

Welcome to BIO-210

Applied software engineering for life sciences

November 10th 2025 – Lecture 9

Prof. Alexander MATHIS

EPFL

Announcements

- Today (Monday at 10am) v4 on testing was due, we will grade it. If you haven't released it yet, please do so asap.
- Monday 16:15 - 17: my office hours at SV 2811

	Date	Topic	Software version	Software releases	Feedback
0	08/09/2025	Python introduction I			
1	15/09/2025	Python introduction II			
2	22/09/2025	Public holiday			
3	29/09/2025	Git & GitHub (+ installation)			
4	06/10/2025	Project introduction	v1		
5	13/10/2025	Functionify	v2	v1	
6	20/10/2025	EPFL fall break			
7	27/10/2025	Visualization & documentation	v3	v2	code review (API)
8	03/11/2025	Unit-tests, functional tests	v4	v3	
9	10/11/2025	Code refactoring	v5	v4	graded (tests)
10	17/11/2025	Profiling & code optimization	v6	v5	code review
11	24/11/2025	Object oriented programming	v7	v6	graded (speed)
12	01/12/2025	Model analysis & project report	v8	v7	code review (OO)
13	08/12/2025	Work on project (no class)			
14	15/12/2025	Wrap up		v8	graded (project)

🧩🤔: Is there something wrong with this doctest?

```
1 def get_dna_complement(seq):
2     """
3     Return the complement of a DNA sequence.
4
5     Example
6     -----
7     >>> get_dna_complement('ATCG')
8     TAGC
9     """
10    complement = {'A': 'T', 'T': 'A', 'C': 'G', 'G': 'C'}
11    return ''.join(complement[base] for base in seq)
```

🧩🤔 : Yes, missing quotes!

```
1 def get_dna_complement(seq):
2     """
3     Return the complement of a DNA sequence (adenine (A) pairs with thymine (T),
4     and cytosine (C) pairs with guanine (G)).
5
6     Example
7     -----
8     >>> get_dna_complement('ATCG')
9     'TAGC'                                     # Missing quotes. It expects a string!
10    """
11    complement = {'A': 'T', 'T': 'A', 'C': 'G', 'G': 'C'}
12    return ''.join(complement[base] for base in seq)
```

🧩🤔 : How do we test floating-point results?

```
1 import numpy as np
2
3 def calculate_std(data):
4     """Calculate standard deviation of data."""
5     return np.std(data)
```

Test that `calculate_std([1, 2, 3, 4, 5])` returns approximately `1.4142135623730951`.

Option 1: Doctest with rounding

```
1 def calculate_std(data):
2     """
3     Calculate standard deviation.
4
5     >>> round(calculate_std([1, 2, 3, 4, 5]), 2)
6     1.41
7     """
8     return np.std(data)
```

Option 2: Pytest with approx

```
1 import pytest
2
3 def test_calculate_std():
4     """Test standard deviation calculation"""
5     result = calculate_std([1, 2, 3, 4, 5])
6     assert result == pytest.approx(1.414, rel=1e-3)
7
8     # Or with absolute tolerance
9     assert result == pytest.approx(1.414, abs=0.001)
```

Pytest's `approx()` is cleaner for floating-point comparisons!

Summary: Doctest vs Pytest

Feature	Doctest	Pytest
Purpose	Documentation + basic testing	Comprehensive testing
Location	Inside docstrings	Separate test files
Best for	Simple examples, user-facing docs	Edge cases, complex scenarios
Syntax	<pre>>>> function(input)</pre>	<pre>assert function(input) == expected</pre>
Floating-point	Use <code>round()</code>	Use <code>pytest.approx()</code>
Exceptions	<pre>Traceback...</pre>	<pre>pytest.raises()</pre>
Multiple cases	Repeat examples	<code>@pytest.mark.parametrize</code> (see later!)

Your README.md file

A README is a text file that introduces and explains your GitHub project.

- should be succinct, but detailed
- notice that it's the face of your project; i.e. if one goes to your repository, the README is rendered below the file list.

Typical content:

- Project title + description
- Installation guide/requirements (see below)
- Examples (for using the code)
- License

A brief overview of licenses

When you create a project, you should also choose a license for your code. A license tells others what they can and cannot do with your code. Unlicensed code is considered "all rights reserved" by default, meaning others cannot use, modify, or distribute it without your permission.

