# Data visualization using Python and Vega-Altair

In this introductory-level workshop, we will learn to produce **reproducible data visualization pipelines using the Python programming language** and the Vega-Altair declarative visualization library.

We will work in Jupyter Notebooks and start with **Python basics**. We will introduce the pandas library for "data wrangling" (reading, writing, sorting, and filtering of data). With pandas, we will be able to read data from Excel sheets and comma-separated values (CSV) files.

Finally, will learn how to produce and share **reproducible plots** using Vega-Altair.

## Who is the course for?

- Somebody starting with Python or curious about Python.
- Somebody who needs to read, process, and plot data for their work or studies and would like to try it out with Python.
- Persons who already use Python for this but want to learn about libraries to simplify common tasks and about how to share their workflow in a reproducible way.

> ⚙ **Preparations**
>
> - No programming language experience needed, we will start from zero and learn the basics together
> - Computer with network access
> - Software install instructions
> - Bring one of your recent plotting tasks or challenges

## What is not taught?

- Version control. Although super useful it is outside of this workshop.
- Python outside a Jupyter Notebook.
- **Running the examples in VS Code or Spyder might not be possible**. Please use Jupyter Notebooks for this course.
- Python sets and tuples are only mentioned.
- File input/output is only used via libraries and doing "own" file-I/O is only part of optional material.

- How to choose the right visualization format for the data at hand.
- Python object oriented design.
- Python packaging.
- NumPy arrays.
- Managing environments and installing Python packages.

# Episode overview

Day 1 morning:

- Software install instructions
- Jupyter Notebooks  `view on` `nbviewer`  `CO` `Open in Colab`  `launch` `binder`
- Python basics
- Plotting with Vega-Altair  `view on` `nbviewer`  `CO` `Open in Colab`  `launch` `binder`

Day 1 afternoon:

- More Plotting with Vega-Altair
- Tidy data and dealing with messy data

Day 2 morning:

- Customizing plots  `view on` `nbviewer`  `CO` `Open in Colab`  `launch` `binder`
- Learning how to adapt existing gallery examples

Day 3 morning:

- How to find help and how to navigate the documentation
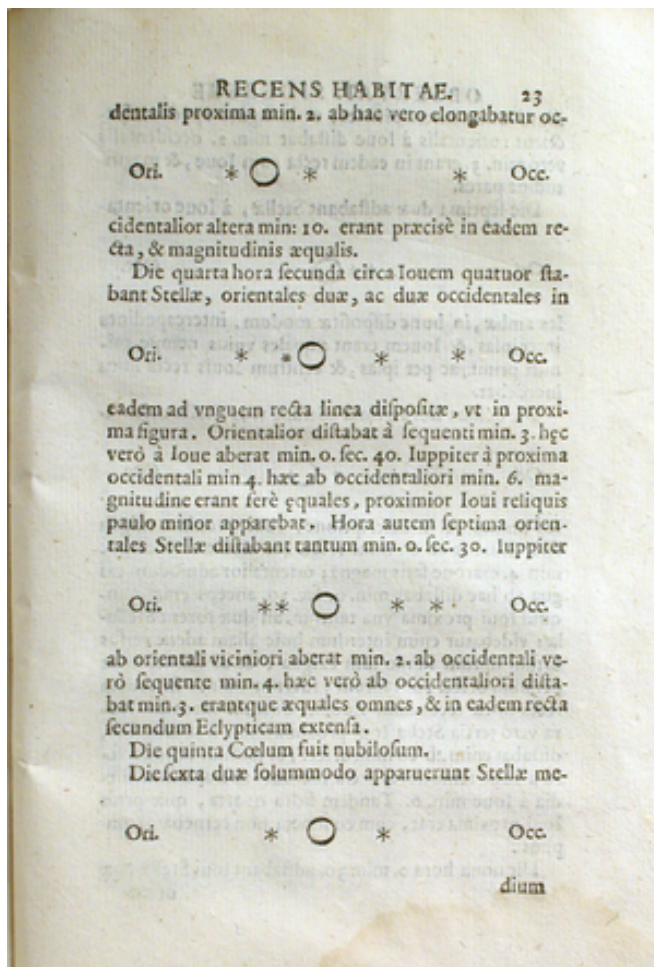- Sharing plots and notebooks

# Jupyter Notebooks

**❶ Objectives**

- Know what it is
- Create a new notebook and save it
- Open existing notebooks from the web
- Be able to create text/markdown cells, code cells, images, and equations
- Know when to use a Jupyter Notebook for a Python project and when perhaps not to
- We will build up this notebook (spoiler alert!)

[this lesson is adapted from https://coderefinery.github.io/jupyter/motivation/]

# Motivation for Jupyter Notebooks

*One of the first notebooks: Galileo's drawings of Jupiter and its Medicean Stars* from *Sidereus Nuncius. Image courtesy of the History of Science Collections, University of Oklahoma Libraries (CC-BY).*

- **Code, text, equations, figures, plots**, etc. are interleaved, creating a *computational narrative*.
- *"an environment in which users execute code, see what happens, modify and repeat in a kind of iterative conversation between researcher and data"*
- The name "Jupyter" derives from Julia+Python+R, but today Jupyter kernels exist for dozens of programming languages.
- Gallery of interesting Jupyter Notebooks.

## Our first notebook

### ✍️ Exercise Jupyter-1: Create a notebook (15 min)

- Open a new notebook (on Windows: open Anaconda Navigator, then launch JupyterLab; on macOS/Linux: you can open JupyterLab from the terminal by typing `jupyter-lab` )
- Rename the notebook
- Create a **markdown cell** with a section title, a short text, an image, and an equation

```
# Title of my notebook

Some text.

![Photo of Galilei's manuscript]
(https://upload.wikimedia.org/wikipedia/commons/b/b3/Galileo_Galilei_%281564_-
_1642%29_-_Serenissimo_Principe_-
_manuscript_with_observations_of_Jupiter_and_four_of_its_moons%2C_1610.png)

$E = mc^2$
```

- Most important shortcut: **Shift + Enter**, to run current cell and create a new one below.
- Create a **code cell** where you define the `arithmetic_mean` function:

```python
def arithmetic_mean(sequence):
    s = 0.0
    for element in sequence:
        s += element
    n = len(sequence)
    return s / n
```

- In a different cell, call the function:

```python
arithmetic_mean([1, 2, 3, 4, 5])
```

- In a new cell, let us try to plot a layered histogram:

```python
# this example is from https://altair-
viz.github.io/gallery/layered_histogram.html

import pandas as pd
import altair as alt
import numpy as np
np.random.seed(42)

# Generating Data
source = pd.DataFrame({
    'Trial A': np.random.normal(0, 0.8, 1000),
    'Trial B': np.random.normal(-2, 1, 1000),
    'Trial C': np.random.normal(3, 2, 1000)
})

alt.Chart(source).transform_fold(
    ['Trial A', 'Trial B', 'Trial C'],
    as_=['Experiment', 'Measurement']
).mark_bar(
    opacity=0.3,
    binSpacing=0
).encode(
    alt.X('Measurement:Q').bin(maxbins=100),
    alt.Y('count()').stack(None),
    alt.Color('Experiment:N')
)
```

## Use cases for notebooks

- Really good for step-by-step recipes (e.g. read data, filter data, do some statistics, plot the results)
- Experimenting with new ideas, testing new libraries/databases
- As an *interactive* development environment for code, data analysis, and visualization
- Keeping track of interactive sessions, like a **digital lab notebook**
- **Supplementary information with published articles**

## Good practices

**Run all cells** or even **Restart Kernel and Run All Cells** before sharing/saving to verify that the results you see on your computer were not due to cells being run out of order.

This can be demonstrated with the following example:

```
numbers = [1, 2, 3, 4, 5]
arithmetic_mean(numbers)
```

We can first split this code into two cells and then re-define `numbers` further down in the notebook. If we run the cells out of order, the result will be different.

## Python basics

**❗ Objectives**

- Knowing what types exist
- Knowing the most common data structures: lists, tuples, dictionaries, and sets
- Creating and using functions
- Knowing what a library is
- Knowing what `import` does
- Being able to "read" an error

### Motivation for Python

- **Free**
- Huge **ecosystem of examples, libraries, and tools**
- Relatively easy to read and understand

- Similar in scope and use cases to R, Julia, and Matlab

## Basic types

```python
# int
num_measurements = 13

# float
some_fraction = 0.25

# string
name = "Bruce Wayne"

# bool
value_is_missing = False
skip_verification = True

# we can print values
print(name)

# and we can do arithmetics with ints and floats
print(5 * num_measurements)
print(1.0 - some_fraction)
```

- Python is **dynamically typed**: We do not have to define that an integer is an `int`, we can use it this way and Python will infer it.
- However, one can use type annotations in Python (see also mypy).
- Now you also know that we can add `# comments` to our code.

## Data structures for collections: lists, dictionaries, sets, and tuples

```python
# lists are good when order is important
scores = [13, 5, 2, 3, 4, 3]

# first element
print(scores[0])

# we can add items to lists
scores.append(4)

# lists can be sorted
scores.sort()
print(scores)

# dictionaries are useful if you want to look up
# elements in a collection by something else than position
experiment = {"location": "Svalbard", "date": "2021-03-23", "num_measurements": 23}

print(experiment["date"])

# we can add items to dictionaries
experiment["instrument"] = "a particular brand"
print(experiment)

if "instrument" in experiment:
    print("yes, the dictionary 'experiment' contains the key 'instrument'")
else:
    print("no, it doesn't")
```

- `Lists` are good when order is important, and it needs to be changed
- `Dictionaries` are mappings key→value.
- `Sets` are useful for unordered collections where you want to make sure that there are no repetitions.
- There are also `tuples` that are similar to lists but their items cannot be modified.

You can put:

- dictionaries inside lists
- lists inside dictionaries
- dictionaries inside dictionaries
- lists inside lists
- tuples inside ...
- ...

## Iterating over collections

Often we wish to iterate over collections.

Iterating over a list:

```
scores = [13, 5, 2, 3, 4, 3]

for score in scores:
    print(score)

# example with f-strings
for score in scores:
    print(f"the score is {score}")
```

We don't have to call the variable inside the for-loop "score". This is up to us. We can do this instead (but is this more understandable for humans?):

```
scores = [13, 5, 2, 3, 4, 3]

for x in scores:
    print(x)
```

Iterating over a dictionary:

```
experiment = {"location": "Svalbard", "date": "2021-03-23", "num_measurements": 23}

for key in experiment:
    print(experiment[key])

# another way to iterate
for (key, value) in experiment.items():
    print(key, value)
```

## Functions

- Functions are like **reusable recipes**. They receive ingredients (input arguments), then inside the function we do/compute something with these arguments, and they return a result.

```
def add(a, b):
    result = a + b
    return result
```

- Together we write a function which sums all elements in a list:

```python
def add_all_elements(sequence):
    """
    This function adds all elements.
    This here is a docstring, a documentation string for a function.
    """
    s = 0.0
    for element in sequence:
        s += element
    return s


measurements = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]

print(add_all_elements(measurements))
```

- We reuse this function to write a function which computes the mean:

```python
def arithmetic_mean(sequence):
    # we are reusing add_all_elements written above
    s = add_all_elements(sequence)
    n = len(sequence)
    return s / n


measurements = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]

mean = arithmetic_mean(measurements)

print(mean)
```

- Functions can call other functions. Functions can also get other functions as input arguments.
- Functions can return more than one thing:

```python
def uppercase_and_lowercase(text):
    u = text.upper()
    l = text.lower()
    return u, l


some_text = "SequenceOfCharacters"
uppercased_text, lowercased_text = uppercase_and_lowercase(some_text)

print(uppercased_text)
print(lowercased_text)
```

**Why functions?** Less repetition but also simplify reading and understanding code.

## Reading error messages

Here we introduce a mistake and we together try to make sense of the traceback:

```
------------------------------------------------------------------------
--
NameError                                Traceback (most recent call las
t)
<ipython-input-10-8e6794a136dc> in <module>
      8 measurements = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
      9
---> 10 mean = arithmetic_mean(measurements)
     11
     12 print(mean)

<ipython-input-10-8e6794a136dc> in arithmetic_mean(sequence)
      1 def arithmetic_mean(sequence):
      2     # we are reusing add_all_elements written above
----> 3     s = add_all_element(sequence)
      4     n = len(sequence)
      5     return s / n

NameError: name 'add_all_element' is not defined
```

*Example error traceback. Can you explain the error?*

## Libraries

We can look at libraries as collections of functions. We can import the libraries/modules and then reuse the functions defined inside these libraries.

Try this:

```
import numpy

measurements = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]

result = numpy.std(measurements)

print(result)
```

This means `numpy` contains a function called `std` which apparently computes the standard deviation (check also its documentation).

Often you see this in tutorials (the module is imported and renamed to a shortcut):

```
import numpy as np

result = np.std([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
```

It is possible to create own modules to collect own functions for reuse.

## Great resources to learn more

# Exercises

## ✍️ Exercise: create a function that computes the standard deviation

- Arithmetic mean:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

- Standard deviation:

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

- In other words the computation is similar but we need to sum over squares of differences and at the end take a square root.
- Take this as a starting point:

```python
# we have written this one together previously
def arithmetic_mean(sequence):
    s = 0.0
    for element in sequence:
        s += element
    n = len(sequence)
    return s / n


def standard_deviation(sequence):
    # here we need to do some work:
    # mean = ?
    # s = ?
    n = len(sequence)
    return (s / n) ** 0.5
```

- If this is the input list:

```python
measurements = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
```

Then the result would be: 2.872...

## ✔ Solution 1 (longer but hopefully easier to understand)

```python
# we have written this one together previously
def arithmetic_mean(sequence):
    s = 0.0
    for element in sequence:
        s += element
    n = len(sequence)
    return s / n


# notice how this function reuses the other
def standard_deviation(sequence):
    mean = arithmetic_mean(sequence)
    s = 0.0
    for element in sequence:
        s += (element - mean) ** 2
    n = len(sequence)
    return (s / n) ** 0.5
```

## ✔ Solution 2 (more compact)

```python
def arithmetic_mean(sequence):
    return sum(sequence) / len(sequence)


def standard_deviation(sequence):
    mean = arithmetic_mean(sequence)
    s = sum([(x - mean) ** 2 for x in sequence])
    n = len(sequence)
    return (s / n) ** 0.5
```

## 🖊 Exercise: working with a dictionary

- We have this dictionary as a starting point:

  ```python
  grades = {"Alice": 80, "Bob": 95}
  ```

- Add the grades of few more (fictious) persons to this dictionary.
- Print the entire dictionary.
- What happens when you add a name which already exists (with a different grade)?
- Print the grade for one particular person only.
- What happens when you try to print the result for a person that wasn't there?
- Try also these:

  ```python
  print(grades.keys())
  print(grades.values())
  print(grades.items())
  ```

## ✔ Solution

We can add more people like this:

```
grades["Craig"] = 56
grades["Dave"] = 28
grades["Eve"] = 75
```

Print the entire dictionary with:

```
print(grades)
```

We get:

```
{'Alice': 80, 'Bob': 95, 'Craig': 56, 'Dave': 28, 'Eve': 75}
```

Adding an entry which already exists updates the entry (please try it).

Printing the result for one particular person:

```
print(grades["Eve"])
```

Printing the result for a person which does not exists, gives a `KeyError`.

