

Digital Skills

Christos Dimitrakakis

December 2, 2022

Contents

1	Introduction	5
1.1	Learning goals:	5
1.2	Administration	6
1.3	Assessment	6
1.4	Data sources	8
1.4.1	Synthetic data	8
1.4.2	UCI machine learning repository	8
1.4.3	Wikipedia and newspaper articles	8
1.4.4	Economics data	8
1.4.5	PET Statistics Hackathon	8
1.4.6	NASA Mars Challenge	8
2	Module 1: Visualisation as models and data summary (4 weeks)	9
2.1	Histograms: model a distribution	9
2.1.1	Bar graph activity	9
2.1.2	Introduction to histograms THEORY:STATISTICS . .	11
2.1.3	General python help CODE:PYTHON	12
2.1.4	A simple program	13
2.1.5	Consoles, Scripts and Notebooks	13
2.1.6	Python variables CODE	14
2.1.7	Python lists CODE	16
2.1.8	Numpy arrays ADVANCED:CODE	17
2.1.9	Python control structures	17
2.1.10	Python functions CODE	18
2.1.11	Pandas and Histograms PLOT:CODE	19
2.1.12	Histograms vs Pie Charts	20

2.2	Time-Series: model the evolution of a system	20
2.2.1	Plotting lines	22
2.2.2	Race times	22
2.2.3	Example: The inclination of Mars	23
2.2.4	Example: Covid	23
2.2.5	Example: Stock market prices	23
2.3	Scatterplots: model a relationship	23
2.3.1	Relationships as functions	24
2.3.2	Relationships as joint distributions	25
2.3.3	Relationships as conditional distribution.	25
2.3.4	Example: Unemployment, GDP	25
3	Module 2: Experiment design (3 weeks)	26
3.1	Data collection and cleaning	26
3.2	Random sampling	26
3.2.1	Survey of political opinions	26
3.2.2	Uniform sampling	30
3.2.3	Expectation	31
3.2.4	Survey of heights	33
3.3	The data science pipeline	33
3.3.1	Salaries of men and women	34
3.3.2	A question of voting	36
4	Module 3: Inference (2 weeks)	37
4.1	Logic and bar graphs	37
4.1.1	Logic and sets	37
4.2	Conditional Probability and the Theorem of Bayes	38
4.2.1	Bayes theorem	38
4.2.2	The cards problem	38
4.3	Bayesian analysis example	39
4.3.1	The murder problem	39
4.3.2	DNA test	39
4.3.3	Everybody is a suspect	40
4.3.4	Explanation	40
4.3.5	Python example	40
4.4	More general problems	41
4.4.1	The \$\$\$-meteorologists problem	41
4.4.2	Types of hypothesis testing problems	41
4.4.3	Pitfalls	42
4.4.4	The covid test problem	42

4.5	Hypothesis testing	43
4.5.1	Estimation intervals	43
4.6	Simulation study	43
	ACTIVITY	
4.7	Introduction: Estimation, Prediction and Decision Making . .	44
4.7.1	Correlation example	44
4.7.2	The three main problems in statistics	44
4.7.3	Types of variables	45
4.7.4	From correlation to causation.	45
4.7.5	Correlation versus Causality	46
4.8	Joint and Conditional Probabilities	46
4.8.1	Correlation: Setting	46
4.8.2	Conditional probabilities	46
4.8.3	Python example: independent RVs	47
4.8.4	Dependent, Discrete x, y	47
4.8.5	Python example: multiple draws	48
4.8.6	Empirical joint probability of x, y	49
4.8.7	Empirical joint probability of x, y	49
4.8.8	Continuous x, y	49
4.8.9	Continuous x , Discrete y	50
4.8.10	Covariance matrix	51
4.8.11	Correlation versus dependence	51
4.9	Models of Causation	52
4.9.1	Causal inference vs the actual cause	52
4.9.2	Confounding variables	52
4.9.3	Instrumental variables	53
4.9.4	Avoiding confounding	54
5	Module 4: Advanced visualisation (2 weeks)	54
5.1	Geographical data	54
5.1.1	Colour maps	54
5.1.2	Contour maps	54
5.2	Text data	55
6	Module 5: Data analysis in practice (2 weeks)	55
6.1	Visualising fMRI data	55
6.2	Visualising GWAS data	55
7	Module 6: Project work and presentations (2 weeks)	55

8	Assignments	55
8.1	Table To Picture	55
8.1.1	Instructions	55
8.1.2	Example	56
8.2	Plot deconstruction	56
8.2.1	Instructions	56
8.2.2	Example	57
8.3	Newspaper article analysis	57
8.4	Open project	57
8.4.1	Project proposal	57
8.4.2	Project Highlight	57
8.4.3	Project presentation	57
8.4.4	Project report	57
9	Background	57
9.1	Notation	57
9.1.1	Sets	58
9.1.2	Analysis	58
9.1.3	Probability	58
9.2	Probability background	59
9.2.1	Randomness CODE:ACTIVITY	59
9.2.2	Uncertainty ACTIVITY	61
9.2.3	Probability space THEORY:PROBABILITY	61
9.2.4	Probability measures THEORY:PROBABILITY	62
9.2.5	Marginalisation	62
9.2.6	Mutually exclusive events DEFINITION:PROBABILITY	63
9.2.7	Independence DEFINITION:PROBABILITY	63
9.2.8	Conditional probability DEFINITION:PROBABILITY	63
9.2.9	Bayes theorem THEORY:PROBABILITY	63
9.2.10	Conditional independence DEFINITION:PROBABILITY	64
9.2.11	Example Distributions	64
9.2.12	Random variables THEORY:PROBABILITY	66
9.2.13	Expectation DEFINITION:PROBABILITY	67
9.2.14	Variance DEFINITION:PROBABILITY	68
9.2.15	Conditional expectation DEFINITION:PROBABILITY	68
10	Graphics types	69
11	Schedule and links to other courses	70

12 Glossary	72
13 References	73
13.1 Some workshops	73
13.2 Books	73

1 Introduction

This is a hands-on course that integrates introductory programming, statistics and data science. Through out this course, we will formulate scientific hypotheses, design experiments, and collect and analyse data visually and through formal models. You are expected to supplement this course with homework, self-study and other courses in descriptive statistics and python programming, but no prior knowledge is assumed. Formal concepts will be introduced through class activities and examples.

The course will focus neither on statistical theory nor on programming. We will only lightly build those skills within the course. Our main aim will be to build scientific and statistical intuition through practical work. However, to achieve this you need to be able to independently pick up theoretical and programming skills. For that reason, it is necessary for you to do outside reading either (ideally) by taking statistics and programming courses in parallel, or through self-study.

1.1 Learning goals:

Graphical comprehension

1. Recognise structural elements in a statistical graph (e.g. axis, symbols, labels) and evaluate the effectiveness (for perception and judgment) and appropriateness (for the type of data) of structural element.
2. Translate relationships reflected in a graph to the data represented.
3. Recognise when one graph is more useful than another and organise/reorganise data to make an alternative representation.
4. Use context to make sense of what is presented in a graph and avoid reading too much into any relationships observed.
5. Express creative thinking via the production of an innovative graphical presentation.

Scientific process

1. Understanding the randomness, variability and uncertainty inherent in a problem.
2. Developing clear statements of the problem/scientific research question; understanding the purpose of the answer.
3. Be able to perform a basic experiment design.
4. Identify sources of bias in data collection and analysis.
5. Ensuring acquisition of high-quality data and not just a lot of numbers.
6. Understanding the process that produced the data, to provide proper context for analysis.
7. Allowing domain knowledge to guide both data collection and analysis.
8. Quantify uncertainty—and knowledge—visually.
9. Realise that all visualisations are model summaries.
10. Be able to write simple python programs for data science workflows.

1.2 Administration

1. Make sure you are registered on IS-Academia
2. Also register on Moodle: this is where the assignments will be
3. Clone this git repository

1.3 Assessment

The assessment is purely through in-class exercises, quizzes and homework assignments. There will be assignments spread over the semester, as well as a group project. The project will be performed in pairs.

For all assignment and the project, the following rubrik is used. Some of the assignments may not involve all parts.

Experiment design. The first stage any project, no matter how small, is the experiment design and analysis. This includes a plan for how to collect data, methodologies for analysing the data, and the development of a pipeline, preferrably in the form of a program, for collecting data and

analysing it. In addition, the experiment design must be reproducible: This can be ensured by running the data collection and analysis pipeline on simulated data, and seeing if the results are as expected.

Computation. Here you must instantiate the experiment design and analysis with concrete computations. For reproducibility, the computations you perform should be independent of the data you actually have. Correctness of the computations is the most important aspect, here. However, you should also take care to document why and how you are doing the computations.

Graphics. This addresses the creation of visualisations of your analysis. It is recommended to do this fully automatically, so that you can simply run your pipeline and get all the results you need. Be sure to quantify uncertainty.

Text. Here you should explain in text what the graphics mean. Point out any interesting things you can see in the visualisation and try to explain it. Do not be overconfident, but quantify uncertainty properly.

Synthesis. Here you should summarise the most important findings from your analysis. Be careful to not over-interpret your results. A lot of results can be imaginary and can be attributed to insufficient data, biased sampling, improper modelling or p -value hacking. Again, be sure to quantify uncertainty.

Skill	Needs Improvement	Basic Level	Advanced Level
Experiment design: Data collection and analysis pipeline	Inappropriate sampling, non-reproducible analysis	Data collection biased or analysis not reproducible.	Unbiased sampling and reproducible experiment design and analysis.
Computation: Perform computations	Computations contain errors and extraneous code	Computations are correct but contain extraneous/unnecessary code.	Computations are correct and properly justified and explained.
Graphics: Communicate findings graphically clearly, precisely and concisely	Inappropriate choice of plots; poorly labelled plots; plots missing.	Plots convey information correctly but lack context for interpretation.	Plots convey information correctly with adequate and appropriate information
Text: Communicate findings clearly, precisely and concisely	Explanation is illogical, incorrect or incoherent.	Explanation is partially correct but incomplete or unconvincing	Explanation is correct, complete and convincing.
Synthesis: Identify key features of the analysis and interpret results	Conclusions are missing, incorrect, or not made based on analysis	Conclusions reasonable, but partially correct or incomplete.	Relevant conclusions explicitly connected to analysis and context.

Pass: All parts must be addressed, the 'default' grade is 75%. 5% is added for every 'advanced' skill and removed for every 'needs improvement'

skill'. Thus the passing grades are 50-100%.

Fail: If not all parts are explicitly addressed, the assignment is failed.

1.4 Data sources

This course will consider the following data sources in order of importance.

1.4.1 Synthetic data

This data is obtained through simulation, and it is useful in order to test whether a particular pipeline is working as intended. In particular, it is a great way to test the performance of a method as you vary the data generation process so that different assumptions are satisfied. This allows you to verify robustness.

1.4.2 UCI machine learning repository

The UCI repository has a large collection of datasets in an easy to access format. These have already been used in many academic papers, and are a good starting point for you to look at real data. All the data is formatted in an easy-to-use some format, but some pre-processing may still be necessary.

1.4.3 Wikipedia and newspaper articles

Wikipedia has many interesting articles, from which you can extra tabular data, as well as more contextual information. It is possible to also discuss newspaper articles. Wikipedia and newspaper articles can be used in the context of some assignments.

1.4.4 Economics data

- FRED: Federal Reserve Economic Data
- OECD: Organisation for Economic Co-operation and Development

1.4.5 PET Statistics Hackathon

<https://petlab.officialstatistics.org/>

1.4.6 NASA Mars Challenge

<https://www.drivendata.org/competitions/97/nasa-mars-gcms/>

2 Module 1: Visualisation as models and data summary (4 weeks)

What is visualisation? It is a way to *summarise data*. It is also a way to view relationships between variables. Visualisation helps us to find patterns and understand the underlying laws behind how the data was generated. This is, in fact, the essence of modelling.

A model is *also* a way of summarising the essential features of the data. A visualisation differs from a model only in one sense: It is easy to interpret visually.

Every data visualisation implicitly assumes a model of the data generating process. This is true for even the simplest visualisations, like histograms. There is no escape from the fact that any visualisation makes a lot of assumptions. We must emphasize what those assumptions are. What happens if they are not true?

