

Contents

Introduction	1
Code-Automation as a Standard in Modern Biosciences	2
How Code Quality Improves Scientific Reproducibility	3
Python as a Programming Language	5
The Potential of Python Data Science Packages for Biomedicine	9
Aims	13
References	14
Appendices	17
A Supplementary Data & Methods	17
A.1 Figures	17
A.2 Tables	18
A.3 Materials & Methods	19
B Documentation of plotastic	20

Introduction

To provide a comprehensive background for the following chapters that focus on the interaction of human mesenchymal stromal cells (hMSCs) with multiple myeloma (MM) cells, this

Code-Automation as a Standard in Modern Biosciences

Beschreibe die Situation. - Big Data in Biosciences - what is big data, examples - Define citable challenges: - reproducibility crisis - lack of tools

In recent years, the biosciences have evolved dramatically, with a marked increase in the volume and complexity of data generated (Yang et al., 2017; Ekmekci et al., 2016). This transformation necessitates robust software tools, many of which require coding skills to use effectively. Here we summarize standard tools used by biosciences today and show their reliance on coding. The author argues that the role of a modern independent researcher is now intertwined with coding skills similar to a role of “precision medicine bioninformatician” (Gómez-López et al., 2019).

Statistical analysis in biosciences has traditionally been reliant on user-friendly tools like *Microsoft Excel* and *GraphPad Prism*. While *Excel* by itself is recognized as limited for complex data analysis (Tanavalee et al., 2016; Incerti et al., 2019), GraphPad Prism offers more advanced statistical models .

However, increasingly demands more sophisticated approaches as data sets grow in size and complexity.

R and Python scripts offer more efficient and versatile solutions, enabling complex analyses with a few lines of code (R Core Team, 2018; Vallat, 2018).

Recognizing this trend, Microsoft has integrated a Python interpreter into *Excel* to computations more accessible within a widely used platform (?).

Next-generation sequencing, such as bulk RNAseq, has become affordable, allowing for larger sample sets during a single PhD project. This technology offers advanced tools that are most efficiently used through scripting in R or Python. In the absence of a dedicated statistician, researchers are compelled to learn coding.

In gene ontology, tools such as Metascape facilitate the integration of vast datasets and outputs multiple useful data visualizations. Metascape also provides multiple excel sheets, containing all results, sometimes in a nested format, which provides even further information that’s adaptable for specific hypotheses, but given the sheer amount of data, is impractical to analyze manually.

since Metascape returns large *Excel* sheets with complex nested information, a researcher without coding skills requires manual work to adapt the results to specific research hypotheses.

its true potential is unlocked only when researchers can manipulate and analyze these data through scripting.

Modern gene ontology tools like Metascope offer powerful graphical user interfaces. However, their effectiveness is only possible through standardizing multiple large datasets.

The output from Metascope, large *Excel* sheets with complex nested information, is more efficiently analyzed through scripting, which is often necessary to adapt metascope results to specific research hypotheses.

Image analysis is another area where coding skills are essential. ImageJ/FIJI, a standard tool in the field, requires scripting for batch processing of multiple images and automating multiple processing steps into a pipeline. While macros can be recorded, understanding the underlying code is necessary for troubleshooting and adapting the macro to new datasets.

In the field of protein structural biology, Pymol is a standard tool that also has a Python command interface.

Similarly, artificial intelligence (AI), a game-changer in biomedicine, primarily uses Python due to its extensive libraries for machine learning and scientific computing. Python is also a standard for integrative biomedicine simulations.

Finally, databases and repositories are essential for storing, retrieving, and sharing data. Researchers need to understand common file formats to adhere to standards that ensure re-usability and interoperability. Scripting helps automate the process of formatting data for submission to these databases.

In conclusion, the integration of coding in bioscience research is not just a trend but a necessity. As the field continues to evolve, the demarcation between biologists and computational scientists blurs, underscoring the importance of coding skills for the next generation of researchers. The ability to code is fast becoming an indispensable asset, as integral to bioscience as traditional laboratory skills.

How Code Quality Improves Scientific Reproducibility

A main reason to write software is to define re-usable instructions for task automation (Narzt et al., 1998). However, the complexity of the code makes it prone to errors and can prevent usage by persons other than the author himself. This is a problem for the general scientific community, as the software is often the only way to reproduce the results of a study (Sandve et al., 2013). Hence, modern journals aim to enforce standards to software development, including software written and used by biological researchers (Smith et al., 2018). Here, we provide a brief overview of the standards utilized by plotastic that to ensure its reliability and reproducibility by the scientific community (Peng, 2011).