The outputs of these three:

```
print(grades.keys())
print(grades.values())
print(grades.items())
```

... are either the only the keys or only the values, or in the case of `items()`, key-value pairs (tuples):

```
dict_keys(['Alice', 'Bob', 'Craig', 'Dave', 'Eve'])
dict_values([80, 95, 56, 28, 75])
dict_items([('Alice', 80), ('Bob', 95), ('Craig', 56), ('Dave', 28), ('Eve', 75)])
```

The exercises below use if-statements.

## ✍️ Optional exercise/ homework: removing duplicates

- This list contains duplicates:

```
measurements = [2, 2, 1, 17, 3, 3, 2, 1, 13, 14, 17, 14, 4]
```

- Write a function which removes duplicates from the list and sorts the list. In this case it would produce:

```
[1, 2, 3, 4, 13, 14, 17]
```

## ✔ Solution 1 (longer but hopefully easier to understand)

The function `sorted` sorts a sequence but it creates a new sequence. This is useful if you need a sorted result without changing the original sequence.

We could have achieved the same result with `list.sort()`.

```python
def remove_duplicates_and_sort(sequence):
    new_sequence = []
    for element in sequence:
        if element not in new_sequence:
            new_sequence.append(element)
    return sorted(new_sequence)
```

## ✔ Solution 2 (more compact)

Converting to set removes duplicates. Then we convert back to list:

```python
def remove_duplicates_and_sort(sequence):
    new_sequence = list(set(sequence))
    return sorted(new_sequence)
```

## ✍️ Optional exercise/ homework: counting how often an item appears

- Back to our list with duplicates:

```
measurements = [2, 2, 1, 17, 3, 3, 2, 1, 13, 14, 17, 14, 4]
```

- Your goal is to write a function which will return a dictionary mapping each number to how often it appears. In this case it would produce:

```
{2: 3, 1: 2, 17: 2, 3: 2, 13: 1, 14: 2, 4: 1}
```

## ✔ Solution 1 (longer but hopefully easier to understand)

```python
def how_often(sequence):
    counts = {}
    for element in sequence:
        if element in counts:
            counts[element] += 1
        else:
            counts[element] = 1
    return counts
```

## ✔ Solution 2 (more compact)

The point of this solution is to show that for such common operations, ready-made functions and objects already exist and is is worth to check out the documentation about the collections module.

```python
from collections import Counter, defaultdict


def how_often_alternative1(sequence):
    return dict(Counter(sequence))


def how_often_alternative2(sequence):
    counts = defaultdict(int)
    for element in sequence:
        counts[element] += 1
    return dict(counts)
```

# Plotting with Vega-Altair

## ❗ Objectives

- Be able to create simple plots with Vega-Altair and tweak them
- Know how to look for help
- Reading data with Pandas from disk or a web resource
- Know how to tweak example plots from a gallery for your own purpose
- We will build up this notebook (spoiler alert!)

## Repeatability/reproducibility

From Claus O. Wilke: "Fundamentals of Data Visualization":

> One thing I have learned over the years is that automation is your friend. I think figures should be autogenerated as part of the data analysis pipeline (which should also be automated), and they should come out of the pipeline ready to be sent to the printer, no manual post-processing needed.

- **Try to minimize manual post-processing**. This could bite you when you need to regenerate 50 figures one day before submission deadline or regenerate a set of figures after the person who created them left the group.
- There is not the one perfect language and **not the one perfect library** for everything.
- Within Python, many libraries exist:
    - Vega-Altair: declarative visualization, statistics built in
    - Matplotlib: probably the most standard and most widely used
    - Seaborn: high-level interface to Matplotlib, statistical functions built in
    - Plotly: interactive graphs
    - Bokeh: also here good for interactivity
    - plotnine: implementation of a grammar of graphics in Python, it is based on ggplot2
    - ggplot: R users will be more at home
    - PyNGL: used in the weather forecast community
    - K3D: Jupyter Notebook extension for 3D visualization
    - Mayavi: 3D scientific data visualization and plotting in Python
    - …
- Two main families of libraries: procedural (e.g. Matplotlib) and declarative (e.g. Vega-Altair).

## Why are we starting with Vega-Altair?

- Concise and powerful
- "Simple, friendly and consistent API" allows us to focus on the data visualization part and get started without too much Python knowledge
- The way it **combines visual channels with data columns** can feel intuitive
- Interfaces very nicely with Pandas
- Easy to change figures
- Good documentation
- Open source
- Makes it easy to save figures in a number of formats (svg, png, html)
- Easy to save interactive visualizations to be used in websites

## Example data: Weather data from two Norwegian cities

We will experiment with some example weather data obtained from Norsk KlimaServiceSenter, Meteorologisk institutt (MET) (CC BY 4.0). The data is in CSV format (comma-separated values) and contains daily and monthly weather data for two cities in Norway: Oslo and Tromsø. You can browse the data here in the lesson repository.

We will use the Pandas library to read the data into a dataframe.

Pandas can read from and write to a large set of formats (overview of input/output functions and formats). We will load a CSV file directly from the web. Instead of using a web URL we could use a local file name instead.

Pandas dataframes are a great data structure for **tabular data** and tabular data turns out to be a great input format for data visualization libraries. Vega-Altair understands Pandas dataframes and can plot them directly.

## Reading data into a dataframe

We can try this together in a notebook: Using Pandas we can **merge, join, concatenate, and compare** dataframes, see https://pandas.pydata.org/pandas-docs/stable/user_guide/merging.html.

Let us try to **concatenate** two dataframes: one for Tromsø weather data (we will now load monthly values) and one for Oslo:

```python
import pandas as pd

url_prefix = "https://raw.githubusercontent.com/coderefinery/data-visualization-python/main/data/"

data_tromso = pd.read_csv(url_prefix + "tromso-monthly.csv")
data_oslo = pd.read_csv(url_prefix + "oslo-monthly.csv")

data_monthly = pd.concat([data_tromso, data_oslo], axis=0)

# let us print the combined result
data_monthly
```

Before plotting the data, there is a problem which we may not see yet: Dates are not in a standard date format (YYYY-MM-DD). We can fix this:
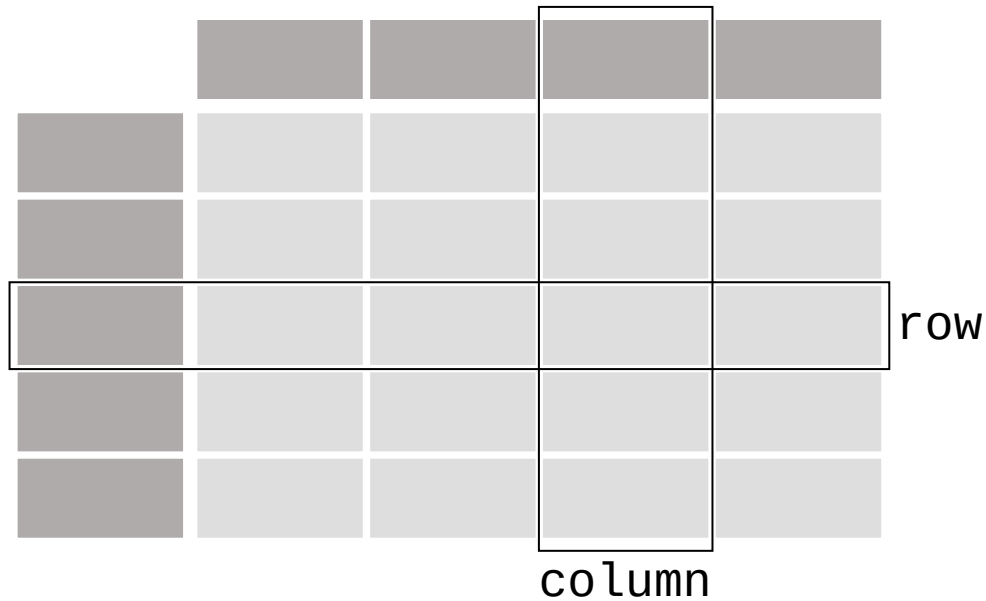
```python
# replace mm.yyyy to date format
data_monthly["date"] = pd.to_datetime(list(data_monthly["date"]), format="%m.%Y")
```

With Pandas it is possible to do a lot more (adjusting missing values, fixing inconsistencies, changing format).

**What is in a dataframe?**

The name pandas is derived from the term "**pan**el **da**ta".

# DataFrame



*A pandas dataframe object is composed of rows and columns.*

Let us explore these together in the notebook (run these in separate cells):

```python
# print an overview of the dataset
data_monthly

# print the first 5 rows
data_monthly.head()

# print the last 5 rows
data_monthly.tail()

# print all column titles - no parentheses here
data_monthly.columns

# show which data types were detected
data_monthly.dtypes

# print table dimensions - no parentheses here
data_monthly.shape

# print one column
data_monthly["max temperature"]

# get some statistics
data_monthly["max temperature"].describe()

# what was the maximum temperature?
data_monthly["max temperature"].max()

# print all rows where max temperature was above 20
data_monthly[data_monthly["max temperature"] > 20.0]
```

## Where to learn more about pandas

Pandas is extremely powerful and there is a lot that can be done and there are great resources to explore more:

- Getting started guide (including tutorials and a 10 minute flash intro)
- 10 minutes to pandas tutorial
- Pandas documentation
- Cheatsheet
- Cookbook
- Data Carpentry lesson "Data Analysis and Visualization in Python for Ecologists" (useful not only for ecologists)

## Plotting the data

Now let's plot the data. We will start with a plot that is not optimal and then we will explore and improve a bit as we go:

```
import altair as alt

alt.Chart(data_monthly).mark_bar().encode(
    x="date",
    y="precipitation",
    color="name",
)
```
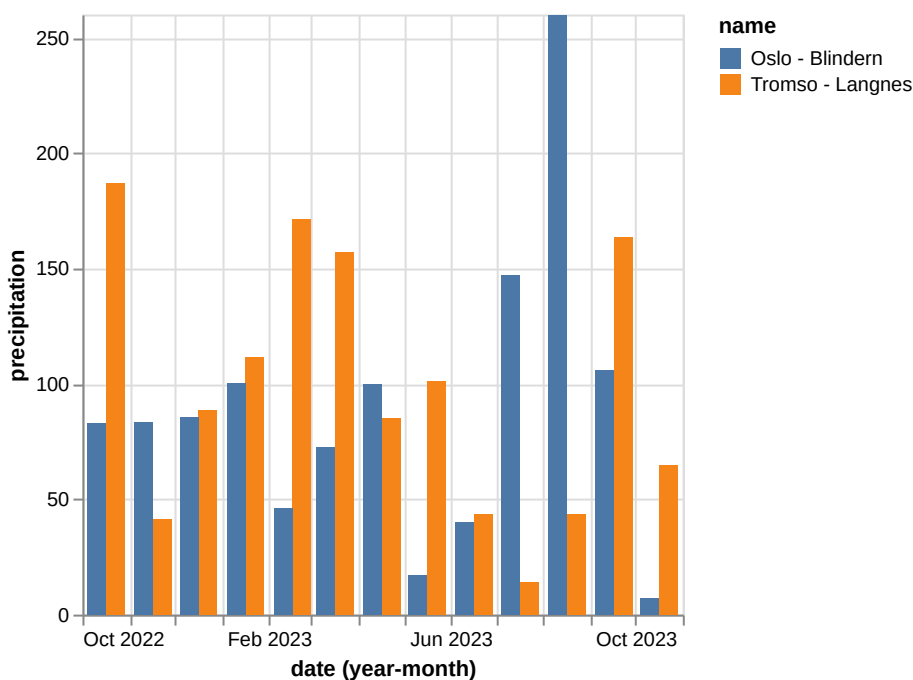


*Monthly precipitation for the cities Oslo and Tromsø over the course of a year.*

> 💬 **Let us pause and explain the code**
>
> - `alt` is a short-hand for `altair` which we imported on top of the notebook
> - `Chart()` is a function defined inside `altair` which takes the data as argument
> - `mark_bar()` is a function that produces bar charts
> - `encode()` is a function which encodes data columns to **visual channels**
>
> Observe how we connect (encode) **visual channels** to data columns:
>
> - x-coordinate with "date"
> - y-coordinate with "precipitation"
> - color with "name" (name of weather station; city)

We can improve the plot by giving Vega-Altair a bit more information that the x-axis is **temporal** (T) and that we would like to see the year and month (yearmonth):

```
alt.Chart(data_monthly).mark_bar().encode(
    x="yearmonth(date):T",
    y="precipitation",
    color="name",
)
```

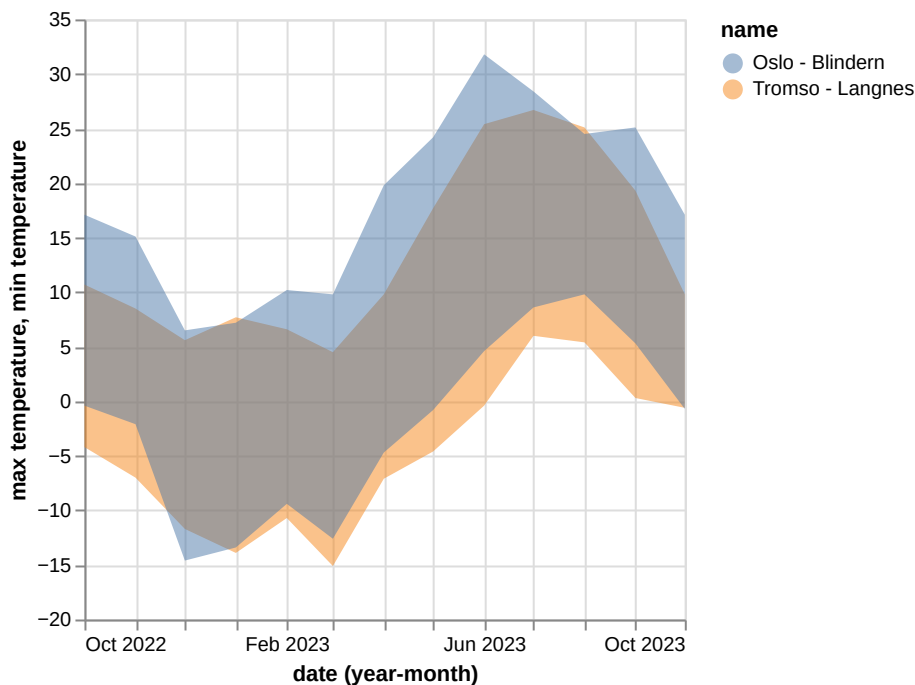Apart from T (temporal), there are other encoding data types:

- Q (quantitative)
- O (ordinal)
- N (nominal)
- T (temporal)
- G (geojson)



*Monthly precipitation for the cities Oslo and Tromsø over the course of a year.*

Let us improve the plot with another one-line change:

```
alt.Chart(data_monthly).mark_bar().encode(
    x="yearmonth(date):T",
    y="precipitation",
    color="name",
    column="name",
)
```

*Monthly precipitation for the cities Oslo and Tromsø over the course of a year with with both cities plotted side by side.*

With another one-line change we can make the bar chart stacked:

```python
alt.Chart(data_monthly).mark_bar().encode(
    x="yearmonth(date):T",
    y="precipitation",
    color="name",
    xOffset="name",
)
```



*Monthly precipitation for the cities Oslo and Tromsø over the course of a year plotted as stacked bar chart.*

This is not publication-quality yet but a really good start!