Some common open-source licenses:

- MIT License: very permissive, allows reuse with minimal restrictions
- GNU General Public License (GPL): requires derivative works to also be open-source under the same license
- Apache License 2.0: allows use, modification, and distribution with some conditions regarding patents and trademarks
- Check out this page for comparing licenses

Naturally for your (private) BIO-210 project, the license is not a big issue, but for real-world projects, it's important to choose an appropriate license to protect your work and clarify how others can use it.

Examples of good READMEs

Check out some example projects. Firstly, typical research codebases:

- two recent machine learning projects from my group: Kinesis and hBehaveMAE






Secondly, popular open-source projects with well-structured READMEs:

- Scipy
- MMore
- Awesome Python
- Python The Algorithms
- Strix

Also check out the Demo project.

How to format a Readme.md?

Formatting READMEs is based on the Markdown language, which enables you to write nicely formatted and visually appealing texts!

- it's simple to learn, just skim this Markdown Guide      ...
- Check out this readme generator
- You can also use Markdown live preview to see how your markdown looks like while writing it.

Git management: .gitignore file

Use a **.gitignore** file to exclude files you do not want to version control (e.g. notebook checkpoints, pycache, files from your IDE, ...). This keeps your repository clean and avoids committing unnecessary files. The .gitignore file should be placed in the root directory of your repository.

- [Docs for gitignore](#)
- [demo-project/.gitignore.](#)

Note: if you already have unwanted files in your repository and create a .gitignore, you need to remove those, as they are already tracked with git.

How to make code reproducible?

- make (software) dependencies and requirements explicit (next slide)
- use version control and store which version produced what results
- use tests / doctests
- do not comment/uncomment sections of your code to control the behavior. This approach makes it error prone, hard to reproduce and difficult to automate. Instead, use if/else statements to control the flow of your program. If needed, use configuration files or command-line arguments to set parameters.
- Leverage scripts to run experiments, which can help in automating and reproducing results.
- You could also consider making different experiment scripts (to reproduce each experiment)

Documenting requirements: Software dependencies

- Software often has many dependencies.
 - Code might run differently with different library versions!
 - Thus, placing dependencies within a contained environment can minimize issues and allow others to run the code just like it runs on your system. It also means that you can run the same code in the future, even if libraries have changed.
 - Different programs might also require different libraries, which can be supported by `environments` . This avoids dependency conflicts between different projects on your system.
 - Common python environments include MiniForge, Anaconda (conda) and virtualenv. You are already familiar with MiniForge from the course setup
- > share an environment/requirement file (e.g. see demo-project)
- > also see hBehaveMAE requirements

Discussion: How should your code be organized?

What components should your project have? How should files and folders be organized?

How should your code be organized?

Code and Basics:

- README (see above, also demo code!) + .gitignore file
- a main function that is clean and easy to run (additionally `experimentX.py` , `exptY.py` , ...)
- a test function with all our doctests/pytests
- scripts containing functions (named reasonably, e.g. `TuringModel.py`) and a `utils.py` script containing plotting etc.

Results:

- a folder (e.g. called 'results') containing saved experiments/outcomes (could have traceable names: e.g. `experiment17_parametersXYZ.png` are the results for `experiment17.py`)
- Remark 1: you could also add jupyter notebooks that integrate documentation (description + figures)
- Remark 2: for BIO-210 just follow the problem sets (to get the highest score!); e.g. `main.py` in v2 should do whatever you're asked to have implemented.

Further reading

- Good enough practices in scientific computing - a very short article!
- Guide for reproducible research by the Alan Turing Institute

Questions?

: How can we make this faster?

```
1 def generate_patterns(num_patterns, pattern_size):
2     patterns = np.zeros([num_patterns, pattern_size])
3     for i in range(num_patterns):
4         for j in range(pattern_size):
5             patterns[i,j]=np.random.choice([-1,1])
6
7     return patterns
```

: Are these good docstrings?

```
1 def generate_pattern(pattern_number,pattern_size):
2     '''
3     Generates an array of (pattern_number) patterns of size (pattern_size)
4
5     Parameters
6     -----
7     pattern_number : int
8         Number of patterns to generate.
9     pattern_size : int
10        Size (number of neurons) in each pattern.
11
12    Returns
13    -----
14    np.ndarray
15        Array of shape (pattern_number, pattern_size) with entries of -1 or 1.
16
17    Examples
18    -----
19    >>> generate_pattern(3,4).shape
20    (3,4)
21    '''
22    return np.random.choice([-1,1],size=(pattern_number,pattern_size))
```

Questions?

