# Lecture 0: Course Overview

Instructor:

Parisa Rashidi

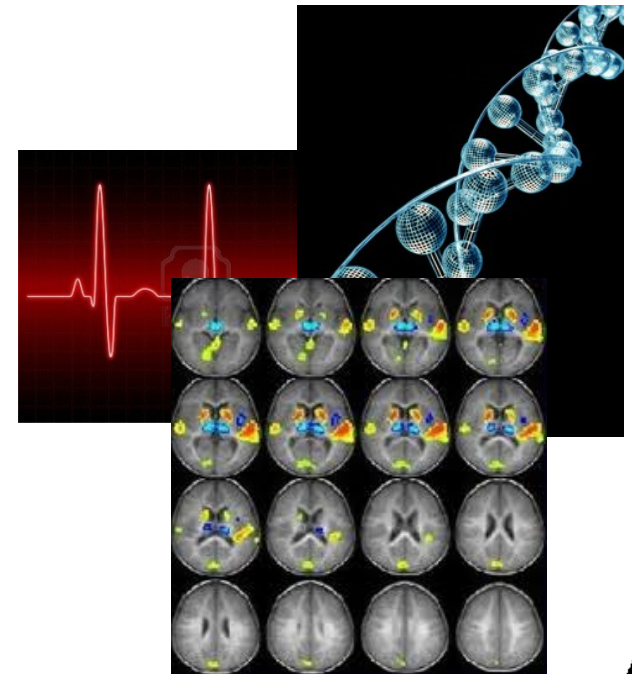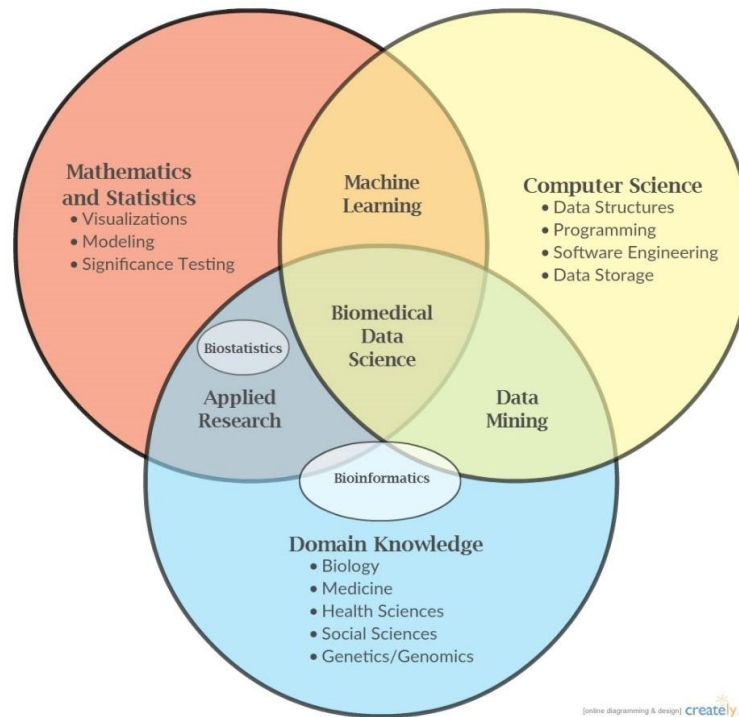FALL 2019

## BIOMEDICAL DATA SCIENCE

# Agenda

- Short Introduction
  - What you will learn in this course
- Logistics
  - What is expected of you in this course

