

# Data Science and Machine Learning in Python

Stephan Weyers

# Topics covered in the online lectures

## Part 1: Data Science

|   | Date                 | Topics covered   |
|---|----------------------|--|
| 1 | Apr 13 <sup>th</sup> | Course introduction<br>Data Science motivation<br>How to use Jupyter Notebook<br>Python types and lists<br>Loops, if/else, functions   |
| 2 | Apr 20 <sup>th</sup> | Python tuples, lists, dictionaries<br>Functions<br>Numpy basics, operations<br>Image processing  |
| 3 | Apr 27 <sup>th</sup> | Pandas Series, DataFrame<br>Pandas basic operations<br>Import/export files   |
| 4 | May 4 <sup>th</sup>  | Principles of data visualization<br>Data cleaning and preparation<br>Join, combine and reshape data  |
| 5 | May 11 <sup>th</sup> | <b>Volkswahl Bund dataset<br/>Johanna Dahlbeck<br/>14:30-15:30 German time</b><br>Data visualization in Python<br>How to write Data Science reports<br>Data aggregation and grouping |

## Part 2: Machine Learning

|    | Date                 | Topics covered   |
|----|----------------------|--|
| 6  | Jun 1 <sup>st</sup>  | Introduction to supervised learning<br>Classification and regression<br>scikit-learn<br>k-Nearest Neighbors<br>Linear Models   |
| 7  | Jun 8 <sup>th</sup>  | Decision trees<br>Random forests and gradient boosting<br>Support vector machines<br>Neural networks                           |
| 8  | Jun 15 <sup>th</sup> | Introduction to unsupervised learning<br>Preprocessing and scaling<br>Dimensionality reduction<br>Principal component analysis |
| 9  | Jun 22 <sup>nd</sup> | k-means clustering<br>Hierarchical clustering<br>DBSCAN  |
| 10 | Jun 29 <sup>th</sup> | Representing data<br>Engineering features  |
| 11 | Jul 6 <sup>th</sup>  | Model evaluation and improvement<br>Text data analysis   |

# Agenda for online lecture 4

| Session     | Topic                        | Mode                         | Materials used      | Time (min) |
|-------------|------------------------------|------------------------------|---------------------|------------|
| 14:30-16:00 | Tutorial Marius Meiners      | Q&A                          |                     | 10         |
|             | Data preparation / cleaning  | Lecture / Q&A                | Lecture slides      | 20         |
|             | Inventory step 1 & 2         | Lecture / Q&A                | Inventory Excel     | 10         |
|             | Inventory step 1             | Lecture / Q&A                | Lecture 04 notebook | 10         |
|             | Inventory step 2             | Team work in break-out rooms | Lecture 04 notebook | 25         |
|             | Inventory step 2             | Presentation of results      | Lecture 04 notebook | 10         |
| 16:10-17:40 | Importance of Visualization  | Lecture / Q&A                | Lecture slides      | 5          |
|             | Visualization: Trustworthy   | Discussion in main room      | Lecture slides      | 15         |
|             | Visualization: Accessibility | Discussion in main room      | Lecture slides      | 15         |
|             | Visualization: Elegance      | Discussion in main room      | Lecture slides      | 15         |
|             | Inventory step 3 & 4         | Lecture / Q&A                | Inventory Excel     | 10         |
|             | Inventory step 3             | Lecture / Q&A                | Lecture 04 notebook | 10         |
|             | Inventory step 4             | Team work in break-out rooms | Lecture 04 notebook | 15         |
| 17:50-19:20 | Inventory step 5 & 6 & 7     | Lecture / Q&A                | Inventory Excel     | 15         |
|             | Inventory step 5             | Lecture / Q&A                | Lecture 04 notebook | 10         |
|             | Inventory step 4 & 6 & 7     | Team work in break-out rooms | Lecture 04 notebook | 40         |
|             | Inventory step 4 & 6 & 7     | Presentation of results      | Lecture 04 notebook | 10         |
|             | Questions W02                | Q&A                          |                     | 20         |
|             | Team grouping, deadlines     | Lecture / Q&A                |                     | 5          |

|                        |             | Description  | Examples  |
|------------------------|-------------|--|---|
| Selection of cases     | Validity    | Restriction to valid cases                               | <ul style="list-style-type: none"> <li>Only active customers</li> <li>No “test-buyers”</li> </ul> |
|                        | Relevance   | Restriction to relevant cases related to target          | <ul style="list-style-type: none"> <li>Only contract customers, no pre-paid customers</li> </ul>  |
|                        | Sample      | Restriction to sample in case of large dataset           | <ul style="list-style-type: none"> <li>Random sampling vs. stratified sampling</li> </ul>         |
|                        | Partition   | Split up into homogeneous segments                       | <ul style="list-style-type: none"> <li>Certain customer segments, e.g. high-value</li> </ul>      |
| Selection of variables | Quality     | Exclude variable of low quality                          | <ul style="list-style-type: none"> <li>Zero variance</li> <li>Many missing values</li> </ul>      |
|                        | Time        | Use variables time-related to target variable            | <ul style="list-style-type: none"> <li>Use only sales data of last 12 months</li> </ul>           |
|                        | Correlation | Exclude variables with high correlation among each other | <ul style="list-style-type: none"> <li>Weight and height</li> <li>Orders and returns</li> </ul>   |
|                        | Relevance   | Exclude variables with no relevance for target           | <ul style="list-style-type: none"> <li>Postal code</li> <li>Insurance ID</li> </ul>               |

| ID | Name      | Income  | Age | Gender  | Target |
|----|-----------|---------|-----|---------|--------|
| 1  | Müller    | 50.570  | 45  | M       | 1      |
| 2  | Meyer     | 19.032  | 21  | M       | 0      |
| 3  | Schneider | 43.452  |     | F       | 0      |
| 4  | Jäger     | 75.976  | 37  | Male    | 1      |
| 5  | Muller    | 50.570  | 45  | M       | J      |
| 6  | Langer    | 47.414  | 140 | Femnale | 1      |
| 7  | Schäfer   | 228.011 | 48  | F       | 0      |
| 8  | Simon     | 36.976  | 27  | M       | 1      |

## Question for discussion

- Which errors can you observe in the dataset?

## Options for missing values

- Delete (e.g. remove all customers without email-address)
- Replace (e.g. by mean values by segment)
- Accept (use missing as new category)

## Options for outliers

- Delete (e.g. remove all customers sales above 1 Mio. EUR)
- Replace (e.g. truncate all sales above 10.000 EUR to 10.000 EUR)
- Accept (maybe outliers are especially interesting for analysis)
- Binning (use categories like 0-20, 20-50, 50-100, 100-500, >500)

## Options for values with errors

- Delete (e.g. remove all customers with sales < 0 EUR)
- Replace (e.g. replace “US” by “United States”)
- Accept (if reasonable, e.g. negative sales = returns)

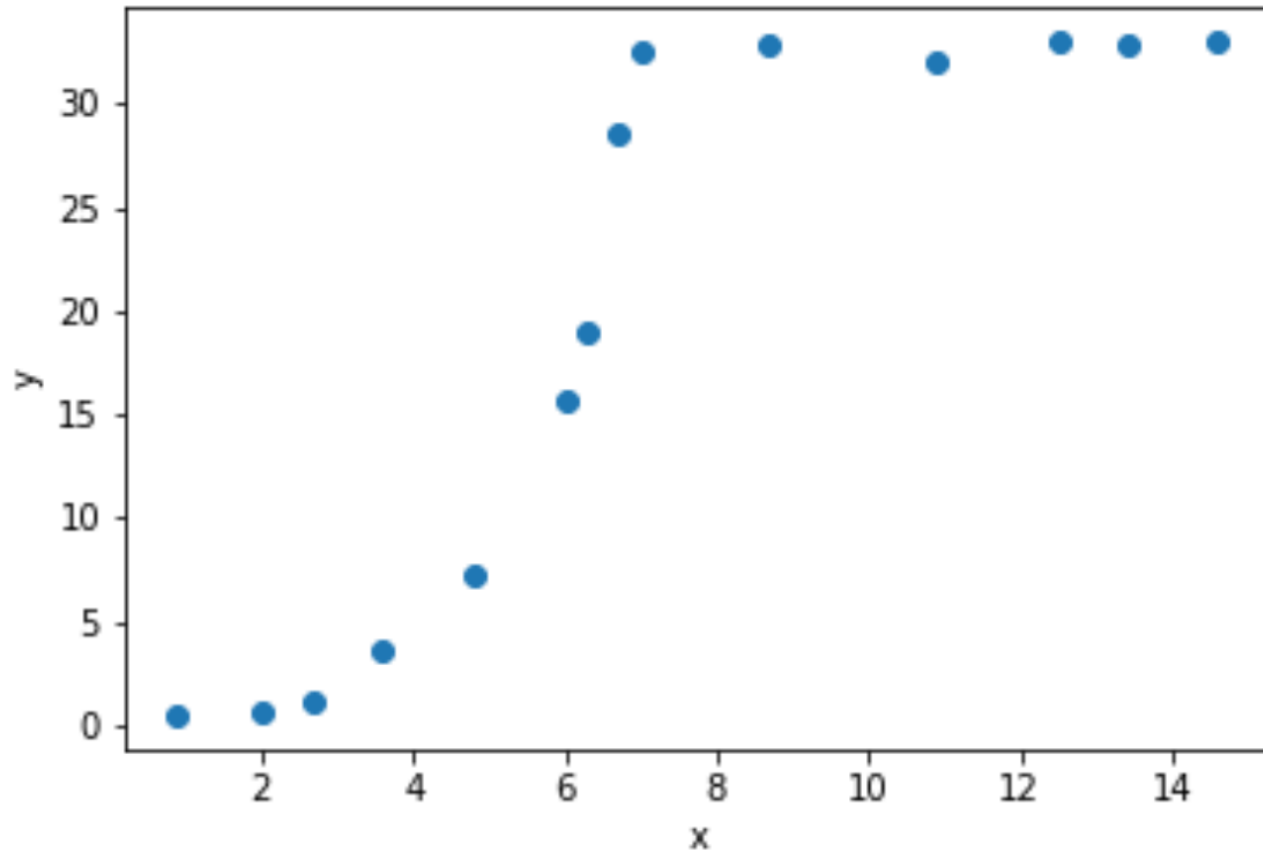
## Example 1 – raw data

| X    | Y    | X    | Y    |
|------|------|------|------|
| 0.9  | 0.5  | 8.7  | 32.3 |
| 2.7  | 1.1  | 4.8  | 7.3  |
| 6.7  | 28.6 | 12.5 | 33.1 |
| 10.9 | 32.8 | 13.4 | 32.9 |
| 6.0  | 15.7 | 2.0  | 0.75 |
| 6.3  | 19   | 3.6  | 3.6  |
| 7.0  | 32.6 | 14.6 | 33   |