Selected function topics

Reminder: Function-related statements and expressions

Statement or expression	Examples
Call expression	<code>myfunc('Sepp1',175,age=22,*rest)</code>
<code>def</code>	<pre>def printer(message): print('Hello'+message)</pre>
<code>return</code>	<pre>def adder(a,b=1,*c): return a+b+c[0]</pre>
<code>global</code>	<pre>x = 'outside' def changer(): global x; x= 'new'</pre>
<code>lambda</code>	<code>func = [lambda x: x**2, lambda x: x**3]</code>

Reminder: Argument-matching modes

Python functions allow highly flexible calling patterns for functions

- *Positionals*: matched left to right
- *Keywords*: matched by argument name; `name = value` syntax
- *Defaults*: specify values for optional arguments (that do not need to be passed)
- ***Varargs collecting*: pass arbitrarily many positional or keyword arguments**

Arbitrary positional arguments collectors (*args)

When this function is called all arguments are assigned to a tuple *args*.

```
1 In [1]: def f(*args): print(args) # Note use of *
2
3 In [2]: f()
4 () # returns a tuple!
5
6 In [3]: f(1)
7 (1,) # returns a tuple!
8
9 In [4]: f(1,2,3,4)
10 (1, 2, 3, 4) # returns a tuple!
11
12 In [5]: f("EPFL","is","great!")
13 ('EPFL', 'is', 'great!') # returns a tuple!
```

Why is *args useful?

For instance, imagine you need to sum all numbers somebody gives you ...

```
1 In [1]: def summation(a,b,c):
2         ...:     return a+b+c
3         ...: summation(1,2,3)
4 Out[1]: 6
5 In [2]: summation(1,2,3,4,5)           # That's too many for your function!
6 -----
7 TypeError                                 Traceback (most recent call last)
8 <ipython-input-2-0851c3d51d8a> in <module>
9 ----> 1 summation(1,2,3,4,5)
10 TypeError: summation() takes 3 positional arguments but 5 were given
11 In [3]: def summation(*args):
12         ...:     Sigma=0
13         ...:     for s in args:
14         ...:         Sigma+=s
15         ...:     return Sigma
16         ...:
17 In [4]: summation(1,2,3,4,5)
18 Out[4]: 15
19 In [5]: summation(213.1,445560123,111)
20 Out[5]: 445560447.1
```

Arbitrary keyword arguments collectors (**kargs)

When this function is called all arguments are assigned to a dictionary *args*.

```
1 In [1]: def f(**kargs): print(kargs)      # Note use of **
2
3 In [2]: f()
4 {}                                         # returns a dict
5
6 In [3]: f(1)
7 -----
8 TypeError                                 Traceback (most recent call last)
9 <ipython-input-3-281ab0a37d7d> in <module>
10 ----> 1 f(1)
11
12 TypeError: f() takes 0 positional arguments but 1 was given
13
14 In [4]: f(x=1,y=2)
15 {'x': 1, 'y': 2}                         # returns a dict
```

Argument matching in Python function

Now we have seen all four styles: *Positionals*, *Keywords*, *Defaults*, and *Varargs collecting*

These methods can also be combined, but note positional arguments come before keyword arguments. Just like default variables come after positional variables.

```
1 def myfun(a,b,*args,**kwargs):  
2     pass
```

Read more in the [docs](#)

: How can we calculate the mean of each column ignoring missing data?

```
1 In [1]: data = np.array([[2,np.nan],[3, 0],[4,1]]) # Notice: missing data
```

: How can we calculate the mean of each column ignoring missing data?

```
1 In [1]: data = np.array([[2,np.nan],[3, 0],[4,1]])      # Notice: missing data
2 In [2]: np.nanmean(data,axis=0)
3 Out[2]: array([3.0,  0.5  ])
```

Let's verify: $(2+3+4)/3 = 3$ and $(0+1)/2 = 0.5$. The denominator is automatically adjusted to the number of non-nan items!

