

# From bioinformatics user to bioinformatics engineer: a report

*Gilderlanio Santana de Araújo*

*gilderlanio@gmail.com*

*Laboratório de Genética Humana e Médica*

*Programa de Pós-Graduação em Genética e Biologia Molecular*

*Instituto de Ciências Biológicas*

*Universidade Federal do Pará, Belém, Brazil.*

## 1 Abstract

2 Teaching computer programming is not a simple task and it is challenging to introduce the concepts of program-  
3 ming in graduate programs of other fields. Little efforts have been made on engaging students in computational  
4 development after programming trainings. An emerging need is to establish subjects of bioinformatics and pro-  
5 gramming languages in genetics and molecular biology graduate programs, when students in these degree programs  
6 are immersed in a sea of genomic and transcriptomic data, which demands proficient computational treatment.  
7 I report an empirical guideline to introduce programming languages and recommend Python as first language  
8 for graduate programs in which students were from genetics and molecular biology backgrounds. Including the  
9 development of programming solutions related to graduate students' research activities may improve program-  
10 ming skills and better engagement. These results suggest that the applied approach leads to enhanced learning  
11 of introductory to autonomy in highly advanced programming concepts by graduate students. This guide should  
12 be extended for other research programs.

## 13 1 Introduction

14 Bioinformatics is an interdisciplinary area that requires a depth knowledge in computational, statisti-  
15 cal/mathematical and life sciences subjects. In 2016, the ISCB Education Committee's Curriculum Task  
16 Force described needs for bioinformatics education and competencies of bioinformatics engineers who  
17 actively develops algorithms and computational systems, as well as of bioinformatics users that explore  
18 computational infrastructures and softwares in specific contexts to work on data analysis of different  
19 sorts, such as population genetics, phylogeny, medical genetics (Mulder *et al.*, 2018; Welch *et al.*, 2016).

20 Initiatives such as The European Bioinformatics Institute (EMBL-EBI - <https://www.ebi.ac.uk/>), the  
21 Pan African bioinformatics network for H3Africa (H3AbioNet - [www.h3abionet.org](http://www.h3abionet.org)), the Bioinformatics  
22 Multidisciplinary Environment (BioME - <https://bioinfo.imd.ufrn.br/>) and the Bioinformatics Interunits  
23 Program of the Federal University of Minas Gerais (<http://www.pgbioinfo.icb.ufmg.br/>) extensively pro-  
24 mote academic education and training in bioinformatics subjects as a way to supply the demand for  
25 bioinformatics engineers and users capable of processing and analyzing the high volume of data from the  
26 'omics' sciences, such as genomics, transcriptomics and proteomics.

27 Despite the number of programs in Bioinformatics, Attwood *et al.* (2017) identified a strong bias  
28 for conducting short courses to fill skill gaps towards bioinformatics analysis and, in the same study,  
29 the authors state that it is more necessary to mix bioinformatics subjects with life sciences programs.  
30 Next generation bioinformatics must be supported by feasible and engaging *curricula* in their graduate  
31 programs. The need for effective programming training is emergent, due to the problems of computational  
32 biology that have required implementation of robust algorithms, which must be able to handle large  
33 volumes of data generated by high-throughput sequencing technologies, for example.

34 Currently, most of life sciences graduate students are submerged in a sea of "omics" data and their  
35 educational programs have shown an inhibition or have adapted their curriculum program to include  
36 computer programming and information technology-related trainings. Teaching algorithms, logics or  
37 programming languages, in many cases, is a challenging task. Interestingly, bioinformatics candidate  
38 students in the North of Brazil are predominantly from different undergraduate backgrounds, in particular  
39 biology, biomedical sciences, medical schools, as further discussed in subsequent sections. In addition,  
40 there is a lack of methodologies aimed at real-time learning that demystifies and motivates the student to  
41 learn programming languages, which can favor the use and application of computational thinking after  
42 the end of the course.