|         | X   | Y    |
|---------|-----|------|
| Median  | 6.5 | 23.8 |
| Mean    | 7.2 | 19.5 |
| STD.DEV | 4.2 | 13.6 |

Correlation = 0.88

## Example 1 – Visualized

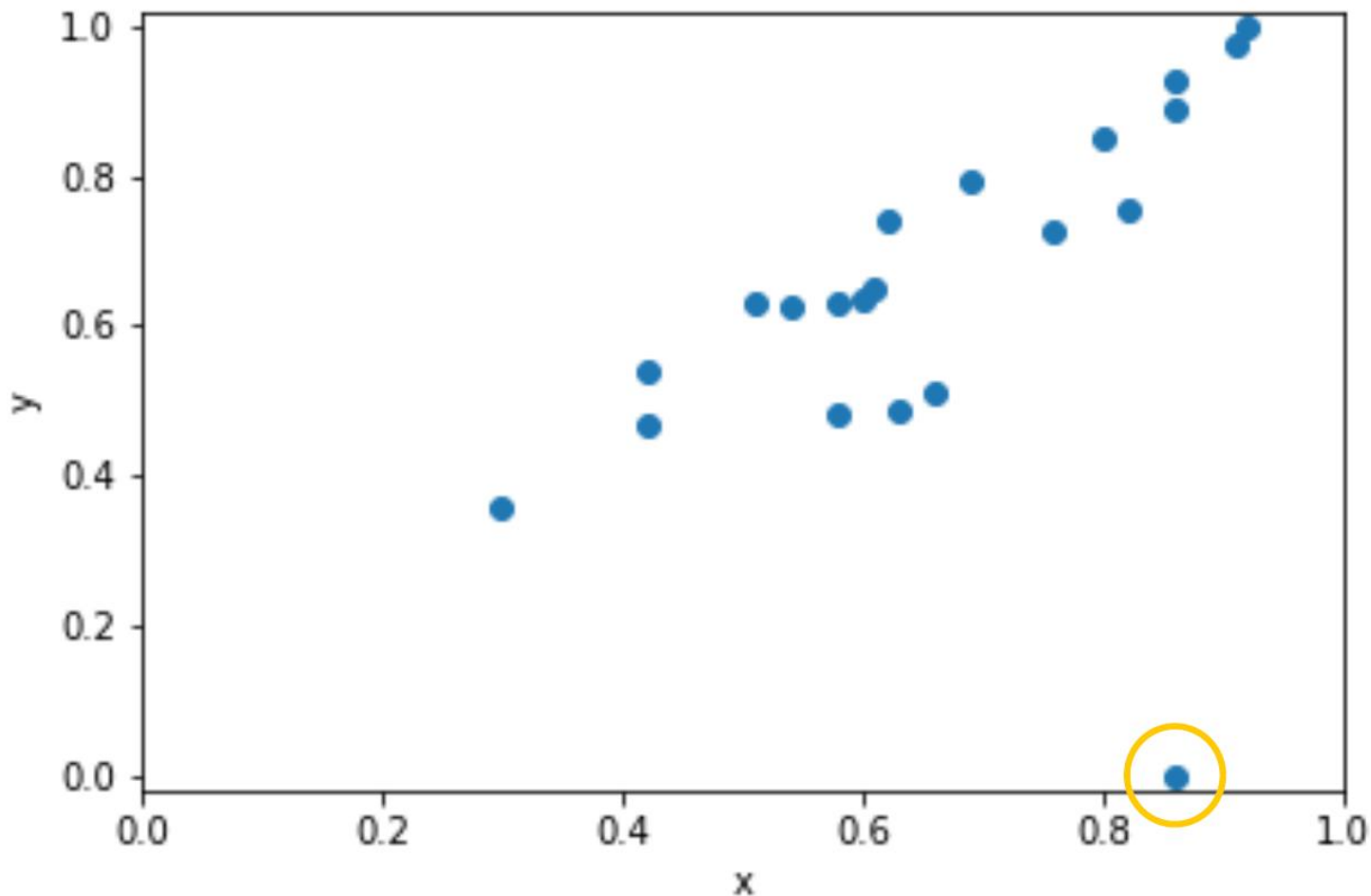




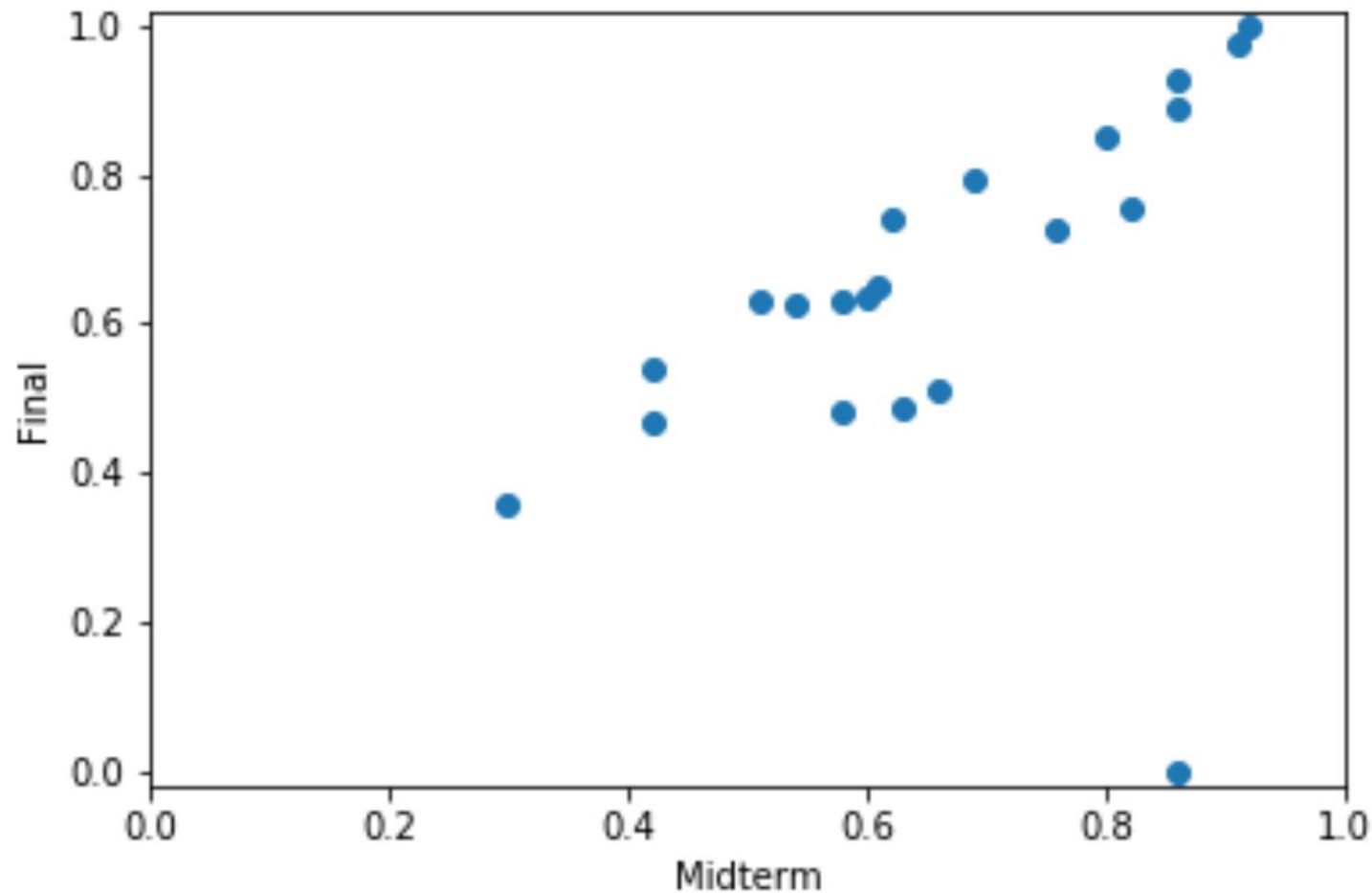
## Example 2 – raw data

| X    | Y       | X    | Y       |
|------|---------|------|---------|
| 0.42 | 0.46750 | 0.66 | 0.51250 |
| 0.54 | 0.62500 | 0.63 | 0.48750 |
| 0.42 | 0.53750 | 0.92 | 1.00000 |
| 0.86 | 0.92750 | 0.86 | 0.88750 |
| 0.60 | 0.63750 | 0.91 | 0.97500 |
| 0.51 | 0.63125 | 0.82 | 0.75625 |
| 0.30 | 0.35625 | 0.86 | 0.00000 |
| 0.61 | 0.65000 | 0.80 | 0.85000 |
| 0.58 | 0.63125 | 0.69 | 0.79375 |
| 0.76 | 0.72500 | 0.62 | 0.74000 |
| 0.58 | 0.48125 |      |         |