Every data visualisation, then, proceeds in three steps:

1. Data transformation
2. Model creation
3. Model visualisation

Parameters. Every model is defined by a number of parameters. This is what is displayed when we visualise data. You can think of the model as the underlying theory, and the visualisation as a way to explain the theory visually.

2.1 Histograms: model a distribution

Histograms are a simple tool for modelling distributions. In their simplest application, they are used to simply count the number of items in distinct bins of a dataset. While typically employed to represent the empirical distribution of one-dimensional variables, they can be generalised to multiple dimensions .

2.1.1 Bar graph activity

1. All students who are male raise their hand
2. All students who are female raise their hand.

3. We count, and draw a bar graph: the number of male, female and other students.
4. We count how many are in the BSc of DS of those
5. We also count how many are taking a programming course
6. We also count how many are taking a maths/stats course
7. What does the graph tell us about:
 - (a) Computer Science students
 - (b) Students at Neuchatel
 - (c) Residents of Neuchatel
 - (d) Other subsets of the population
8. Can we make a similar graph when measuring a continuous variable?

1. Definition of a bar graph

THEORY

Consider a set of k categories $C = \{1, \dots, k\}$. Every individual belongs in one category. This can be defined through a function f assigning categories to individuals.

In particular, imagine a dataset of individuals D . There exists a membership function

$$f : D \rightarrow C$$

so that $f(x)$ is the category to which individual $x \in D$ belongs.

A bar graph is a visual representation of the following count vector,

$$n_c(D) = \sum_{i \in C} \mathbb{I}\{f(i) = c\},$$

where \mathbb{I} is an indicator function:

$$\mathbb{I}\{A\} = \begin{cases} 1, & A \text{ is true} \\ 0, & A \text{ is false,} \end{cases}$$

so that the height of the c -th bar is proportional to n_c .

2.1.2 Introduction to histograms

THEORY:STATISTICS

Assume data is in \mathbb{R} . Then split the real line into intervals $[a_i, b_i)$. For a given dataset D , for each interval i , count the amount of data $n_i(D)$ in the interval. We can also normalise to obtain $p_i(D) = n_i(D) / \sum_j n_j(D)$

More generally, a (counting) histogram is defined as a collection of disjoint sets called **bins**

$$\{A_i | i = 1, \dots, k\}$$

with associated counts n_i , so that, given some data D ,

$$n_i(D) = \sum_{x \in D} \mathbb{I}[x \in A_i],$$

where n_i is the number of datapoints in A_i . Typically $A_i \subset \mathbb{R}$.

We can use the histogram as the model of a distribution. For that, we use the relative frequency of points in each bin: $p_i(D) = n_i(D) / \sum_j n_j(D)$. The selection of bins influences the model.

See also: <https://en.wikipedia.org/wiki/Histogram>

1. Histogram activity

ACTIVITY

- (a) Introduce the concept of a histogram on the board.
- (b) Split the students in two groups.
- (c) Have each group collect the height of every student.
- (d) How can we summarise the data of each group?
- (e) Now the students will individually draw a histogram from the data of their group.
- (f) Show two different histograms from two people in the same group. Why are they different? Discuss in pairs and then in class.
- (g) Now show a histogram from a person in another group. Why are the histograms in the two groups different? Discuss.
- (h) Collect the data of all students in the online excel file.
- (i) Now we shall plot a histogram of the students using the sheet. How does that differ?

[If there are not enough students, the exercise can be performed by adding random numbers using dice]

2. Measuring a discrete distribution

ACTIVITY

- (a) Toss a coin 10 times and record each one of the results, e.g. $\{0,1,1,0,0,0,0,1,1,1\}$.

- (b) Count the number of times it comes heads or tails.
- (c) We then summarise the result.

Let us denote the number of times you have heads by $N(x = k)$. It should be approximately true that $N(x = k) \approx 5$, however, this may not be true for everybody.

We can visualise this by plotting bars or lines, whose height is proportional to $N(x = k)$.

We typically assume that individual coin tosses are generated from The Bernoulli distribution. This means that the probability of heads or tails is fixed, and does not depend on the result of the previous tosses. Why might not that be the case?

If individual tosses are Bernoulli, then the distribution of the number of heads (or tails) is a binomial distribution.

We will now show how to achieve the same results programmatically.

2.1.3 General python help

CODE:PYTHON

To help yourself understand python, you can always take a look at the documentation in English or French. To start with, check out the Tutorial. Then use the library reference for advanced usage.

Sometimes it is quicker to just use the **help** command in the python console or the **?** command in the jupyter notebook.

Python can be used as a simple calculator

```
$ python3
Python 3.8.10 (default, Jun 22 2022, 20:18:18)
[GCC 9.4.0] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> 1 + 1
2
>>> 2 * 3 + 1
7
>>> 2 * (3 + 1)
8
>>> exit()
```

This interactive console is the most usual way of playing with python in the beginning, but it is not useful in general.

2.1.4 A simple program

Python programs are executed one statement at a time. Statements are separated by newlines. Anything appearing after a `#` symbol is not executed.

```
print("Hello world") # first statement, with a comment
print("Goodbye, world.") # second statement!
# print("This is not printed") - it is a comment, you see
```

Python programs can be generate text, write and read from files, access the internet, generate and display or save plots to disc, play and record music, record images from a camera....

Before we actually run the program, one of you can play the role of the python **interpreter**. The interpreter goes through each line of the program, interprets it and executes it.

When we execute a program in the console, we are assuming the role of the interpreter that steps through the program.

Run the above statements in the python interpreter

```
$ python3
Python 3.8.10 (default, Jun 22 2022, 20:18:18)
[GCC 9.4.0] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> print ("Hello world")
Hello world
>>> print ("Goodbye, world")
Goodbye, world
>>> exit()
$
```

The statements `print()` and `exit()` are called **functions**. Function names must always be used with parenthesis. The contents of the parentheses are called **arguments**. Changing the arguments to a function has a different effect.

We can now try and save the above commands in a file called "python-test.py" and execute it via

```
python3 pythontest.py
```

2.1.5 Consoles, Scripts and Notebooks

Console input is used when you want to have a purely interactive session to test something. In console mode, the interpreter executes each line as you enter it.

Script files are used to save your work and re-run it. They also allow you to build complex programs from multiple files, where each file has a different functionality. In script mode, python acts as though you were entering each line one-by-one. It reads each line and executes in turn.

Notebooks are something in between. They are script files with interactive output, and are very useful for rapid development and testing. They also save their state and output in between runs, so they help to document your code. We will use them a lot in class. There are two methods to use notebooks:

- Locally through e.g. jupyter-lab or jupyter-notebook
- Online through <https://colab.research.google.com/> <https://replit.com> or <https://noto.epfl.ch>

Most of the time, you want to be saving the code you write in a script file and executing it, instead of using the console. However, sometimes the interactivity of the console is helpful. This is when notebooks are used. After you are done developing something with the notebook, you can then extract what you need in a simple python script.

2.1.6 Python variables

CODE

The python interpreter has a **state**. This includes the contents of a memory where variables are stored, and the current location of the code pointer, that is which line will be executed next.

Variables are alphanumeric references to simple or complex objects. Possible variable names:

- X
- NumberOfApples
- salary
- scratch_{variable}
- y2

A variable can be assigned a value with the = operator.

```
x = 2 # this gives the numeric value of '2' to the variable
```

Variables must be defined before they can be used for the first time

```
>>> x
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'x' is not defined
>>> x = 2
>>> x # typing the name of a variable in the console gives you its value
2
```

Numerical Python variables are very simple entities. Let us go through this is easy program for a warm-up.

- `x=value`; assigns a value to a variable named `x`
- `print()`; displays something in the terminal

```
x = 1 # a variable
y = 2 # another variable
print(x+y) # print the value of this variable sum
x = y # assignment operation: now x has the same value as y
print(x) #what would this value be?
y = 3
print(x) #is x changed?
```

1. Possible confusion point: assignment operator and math equations

The assignment operator `=` is **not** a mathematical equation. For example, in mathematics I may write

$$x = y + 1 \tag{1}$$

$$x = 5 \tag{2}$$

This is a system of equations, which can be solved to obtain $5 = y + 1$ and so $y = 4$. This is **not** what the assignment operator means.

Consequently, the following program will fail with an error, as `y` is not defined

```
x = y + 1
x = 5
```

In the following program, the value of `y` will remain `-1` after the program ends

```

y = -1 # y = -1, x is not defined
x = y + 1 # now y = -1, x = 0
x = 5 # now y = -1, x = 5

```

In fact, writing the above as a system of equations makes no sense, as $x = y + 1$ cannot be true if $x = 5, y = -1$. Replacing, we obtain $5 = 0$, which is false.

So, while math-like notation is used in programming, its meaning is not really the same as in mathematics, most of the time.

2.1.7 Python lists

CODE

A slightly more complex object are python lists. A list can contain anything, and is so very flexible. It can contain numbers, strings, or arbitrary 'objects'.

Now check out the first part of the Histogram example. For that we need one line of setup so we can plot stuff.

```

import matplotlib.pyplot as plt # this is used for plotting
X=[0,1,0,0,0,0,1]; # list of coin tosses
plt.hist(X) # plot a histogram - this automatically splits everything into bins

```

In reality, the histogram function creates a so-called bar plot

```

import matplotlib.pyplot as plt # this is used for plotting
X=[0,1,0,0,0,0,1]; # list of coin tosses
plt.bar(["heads","tails"], [sum(X), len(X) - sum(X)]) # do a bar plot!

```

The following source creates a list of four numbers and returns one element. Things to unpack here:

- `x[i]` returns the **(i+1)-th** element: we start counting **from 0**
- the **return** statement sends a value back to the whatever started the python program: in this case this .org file.

```

x = [1, 2, 3, 4]
return x[3] # returns the last element of the list

```

The following program assigns arrays-values to variables. Now x, y are both lists.


```

x = [1, 2, 3, 4]
y = [-1, -2]
x = y # assignment operation: now x is just a different name for y
y[0] = 1 # modify the 0th element of y
return x # what would the value of x be?

```

Lists are different in one respect: when we assign one list name to another, this does not copy any data. Both names refer to the same data. Consequently, if we change the data, it changes for both variable names. The way to avoid that is to use the **copy()** function.

```

x = [1, 2, 3, 4]
y = [-1, -2]
x = y.copy() # copy operation: now x has a copy of y's data
y[0] = 1 # modify the 0th element of y
return x # what would the value of x be?

```

2.1.8 Numpy arrays

ADVANCED:CODE

Because lists are very flexible, they are a bit slow. A special type of object, an array, is used to handle lists of numbers. This is not defined in basic python, but only in one module called *numpy*. Even though basic Python has only a few commands, it has many modules that extend the language to perform complex tasks without having to code everything from scratch.

```

import numpy as np
x = np.array([1, 2, 3, 4])
y = np.array([-1, -2])
x = y # assignment operation:
y[0] = 1
return x

```

2.1.9 Python control structures

Sometimes we want to repeat some code. For example, we have two matrices of X, Y values and we wish to plot them:

```

import matplotlib.pyplot as plt
import numpy as np
X = np.random.uniform([10,128])
Y = X + np.random.uniform([10,128])
plt.plot(X[0], Y[0])

```

```
plt.plot(X[1], Y[1])

#.... etc - to avoid repetition we can use this:
for t in range(10): # this defines the variable t and it cycles it through the values 0 to 9
    # start of repeated block
    plt.plot(X[t], Y[t])
# end of repeated block - blocks are identified by indentation

# we can also loop through a specific list of values
for t in [1, 2, -1]:
    print(t)
# this should output 1, 2, -1
```

2.1.10 Python functions

CODE

Sometimes we want to repeat a complex bit of code in different places. So a loop won't do. The way to do that is to use a function:

```
def function_name(first_argument, second_argument): # there can be zero or more arguments
    return first_argument + second_argument # this function just returns the sum of its arguments
```

1. Function scope ADVANCED:CODE Whenever code is executed inside a function, the variables created there are only valid within the function. The function arguments are also effectively new variables. To see this, consider the following example.

```
def example_function(argument):
    # The following does not necessarily modify the original variable
    # passed. It depends on the effect of 'argument =
    # original_variable'. If it copies the value, then the original
    # variable remains the same. If it merely acts as a reference (as
    # is the case with lists and arrays) then a modification happens.
    argument += 1
    # this variable should not be visible outside the function
    hidden_variable = 2
    # variables defined outside the function are still readable!
    print(outside_variable)
    # but, we cannot affect the variables outside the
    # function. Otherwise there would be a mess.
    another_variable = 0
    # for that reason, functions should only use the arguments passed
```

```

        return argument
test = 100
outside_variable = "I am defined outside the function"
another_variable = -1 # so am I
foo = example_function(test) # line to unpack
print("test:", test)
print("foo:", foo)
print("outside_variable:", outside_variable)
#print("hidden_variable:", hidden_variable) # it complains of 'hidden_variable' not found

```

Let us unpack what happens. When we write `foo = example_function(test)`, what happens is as follows.

```

argument = test # create a new variable: all other variables are now hidden from scope
argument += 1 # execute the function's code block
foo = argument # apply the return operator to 'foo = example_function()'

```

2.1.11 Pandas and Histograms

PLOT:CODE

For this, we work on the Histogram example.