43 Here, I report the experience of teaching programming languages to students of the Graduate Program  
44 in Genetics and Molecular Biology (PPGBM - <http://ppgbm.propesp.ufpa.br/>) at the Federal University  
45 of Pará, as well as the description of an empirical method for keeping these students engaged in the  
46 development of their computational solutions, after the training time. I proposed a course entitled  
47 "Programming for Bioinformatics with Python", in which I included core topics of Python, as well as  
48 empirical software engineering tasks, such as: definition of the scope of the candidates' research project  
49 and functional requirement elucidation. I adopted "divide and conquer" paradigm to script programming,  
50 pattern recognition of their data and processes, which suits well biological application development.  
51 Practical, interactive, and personalized activities in the context of each candidate in a real-time way  
52 improves consolidating concepts of programming languages and autonomy for life science candidates,  
53 which, in practice, put them in the path to become bioinformatics engineers.

## 54 2 Methods

### 55 2.1 Python for Bioinformatics

56 Python is a hybrid programming language that allows scripting in a functional and object-oriented  
57 paradigm (<https://www.python.org/doc/>). By its resources, Python is considered a production-ready  
58 language, provides clear syntax and semantics, taking advantages of mandatory code indentation, which  
59 improves readability and refactoring. Python prioritizes the developer experience, making many software  
60 engineers choose it as a programming language, based on their potential for productivity, learning curve  
61 cost, and computational support.

62 Python can be applied for general purposes, and have increased its use in bioinformatics. A survey  
63 of programming languages for bioinformatics was made on GitHub and points out the massive and  
64 predominant use of Python for genomic analysis (Suarez *et al.*, 2018).

65 Python improves productivity on the implementation of pipelines and integration of scientific work-  
66 flows. Python shows abundance of statistical libraries and supports mathematical computation with high  
67 potential for data science. Lately, some libraries in Python have been implemented with functions for  
68 large scale statistics, machine learning and high quality graphical representation of data. These are new  
69 features in Python, which aggregates some immeasurable features of the R language context, as for data  
70 visualization and data manipulation in dataframes structures.

## 71 2.2 Course structure and execution period

72 The course “Programming for Bioinformatics with Python” was designed to provide a basic knowledge in  
73 high-level programming languages for graduate students in genetics and molecular biology, at the Federal  
74 University of Pará.

75 This curricular component is a way of improving computational skills, encouraging and arousing  
76 autonomy in algorithm implementation for processing and analyzing biological data within the context  
77 of a master, doctoral or post-doctoral research project of candidates with little or no knowledge of  
78 programming languages.

79 The course was offered twice, from March to April 2019, and for the same months in 2020. The first  
80 class was formed by 16 students with little or no programming experience, while the second class was  
81 formed by 10 students with a similar background. In both periods, the programming classes had a 45-hour  
82 workload, with 15 classes of three hours. As of March 18, 2020, due to quarantine restrictions caused by  
83 the pandemics of COVID-19 in Brazil, all classes and monitoring of students’ project development in  
84 the second course started to be conducted remotely and this was maintained until the end of the course.

## 85 2.3 Bibliography and Integrated development environment (IDE)

86 The course was conducted with theoretical and essentially practical classes based on the computational  
87 biology literature. I adopted three reference textbooks: Bioinformatics Algorithms: Design and Imple-  
88 mentation in Python (Rocha and Ferreira, 2018), Learning python: Powerful object-oriented program-  
89 ming (Lutz, 2013) and the Introdução a Programação para Bioinformática com Biopython, which is only  
90 available in Brazilian Portuguese (Mariano *et al.*, 2015).

91 Differently from some courses, that use terminal or simple text-based editors, we conducted pro-  
92 gramming practical classes using the PyCharm Community (<https://www.jetbrains.com/pt-br/>), which  
93 is a professional integrated development environment (IDE) dedicated to improve Python programming.  
94 Linux Ubuntu environment was used for running scripts by command line.

## 95 2.4 Computer programming and research topics

96 I designed the course structure with 15 classes, 3 hours each, which included Python programming and  
97 research activities in parallel (see Figure 1). A total of 10 theoretical-practical lessons were performed,  
98 covering the following topics: introduction to computer programming; aspects of programming with  
99 Python (v3.0); data types; logical and arithmetic operations; data structures (list, dictionaries, tuples  
100 and sets); manipulation of strings; conditional and iterative structures; built-in functions, parameters,  
101 implementing new and reusing; using the command line and importing system libraries; file manipulation  
102 (creating, merging and writing raw data files). In the last practical lessons of the course, we explored  
103 libraries, such as BioPython (<https://biopython.org/>), Pandas (<https://pandas.pydata.org/>) and Seaborn  
104 (<https://seaborn.pydata.org/>), that are proposed to assist routines for data manipulation and visualization.