Let us try one more example where we can nicely see how Vega-Altair is able to map visual channels to data columns:

```
alt.Chart(data_monthly).mark_area(opacity=0.5).encode(
    x="yearmonth(date):T",
    y="max temperature",
    y2="min temperature",
    color="name",
)
```
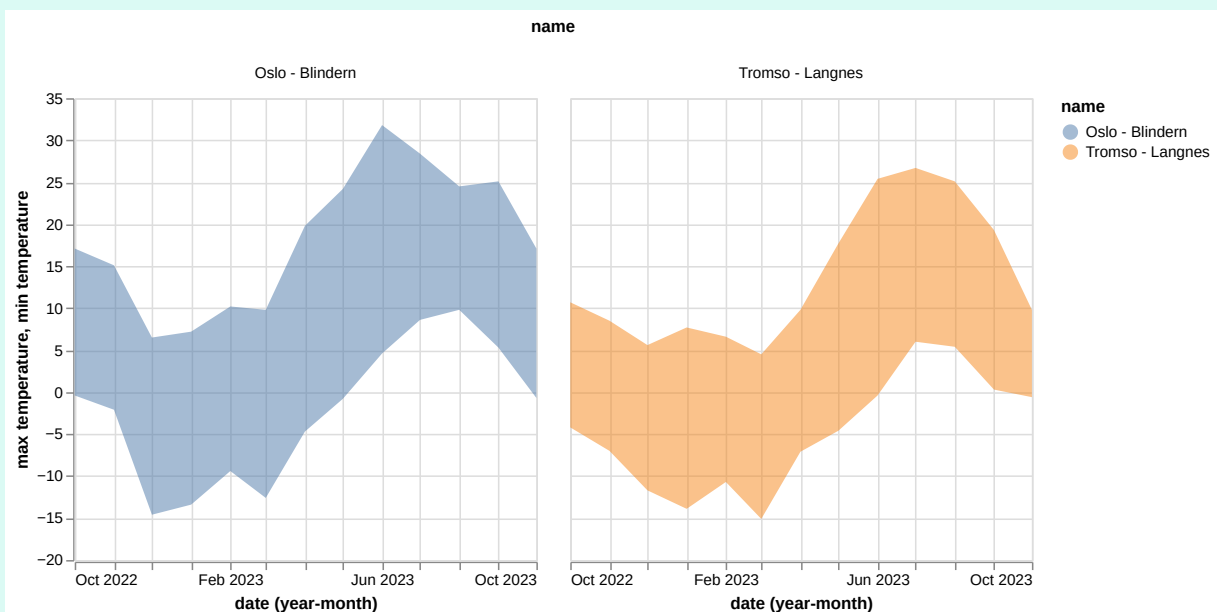


*Monthly temperature ranges for two cities in Norway.*

## Exercise: Using visual channels to re-arrange plots

🖊 **Exercise Plotting-1: Using visual channels to re-arrange plots**

1. Try to reproduce the above plots if they are not already in your notebook.
2. Above we have plotted the monthly precipitation for two cities side by side using a stacked plot. Try to arrive at the following plot where months are along the y-axis and the precipitation amount is along the x-axis:

3. Ask the Internet or AI how to change the axis title from "precipitation" to "Precipitation (mm)".

4. Modify the temperature range plot to show the temperature ranges for the two cities side by side like this:



## ✔ Solution

1. Copy-paste code blocks from above.
2. Basically we switched x and y:

```
alt.Chart(data_monthly).mark_bar().encode(
    y="yearmonth(date):T",
    x="precipitation",
    color="name",
    yOffset="name",
)
```

3. This can be done with the following modification:

```
alt.Chart(data_monthly).mark_bar().encode(
    y="yearmonth(date):T",
    x=alt.X("precipitation").title("Precipitation (mm)"),
    color="name",
    yOffset="name",
)
```

4. We added one line:

```
alt.Chart(data_monthly).mark_area(opacity=0.5).encode(
    x="yearmonth(date):T",
    y="max temperature",
    y2="min temperature",
    color="name",
    column="name",
)
```

## More fun with visual channels

Now we will try to **plot the daily data and look at snow depths**. We first read and concatenate two datasets:

```
url_prefix = "https://raw.githubusercontent.com/coderefinery/data-visualization-
python/main/data/"

data_tromso = pd.read_csv(url_prefix + "tromso-daily.csv")
data_oslo = pd.read_csv(url_prefix + "oslo-daily.csv")

data_daily = pd.concat([data_tromso, data_oslo], axis=0)
```

We adjust the data a bit:

```
# replace dd.mm.yyyy to date format
data_daily["date"] = pd.to_datetime(list(data_daily["date"]), format="%d.%m.%Y")

# we are here only interested in the range december to may
data_daily = data_daily[
    (data_daily["date"] > "2022-12-01") & (data_daily["date"] < "2023-05-01")
]
```

Now we can plot the snow depths for the months December to May for the two cities:

```
alt.Chart(data_daily).mark_bar().encode(
    x="date",
    y="snow depth",
    column="name",
)
```
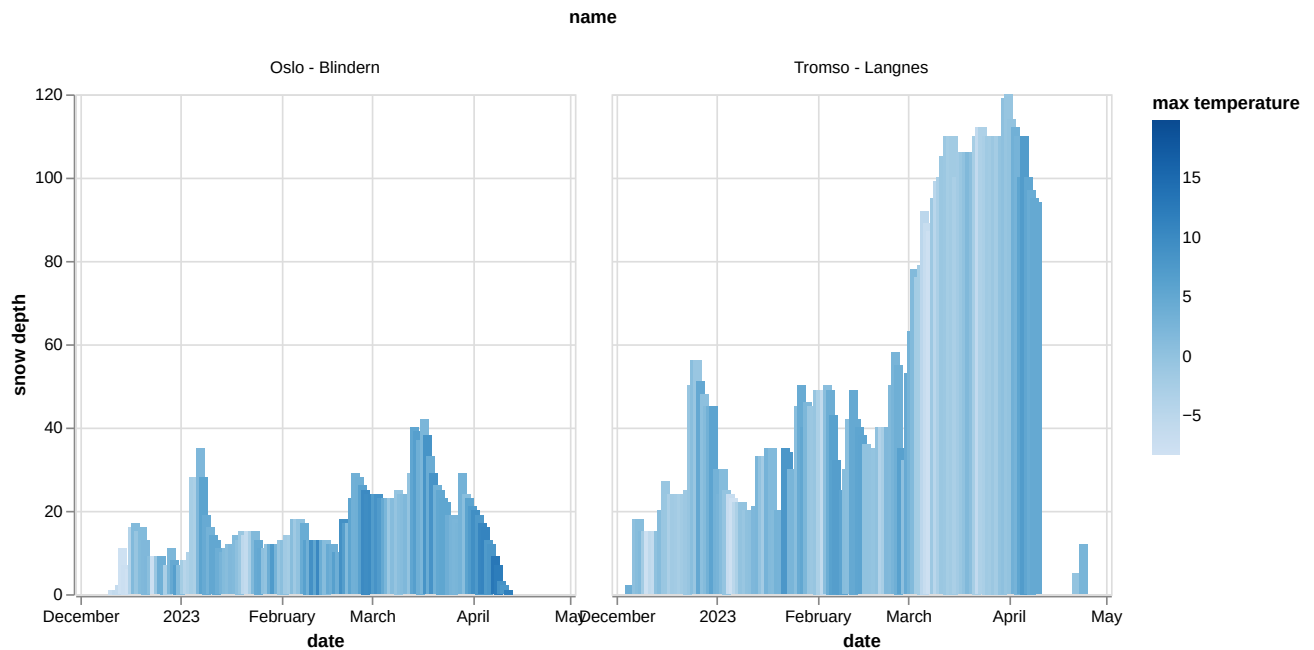


*Snow depth (in cm) for the months December 2022 to May 2023 for two cities in Norway.*

What happens if we try to color the plot by the "max temperature" values?

```
alt.Chart(data_daily).mark_bar().encode(
    x="date",
    y="snow depth",
    color="max temperature",
    column="name",
)
```
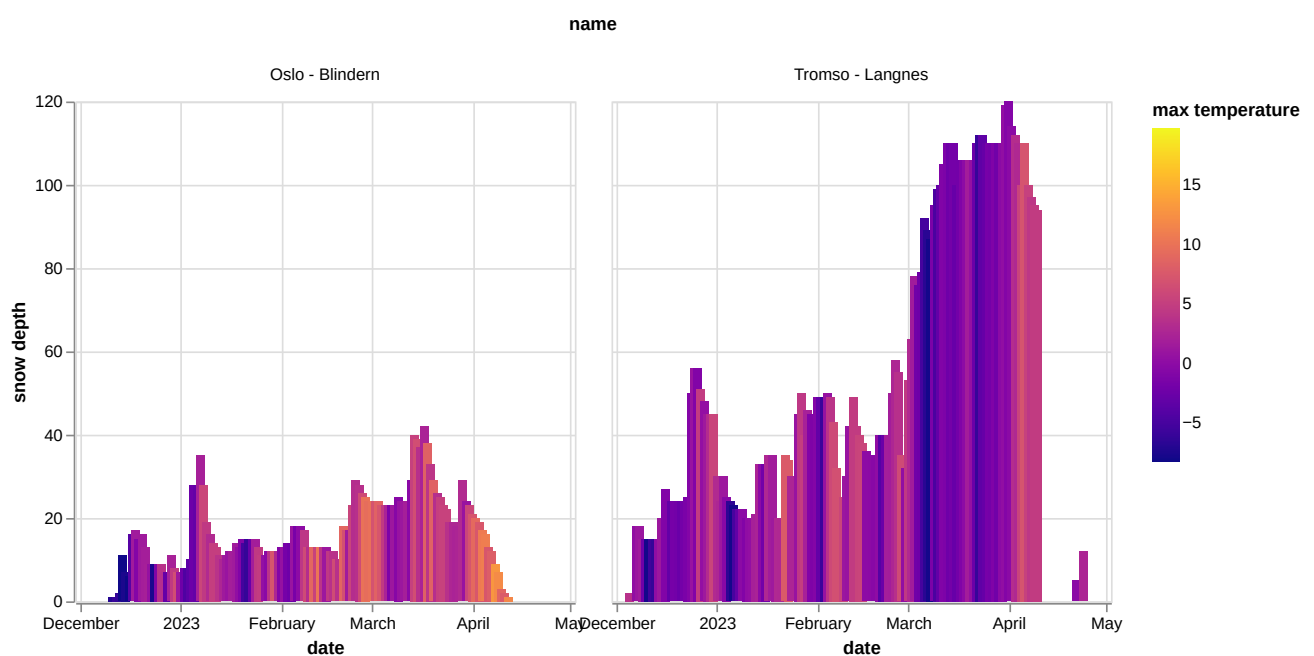
The result looks neat:

*Snow depth (in cm) for the months December 2022 to May 2023 for two cities in Norway. Colored by daily max temperature.*

We can change the color scheme ([available color schemes](#)):

```
alt.Chart(data_daily).mark_bar().encode(
    x="date",
    y="snow depth",
    color=alt.Color("max temperature").scale(scheme="plasma"),
    column="name",
)
```
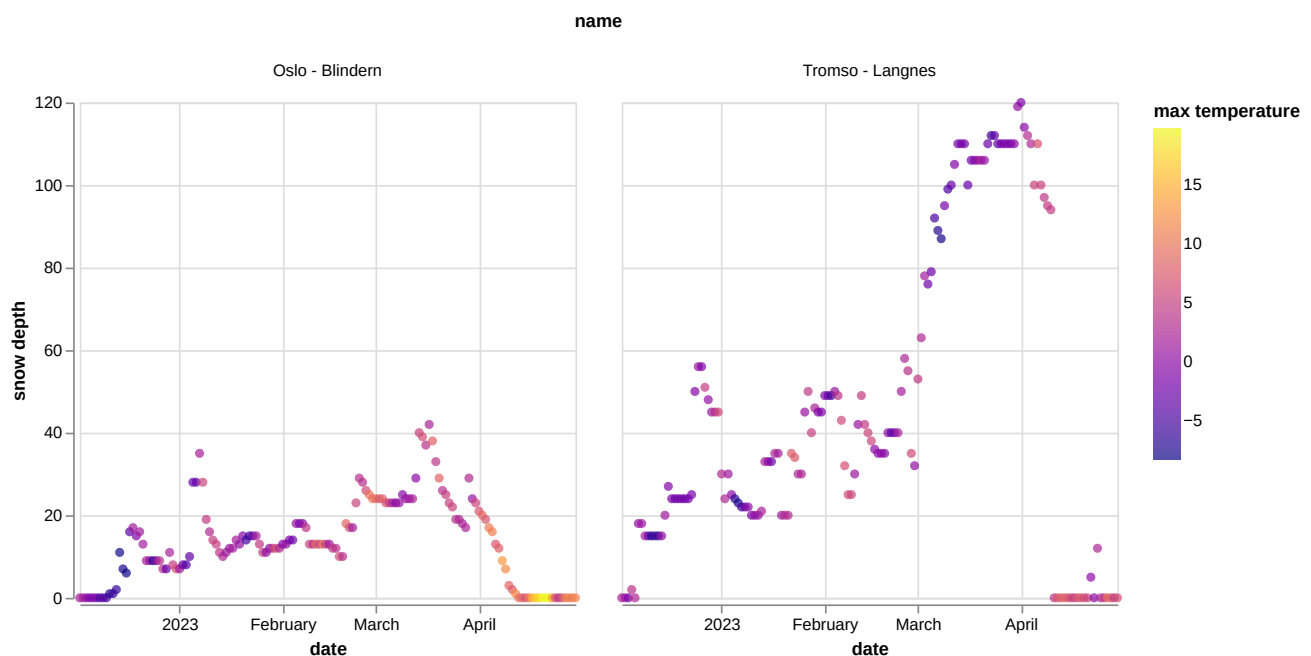
With the following result:



*Snow depth (in cm) for the months December 2022 to May 2023 for two cities in Norway. Colored by daily max temperature. Warmer days are often followed by reduced snow depth.*

Let's try one more change to show that we can experiment with different plot types by changing `mark_bar()` to something else, in this case `mark_circle()`:

```
alt.Chart(data_daily).mark_circle().encode(
    x="date",
    y="snow depth",
    color=alt.Color("max temperature").scale(scheme="plasma"),
    column="name",
)
```



*Snow depth (in cm) for the months December 2022 to May 2023 for two cities in Norway. Colored by daily max temperature. Warmer days are often followed by reduced snow depth.*

## Themes

In Vega-Altair you can change the theme and select from a long list of themes. On top of your notebook try to add:

```
alt.themes.enable('dark')
```

Then re-run all cells. Later you can try some other themes such as:

- `fivethirtyeight`
- `latimes`
- `urbaninstitute`

You can even define your own themes!