Modern software development is a long-term commitment of maintaining and improving

code after initial release (Boswell & Foucher, 2011). Hence, it is good practice to write the software such that it is scalable, maintainable and usable. Scalability or, to be precise, structural scalability means that the software can easily be expanded with new features without major modifications to its architecture (Bondi, 2000). This is achieved by writing the software in a modular fashion, where each module is responsible for a single function. Maintainability means that the software can easily be fixed from bugs and adapted to new requirements (Kazman et al., 2020). This is achieved by writing the code in a clear and readable manner, and by writing tests that ensure that the code works as expected (Boswell & Foucher, 2011). Usability is hard to define (Brooke, 1996), yet one can consider a software as usable if the commands have intuitive names and if the software’s manual, termed “documentation”, is up-to-date and easy to understand for new users with minimal coding experience. A software package that has not received an update for a long time (approx. one year) could be considered abandoned. Abandoned software is unlikely to be fully functional, since it relies on other software (dependencies) that has changed in functionality or introduce bugs that were not expected by the developers of all dependencies. Together, software that’s scalable, maintainable and usable requires continuous changes to its codebase. There are best practices that standardize the continuous change of the codebase, including version control, continuous integration (often referred to as CI), and software testing.

Version control is a system that records changes to the codebase line by line, allowing the documentation of the history of the codebase, including who made which changes and when. This is required to isolate new and experimental features into newer versions and away from the stable version that’s known to work. The most popular version control system is Git, which is considered the industry standard for software development (Chacon & Straub, 2024). Git can use GitHub.com as a platform to store and host codebases in the form of software repositories. GitHub’s most famous feature is called “pull request”. A pull request is a request from anyone registered on GitHub to include their changes to the codebase (as in “please pull this into your main code”). One could see pull requests as the identifying feature of the open source community, since it exposes the codebase to potentially thousands of independent developers, reaching a workforce that is impossible to achieve with closed source models used by paid software companies.

Continuous integration (CI) is a software development practice in which developers integrate code changes into a shared repository several times a day (Duvall et al., 2007). Each integration triggers the test suite, aiming to detect errors as soon as possible. The test suite includes building the software, setting up an environment for the software to run and then executing the programmed tests, ensuring that the software runs as a whole. Continuous integration is often used together with software branches. Branches are independent copies of the codebase that are meant to be merged back into the original code once the changes are finished. Since branches

accumulate multiple changes over time, this can lead to minor incompatibilities between the branches of all developers (integration conflicts), which is something that CI helps to prevent.

Continuous integration especially relies on a thorough software testing suite. Software testing is the practice of writing code that checks if the codebase works as expected (Myers et al., 2011). The main type of software testing is unit testing, which tests the smallest units of the codebase (functions and classes) in isolation (Listing 1).

Listing 1: Example of an arbitrary Python function and its respective unit test function. The first function simply returns the number 5. The second function tests if the first function indeed returns the number 5. The test function is named with the prefix “test_” and is placed in a file that ends with the suffix “_test.py”. The test function is executed by the testing framework pytest. Note that code after “#” is considered a comment and won’t be executed.

```
1 # Define a function called "give_me_five" that returns the number 5
2 def give_me_five():
3     return 5
4 # Define a test function asserting that "give_me_five" returns 5
5 def test_give_me_five():
6     assert give_me_five() == 5
```

The quality of the software testing suite is measured by the code coverage, the precision of the tests, and the number of test-cases that are checked. The code coverage is the percentage of the codebase that is called by the testing functions, which should be as close to 100% as possible, although it does not measure how well the code is tested. The precision of the test is not a measurable quantity, but it represents if the tests truly checks if the code works as expected. The number of test-cases is the number of different scenarios that are checked by the testing functions, for example testing every possible option or combinations of options for functions that have multiple options. The most popular software testing framework for Python is pytest, which is utilized by plotastic (Krekel et al., 2004).

Together, the standards of software development, including version control, continuous integration, and software testing, ensure that the software is scalable, maintainable, and usable. This is especially important for software that is used by the scientific community, as it ensures that the software is working as expected at defined versions years after publishing scientific results.