105 In parallel to programming Python topics, the course was designed to execute students’ research  
106 tasks. First, at the beginning of the course, all students were asked to provide a summary of their  
107 graduate projects. Second, problems in computational biology must be identified within the scope of  
108 each research project. Then, we were able to draw an overview of computational solutions, elicit and  
109 select functional requirements for bioinformatics problems, considering the time and scope. We adopted  
110 a divide-and-conquer strategy to implement solutions in the course time.

111 Five classes were used to coordinate and supervise the development of computational solutions in can-  
112 didates’ research context. The project development included four steps: a) initialization of bioinformatics

113 project development after lessons covering BioPython, Pandas and Seaborn; b) monitoring project de-  
114 velopment, which includes solving programming questions; c) ending project development and preparing  
115 presentations; and, finally, d) in-person or remote presentation of final projects to classmates.

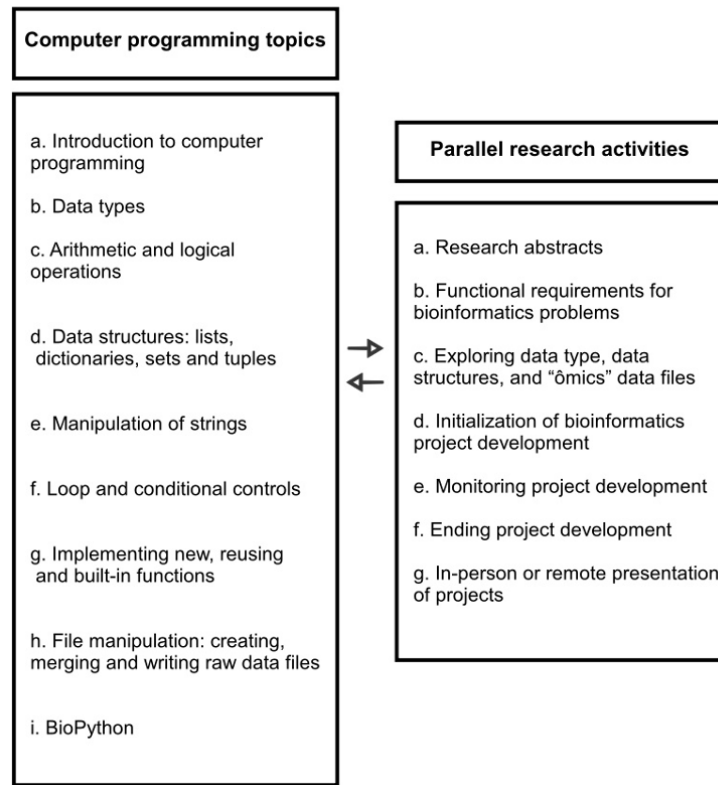


Fig. 1: Outline of the topics covered in the discipline of Programming for Bioinformatics with Python (left) and parallel research activities (right). Research activities were performed to engage the students with Python programming in face of their graduate research project.

116 In the classes on data types and data structures, we explored how "omics" data could be modeled  
117 and loaded in memory. Many research routines in genomics and transcriptomics were discussed, as well  
118 as how to use Python resources to model and implement solutions. For example, a DNA sequence can  
119 be simply represented by a string, or even by a more robust object like *Bio.Seq* from BioPython library.  
120 The *Seq* object provides methods similar to those implemented for strings, such as count, find, split  
121 and strip. In addition, the *Seq* object has an alphabet as an attribute, which can be instantiated from  
122 the *Bio.Alphabet* class, allowing you to build objects with a generic DNA or a standard alphabet of  
123 International Union of Pure and Applied Chemistry (IUPAC, <https://iupac.org/>). The *Seq* object also  
124 provides specific methods for manipulating DNA sequences, such as to obtain the complement, reverse  
125 complement and RNA sequence by transcription. On a large scale, the *Bio.SeqIO* object allows us to  
126 read sequences stored in *.fasta* files, which is the standard file to represent both nucleotide and peptide  
127 sequences.

## 128 2.5 Continuous practical lessons

129 We adopted a continuously and individually assessment through problem-solving and mainly development  
130 of projects in Python language. The final projects were of mandatory public presentation in the case  
131 of the first course, and remotely via Google Hangouts, in the case of the second course, as previously  
132 mentioned.