Note:

```
1 In [3]: data.mean(axis=0)
2 Out[3]: array([3.,   nan])
```

Entries of `np.nan` is very useful if you have missing data in some table (dataframe/array). The `np.nanX` functions then make sure that missing data are handled correctly.

Anonymous functions: lambda

Python has an expression `lambda` to generate function objects (name from lambda calculus).

```
1 lambda arg1, arg2, ... argN: expression using args
```

Functions returned by running `lambda` expressions are like those created and assigned by `def`.

```
1 In [1]: def f(x,y,z): return x*y*z
2 In [2]: f(1,2,3)
3 Out[2]: 6
4 In [3]: f=lambda x,y,z: x*y*z      # explicitly assign
5 In [4]: f(1,2,3)
6 Out[4]: 6
7 In [5]: f=lambda x,y=2,z=2: x*y*z  # defaults work similarly
8 In [6]: f(1)
9 Out[6]: 4
```

Differences of `lambda` and `def`

- `lambda` is an expression not a statement. Thus, `lambda` can appear in different places!
- `lambda`'s body is a single expression not a block of statements.

Expression vs. Statement

Aspect	Expression	Statement
Returns	Evaluates to a value	Performs an action
Can assign?	✓ Yes: <code>x = expression</code>	✗ No: <code>x = statement</code>
Use inline?	✓ Yes: <code>print(expression)</code>	✗ No: <code>print(statement)</code>
Examples	<code>5 + 3</code> , <code>lambda x: x**2</code> , <code>x > 5</code>	<code>x = 5</code> , <code>def f():</code> , <code>if x > 5:</code>

Key insight: `lambda` is powerful because it's an *expression* - unlike `def` you can use it anywhere a value is needed! This is particularly helpful with other tools...

Mapping functions over iterables: map

```
1 In [1]: data=[0,123,1224,412.23]
2 In [2]: map?
3 Init signature: map(self, /, *args, **kwargs)
4 Docstring:
5 map(func, *iterables) --> map object
6
7 Make an iterator that computes the function using arguments from
8 each of the iterables. Stops when the shortest iterable is exhausted.
9 Type:          type
10 Subclasses:
11 In [3]: map(lambda x: x**17-13.2,data)          # it is an iterator!
12 Out[3]: <map at 0x7fa1515be9a0>
13 In [4]: list(map(lambda x: x**17-13.2,data))    # convert to list for displaying
14 Out[4]:
15 [-13.2, 3.3758791744665375e+35, 3.1065911647383047e+52, 2.866639650957827e+44]
16 # Remember pow(base,exp) = base ** exp
17 In [5]: list(map(pow,[0,1,2,3,4],[2,2,2,3]))    # Notice *args in use!
18 Out[5]: [0, 1, 4, 27]
```

Selecting items in iterables: filter

`filter` allows the selection of iterable's items based on a function.

```
1 In [1]: list(range(-10,10))
2 Out[1]: [-10, -9, -8, -7, -6, -5, -4, -3, -2, -1, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
3
4 In [2]: list(filter((lambda x: x>3.4), range(-10,10)))
5 Out[2]: [4, 5, 6, 7, 8, 9]
```

Combining items in iterables: reduce

`reduce` allows cumulatively applying a function!

```
1 In [1]: from functools import reduce
2
3 In [2]: reduce(lambda x,y: x+y, [1,2,3,4,5])
4 Out[2]: 15
5
6 In [3]: reduce(lambda x,y: x*y, [1,2,3,4,5])
7 Out[3]: 120
8
9 In [4]: reduce?
10 Docstring:
11 reduce(function, sequence[, initial]) -> value
12
13 Apply a function of two arguments cumulatively to the items of a sequence,
14 from left to right, so as to reduce the sequence to a single value.
15 For example, reduce(lambda x, y: x+y, [1, 2, 3, 4, 5]) calculates
16 (((1+2)+3)+4)+5). If initial is present, it is placed before the items
17 of the sequence in the calculation, and serves as a default when the
18 sequence is empty.
19 Type:          builtin_function_or_method
```

: What does this reduce return?

```
1 from functools import reduce
2 numbers = [2, 3, 2]
3
4 result = reduce(lambda x, y: x ** y, numbers)
5 print(result)
```

Let's trace it step by step:

1. $x=2, y=3 \rightarrow 2 ** 3 = 8$
2. $x=8, y=2 \rightarrow 8 ** 2 = 64$

Thus, the answer is: `64`

What is a Decorator?

