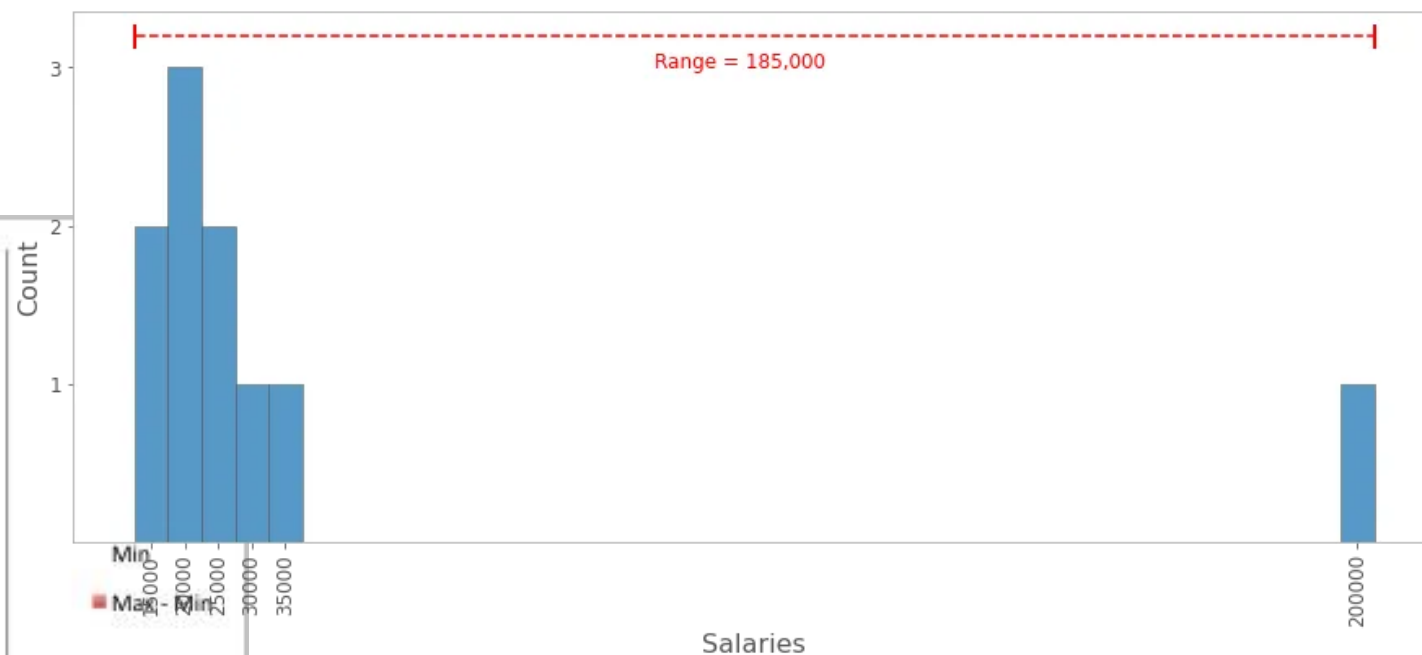
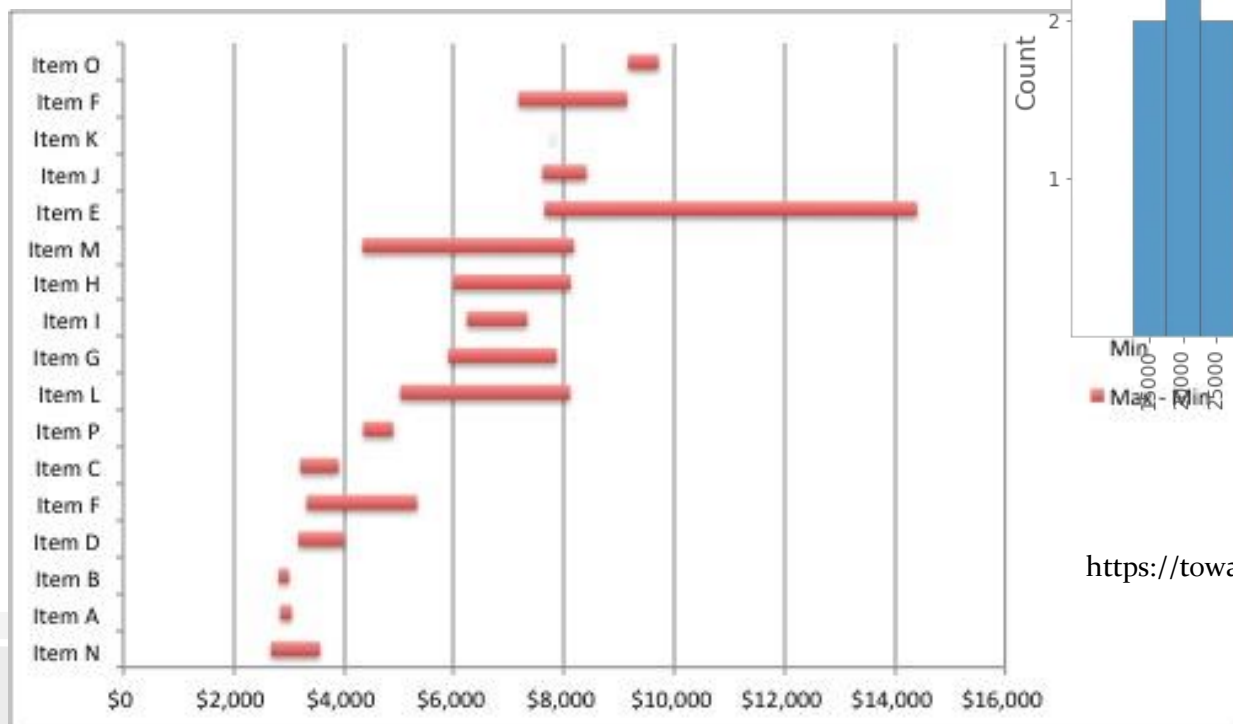




# Descriptive statistics: Spread

Measures to describe the **spread/dispersion** of the set

- **Range** = max – min



<https://towardsdatascience.com/statistics-02-measuring-and-visualizing-the-spread-of-data-2fc31d928830>

# Descriptive statistics: Spread

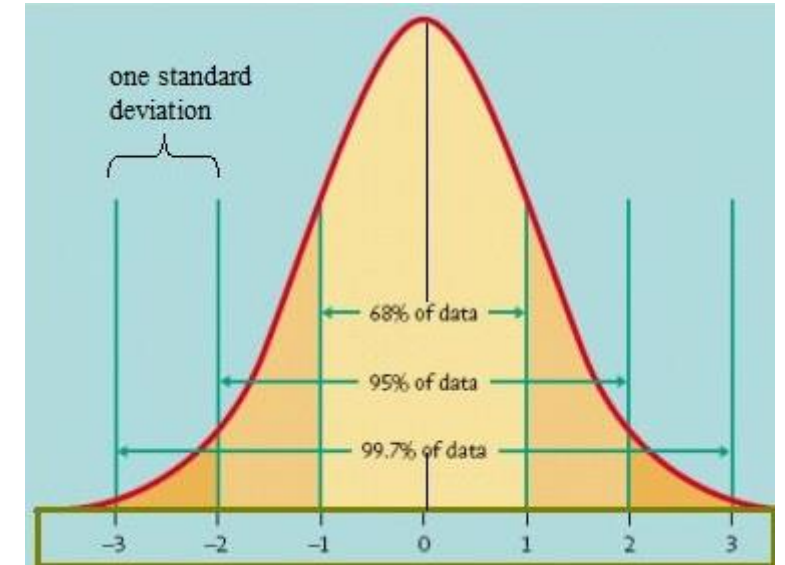
Measures to describe the **spread/dispersion** of the set