133 For example, a simple solution for calculating the genotype frequencies based on allele frequencies  
134 was requested, given that it is a routine task in the context of population genetics studies. The imple-  
135 mented algorithm must receive the frequencies of two alleles (A and B) of a genomic variant and print  
136 the genotypic frequencies at the terminal. Starting from this task, students were required to have a  
137 simple program, which would receive allele frequencies as input, and the system should inform genotype  
138 frequencies (AA, AB, BB).

139 A possible solution to this problem was implemented in Listing 1. In this example, we explored pro-  
140 gramming concepts of using Python libraries, reading and processing data from shell terminal, converting  
141 data types, conditional structures, arithmetic operations and user interaction.

```
142 1 import sys
143 2
144 3 # Capture from command line the frequencies of allele A and allele B.
145 4 frq_allele_a = float(sys.argv[1])
146 5 frq_allele_b = float(sys.argv[2])
147 6
148 7 # Validation of allele frequencies.
149 8 if (frq_allele_a + frq_allele_b == 1):
150 9
151 10     frq_aa = frq_allele_a ** 2
152 11     frq_ab = 2 * frq_allele_a * frq_allele_b
153 12     frq_bb = frq_allele_b ** 2
154 13     print("Freq. AA: ", frq_aa)
155 14     print("Freq. AB: ", frq_ab)
156 15     print("Freq. BB: ", frq_bb)
157 16 else:
158 17     print("Allele frequencies must be equals 1.")
```

Listing 1: Code example to calculate genotypic frequencies of a population from allele frequencies.

## 159 2.6 Course structure evaluation

160 For course evaluation purposes, a Google Forms was created to question students of both courses about  
161 basic training, the context in which the students developed or still develop their research project, which  
162 packages they used, as well as whether the student continued to work with Python. The questions that  
163 formed the research were defined as follows:

- 164 • Enter your graduation course.
- 165 • Have you continued to develop analyses using Python?
- 166 • In what context/theme did you develop the course project in Python?
- 167 • Which libraries, packages or modules do you use in your analyses?

168 We performed a qualitative analysis based on the students' discursive responses to assess the impact  
169 of the proposed course in their research after the end of the training period, as well as to understand the  
170 needs of genetic and molecular biology graduate programs regarding the development of bioinformatics  
171 tools and future training.

### 172 3 Results

173 The questionnaire via Google Forms was answered by 13 students distributed from Graduate Program  
174 in Genetics and Molecular Biology and from Graduate Program in Oncology and Medical Sciences, both  
175 at the Federal University of Pará. The students presented different undergraduate degree backgrounds,  
176 being predominant the education in Biology ( $n = 6$ ). Only one of the students reported a previous  
177 training in a field related to technology and data processing. The entire distribution of students by area  
178 is shown in Figure 2A.

179 In the form, there was a question on whether these students have developed their analysis using Python  
180 since the end of the course. On that matter, a group of 84 % of the students reported a continuous use  
181 of Python libraries for their research activities, while only two answered that they have not been using  
182 Python.

