

Syllabus for GEN220: High Throughput Biological Data Processing

Course Description

This course focuses on computational skills for processing data using programming language Python and UNIX environment. No prior programming experience is required, but some basic computer skills will be useful.

With the advancement of high throughput data generation methods, a major challenge that graduate students in life sciences have to face today is to analyze large amount of biological data. The objective of this course is to provide an opportunity for graduate students with no computer science background to learn the basic skills of handling high throughput biological data. It covers the Linux/Unix environment and the importance of the command line interface; the Python programming language; program design, implementation, and testing; BioPython; Strategies for analyzing genome resequencing, RNASeq, sequencing data. Students build hands-on skills by analyzing real high throughput biological data through homework assignments and team projects.

Units: 3

Instructor: Jason Stajich (jason.stajich@ucr.edu)

Time and location: W 2:10-4:00PM, F 2:10-3:00PM, ULB104

Office Hours: By Appointment, 1207K Genomics

Prerequisites

- Coursework in genetics or molecular biology or permission of instructor

Resources

None of these texts are required for completion of the course but they will provide a great deal of helpful background and examples that will improve your ability to master UNIX or Programming in Python.

1. *Bioinformatics Data Skills: Reproducible and Robust Research with Open Source Tools*. Vince Buffalo. 2015 O'Reilly & Associates. Available from [O'Reilly and Associates](#), [Amazon](#)
2. *Unix and Perl to the Rescue: A Primer*. Keith Bradnam and Ian Korf. [Unix and Perl Primer for Biologists](#)
3. *Unix and Perl to the rescue!* Bradnam and Korf. [Amazon](#)

4. [Rosalind](#) - An online platform to learn bioinformatics and programming in Python.
5. Software Carpentry - <https://software-carpentry.org/> and Data Carpentry - <http://www.datacarpentry.org/>.
6. Berk Ekmekci, Charles E. McAnany, Cameron Mura. An Introduction to Programming for Bioscientists: A Python-Based Primer. PLoS Comp Bio. DOI: [10.1371/journal.pcbi.1004867](https://doi.org/10.1371/journal.pcbi.1004867)

Grading

- Programming Homework assignments (5 in total): 50%
- Team project: 50%

Homework

- There will be a programming assignment every two weeks. Programming assignments must be prepared along with any necessary input files or documentation to demonstrate program usage.
- Code should be runnable as turned in. You will deposit your code in your github repository or if not possible, by iLearn. You can make one private personal repository to deposit and should organize a folder for each homework assignment (e.g. hw1, hw2, hw3).
- Homework is due BEFORE class on the Wednesdays it is due. The next homework will be posted on Friday at latest.

Projects

- Topics to be selected from a set of choices or the team's choosing (with approval from instructor). Selection of topics will occur by end of October.
- Project teams will be 2-3 individuals working together.
- A presentation will be made by each team - last day(s) of class.
- A final report with the details will be turned in by the group.
- The report needs to detail what each person's contribution is to the project.

Schedule

Date	Day	Lecture Topic	Notes
Sep-29	F	Course Outline / UNIX I: Introduction	
Oct-4	W	UNIX II: Running programs, capturing output	
Oct-6	F	UNIX III: Tools for data processing	
Oct-11	W	Logging into Biocluster. Running Jobs	Guest Lecture
Oct-13	F	Python 1 - Variables, running, cmdline, strings, math	
Oct-18	W	Python 2 - Logic, loops, lists, iterator	Guest Lecture
Oct-20	F	Python 3 - I/O reading/writing files, directories	
Oct-25	W	Python 4 - Dictionaries, Arrays, functions	
Oct-27	F	Python 5 - Libraries, packages, BioPython	
Nov-1	W	Python 6 - Structured data (CSV, XML, GFF, BED)	
Nov-3	F	Data Plotting and R graphics	
Nov-8	W	Bioinformatics 1 - BLAST, cmdline & automation	
Nov-10	F	Veteran's Day Holiday: No class	
Nov-15	W	Bioinformatics 2 - short reads, alignment, stats, SNPs	
Nov-17	F	Bioinformatics 3 - Genome Assembly	
Nov-22	W	<i>TBD - No Class?</i>	
<i>Nov-24</i>	<i>F</i>	<i>Thanksgiving Holiday - no class</i>	
Nov-29	W	Bioinformatics 5 - SNP and variant discovery	
Dec-1	F	Bioinformatics 6 - Phylogenetic trees	
Dec-6	W	TBD	Guest Lecture
Dec-8	F	Presentations	
Exam week	TBD	Additional presentations TBD	