## Example 2 – Visualized



## Example 2 – Visualized

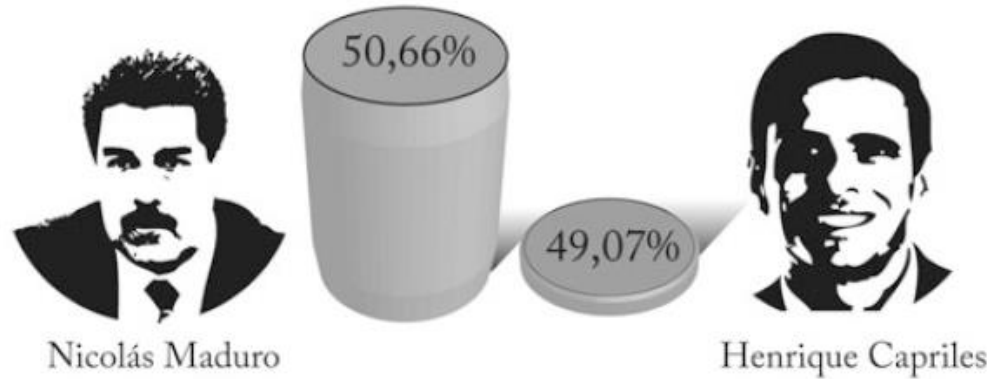


## Good data visualization

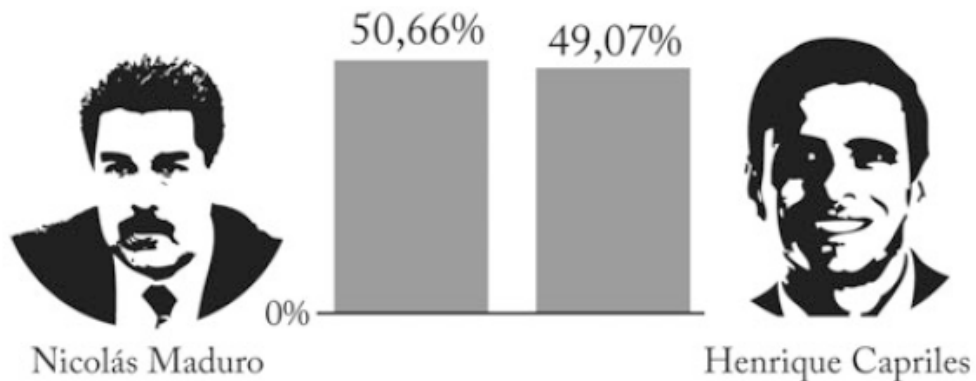
- Trustworthy
- Accessible
- Elegant



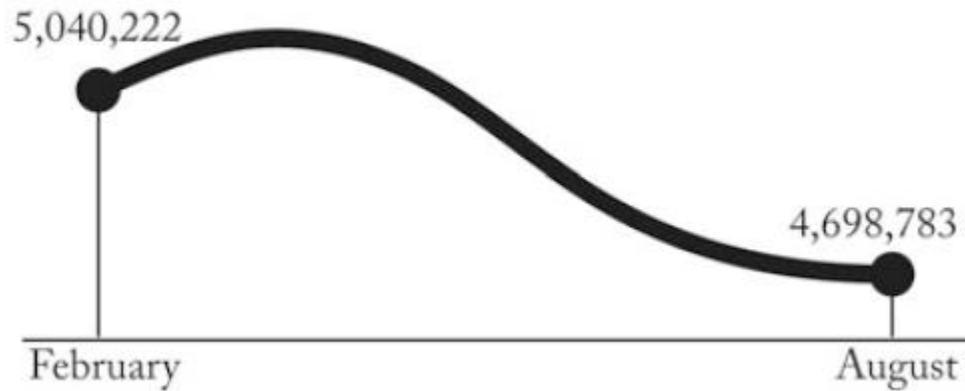
## PRESIDENTIAL ELECTIONS, 2013



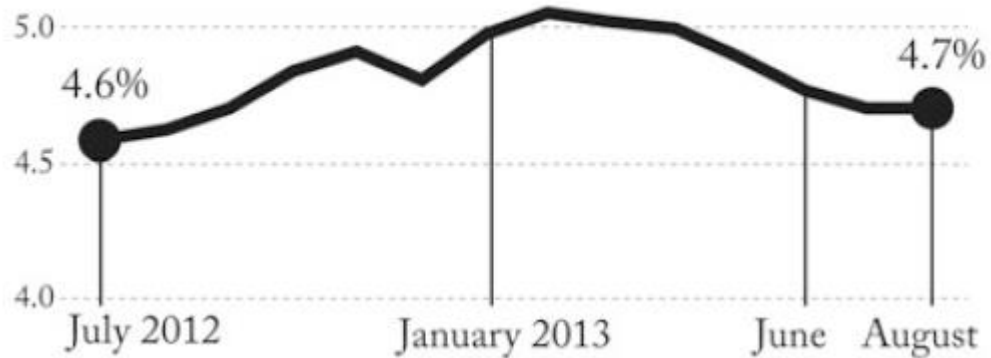
## PRESIDENTIAL ELECTIONS, 2013



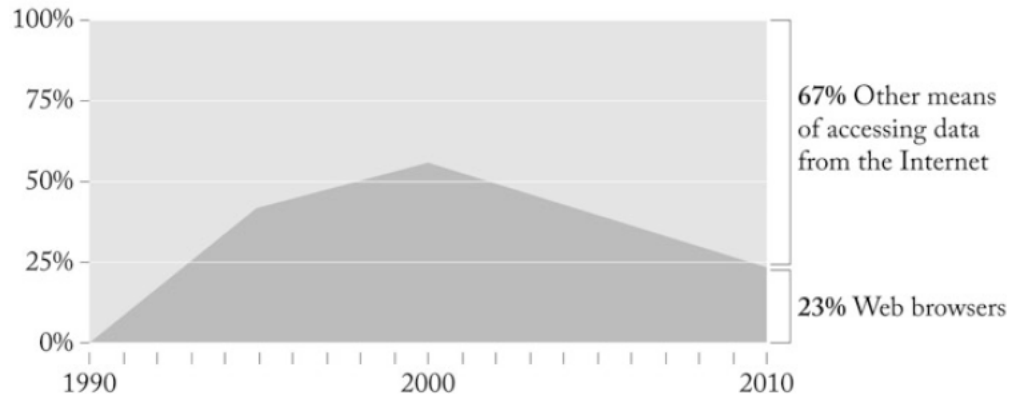
UNEMPLOYMENT 2013



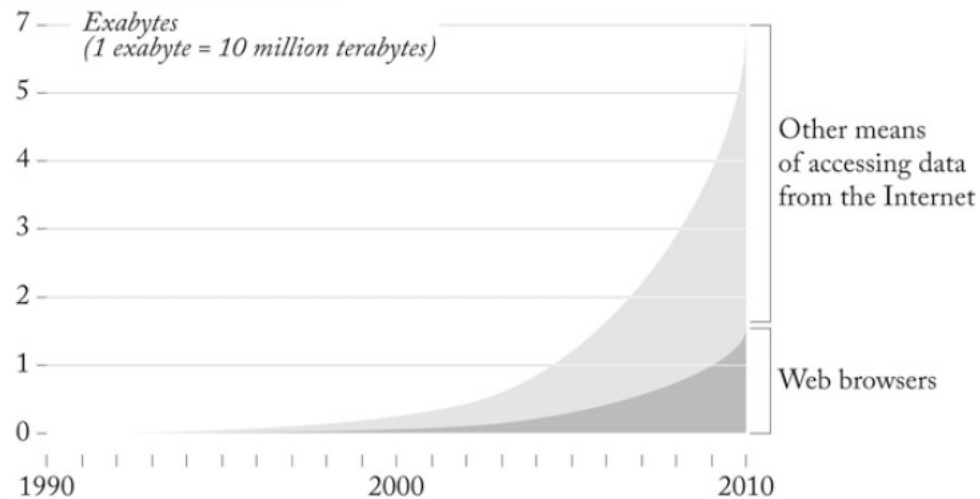
UNEMPLOYMENT



Internet Traffic in the U.S.