Pandas is a module for simple and efficient data I/O processing and visualisation. The following code snippet demonstrates a couple of features.

```

import pandas as pd # we need to load a library first
# loading data into pandas creates a data frame df
df['column-name'] # selects a column
df.hist() # creates a plot with many histograms

```

1. Coin example

ACTIVITY:PLOT

Plotting is also possible through the matplotlib. This is the module that pandas uses to plot stuff. It just has a simpler interface for doing so. But if you want to create custom plots, matplotlib is what you need to use.

```

X = [1, 0, 1, 0, 1, 1, 0, 1, 0] # a sequence of coin tosses.
import matplotlib.pyplot as plt # python has no default plot function, we must import
plt.hist(X) # this function plots the histogram

```

Each one of you should predict the result of a number of coin tosses. Let us do a histogram of the predictions. This is a binomial distribution.

- (a) The students record their data in the shared spreadsheet
- (b) Firstly, plot the histogram of the data with default settings.
- (c) What is the eff

Let us look at the student data: see `src/histograms/heights.ipynb`

2. Heights example

ACTIVITY

```
import pandas as pd
X = pd.read_csv("class-data.csv") # read the data into a DataFrame
X['Height (cm)'].hist() #directly plot the histogram
```

2.1.12 Histograms vs Pie Charts

While histograms are good visualisations of distributions on the real line, distributions over a discrete set of possible values are best-represented by a pie-chart. This especially if there is no relation between the different values. As an example, if the values are distinct categories, there is no particular reason to order them on an axis.

- What are the advantages and disadvantages of pie charts and histograms?

	Histogram	Pie Chart
To show proportions		
For more categories		
To compare relative size		
For real-valued data		

- Why is a 3D pie chart never a good idea?

```
#+BEGIN_SRC: python plt.pie(counts) # plot counts #+END_SRC
```

2.2 Time-Series: model the evolution of a system

A time series x_1, \dots, x_t is simply a sequence of variables. We typically assume that this is random. How can we capture this dependency between variables? Does the value of x_t depend only on the value of x_{t-1} ? On all the previous values? Only on the time index t ?

Frequently, sequential observations of a variable x_t are in fact noisy measurements of the true variable of interest, y_t , which we never observe. As an

example, consider covid infections. There is a true, underlying, number of infections, but we only ever measure the number of positive cases detected in a day.

Generally, there are three tasks associate with time series modelling, always given data up to this point, i.e. x_1, \dots, x_t .

1. Smoothing: What has happened in the past? Here we estimate y_{t-k} for $k > 0$.
2. Filtering: What is the current situation? To solve this problem we must estimate y_t .
3. Prediction: What will happen in the future? This involves predicting y_{t+k} for some $k > 0$.

These problems are all related and can be formalised in a statistical manner, and there are multiple algorithms that can be used to solve each problem. When $x_t = y_t$, then smoothing and filtering are trivial, but prediction is still an important problem. We focus here on a simple linear transformation such as the moving average as a basic solution method.

Smoothing For smoothing, a moving average filter is typically sufficient whenever $\mathbb{E}(x_t) = y_t$, i.e. x is just a zero-mean noisy measurement of y . Then we can construct the estimator

$$\hat{y}_t = \frac{1}{2n+1} \sum_{k=t-n}^{t+n} x_k.$$

Filtering When we wish to filter, at best we can take the moving average from the past n observations. If n is very large, then there is a corresponding delay between our filtering and the final prediction.

$$\hat{y}_t = \frac{1}{n+1} \sum_{k=t-n}^t x_k.$$

The only way to remove the lag is to perform a more complex transformation of the original data. To see this, consider the problem of prediction.

Prediction Prediction means estimating something in the future. This task is never trivial, even with perfect observations, i.e. when $x_t = y_t$. In this setting moving averages do not make sense. A simple idea is to *assume a linear trend*, e.g. that $y_{t+1} - y_t = y_t - y_{t-1}$. By re-arranging terms, we have that $y_{t+1} = 2y_t - y_{t-1}$. This gives us the estimator:

$$\hat{y}_{t+1} = 2x_t - x_{t-1}$$

2.2.1 Plotting lines

Here is a simple example of line plotting.

```
import numpy as np
X = [1, 2, 3, 4, 5, 4, 3, 2, 1] # define a small number of points
import matplotlib.pyplot as plt # import the plotting library
plt.plot(x) # perform a standard, simple plot
plt.savefig(f)
return f
```

What are such plots useful for?

2.2.2 Race times

ACTIVITY

https://en.wikipedia.org/wiki/1500_metres_world_record_progression

Wikipedia has a table that shows the progression of 1500m world records.

1. Let us first show the records up to 1950 .
 2. Try and predict the progression of world records on the board.
 3. Let us now look at the actual graph. Is it what you expected?
 4. How do you expect the progression to continue after 2020?
 5. How do you explain this progression? Can you find data to validate or refute your explanation?
1. Scraping tables example :example:data-collection:

```
import pandas
tables=pandas.read_html("URL") # read a table
# convert date-string:
dt = datetime.datetime.strptime(string, '%Y-%m-%d').year
# string manipulation
string.replace("+", "0") # replaces a + with a 0
string.split(":") # splits a string into multiple strings
# data formats
float("12.2"); # converts a number into a float
```

2.2.3 Example: The inclination of Mars

EXAMPLE

1. Plot Mars data
2. Show orbits
3. 3-body system, chaos and randomness

2.2.4 Example: Covid

EXAMPLE

1. Plot covid data.
2. Smooth the data: moving average plots
3. Try and estimate past, current and future infections with simple tools.
4. Discuss: Are those simple tools sufficient? Is our visualisation consistent? Do we need something further?

2.2.5 Example: Stock market prices

See: Trading Economics

2.3 Scatterplots: model a relationship

Let us start with an example where we just have three variables. We can plot the relationship between any two of them.

```
X=[1, 2, 3, 4, 10, 6]
Y=[5, 2, 5, 3, 1, 2]
Z=[0, 1, 0, 1, 0, 1]
import matplotlib.pyplot as plt
plt.scatter(X,Y)
```

Variables are frequently in some array instead.

```
import numpy as np
n_data = 10
n_features = 3
data = np.random.uniform(size=[n_data, n_features]) # create some random data
plt.scatter(data[:,0], data[:,1]) #plot the first against the second column
# We can always take a 'slice' of the data:
data[:, [1,2]] # get columns 1 and 2 and all the rows
data[1:10, [0,2]] # get columns 0 and 2 and all rows 1-10
```

```
## : means everything
## a:b means everything from a to b
## [a,b,c] means a, b and c.
```

In dataframes, it we can deal with multiple variables by name

```
import pandas as pd
df = pd.DataFrame(data, columns = ["Alcohol", "Caffeine", "Sugar"])
plt.scatter(df["Alcohol"], df["Caffeine"]) #plot the first against the second column
# getting slices is also possible in pandas dataframes, just slightly different:
df.loc[:, 'Alcohol'] # get column Alcohol
df.loc[:, ['Alcohol', 'Caffeine']] # get column Alcohol and Caffeine
```

2.3.1 Relationships as functions

A lot of relationships between two variables $x \in X$, $y \in Y$, can be described through some deterministic function $f : X \rightarrow Y$, i.e.

$$y = f(x).$$

If the relationship is one-to-one, then there exists an inverse function $f^{-1} : Y \rightarrow X$, so that

$$x = f^{-1}(y),$$

with

$$x = f^{-1}[f(x)].$$

Sometimes, however, the relationship between the two variables is not deterministic, that is the value of x does not uniquely determine the value of y , or the converse may occur... or both.

1. Physical relations

EXAMPLE

Many equations in physics relate two quantities. For example, there is the equation relating current I , voltage V and resistance R :

$$V = IR.$$

This relation can be inverted to obtain

$$R = V/I, \quad I = V/R.$$

Let us say that the resistance R is fixed. By altering the voltage (e.g. by adding more batteries to a circuit) we can see that the current increases.

2.3.2 Relationships as joint distributions

The simplest way to model a stochastic relationship between two variables is to model the joint distribution $P(X, Y)$. Consider the example of heights versus weights: We can expect that the taller a person is, the heavier they will be. However, their weight will depend on the mass of their muscle and adipose tissue. These in turn depend on their age, sex, genetics and lifetime calorie expenditure and intake.

1. Weight and height distribution

EXAMPLE

Here we can plot the number of people having a certain height and weight combination. This can be done with a colour-map. This is not much different than a normal histogram - and is called `hist2d` in `pyplot`:

```
X = 100 / (1 + np.exp(-np.random.normal(size=100))) + 125
Y = X * (1 + 0.1*abs(np.random.normal(size=100))) - 100
```

2.3.3 Relationships as conditional distribution.

Here we model the distribution of one variable given a fixed value for the other, e.g. $P(X|Y)$. The simplest thing is to only try to model the expected value $E[X|Y]$. How can we do this?

Method 1: polynomial fitting!

```
# returns the best fitting line to the data
a, b = np.polyfit(data_x, data_y, 1)
# Why is this the 'best' line? Because it minimises the total squared
# error between the predicted value and the actual ones.
# We can plot the line by this simple linear equation
ax = np.linspace(0,1)
plt.plot(ax, a * ax + b)
```

2.3.4 Example: Unemployment, GDP

Get some data financial data from FRED. This is time-series data. Can we actually make sense out of it in terms of correlations? Explore.

First, the unemployment rate: <https://fred.stlouisfed.org/series/UNRATE> Then, the GDP: <https://fred.stlouisfed.org/series/GDPC1> This has two different data frames.

```
# read the files
import pandas as pd
```

```
ur=pd.read_csv("UNRATE.csv")
gdp=pd.read_csv("GDPC1.csv")
# the date ranges are different, so we must try to merge them (inner join!)
merged = ur.merge(gdp)
```

3 Module 2: Experiment design (3 weeks)

3.1 Data collection and cleaning

ACTIVITY

3.2 Random sampling

In this module, we will perform the following activities:

1. Uniformly random sampling. How can we perform it?
2. Biased sampling, correcting for effects.
3. Importance sampling.

3.2.1 Survey of political opinions

ACTIVITY

You each support one of the following political groups:

- R: Red.
- B: Blue.

In this exercise, we will try to measure the support for different political parties.

1. Fixed affiliations

ACTIVITY

- (a) Deal cards so that 40% of the students are red, and 60% are blue.
- (b) Now sample the population in the following manner:
Everybody throws a die. Those with a value 5 or greater are part of our sample, The sample ω comes to the board.
The set of all possible samples is called the universe Ω . It is possible that we sample everybody, or nobody, or only the boys, or only the girls, or any combination.
Mathematically, we can say that $\Omega = \{0,1\}^n$, where n is the number of students, and $\omega_i = 1$ if a student has been selected.
- (c) Write a tick mark in the box saying "Red" or "Blue".

- (d) We then measure the proportion of red and blue votes. This proportion is the random variable!

We can repeat the same procedure with a different sampling method:

- (a) Assign a number to each student.
- (b) Cast a die and see which student it corresponds to.
- (c) The student tells me their vote.
- (d) I repeat.

2. Simulated sampling from the larger population ACTIVITY

Here we assign a random affiliation to each one of you. Throw a die for your political affiliation:

Die	Party
0	Red
1	
2	
3	
4	
5	
6	Blue
7	
8	
9	

We now count the number of people having different affiliations. This are your true voting affiliations. If you were to vote, then you would vote for these specific parties.

Here, the underlying random space is the combined dice throws of all the class, and the random variable of interest is the number of votes for each candidate.