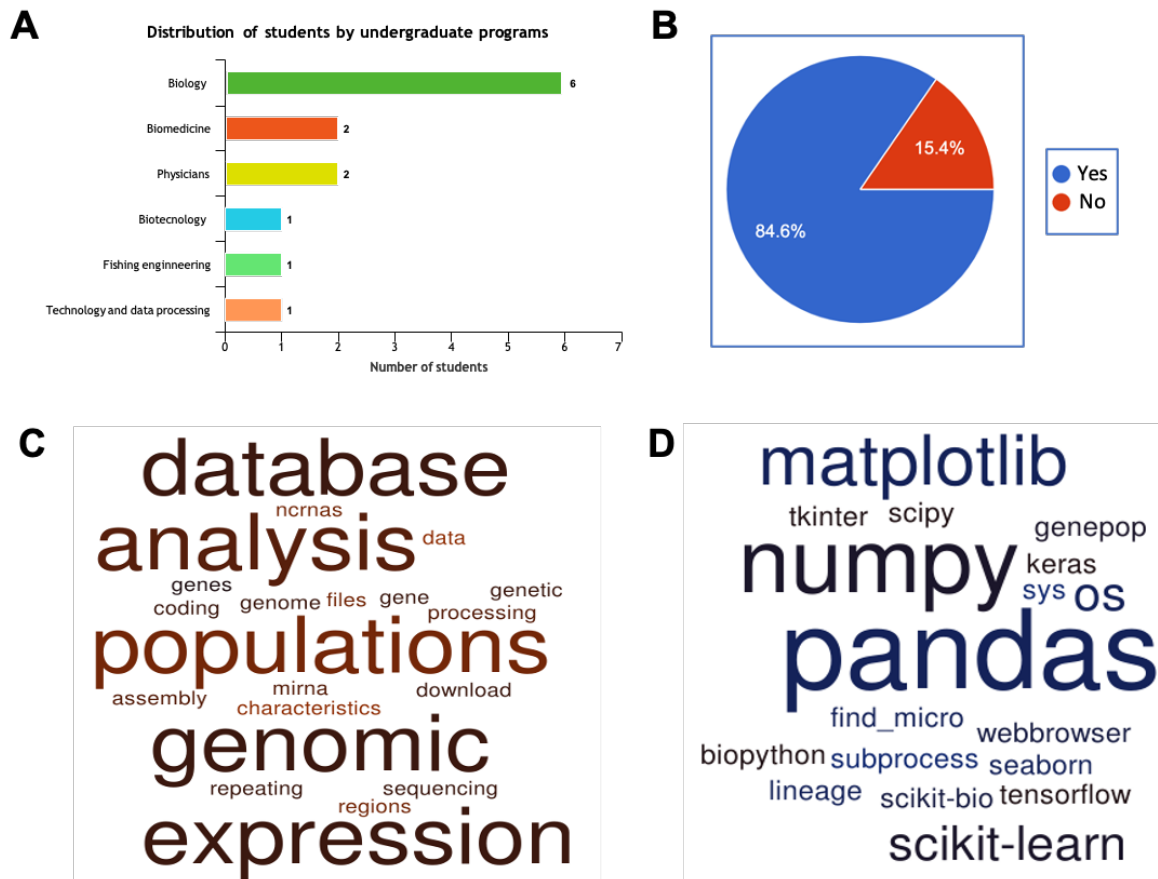


Fig. 2: Statistical summary of the Python discipline profile. (A) Distribution of students by undergraduate degree. (B) Proportion of students who continue to develop their Python projects after Python training. (C) Word cloud of the biological contexts, regarding the developed projects. (D) Python libraries used for the development of students' graduate projects.

183 From the students' discursive response, I analyzed the context in which the projects were developed  
184 and which libraries in Python were used to solve their bioinformatics research problems. In Figure ?? C,  
185 a word cloud summarizes several scopes related to biological fields and molecular and genetic elements.  
186 All candidates have developed scripts to process Next-Generation Sequencing (NGS) data, being mostly  
187 studies related to population genetics and data integration of coding and non-coding RNA databases.

188 In the context of each research projects, most candidates have implemented routines for extraction,  
189 loading and transformation of biological data for analysis. Data processing is a task that cost most of  
190 the time in implementing computational pipelines, for example input data for statistical tests, such as  
191 Fisher Test and Linear and Logistic regressions for studies of population genetics and gene expression,  
192 that require quality controls and format conversions. Other interesting analysis was a script developed to  
193 explore repeat elements in DNA sequences, which are short sequence of characters with 3, 4, 5 in length,  
194 that was screened in genome sequences of fishes. Large scale data were used by these students, such as  
195 databases that included 1000 Genomes Project (Consortium, 2010), Geuvadis (Lappalainen *et al.*, 2013),  
196 circBase (Glažar *et al.*, 2014), mirTarbase (Chou *et al.*, 2017), pirBase (Wang *et al.*, 2018) and The  
197 Cancer Genome Atlas Program (Weinstein *et al.*, 2013).

198 We found a predominance of data science-related libraries (see word cloud in Figure 2D) for genomic  
199 and transcriptomic analysis. More sophisticated analytic scripts were generated through the application  
200 of machine learning algorithms. In Python, several specialized libraries were adopted, such as BioPython,  
201 GenePop, Pandas, Seaborn, Scikit-learn and scikit-image. Based on these findings, I strong recommend  
202 data science-related courses for graduate students in Genetics and Molecular Biology in addition to  
203 Biostatistics.