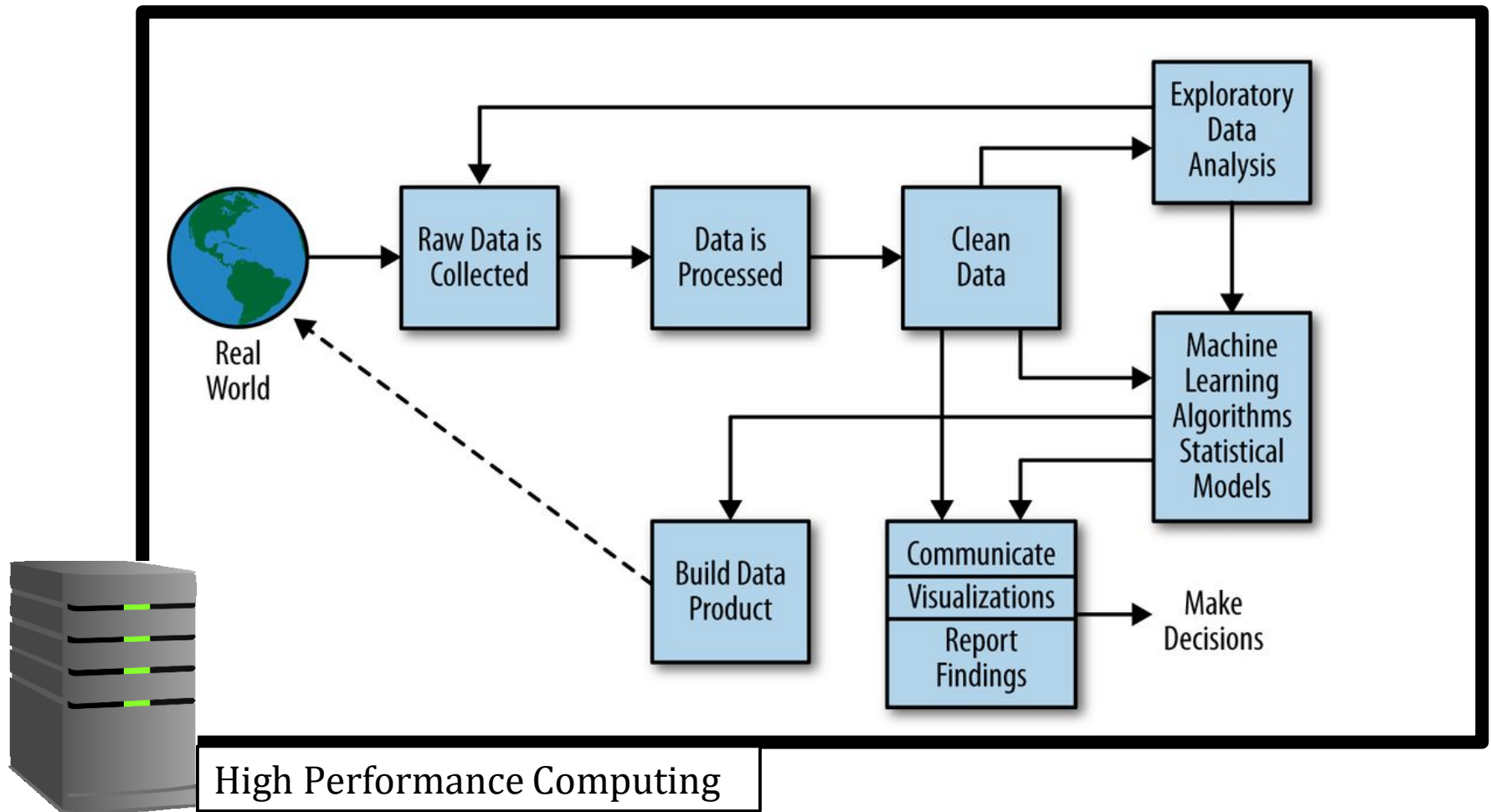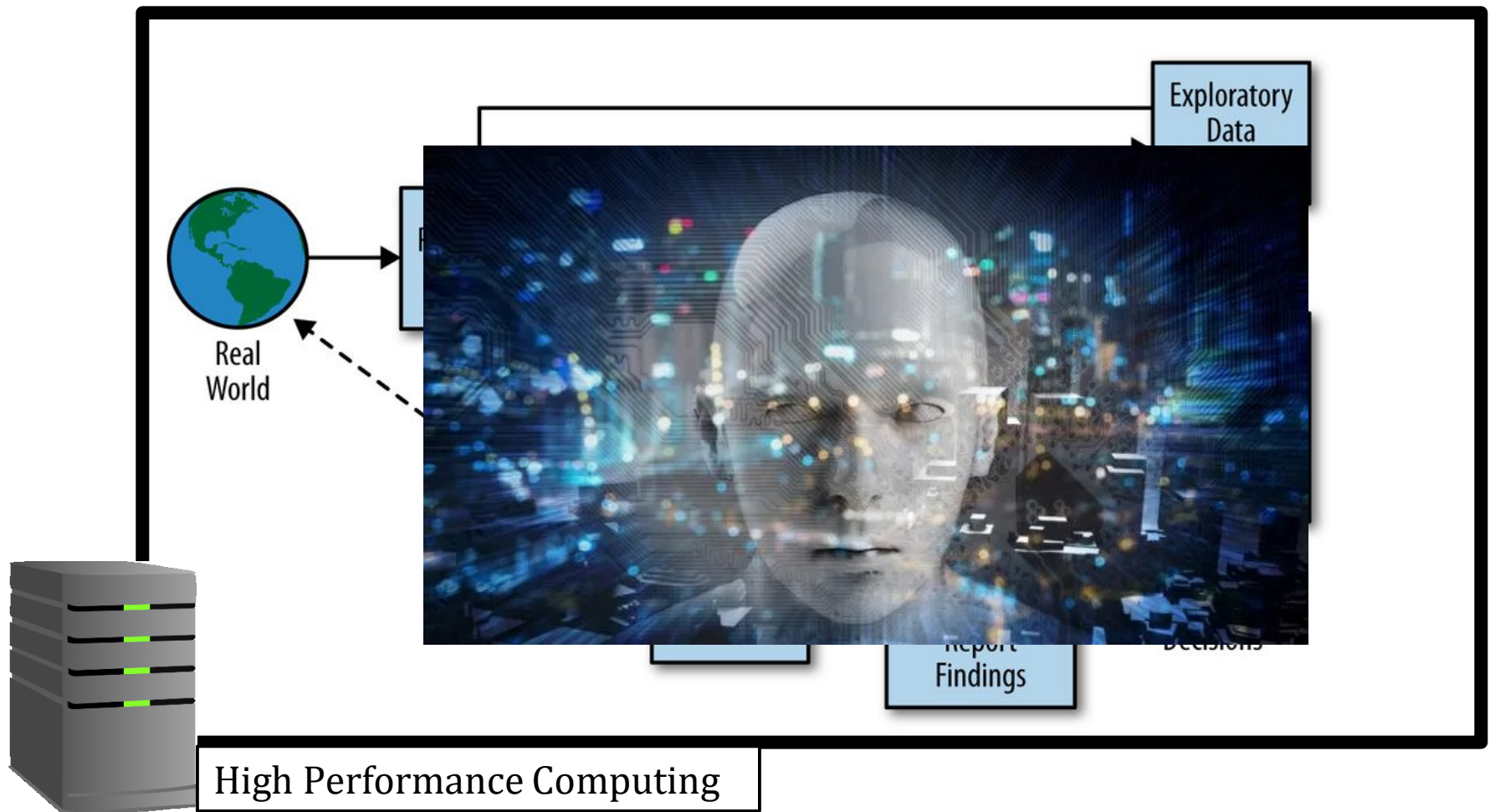# Introduction