- Standard deviation  $\sigma$
- variance =  $\sigma^2$

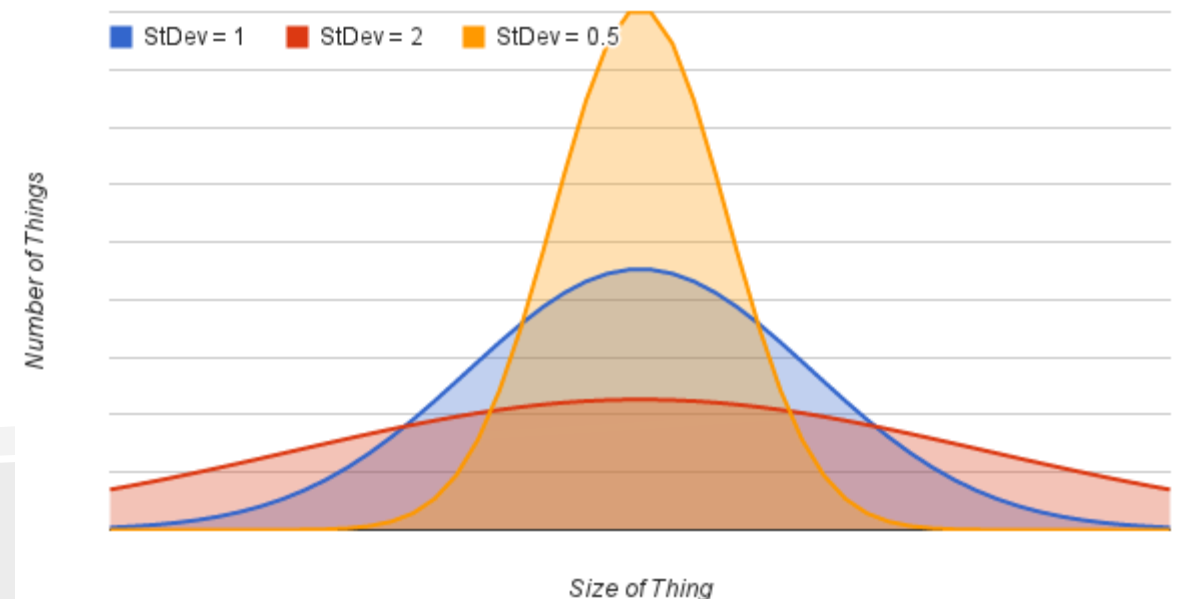
Only for  
numerical  
data

$$\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

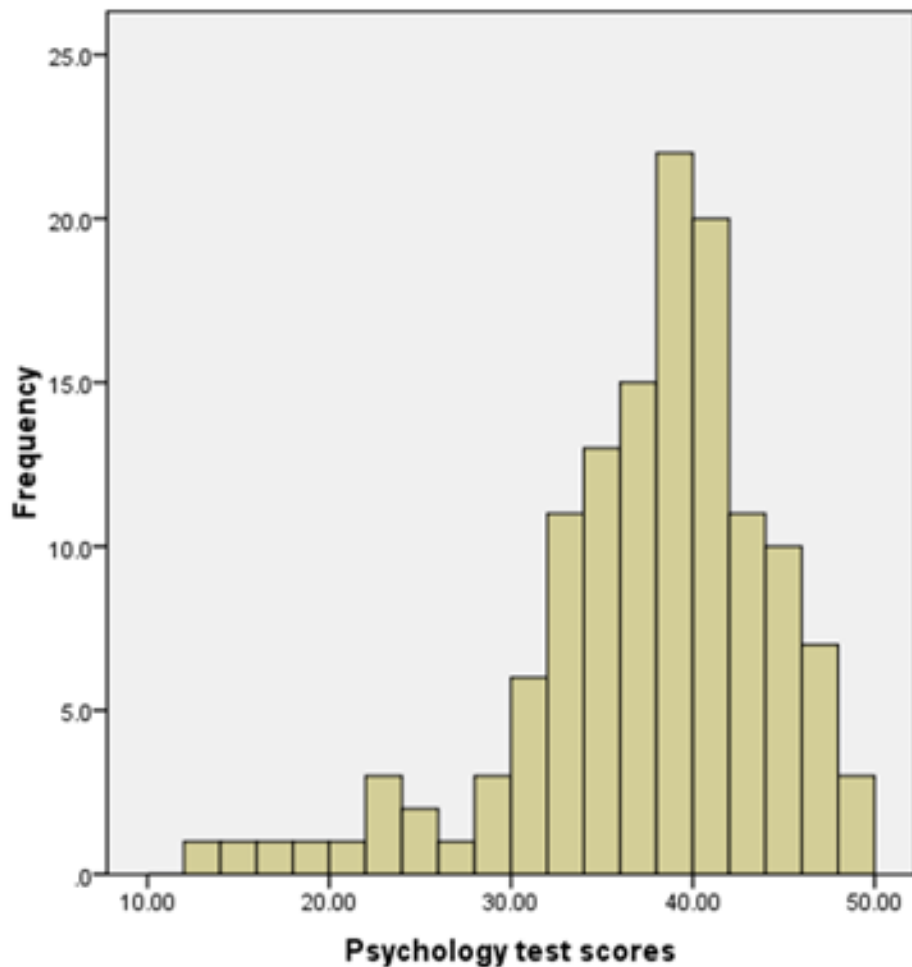
‘distance’ to the mean value



Normal Distributions with Different Standard Deviations

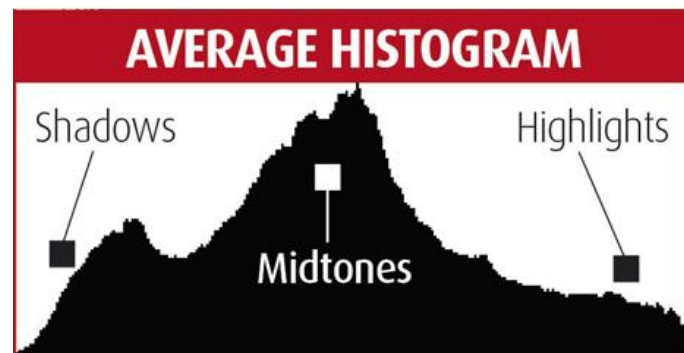


# Basic Charts – 1. Histogram



Mean =37.0227  
Std. Dev. =6.82549  
N =132

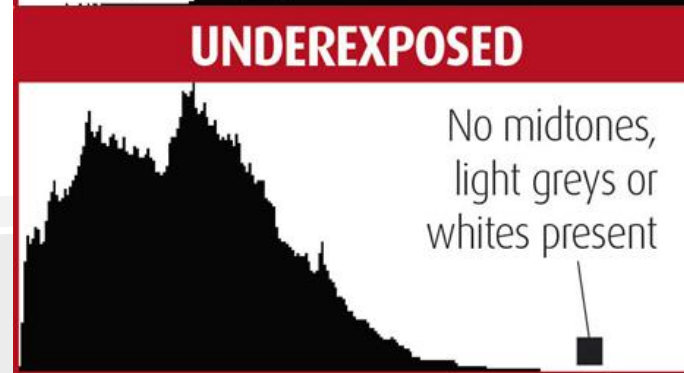
1 nummmerical  
variable



correct exposure

=

the peak of the black mountain  
sitting halfway between shadows  
and highlights

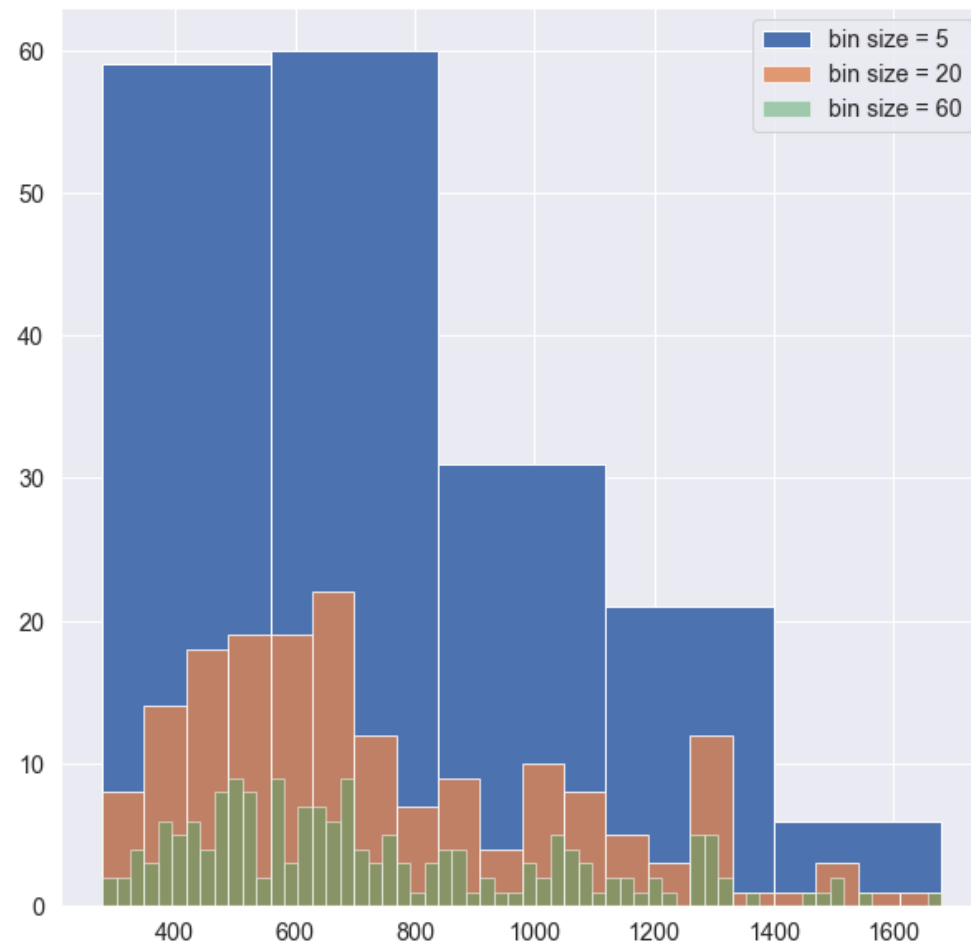




# Basic Charts – 1. Histogram

Bin size:

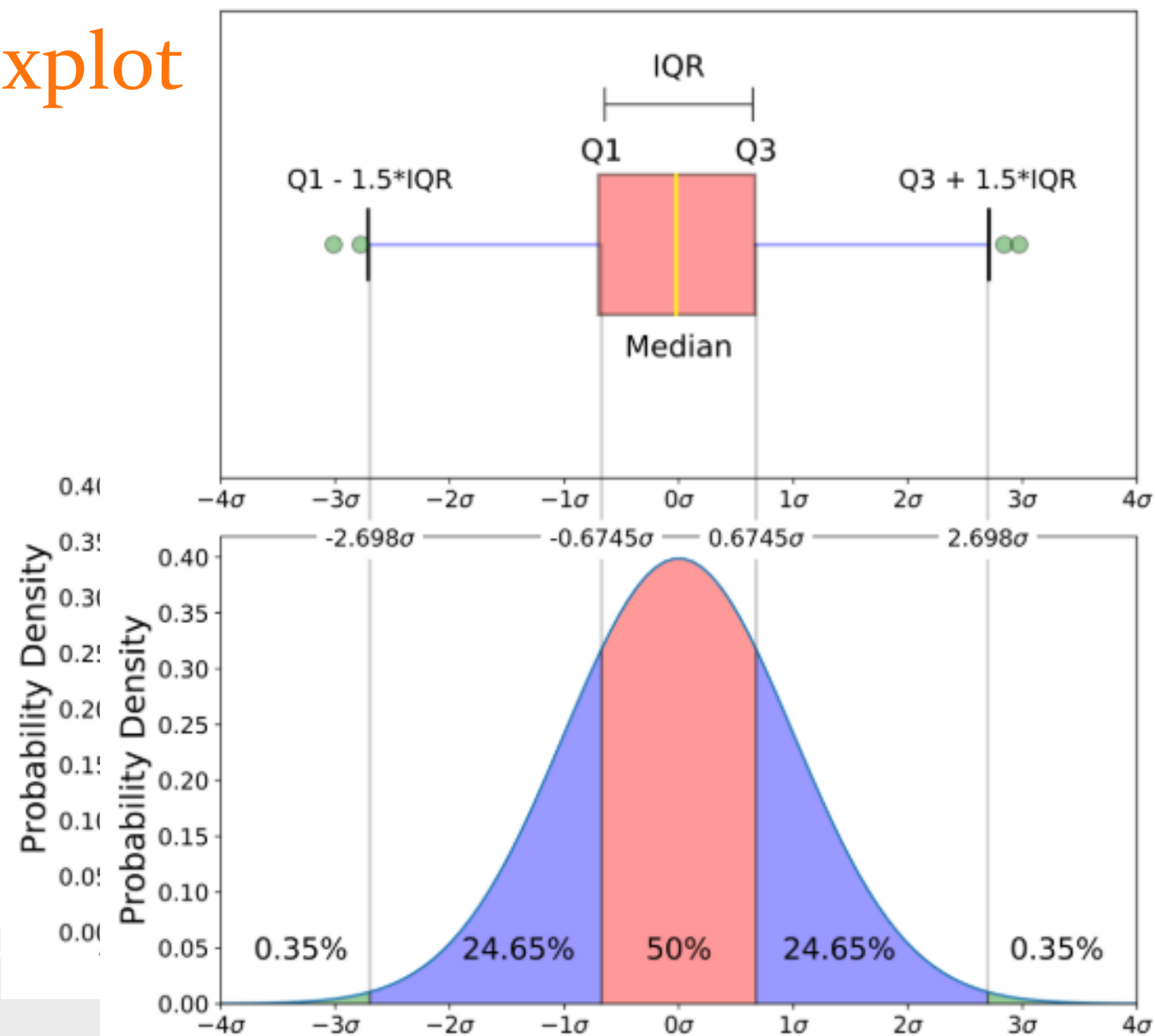
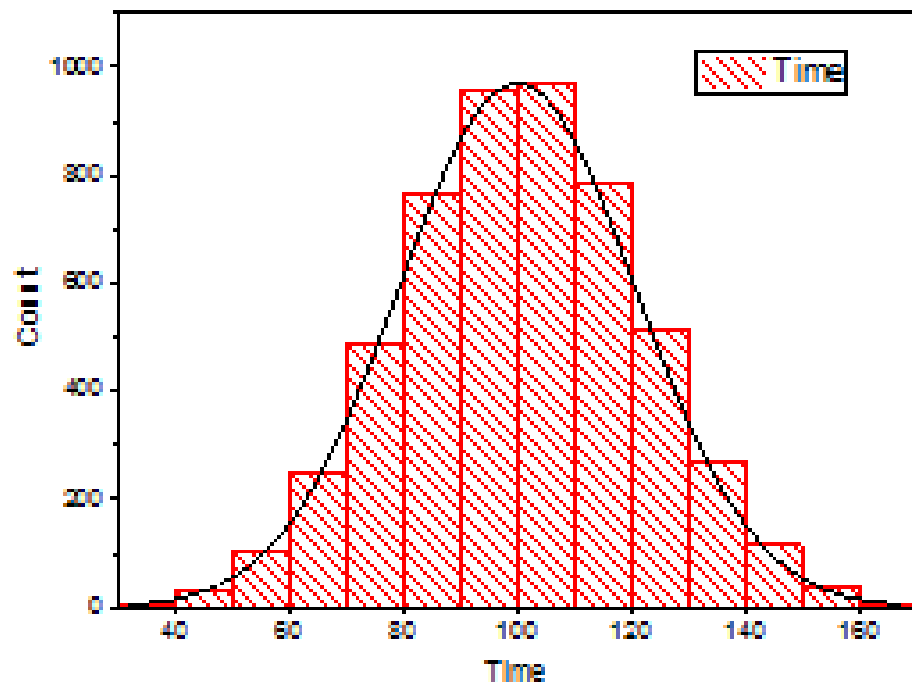
- Too small: no grouping/aggregation effect
- Too big: no differentiation



```
#(bin) size does matter!
```

```
plt.hist([winedata['proline']], bins= 5, alpha=1, label='bin size = 5')  
plt.hist([winedata['proline']], bins= 20, alpha=0.8, label='bin size = 20')  
plt.hist([winedata['proline']], bins= 60, alpha=0.5, label='bin size 60')  
plt.legend(loc='upper right')  
plt.show()
```