Python as a Programming Language

Here, we provide a general overview of the Python programming language, explaining terms like “*type*”, “*method*”, etc., in order to prepare readers without prior programming experience for the following chapters. We also describe the design principles of Python to lay out the key concepts that differentiate Python compared to other programming languages. A more detailed tutorial on Python that’s specialized for bioscientists is found in Ekmekci et al. 2016

Listing 2: Example of readable Python code. This one-line code returns the words (string) 'Hello, World!' when executed. The command is straightforward and easy to understand.

```
1 print("Hello, World!")
2 # Output: Hello, World!
```

Languages such as Python are considered “*high-level*”, which means that it is designed to be easy to read and write, but also independent of hardware by hiding (“*abstracting*”) underlying details (*The Python Language Reference*, 2024). A key principle of Python is the emphasis on implementing a syntax that is concise and close to human language (Listing 2, Listing 3).

Listing 3: Example of less readable code written in the low-level programming language C. This code is doing exactly the same as the Python code in Listing 2. The command is harder to understand because more steps are needed to access the same functionality, including the definition of a function

```
1 #include <stdio.h>
2 int main() {
3     printf("Hello, World!");
4     return 0;
5 }
6 // Output: Hello, World!
```

Furthermore, Python is an *interpreted* language, which means that the code is executed line by line. This makes coding easier because the programmer can see the results of the code immediately after writing it, and error messages point to the exact line where the error occurred. This is in contrast to *compiled* languages, where the code has to be compiled into machine code before it can be executed. The advantage of compiled languages is that the code runs faster, because the machine code is optimized for the hardware.

Python automates tasks that would otherwise require an advanced understanding of computer hardware, like the need for manual allocation of memory space. This is achieved by using a technique called “*garbage collection*”, which automatically frees memory space that is no longer needed by the program. This is a feature that is not present in low-level programming languages like C or C++, that were designed to maximize control over hardware.

Another hallmark of Python is its *dynamic typing system*. In Python the type is inferred automatically during code execution (Listing 4). This is in contrast to *statically* typed languages like C, where the type of a variable has to be declared explicitly and cannot be changed during code execution (Listing 5) (?).

Dynamic typing makes Python a very beginner-friendly language, since one does not have to keep track of the type of each variable. However, this also makes Python a slower language, because the interpreter has to check the type of each variable during code execution. Also, developing code with dynamic typing systems is prone to introducing bugs (“type errors”), because it allows unexperienced developers to convert variables from one type to another without noticing, leading to unexpected behavior. Hence, larger Python projects require disciplined

Listing 4: Example of dynamic typing in Python. The variable “a” is assigned the value 5, which is of type integer. The variable “a” is then assigned the value “Hello, World!”, which is of type string. Python allows dynamic re-assignment of variables with different types. Note that code after “#” is considered a comment and won’t be executed.

```
1 a = 5 # Type integer
2 a = 5.0 # Type float
3 a = 'Hello, World!' # Type string
4 a = True # Type boolean
5 a = False # Type boolean
6 a = [1, 2, 3] # Type list of integers
7 a = {'name': 'Regina'} # Type dictionary
```

Listing 5: Example of static typing in C. The variable “a” is declared as an integer (int), and can only store integers. The variable “a” is then assigned the value 5, which is an integer. The variable “a” is then assigned the value ‘Hello, World!’, which is a string. This results in a compilation error, because the variable “a” can only store integers. Note that code after “//” is considered a comment and won’t be executed.

```
1 int a; // Declare type as integer
2 a = 5;
3 a = 'Hello, World!'; // Compilation error!
```

adherence to programming conventions. One such convention is *type hinting*, which is a way to explicitly note the type of a variable. Type hinting does not have an effect on the code, but it makes the code more readable and understandable for other developers, and allows for development environments to detect type errors before execution (Listing 6) (van Rossum et al., 2014).

Listing 6: Example of type hints used in Python. Explicitly stating the type of the variable is optional and does not change the behavior of the code as shown in Listing 4.

```
1 a: int = 5
2 a: str = 'Hello, World!'
```

Python supports both functional and object-oriented programming paradigms. In functional programming, the code is written in a way that the program is a sequence of function calls, where each function call returns a value that is used in the next function call (Listing 7). This approach is useful when multiple actions have to be performed on the same data and the structure of the data is relatively simple, for example a string of a gene sequence.

When the data itself gains in complexity, for example when storing not just the gene sequence, but also the promotor sequence, an object-oriented approach is more suitable (Listing 8). Object-oriented programming is a programming paradigm that uses objects and classes. An object is a collection of both data and functions, and a class is a blueprint for creating objects. The data of an object is stored as attributes. Functions that are associated with an object are called methods.