Given the number of people in the course, what is the expected number of votes for each party?

3. Connecting probabilities of outcomes to probabilities PROBABILITY:THEORY:EXAMPLE

From a probability perspective, we can think

of the die as having random outcomes in $\Omega = [9]$. The random variable is the party vote $v : \Omega \rightarrow \{R, G, B\}$ with

$$v(\omega) = \begin{cases} R, & \omega \in \{0, 1, 2, 3\} \\ B, & \omega \in \{4, 5, 6, 7, 8, 9\} \end{cases}$$

Let us assume that the die has a uniform distribution P so that $P(\omega) = 1/10$ for all $\omega \in \Omega$ and $P(S) = |S|/10$ for all subsets $S \subset \Omega$. What is then the probability that somebody supports the Red party?

The probability that $v = R$, which we write informally as $\mathbb{P}(v = R)$, is simply the probability of all ω that lead to R , that is:

$$\mathbb{P}(v = R) = P(\{\omega : v(\omega) = R\}) = P(\{0, 1\}) = 2.$$

4. Expectations

PROBABILITY:THEORY:EXAMPLE

If we know the probability that a randomly chosen voter will vote for the i -party, then what is the probability of different numbers of votes? What is the expected number of votes for each party?

First, let us consider another random variable:

- n_i : the number of votes cast for each party i .

This total number of votes depends on the party affiliation of each voter. Let

- ω_t be the random die of person t .
- $v_t = v(\omega_t)$ is then the party affiliation of person t .

We collect the random die into one big vector

$$\omega = (\omega_1, \dots, \omega_t, \dots, \omega_T), \quad \omega \in \Omega^T, \omega_t \in \Omega$$

Then the total number of votes for party i is simply

$$n_i(\omega) = \sum_{t=1}^T \mathbb{I}\{f_t = i\} = \sum_{t=1}^T \mathbb{I}\{v(\omega_t) = i\}$$

Clearly, if there is only one voter, the expected number of votes for each party i is simply: $\mathbb{P}(v = i)$. If we had T voters then this should be $n_i = T \mathbb{P}(v = i)$. We can verify this by writing out the expectation.

$$\mathbb{E}_P[n_i] = \sum_{\omega} P^T(\omega) n_i(\omega)$$

Here P^T is the combined distribution of all persons dice. The main assumption we must make is that each person's die is independent of everybody else's. Then

$$P^T(\omega) = \prod_{t=1}^T P(\omega_t).$$

Consequently, the expectation becomes

$$\begin{aligned} \mathbb{E}_P[n_i] &= \sum_{\omega \in \Omega^T} P^T(\omega) n_i(\omega) \\ &= \sum_{\omega \in \Omega^T} P^T(\omega) \sum_{j=1}^T \mathbb{I}\{v_j = i\} \\ &= \sum_{\omega \in \Omega^T} \sum_{j=1}^T P^T(\omega) \mathbb{I}\{v_j = i\} \\ &= \sum_{j=1}^T \sum_{\omega \in \Omega^T} [P^T(\omega) \mathbb{I}\{v_j = i\}] \\ &= \sum_{j=1}^T P^T(\{\omega : v(\omega_j) = i\}) \\ &= \sum_{j=1}^T P(\{\omega_j : v(\omega_j) = i\}) \\ &= \sum_{j=1}^T P_v(i) \\ &= TP_v(i). \end{aligned}$$

The independence assumption is used to show that $P^T(\{\omega : v(\omega_j) = i\}) = P(\{\omega_j : v(\omega_j) = i\})$. The other steps do not require any assumptions.

5. The probability measure induced by a random variable f PROBABILITY:THEORY:ADVANCED:EXAMPLE

We already define a probability measure P on the outcomes Ω . Since f is a function on Ω , we can also define an appropriate probability measure P_f for the outputs of f . In particular, for any subset $S \subset \{R, G, B\}$, we define

$$P_f(B) = P(\{\omega : f(\omega) \in B\}).$$

We can this identify the informal probability $\mathbb{P}(f = R)$ with the measure $P_f(R)$. Here, even though R is an element, and not a set, we abuse notation. Normally we would write $P_f(\{R\})$ to denote the set consisting only of R .

6. Random sampling

ACTIVITY

Each one of you throws a second die and records the outcome. We now have

Die	Party	Response
0	Red	Not Reachable
1-2		Refuse
3-9		Green
0-4	Blue	Not Reachable
5		Refuse
6-9		Blue

- Make a histogram / bar-char pie plot on the number of votes

3.2.2 Uniform sampling

CODE

In uniform sampling, the probability of all outcomes is the same.

Sampling with replacement In sampling with replacement, each outcome can appear more than once.

```
import numpy as np
population_size = 10 # say we have 10 people we want to sample from
n_samples = 5 # say we want to take 5 samples from the population the
# following will give us a sample drawn with replacement: It doesn't
# matter who we selected before, the next one will be randomly
# selected independently of the previous selection
sample = np.random.choice(population_size, size = n_samples)
return sample
```

The following code has the same effect

```

import numpy as np
population_size = 10 # say we have 10 people we want to sample from
n_samples = 5 # say we want to take 5 samples from the populatoin
sample = np.zeros(n_samples)
for i in range(n_samples):
    sample[i] = np.random.choice(population_size)
return sample

```

Sampling without replacement In sampling without replacement, each outcome can appear at most once.

```

import numpy as np
population_size = 10 # say we have 10 people we want to sample from
n_samples = 5 # say we want to take 5 samples from the population the
# following will give us a sample drawn with replacement: It doesn't
# matter who we selected before, the next one will be randomly
# selected independently of the previous selection
sample = np.random.choice(population_size, size = n_samples, replace='False')
return sample

```

3.2.3 Expectation

THEORY:PROBABILITY

Recall that a random variable f is a function $f : \Omega \rightarrow \mathbb{R}$. The expectation of a random variable with underlying distribution $P(\omega)$ is simply

$$\mathbb{E}_P[f] \triangleq \sum_{\omega \in \Omega} f(\omega)P(\omega).$$

There is nothing random about the variable itself, it is only the random input that makes its value random.

In the following example, $\omega \in \{0, 1, 2, 3\}$, with $P(\omega)$, specified by the vector $P[]$

z1

```

import numpy as np
# Let us define the space Omega
Omega = np.array([0, 1, 2, 3])

# Let us define a vector P so that P[omega] is the probability of omega
# e.g the probability that omega = 0 is 0.1, and that omega = 3 is 0.5.
P = np.array([0.1, 0.2, 0.3, 0.4])

```

```

# Here is our random variable f(). It is just a function of omega. There
# is nothing random about it, only omega is random!
def f(omega):
    return omega * omega

# Let us generate a random omega:
random_outcome = np.random.choice(0mega, p = P)
random_variable_value = f(random_outcome)
print("omega:", random_outcome,
      "f:", random_variable_value)

# We can also easily calculate the expectation of the random variable
# through the dot product.
print ("Expected value:", np.dot(P, f(0mega)))

```

1. What is the expectation of this variable?

EXERCISE

In our case, $f(\omega) = \omega^2$. Let us first consider a discrete Ω :

Let $\omega \in \{0, 1, 2\}$ and $P(1) = 0.3$ and $P(2) = 0.2$. Then

$$\begin{aligned}
 \mathbb{E}_P(f) &= \sum_{\omega} P(\omega) f(\omega) \\
 &= P(0)f(0) + P(1)f(1) + P(2)f(2) \\
 &= 0.5 \times 0^2 + 0.3 \times 1^2 + 0.2 \times 2^2 \\
 &= 0.3 + 0.8 = 1.1.
 \end{aligned}$$

Let us now consider a continuous Ω with probability measure $P(A)$ for all (measurable) subsets A of Ω . This can be defined through the probability density $p(\omega)$,

$$P(A) = \int_A p(\omega) d\omega.$$

The corresponding expectation of f is then given by

$$\mathbb{E}_P(f) = \int_{\Omega} f(\omega) p(\omega) d\omega$$

For our specific example, let us choose p to be the uniform distribution on $[0, 2]$. Then $p(\omega) = 1/2$ and

$$\mathbb{E}_P(f) = \int_0^2 \omega^2 / 2 d\omega = [x^3 / 6]_0^2 = 2^3 / 6 = 8 / 6 = 4 / 3.$$

2. Centime exercise

ACTIVITY:ADVANCED

A jar with coins is passed around the class.

- (a) The students are asked to guess how many coins it contains.
- (b) The students agree on a 50% confidence interval.
- (c) The students fit a normal distribution on this interval $[\mu - \frac{2}{3}\sigma, \mu + \frac{2}{3}\sigma]$.
- (d) Is this normal distribution a good choice? Are you 90\
- (e) Is a normal distribution generally appropriate?
- (f) Puzzle: Guess how many coins there are. If correct, then the class will share the money. If not, they will get nothing. What is the correct guess?

(If students have trouble with this, try with small numbers of coins and finite number of possibilities - demonstrate by playing the guessing game repeatedly)

3.2.4 Survey of heights

ACTIVITY

- 1. First, randomly select a student. How? Everybody gets a different number. Then I throw a die until I get a single student matching this number. The student is my sample ω .
- 2. The student comes to the board and measures his height. This height is a random variable $h(\omega)$. Each student has their own fixed height, but which student I select is random. Thus, the height that I measure is random.
- 3. Repeat the experiment once more.
- 4. Now we randomly select multiple students. How? Each student throws a die. If the die is >4 , then they are selected. All the students come to the board. They are our sample ω .
- 5. We now write the student heights, and average. The average height of the sample is our random variable.

3.3 The data science pipeline

The experimental pipeline has a number of different components.

1. Formulating the problem.
2. Deciding what type of data is needed.
3. Choosing the model and visualisation needed.
4. Designing the experimental protocol.
5. Generating data confirming to our assumptions.
6. Testing the protocol on synthetic data. Is it working as expected?
7. Putting the protocol through on real data.

3.3.1 Salaries of men and women

EXAMPLE

You want to measure if men and women have different salaries, as well as the reasons why. In particular, you want to check if men are paid more than women on average, and if this can be explained purely by age.

1. Formulate the problem

Let us define three variables, gender g , age a and salary s . Our possible hypotheses can be expressed mathematically as follows:

For background, check out Hypotheses about the expected salary relate to conditional expectation

- $E(s|g = m) = E(s|g = f)$. Men and women have the same salary in expectation
- $E(s|g = m) > E(s|g = f)$. Women have a lower salary in expectation. But how much lower?

(Note that a stronger hypothesis would involve the conditional probabilities rather than the expectation)

If we find a difference, then we may want to **explain** it. If the age is a **sufficient explanation** for the salary, then the salary is conditionally independent of the gender given the age:

- $E(s|a, g) = E(s|a)$. The salary is independent of the gender, given the age.

2. What type of data do we want?

We need salaries, age and gender.

3. Model and visualisation

Since we are looking at means, then it is maybe enough to look at averages. Averages are single numbers, so maybe a bar plot is enough.

4. Experimental protocol

- (a) Collect data.
- (b) Plot average for men and women. Use the simulation to find the right type of plot.
- (c) If average salary is very different (how different?) then say women are paid less than men. Use the simulation to get an idea.

5. Generate data according to our assumptions

We can have three different tests:

- (a) Generate everything independently

```
import numpy as np
# Generate data
def identical_populations(n_samples):
    gender = np.random.choice(samples = n_samples)
    age = 18 + 60*np.random.beta(3,4, size = n_samples)
    salary = 100* np.random.exponential(100, size = n_samples)
    return age, gender, salary
```

- (a) Make the salary depend on the gender

```
def different_populations(n_samples):
    gender = np.random.choice(samples = n_samples)
    age = 18 + 60*np.random.beta(3,4, size = n_samples)
    salary = (100 + gender*10)* np.random.exponential(100, size = n_samples)
    return age, gender, salary
```

- (a) Make the salary depend on the age, but have fewer women working after some age.

```
def age_effect_populations(n_samples):
    age = 18 + 60*np.random.beta(3,4, size = n_samples)
    q = (age - 18) / 60
```

```

p = q * 1 + (1 - q)*0.5
gender = np.random.choice(2, p = [p, 1-p], samples = n_samples)
salary = (100 + age)* np.random.exponential(100, size = n_samples)
return age, gender, salary

```

See Salary example notebook.