## Basic Charts – 2. Boxplot



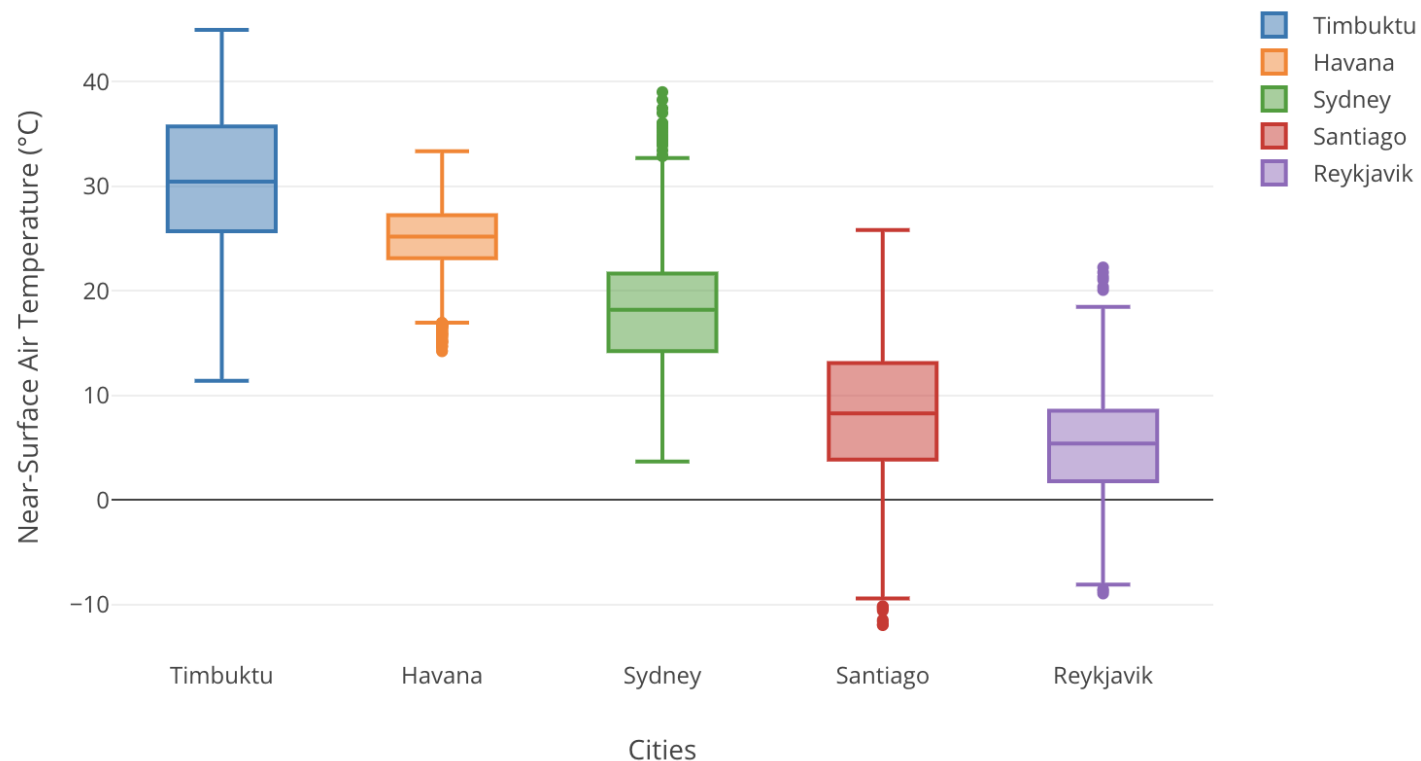


## Basic Charts – 2. Boxplot(s)

Boxplots can show a lot of information in 1 graph

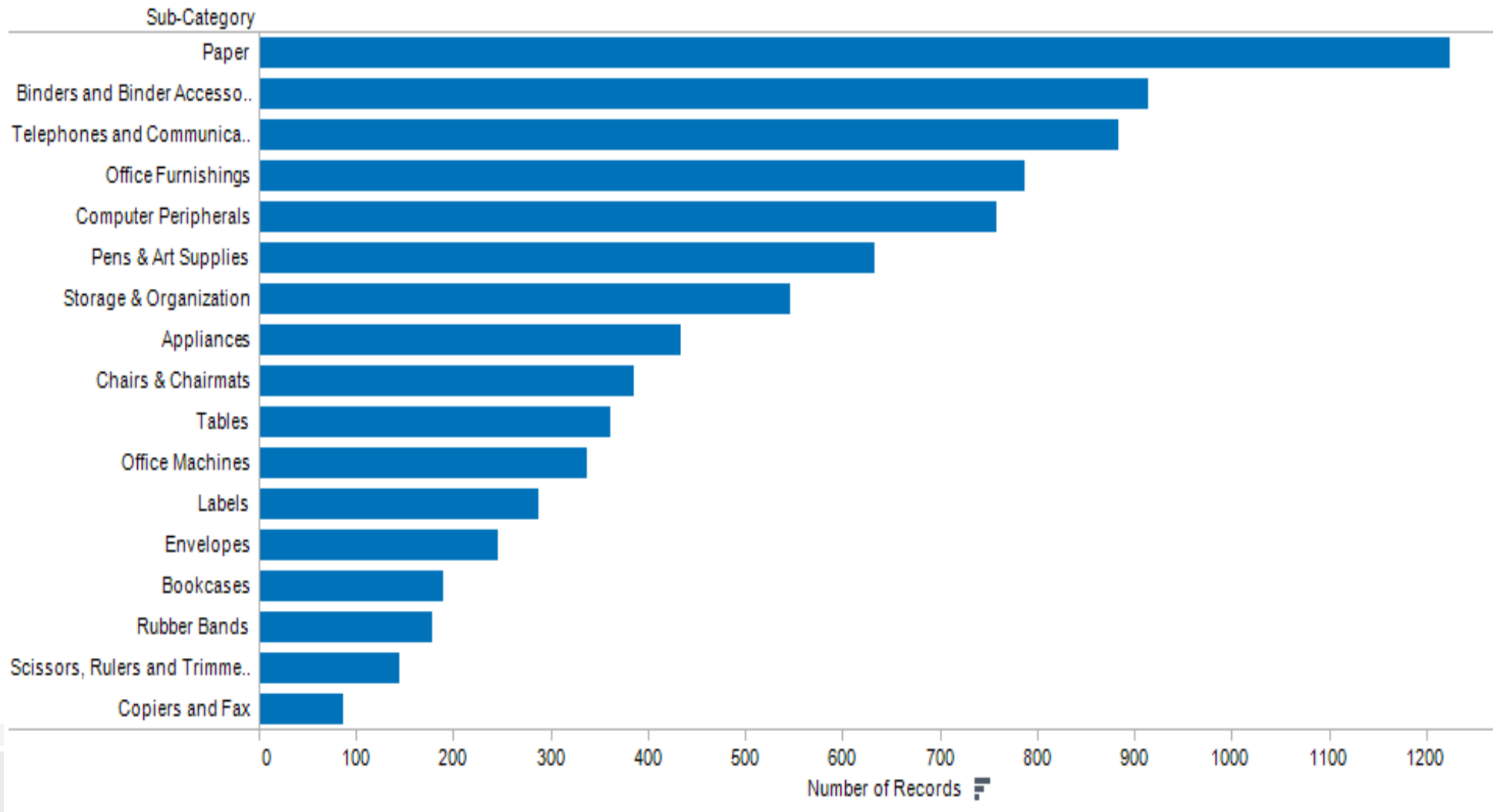
Make sure you can explain what you see (and what strikes you)