## Exercise: Anscombe's quartet

Save the following data as `example.csv` (you can do this directly from JupyterLab; this data is the Anscombe's quartet):

```
dataset,x,y
I,10.0,8.04
I,8.0,6.95
I,13.0,7.58
I,9.0,8.81
I,11.0,8.33
I,14.0,9.96
I,6.0,7.24
I,4.0,4.26
I,12.0,10.84
I,7.0,4.82
I,5.0,5.68
II,10.0,9.14
II,8.0,8.14
II,13.0,8.74
II,9.0,8.77
II,11.0,9.26
II,14.0,8.1
II,6.0,6.13
II,4.0,3.1
II,12.0,9.13
II,7.0,7.26
II,5.0,4.74
III,10.0,7.46
III,8.0,6.77
III,13.0,12.74
III,9.0,7.11
III,11.0,7.81
III,14.0,8.84
III,6.0,6.08
III,4.0,5.39
III,12.0,8.15
III,7.0,6.42
III,5.0,5.73
IV,8.0,6.58
IV,8.0,5.76
IV,8.0,7.71
IV,8.0,8.84
IV,8.0,8.47
IV,8.0,7.04
IV,8.0,5.25
IV,19.0,12.5
IV,8.0,5.56
IV,8.0,7.91
IV,8.0,6.89
```

## 🧹 Exercise Plotting-2: Read and plot a CSV file

- Save the above CSV file to disk as `example.csv` in the same folder where you run JupyterLab. We recommend to create the file in the JupyterLab interface.
- Plot the data using `mark_point`.
- Your goal is to arrive at four plots for the four data sets, all side by side.
- If you have time, try to customize the plot.

```python
# we don't need to import again but just in case you started here
import pandas as pd

data = pd.read_csv("example.csv")

alt.Chart(data).mark_point().encode(
    x="x",
    y="y",
    color="dataset",
    column="dataset",
    )
```

Here is a more advanced example where the four plots are arranged in a 2 x 2 grid:

```python
def create_chart(data, number):
    chart = (
        alt.Chart(data)
        .transform_filter(alt.datum.dataset == number)
        .mark_point()
        .encode(x="x", y="y")
    )
    return chart


chart1 = create_chart(data_example, "I")
chart2 = create_chart(data_example, "II")
chart3 = create_chart(data_example, "III")
chart4 = create_chart(data_example, "IV")

chart = alt.vconcat(
    alt.hconcat(chart1, chart2),
    alt.hconcat(chart3, chart4),
)

chart.display()
```

- Browse a number of example galleries to help you choose the library that fits best your work/style.
- Minimize manual post-processing and try to script all steps.
- CSV (comma-separated values) files are often a good format to store the data that we wish to plot.
- Read the data into a Pandas dataframe and then plot it with Vega-Altair where you connect data columns to visual channels.

# Tidy data and dealing with messy data

In the previous episode we read data from nicely formatted "plain" text files. But sometimes the data is in a spreadsheet or in less nicely formatted text files. In this episode we will discuss strategies for how to work with these.

## Importing data from spreadsheets

We will create a spreadsheet with the following content (only columns A and B; the actual content does not have to be exactly the same):

| | A | B | C | D |
|---|---|---|---|---|
| 1 | weekday | number of coffees | | |
| 2 | monday | 3 | | |
| 3 | tuesday | 2 | | |
| 4 | wednesday | 3 | | some side note |
| 5 | thursday | 4 | | and some color |
| 6 | friday | 2 | | |
| 7 | saturday | 3 | | |
| 8 | sunday | 3 | | |
| 9 | | | | |

*Example spreadsheet with a side note.*

Copy this also to the second sheet and for demonstration purpose add some side-notes to the second sheet and also color one or two cells (some people like to give some meaning to cells using color).

Save the spreadsheet as `experiment.xls`.

Now we will together try to read and inspect both sheets in the Jupyter Notebook:
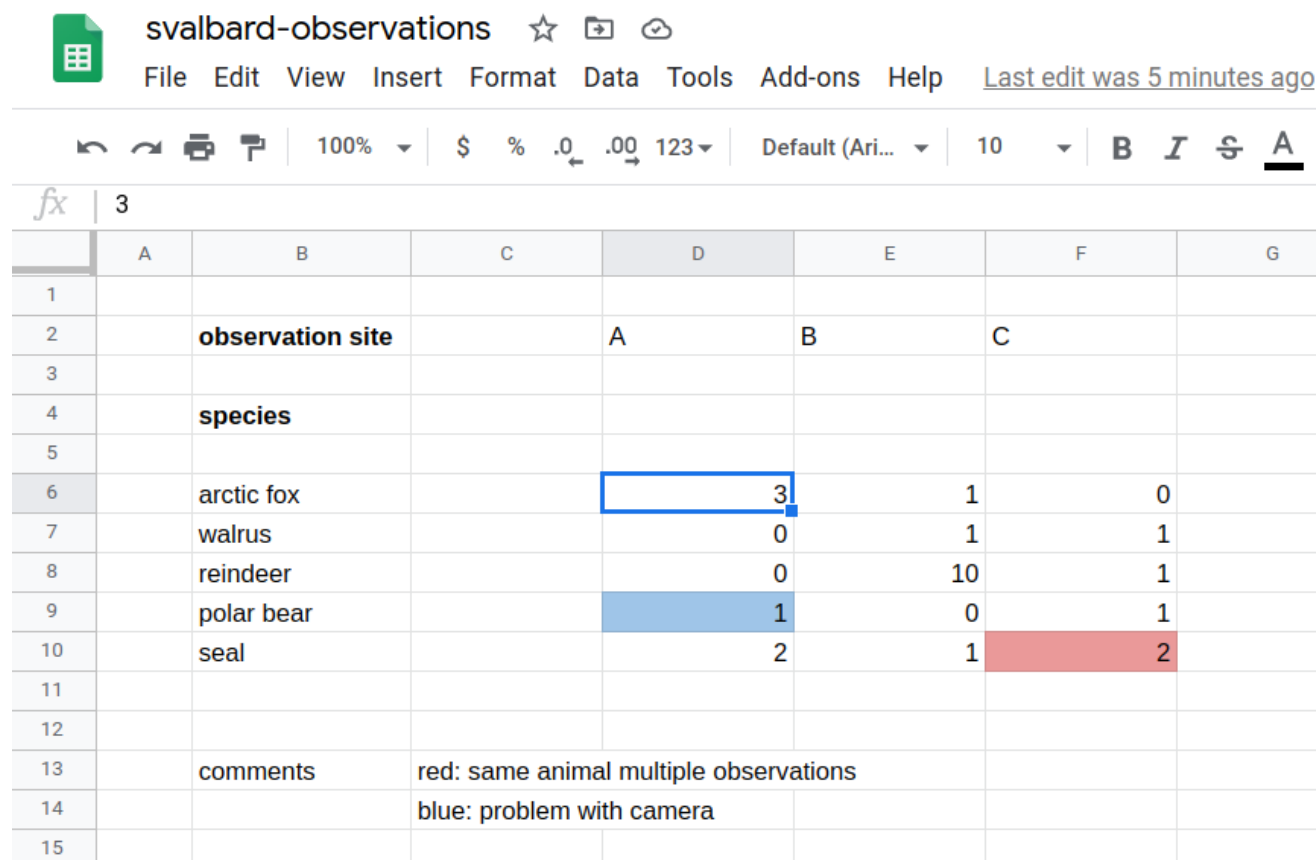
```python
import pandas as pd

data = pd.read_excel('experiment.xls', sheet_name="Sheet1")
data

data = pd.read_excel('experiment.xls', sheet_name="Sheet2")
data
```

## Tidy data



*Example spreadsheet (this is a phantasy dataset, apologies to biology students/researchers - this is not my domain).*

How should we arrange the data?

| Species | Observation sites |
|---|---|
| arctic fox | A, B |
| walrus | B, C |
| reindeer | B, C |
| polar bear | A, C |
| seal | A, B, C |

*Attempt 1: Not great since we need to somehow divide at the comma. How should we deal with multiple sightings?*

| Species | Observation site A | Observation site B | Observation site C |
|---|---|---|---|
| arctic fox | 3 | 1 | 0 |
| walrus | 0 | 1 | 1 |
| reindeer | 0 | 10 | 1 |
| polar bear | 1 | 0 | 1 |
| seal | 2 | 1 | 2 |

*Attempt 2: Adding observation sites will force us to add columns.*

| Species | arctic fox | walrus | reindeer | polar bear | seal |
|---|---|---|---|---|---|
| **Observation site A** | 3 | 0 | 0 | 1 | 2 |
| **Observation site B** | 1 | 1 | 10 | 0 | 1 |
| **Observation site C** | 0 | 1 | 1 | 1 | 2 |

*Attempt 3: Adding species will force us to add columns.*

| Species | Observation site | Number of sightings |
|---|---|---|
| arctic fox | A | 3 |
| arctic fox | B | 1 |
| walrus | B | 1 |
| walrus | C | 1 |
| reindeer | B | 10 |
| reindeer | C | 1 |
| polar bear | A | 1 |
| polar bear | C | 1 |
| seal | A | 2 |
| seal | B | 1 |
| seal | C | 2 |

*Tidy data format: Columns are variables, rows are observations/measurements. Easy to add new species and sites.*

> ❶ **Tidy data format**
>
> - Hadley Wickham: Tidy Data

- Columns are variables
- Rows are observations/measurements
- "Long form"
- Order does not matter
- **Easy to extend** with more species and more sites without modifying the code
- **Structure for storing data** - this does not mean that this is ideal for tables in presentations or publications
- It is possible to convert between wide form and long form and back (e.g. using `pandas.melt` or `pandas.pivot` ), see this example notebook

## Use a standard format

```
Species,Observation site,Number of sightings
arctic fox,A,3
arctic fox,B,1
walrus,B,1
walrus,C,1
reindeer,B,10
reindeer,C,1
polar bear,A,1
polar bear,C,1
seal,A,2
seal,B,1
seal,C,2
```

- **Use a format that is standard in your community, don't invent your own**
- CSV is often a good choice since most visualization tools can read CSV data

There are many more formats (adapted after Python for Scientific Computing):

| Name: | Human readable: | Space efficiency: | Arbitrary data: | Tidy data: | Array data: | Long term storage/sharing: |
|---|---|---|---|---|---|---|
| CSV | ✅ | ❌ | ❌ | ✅ | 🟨 | ✅ |
| Feather | ❌ | ✅ | ❌ | ✅ | ❌ | ❌ |
| Parquet | ❌ | ✅ | 🟨 | ✅ | 🟨 | ✅ |
| NPY | ❌ | 🟨 | ❌ | ❌ | ✅ | ❌ |
| HDF5 | ❌ | ✅ | ❌ | ❌ | ✅ | ✅ |
| NetCDF | ❌ | ✅ | ❌ | ❌ | ✅ | ✅ |
| JSON | ✅ | ❌ | 🟨 | ❌ | ❌ | ✅ |
| GeoJSON | ✅ | ❌ | 🟨 | ❌ | ❌ | ✅ |
| Excel | ❌ | ❌ | ❌ | 🟨 | ❌ | 🟨 |
| Graph formats | 🟨 | 🟨 | ❌ | ❌ | ❌ | ✅ |

| Name: | Human readable: | Space efficiency: | Arbitrary data: | Tidy data: | Array data: | Long term storage/sharing: |
|---|---|---|---|---|---|---|
| SQL | ❌ | 🟨 | ❌ | ❌ | ❌ | ❌ |

> **❗ Note**
>
> - ✅ : Good
> - 🟨 : Ok / depends on a case
> - ❌ : Bad

## Data cleaning

Often we want to visualize data sets with inconsistent or missing entries:

```
Date,Organization,Number of participants
2020-09-27,UiT,20
Oct 10 2020,UiT Norges arktiske universitet,15
"Nov. 11, 2020",UiT The Arctic University of Norway,40
2020-12-12,UiT The Arctic University of Norway,-
```

Data cleaning is a bit outside the scope of this course (although we have done some of this in the pandas episode) but still good to know:

- There are tools to clean and merge inconsistent data sets (e.g. OpenRefine, see also this Data Carpentry lesson)
- This does not have to be done manually

## Customizing plots

> **❗ Objectives**
>
> - Know where to look to find out how to tweak plots
> - Start with a relatively simple example and build up more and more features
> - See the process of going from a raw plot towards a publication-ready plot
> - We will build up this notebook (spoiler alert!)

### Loading and plotting a dataset

In this lesson will work with one of the Gapminder datasets.

Let us together read and plot the data and then we explain what is happening and we will improve the figure together. First we read and inspect the data:
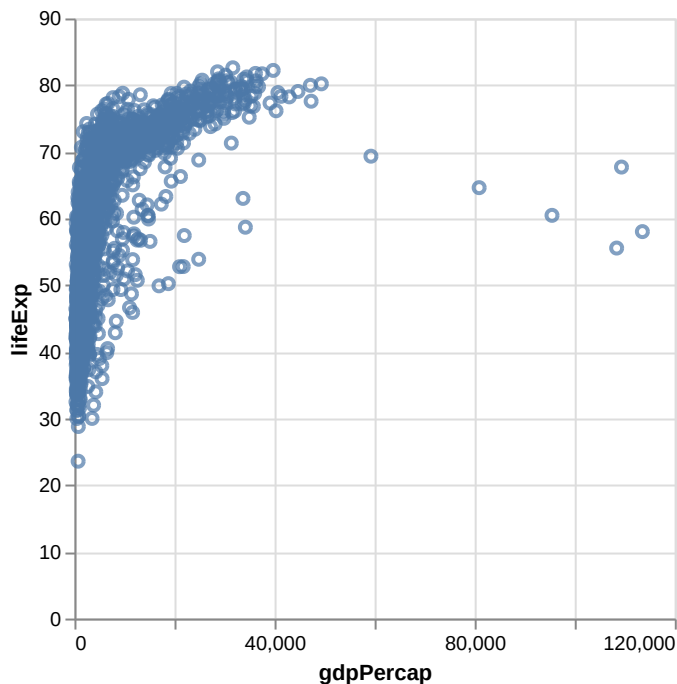
```
# import necessary libraries
import altair as alt
import pandas as pd

# read the data
url_prefix = "https://raw.githubusercontent.com/plotly/datasets/master/"
data = pd.read_csv(url_prefix + "gapminder_with_codes.csv")

# print overview of the dataset
data
```

With very few lines we can get the first raw plot:

```
alt.Chart(data).mark_point().encode(
    x="gdpPercap",
    y="lifeExp",
)
```



*First raw plot with all countries and all years.*

Observe how we connect (encode) **visual channels** to data columns:

- x-coordinate with "gdpPercap"
- y-coordinate with "lifeExp"

The following code would have the same effect but the above version might be easier to read:

```
alt.Chart(data).mark_point().encode(x="gdpPercap", y="lifeExp")
```

## Filtering data

In Vega-Altair we can chain functions. Let us add two more functions: The first will apply a filter, the second will make the plot interactive:

```
alt.Chart(data).mark_point().encode(
    x="gdpPercap",
    y="lifeExp",
).transform_filter(alt.datum.year == 2007).interactive()
```
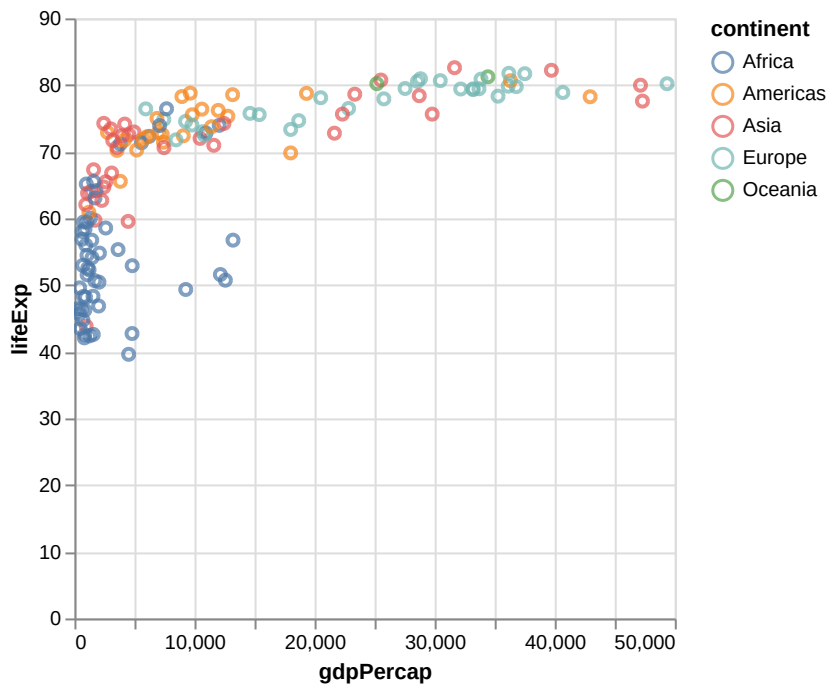


*Now we only keep the year 2007.*

Alternatively, we could have filtered the data before plotting using pandas.