A **decorator** is a function that modifies or enhances another function without changing its code.

Decorators wrap a function to add extra behavior. We'll see an example next!

Parametrized Testing with Pytest

`@pytest.mark.parametrize` is a decorator. It allows you to run the same test function multiple times with different inputs.

Benefits:

- Avoid repetitive code (no need to write `assert` for each case)
- Each parameter set runs as a separate test
- Easy to add new test cases
- Clear reporting of which specific inputs fail

Generic syntax:

```
1 @pytest.mark.parametrize("argument_names", [  
2     (value1_arg1, value1_arg2, ...),  
3     (value2_arg1, value2_arg2, ...),  
4 ])  
5 def test_function(argument_names):  
6     # test code using the arguments
```

Pytest with multiple test cases (parametrize)

```
1 import pytest
2 import numpy as np
3
4 @pytest.mark.parametrize("data,expected", [
5     ([1, 2, 3, 4, 5], 1.4142135623730951),
6     ([10, 10, 10, 10], 0.0),
7     ([0, 0, 0, 0], 0.0),
8     ([1, 5], 2.0),
9     ([-1, 1], 1.0),
10    ([100, 200, 300], 81.64965809277260),
11 ])
12 def test_calculate_std(data, expected):
13     """Test standard deviation calculation with multiple inputs"""
14     result = calculate_std(data)
15     assert result == pytest.approx(expected, rel=1e-6)
```

Imperative vs. functional programming

```
1 In [1]: input = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
2     ...: output = 0
3     ...: for i in range(len(input)):
4     ...:     if input[i] % 2 == 0:
5     ...:         output += input[i]*10
6     ...:
7
8 In [2]: output
9 Out[2]: 300
```

Imperative programming example:

```
In [3]: reduce(lambda x,y: x+y, map(lambda i: 10*i, \
                                     filter(lambda i: i%2 == 0, range(1,11))))
Out[3]: 300

In [4]: np.sum(np.arange(0,11,2)*10)      # Numpy implementation
Out[4]: 300
```

For some theory, check out: [Can programming be liberated from the von Neumann style?: a functional style and its algebra of programs](#) by John Backus.

Scipy library: scientific computing

Important tools in the scipy library:

- numerical integration `scipy.integrate`
- optimization `scipy.optimize`
- Fourier transforms `scipy.fft`

For a full list, see the [docs](#) and the [cookbook](#).

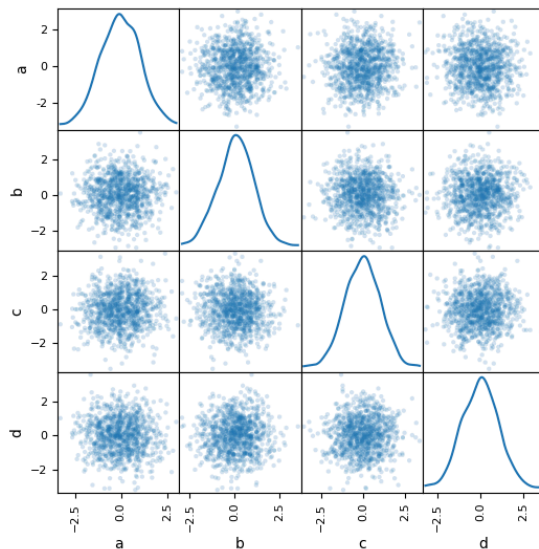
Pandas: data structures and analysis

High-level data analysis package

```
1 >>> import pandas as pd
2 >>> import numpy as np
3 # Defining a dataframe (df):
4 >>> df = pd.DataFrame(np.random.randn(1000, 4), columns=["a", "b", "c", "d"])
5 >>> print(df.head())           # prints the first 5 rows
6 a          b          c          d
7 0 -0.182791 -0.768629 -1.381591  0.035229
8 1  1.210803  0.487254 -0.087025  0.253478
9 2 -0.371740  1.092439 -0.829110  0.518891
10 3 -1.364988 -2.046488  0.172973 -2.117577
11 4  1.369287 -0.413473 -2.047923 -1.240338
12 >>> print(df.mean(axis=0))
13 a    -0.016737
14 b     0.000806
15 c    -0.049374
16 d    -0.052093
17 dtype: float64
```