Listing 7: Example of functional programming in Python. The code defines a function called “find_restriction_site” that finds the position of a restriction site in a gene. The function “cut” uses the function “find_restriction_site” to cut the gene at the restriction site.

```
1 def find_restriction_site(gene: str):
2     return gene.find('GCGC')
3
4 def cut(gene: str):
5     position = find_restriction_site(gene)
6     return gene[position:]
7
8 gene1 = 'TGAGCTGAGCTGATGCGCTATATTTAGGCG'
9 gene1_cut = cut(gene1)
10 print(gene1_cut)
11 # Output: GCGCTATATTTAGGCG
```

Listing 8: Example of object oriented programming in Python. The class is called “Gene” and has four methods, “__init__”, “find_promotor”, “find_restriction_site” and “cut”. The method “__init__” is called when creating (“initializing”) an object, which fills the object with user-defined data. The parameter “self” is used to reference the object itself internally. “find_promotor” is a method that finds the position of the promotor in the gene and is called during object initialization.

```
1 class Gene:
2     def __init__(self, sequence: str):
3         self.sequence: str = sequence # Save sequence as attribute
4         self.promotor: str = self.find_promotor()
5     def find_promotor(self):
6         return self.sequence.find('TATA')
7     def find_restriction_site(self):
8         return self.sequence.find('GCGC')
9     def cut(self):
10        position = self.find_restriction_site()
11        return self.sequence[position:]
12
13 gene1 = Gene(sequence='TGAGCTGAGCTGATGCGCTATATTTAGGCG') # Create object
14 gene1_cut = gene1.cut() # Call the method cut
15 print(gene1_cut) # Show result
16 # Output: GCGCTATATTTAGGCG
```

A major benefit of using an object oriented versus a functional approach is that the data itself is programmable, enabling the programmer to define the behavior of the data itself through methods. This is achieved by using the keyword “self” to reference the object itself inside the class. For example, one could extend the class “Gene” with a method that finds the promotor of the gene and stores it as an attribute (Listing 8).

When designing software, both functional and object oriented programming can be used together, where object oriented programming is often used to design the program’s overall architecture, and functional programming is used to implement the algorithms of the program’s features. This allows for scalability of the software, as every single class is extended through the

addition of new methods. Furthermore, classes can be expanded in their functionalities through inheritance (Listing 9).

Listing 9: Example of inheritance in Python. The class “mRNA” inherits from the class “Gene”. The class “mRNA” has two methods, “__init__” and “find_stopcodon”. The method “find_stopcodon” finds the position of stop codons.

```
1 # Define a class called mRNA inheriting from the class Gene
2 class mRNA(Gene):
3     def __init__(self, sequence: str):
4         super().__init__(sequence) # Get attributes from parent class
5         self.sequence.replace('T', 'U') # Replace thymine with uracil
6     def find_stopcodons(self):
7         return self.sequence.find('UGA')
8
9 mrna1 = mRNA(sequence='TGAGCTGAGCTGATGCGCTATATTTAGGCG') # Create object
10 print(mrna1.find_stopcodons()) # Call the method translate
11 # Output: [0, 5, 10]
```

Inheritance is a feature of object-oriented programming that allows a class to access every attribute and method of a parent class. For example, one could extend the class “Gene” with a class “mRNA”, by writing a class “mRNA” that inherits from the class “Gene”.

Together, Python is not just beginner-friendly, but also well respected for its ease in development, which is why it is widely used in professional settings for web development, data analysis, machine learning, biosciences and more (Ekmekci et al., 2016; Rayhan & Gross, 2023).

The Potential of Python Data Science Packages for Biomedicine

Python includes a vast number of built-in packages used for basic data-types, software development, simple math operations, etc., (*The Python Language Reference*, 2024). Still, Python relies on packages developed by its users to provide specialized tools for data analysis. A Python package consists of multiple Python *modules*, where each module is a text-file with a .py ending containing Python code. Famous examples of such packages are *pytorch* and *tensorflow*, that are used to build models of artificial intelligence, including *ChatGPT* (Paszke et al., 2019; Abadi et al., 2016; Radford et al., 2019). Here, we outlay the most important packages used for *plotastic* in Chapter 2 and present examples how these packages are utilized in modern biomedical research.