## Using color as additional channel

A very neat feature of Vega-Altair is that it is easy to add and modify visual channels. Let us try to add one more so that we do something with the "continent" data column:

```
alt.Chart(data).mark_point().encode(
    x="gdpPercap",
    y="lifeExp",
    color="continent",
).transform_filter(alt.datum.year == 2007).interactive()
```
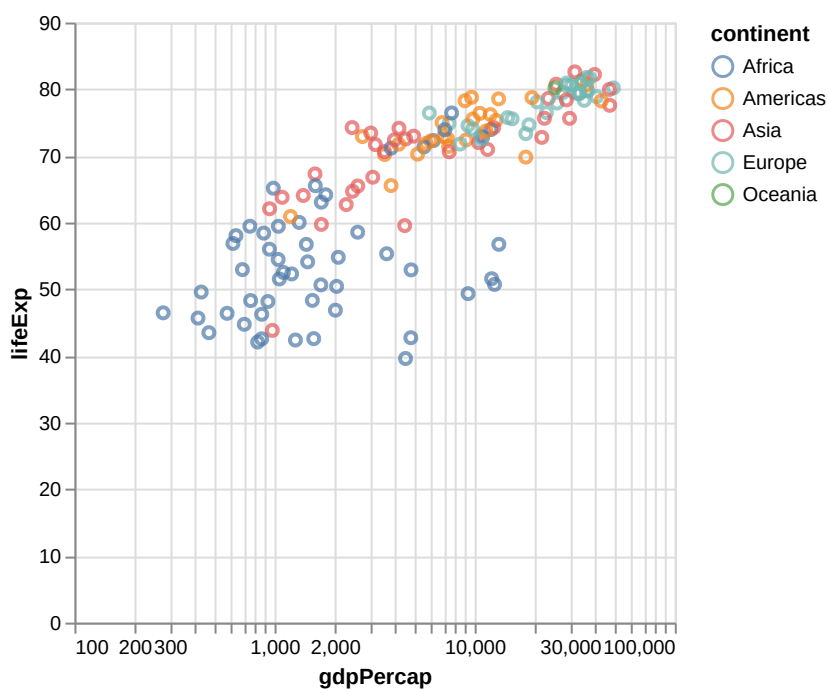
*Using different colors for different continents.*

## Changing to log scale

For this data set we will get a better insight when switching the x-axis from linear to log scale:
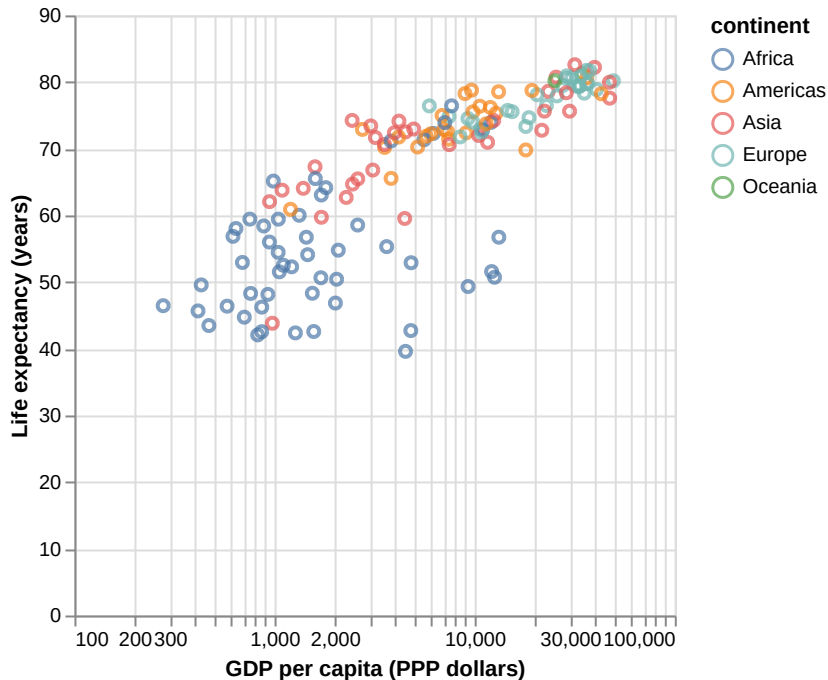
```python
alt.Chart(data).mark_point().encode(
    x=alt.X("gdpPercap").scale(type="log"),
    y=alt.Y("lifeExp"),
    color="continent",
).transform_filter(alt.datum.year == 2007).interactive()
```



*Changing the x axis to log scale.*

## Improving axis titles

```
alt.Chart(data).mark_point().encode(
    x=alt.X("gdpPercap").scale(type="log").title("GDP per capita (PPP dollars)"),
    y=alt.Y("lifeExp").title("Life expectancy (years)"),
    color="continent",
).transform_filter(alt.datum.year == 2007).interactive()
```



*Improving the axis titles.*

## Faceted charts

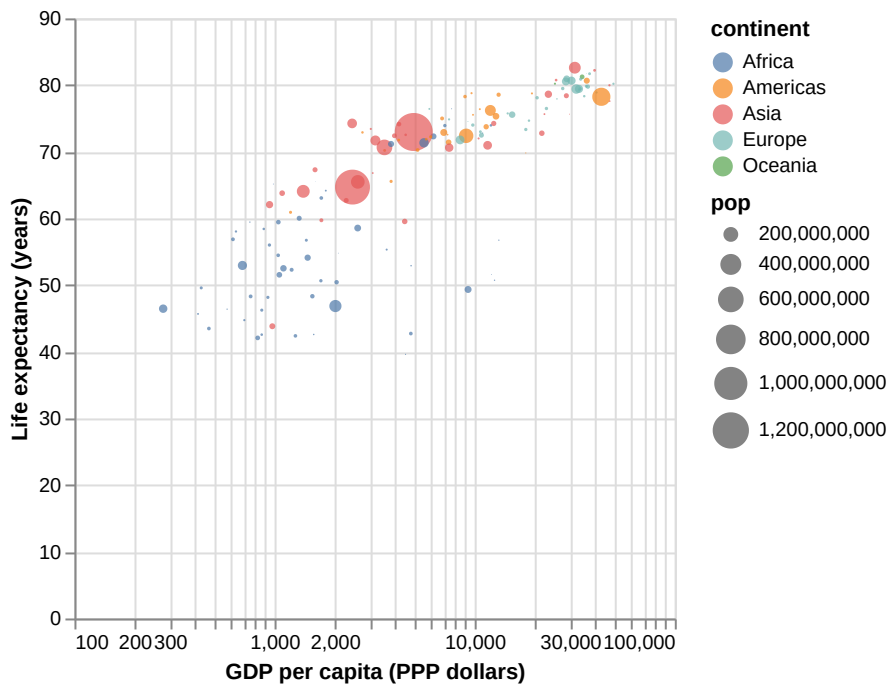To see what faceted charts are and how easy it is to do this, add the following line:

```
alt.Chart(data).mark_point().encode(
    x=alt.X("gdpPercap").scale(type="log").title("GDP per capita (PPP dollars)"),
    y=alt.Y("lifeExp").title("Life expectancy (years)"),
    color="continent",
    row="continent",
).transform_filter(alt.datum.year == 2007).interactive()
```

Guess what happens when you change `row="continent"` to `column="continent"` ?

## Changing from points to circles

Let us add one more visual channel, mapping size of the circle to the population size of a country:

```
alt.Chart(data).mark_circle().encode(
    x=alt.X("gdpPercap").scale(type="log").title("GDP per capita (PPP dollars)"),
    y=alt.Y("lifeExp").title("Life expectancy (years)"),
    color="continent",
    size="pop",
).transform_filter(alt.datum.year == 2007).interactive()
```



*Circle sizes are proportional to population sizes.*

## Title and axis values

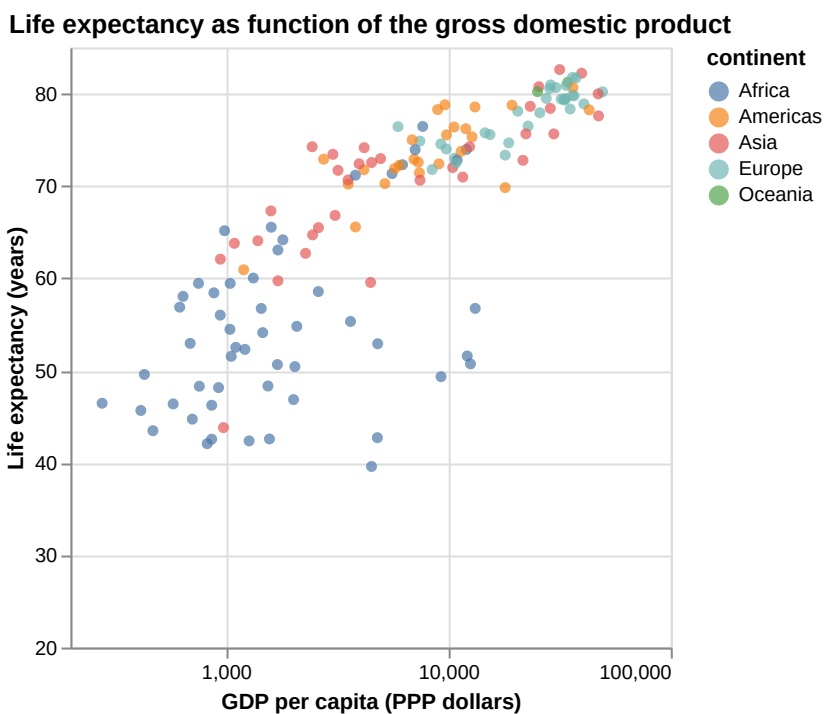In the next step we modify a number of things:

- We go back to the version where all circles have the same size
- Add figure title
- Modify axis domains to "zoom into" the interesting part of the plot
- Set axis values
- Change from `mark_point()` to `mark_circle()`
- Invoke `interactive()` in a separate step

```
chart = (
    alt.Chart(
        data,
        title=alt.Title("Life expectancy as function of the gross domestic product"),
    )
    .mark_circle()
    .encode(
        x=alt.X("gdpPercap", axis=alt.Axis(values=[100, 1000, 10000, 100000]))
        .scale(type="log", domain=(200, 100000))
        .title("GDP per capita (PPP dollars)"),
        y=alt.Y("lifeExp", axis=alt.Axis(values=[20, 30, 40, 50, 60, 70, 80]))
        .title("Life expectancy (years)")
        .scale(domain=(20, 85)),
        color="continent",
    )
    .transform_filter(alt.datum.year == 2007)
)

chart.interactive()
```



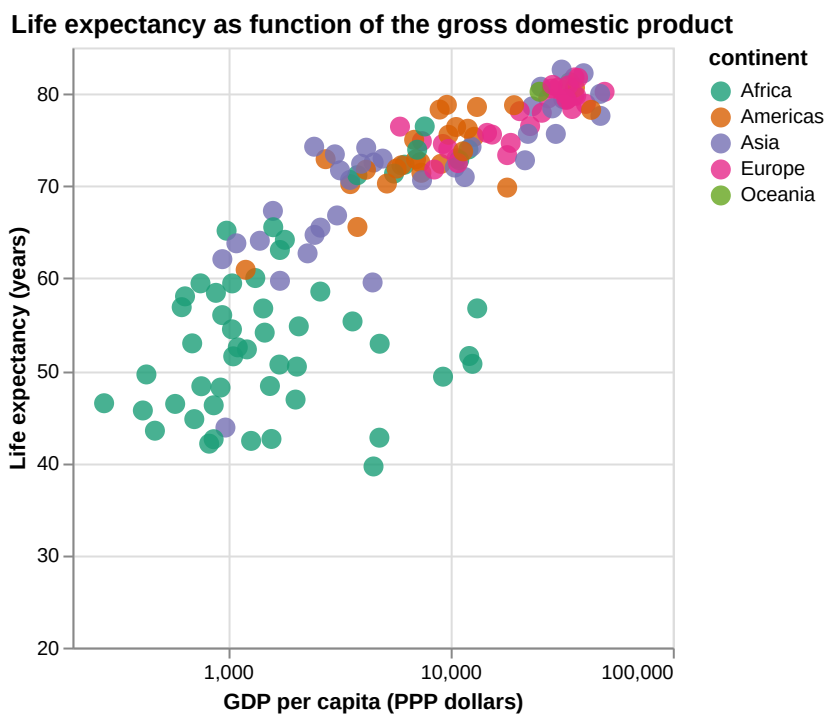*The plot is starting to look better!*

## Colors

In the next step we change the color scheme ([list of all schemes](#)), make the circles larger and slightly transparent:

```
chart = (
    alt.Chart(
        data,
        title=alt.Title("Life expectancy as function of the gross domestic product"),
    )
    .mark_circle(opacity=0.8, size=100.0)
    .encode(
        x=alt.X("gdpPercap", axis=alt.Axis(values=[100, 1000, 10000, 100000]))
        .scale(type="log", domain=(200, 100000))
        .title("GDP per capita (PPP dollars)"),
        y=alt.Y("lifeExp", axis=alt.Axis(values=[20, 30, 40, 50, 60, 70, 80]))
        .title("Life expectancy (years)")
        .scale(domain=(20, 85)),
        color=alt.Color("continent").scale(scheme="dark2"),
    )
    .transform_filter(alt.datum.year == 2007)
)

chart.interactive()
```



*The plot after adjusting circles and colors.*
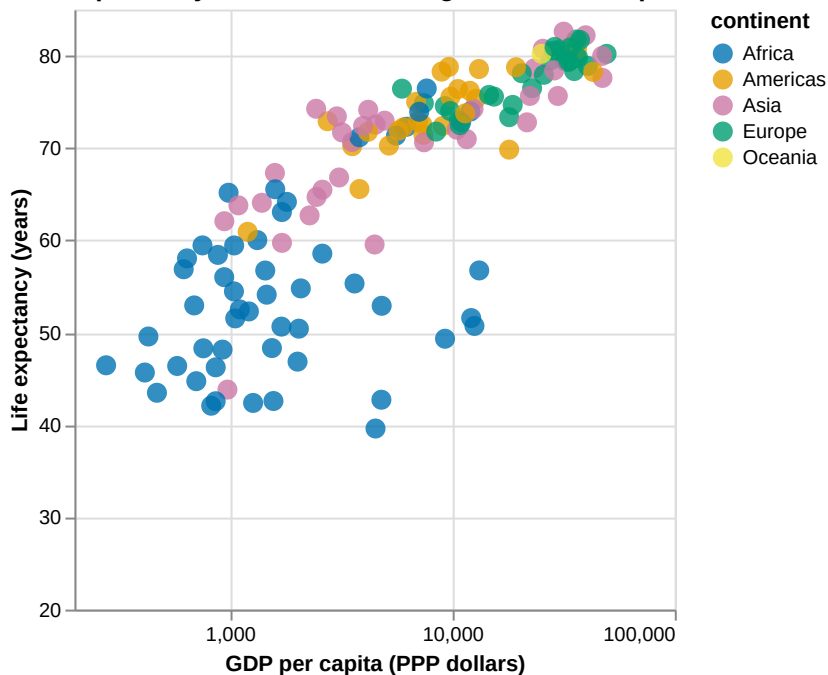
We can also define own colors:

```python
okabe_ito = [
    "#0072b2",
    "#e69f00",
    "#cc79a7",
    "#009e73",
    "#f0e442",
    "#000000",
    "#d55e00",
    "#56b4e9",
]
```

```python
chart = (
    alt.Chart(
        data,
        title=alt.Title("Life expectancy as function of the gross domestic product"),
    )
    .mark_circle(opacity=0.8, size=100.0)
    .encode(
        x=alt.X("gdpPercap", axis=alt.Axis(values=[100, 1000, 10000, 100000]))
        .scale(type="log", domain=(200, 100000))
        .title("GDP per capita (PPP dollars)"),
        y=alt.Y("lifeExp", axis=alt.Axis(values=[20, 30, 40, 50, 60, 70, 80]))
        .title("Life expectancy (years)")
        .scale(domain=(20, 85)),
        color=alt.Color("continent").scale(range=okabe_ito),
    )
    .transform_filter(alt.datum.year == 2007)
)

chart.interactive()
```



*Adjusting colors to those recommended by Okabe and Ito.*

💬 Why these colors?