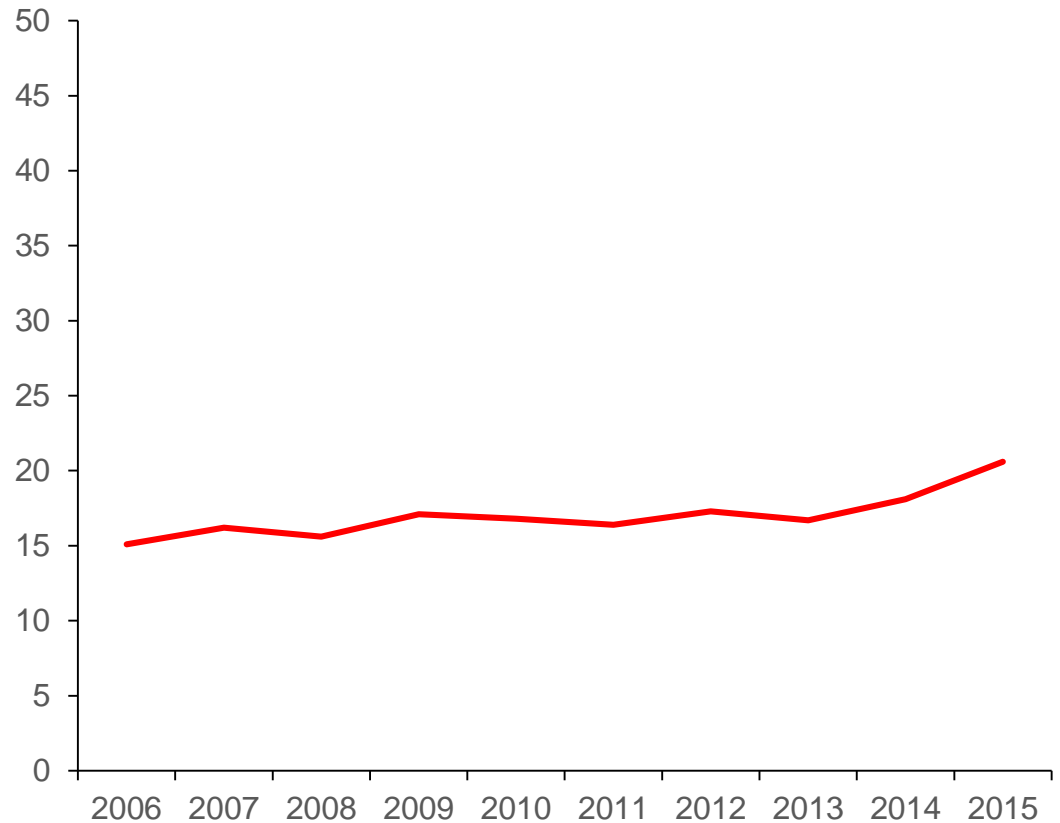
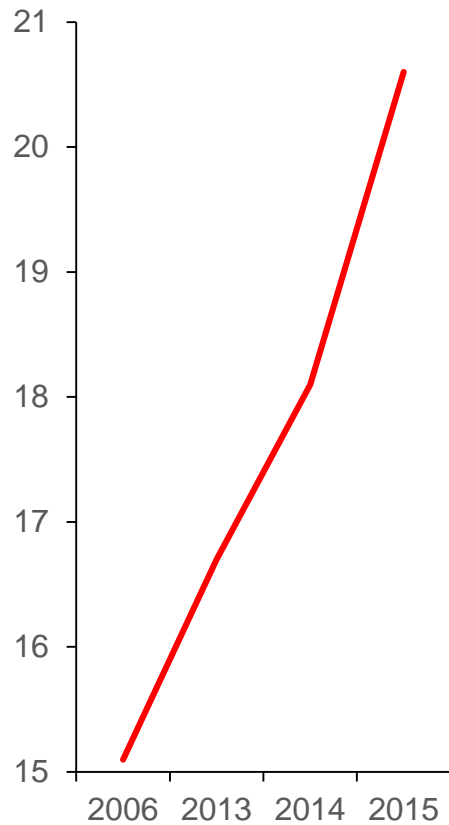


Internet Traffic in the U.S.

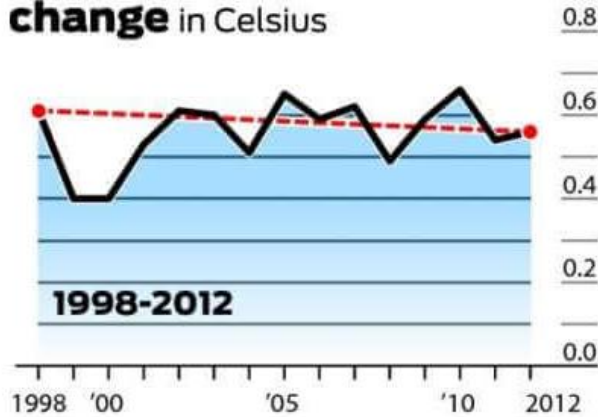




## Proportion of women on supervisory and management boards



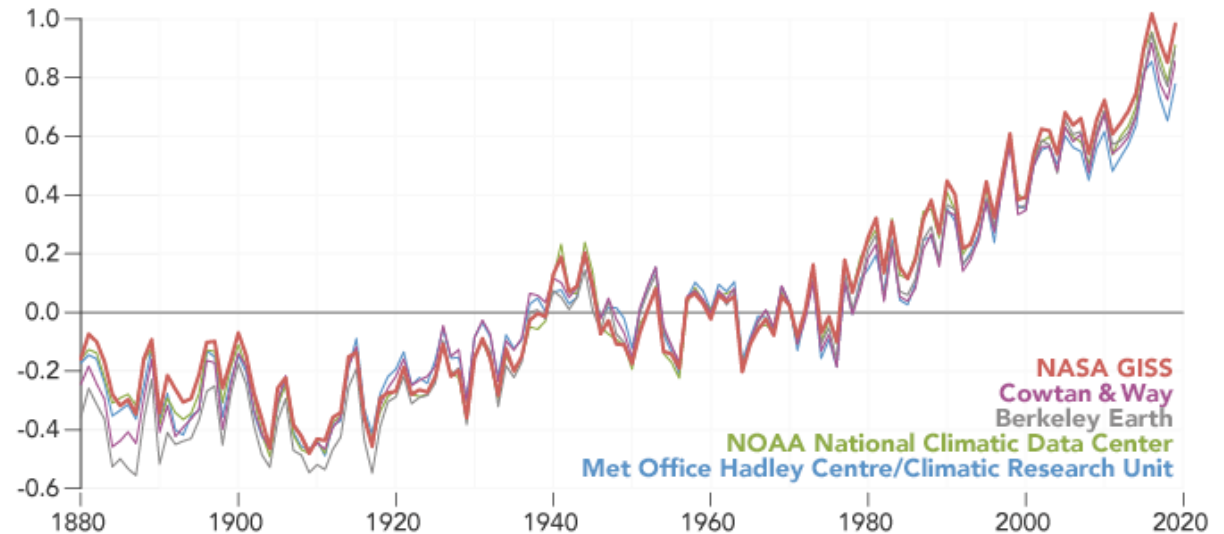
## Global air temperature change in Celsius