If we estimate the means, we see there is a lot of variability. How can we fix that? One idea is to perform a lot of simulations and note how much variability we do have. Another is to use a bootstrap sample.

3.3.2 A question of voting

EXAMPLE

An election is coming in 100 days. A political party supporting one of the options [assume it's just a yes/no vote] in the election gives you 10,000 CHF to spend over these 100 days so as to measure the mood of the population. They can use that information to increase their chances of success. How should you do the study?

1. Formulate the problem ACTIVITY Who is going to use our visualisation? How are they going to use it? Let us brainstorm a little bit about this.
2. What type of data do we want? ACTIVITY What do we want to know about the people we collect the data from?
3. Choose the model and visualisation needed ACTIVITY What can we assume about people's opinions? How about their responses? Are they truthful? What is the most useful visualisation?

After we get the data, we need to analyse it. A simple model would be a moving average of polls over time. Would that work?

4. Design the experimental protocol ACTIVITY
 Assume we need to pay 1 CHF for every time we ask a poll question, and we have a budget of 10,000 CHF. Then, how should we ask questions, assuming the election is in 100 days from now?
 (a) Ask 10,000 people now. (b) Ask 10,000 people one day before the election. (c) Ask 10,000 people 50 days from now. (d) Ask 100 people every day? (e) Ask 1,000 people every 10 days.

5. Generate data according to our assumptions

We can assume a simple model of the electorate here... what should it be? Maybe their opinion changes over time. Maybe some people are not responding, or not reachable. Which one corresponds to our assumptions?

6. Test the protocol

After testing the whole pipeline, we can see if it actually works as intended.

4 Module 3: Inference (2 weeks)

4.1 Logic and bar graphs

THEORY

How many people are taking CS? How many are in other courses? Let us consider the following statements. Each statement is either true or false. If a statement is true, it has the value 1. It consequently has the value 0 if it is false. For each student i , we define:

- m_i : the student is male
- f_i : the student is female
- s_i : the student is in the faculty of science
- l_i : the student is in the faculty of law
- e_i : the student is in the faculty of economics

Clearly, if s_i is true then l_i, e_i are false.

4.1.1 Logic and sets

THEORY

We can associate events with sets in a universe. The following laws apply:

1. There is a universe Ω of possible outcomes
2. Consider an individual outcome $\omega \in \Omega$:

(a) If $\omega \in A$, then we say that A is **true**. (b) If $\omega \notin A$, then we say that A is **false**.

1. Let $\neg A$, read 'not A', be the complement $\Omega \setminus A$: If A is true, then $\neg A$ is false and vice-versa.

2. B is a subset of A (i.e. $B \subset A$) iff B **implies** A .
3. If A and B are disjoint, i.e. $A \cap B = \emptyset$ (they have an empty intersection) then A, B are **mutually exclusive**. That means that it is impossible for A, B to both be true.

4.2 Conditional Probability and the Theorem of Bayes

4.2.1 Bayes theorem

THEORY:PROBABILITY

- Recall the definition of Conditional probability:

$$P(A|B) = P(A \cap B)/P(B)$$

i.e. the probability of A happening if B happens.

- It is also true that:

$$P(B|A) = P(A \cap B)/P(A)$$

- Combining the two equations, reverse the conditioning:

$$P(A|B) = P(B|A)P(A)/P(B)$$

- So we can reverse the order of conditioning, i.e. relate to the probability of A given B to that of B given A .

4.2.2 The cards problem

1. Print out a number of cards, with either $[A|A]$, $[A|B]$ or $[B|B]$ on their sides.
2. If you have an A , what is the probability of an A on the other side?
3. Have the students perform the experiment with:
 - (a) Draw a random card.
 - (b) Count the number of people with A .
 - (c) What is the probability that somebody with an A on one side will have an A on the other?
 - (d) Half of the people should have an A ?
1. The prior and posterior probabilities

A	A	2/6	A observed	2/3
A	B	1/6	A observed	1/3
B	A	1/6		
B	B	2/6		

4.3 Bayesian analysis example

4.3.1 The murder problem

- Somebody saw somebody matching their description and he was found in the neighbourhood. There is no other evidence.
- There are two possibilities:
 - H_0 : They are innocent.
 - H_1 : They are guilty.

What is your belief that they have committed the crime?

1. Prior elicitation

- All those that think the accused is guilty, raise their hand.
- Divide by the number of people in class
- Let us call this $P(H_1)$.
- This is a purely subjective measure!

4.3.2 DNA test

- Let us now do a DNA test on the suspect

1. DNA test properties

- D : Test is positive
- $P(D|H_0) = 1\%$: False positive rate
- $P(D|H_1) = 100\%$: True positive rate

2. Run the test

- The result is either positive or negative ($\neg D$).
- What is your belief **now** that the suspect is guilty?

4.3.3 Everybody is a suspect

- Run a DNA test on everybody.
- What is different from before?
- Who has a positive test?
- What is your belief that the people with the positive test are guilty?

4.3.4 Explanation

- Prior: $P(H_i)$.
- Likelihood $P(D|H_i)$.
- Posterior: $P(H_i|D) = P(D \cap H_i)/P(D) = P(D|H_i)P(H_i)/P(D)$
- Marginal probability: $P(D) = P(D|H_0)P(H_0) + P(D|H_1)P(H_1)$
- Posterior: $P(H_0|D) = \frac{P(D|H_0)P(H_0)}{P(D|H_0)P(H_0) + P(D|H_1)P(H_1)}$
- Assuming $P(D|H_1) = 1$, and setting $P(H_0) = q$, this gives

$$P(H_0|D) = \frac{0.1q}{0.1q + 1 - q} = \frac{q}{10 - 9q}$$

- The posterior can always be updated with more data!

4.3.5 Python example

```
# the input to the function is the prior, the likelihood function, and posteriors
# Input:
# - prior for hypothesis 0 (scalar)
# - data (single data point)
# - likelihood[data][hypothesis] array unction
# Returns:
# - posterior for the data point (if multiple points are given, the calculation is rep
def get_posterior(prior, data, likelihood):
    marginal = prior * likelihood[data][0] + (1 - prior) * likelihood[data][1]
    posterior = prior * likelihood[data][0] / marginal
    return posterior

import numpy as np
```



```

prior = 0.9
likelihood = np.zeros([2, 2])
# pr of negative test if not a match
likelihood[0][0] = 0.9
# pr of positive test if not a match
likelihood[1][0] = 0.1
# pr of negative test if a match
likelihood[0][1] = 0
# pr of positive test if a match
likelihood[1][1] = 1
data = 1
return get_posterior(prior, data, likelihood)

```

4.4 More general problems

4.4.1 The k -meteorologists problem

We have k meteorological stations.

1. Predictions and outcomes

Station	M	T	W	T
MeteoSuisse	25%	20%	10%	5%
Wunderground	30%	50%	20%	10%
AccuWeather	90%	70%	10%	0%
Rain	Y	N	N	Y

2. H_i : The i -th station's model is correct

- $P(H_i)$: prior
- $P(D|H_i)$: likelihood according to station i
- $P(H_i|D)$: posterior

4.4.2 Types of hypothesis testing problems

1. Simple Hypothesis Test Example: DNA evidence, Covid tests
 - Two hypotheses H_0, H_1
 - $P(D|H_i)$ is defined for all i
2. Multiple Hypotheses Test Example: Model selection

- H_i : One of many mutually exclusive models
 - $P(D|H_i)$ is defined for all i
3. Null Hypothesis Test Example: Are men's and women's heights the same?
- H_0 : The 'null' hypothesis
 - $P(D|H_0)$ is defined
 - The alternative is **undefined**

4.4.3 Pitfalls

1. Problem definition
 - Defining the models $P(D|H_i)$ incorrectly.
 - Using an "unreasonable" prior $P(H_i)$
2. The garden of many paths
 - Having a huge hypothesis space
 - Selecting the relevant hypothesis after seeing the data

4.4.4 The covid test problem

ACTIVITY

10% of the class has covid, i.e. $P(\text{covid}) = 0.1$. Each one of you performs a covid test. If you have covid, the test is correct 80% of the time, i.e. $P(\text{positive} | \text{covid}) = 0.8$. Conversely, if you do not have covid, there is still a 10% chance of a positive test, with $P(\text{positive} | \text{not-covid}) = 0.1$

How likely is it that you have covid if your test is positive or negative, i.e. $P(\text{covid} | \text{positive})$, vs. $P(\text{covid} | \text{negative})$?

First of all, each one of you should independently generate a uniform random number between 1 and 10. For that, you can each throw a die, and record the outcome.

Then you throw a second die, and record that as well.

I will now pass over the tables and tell each one of you if they have a positive test.

Now, everybody with a positive test raises their hand. I expect it to be slightly more than 10% (but it depends).

4.5 Hypothesis testing

Hypothesis testing, or model selection, is a general problem in machine learning and statistics.

4.5.1 Estimation intervals

Intervals are a generalisation of hypothesis tests. We want to know if a given unknown parameter is within some range. Let us give an example.

1. Estimating an expectation We might have data x_1, \dots, x_T , with $x_t \sim \text{Ber}(\theta)$ from a Bernoulli distribution with parameter θ . While we can calculate the mean from the samples, what can we actually say about the unknown expectation $\mathbb{E}[x_i] = \theta$?
2. Bayesian credible interval The Bayesian idea is intuitive. Through the posterior, we can calculate the probability that the true parameter is in any interval $S \subset [0, 1]$:

$$P(\theta \in S | x_1, \dots, x_T)$$

This idea works for any kind of estimation.

3. Confidence interval Confidence intervals work the other way round. First we construct an algorithm for obtaining an estimate of the unknown parameter, e.g. the empirical mean estimate $\mu(x_1, \dots, x_T)$. The confidence interval tells us the probability that our estimate is further away than ϵ from the expectation, as an analytic function of ϵ and the amount of data n

$$P(|\mu - \theta| > \epsilon | \theta) \leq \delta(\epsilon, n)$$

4. Difference between confidence intervals and credible intervals.

In confidence intervals, we condition on the unknown parameter θ , but the function δ is constructed so that it is independent of θ , since that is unknown.

4.6 Simulation study

ACTIVITY

For a simple visualisation problem, vary parameter values and simulate thousands of times under each set of conditions. Summarise your findings graphically.

4.7 Introduction: Estimation, Prediction and Decision Making

4.7.1 Correlation example

STOP GLOBAL WARMING: BECOME A PIRATE

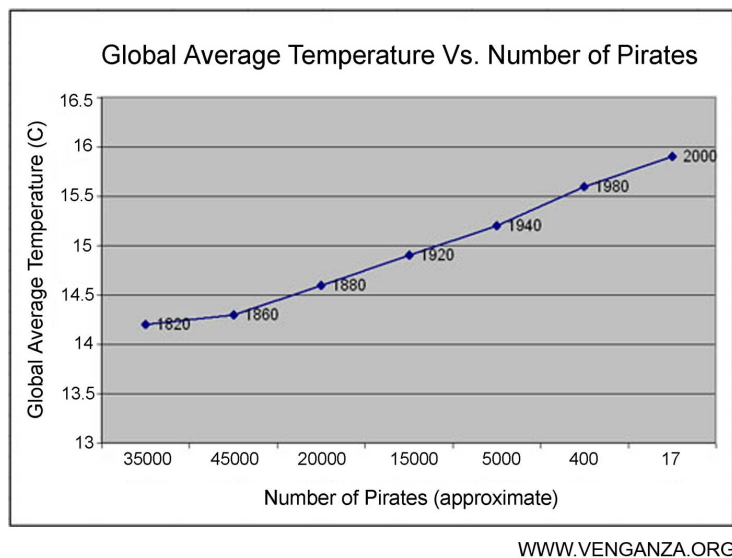


Figure 1: Pirates cool the earth?

:BEAMER_{env}: ignoreheading

4.7.2 The three main problems in statistics

1. Estimation

- How many covid patients are there now?
- Which is the best meteorological station?
- What is the right model for planetary motion?

2. Prediction

- How many covid patients will there be next week?
- Will it rain tomorrow?

- When is the next lunar eclipse?

3. Decision making

- What covid measures should I take?
- Should I take the umbrella?
- What's the best way to reach the next eclipse?