Interactive Python: The standard Python interface is insufficient for data science, because it lacks the tools to quickly and conveniently visualize and explore data. IPython can be understood as an enhanced version of the standard Python interpreter, designed to improve the interactivity of Python code execution (Perez & Granger, 2007). IPython introduces features like dynamic type introspection and an interactive shell that offers rich media support. This

functionality is akin to what *MATLAB* and *RStudio* provide through their advanced graphical user interfaces and extensive debugging tools. This is particularly advantageous in the field of biomedicine, where visualizing data trends and patterns can lead to significant insights; however, IPython is most often utilized in the form of *Jupyter Notebooks*.

Jupyter: Jupyter is an evolution of IPython, introducing the *Jupyter notebook* format, having a file-ending `.ipynb` (Kluyver et al., 2016). Jupyter Notebooks are documents that combine both code and text structured as *code cells* and *markdown cells*, respectively. Markdown cells allow the author to provide additional information with text formatting, for example structuring the document with headings and subheadings, adding hyperlinks, images and mathematical formulas. Code cells can be executed individually, displaying the output directly below the cell. This allows for an interactive exploration of data, but also makes Jupyter Notebooks a very human-readable format that outlays data analysis in a clear manner with precise and reproducible documentation of all data processing steps. A major benefit of Jupyter Notebooks are interchangeable *Kernels*, allowing the execution of code in different programming languages, such as R, Julia, and C++ (Giorgi et al., 2022). Today, Jupyter Notebooks have become a standard format compatible with collaborative platforms like *Google Colab* and *JupyterLab*, but also professional software development tools like *VS Code*, and *PyCharm*. For biomedical research, Jupyter Notebooks hold great potential to improve reproducibility, as they provide a standardized format to present data analyses, and are found in the supplemental of modern publications of both bioinformatics and wet-lab research (Taskiran et al., 2024; Bosch-Queralt et al., 2022; Howe & Chain, 2015).

NumPy: Central processing units (CPU) usually execute one instruction on one data point at a time. For manipulating tabular data, this is inefficient as the same instruction must be repeatedly loaded for every data point. NumPy accelerates the mathematical capabilities of Python by enabling large-scale operations on multi-dimensional arrays and matrices with high efficiency (Harris et al., 2020). One key feature of NumPy is the implementation of “vectorization” or SIMD (Single Instruction, Multiple Data) instructions. SIMD allows multiple data points to be processed simultaneously, significantly speeding up operations that are inherently parallelizable, such as matrix addition or multiplication. NumPy’s syntax and functional approach to array manipulation have set a standard for matrix computation, influencing the design of advanced AI frameworks such as PyTorch and `mlx` (Paszke et al., 2019; *ML-Explore/MLx*, 2024), which mirrors several of NumPy’s functionalities to facilitate ease of use for those familiar with NumPy. This standardization has made NumPy an attractive tool not only in genomics (Ding et al., 2023), but also for modern clinical applications, such as imaging technologies and augmented-reality in surgery (Thompson et al., 2020).

Pandas: Tables are the most common way to store experimental results. Pandas extends

Python with a tabular datatype, called `DataFrame`, which allows for easy data manipulation with integrated indexing (McKinney, 2011). The intuitive interface of Pandas can be likened to *Microsoft Excel*; however, it is vastly more powerful due to its speed, functionality, and ability to handle larger datasets by running efficient numpy vectorization in the background. Unlike *Excel*, Pandas enables automation by summarizing processing commands into scripts, documenting each step, and ensuring reproducibility. Pandas is used in biomedicine for data wrangling, data cleaning, and data analysis, as it allows for the integration of multiple data sources into a single table (Santos et al., 2020).

matplotlib: `matplotlib` is a plotting library that provides a wide range of static, animated, and interactive plots and graphs (Hunter, 2007). It serves as the foundation for many other visualization tools and is particularly valued for its flexibility and customization options. Researchers in biomedicine use `matplotlib` to create a variety of graphs (like histograms or scatter plots), which are essential for preliminary data analysis and checking data distributions.

seaborn: While `matplotlib` is valued for its flexibility, it can be cumbersome to use for complex visualizations, e.g. by using different syntaxes for different types of plots. `seaborn` builds on `matplotlib` by integrating closely with Pandas data structures and providing a high-level interface for drawing attractive and informative statistical graphics (Waskom, 2021). It simplifies the creation of complex visualizations involving multidimensional data, making it easier to reveal patterns and relationships via color encoding, faceting, and automated statistical fits. This is particularly useful in biomedical research for visualizing and understanding complex datasets, such as large quantities of protein data (Weiss, 2022). `seaborn` could indirectly contribute to improving reproducibility in biomedical research by making visualizations of complex data very accessible through an easy and standardized syntax.