Box plots

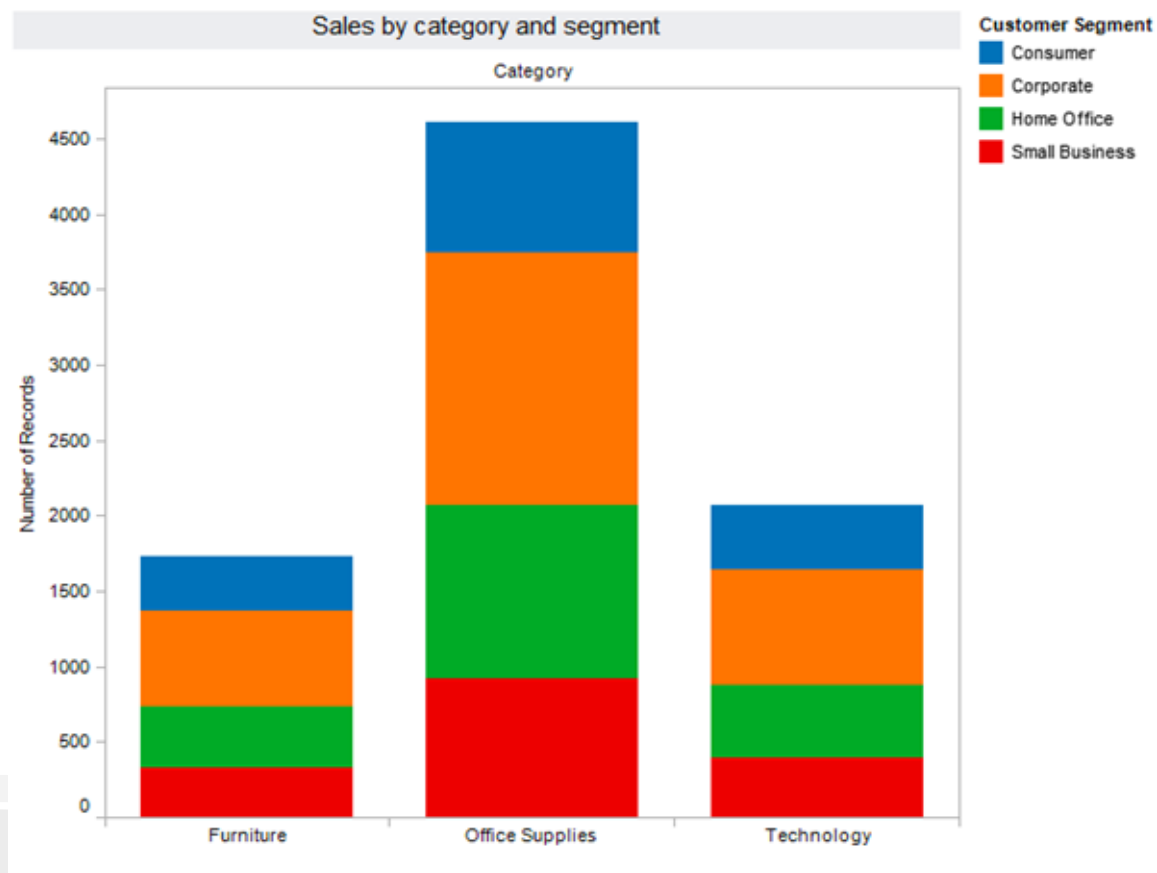


1 category, 1 numerical

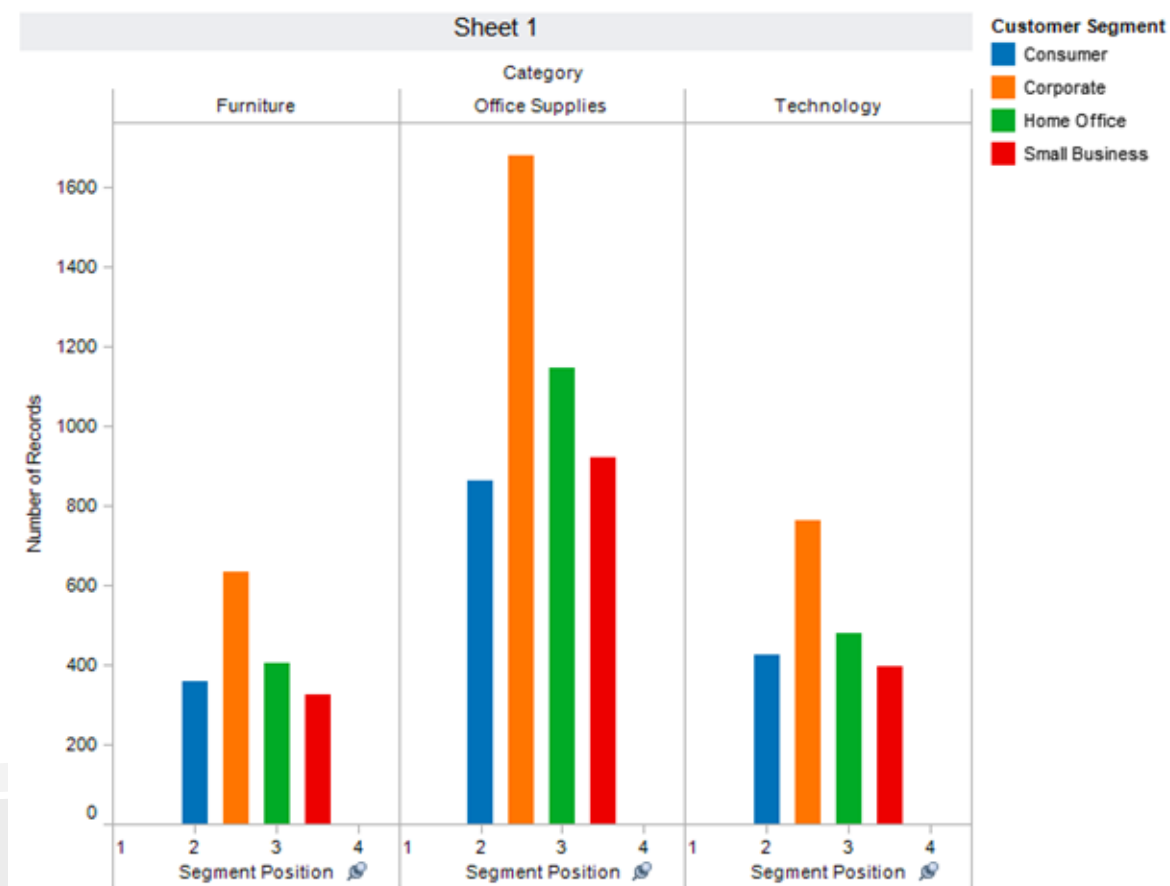
## Basic Charts – 3. Bar Chart



... stacked and grouped



2 categorical keys, 1 quant. value

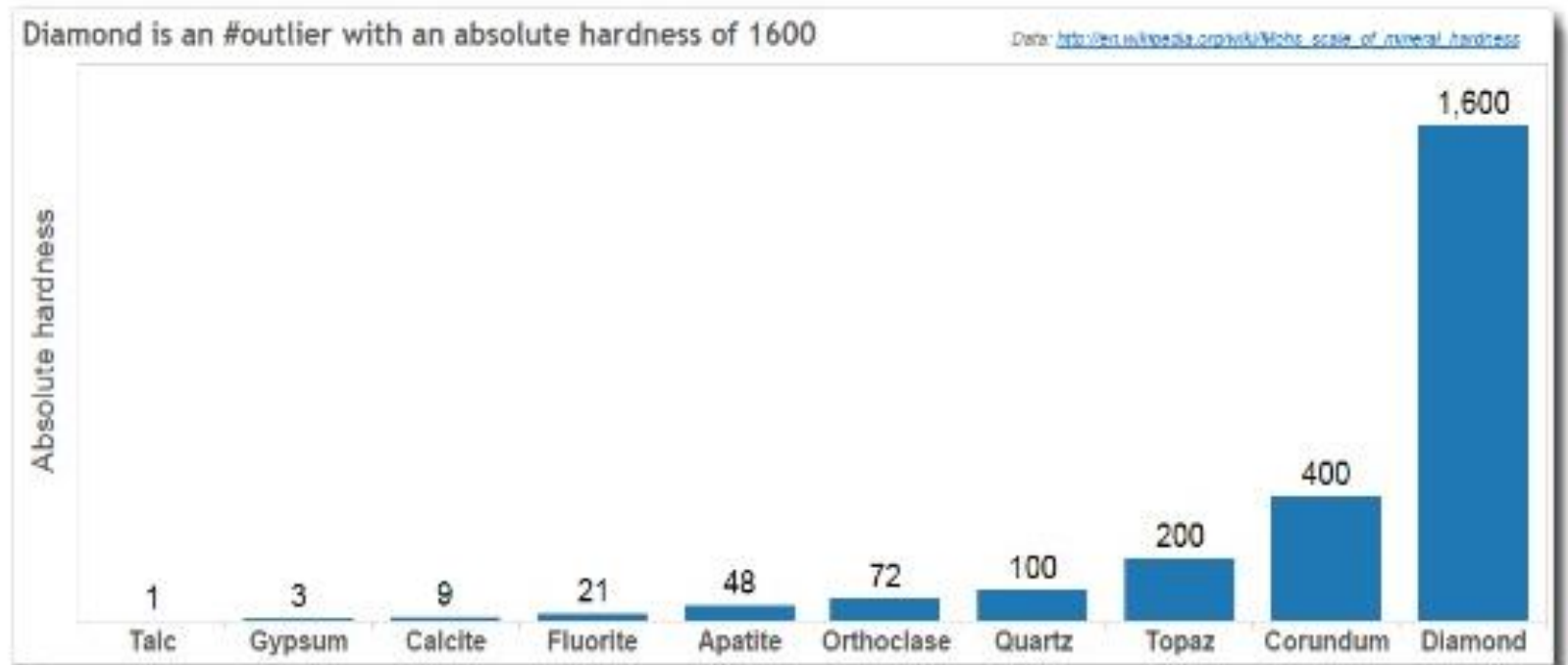


# ... insights through Outliers and Comparisons

Insights about how **special** something is

or **larger/smaller** than something else

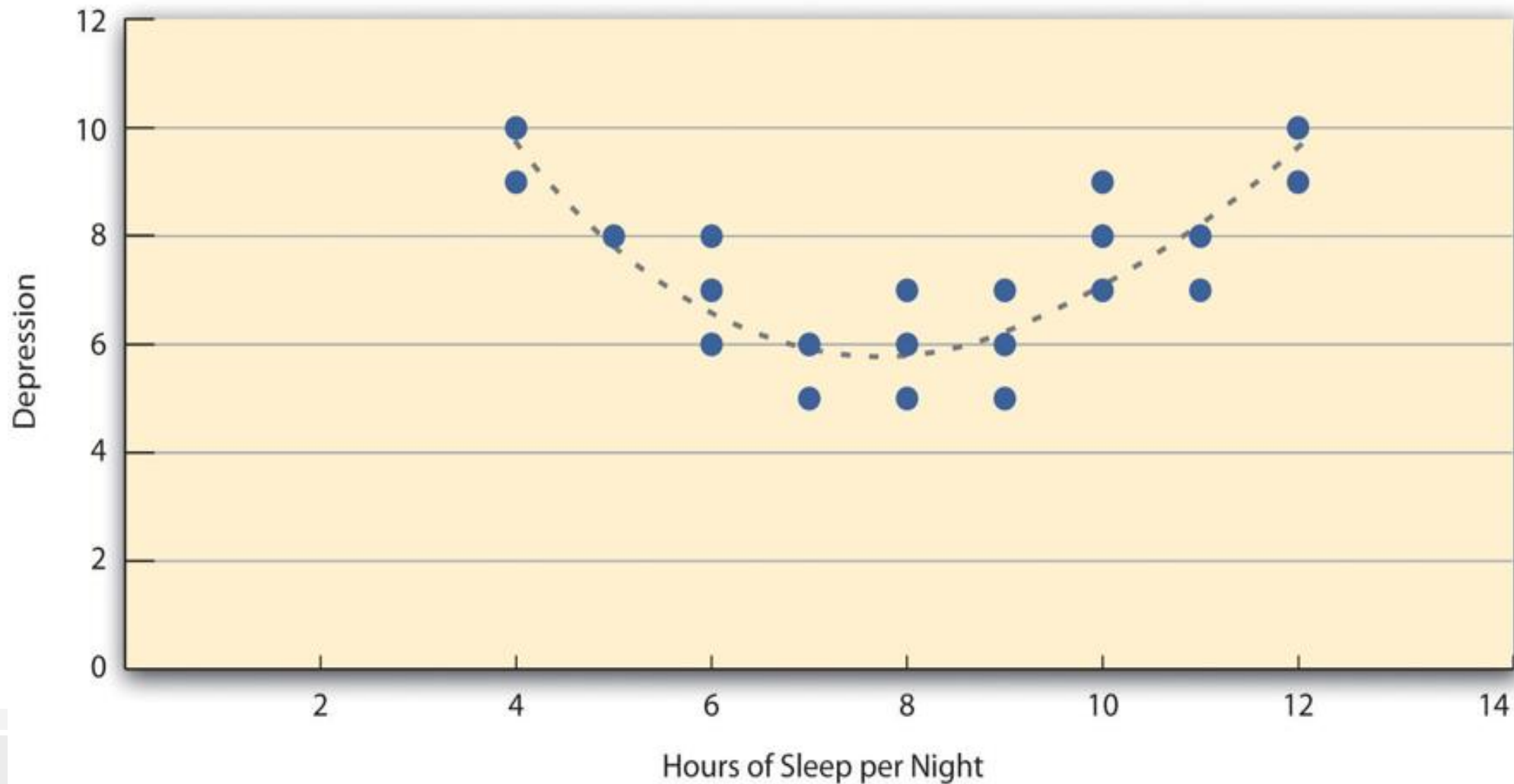
Which mineral is an #outlier? Why, Diamond: 1,600 on Moh's Scale of Hardness.





2 numerical variables

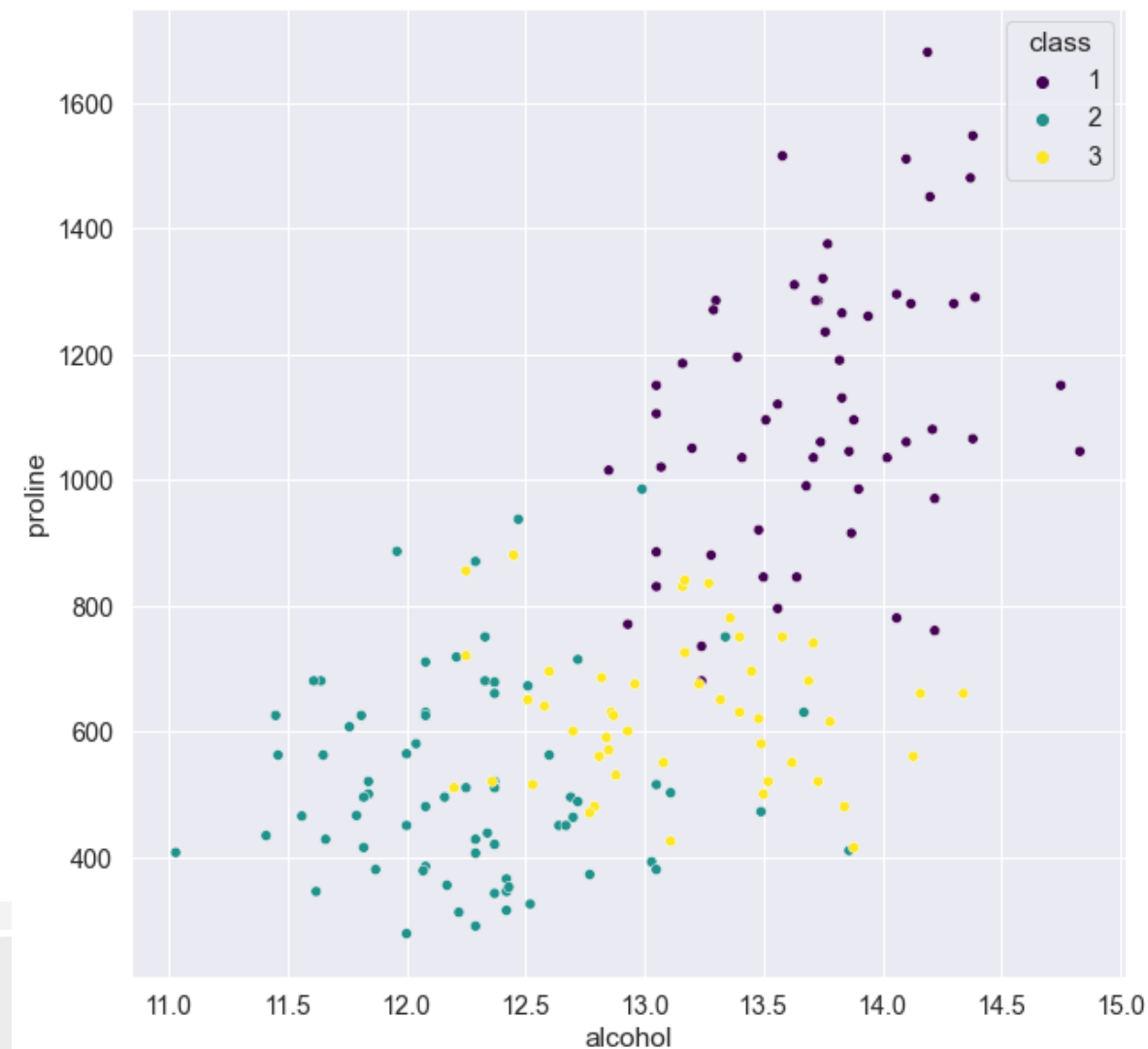
## Basic Charts – 4. Scatter plot



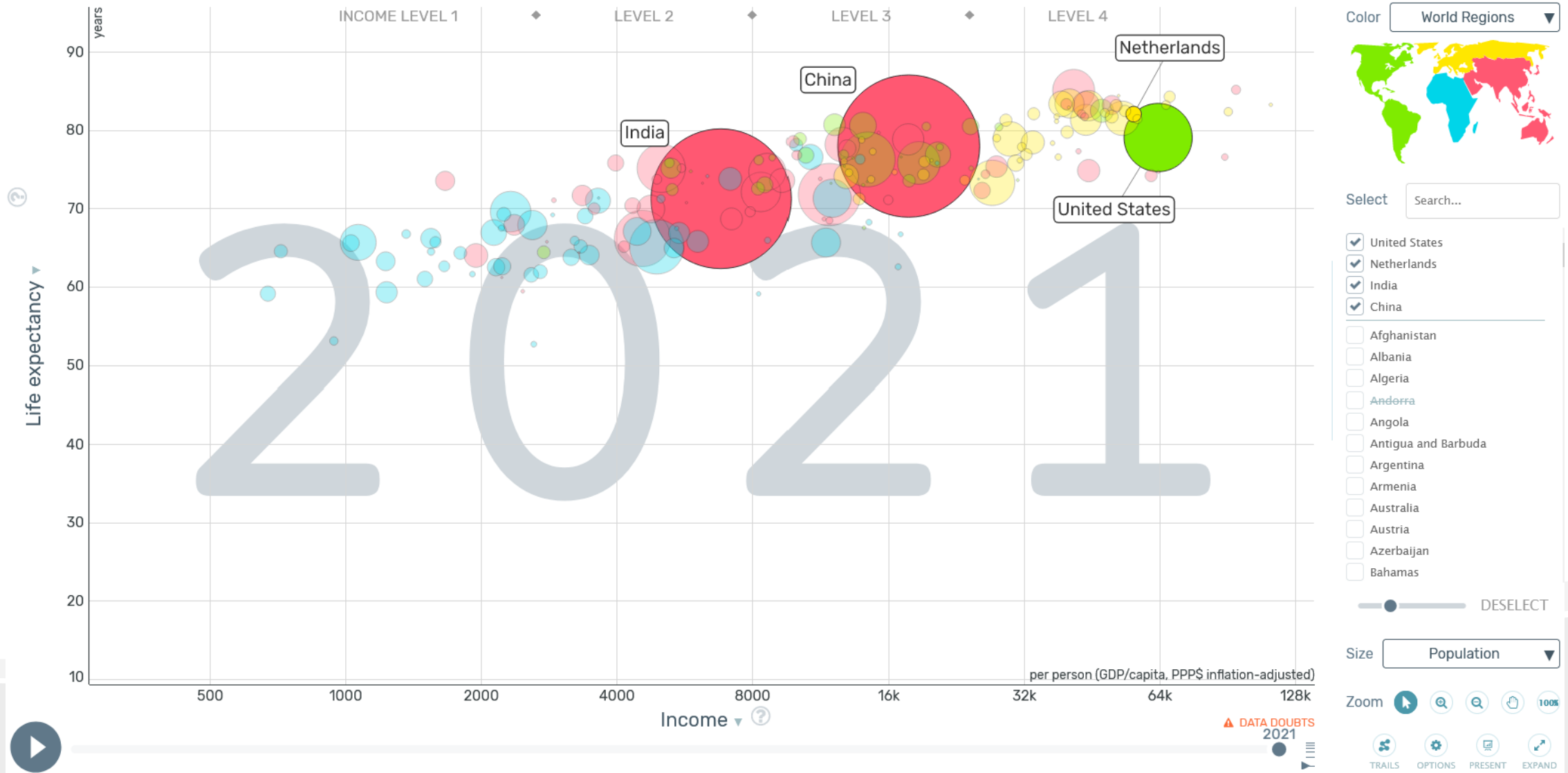
## Basic Charts – 4. Scatter plot

```
#plot 2 features as data points (a 3rd feature as hue)
```