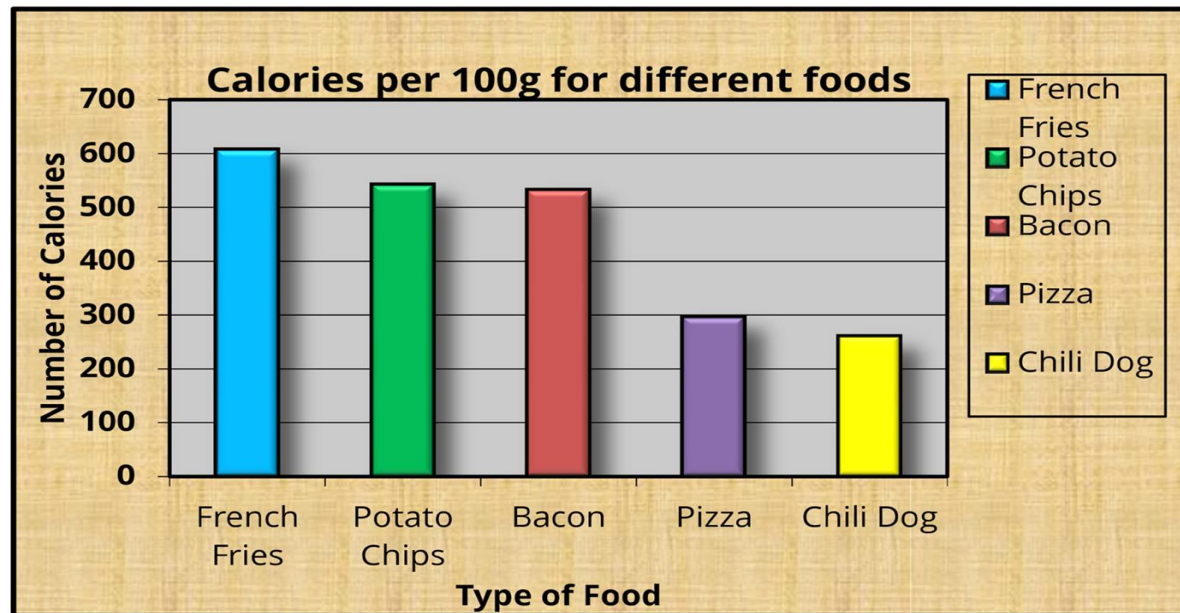
Source: NASA

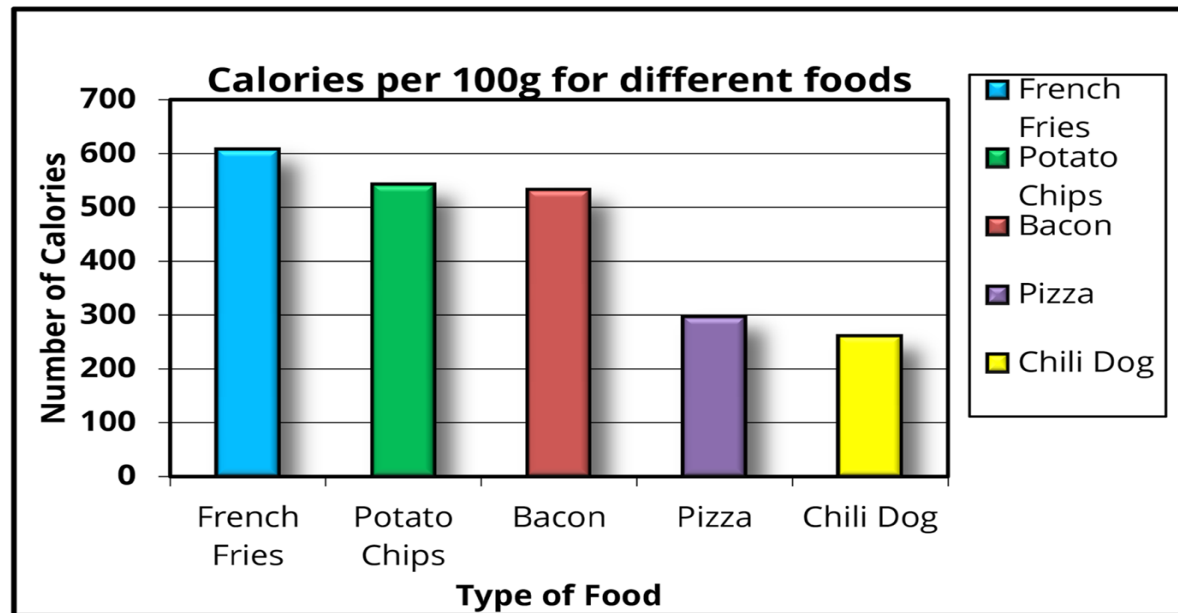
Source:  
<https://www.datapine.com/blog/misleading-statistics-and-data/>

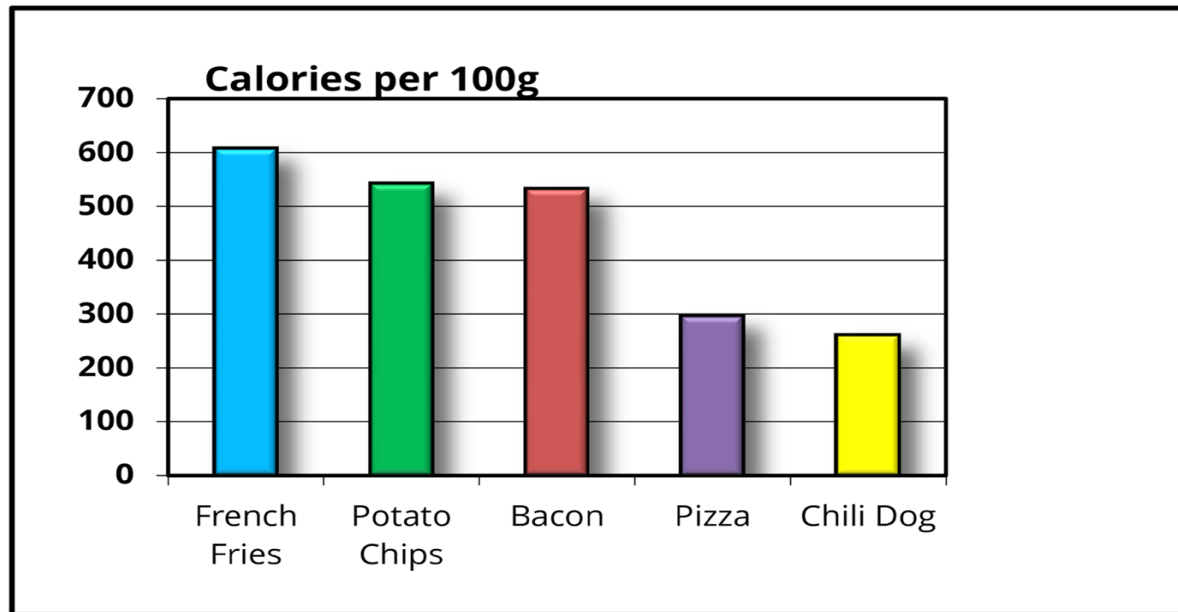
## A World of Agreement: Temperatures are Rising Global Temperature Anomaly (relative to 1951-1980, °C)

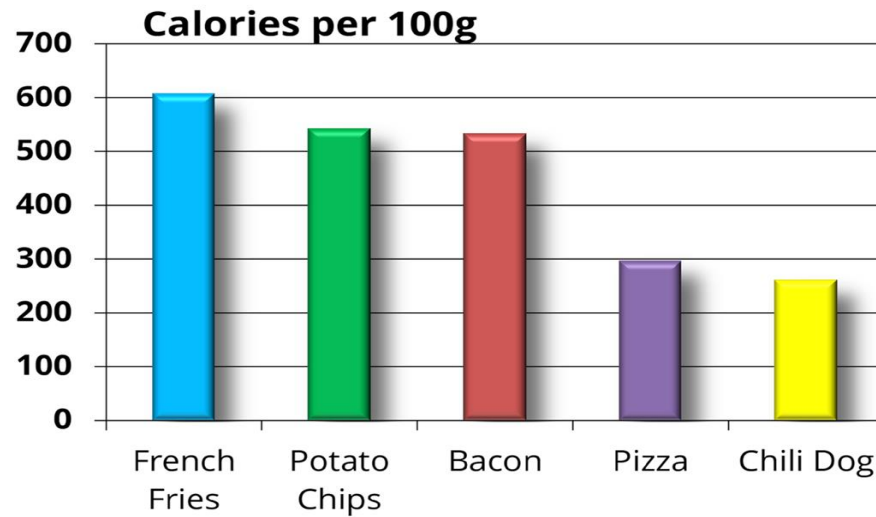


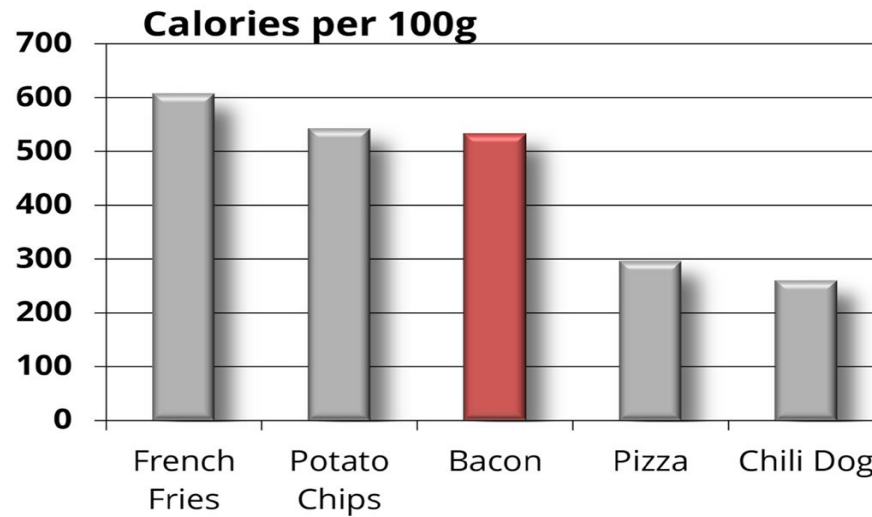
Source: <https://earthobservatory.nasa.gov/world-of-change/global-temperatures>