Pingouin: Integrating both data visualization and statistical analysis is beneficial for researchers who wish to conduct advanced statistical analysis without switching between different software environments. `Pingouin` is designed to be a user-friendly statistical tool that offers a straightforward syntax for performing statistical tests, which are commonly implemented in R (Vallat, 2018). Unlike R, `Pingouin` integrates seamlessly within the Python ecosystem, which allows combining data manipulation, analysis, and visualization all in one platform. This improves reproducibility by reducing the number of software tools required to analyze data. Despite its potential to streamline the data analysis process, `Pingouin` has not been widely adopted by biomedical research, yet. One example of a study that utilized `Pingouin` is the work of Kelly et al. (2023) in the field of Patient Public Involvement (PPI), producing an ethical matrix that allows for the inclusion of stakeholder opinion in medical research design. This lack of `Pingouin`'s adoption in biomedicine could be due to recent development and the dominance of R in the field. However, since Python offers multiple benefits over R in ease of use, software

development, runtime performance and integration with other tools (like including performant C++ code), Pingouin is an attractive standard for future statistical analyses in biomedicine (Gorelick & Ozsvald, 2020).

Together, these python packages form the backbone of modern data analysis in Python, often times combining software from different languages to accelerate certain features, while retaining the ease of use and readability that Python is known for. This is particularly advantageous in the field of biomedicine, where the requirements of modern data analysis are often complex and require a high degree of flexibility and customization.

Aims

This project defines these aims:

- Characterize the interaction between myeloma cells and mesenchymal stromal cells
- Aim 2
- Aim 3

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... Zheng, X. (2016, March). *TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems* (No. arXiv:1603.04467). arXiv. Retrieved 2024-03-07, from <http://arxiv.org/abs/1603.04467> doi: 10.48550/arXiv.1603.04467
- Bondi, A. B. (2000, September). Characteristics of scalability and their impact on performance. In *Proceedings of the 2nd international workshop on Software and performance* (pp. 195–203). New York, NY, USA: Association for Computing Machinery. Retrieved 2024-03-07, from <https://dl.acm.org/doi/10.1145/350391.350432> doi: 10.1145/350391.350432
- Bosch-Queralt, M., Tiwari, V., Damkou, A., Vaculčíaková, L., Alexopoulos, I., & Simons, M. (2022, March). A fluorescence microscopy-based protocol for volumetric measurement of lysolecithin lesion-associated de- and re-myelination in mouse brain. *STAR protocols*, 3(1), 101141. doi: 10.1016/j.xpro.2022.101141
- Boswell, D., & Foucher, T. (2011). *The Art of Readable Code: Simple and Practical Techniques for Writing Better Code*. "O'Reilly Media, Inc."
- Brooke, J. (1996, January). SUS – a quick and dirty usability scale. In (pp. 189–194).
- Chacon, S., & Straub, B. (2024, March). *Git - Book*. Retrieved 2024-03-07, from <https://git-scm.com/book/de/v2>
- Ding, W., Goldberg, D., & Zhou, W. (2023, August). PyComplexHeatmap: A Python package to visualize multimodal genomics data. *iMeta*, 2(3), e115. doi: 10.1002/imt2.115
- Duvall, P., Matyas, S., & Glover, A. (2007). *Continuous integration: Improving software quality and reducing risk*. Pearson Education. Retrieved from <https://books.google.de/books?id=PV9qfEdv9L0C>
- Ekmekci, B., McAnany, C. E., & Mura, C. (2016, July). An Introduction to Programming for Bioscientists: A Python-Based Primer. *PLOS Computational Biology*, 12(6), e1004867. Retrieved 2024-03-10, from <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1004867> doi: 10.1371/journal.pcbi.1004867
- Giorgi, F. M., Ceraolo, C., & Mercatelli, D. (2022, April). The R Language: An Engine for Bioinformatics and Data Science. *Life*, 12(5), 648. Retrieved 2024-04-21, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9148156/> doi: 10.3390/life12050648
- Gómez-López, G., Dopazo, J., Cigudosa, J. C., Valencia, A., & Al-Shahrour, F. (2019, May). Precision medicine needs pioneering clinical bioinformaticians. *Briefings in Bioinformatics*, 20(3), 752–766. doi: 10.1093/bib/bbx144
- Gorelick, M., & Ozsvald, I. (2020). *High Performance Python: Practical Performant Programming for Humans*. "O'Reilly Media, Inc."
- Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., ... Oliphant, T. E. (2020, September). Array programming with NumPy. *Nature*, 585(7825), 357–362. Retrieved 2023-08-09, from <https://www.nature.com/articles/s41586-020-2649-2> doi: 10.1038/s41586-020-2649-2
- Howe, A., & Chain, P. S. G. (2015). Challenges and opportunities in understanding microbial communities with metagenome assembly (accompanied by IPython Notebook tutorial). *Frontiers in Microbiology*, 6, 678. doi: 10.3389/fmicb.2015.00678
- Hunter, J. D. (2007, May). Matplotlib: A 2D Graphics Environment. *Computing in Science & Engineering*, 9(3), 90–95. Retrieved 2023-11-15, from <https://ieeexplore.ieee.org/document/4160265> doi: 10.1109/MCSE.2007.55
- Incerti, D., Thom, H., Baio, G., & Jansen, J. P. (2019, May). R You Still Using Excel? The Advantages of Modern Software Tools for Health Technology Assessment. *Value in Health*, 22(5), 575–579. Retrieved