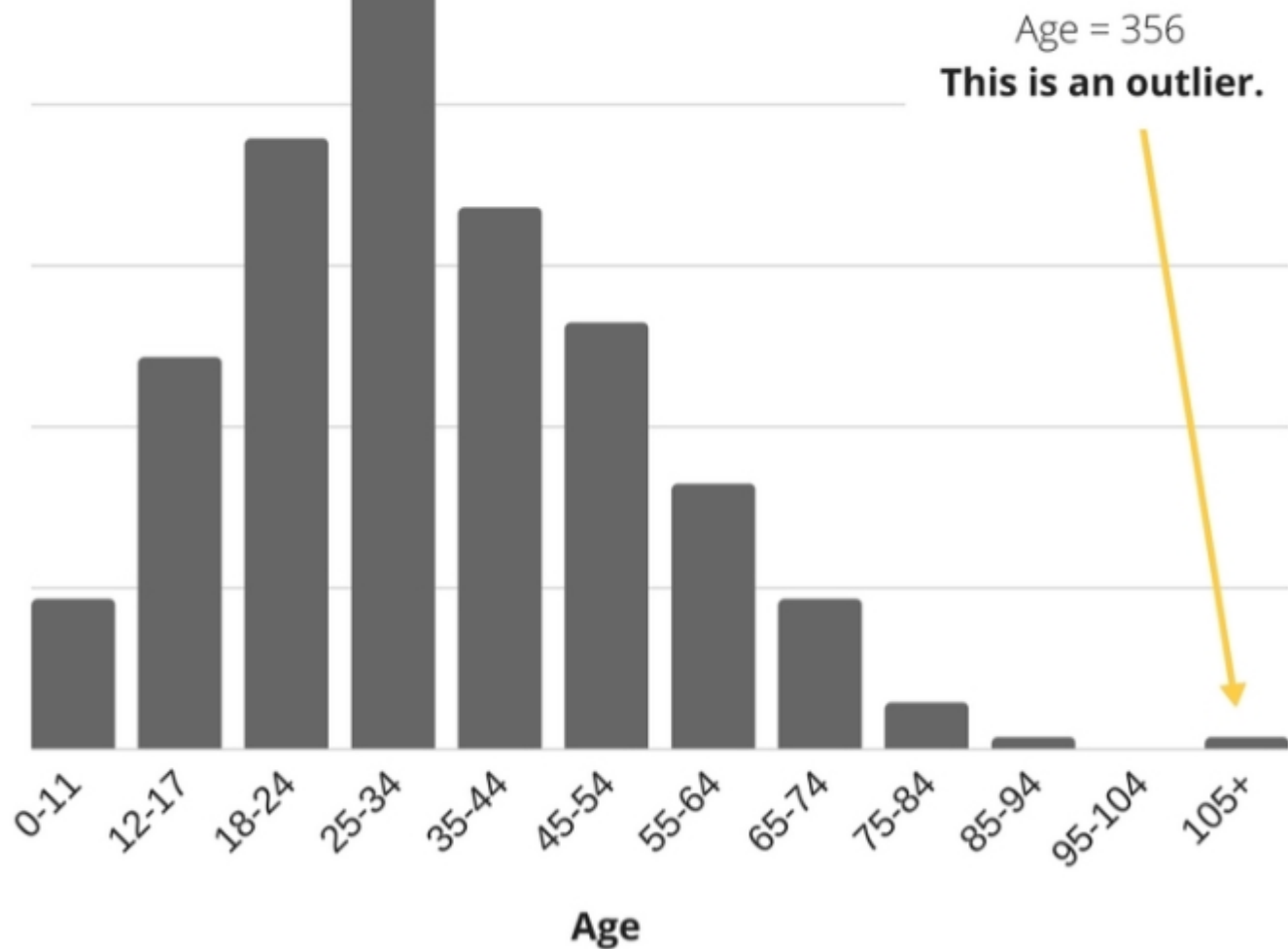
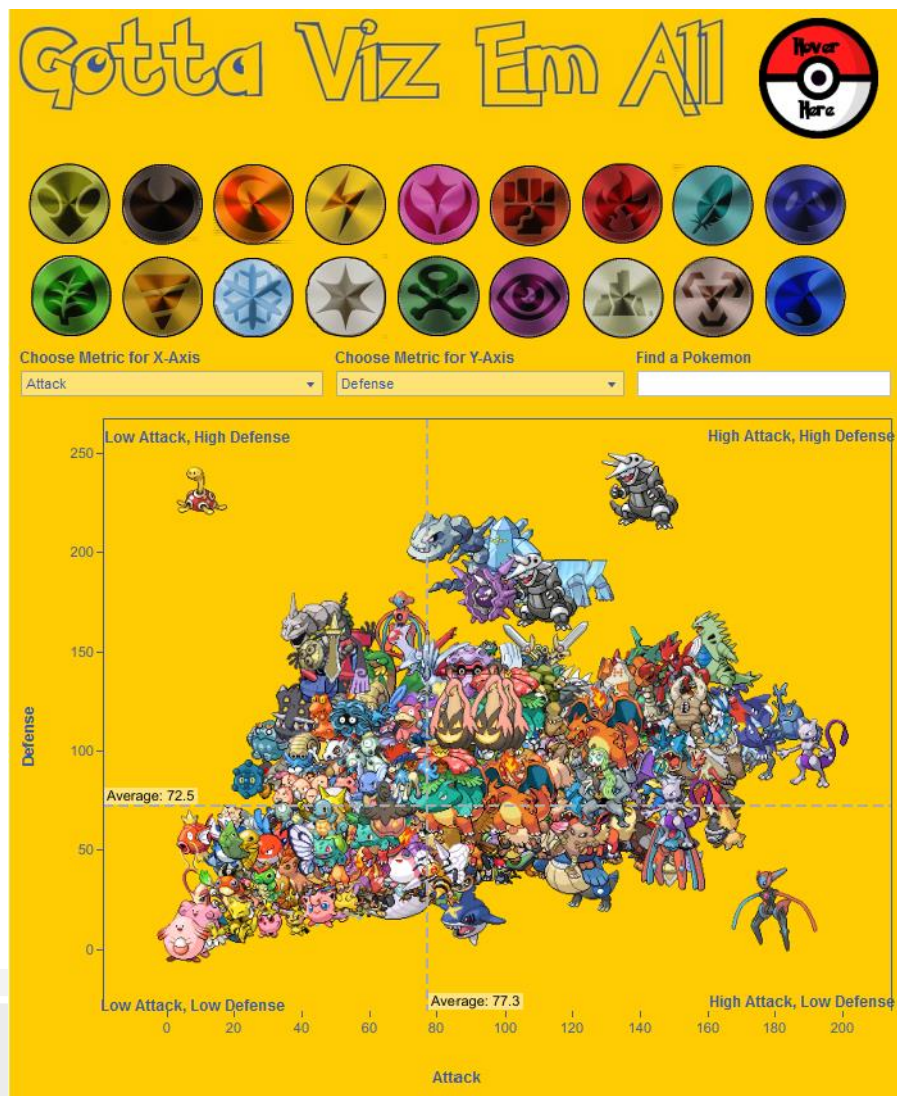
```
sns.scatterplot(data=winedata, x='alcohol',  
y='proline', hue='class', palette='viridis')
```