204 As for reasearch topics, different themes were reported. An example was an analysis implemented in  
205 Python that aimed to perform expression quantitative trait analysis, in order to investigate the influence  
206 of polymorphisms on gene expression. This analysis was performed regarding genomic and transcriptomic  
207 data for African and European populations. Data were extracted from gEUVADIS and sample sequencing  
208 from the 1000 Genomes Project(Consortium, 2010; Lappalainen *et al.*, 2013). Initially, the *LabelEncoder*  
209 function of the *sklearn.preprocessing* package was used to encode the genotypes AA, AB, and BB in 0,  
210 1, and 2, respectively. Then, the function *snphwe* library supports us to evaluate the Hardy-Weinberg  
211 Equilibrium (HWE) for each polymorphism. The *shapiro* function from *scipy.stats* was used to assess  
212 the normality of gene expression data. For graphical representations, histograms were used to draw the  
213 distribution of gene expression and boxplots were used to represent the regressions. For this purpose, the  
214 following libraries were used: *matplotlib.pyplot*, *seaborn*, *numpy* and *statannot*.

## 215 4 Discussion

216 Here, I discuss and report the experience of teaching high-level programming languages, like Python,  
217 adding in parallel the active and analytic activities for the development of computational solutions,  
218 implemented by graduate students, with a focus on developing scripts for their research projects. We  
219 designed the training to cover the fundamental aspects of programming languages and research activities  
220 as differential on engaging students on scripting.

221 A high heterogeneity of 'omics' data in the candidate's research projects was notable. The greater the  
222 number of students, the greater the level of heterogeneity in research and development projects. Then, I  
223 carried out parallel data description and modeling activities in Python, as a way of better comprehending  
224 data types, data structures and file formats, which allowed a better perception to students regarding the  
225 manipulation of their data and projects.

226 Harmony was aimed between the research activities of graduate students and the development of their  
227 own solutions for their research projects. Remaining to develop solutions is no occasion. In most of the

228 reported methodologies, only the theoretical-practical basis in programming languages is explored, being  
229 necessary to instigate candidates' involvement spirit, that is, to make them less passive in the learning  
230 process and often dependent on computer specialists. This continuous engagement in the development  
231 of computational solutions is probably a reflection of the rapid and satisfactory return of parallel and  
232 applied activities of programming and own research activities. In this context, it is recommended to  
233 merge research activities in programming language subjects for graduate students of life sciences.

234 Notably, implementing solutions for student's research projects and obtaining quick solutions for their  
235 problems enhances their interest and curiosity to implement and use other Python resources not explored  
236 in the course. This is an aspect that highlights a level of autonomy achieved by students on developing  
237 their own solutions.

238 With the described approach, new perspectives on training graduate students were conceived, with  
239 subjects related to programming languages. These students are now able to deal with bioinformatics  
240 problems that require analysis of large scale data, such as genome sequences and transcriptomic data.  
241 The course methodology consequently demystifies the use of programming languages and presents itself  
242 as a unique opportunity for the application of computer knowledge, to achieve quick solutions.

243 Thus, I believe that the present work may contribute with ideas in the practical teaching of program-  
244 ming languages in the "omics" era, being a facilitator in the construction of knowledge in life sciences  
245 undergraduate and graduate programs.

246 I encourage the development of technical skills with professional tools such as PyCharm, qualifying  
247 the student to enter the industry market. In addition, this report reinforces the approaches of adapt-  
248 ing bioinformatics *curricula* for data science subjects, whereas mentioned techniques and methods are  
249 common in several research contents.

## 250 5 Conclusion

251 Genomics and transcriptomics are two research areas of constant application of data science methods  
252 and techniques to perform analysis on large data volumes. Python programming language have stood  
253 out in the scientific and industrial environment as support languages for building solutions. Among  
254 the "omics" sciences, this language has reached out by providing several tools for general purposes.  
255 Modules like BioPython, scikit-bio, scikit-learn, Pandas and seaborn have been used successfully by most  
256 graduate students for solving statistical problems with machine learning and functions for large-scale  
257 data manipulation.

258 Considering that only providing the essentials of programming languages might not give satisfactory  
259 results, adopting real-time computational development tasks to solve problems in each student's research  
260 context entails the engagement in the development of scripts that automate their daily tasks in labora-  
261 tories. This aspect has been the motivational element to make the students have a real perception of the  
262 applicability of their own scripts or computational pipelines. This fact corroborates the high percentage  
263 of students who still use Python after the end of the course.