- 2024-03-11, from <https://www.sciencedirect.com/science/article/pii/S1098301519300506> doi: 10.1016/j.jval.2019.01.003
- Kazman, R., Bianco, P., Ivers, J., & Klein, J. (2020, December). *Maintainability* (Report). Carnegie Mellon University. Retrieved 2024-03-07, from <https://kilthub.cmu.edu/articles/report/Maintainability/12954908/1> doi: 10.1184/R1/12954908.v1
- Kelly, B. S., Kirwan, A., Quinn, M. S., Kelly, A. M., Mathur, P., Lawlor, A., & Killeen, R. P. (2023, May). The ethical matrix as a method for involving people living with disease and the wider public (PPI) in near-term artificial intelligence research. *Radiography (London, England: 1995)*, 29 Suppl 1, S103-S111. doi: 10.1016/j.radi.2023.03.009
- Kluyver, T., Ragan-Kelley, B., Pérez, F., Granger, B., Bussonnier, M., Frederic, J., ... Jupyter Development Team (2016). *Jupyter Notebooks—a publishing format for reproducible computational workflows*. Retrieved 2024-04-20, from <https://ui.adsabs.harvard.edu/abs/2016ppap.book...87K> doi: 10.3233/978-1-61499-649-1-87
- Krekel, H., Oliveira, B., Pfannschmidt, R., Bruynooghe, F., Laughner, B., & Bruhin, F. (2004). *Pytest*. Retrieved from <https://github.com/pytest-dev/pytest>
- McKinney, W. (2011, January). Pandas: A Foundational Python Library for Data Analysis and Statistics. *Python High Performance Science Computer*.
- ML-explore/mlx*. (2024, April). ml-explore. Retrieved 2024-04-21, from <https://github.com/ml-explore/mlx>
- Myers, G. J., Sandler, C., & Badgett, T. (2011). *The art of software testing* (3rd ed.). Wiley Publishing. Retrieved from <https://malenezi.github.io/malenezi/SE401/Books/114-the-art-of-software-testing-3-edition.pdf>
- Narzt, W., Pichler, J., Pirklbauer, K., & Zwintz, M. (1998, January). A Reusability Concept for Process Automation Software..
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., ... Chintala, S. (2019, December). *PyTorch: An Imperative Style, High-Performance Deep Learning Library* (No. arXiv:1912.01703). arXiv. Retrieved 2024-03-07, from <http://arxiv.org/abs/1912.01703> doi: 10.48550/arXiv.1912.01703
- Peng, R. D. (2011, December). Reproducible Research in Computational Science. *Science*, 334(6060), 1226–1227. Retrieved 2024-03-18, from <https://www.science.org/doi/10.1126/science.1213847> doi: 10.1126/science.1213847
- Perez, F., & Granger, B. E. (2007, May). IPython: A System for Interactive Scientific Computing. *Computing in Science & Engineering*, 9(3), 21–29. Retrieved 2024-04-20, from <https://ieeexplore.ieee.org/document/4160251> doi: 10.1109/MCSE.2007.53
- The Python Language Reference*. (2024). Retrieved 2024-03-07, from <https://docs.python.org/3/reference/index.html>
- R Core Team. (2018). *R: A language and environment for statistical computing* [Manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language Models are Unsupervised Multitask Learners.. Retrieved 2024-03-07, from <https://www.semanticscholar.org/paper/Language-Models-are-Unsupervised-Multitask-Learners-Radford-Wu/9405cc0d6169988371b2755e573cc28650d14dfe>
- Rayhan, A., & Gross, D. (2023). *The Rise of Python: A Survey of Recent Research*. doi: 10.13140/RG.2.2.27388.92809
- Sandve, G. K., Nekrutenko, A., Taylor, J., & Hovig, E. (2013, October). Ten Simple Rules for Reproducible Computational Research. *PLoS Computational Biology*, 9(10), e1003285. Retrieved 2024-03-07, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3812051/> doi: 10.1371/journal.pcbi.1003285
- Santos, B. S., Silva, I., Ribeiro-Dantas, M. d. C., Alves, G., Endo, P. T., & Lima, L. (2020, October). COVID-19: A scholarly production dataset report for research analysis. *Data in Brief*, 32, 106178. doi: 10.1016/