## Basic Charts – 4. Scatter plot

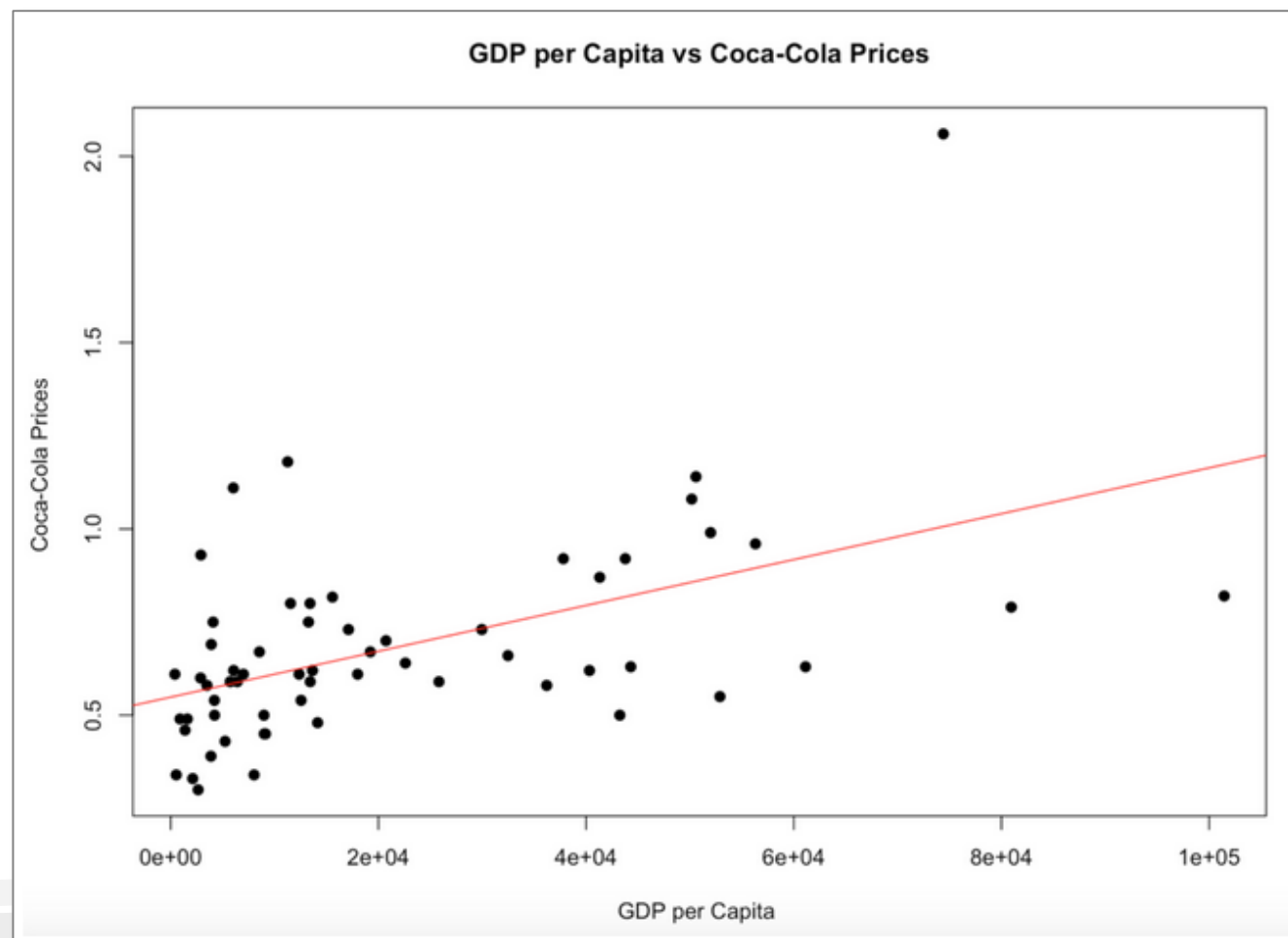


# ... insights through Outliers



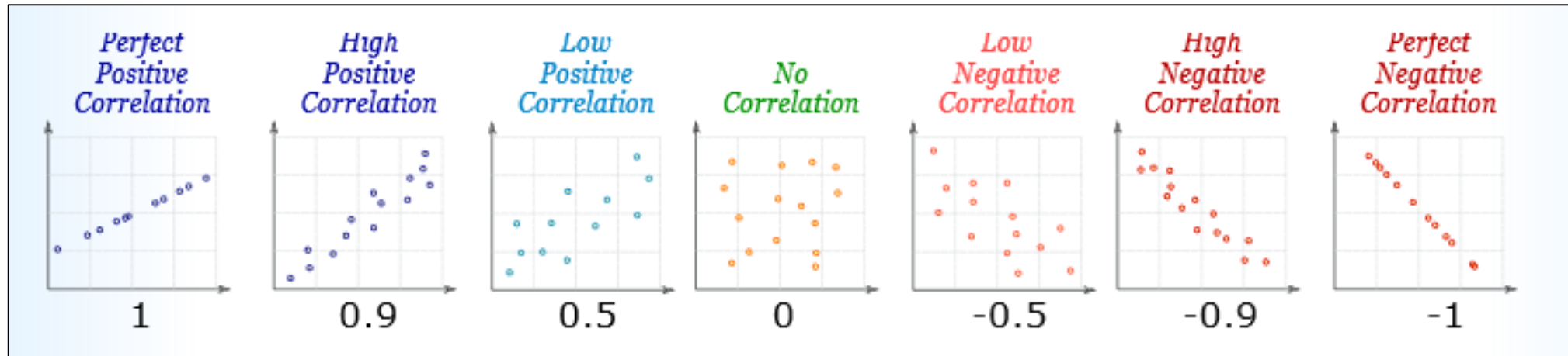
## ... insights through Correlations

Insights about how strongly numerical features are **linked** and therefore possibly influence each other.



## ... insights through Correlations

also as a numerical value, most commonly used is the **correlation coefficient R**



Try this fun game: <http://guessthecorrelation.com/>



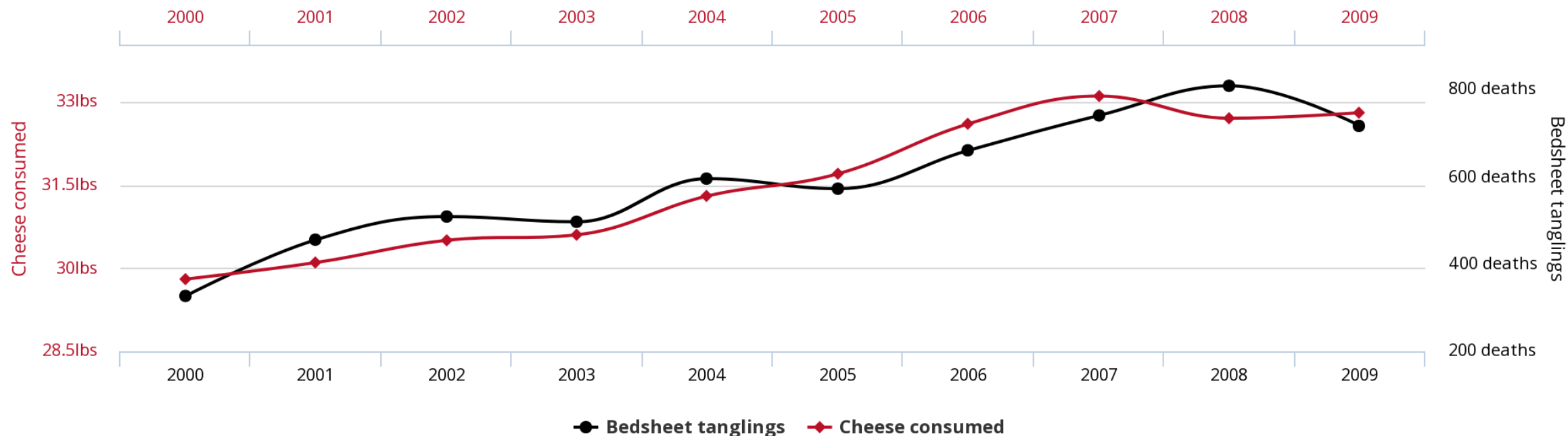


## ... insights through **Correlations**

**Per capita cheese consumption**

correlates with

**Number of people who died by becoming tangled in their bedsheets**



tylervigen.com

**Correlation does not always mean a cause-effect relation!**

See many more interesting examples here: <http://www.tylervigen.com/spurious-correlations>



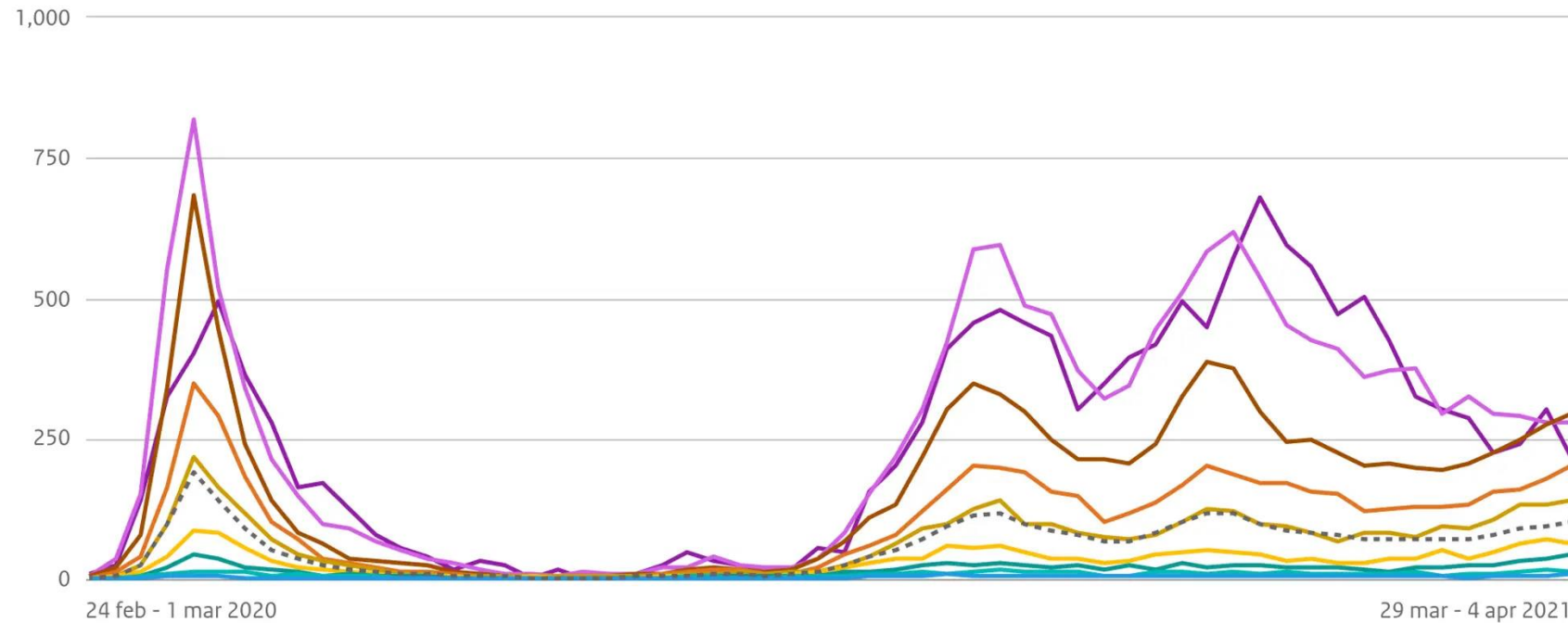
## Basic Charts – 5. Line chart

2 numerical variables, usually  $x = \text{time}$

Number of hospital admissions per age group over time

— 0-19 — 20-29 — 30-39 — 40-49 — 50-59 — 60-69 — 70-79 — 80-89 — 90+ .... All age groups

