

## Human Genetics 551: Computational Genomics Winter 2018

### Course description

This course provides an introduction to the theory and practice of genomic data analysis including DNA sequence mapping, variant calling, gene expression analysis, and functional genomics. This course will introduce basic programming concepts using Python. A final course project will require applying standard genomics tools to existing data sets. This course is appropriate for students with no background in computer programming.

Students are required to bring their own laptop computers to class. Students should have the ability to install software on their computer.

### Course Objectives

1. Learn the fundamentals of Python programming for manipulating genomics data
2. Develop an understanding of the computational principles underlying common genomics analyses
3. Apply computational tools to address important questions in genomics

### Suggested Text Books and Other Resources

Learning Python by Mark Lutz (O'Reilly publishing). We will focus on Python v 2.x (not 3.x).

Bioinformatics Data Skills by Vince Buffalo

Bioinformatics Programming Using Python by Mitchell Model

<http://www.ee.surrey.ac.uk/Teaching/Unix/index.html>

[https://github.com/blahah/command\\_line\\_bootcamp](https://github.com/blahah/command_line_bootcamp)

<http://www.openbookproject.net/thinkcs/python/english2e/index.html>

<http://rosalind.info/problems/locations/>

<http://opentechschoo.github.io/python-data-intro/>

Additional material including papers and web-resources will be made available using Canvas.

**Start Date:** January 4, 2018

**Time:** 10:10am-11:30am (Tuesdays and Thursdays)

*An optional practicum will also be scheduled based on student availability.*

**Location:** Taubman Library 5000

**Credits:** 3

**Instructors:** Jeff Kidd, Ph.D.                      Jacob Kitzman, Ph.D.

**Office:** 3726 MS II                                      4811 MS II

**Phone:** 734-764-6922                              734-764-9587

**Email:** [jmkidd@umich.edu](mailto:jmkidd@umich.edu)                      [kitzmanj@umich.edu](mailto:kitzmanj@umich.edu)

**Conflict of Interest:** The instructors, Jeff Kidd and Jacob Kitzman, have no industry relationships.

### Grading

Course grades will be based on programming/analysis assignments (65%), a final course project and presentation (25%) and active and thoughtful participation during lectures and discussions (10%).

**Lectures** (Note assignment dates are subject to change -- refer to Canvas for uptodate information):

**1. 1/4 (Th)**

Introduction to Computational Genomics

Motivation for learning computation, programs vs data, manipulating text files, file transfer, and using unix.

Reading:

Sean Eddy: High Throughput Sequencing for Neuroscience

<http://cryptogenomicon.org/2014/11/01/high-throughput-sequencing-for-neuroscience/>

Noble, WS *Plos Comp Biol* 2009

<http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1000424>

Introduction to UNIX commands

<https://swcarpentry.github.io/shell-novice/01-intro/>

<https://swcarpentry.github.io/shell-novice/02-filedir/>

<https://swcarpentry.github.io/shell-novice/03-create/>

<https://swcarpentry.github.io/shell-novice/04-pipefilter/>

Additional Unix tutorial (with built-in interface):

[http://rik.smith-unna.com/command\\_line\\_bootcamp](http://rik.smith-unna.com/command_line_bootcamp)

Unix quick reference sheet:

<https://www.rain.org/~mkummel/unix.html>

**2. 1/9 (Tu)**

Intro to Python I

iPython notebooks, visualization, configuration, interactive vs saved-file computing. Variables, objects, concept of immutability. Numbers and strings

Reading:

*How to Think Like a Computer Scientist: Learning with Python.*

Chapter 1: The way of the program

<http://www.openbookproject.net/thinkcs/python/english2e/ch01.html>

Chapter 2: Variables, expressions, and statements.

<http://www.openbookproject.net/thinkcs/python/english2e/ch02.html>

Assignment:

Python programming assignment #1

**3. 1/11 (Th)**

Computational Topic: Intro to Python II

Flow control: if, elif, else, loops.

Data structures: lists and tuples, and list comprehensions.

Mutable vs immutable variables

Reading:

*How to Think Like a Computer Scientist: Learning with Python.*

Chapter 4: Conditionals

Chapter 6: Iteration

Chapter 7: Strings

Chapter 9: Lists

<http://www.openbookproject.net/thinkcs/python/english2e/>

Optionally also see:

Python Tutorial sections 3.1.4, 4.1-4.4, 5.1-5.3

<https://docs.python.org/2/tutorial/>

Assignment: Python exercise

#### 4. 1/16 (Tu)

Genomics Topic: Storing genomics coordinate data

coordinate intervals, 0-based vs 1-based, UCSC-bed format, bedtools

Computational Topic: Introduction to Unix computing

Reading:

Introduction to UNIX commands

<https://swcarpentry.github.io/shell-novice/01-intro/>

<https://swcarpentry.github.io/shell-novice/02-filedir/>

<https://swcarpentry.github.io/shell-novice/03-create/>

<https://swcarpentry.github.io/shell-novice/04-pipefilter/>

Additional Unix tutorial (with built-in interface):

[http://rik.smith-unna.com/command\\_line\\_bootcamp](http://rik.smith-unna.com/command_line_bootcamp)

Unix quick reference sheet:

<https://www.rain.org/~mkummel/unix.html>

Portions of Bedtools PDF manual ( particularly Chapter 5, posted on Canvas)

Assignment:

Unix and bedtools exercise

#### 5. 1/18 (Th)

Genomics Topic: \*-seq

Assays for variation and function using sequencing as a read-out

Computational Topic: Intro to Python III

loops, iterables, methods, reading and writing files

Reading:

Shendure and Ji, *Nat Biotech*, 2008

<http://www.nature.com/nbt/journal/v26/n10/pdf/nbt1486.pdf>

Shendure and Lieberman Aiden, *Nat Biotech*, 2012

<http://www.nature.com/nbt/journal/v30/n11/pdf/nbt.2421.pdf>

Python file I/O

[http://www.tutorialspoint.com/python/python\\_files\\_io.htm](http://www.tutorialspoint.com/python/python_files_io.htm)

Assignment: Python programming assignment - file and string manipulation

## 6. 1/23 (Tu)

Genomics Topic: Illumina Fastq Files

Sequence representation and quality scores, storing of paired-end reads

Computational Topic: Intro to Python IV

Functions, classes, dictionaries and sets

Reading:

FASTQ Format (Wikipedia)

[https://en.wikipedia.org/wiki/FASTQ\\_format](https://en.wikipedia.org/wiki/FASTQ_format)

Python:

*How to Think Like a Computer Scientist: Learning with Python.*

Chapter 3: Functions

<http://www.openbookproject.net/thinkcs/python/english2e/ch03.html>

Chapter 5: Fruitful functions

<http://www.openbookproject.net/thinkcs/python/english2e/ch05.html>

Chapter 12: Dictionaries

<http://www.openbookproject.net/thinkcs/python/english2e/ch12.html>

See also dictionaries documentation

<https://docs.python.org/2/tutorial/datastructures.html#dictionaries>

Assignment:

fastq file analysis exercise

## 7. 1/25 (Th)

Genomics Topic: Sequence alignment and Dynamic Programming

Local vs global alignments, Burrows-Wheeler transform

Reading:

Eddy, *Nat Biotech*, 2014

<http://www.nature.com/nbt/journal/v22/n7/full/nbt0704-909.html>

Assignment:

No assignment

- 8. 1/30 (Tu)**      Genomics Topic: Significance of sequence matches  
Dynamic programming, continued. Local vs global alignments.
- Computational Topic: Intro to Python V  
Basic python file I/O with tabular data  
Reading in a matrix to a dictionary
- Assignment  
Dynamic programming and sequence alignment (by hand)
- 9. 2/1 (Th)**      Genomics Topic: Assessing Significance  
Basic parametric and nonparametric statistical tests
- Computational Topic: Statistics in Python I  
Array manipulation and numeric computing with numpy/scipy. Random number generation for simulations, and statistical testing in python.
- Reading:  
*Points of Significance* series from *Nature Methods*:  
Importance of being uncertain  
<http://www.nature.com/nmeth/journal/v10/n9/full/nmeth.2613.html>  
and  
Non-parametric tests  
<http://www.nature.com/nmeth/journal/v11/n5/full/nmeth.2937.html>
- The rest of this excellent series is at:  
<http://www.nature.com/collections/qghhqm/pointsofsignificance>
- More lengthy: *A biologist's guide to statistical thinking and analysis*  
<http://www.ncbi.nlm.nih.gov/books/NBK153593/>
- 10. 2/6 (Tu)**      Computational Topic I: Statistics in Python (continued)
- Computational Topic II: Python data visualization I  
Matplotlib – basic plot construction
- Assignment:  
Python statistics and plotting assignment
- 11. 2/8 (Th)**      Introduction to Cluster Computing  
FLUX, process management, qsub
- Assignment:  
Cluster computing exercise