- j.dib.2020.106178
- Smith, A. M., Niemeyer, K. E., Katz, D. S., Barba, L. A., Githinji, G., Gymrek, M., ... Vanderplas, J. T. (2018). Journal of Open Source Software (JOSS): Design and first-year review. *PeerJ Preprints*, 4, e147. doi: 10.7717/peerj-cs.147
- Tanavalee, C., Luksanapruksa, P., & Singhatanadgige, W. (2016, June). Limitations of Using Microsoft Excel Version 2016 (MS Excel 2016) for Statistical Analysis for Medical Research. *Clinical Spine Surgery*, 29(5), 203. Retrieved 2024-03-11, from https://journals.lww.com/jspinaldisorders/fulltext/2016/06000/limitations_of_using_microsoft_excel_version_2016.5.aspx doi: 10.1097/BSD.0000000000000382
- Taskiran, I. I., Spanier, K. I., Dickmanken, H., Kempynck, N., Pančíková, A., Ekşi, E. C., ... Aerts, S. (2024, February). Cell-type-directed design of synthetic enhancers. *Nature*, 626(7997), 212–220. Retrieved 2024-04-21, from <https://www.nature.com/articles/s41586-023-06936-2> doi: 10.1038/s41586-023-06936-2
- Thompson, S., Dowrick, T., Ahmad, M., Xiao, G., Koo, B., Bonmati, E., ... Clarkson, M. J. (2020, July). SciKit-Surgery: Compact libraries for surgical navigation. *International Journal of Computer Assisted Radiology and Surgery*, 15(7), 1075–1084. doi: 10.1007/s11548-020-02180-5
- Vallat, R. (2018, November). Pingouin: Statistics in Python. *Journal of Open Source Software*, 3(31), 1026. Retrieved 2023-05-29, from <https://joss.theoj.org/papers/10.21105/joss.01026> doi: 10.21105/joss.01026
- van Rossum, G., Lehtosalo, J., & Langa, L. (2014). *PEP 484 – Type Hints* / [peps.python.org](https://peps.python.org/pep-0484/). Retrieved 2024-03-08, from <https://peps.python.org/pep-0484/>
- Waskom, M. L. (2021, April). Seaborn: Statistical data visualization. *Journal of Open Source Software*, 6(60), 3021. Retrieved 2023-03-26, from <https://joss.theoj.org/papers/10.21105/joss.03021> doi: 10.21105/joss.03021
- Weiss, C. J. (2022, September). Visualizing protein big data using Python and Jupyter notebooks. *Biochemistry and Molecular Biology Education: A Bimonthly Publication of the International Union of Biochemistry and Molecular Biology*, 50(5), 431–436. doi: 10.1002/bmb.21621
- Yang, A., Troup, M., & Ho, J. W. (2017, July). Scalability and Validation of Big Data Bioinformatics Software. *Computational and Structural Biotechnology Journal*, 15, 379–386. Retrieved 2024-03-07, from <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5537105/> doi: 10.1016/j.csbj.2017.07.002

Appendices

A Supplementary Data & Methods

A.1 Figures

A.2 Tables

A.3 Materials & Methods

B Documentation of `plotastic`