Total per week per 1,000,000 people



Source: NICE via RIVM



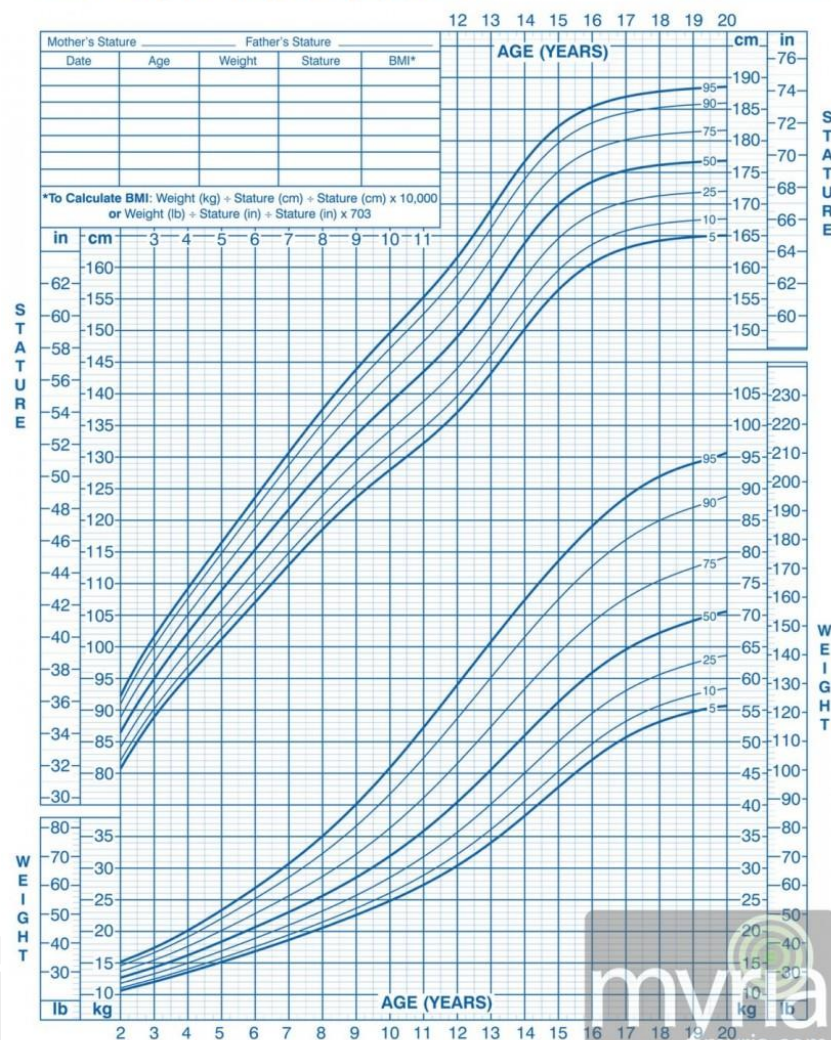
# Basic Charts – 5. Line chart

2 to 20 years: Boys

Stature-for-age and Weight-for-age percentiles

NAME \_\_\_\_\_

RECORD # \_\_\_\_\_



2 numerical variables, usually  $X = \text{time}$

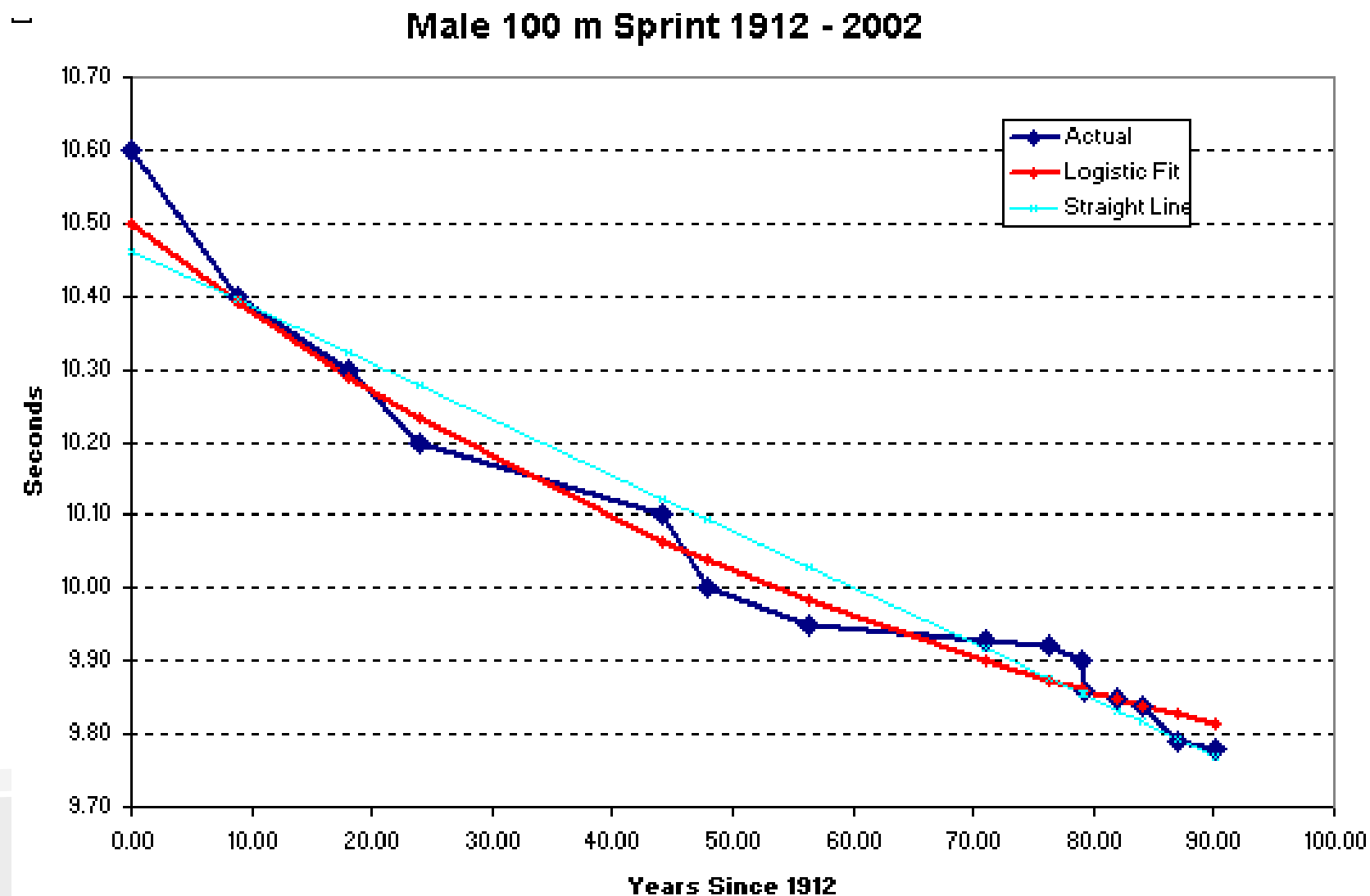
Life is like an  
ECG diagram



If it goes smoothly,  
you are dead.