4.7.3 Types of variables

1. Observed variables

- Telescope images
- Survey data from a health authority
- Radar measurements in an AV

2. Latent variables

- The mass, position and velocity of Mars
- Actual number of covid patients
- The location and speed of nearby vehicles

3. Decision variables

- Can be observed or latent
- Selected arbitrarily by humans or machines

4.7.4 From correlation to causation.

1. Statistical relation.

- Is smoking related to lung cancer?

2. Statistical causation.

- Does smoking cause lung cancer?
- Will I get lung cancer if I smoke?

3. Actual causes.

- Did I get lung cancer because I smoked?

4.7.5 Correlation versus Causality

1. Correlation and dependence
 - Does **knowing** the value of one variable give us **information** about another variables?
2. Causality
 - Does **changing** the value of one variable **affect** other variables?

4.8 Joint and Conditional Probabilities

4.8.1 Correlation: Setting

1. Outline
 - Random variables x, y taking values in \mathcal{X}, Y .
 - We write $\mathbb{P}(x, y)$ to informally mean their joint distribution.
2. Formal definition
 - Underlying probability space P, Ω, Σ with
 - Outcome space Ω
 - Event space Σ , so that $A \in \Sigma$ are subsets of Ω .
 - Probability measure $P : \Sigma \rightarrow \Omega$.
 - RVs $x : \Omega \rightarrow \mathcal{X}, y : \Omega \rightarrow \mathcal{Y}$
 - Joint measure $P_{x,y}(S_x, S_y) \triangleq P(\{\omega : x(\omega) \in S_x, y(\omega) \in S_y\})$

4.8.2 Conditional probabilities

1. Definition The conditional distribution of x given y is:

$$\mathbb{P}(x|y) = \mathbb{P}(x, y) / \mathbb{P}(y)$$

Thus, for every value of y we get a different distribution for x .

2. Recall definition for events This has the same form.

$$P(A|B) = P(A \cap B) / P(B).$$

4.8.3 Python example: independent RVs

Let us begin with a simple example of two variables which are independent Bernoulli.

1. Bernoulli-distributed $x, y \in \{0, 1\}$

- $\mathbb{P}(x = 1) = \theta$
- $\mathbb{P}(y = 1) = v$

The independence can also be seen in the structure of the program: x and y are generated by two separate calls to `random.choice()`

2. One draw of x, y

```
import numpy as np
theta = 0.6
v = 0.8
x = np.random.choice(2, p = [1 - theta, theta])
y = np.random.choice(2, p = [1 - v, v])
```

4.8.4 Dependent, Discrete x, y

Now let us look at dependent variables. Here x is generated first. Then y has a value which depends on what value of x has been generated.

1. Bernoulli-distributed $x, y \in \{0, 1\}$

- $\mathbb{P}(x = 1) = \theta$
- $\mathbb{P}(y = 1|x = 0) = v_0$
- $\mathbb{P}(y = 1|x = 1) = v_1$

We can also calculate $\mathbb{P}(y = 1)$ by using the marginalisation rule:

$$\mathbb{P}(y = 1) = \mathbb{P}(y = 1|x = 0)\mathbb{P}(x = 0) + \mathbb{P}(y = 1|x = 1)\mathbb{P}(x = 1) = v_0(1-\theta) + v_1\theta.$$

Now let us how to do this programmatically. Here, x is an explicit input to the function that generates y :

2. One draw of x, y

```

import numpy as np
theta = 0.6
v = np.zeros(2)
v[0] = 0.4
v[1] = 0.8
x = np.random.choice(2, p = [1 - theta, theta])
y = np.random.choice(2, p = [1 - v[x], v[x]])
return x,y

```

4.8.5 Python example: multiple draws

Doing this for multiple draws is a bit more complicated, but it can be done. We just need to use the set-builder notation to generate one y for each x generated.

```

import numpy as np
n = 10000
theta = 0.6
v = np.zeros(2)
v[0] = 0.4
v[1] = 0.8
x = np.random.choice(2, p = [1 - theta, theta], size = n)
y = np.array([np.random.choice(2, p = [1 - v[x_t], v[x_t]]) for x_t in x])
import matplotlib.pyplot as plt
A = np.zeros([2,2])

for i in range(2):
    for j in range(2):
        A[i,j] = sum((x==i) & (y==j))

plt.imshow(A)
plt.savefig("correlated-binary.png")
plt.show()
return A

```

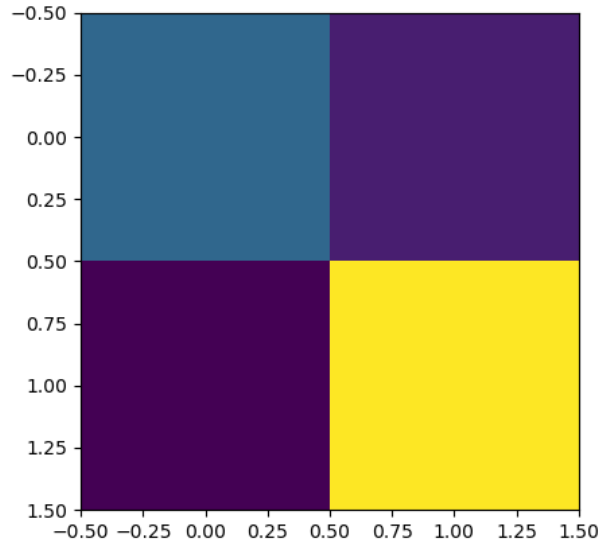



Figure 2: Here $x \sim \text{Bern}(0.8)$ and $y \sim \text{Bern}(0.9x)$.

4.8.6 Empirical joint probability of x, y

4.8.7 Empirical joint probability of x, y

4.8.8 Continuous x, y

This is the typical structure of regression problems. Let us start with a simple normal distribution for both x, y :

1. Normal-distributed x, y
 - $x \sim \text{Normal}(0, 1)$.
 - $y|x \sim \text{Normal}(x, 1)$.

The above notation means that:

- x has a normal distribution with mean 0 and variance 1.
 - Given the value of x , y has a normal distribution with mean x and variance 1.
2. One draw from x, y This dependence is made explicit in the following code.

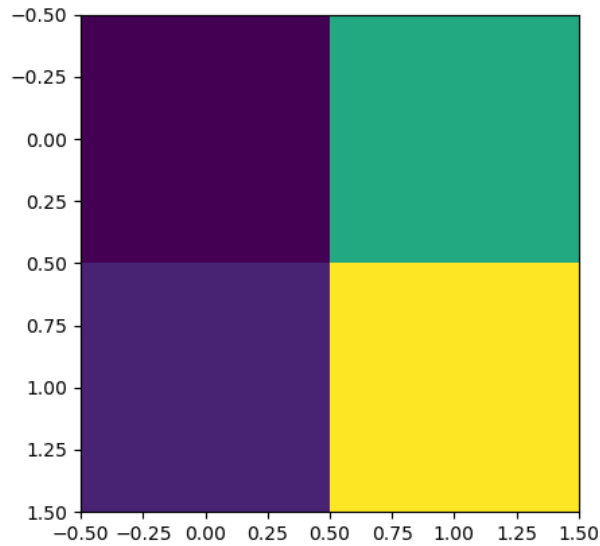


Figure 3: Here $x \sim \text{Bern}(0.8)$ and $y \sim \text{Bern}(0.1)$.

```
import numpy as np
theta = 0.8
x = np.random.normal(0, 1)
y = np.random.normal(x, 1)
return x,y
```

4.8.9 Continuous x , Discrete y

This is the typical structure of classification problems.

1. Normal-distributed x , Bernoulli-distributed y

- $y \sim \text{Bernoulli}(0.6)$
- $x|y \sim 160 + \text{Normal}(10 * y, 1)$.

The above notation means that:

- y has a Bernoulli distribution with mean 0.6.
- Given the value of y , x has a normal distribution with mean $160 + 10x$ and variance 1.

2. One draw from x, y

```
import numpy as np
y = np.random.choice(2, p = [0.4, 0.6])
x = np.random.normal(x, 1)
return x,y
```

4.8.10 Covariance matrix

- Consider a collection of RVs x_1, \dots, x_n .
 - The joint distribution is a complicated object.
 - Visualised with scatterplots (x_i, x_j) , e.g. `sns.pairplot()`
1. Covariance matrix C Instead, we can calculate the covariance matrix

$$C_{ij} = \frac{\mathbb{E}\{[x_i - \mathbb{E}(x_i)][x_j - \mathbb{E}(x_j)]\}}{\sqrt{\mathbb{E}\{[x_i - \mathbb{E}(x_i)]^2\} \mathbb{E}\{[x_j - \mathbb{E}(x_j)]^2\}}}$$

2. Approximating the covariance matrix
 - Assuming data $x(t)$ with components $x_i(t)$:
 - $C_{ij} \approx \frac{1}{T} \sum_{t=1}^T [x_i(t) - \mu_i][x_j(t) - \mu_j] / \sigma_i \sigma_j$.
 - μ_i : (empirical) mean of x_i
 - σ_i : (empirical) standard deviation of x_i

4.8.11 Correlation versus dependence

1. Dependent random variables
 - x, y are independent if $\mathbb{P}(x, y) = \mathbb{P}(x) \mathbb{P}(y)$
 - equivalently, if $\mathbb{P}(x|y) = \mathbb{P}(x)$
 - x, y are dependent if they are not independent.
2. Correlated random variables
 - x, y are uncorrelated if $\mathbb{E}(x, y) = \mathbb{E}(x) \mathbb{E}(y)$
 - Equivalently, if $\mathbb{E}(x|y) = \mathbb{E}(x)$
 - x, y are correlated if $\mathbb{E}(x, y) \neq \mathbb{E}(x) \mathbb{E}(y)$
3. Theorem
 - If x, y are correlated then they are dependent.
 - If x, y are independent then they are uncorrelated.

4.9 Models of Causation

4.9.1 Causal inference vs the actual cause

1. Causal inference
 - Can aspirine cure headaches?
 - Does smoking cause lung cancer?
 - Or do cancer patients become smokers?
 - Or is there a third factor causing both?
2. The actual cause
 - Did aspirin cure **my** headache?
 - Did smoking cause **my** cancer?
3. Applications
 - Causal inference useful in a scientific setting.
 - Reliable methods for causal inference exist.
 - Actual causes useful in a legal setting.
 - No reliable method or definition exists for determining actual causes.

4.9.2 Confounding variables

1. Arrival at work
 - Tom and Fatima both work in Lausanne.
 - Whenever Tom is late to work, so is Fatima.
 - When this happens, there is also a traffic jam.
2. Kidney stone treatment
 - Treatment A is effective 90% of the time
 - Treatment B is effective 50% of the time.
 - Why is that? (see src/causality/confounders)x

Treatment	Small Stones	Large Stones
A	80%	30%
B	90%	60%

We see here that treatment B is always better, but Large Stones are harder to treat. Let us now consider the treatment policy:

Treatment	Small Stones	Large Stones
A	80%	20%
B	20%	80%

The treatment policy mostly assigns treatment A when patients have small stones, and B for large stones. Hence, as treatment B is applied to harder cases, it might look worse on average. In particular, if we calculate the average treatment effect, by marginalising over stones, we have

Treatment	Average Effect
A	0.7
B	0.66

So, even though the second treatment better both for large and small stones, if we ignore the stone size, it looks worse.

Ignoring the stone size makes it a **confounder**, because it was used in the policy, but we ignore it. It wouldn't matter if stone size was not used to select the treatment, which is why we typically use randomised trials to study the effect of treatments.

4.9.3 Instrumental variables

1. Setting
 - a : Treatment
 - x : Patient characteristics
 - y : Outcome
2. Confounding
 - x is ignored
 - a depends on x .
 - y depends on a, x .
3. Using x as an instrument
 - x is taken into account
 - We model $\mathbb{P}(y|x, a)$, $\mathbb{P}(a|x)$

4.9.4 Avoiding confounding

1. Randomised studies
 - Treatment is random
 - Hence dependence to x is cut.
2. Selecting an instrument
 - Must be an informed choice
 - A complex instrument can lead to spurious correlations
3. Explicitly model policies π
 - Treatment policy π with $\pi(a|x, \pi)$.
 - Outcomes depend on policy π

5 Module 4: Advanced visualisation (2 weeks)

5.1 Geographical data

https://scikit-learn.org/stable/auto_examples/neighbors/plot_species_kde.html

5.1.1 Colour maps

1. Colour as a continuous variable.
2. Colour as a discrete variable.
3. Colour perception and interpretation.

5.1.2 Contour maps

1. Geographical contour
2. Density plots

5.2 Text data

6 Module 5: Data analysis in practice (2 weeks)

6.1 Visualising fMRI data

6.2 Visualising GWAS data

7 Module 6: Project work and presentations (2 weeks)

8 Assignments

The course contains assignments and a project. The instructions for each assignment are given below. The assignments are largely done in class, but completed at home.