## More tweaking towards a publication-ready figure

Let us add a subtitle and adjust sizing and positioning:

```python
chart = (
    alt.Chart(
        data,
        title=alt.Title(
            "Life expectancy as function of the gross domestic product",
            subtitle=[
                "Gross domestic product (GDP) per capita measures the value of everything",
                "produced in a country during a year, divided by the number of people.",
                "The unit is in purchasing power parities (PPP dollars), fixed to 2017 prices.",
                "Data is adjusted for inflation and differences in the cost of living between countries.",
            ],
        ),
    )
    .mark_circle(opacity=0.8, size=100.0)
    .encode(
        x=alt.X("gdpPercap", axis=alt.Axis(values=[100, 1000, 10000, 100000]))
        .scale(type="log", domain=(200, 100000))
        .title("GDP per capita (PPP dollars)"),
        y=alt.Y("lifeExp", axis=alt.Axis(values=[20, 30, 40, 50, 60, 70, 80]))
        .title("Life expectancy (years)")
        .scale(domain=(20, 85)),
        color=alt.Color("continent").scale(range=okabe_ito),
    )
    .transform_filter(alt.datum.year == 2007)
)

chart = chart.configure_axis(labelFontSize=20, titleFontSize=20)

chart = chart.properties(width=600, height=500)

chart = chart.configure_title(
    fontSize=20,
    subtitleFontSize=20,
    anchor="start",
    orient="bottom",
    offset=20,
    subtitleColor="gray",
)

chart = chart.configure_legend(
    titleFontSize=20,
    labelFontSize=20,
    padding=10,
)

chart.interactive()
```
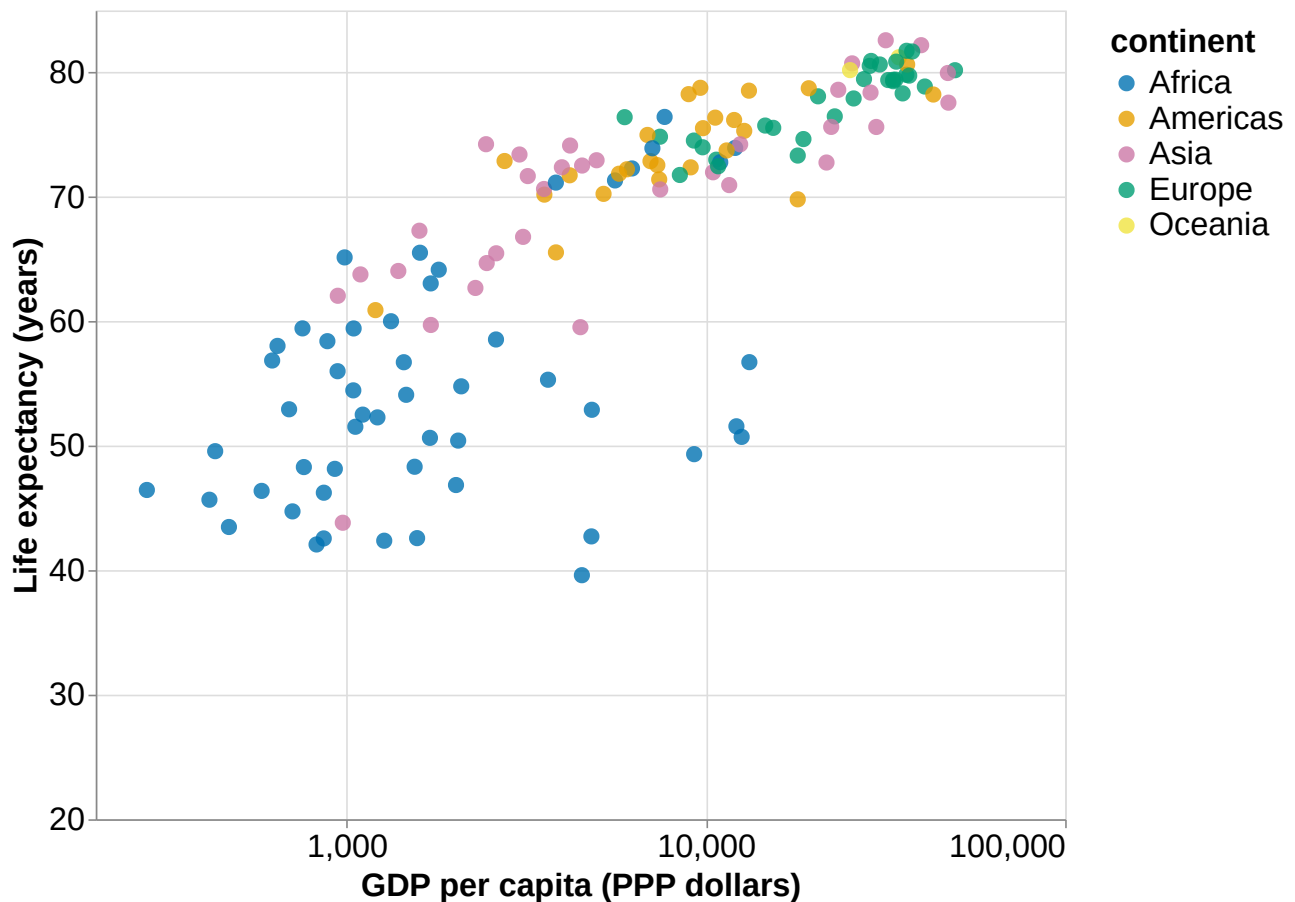
**Life expectancy as function of the gross domestic product**
Gross domestic product (GDP) per capita measures the value of everything
produced in a country during a year, divided by the number of people.
The unit is in purchasing power parities (PPP dollars), fixed to 2017 prices.
Data is adjusted for inflation and differences in the cost of living between countries.

## Interactive charts

With not too many changes we can make the chart interactive and add a slider for the year
(please try this in this notebook):

```
year_slider = alt.binding_range(min=1952, max=2007, step=5, name="Year")
slider_selection = alt.selection_point(bind=year_slider, fields=["year"], value=2007)

chart = (
    alt.Chart(
        data,
        title=alt.Title(
            "How life expectancy and gross domestic product evolved over time",
            subtitle=[
                "Gross domestic product (GDP) per capita measures the value of
everything",
                "produced in a country during a year, divided by the number of
people.",
                "The unit is in purchasing power parities (PPP dollars), fixed to 2017
prices.",
                "Data is adjusted for inflation and differences in the cost of living
between countries.",
            ],
        ),
    )
    .mark_circle(opacity=0.8, size=100.0)
    .encode(
        x=alt.X("gdpPercap", axis=alt.Axis(values=[100, 1000, 10000, 100000]))
        .scale(type="log", domain=(200, 100000))
        .title("GDP per capita (PPP dollars)"),
        y=alt.Y("lifeExp", axis=alt.Axis(values=[20, 30, 40, 50, 60, 70, 80]))
        .title("Life expectancy (years)")
        .scale(domain=(20, 85)),
        color=alt.Color("continent").scale(range=okabe_ito),
    )
    .add_params(slider_selection)
    .transform_filter(slider_selection)
)

chart = chart.configure_axis(labelFontSize=20, titleFontSize=20)

chart = chart.properties(width=600, height=500)

chart = chart.configure_title(
    fontSize=20,
    subtitleFontSize=20,
    anchor="start",
    orient="bottom",
    offset=20,
    subtitleColor="gray",
)

chart = chart.configure_legend(
    titleFontSize=20,
    labelFontSize=20,
    padding=10,
)

chart.interactive()
```

## Adding more annotation

With few more lines we can add extra annotation that can help to highlight certain aspects
of the plot and to tell a story:

```python
year_slider = alt.binding_range(min=1952, max=2007, step=5, name="Year")
slider_selection = alt.selection_point(bind=year_slider, fields=["year"], value=2007)

chart = (
    alt.Chart(
        data,
        title=alt.Title(
            "How life expectancy and gross domestic product evolved over time",
            subtitle=[
                "Gross domestic product (GDP) per capita measures the value of
everything",
                "produced in a country during a year, divided by the number of
people.",
                "The unit is in purchasing power parities (PPP dollars), fixed to 2017
prices.",
                "Data is adjusted for inflation and differences in the cost of living
between countries.",
            ],
        ),
    )
    .mark_circle(opacity=0.8, size=100.0)
    .encode(
        x=alt.X("gdpPercap", axis=alt.Axis(values=[100, 1000, 10000, 100000]))
        .scale(type="log", domain=(200, 100000))
        .title("GDP per capita (PPP dollars)"),
        y=alt.Y("lifeExp", axis=alt.Axis(values=[20, 30, 40, 50, 60, 70, 80]))
        .title("Life expectancy (years)")
        .scale(domain=(20, 85)),
        color=alt.Color("continent").scale(range=okabe_ito),
    )
    .add_params(slider_selection)
    .transform_filter(slider_selection)
)

annotation = (
    alt.Chart(data)
    .encode(
        x="gdpPercap",
        y="lifeExp",
        text="country",
        color=alt.value("black"),
    )
    .transform_filter((slider_selection) & (alt.datum.country == "Norway"))
)

chart = (
    chart
    + annotation.mark_point(size=100.0)
    + annotation.mark_text(size=15, xOffset=10, align="left", baseline="middle")
)

chart = chart.configure_axis(labelFontSize=20, titleFontSize=20)

chart = chart.properties(width=600, height=500)

chart = chart.configure_title(
    fontSize=20,
    subtitleFontSize=20,
    anchor="start",
    orient="bottom",
    offset=20,
    subtitleColor="gray",
)
```
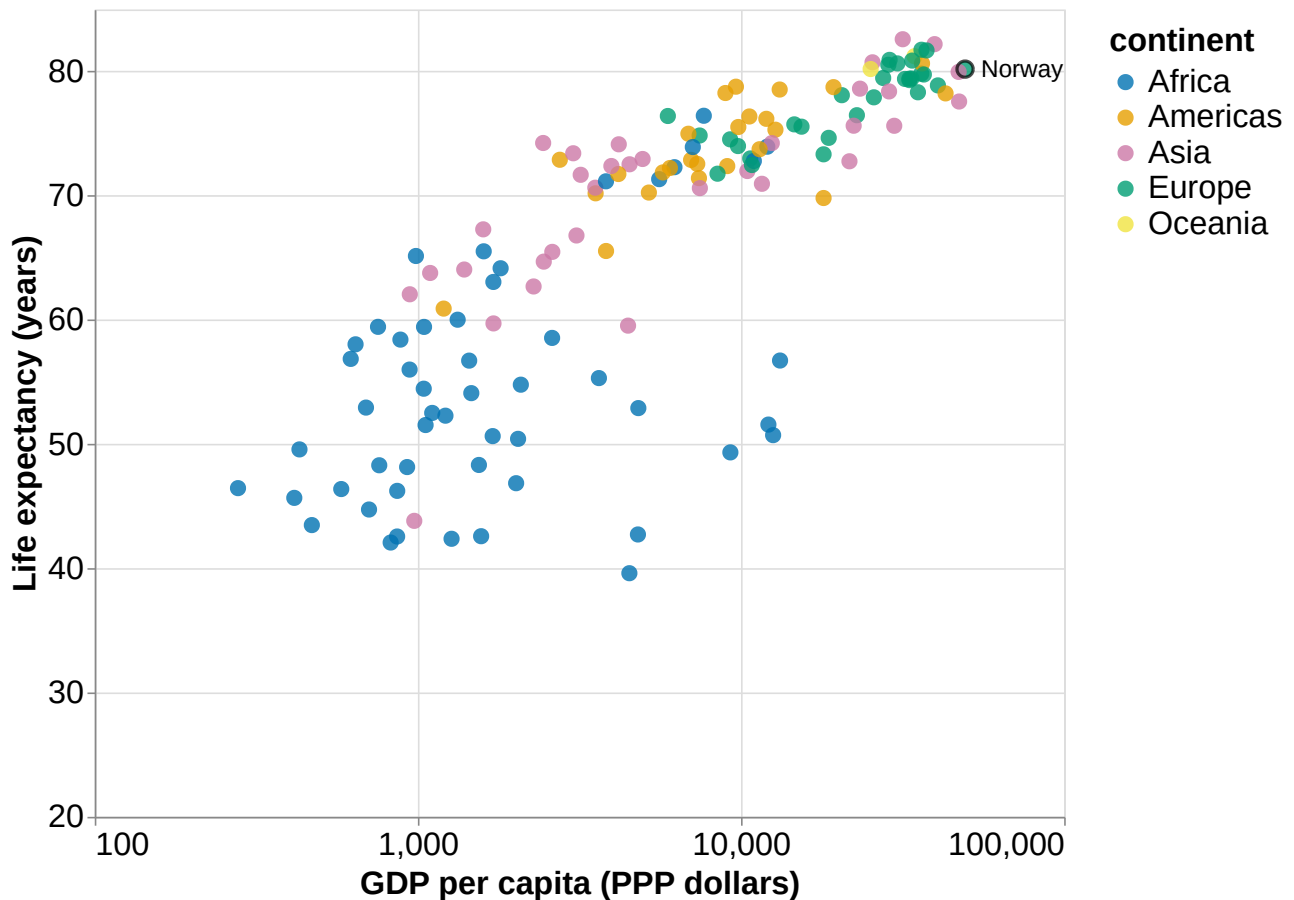
```
chart = chart.configure_legend(
    titleFontSize=20,
    labelFontSize=20,
    padding=10,
)

chart.interactive()
```



**How life expectancy and gross domestic product evolved over time**
Gross domestic product (GDP) per capita measures the value of everything
produced in a country during a year, divided by the number of people.
The unit is in purchasing power parities (PPP dollars), fixed to 2017 prices.
Data is adjusted for inflation and differences in the cost of living between countries.