264 In this way, I believe that this report contributes to consolidate new teaching methodologies, including  
265 applied classes of high-level programming languages like Python for bioinformatics in the era of "omics"  
266 sciences.

## 267 Acknowledgements

268 Thanks to all graduate students that answered the online questionnaire and remain doing science and  
269 scripting in Python, even in pandemic situations. The outcomes of that questionnaire provided a helpful  
270 feedback to improve programming language training.



## 271 References

- 272 Attwood, T. K., Blackford, S., Brazas, M. D., Davies, A., and Schneider, M. V. (2017). A global  
273 perspective on evolving bioinformatics and data science training needs. *Briefings in Bioinformatics*,  
274 **20**(2), 398–404.
- 275 Chou, C.-H., Shrestha, S., Yang, C.-D., Chang, N.-W., Lin, Y.-L., Liao, K.-W., Huang, W.-C., Sun,  
276 T.-H., Tu, S.-J., Lee, W.-H., *et al.* (2017). Mirtarbase update 2018: A resource for experimentally  
277 validated microrna-target interactions. *Nucleic acids research*, **46**(D1), D296–D302.
- 278 Consortium, . G. P. (2010). A map of human genome variation from population-scale sequencing. *Nature*,  
279 **467**(7319), 1061.
- 280 Glázar, P., Papavasileiou, P., and Rajewsky, N. (2014). Circbase: A database for circular rnas. *Rna*,  
281 **20**(11), 1666–1670.
- 282 Lappalainen, T., Sammeth, M., Friedländer, M. R., Ac't Hoen, P., Monlong, J., Rivas, M. A., Gonzalez-  
283 Porta, M., Kurbatova, N., Griebel, T., Ferreira, P. G., *et al.* (2013). Transcriptome and genome  
284 sequencing uncovers functional variation in humans. *Nature*, **501**(7468), 506–511.
- 285 Lutz, M. (2013). *Learning python: Powerful object-oriented programming.* " O'Reilly Media, Inc."
- 286 Mariano, D. D. C., Barroso, J. R. P., Correia, T. S., and Melo-Minardi, R. C. (2015). *Introdução à*  
287 *Programação para Bioinformática com Biopython*, volume 1. CreateSpace Independent Publishing  
288 Platform, 3 edition.
- 289 Mulder, N., Schwartz, R., Brazas, M. D., Brooksbank, C., Gaeta, B., Morgan, S. L., Pauley, M. A.,  
290 Rosenwald, A., Rustici, G., Sierk, M., *et al.* (2018). The development and application of bioinformatics  
291 core competencies to improve bioinformatics training and education. *PLoS computational biology*,  
292 **14**(2), e1005772.
- 293 Rocha, M. and Ferreira, P. G. (2018). *Bioinformatics Algorithms: Design and Implementation in Python*.  
294 Academic Press.
- 295 Suarez, C. G. H., Burbano, M. E. G., Guerrero, V. A. B., and Tovar, P. A. M. (2018). Bioinformatics  
296 software for genomic: a systematic review on github. Technical report, PeerJ Preprints.
- 297 Wang, J., Zhang, P., Lu, Y., Li, Y., Zheng, Y., Kan, Y., Chen, R., and He, S. (2018). Pirbase: A  
298 comprehensive database of pirna sequences. *Nucleic acids research*, **47**(D1), D175–D180.
- 299 Weinstein, J. N., Collisson, E. A., Mills, G. B., Shaw, K. R. M., Ozenberger, B. A., Ellrott, K., Shmule-  
300 vich, I., Sander, C., Stuart, J. M., Network, C. G. A. R., *et al.* (2013). The cancer genome atlas  
301 pan-cancer analysis project. *Nature genetics*, **45**(10), 1113.
- 302 Welch, L., Brooksbank, C., Schwartz, R., Morgan, S. L., Gaeta, B., Kilpatrick, A. M., Mietchen, D.,  
303 Moore, B. L., Mulder, N., Pauley, M., *et al.* (2016). Applying, evaluating and refining bioinformatics  
304 core competencies (an update from the curriculum task force of iscb's education committee). *PLoS*  
305 *computational biology*, **12**(5), e1004943.