8.1 Table To Picture

TLDR: Find a table in wikipedia on a topic of interest, and convert the table into a graph.

The purpose of this assignment is for you to create a graphic that demonstrates something interesting about a data table found on the web. Please provide as precise and concise answers as possible. This assignment is **graded**. You are *encouraged* to discuss the assignment with other students in a group. However each student must prepare their own *individual* report. Please submit your answers on *moodle*

8.1.1 Instructions

In this exercise, you must create a plot from an existing dataset and write a short report. Use the following steps as a guideline.

1. Find a data table on the internet.
2. Write a short description of the data on the table.
3. Create one or two plots of your choice, summarising the data on the table.
4. Explain what the graph shows about the data.
5. Try and draw some conclusions or generalisations from the graph. Does it make logical sense?

8.1.2 Example

This is an example of this exercise for a dataset we already saw in class.

1. Use the world records data for 100m/400m/1500m or some other distance.
2. Explain what the records show.
3. Show how the world record changes over time for men and women, with different colours. Be sure to plot records with the x axis showing time.
4. Draw regression lines over the world record graph: the records reduce over time.
5. It is not logically possible to expect the times to reduce linearly with time! Is there a fundamental limit? How can the data be best explained? Is human performance constant over time, and records are falling due to random chance? Is human performance slowly increasing over time?

8.2 Plot deconstruction

8.2.1 Instructions

Find an interesting plot from a web page on e.g. wikipedia. Try to identify some problem with the plot. To help you, ask yourself the following questions:

- Is the plot type appropriate?
- Is the data correct?
- Does the plot convey an appropriate message?
- Is there more data somewhere that you could combine with the original to obtain a better picture?

After you have identified problems with the plot, data sources, or missing data, create a new plot, along with an explanation of how you addressed the original plot's deficiencies.

8.2.2 Example

8.3 Newspaper article analysis

In this assignment you will read a newspaper article with some statistics and visualisations, and try to interpret what it says. You must study the article critically. Are the conclusions supported by the data? Does the methodology make sense? Find primary sources that confirm or challenge the article to obtain a more rounded picture.

Here is a list of possible articles you can use. Feel free to suggest your own article and add it to the list.

https://docs.google.com/spreadsheets/d/1QKj_L9f0UIH80qgs2kcjc8AU1eKZzsT0cYFSU06HJBYedit#gid=0

8.4 Open project

8.4.1 Project proposal

Propose a problem to solve, including:

- Hypotheses to test
- How to collect data
- How to analyse the collected data

8.4.2 Project Highlight

After you have started your project, each one of the project members presents a preliminary plot and explains it (5 minutes).

8.4.3 Project presentation

8.4.4 Project report

The completed project should include a report written by both students in the team. This should address the points in the Assessment description.

9 Background

9.1 Notation

For convenience, I include necessary mathematical notation

9.1.1 Sets

- \mathbb{R} : Real numbers
- \mathbb{R}^d : d-dimensional Euclidean space
- \emptyset : The empty set
- $A \subset B$: A is a subset of B.
- $A \cap B$: The intersection of A and B
- $A \cup B$: The union of A and B
- $A \setminus B$: Removing B from A
- Ω : The "universe"
- $A^c = \Omega \setminus A$: The complement of a set.
- $\{x|f(x) = 0\}$: The set of x so that $f(x) = 0$.

9.1.2 Analysis

- $\mathbb{I}\{x \in A\}$: indicator function (takes the value 1 if $x \in A$, 0 otherwise)
- $\sum_{x \in X} f(x) = f(x_1) + \dots + f(x_n)$, with $X = \{x_1, \dots, x_n\}$
- $d/dx f(x)$: derivative of f
- $\partial/\partial x f(x, y)$: partial derivative of f
- $\nabla_x = (\partial/\partial x_1, \dots, \partial/\partial x_n)$, vector of partial derivatives.

9.1.3 Probability

- \mathbb{P} : Probability (informally generally)
- \mathbb{E} : Probability
- P : A probability measure
- p : A probability density
- $P(A|B) = P(A \cap B)/P(B)$. Conditional probability, $A, B \subset \Omega$.
- θ : Parameter

- Θ : Parameter set
- $\{P_\theta | \theta \in \Theta\}$: A family of parametrised models
- $\mathbb{P}(x|y)$ conditional probability for random variables x, y (generally)

9.2 Probability background

The theory of probability is used to mathematically define processes with uncertain outcomes. The set of all possible outcomes depends on the process. For example, if we throw a die, this is the set of all possible ways, locations etc that the die can fall and land. However, we may only be interested in two events: whether the die lands showing a '6' or not. Formally, an event A is a subset of all the possible outcomes Ω . In our example, A can be the set of all ways in which the die can land so that its top shows "6". A probability measure P simply assigns a number between 0 and 1 to every subset A we might be interested in. This can be thought of as the area of A . Different probability measures P assign different areas to different sets.

For some more technical details, see Probability space.

See also:

- <https://en.wikipedia.org/wiki/Probability>

9.2.1 Randomness

CODE:ACTIVITY

Random algorithms using coins.

```
y = 0 # y is a variable, with the value zero currently
import numpy as np # this library has many useful functions
x = np.random.choice(100) # x takes values 'randomly'. It is a 'random variable'.
return x # let's see what value it takes
```

Uncertainty vs randomness: coin-flipping experiment

1. Everybody flips a coin 10 times.
2. Record each throw with 0, 1 in this spreadsheet: https://docs.google.com/spreadsheets/d/1E4bs05HnKXf1GZe4g3v6RLnHsj-YcaWg3Qe_RQyfhHU/edit?usp=sharing
3. Then record how you threw the coin and what coin it was.
4. Discuss if the coin is really random.

5. What is the distribution of coin throws for the first throw?
6. What is the distribution of recorded coin biases? Why do some coins appear more biased than others?
7. Does it make sense to aggregate all the results? What does that assume?

In the context of experiment design and data analysis, it is very common to have conditions like those in this example. Even though we wish there was such a thing as the 'repeated experiment', in practice it never is repeated. There is always some varying factor.

Pseudo-random numbers

Let us now repeat the experiment with data generated via a computer.

```
# here is a default way to generate 'random' numbers
import random
X = random.choices([0, 1], k=10) # uniformly choose 10 times between 0 and 1.
plt.hist(X) # everytime we run these commands, we get a different proportion
```

This python code is completely deterministic. A complicated calculation is used to generate the next 'random' number from the previous one. Consider this example:

```
import random
seed(5) #this sets the 'state' of the random number generating machine
print(random.uniform(0,1)) # the random number is a function of the state
print(random.uniform(0,1)) # the state changes after we generate a new number
print(random.uniform(0,1))
seed(5) # when we reset the state, we get the same sequence of numbers
print(random.uniform(0,1)) #
print(random.uniform(0,1))
print(random.uniform(0,1))
```

For cryptographically strong random numbers you need to use the secrets module:

```
import secrets
secrets.choice(range(100))
```

Physical sources of randomness

Let's go back to throwing coins now. Coins are completely deterministic. Whenever we have a specific coin to throw in the air, there are two things

we do not know. The first is which side the coin will land on. Why is that? The second is uncertainty about the coin bias: is the probability of landing heads exactly 50%? How can we quantify this? What does it depend on? Discuss in class.

What physical source of randomness can we use instead of coins?

9.2.2 Uncertainty

ACTIVITY

Probability is not only used to model random events. In fact, almost nothing can be said to be really random, unless we go into quantum physics. Even a die thrown in the air follows precise mechanical laws. Given enough information, it is possible to accurately predict the outcome of a throw.

For that reason, probability is best thought of as a way to model any residual uncertainty we have about an event. Then the probability of an event is simply a subjective measure of the likelihood.

While probability offers a nice mathematical formulation of uncertainty, when this uncertainty is subjective, the question arises: how can we elicit precise probabilities about uncertain events from individuals? Here is an example.

1. The number of immigrants ACTIVITY Consider the following question: how many immigrants live in Switzerland?

- (a) In-class discussion: what do we mean by that?
- (b) Now everybody can make a guess and record it on this form:
<https://moodle.unine.ch/mod/evoting/view.php?id=295622>

What does this distribution mean? Can we use it as an estimate of uncertainty?

- (a) Now let us create some confidence intervals. The procedure is as follows. Let us take a first guess at an interval, (say 5-10%) and ask: (a) Are you willing to take an even bet that the true number is between [5-10%]?

9.2.3 Probability space

THEORY:PROBABILITY

In probability theory, we typically define the set of all possible events that we care about as the algebra Σ , so that any possible event $A \in \Sigma$ and so that $A \subset \Omega$. The algebra has the property that it is closed under union and complement, that is:

1. If $A, B \in \Sigma$ then $A \cup B \in \Sigma$
2. If $A \in \Sigma$ then $\neg A \in \Sigma$.

Here, $\neg A \triangleq \Omega \setminus A$, i.e. the subset of Ω not containing A .

Together with a probability measure P , the tuple (Ω, Σ, P) defines a probability space. Simply put, $P(A)$ is the probability that event A happens.

9.2.4 Probability measures

THEORY:PROBABILITY

A probability measure P is a function from sets to the interval $[0, 1]$. Measuring the probability of a set is technically the same as measuring the area of a region, or the number of items in a given region. Formally, for a probability measure is defined on:

- A **universe** Ω of outcomes
- The **algebra** Σ of subsets of Ω (which we can think of as all the 'events' of interest) so that:

(a) If $A \in \Sigma$, then $A \subset \Omega$ (b) If $A, B \in \Sigma$ then $A \cup B \in \Sigma$. (b) If $A \in \Sigma$ then $\Omega \setminus A \in \Sigma$.

The axioms of probability A probability measure $P : \Sigma \rightarrow [0, 1]$ on Ω satisfies the following axioms:

1. $P(\Omega) = 1$.
2. If $A \cap B = \emptyset$ then $P(A \cup B) = P(A) + P(B)$.

From these, it also follows that $P(\emptyset) = 0$.

See also: https://en.wikipedia.org/wiki/Probability_measure

9.2.5 Marginalisation

If B_1, \dots, B_n is a partition of Ω , i.e. a collection of sets so that $B_i \cap B_j = \emptyset$ are disjoint and $\bigcup_i B_i = \Omega$ cover all of Ω , then we have that

$$P(A \cap B) = \sum_i P(A \cap B_i),$$

because $(A \cap B_i)$ are disjoint and $\bigcup_i (A \cap B_i) = A$.

9.2.6 Mutually exclusive events

DEFINITION:PROBABILITY

Two events A, B are mutually exclusive if $A \cap B = \emptyset$.

This means that there is no random outcome ω that is in both of them, so they can never be true at the same time.

9.2.7 Independence

DEFINITION:PROBABILITY

Two events A, B are said to be independent iff $P(A \cap B) = P(A)P(B)$.

Intuitively, this means that knowing if A happened tells us nothing about whether B happened.

9.2.8 Conditional probability

DEFINITION:PROBABILITY

For any events $A, B \subset \Omega$ and any probability on Ω , we define the conditional probability that A is true if B is true as follows:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

In other words, this is the probability that proportion of times that A is true when B is true. It is akin to restrict the universe of outcomes to B . We then count: how many times is A true?

We can also express independence in terms of conditional probability. A and B are independent if:

$$P(A|B) = P(A),$$

or, if $P(A) = 0$:

$$P(B|A) = P(B).$$

9.2.9 Bayes theorem

THEORY:PROBABILITY

We can use the conditional probability definition to relate $P(A|B)$ to $P(B|A)$.

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)},$$

since $P(B|A) = P(A \cap B)/P(A)$.

This is most useful for statistical inference, where we think of $P(A)$ the prior, $P(A|B)$ as the posterior, and $P(B|A)$ as the evidence likelihood.