12. 2/13 (Tu) Genomics Topic: The multiple testing problem  
p-values, False discovery rate, q-values.
- Reading:  
Storey and Tibshirani, *PNAS*, 2003  
<http://www.pnas.org/content/100/16/9440.full>  
Online plotting reference
13. 2/15 (Th) Genomics Topic: Sequence Analysis I  
Properties of Illumina data, introduction to mapping reads and the BAM format
- Computational Topic: BWA and samtools  
Manipulating BAM files using samtools and pysam
- Reading:  
Li et al *Bioinformatics*, 2009  
<http://bioinformatics.oxfordjournals.org/content/25/16/2078.long>
- SAM/BAM specification: <http://samtools.github.io/hts-specs/SAMv1.pdf>
14. 2/20 (Tu) Genomics Topic: Sequence Analysis II  
Theory and implementation of short read alignment, including Burrows-Wheeler read mapping.
- Computational Topic: Picard tools  
Read mapping with bwa on FLUX
- Reading:  
Trapnell and Salzberg, *Nature Biotech*, 2009  
<http://www.nature.com/nbt/journal/v27/n5/abs/nbt0509-455.html>
15. 2/22 (Th) Genomics Topic: Sequence Analysis III  
Variant calling, statistical models for joint vs individual calling. Considering different scenarios, somatic vs germline
- Computational Topic: GATK and samtools  
Variant calling models
- Reading:  
Li, *Bioinformatics*, 2011  
<http://bioinformatics.oxfordjournals.org/content/27/21/2987.long>
- DePristo et al., *Nature Genetics*, 2011  
<http://www.nature.com/ng/journal/v43/n5/full/ng.806.html>
- 2/27 (Tu) **No class – “Spring” Break**  
3/1 (Th) **No class – “Spring” Break**

**16. 3/6 (Tu)**Genomics Topic: Sequence Analysis IV

Selecting variants, hard filters vs variant quality calibration., VCF format.

## Reading:

DePristo et al., *Nature Genetics*, 2011

<http://www.nature.com/ng/journal/v43/n5/full/ng.806.html>

Danecek et al. *Bioinformatics*, 2011

<http://bioinformatics.oxfordjournals.org/content/27/15/2156>

VCF format specification:

<https://samtools.github.io/hts-specs/VCFv4.2.pdf>

**17. 3/8 (Th)**Genomics Topic: RNA-seq (I)

Experimental strategies for expression profiling using sequencing. Algorithmic aspects of RNA read mapping, isoform detection, and quantification at the gene, isoform, and exon levels. Normalization.

Computational Topic: RNA-seq read mapping

Alignment vs alignment-free methods.

## Reading:

Oshlack et al., *Genome Biol*, 2010

<https://genomebiology.biomedcentral.com/articles/10.1186/gb-2010-11-12-2>

[20](#)

Conesa et al., *Genome Biol*, 2016

<https://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-088>

[1-8](#)

**18. 3/13 (Tu)**Genomics Topic: RNA-Seq II

Measures of abundance (FPKM, TPM). Expression quantitative trait locus (QTL) statistical framework for studies of genetic variants' impact upon expression.

Computational Topic: Bowtie, STAR, and alignment-free methods

## Reading:

Lior Pachter, "Estimating number of transcripts from RNA-Seq measurements (and why I believe in paywall)"

<https://liorpachter.wordpress.com/2014/04/30/estimating-number-of-transcripts-from-rna-seq-measurements-and-why-i-believe-in-paywall/>

Trapnell et al. *Nature Biotech*, 2013

<http://www.nature.com/nbt/journal/v31/n1/full/nbt.2450.html>

## Assignment:

## eQTL analysis exercise

- 19. 3/15 (Th)**      Python Data Visualization II - Advanced Topics (Final Project Teams Due)  
Pandas and Advanced Plotting  
Reading:  
    To be determined
- 20. 3/20 (Tu)**      Clustering (Ryan Mills, Ph.D, guest-lecture) (Final Project Ideas Due)  
Algorithms and applications for grouping data based on similarity
- 21. 3/22 (Th)**      Genomics Topic: Functional Genomics I  
Chromatin profiling data, chip-seq, Peak calling  
Reading:  
    To be determined
- 22. 3/27 (Tu)**      Genomics Topic: Functional Genomics II  
Chromatin conformation assays, concepts and processing  
Reading:  
    To be determined
- 23. 3/29 (Th)**      HMMs and Segmentation  
Reading:  
    To be determined  
Assignment:  
    Work on final projects
- 24. 4/3 (Tu)**      Machine Learning and Statistical Inference  
SVMs, Decision Tree/Random Forest, Neural Nets  
ROC curves, Cross validation, and the “overfitting” problem  
Reading:  
    To be determined  
Assignment:  
    Work on final projects
- 25. 4/5 (Th)**      Version Control, gitHub, and Reproducible Analyses  
Reading:  
    To be determined  
Assignment:  
    Work on final projects



26. 4/10 (Tu) Final Project Lab
27. 4/12 (Th) Final Project Presentations
28. 4/17 (Tu) Final Project Presentations