## Saving the chart as web page

You can save the chart as a web site and try to open it in a separate browser tab and put it on
your home page or research group website:

```
chart.save("chart.html")
```

## Learning how to adapt existing gallery examples

In this exercise we can try to adapt existing scripts to either **tweak how the plot looks** or to
**modify the input data.** This is very close to real life: there are so many options and
possibilities and it is almost impossible to remember everything so this strategy is useful to
practice:

- Select an example that is close to what you have in mind
- Being able to adapt it to your needs
- Being able to search for help

**This is a great exercise which is very close to real life.**

- Browse the Vega-Altair example gallery.
- Select one example that is close to your current/recent visualization project or simply interests you.
- First try to reproduce this example, as-is, in the Jupyter Notebook.
- **If you get the error "ModuleNotFoundError: No module named 'vega_datasets'", then try one of these examples:** (they do not need the "vega_datasets" module)
  - Slider cutoff (**below you can find a walk-through for this example**)
  - Multi-Line tooltip
  - Heatmap
  - Layered histogram
- Then try to print out the data that is used in this example just before the call of the plotting function to learn about its structure. Consider writing the data to file before changing it.
- Then try to modify the data a bit.
- If you have time, try to feed it different, simplified data. **This will be key for adapting the examples to your projects.**

✔ Example walk-through for the slider cutoff example

In this walk-through I imagine browsing: https://altair-viz.github.io/gallery/index.html

Then this example caught my eye: https://altair-viz.github.io/gallery/slider_cutoff.html

I then copy-paste the example code into a notebook and try to run it and I get the same result.

Next, there is a lot of code that I don't (need to) understand yet but my eyes are trying to find `alt.Chart` which tells me that the data must be the "df" in `alt.Chart(df)`:

```python
import altair as alt
import pandas as pd
import numpy as np

rand = np.random.RandomState(42)

df = pd.DataFrame({
    'xval': range(100),
    'yval': rand.randn(100).cumsum()
})

slider = alt.binding_range(min=0, max=100, step=1)
cutoff = alt.param(bind=slider, value=50)

alt.Chart(df).mark_point().encode(
    x='xval',
    y='yval',
    color=alt.condition(
        alt.datum.xval < cutoff,
        alt.value('red'), alt.value('blue')
    )
).add_params(
    cutoff
)
```

My next step will be to print out the data `df` just before the call to `alt.Chart`:

```python
import altair as alt
import pandas as pd
import numpy as np

rand = np.random.RandomState(42)

df = pd.DataFrame({
    'xval': range(100),
    'yval': rand.randn(100).cumsum()
})

slider = alt.binding_range(min=0, max=100, step=1)
cutoff = alt.param(bind=slider, value=50)

print(df)

alt.Chart(df).mark_point().encode(
    x='xval',
    y='yval',
    color=alt.condition(
        alt.datum.xval < cutoff,
        alt.value('red'), alt.value('blue')
    )
).add_params(
    cutoff
    )
```

The print reveals that `df` is a dataframe which contains x and y values:

```
      xval       yval
0        0    0.496714
1        1    0.358450
2        2    1.006138
3        3    2.529168
4        4    2.295015
..     ...         ...
95      95  -10.712354
96      96  -10.416233
97      97  -10.155178
98      98  -10.150065
99      99  -10.384652

[100 rows x 2 columns]
```

The next thing that often helps me is to save the data to a comma-separated values (CSV) file:

```python
import pandas as pd

df.to_csv("data.csv", index=False)
```

I then open the file in an editor and see that it contains 100 rows:

```
xval,yval
0,0.4967141530112327
1,0.358449851840048
2,1.0061383899407406
3,2.5291682463487657
4,2.2950148716254297
5,2.060877914676249
6,3.6400907301836405
7,4.407525459336549
8,3.938051073401597
9,4.4806111169875615
...
```

Saving the data to file often helps me to see the structure of the data and now I am in a position to replace this with my own data. I create a file called "mydata.csv" and there I use the maximum temperatures for months 1-10 from the Tromso monthly data which we used further up:

```
xval,yval
01,7.7
02,6.6
03,4.5
04,9.8
05,17.7
06,25.4
07,26.7
08,25.1
09,19.3
10,9.8
```

In the notebook I then verify that the reading of the data works:

```
mydata = pd.read_csv("mydata.csv")

mydata
```

Now I can replace the example with my own data (note how I now can comment out some code that I don't need any longer):

```python
import altair as alt
import pandas as pd
# import numpy as np

# rand = np.random.RandomState(42)

# df = pd.DataFrame({
#     'xval': range(100),
#     'yval': rand.randn(100).cumsum()
# })

slider = alt.binding_range(min=0, max=100, step=1)
cutoff = alt.param(bind=slider, value=50)

# print(df)
df = pd.read_csv("mydata.csv")

alt.Chart(df).mark_point().encode(
    x='xval',
    y='yval',
    color=alt.condition(
        alt.datum.xval < cutoff,
        alt.value('red'), alt.value('blue')
    )
).add_params(
    cutoff
    )
```

Seems to work! I then make few more adjustments (I want the slider to work on the y-axis and have a more reasonable default):

```python
import altair as alt
import pandas as pd

slider = alt.binding_range(min=0, max=30, step=1)
cutoff = alt.param(bind=slider, value=15)

df = pd.read_csv("mydata.csv")

alt.Chart(df).mark_point().encode(
    x='xval',
    y='yval',
    color=alt.condition(
        alt.datum.yval < cutoff,
        alt.value('red'), alt.value('blue')
    )
).add_params(
    cutoff
    )
```

My next steps would then be to change axis titles, display the month names, add a legend, and refine from here.

# Sharing plots and notebooks

**❶ Objectives**

- Know about good practices for notebooks to make them reusable
- Have a recipe to share a dynamic and reproducible visualization pipeline

[this lesson is adapted after https://coderefinery.github.io/jupyter/sharing/]

## Document dependencies

If you import libraries into your notebook, note down their versions.

In Python, it is customary to do this either in a `requirements.txt` file (example):

```
jupyterlab
altair == 5.5.0
vega_datasets
pandas == 2.2.3
numpy == 2.1.2
```

... or in an `environment.yml` file (example):

```yaml
name: data-viz
channels:
  - conda-forge
dependencies:
  - python <= 3.12
  - jupyterlab
  - altair-all = 5.5.0
  - vega_datasets
  - pandas = 2.2.3
  - numpy = 2.1.2
```

By the way, this is almost the same `environment.yml` file that we used to install the local software environment in the Software install instructions (the latter did not pin versions).

Place either `requirements.txt` or `environment.yml` in the same folder as the notebook(s).

This is not only useful for people who will try to rerun this in future, it is also understood by some tools (e.g. Binder) which we will see later.

## Different ways to share a Vega-Altair plot

- Save it in SVG format (vector graphics, "maximum resolution")
- Save it in PNG format (raster graphics)
- Share it as notebook (more about it below)
- Save it a web page with `chart.save("chart.html")` and share the HTML file
- You can also get a shareable URL to a chart (example)
- With **sensitive data**, you need to be careful with sharing (see next section)

## Vega-Altair and notebooks containing sensitive data

If you plot **sensitive data** in a notebook with Vega-Altair, you need to be careful.

The author of Vega-Altair provided a good summary in this GitHub comment:

> "Standard Altair rendering requires the entire dataset to be accessible to the viewer's browser: this is a fundamental design decision in Vega/Vega-Lite, in which a chart is equivalent to a dataset plus a specification of how to render it. In general, you should assume that the entire contents of any dataframe you pass to the alt.Chart() object will be saved in the notebook and be inspectable by the viewer."

> "One way to get around this would be to render the chart server-side, export a PNG, and display this png instead of the live chart. Incidentally, in the Jupyter notebook you can do this by running:"

```python
alt.renderers.enable('png')
```

"This sets up Altair such that charts will be rendered to PNG within the kernel, and only that PNG rendering will be embedded in the notebook. Note this requires some extra dependencies, described here."

"But even here, I wouldn't call your data "private" (for example, if you save a scatter plot to PNG, a user can straightforwardly read the data values off the chart!) So this makes me think you're actually doing some sort of aggregation of your data before plotting (e.g. showing a histogram). If this is the case, I would suggest doing those aggregations outside of Altair using e.g. pandas, and then passing the aggregated dataset to the chart. Then you get the normal interactive display of the Altair chart, and your data is just as private as it would have been in the equivalent static rendering – the user can only see the aggregated values you supplied to the chart."

## Different ways to share a notebook

We need to learn how to share notebooks. At the minimum we need to share them with our future selves (backup and reproducibility).

- You can enter a URL, GitHub repo or username, or GIST ID in nbviewer and view a rendered Jupyter notebook
- Read the Docs can render Jupyter Notebooks via the nbsphinx package
- Binder creates live notebooks based on a GitHub repository
- EGI Notebooks (see also https://egi-notebooks.readthedocs.io)
- JupyterLab supports sharing and collaborative editing of notebooks via Google Drive. Recently it also added support for Shared editing with collaborative notebook model.
- JupyterLite creates a Jupyterlab environment in the browser and can be hosted as a GitHub page.
- Notedown, Jupinx and DocOnce can take Markdown or Sphinx files and generate Jupyter Notebooks
- Voilà allows you to convert a Jupyter Notebook into an interactive dashboard
- The `jupyter nbconvert` tool can convert a ( `.ipynb` ) notebook file to:
  - python code ( `.py` file)
  - an HTML file
  - a LaTeX file
  - a PDF file
  - a slide-show in the browser

The following platforms can be used free of charge but have **paid subscriptions** for faster access to cloud resources:

- CoCalc (formerly SageMathCloud) allows collaborative editing of notebooks in the cloud
- Google Colab lets you work on notebooks in the cloud, and you can read and write to notebook files on Drive
- Microsoft Azure Notebooks also offers free notebooks in the cloud
- Deepnote allows real-time collaboration

## Sharing dynamic notebooks using Binder

**✍️ Exercise/demo: Making your notebooks reproducible by anyone (15 min)**

Instructor demonstrates this:

- Instructor creates a GitHub repository.
- Uploads a notebook file that we created in earlier episodes.
- Then we look at the statically rendered version of the notebook on GitHub and also nbviewer.
- Add a file `requirements.txt` which contains:

```
altair == 5.5.0
vega_datasets
pandas == 2.2.3
numpy == 2.1.2
```

- Visit https://mybinder.org:



- Check that your notebook repository now has a "launch binder" badge in your `README.md` file on GitHub.
- Try clicking the button and see how your repository is launched on Binder (can take a minute or two). Your notebooks can now be explored and executed in the cloud.
- Enjoy being fully reproducible!

Also please see how we share the notebooks from this lesson in the Episode overview.

## How to get a digital object identifier (DOI)

- [Zenodo](#) is a great service to get a [DOI](#) for a notebook (but **first practice** with the [Zenodo sandbox](#)).
- [Binder](#) can also run notebooks from Zenodo.
- In the supporting information of your paper you can refer to its DOI.

# Software install instructions

[this page is adapted from [https://aaltoscicomp.github.io/python-for-scicomp/installation/](https://aaltoscicomp.github.io/python-for-scicomp/installation/)]

## Choosing an installation method

For this course we will install an isolated environment with following dependencies:

```
name: data-viz
channels:
  - conda-forge
dependencies:
  - python <= 3.12
  - jupyterlab
  - altair-all
  - vega_datasets
  - pandas
  - numpy
```

If you are used to installing packages in Python and know what to do with the above `environment.yml` file, please follow your own preferred installation method.

**If you are new to Python or unsure** how to create isolated environments in Python from files like the `environment.yml` above, please follow the instructions below.

> 💬 **There are many choices and we try to suggest a good compromise**
>
> There are very many ways to install Python and packages with pros and cons and in addition there are several operating systems with their own quirks. This can be a huge challenge for beginners to navigate. It can also difficult for instructors to give recommendations for something which will work everywhere and which everybody will like.
>
> Below we will recommend **Miniforge** since it is free, open source, general, available on all operating systems, and provides a good basis for reproducible environments. However, it does not provide a graphical user interface during installation. This means that every time we want to start a JupyterLab session, we will have to go through the command line.

> ❶ **Python, conda, anaconda, miniforge, etc?**
>
> Unfortunately there are many options and a lot of jargon. Here is a crash course:

- **Python** is a programming language very commonly used in science, it's the topic of this course.
- **Conda** is a package manager: it allows distributing and installing packages, and is designed for complex scientific code.
- **Mamba** is a re-implementation of Conda to be much faster with resolving dependencies and installing things.
- An **Environment** is a self-contained collections of packages which can be installed separately from others. They are used so each project can install what it needs without affecting others.
- **Anaconda** is a commercial distribution of Python+Conda+many packages that all work together. It used to be freely usable for research, but since ~2023-2024 it's more limited. Thus, we don't recommend it (even though it has a nice graphical user interface).
- **conda-forge** is another channel of distributing packages that is maintained by the community, and thus can be used by anyone. (Anaconda's parent company also hosts conda-forge packages)
- **miniforge** is a distribution of conda pre-configured for conda-forge. It operates via the command line.
- **miniconda** is a distribution of conda pre-configured to use the Anaconda channels.

## Installing Python via Miniforge

Follow the instructions on the miniforge web page. This installs the base, and from here other packages can be installed.

## Installing and activating the software environment

First we will start Python in a way that activates conda/mamba. Then we will install the software environment from this environment.yml file.

An **environment** is a self-contained set of extra libraries - different projects can use different environments to not interfere with each other. This environment will have all of the software needed for this particular course.

We will call the environment `data-viz`.

| Windows | Linux / MacOS |
|---------|---------------|

Use the "Miniforge Prompt" to start Miniforge. This will set up everything so that `conda` and `mamba` are available. Then type (without the `$`):

```
$ mamba env create -n data-viz -f
https://raw.githubusercontent.com/coderefinery/data-visualization-
python/main/software/environment.yml
```

## Starting JupyterLab

Every time we want to start a JupyterLab session, we will have to go through the command line and first activate the `data-viz` environment.

| **Windows** | Linux / MacOS |
|---|---|

Start the Miniforge Prompt. Then type (without the `$`):

```
$ conda activate data-viz
$ jupyter-lab
```