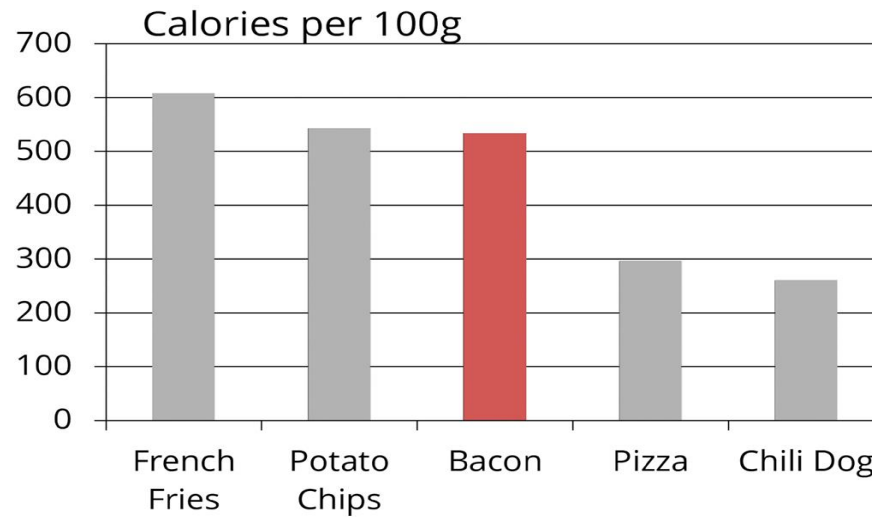




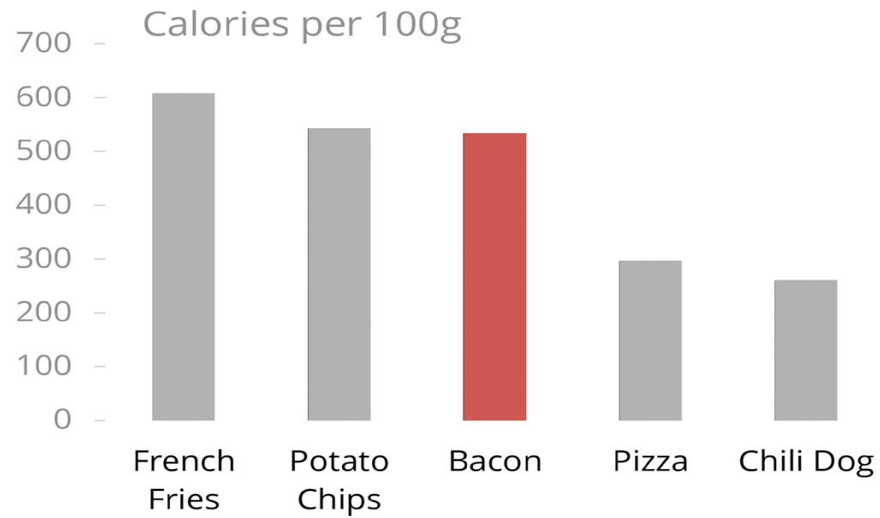


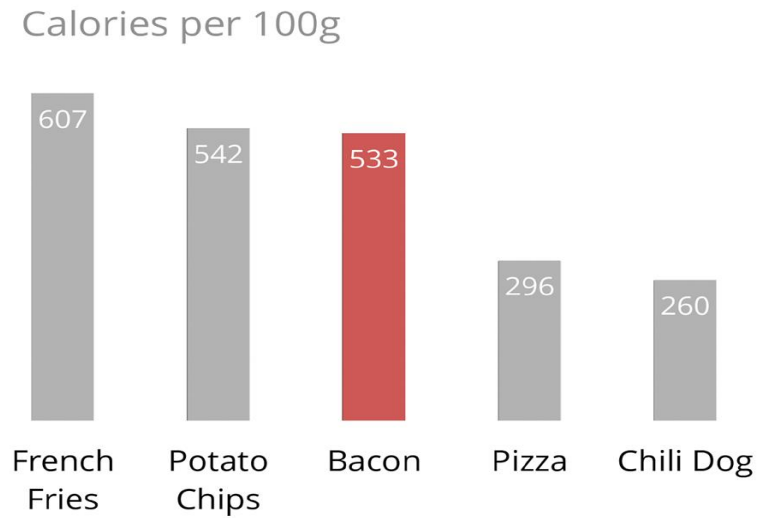




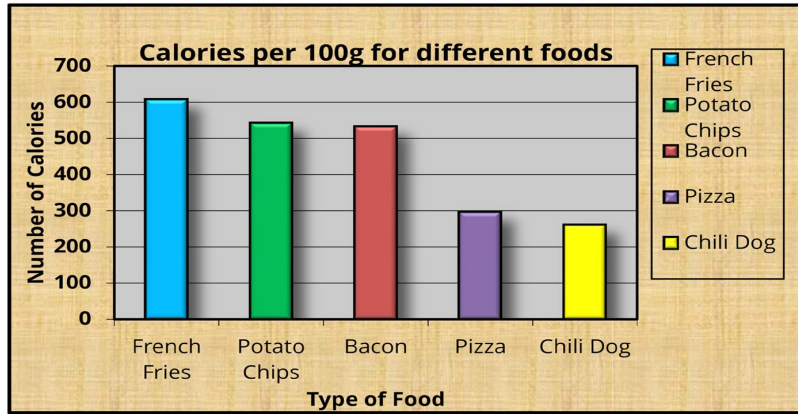




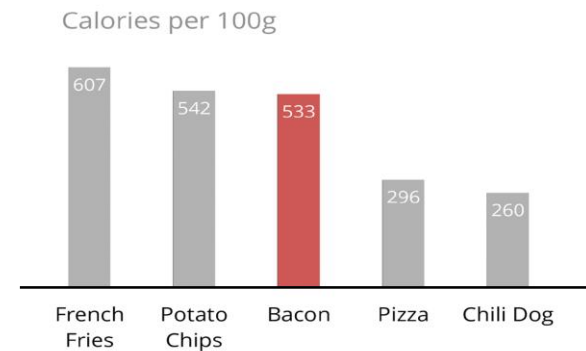




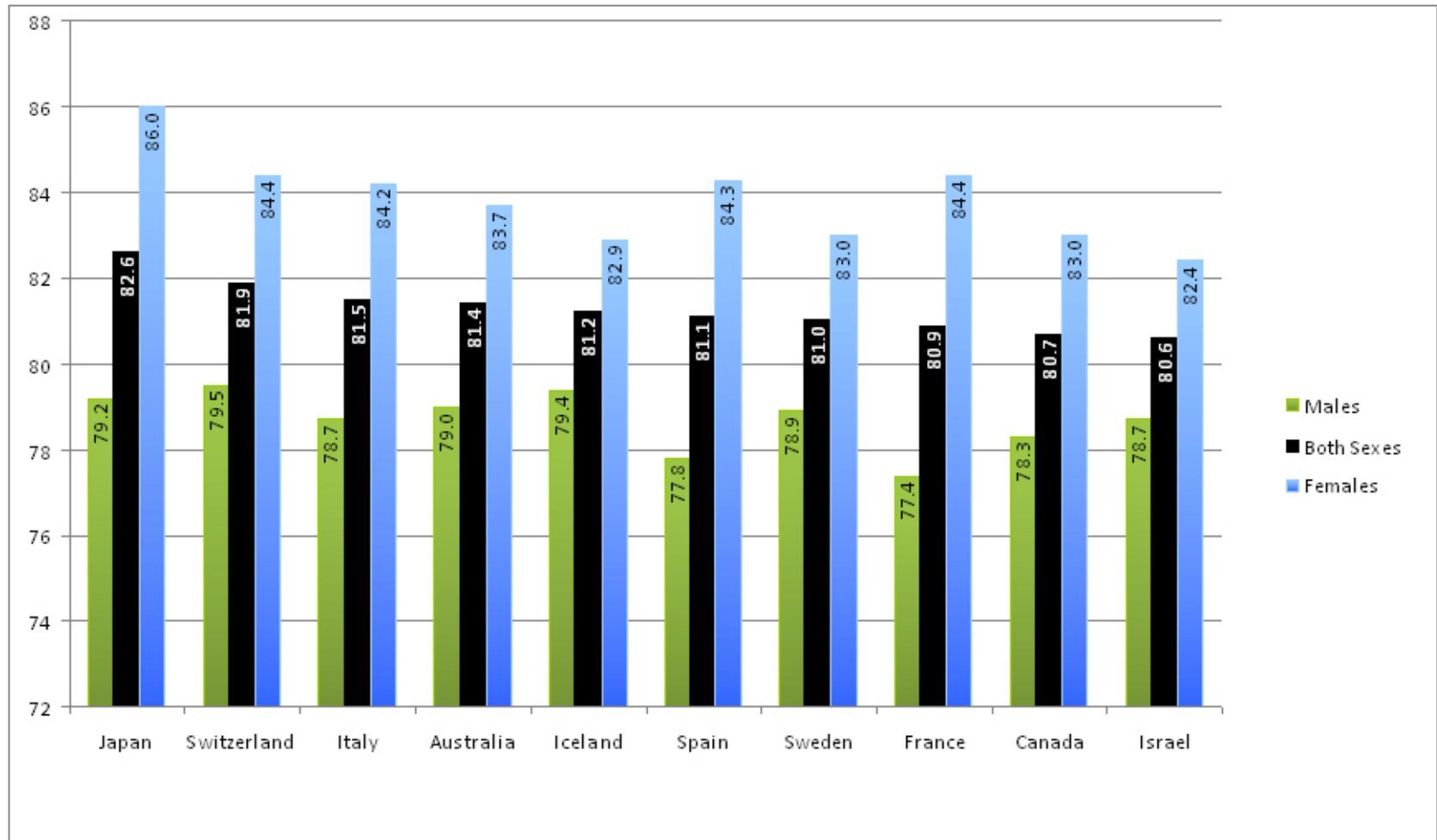
## Before



## After

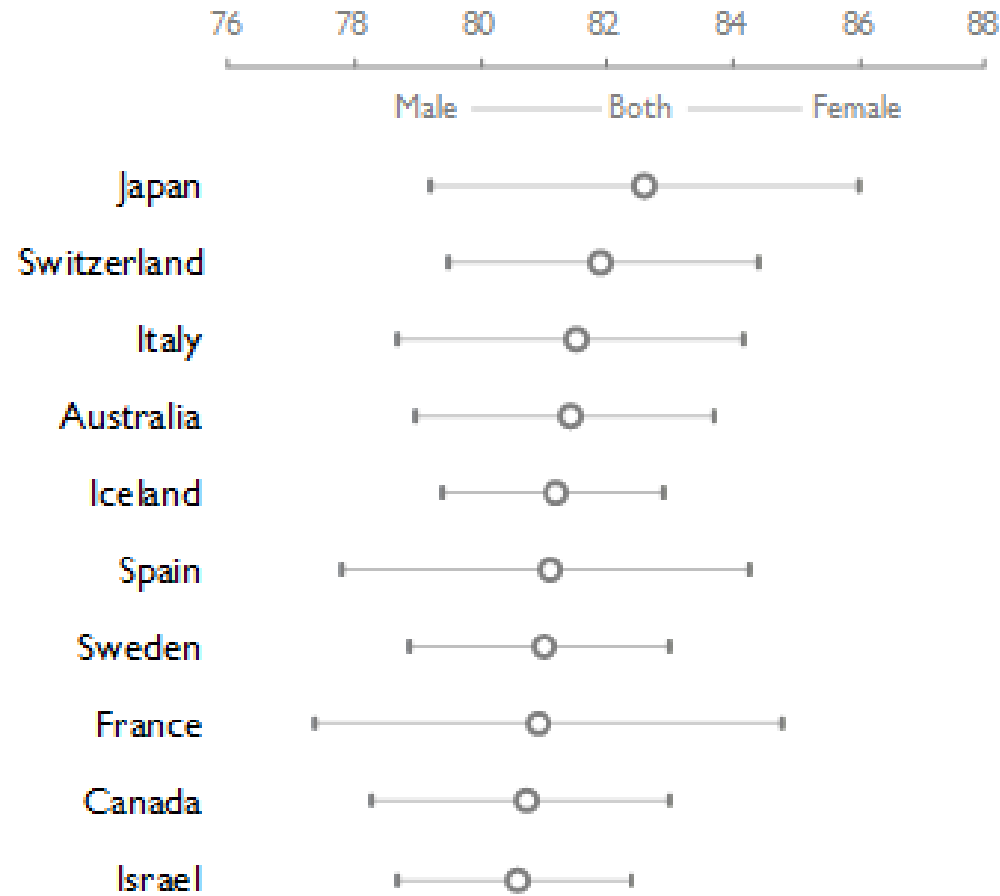


# Accessibility: Too Many Bars – Original Chart



## Life Expectancy at Birth

Top Ten OECD Countries 2010

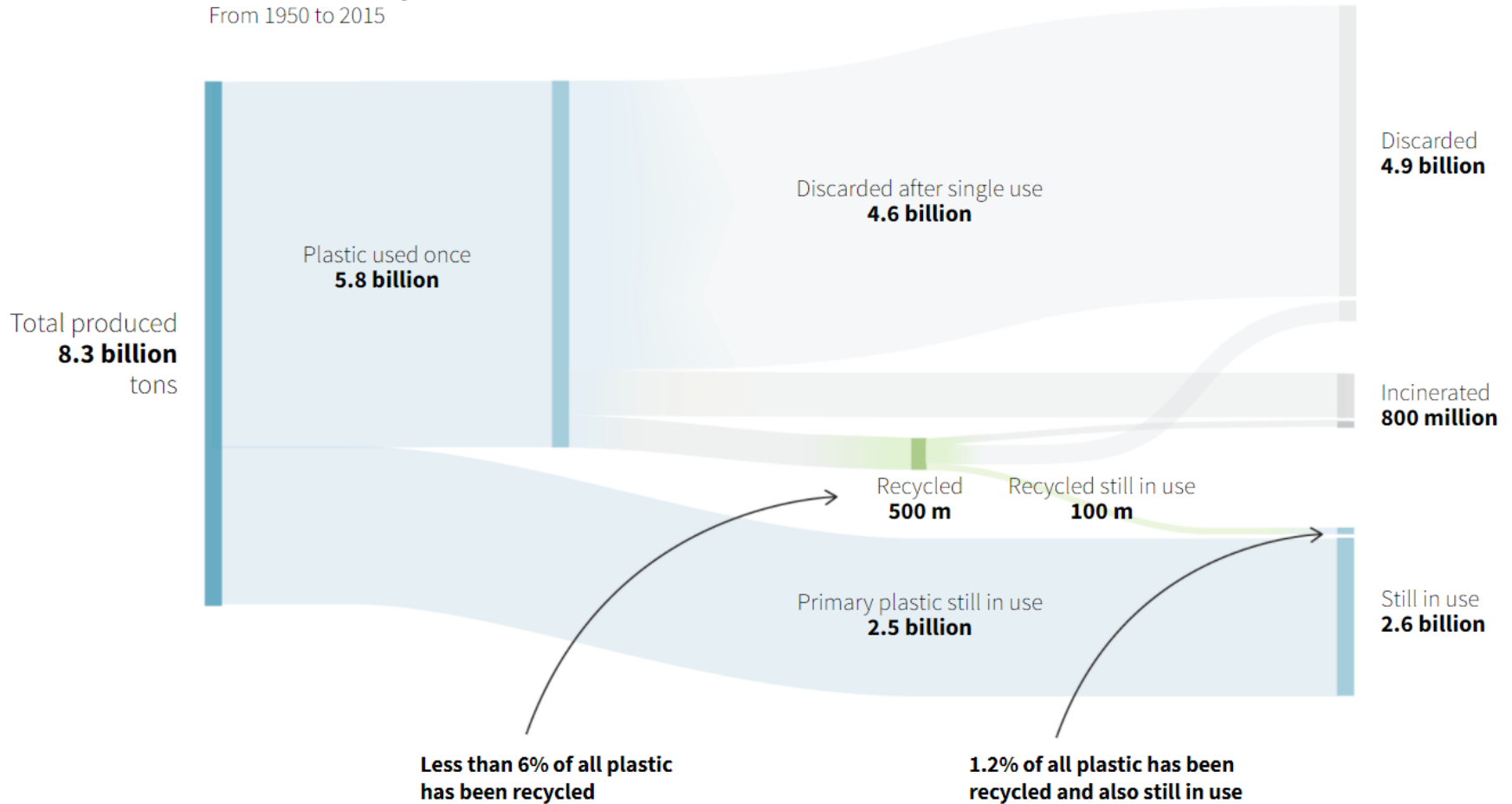


# Elegance: Cholera Outbreak London 1854

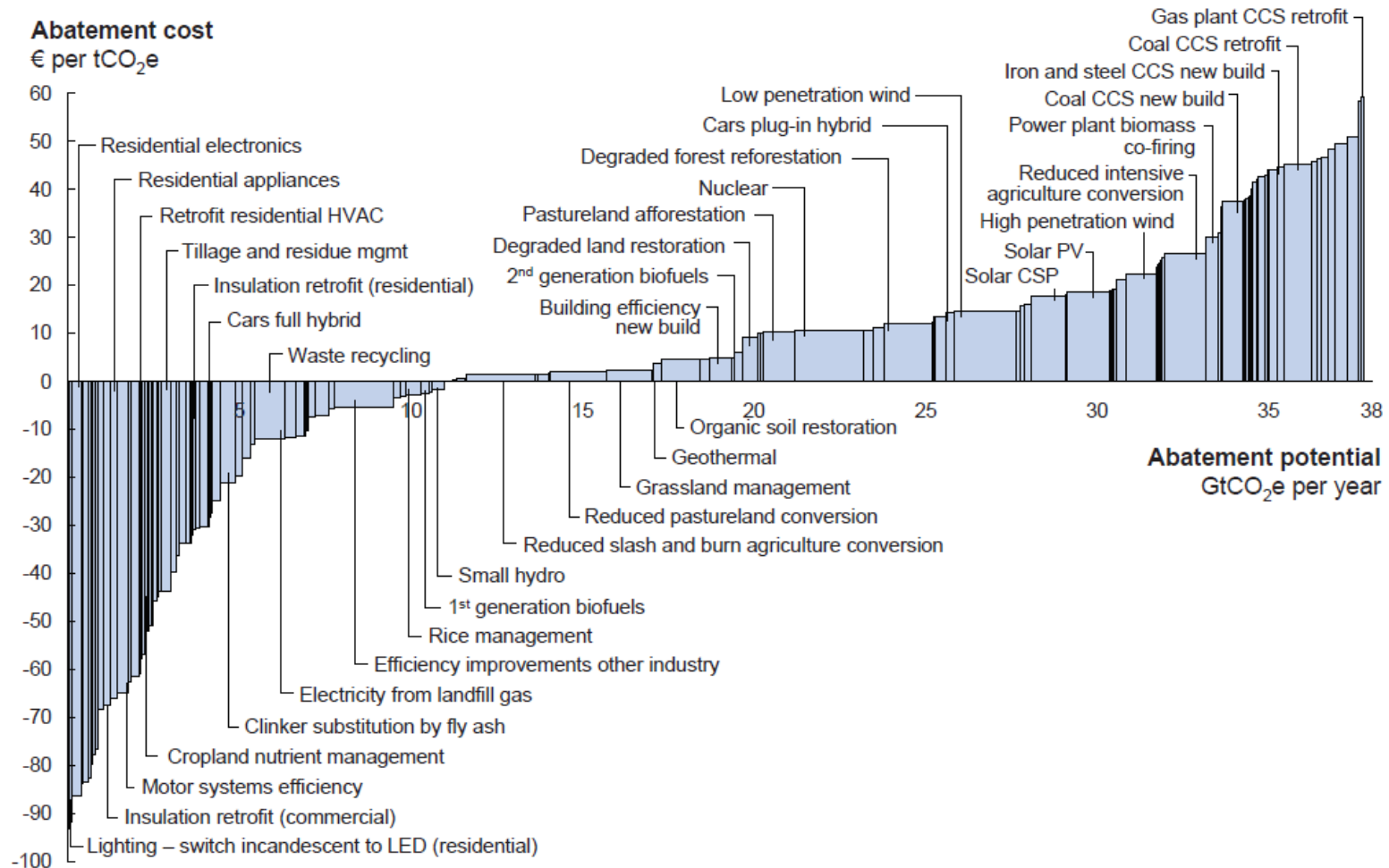


## The fate of all plastic

From 1950 to 2015



## Global GHG abatement cost curve beyond business-as-usual – 2030

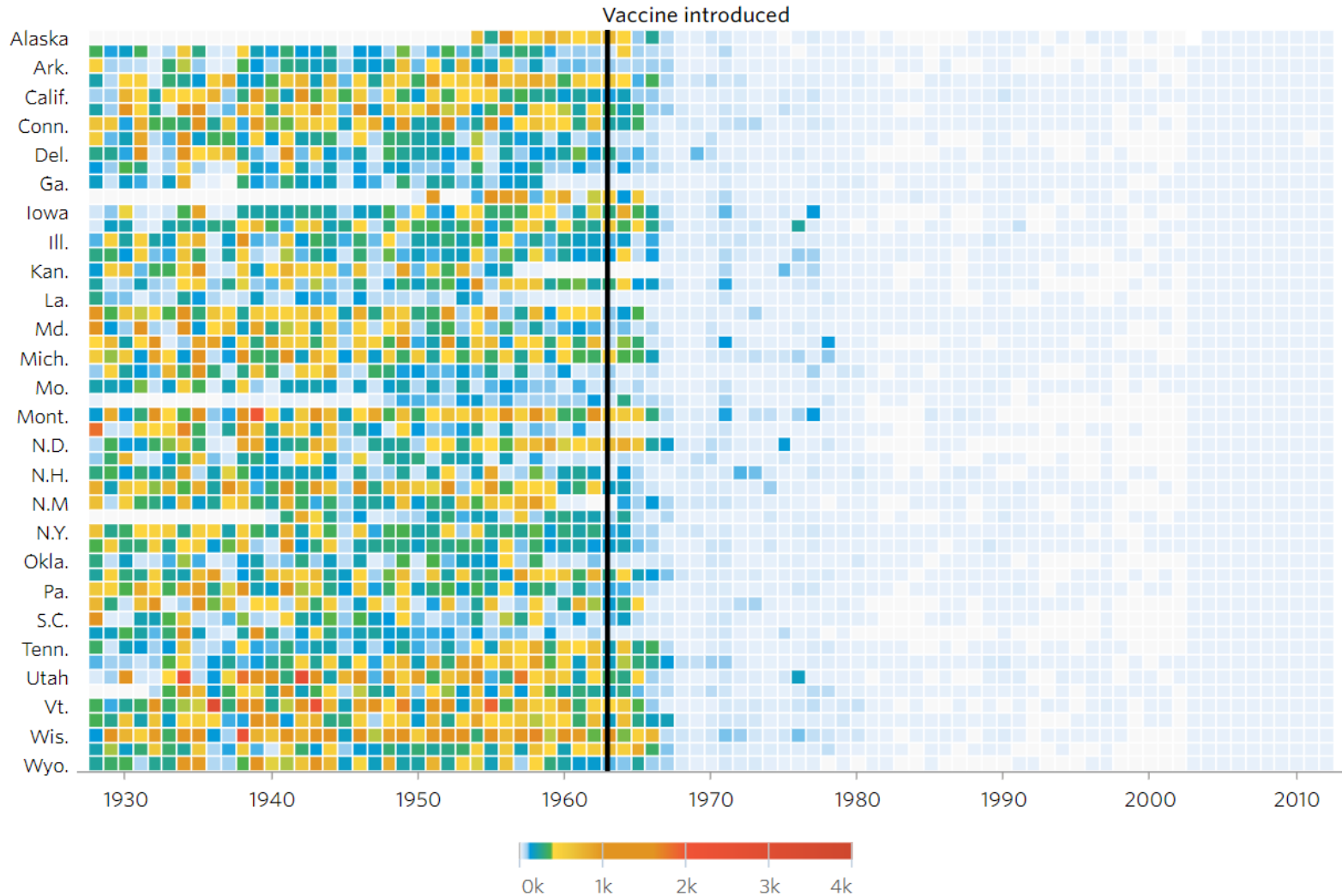


Note: The curve presents an estimate of the maximum potential of all technical GHG abatement measures below €60 per tCO<sub>2</sub>e if each lever was pursued aggressively. It is not a forecast of what role different abatement measures and technologies will play.  
Source: Global GHG Abatement Cost Curve v2.0