Pandas: data structures and analysis

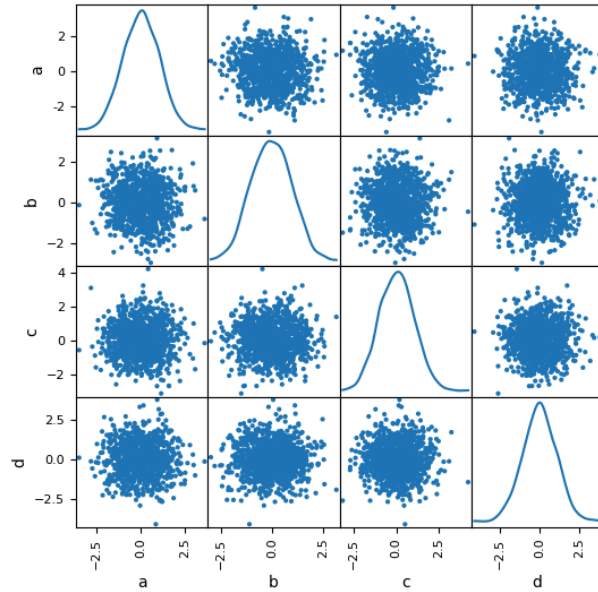
```
>>> from pandas.plotting import scatter_matrix
>>> scatter_matrix(df, alpha=0.2, figsize=(6, 6), diagonal="kde");
```




Note: Matplotlib (and thus pandas) allows you to regulate the transparency of a graph plot using the alpha attribute. By default, alpha=1 (not transparent).

Same code without transparency

```
>>> scatter_matrix(df, alpha=1, figsize=(6, 6), diagonal="kde");
```



Scikit learn: Machine learning in python

 [Install](#) [User Guide](#) [API](#) [Examples](#) [Community](#) [More](#) 🔍 📄 🌐 1.7.2 (stable) ▾

scikit-learn
Machine Learning in Python

[Getting Started](#) [Release Highlights for 1.7](#)

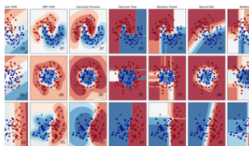
- Simple and efficient tools for predictive data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying which category an object belongs to.

Applications: Spam detection, image recognition.

Algorithms: [Gradient boosting](#), [nearest neighbors](#), [random forest](#), [logistic regression](#), and [more...](#)



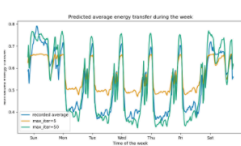
Examples

Regression

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, stock prices.

Algorithms: [Gradient boosting](#), [nearest neighbors](#), [random forest](#), [ridge](#), and [more...](#)



Examples

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, grouping experiment outcomes.

Algorithms: [k-Means](#), [HDBSCAN](#), [hierarchical clustering](#), and [more...](#)



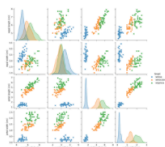
Examples

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, increased efficiency.

Algorithms: [PCA](#), [feature selection](#), [non-negative matrix factorization](#), and [more...](#)

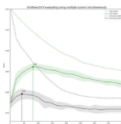


Model selection

Comparing, validating and choosing parameters and models.

Applications: Improved accuracy via parameter tuning.

Algorithms: [Grid search](#), [cross validation](#), [metrics](#), and [more...](#)

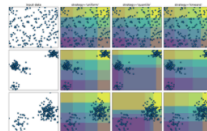


Preprocessing

Feature extraction and normalization.

Applications: Transforming input data such as text for use with machine learning algorithms.

Algorithms: [Preprocessing](#), [feature extraction](#), and [more...](#)



Questions?

Today's summary

- .gitignore, readme.md, markup, reproducible environments, project structure
- functional tools: `lambda` , `map` , `reduce` , `filter`
- important libraries: `scipy`, `pandas`, and `scikit-learn`

Try out the commands in the python shell/notebooks! Practice is key.

After lunch:

- Monday 13 - 15: exercises working on v5 of your project
- Monday 16:15 - 17: my office hours at SV 2811