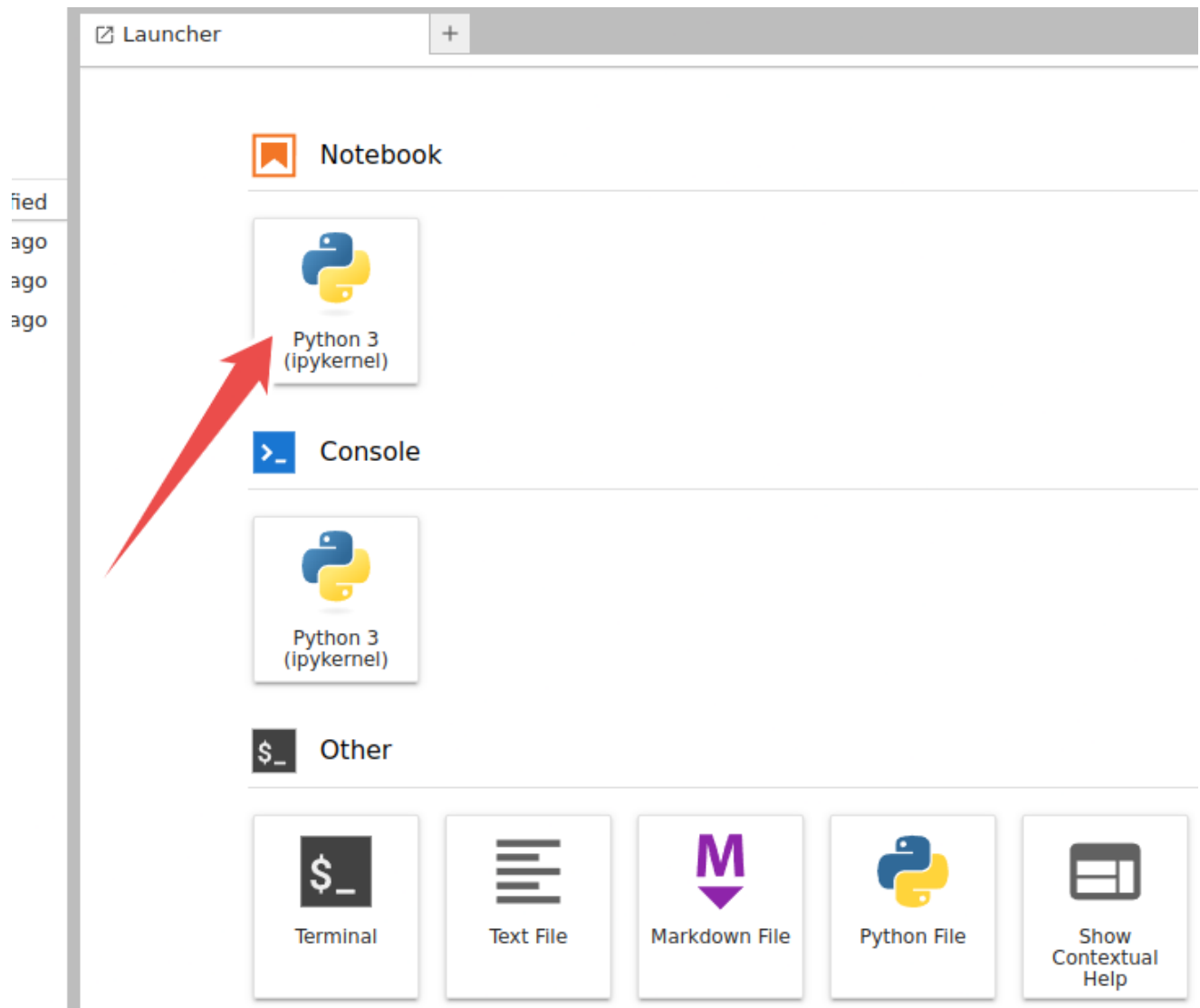## Removing the software environment

| **Windows** | Linux / MacOS |
|---|---|

In the Miniforge Prompt, type (without the `$`):

```
$ conda env list
$ conda env remove --name data-viz
$ conda env list
```
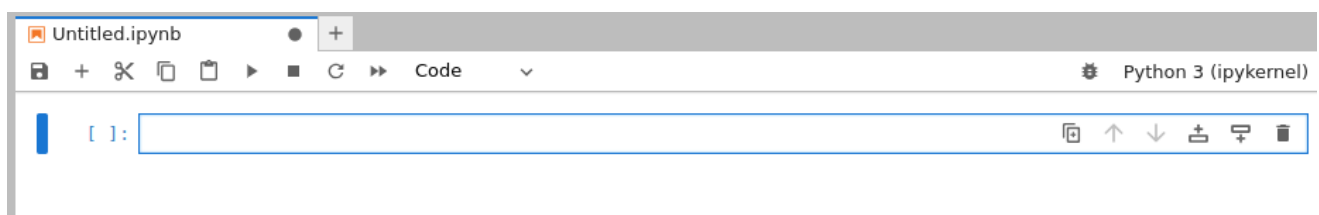
## How to verify your installation

**Start JupyterLab** (as described above). It will hopefully open up your browser and look like this:

*JupyterLab opened in the browser. Click on the Python 3 tile.*

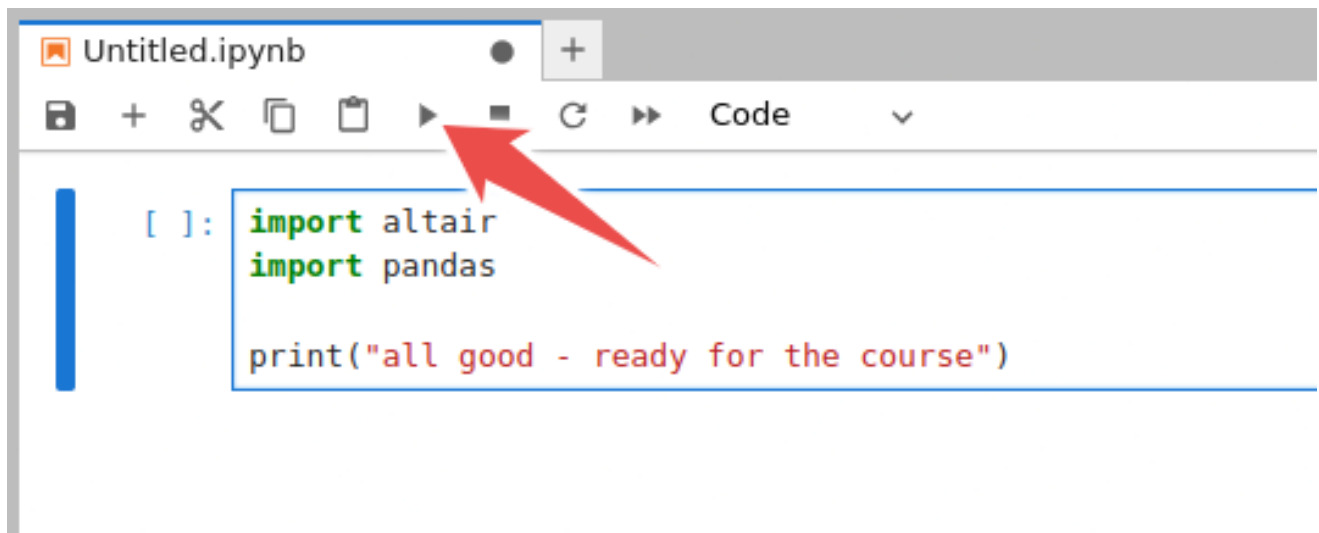Once you clicked the Python 3 tile it should look like this:



*Python 3 notebook started.*

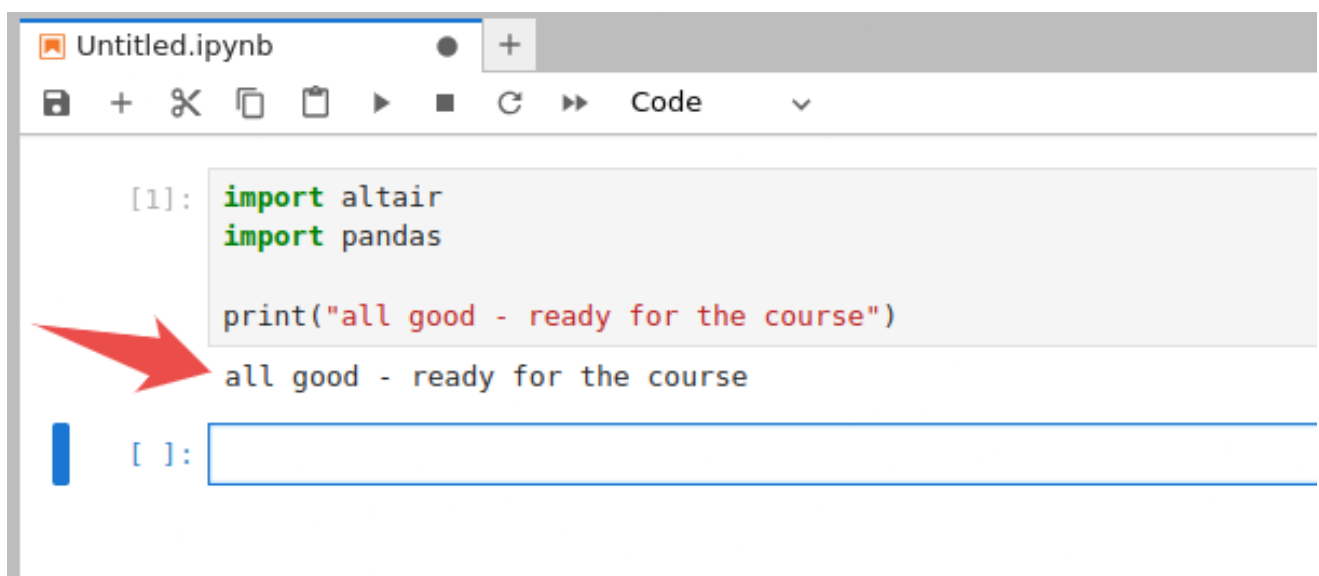Into that blue "cell" please type the following:

```python
import altair
import pandas

print("all good - ready for the course")
```

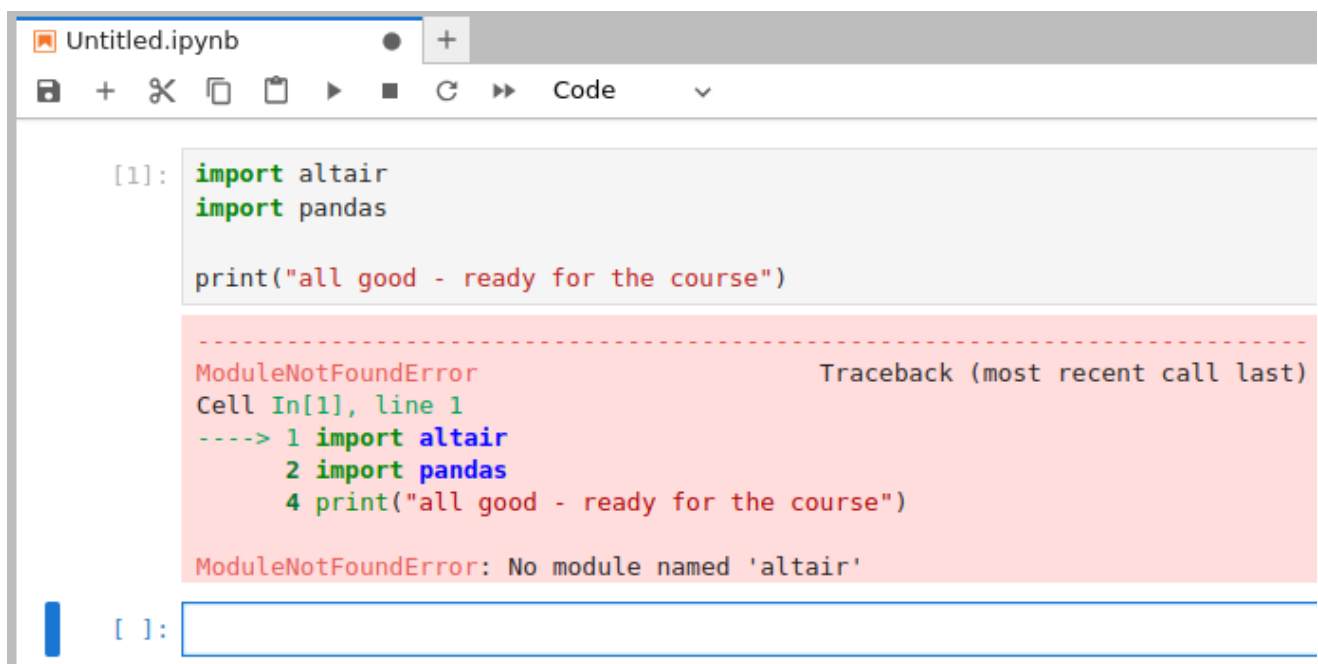*Please copy these lines and click on the "play"/"run" icon.*

This is how it should look:



*Screenshot after successful import.*

If this worked, you are all set and can close JupyterLab (no need to save these changes).

This is how it **should not** look:

*Error: required packages could not be found.*

## Going through the course using Google Colab

It is possible to go through the workshop examples using Google Colab.

However you may need the following cell on top of your notebook:

```python
import altair as alt

# this is here for google colab to update altair
if not alt.__version__.startswith("5"):
    %pip install altair==5.5.0
```

And then you may need to click on "Runtime" -> "Restart session and run all".

## Reading data in custom format

From earlier episodes we know that storing data in standard formats can be convenient and we know about the tidy data format but **sometimes we don't have control over this**: we may get a text file from another program or an instrument and now we are expected to deal with it.

As an example, we got two text files from two different instruments. We wish to extract all frequencies and all intensities into two lists.

This is one, `example1.txt` :

```
result from instrument: R2-D2
-----------------------------------------

  measurement    frequency    intensity

*      1             0.01         0.01
*      2             0.02         0.02
*      3             0.03         0.01
*      4             0.04         0.10
*      5             0.05         0.20
*      6             0.06         0.12
*      7             0.07         0.07
*      8             0.08         0.02
*      9             0.09         0.01
*     10             0.10         0.01


==========================================

timestamp: Sat Mar 27 03:30:34 PM CET 2021
```

And here is the other one, `example2.txt` :

```
result from instrument: C-3PO
==============================

  numbers we want: 10
       0.01    0.01
       0.02    0.02
       0.03    0.01
       0.04    0.10
       0.05    0.20
       0.06    0.12
       0.07    0.07
       0.08    0.02
       0.09    0.01
       0.10    0.01

  unrelated numbers:
       1.23    4.56
       7.89    0.12
```

At the end we want the code to produce these two lists:

```
frequencies = [0.01, 0.02, 0.03, 0.04, 0.05, 0.06, 0.07, 0.08, 0.09, 0.1]
intensities = [0.01, 0.02, 0.01, 0.1, 0.2, 0.12, 0.07, 0.02, 0.01, 0.01]
```

💬 **Discussion**

Before we discuss possible solutions, it can be a good exercise to **describe in words** what the code should do to extract the numbers.

Below you find possible solutions for these two examples. These are designed to be somehow understandable (which does not mean that the code is trivial) and are not as short and as elegant as they could be. One way to make them more elegant would be to use regular expressions.

One solution for `example1.txt` :

```python
def read_data(file_name):
    # we start with empty lists
    frequencies = []
    intensities = []

    # open the file with read permissions
    with open(file_name, "r") as f:
        # iterate over all lines in the file object "f"
        for line in f:
            # we are only interested in lines that start with "*"
            if line.startswith("*"):
                # we split the line on "whitespace"
                # the "_" means we are not interested in the first two entries
                _, _, frequency, intensity = line.split()

                # add the new numbers to the lists
                # we convert from string to float
                frequencies.append(float(frequency))
                intensities.append(float(intensity))

    # function returns the result
    return frequencies, intensities


# here we are outside the function
frequencies, intensities = read_data("example1.txt")

print(frequencies)
print(intensities)
```

One solution for `example2.txt` .

```python
def read_data(file_name):
    # we start with empty lists
    frequencies = []
    intensities = []

    # open the file with read permissions
    with open(file_name, "r") as f:
        # iterate over all lines in the file object "f"
        for line in f:
            # we are only interested in lines that start with "*"
            if "numbers we want" in line:
                # we split the line on "whitespace"
                words = line.split()
                # we are only interested in the last element
                # -1 means last element in that list
                last_element = words[-1]
                # convert from string to int
                num_measurements = int(last_element)

                # now we know the next 10 lines are the interesting ones
                # range(10) produces a list with 10 elements: [0, 1, 2, 3, 4, 5, 6, 7,
8, 9]
                # we use it to iterate exactly 10 times
                for _ in range(num_measurements):
                    # this advances f by one and we get one line at a time
                    next_line = next(f)
                    words = next_line.split()
                    frequency = float(words[0])
                    intensity = float(words[1])
                    # add the new numbers to the lists
                    frequencies.append(frequency)
                    intensities.append(intensity)

    # function returns the result
    return frequencies, intensities


# here we are outside the function
frequencies, intensities = read_data("example2.txt")

print(frequencies)
print(intensities)
```

# Credit

When preparing this lesson, we have reused these resources:

- https://aaltoscicomp.github.io/python-for-scicomp/
- https://datacarpentry.org/python-ecology-lesson/
- https://swcarpentry.github.io/python-novice-inflammation/
- https://coderefinery.github.io/jupyter/