It is also useful to write this rule when A_1, \dots, A_n is a partition of Ω . Then $P(B) = \sum_i P(B \cap A_i) = \sum_i P(B|A_i)P(A_i)$ and so:

$$P(A|B) = \frac{P(B|A)P(A)}{\sum_i P(B|A_i)P(A_i)},$$

9.2.10 Conditional independence

DEFINITION:PROBABILITY

Two events A, B are said to be independent given another event C iff

$$P(A \cap B|C) = P(A|C)P(B|C)$$

This can also be written as

$$P(A|C, B) = P(A|C)$$

9.2.11 Example Distributions

We focus on distributions where there is a finite number of possible outcomes, and hence a finite number of possible events that we might care about. All such distributions are characterised by one or more *parameters*. The simplest such distribution is a distribution on only two outcomes, the family of Bernoulli distributions.

1. The Bernoulli distribution.

DEFINITION:PROBABILITY

Let us start with a simple example, the Bernoulli distribution with parameter $\theta \in [0, 1]$. This is the distribution over two outcomes $\{0, 1\}$, so that if x is a Bernoulli random variable, then:

$$\mathbb{P}(x = 1) = \theta, \quad \mathbb{P}(x = 0) = 1 - \theta.$$

It is typical to think the distribution of heads and tails of a coin as being Bernoulli, with parameter $\theta = 1/2$.

Probability space Formally, if the underlying probability space is (Ω, Σ, P) , with random outcomes $\omega \in \Omega$ and random variable $x : \Omega \rightarrow \{0, 1\}$ then

$$\mathbb{P}(x = 1) = P(\{\omega : x(\omega) = 1\}).$$

See also: https://en.wikipedia.org/wiki/Bernoulli_distribution

2. Binomial distribution

DEFINITION:PROBABILITY

If we repeat a Bernoulli trial, we can also count the number of times the coin comes heads. The distribution of the counts is called the Binomial distribution. If y is a binomial random variable for n throws with parameter θ , then we can write it as the sum of n Bernoulli random variables x_1, \dots, x_n , i.e.:

$$y = \sum_{t=1}^n x_t.$$

The probability of k heads after n throws is given by the formula:

$$\mathbb{P}(y = k) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}$$

, where $\binom{n}{k}$ is the binomial coefficient.

See also: https://en.wikipedia.org/wiki/Binomial_distribution

3. The Categorical/Multinomial distribution

A multinomial distribution is an extension of the Bernoulli and binomial distributions to $m \geq 2$ outcomes.

Categorical distributions Let us start with one trial, e.g. a single throw of a die. We can model this dice throw as the distribution where the probability that the die lands with its k -th face on top is θ_k ,

$$\mathbb{P}(x = k) = \theta_k.$$

Thus, this distribution is parametrised by the vector $\theta = (\theta_1, \dots, \theta_m)$. A random variable $x : \Omega \rightarrow \{1, \dots, m\}$ obeying this distribution is called multinomial.

If the underlying probability space is (Ω, Σ, P) , then

$$\mathbb{P}(x = k) = P(\{\omega : x(\omega) = k\})$$

See also: https://en.wikipedia.org/wiki/Multinomial_distribution

4. Uniform distributions

A special case of binomial and multinomial distributions is the uniform distribution. This is defined as follows.

Let $|A|$ be the size of a set A . Then a distribution P is uniform if it obeys

$$P(A) = \frac{|A|}{|\Omega|}.$$

This definition applies to continuous distributions as well. A standard example is the uniform distribution on the interval $[0, 1)$. Then the probability that we obtain an outcome in the set $[0, p)$ is always equal to p , i.e.

$$P([0, p)) = \mathbb{P}(\omega \in [0, p)) = p.$$

9.2.12 Random variables

THEORY:PROBABILITY

A real-valued random variable $f : \Omega \rightarrow \mathbb{R}$ is simply a function from the outcomes to the real numbers. Even though it is a fixed function, its values are random, because the actual value $\omega \in \Omega$ that will be used to calculate its value $f(\omega)$ is random.

Random variables can be easily generalised to other domains than the real numbers.

See also: https://en.wikipedia.org/wiki/Random_variable En français: https://fr.wikipedia.org/wiki/Variable_al%C3%A9atoire

1. Example random variable

EXAMPLE

Take a 10-sided die. The outcomes of the die represent the space Ω . We can now create a Bernoulli variable from the die. For example, set $x(\omega) = 1$ when $\omega \in \{1, 5\}$ and $x(\omega) = 0$ when $\omega \in \{6, 10\}$.

What is the distribution of x ? Have you seen it before?

2. Generating a Bernoulli random variable.

ACTIVITY:CODE

In python, there are procedures for generating data from many types of random variables. However, all such methods for generating random numbers are not truly random. They rely on something called a pseudorandom number generator. The values output by this generator are then *transformed* so as to become similar to the random variable we want.

- (a) Let us start by throwing a 10-sided die. How do you generate a Bernoulli random variable with $\theta = 0.6$ from the outcomes of the dice throw?
- (b) Now let us consider the following python code, which generates values from a Bernoulli random variable.

```
import numpy as np
return np.random.choice(2,p=[0.6, 0.4])
```

Let us replicate it through a distribution of uniform random variables in $[0, 1]$.

```
import numpy as np
x = np.random.uniform() # returns a uniformly generated number in [0,1).
if x < 0.6:
```

```

    return 1
else:
    return 0

```

9.2.13 Expectation

DEFINITION:PROBABILITY

The expectation of a random variable $f : \Omega \rightarrow \mathbb{R}$ with respect to a probability P on Ω is defined as follows, for finite Ω .

$$\mathbb{E}_P(f) = \sum_{\omega \in \Omega} f(\omega)P(\omega).$$

When P is clear from the context, we can just write $\mathbb{E}(f)$.

1. Expectation of general random variables AD-
VANCED:PROBABILITY:THEORY For the general case, we define it in terms of integrals:

$$\mathbb{E}_P(f) = \int_{\Omega} f(\omega) dP(\omega).$$

2. Averages and expected values

If we have obtained a sequence of samples $\omega_1, \dots, \omega_n$ from $P(\omega)$, and calculated the average

$$M_n(\omega_1, \dots, \omega_n) = \frac{1}{n} \sum_{i=1}^n f(\omega_i).$$

then the average is approximately equal to the expectation:

$$\mathbb{E}_P(f) \approx \frac{1}{n} \sum_{i=1}^n f(\omega_i).$$

To see this, first note that $\mathbb{E}_P(f(\omega_i)) = \mathbb{E}_P(f)$, so the expected value of the average is equal to the expected value of f .

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n f(\omega_i)\right] = \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n f(\omega_i)\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(\omega_i)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}(f) = \mathbb{E}(f).$$

In fact, the average M_n is another random variable, defined on Ω^n . We say that M_n **concentrates** around the expected value, with the property that

$$M_n = \mathbb{E}(f) \mp O(1/\sqrt{n}).$$

Informally, we should expect our error to reduce proportionally to $1/\sqrt{n}$ the number of samples we have.

9.2.14 Variance

DEFINITION:PROBABILITY

The variance of a random variable f is defined as

$$\sigma_f^2 \triangleq \mathbb{E}\{[f - \mathbb{E}(f)]^2\}.$$

We can rewrite the variance as follows:

$$\sigma_f^2 = \mathbb{E}\{[f - \mathbb{E}(f)]^2\} = \mathbb{E}[f^2 - 2f \mathbb{E}(f) + \mathbb{E}(f)^2] = \mathbb{E}[f^2] - 2\mathbb{E}(f) \mathbb{E}(f) + \mathbb{E}(f)^2 = \mathbb{E}[f^2] - \mathbb{E}(f)^2.$$

The variance is another example of a moment of a random variable. More complicated moments, such as skew and kurtosis exist.

9.2.15 Conditional expectation

DEFINITION:PROBABILITY

We can define the conditional expectation of a random variable through conditional probabilities.

In particular, let some variable $f : \Omega \rightarrow \mathbb{R}$ and some event $B \subset \Omega$. Then we define

$$\mathbb{E}(f|B) = \sum_{\omega \in \Omega} f(\omega)P(\omega|B).$$

We can even re-write this in a more intuitive way. Note that if $\omega \notin B$ then $P(\omega|B) = 0$. So, we have

$$\mathbb{E}(f|B) = \sum_{\omega \in B} f(\omega)P(\omega|B) = \sum_{\omega \in B} f(\omega)P(\omega)/P(B).$$

1. Group comparison

ACTIVITY

Let us say we wish to compare the distribution among multiple groups. Perhaps we wish to compare the number of students who achieve a certain test like the table below.

	Success	Success
School	Male	Female
A	62%	82%
B	63%	68%
C	37%	34%
D	33%	35%
E	28%	24%
F	6%	7%
Average	45%	38%

- (a) Let us plot the success rate for females and males over different schools.
- (b) Does this show a bias? What information is missing?
- (c) Let us combine these two plots into one plot now for this we need to use the following code:

```
w=0.5
X = np.linspace(1,6*w,6)
M = ...
F = ...
plt.bar(X, M, align=edge, width=w/2)
plt.bar(X + w, F, align=edge, width=w/2)
```

10 Graphics types

1. Histogram and 3D extensions
2. Density curve
3. Scatterplot
4. Smooth scatterplot
5. Violin plot
6. Line Plot
7. Confidence Intervals
8. Geographical/topological maps
9. Network graphs
10. Word cloud

See also: Catalogue of data visualisation

11 Schedule and links to other courses

The schedule of this and the other courses is in flux, but I do not expect it to change very much. In any case, the course will operate independently of the other courses. You should expect to cover the same topic more than once. However, this course will not focus on either statistical theory or programming.

Week	Statistics	Programming	In-Course	Homework
1 23 Sep	Course intro	Python intro	Histograms Randomness	Math score
2 30 Sep	R Intro Data manipulation Histograms Scatterplots Boxplots Variable types Mosaic plots Functions	Data types	Uncertainty Discrete Variables Continuous Variables	Form groups
3 7 Oct	Quantifying Variability Distribution Density function Histograms Skewness Quantiles	Control	Time-Series Linear functions Stock market prices Crime statistics S&P index World Records	Form groups
4 14 Oct	Qualitative vars in R Discrete vars in R	Structures	Scatterplots Unemployment	Proj. Proposal
5 21 Oct	Continous RV	Functions	Data Cleaning	Table2Picture
6 28 Oct	Continuous RV	Complements	Experiment design Random Sampling Undercounting Represnetative samples	Deconstruction
7 4 Nov	Continuous RV	Classes	Expectations	NewsPaper
8 11 Nov	Dependencies. Joint distribution. Conditional distribution.	Objects	Bayesian Inference	Project Highlight
9 18 Nov	Moments	Errors	Pipelines Simulation studies	Project
10 25 Nov	Covariance Correlation Scatterplots	Iterators	Hypothesis testing The garden of many paths	Project
11 2 Dec	Prices, returns	FP	Examples, project work	Project presentation
12 9 Dec	Conditional expectations		Examples, project work	Project report
13 16 Dec			Project presentations	Project report

12 Glossary

- Area: Air
- Die (Dice): Dé(s)
- Expectation: Espérance
- Experiment: événement
- Histogram: Histogramme?
- Pie chart: Diagramme circulaire
- Bar chart: Diagramme à barres
- Scatterplot: Nuage de points
- Randomness: Hasard
- Uncertainty: Incertitude
- Probability: Probabilité
- Stochastic: stochastique
- Random: Aléatoire
- Random Variable: Variable aléatoire
- Sample Space: see Universe
- Sample (v): Échantillonner
- Sample (n): Échantillon
- Sampling: Échantillonnage
- Set: Ensemble
- Subset: Sous-ensemble
- Superset: Sur-ensemble
- Survey: Sondage
- Universe: Univers
- σ -algebra: tribu, σ -algèbre

13 References

Python help: Use python's `!help!` function whenever you can.

```
help(print)
```

13.1 Some workshops

- <https://vishub.net/bach>
- <https://www.dagstuhl.de/en/program/calendar/semhp/?semnr=22331>
- <https://www.dagstuhl.de/en/program/calendar/semhp/?semnr=22351>

13.2 Books

- [CA] The Truthful Art. Cairo, Alberto :book:
- [GJ] Data Science Par La Pratique (Data Science from Scratch). Grus, Joel (Ch3: Visualisation, Ch6: Probability, Ch7: Inference, Ch9: Data collection, Ch10: Exploration, Ch14: Regression, Ch15: Regression+Bootstrap)