

M21-550 Course Syllabus

Fall 2019

Course Title:	Introduction to Bioinformatics										
Course Master:	C. Charles Gu, Ph.D.										
Teaching Assistant:	Xiaofu (Sissi) Dong, James Perez, Will Yang, Ph.D.										
Schedule:	Tue & Thu, 2:00-3:30pm.										
Place:	MSIBS classroom & Computer Lab (Room 502 & 501, 5 th floor, Becker's Library)										
Format:	Small group lecture and extensive computer labs using real-world data										
Grade Criteria:	Numerical score based on: <table><tr><td>Quizzes</td><td>20%</td></tr><tr><td>Lab assignments</td><td>40%</td></tr><tr><td>Exam-I & II</td><td>10% & 10%</td></tr><tr><td>Final Project</td><td>20%</td></tr><tr><td>Classroom participation</td><td>up to extra 5%</td></tr></table>	Quizzes	20%	Lab assignments	40%	Exam-I & II	10% & 10%	Final Project	20%	Classroom participation	up to extra 5%
Quizzes	20%										
Lab assignments	40%										
Exam-I & II	10% & 10%										
Final Project	20%										
Classroom participation	up to extra 5%										

Reference Text:

1. Andreas Baxevanis and Francis Ouellett (eds.), *Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins*, 3rd Edition, Wiley & Sons, 2005.
2. Robert Gentleman, Vincent Carey, Wolfgang Huber, Rafael Irizarry, Sandrine Dudoit (eds.), *Bioinformatics and Computational Biology Solutions Using R and Bioconductor*, Springer, 2005.
3. Malcolm Campbell and Laurie Heyer, *Discovering Genomics, Proteomics and Bioinformatics* (2nd Edition), Benjamin Cummings, 2007.
4. Richard Durbin, Sean Eddy, Anders Krogh, and Graeme Mitchison, *Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*, Cambridge University Press, 1998.

Other Readings:

Additional reading materials and papers will be assigned.

Description: This course provides a general introduction to basic concepts, methodologies and applications of bioinformatics. Students will learn online bioinformatics databases & data mining tools, and acquire understanding of mathematical/computational algorithms in classic sequence analysis (alignment, gene finding, and hidden Markov models), next-generation sequence analysis (short-read data format and processing, variant calling algorithms), and in selected topics of high-dimensional analysis in the context of mining high-throughput biomedical experiments data such as gene expression microarrays (data QC & normalization, univariate & multivariate differential expression analysis). Students will become familiar with popular bioinformatics software, online tools, and specialized R/BioConductor packages. We will discuss methods for high-dimensional data analysis including classification and clustering analysis, principal component analysis (PCA), statistical/machine learning, and Bayesian inference. There also will be seasonal additional lectures on topics such as proteomics and applications of bioinformatics to real studies of complex diseases.

As an important component of this course, students will conduct hands-on computer labs to learn basics of online bioinformatics databases and tools, and to practice computer programming. The labs require using the statistical computing environment R (i.e., the R primer is required) though introduction to BioConductor basics will be provided. Students will use specialized software and R packages to accomplish tasks including designing experiments, low-level analysis of expression levels, univariate differential expression analysis, and various multivariate analysis techniques taught in class. A variety of software will be used for NGS data analysis covering alignment, variants calls, differential analysis, and visualization of results. Through the lab exercises, students will learn how state-of-the-art computational tools are applied to solving bioinformatics problems in real studies of human diseases.

Schedule of Lectures:

	Tuesday	Thursday
WK-01 8/27,29	Course overview: summary of topics, grading policy, Genomics Primer: genes & genome, heredity & Mendelian inheritance, mutations & disease, omics & bioinformatics	QZ-2; Lab-01: R refresher; install <i>Bioconductor</i> ; allele counts & frequency (<i>genetics</i>), transcription & translation etc. (<i>Biostrings, seqinr</i>)
WK-02 9/3,5	Sequence analysis I: homology detection; pair-wise sequence alignment, alignment score statistics, BLAST search	QZ-2; Lab-02: online databases & tools, BLAST
WK-03 9/10,12	Sequence Analysis II: multiple alignment, Hidden Markov Models (HMM) & applications to motif and gene finding	QZ-3; Lab-03: seq alignment, database/R tools; HMM motif/gene prediction
WK-04 9/17,19	DNA variants: genetic variations & diseases; DNA variants, copy number variations (CNVs), CNVs detection & analysis; impact of SVs	QZ-4; Lab-04: R tools for CNV (<i>DNAcopy, ReadDepth</i> , etc.)
WK-05 9/24,26	NGS data analysis I: next generation sequencing (NGS), sequencing technology & platforms, short reads alignment algorithms & software	QZ-5; Lab-05: BWA, Bowtie; R tools for exploring sequencing data
WK-06 10/1,3	NGS data analysis II: RNA-seq based expression analysis, count-based normalization, transcript-level differential expression	QZ-6; Lab-06: R tools for RNA-seq analysis (<i>edgeR, DESeq2</i>); TopHat, Cufflinks
WK-07 10/8,10	NGS data analysis III: NGS file formats, SAM/BAM; SAMtools; variants calling, VCF format; QC of NGS reads; NGS study design and statistical issues	QZ-7; Lab-07: SAMtools, more R tools for exploring sequencing data
WK-08 10/15,17	Fall Break - No Class	Exam-I
WK-09 10/22,24	High-dimensional analysis I: microarray & transcriptomics; high-dimensional analysis, bias & normalization; robust statistics; differential expression, false discovery rate, resampling statistics	QZ-8; Lab-08: microarray data preprocessing (<i>ReadAffy</i>); differential expression analysis (<i>siggenes</i>);
WK-10 10/29,31	High-dimensional analysis II: multivariate analysis & dimension reduction, PCA (principal component analysis); dissimilarity measures, clustering analysis (hierarchical, k-means)	QZ-9; Lab-09: R tools for PCA (<i>prcomp</i>), clustering analysis (<i>heatmap, hclust</i>)
WK-11 11/5,7	High-dimensional analysis III: statistical learning (supervised/unsupervised), model/feature selection; classification analysis, LDA (linear discriminant analysis), Naïve Bayes classifier	QZ-10; Lab-10: R tools for LDA (<i>lda</i>), naive Bayes classification (<i>naiveBayes</i>)
WK-12 11/12,14	High-dimensional analysis IV: advanced machine learning, support vector machine, artificial neural networks, deep learning in genomics	QZ-11; Lab-11: R tools for advanced ML
WK-13 11/19,21	Bioinformatics in medicine: genetic circuits, environments & interactions; EHR & bigdata; precision medicine & public health initiatives	QZ-12; Lab-12: R tools for network & systems analysis (<i>igraph, etc.</i>)
WK-14 11/26,28	Review	Thanksgiving Break - no class
WK-15 12/3,5	Exam-II	Student Presentation Practice
WK-12 12/12,14	Student Presentation I	Student Presentation II