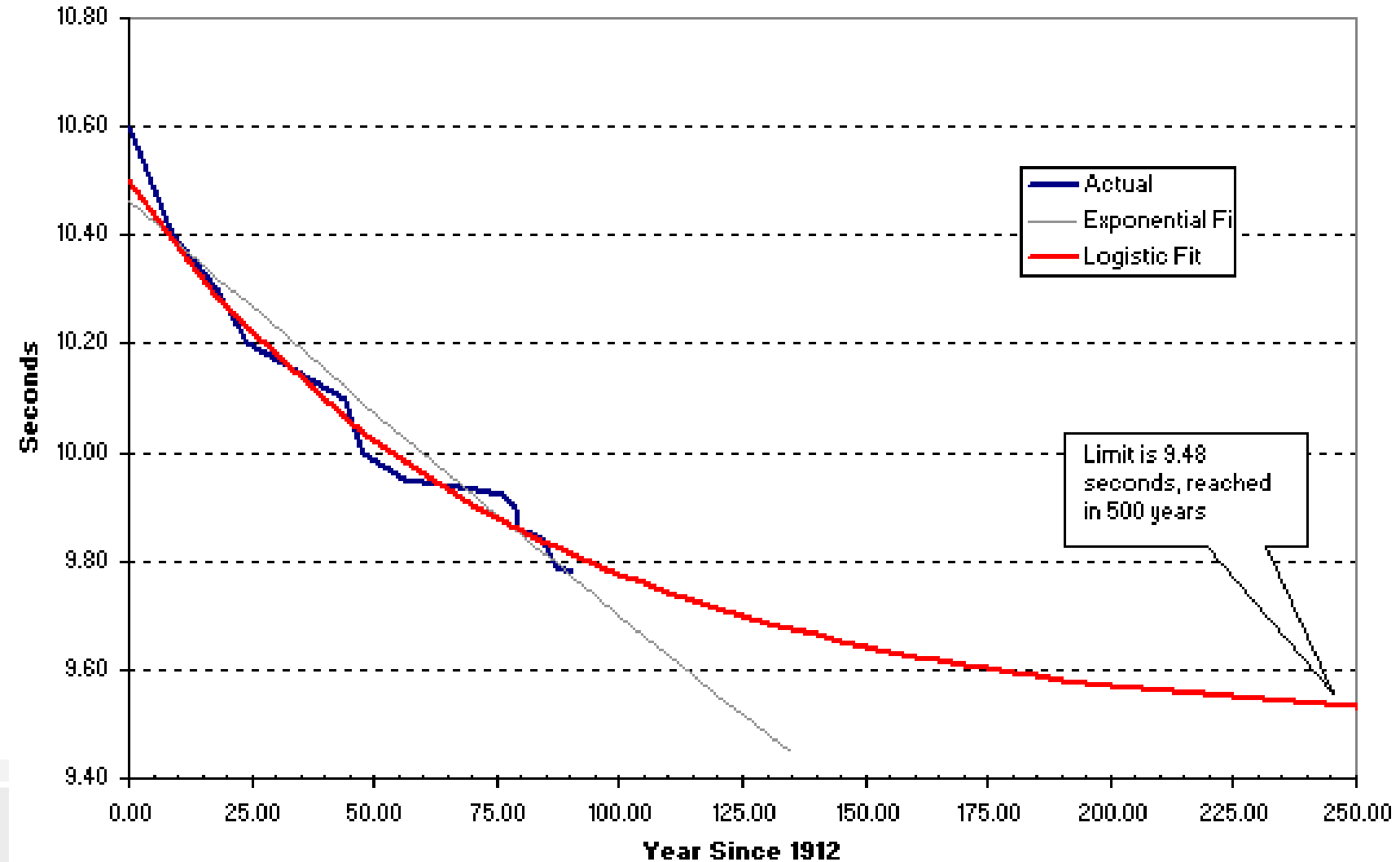
# ... insights through Time: Trends



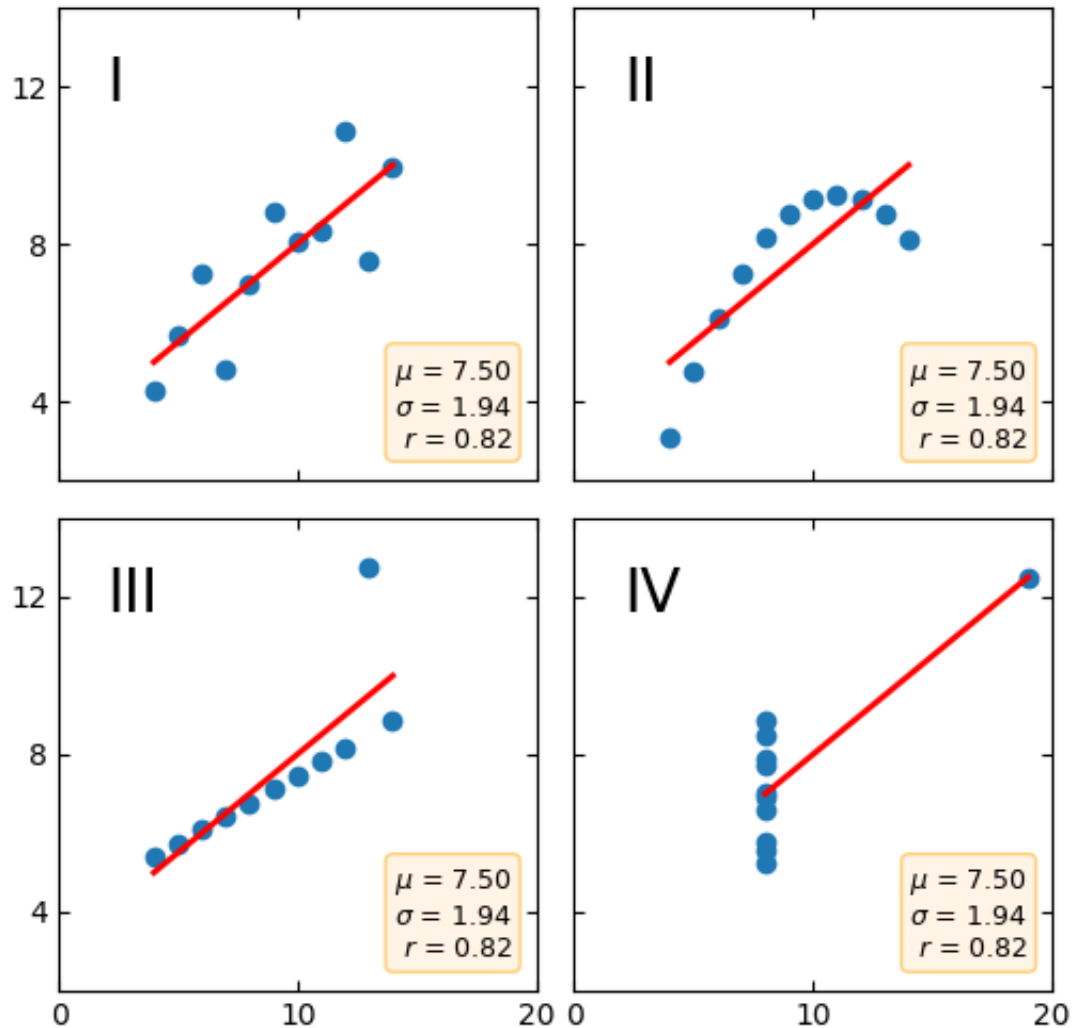


# ... insights through Time: Trends

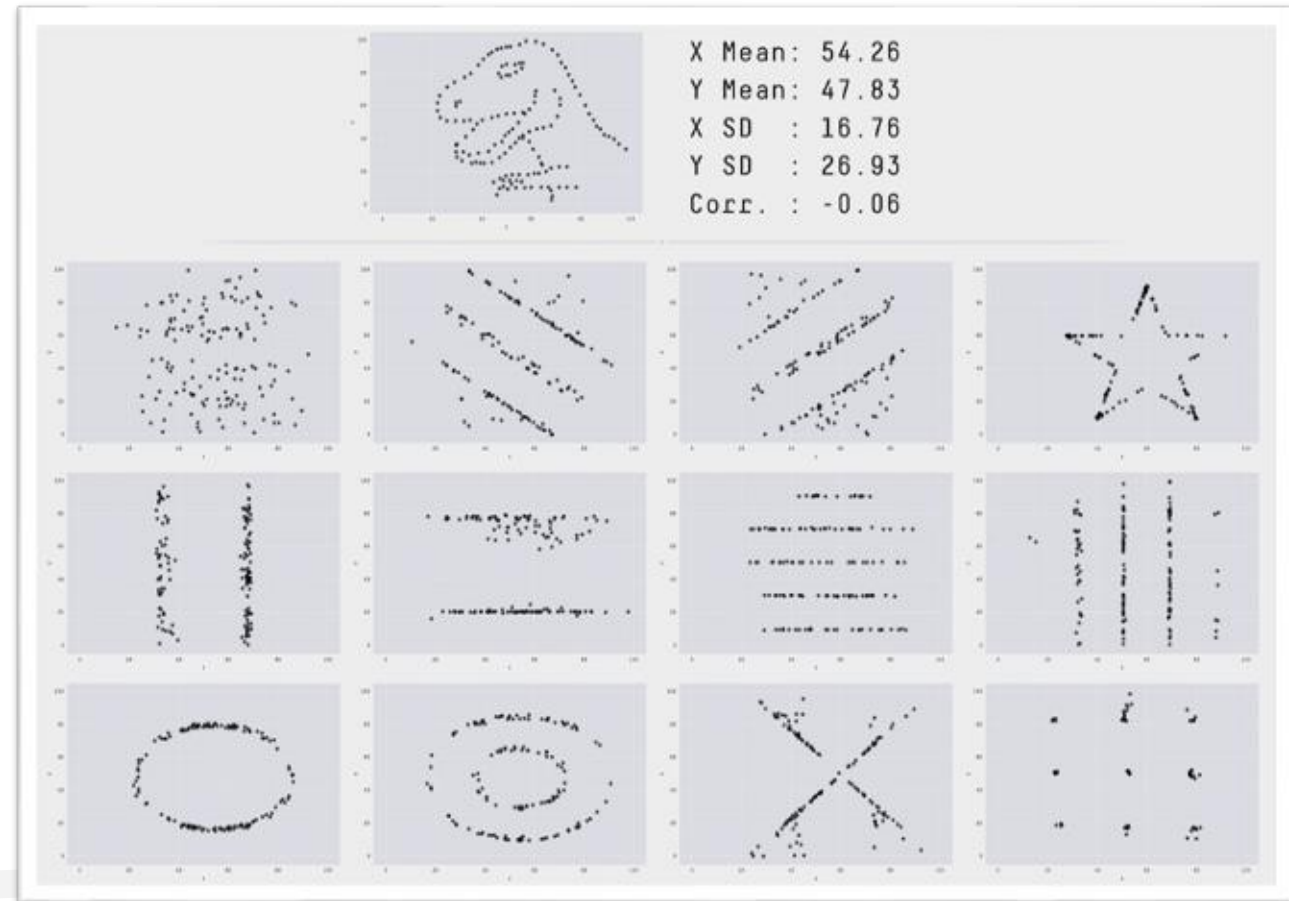
Male 100 m Sprint Prediction



# Why is data visualization important?



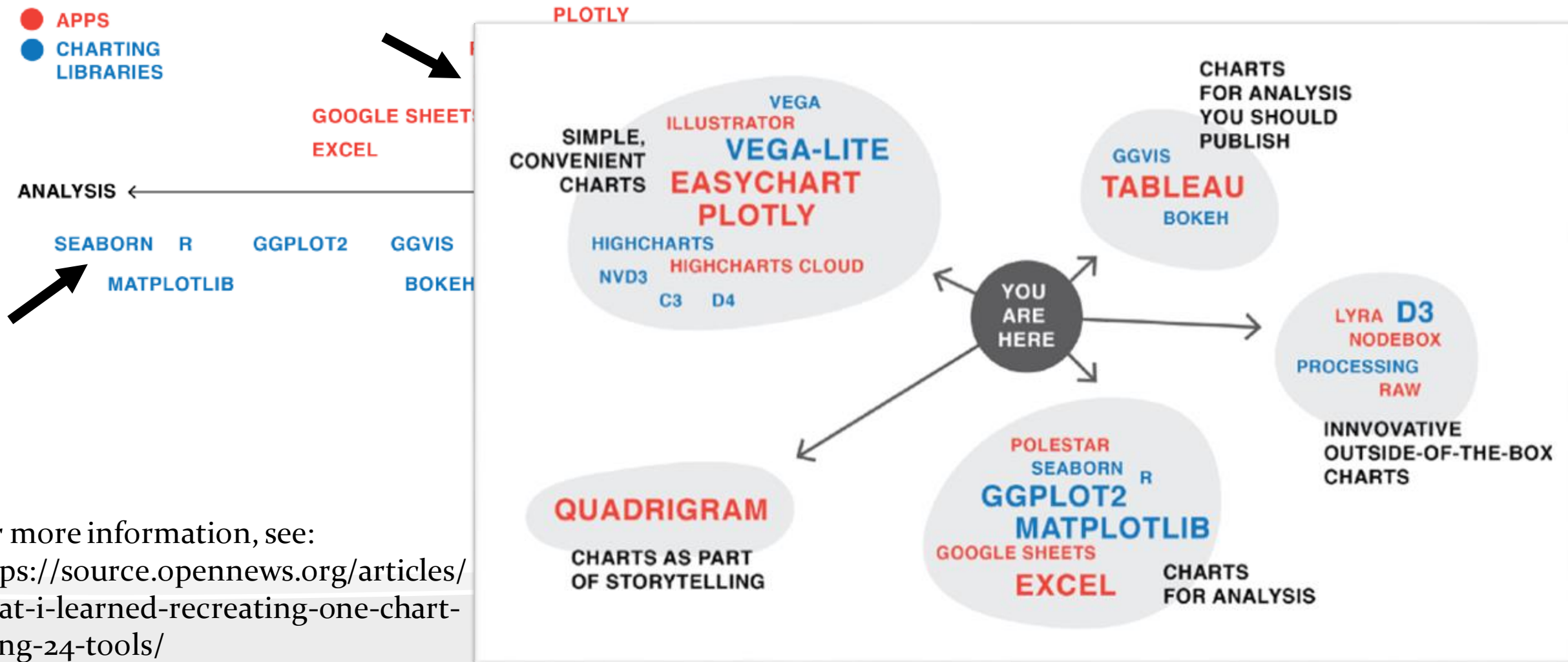
Source: [https://matplotlib.org/3.2.1/gallery/specialty\\_plots/anscombe.html](https://matplotlib.org/3.2.1/gallery/specialty_plots/anscombe.html)



Source: <https://www.autodesk.com/research/publications/same-stats-different-graphs>



# Popular visualization tools





## Tableau

and/or

## Python

- Simple for the exploration phase.
  - Allows also for publishing online interactive visualizations with storytelling elements.
  - Allows for connecting with Big Data sources such as Spark and MongoDB.
- Complete control and reproducibility of datasets
  - Active community
  - Integration with complex libraries like machine learning
  - Can serve web based visualizations (via, for instance, Dash by Plotly)



# Tools for data exploration



- reproducibility of steps is key; that dismisses traditional Excel
- 2 classes of tools
  - end-user tools (GUI based), e.g. Tableau, KNIME, Power BI
  - scripting tools: SQL, Python, R, Julia, JavaScript, ...

## **Typical problems you can run into when using data prep tools**

- data import decisions are usually based on the first  $n$  ( $n < 1000$ ) rows
- dates, dates, dates
- discerning between missing values and empty values/strings
- adaptors for specific data sources / data storage solutions (HTML, XML, databases, filesystems, applications)
- level of documentability: Notebooks shine in this respect



# pandas - what should you know

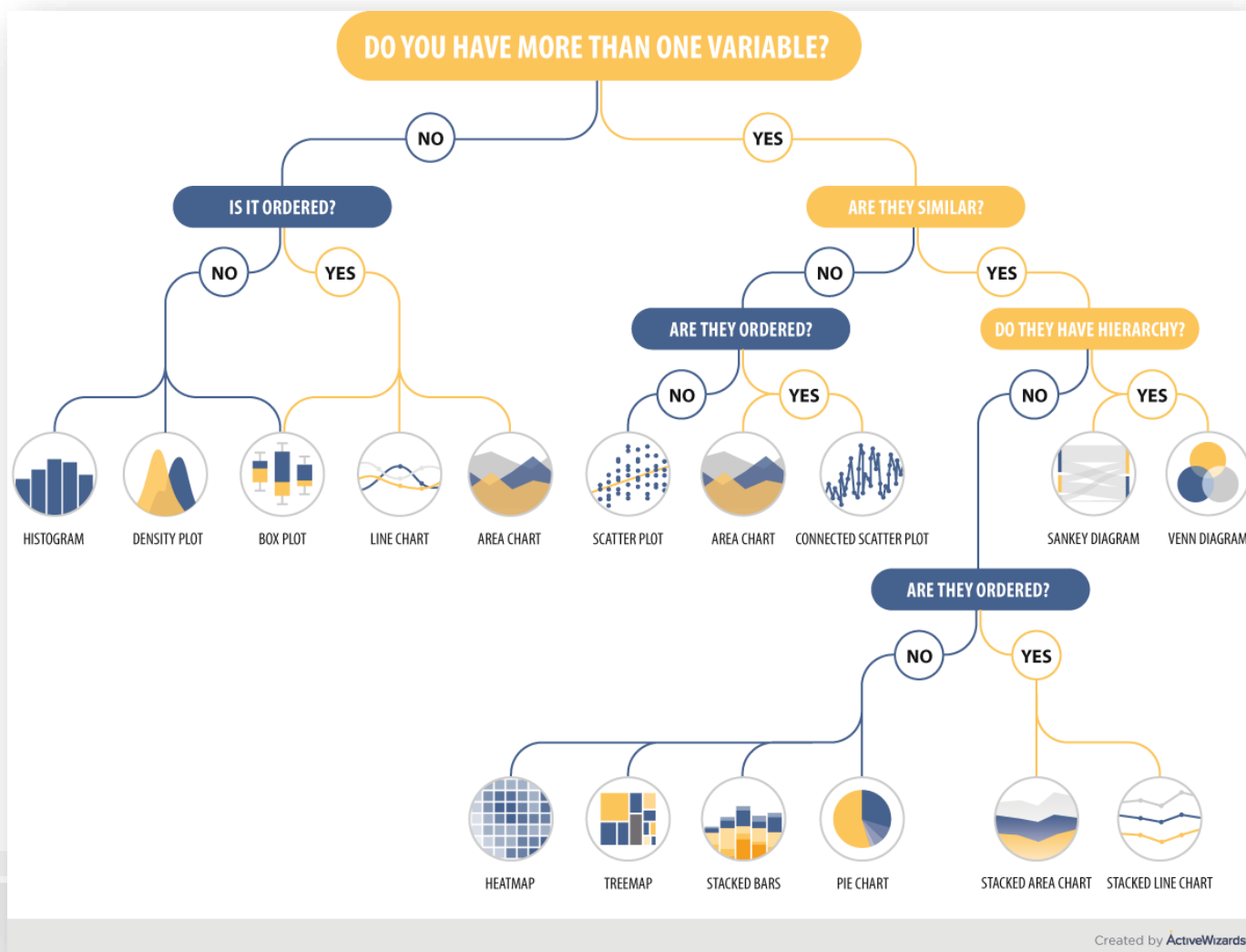
- inspired by R data frames, built on NumPy (and Matplotlib)
- there's more than one way you can do it. That makes using it convenient and difficult at the same time
- no smooth learning curve, although you can get initial work done with a small subset
- it takes quite some time (and necessary mileage) to become a pandas expert
- Numpy is stable, pandas is improved and extended at a rapid rate
- pandas: almost anything regarding data preparation can be done with it
- DataFrame size is bound to available, physical memory:
  - no fancy storage solutions in core pandas
  - a lot of fancy stuff is recently released or in the making (parquet, Arrow)



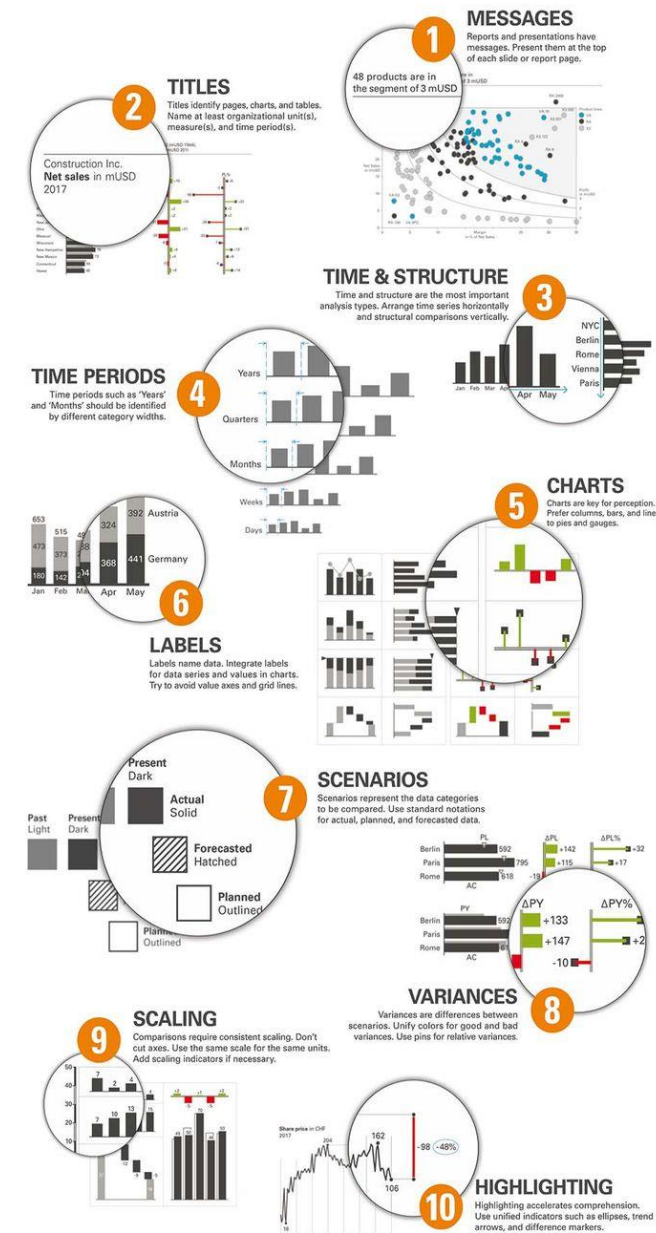
# pandas - basic (necessary) subset

- importing data
  - read from an Excel or csv file
  - read from a Web API
- transforming data
  - handling missing values (NaN, None)
  - mutating: adding columns based on values in other columns
  - subsetting and filtering
  - grouping and aggregating
- presenting data
  - using pandas' built in plot functionality
  - using Seaborn
  - using Matplotlib (only limited subset)

# Useful standards, guidelines



Created by ActiveWizards





## Useful resources

- A comprehensive [Chart Guide](#) poster (+ [large hi-res](#) version)
- <https://elitedatascience.com/python-seaborn-tutorial> - short tutorial on the **Seaborn** library, a fast way of making beautiful charts in Jupyter Notebooks
- <http://seaborn.pydata.org/examples/index.html> - fancy chart examples made with Seaborn, with code samples





# Go prepare your ingredients!

Check the exercise(s) on Canvas