# Trustworthy? Accessible? Elegant?

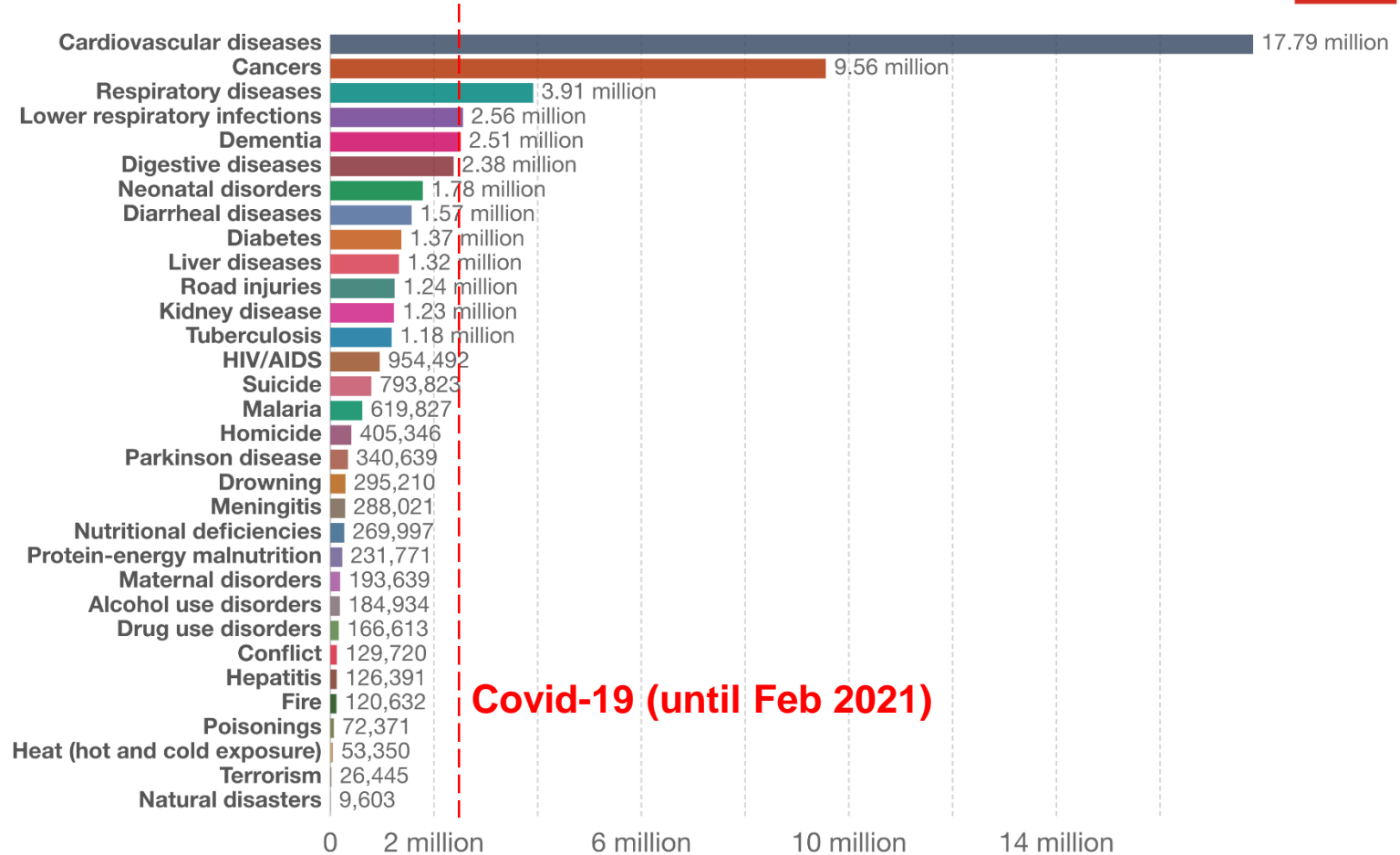
Number of measles cases per 100,000 people



# Trustworthy? Accessible? Elegant?

## Number of deaths by cause, World, 2017

Our World  
in Data



Source: IHME, Global Burden of Disease

OurWorldInData.org/causes-of-death • CC BY

# Deadlines for Submission and Distribution of Grading

| Student task   | Deliverables               | Deadline             | Work       | Share of grade |
|----------------|----------------------------|----------------------|------------|----------------|
| W01 Assignment | Code and results           | Apr 26 <sup>th</sup> | Team A     | 5.0%           |
| W02 Case Study | Code / presentation slides | May 15 <sup>th</sup> | Team B     | 18.0%          |
| W02 Case Study | Peer review*               | May 24 <sup>th</sup> | Individual | 2.0%           |
| W03 Assignment | Code and results           | May 22 <sup>nd</sup> | Team B     | 5.0%           |
| W04 Assignment | Code and results           | Jun 7 <sup>th</sup>  | Team C     | 10.0%          |
| W05 Assignment | Code and results           | Jun 14 <sup>th</sup> | Team D     | 7.0%           |
| W06 Assignment | Code and results           | Jun 28 <sup>th</sup> | Team D     | 13.0%          |
| W07 Case Study | Code / presentation slides | Jul 12 <sup>th</sup> | Team D     | 22.0%          |
| W07 Case Study | Peer review*               | Jul 26 <sup>th</sup> | Individual | 3.0%           |
| DataCamp 1     | Finish course              | May 9 <sup>th</sup>  | Individual | 2.5%           |
| DataCamp 2     | Finish course              | May 30 <sup>th</sup> | Individual | 2.5%           |
| DataCamp 3     | Finish course              | Jun 20 <sup>th</sup> | Individual | 2.5%           |
| DataCamp 4     | Finish course              | Jul 11 <sup>th</sup> | Individual | 2.5%           |

\* Peer review is mandatory. Quality of peer review itself is graded. Not providing peer review at all would result in high point deduction