# Biomedical Data Science

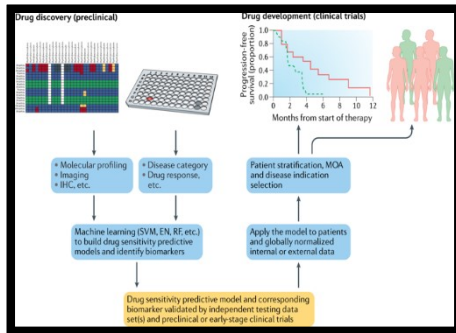- Data science techniques applied to biomedical science problems

# What is Data Science?



High Performance Computing

# Beware of the AI Hype



High Performance Computing

# Why Biomedical Data Science?

# Why Study Data Science?
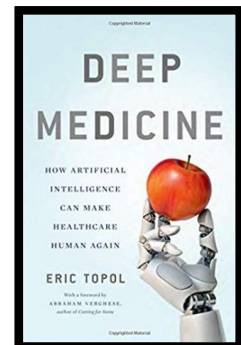


**ML FOR DRUG DISCOVERY**
*(NATURE REVIEWS-DRUG DISCOVERY, 2019)*

**MODELLING TRANSCRIPTION FACTOR BINDING SITES**
*(NATURE REVIEW -GENTICS, 2019)*

Deep Learning at the Cellular Level

"Image-Free" Microscopy
E. Christiansen, *Cell*, 2018

"Ghost Cytometry"
S. Ota, *Science* 2018
N. Nitta, *Cell*, 2018

Subcellular Machine Vision
G. Gut, *Science* 3 August 2018

**AI.GOOGLE/HEALTHCARE**

# Current Healthcare Themes (1): Electronic Health Records (EHR)

- **Moving beyond paper-based records**
  - Regional and national health data integration facilitating
    - Patient care
    - Administrative & financial
    - Research
    - Scholarly information
    - Office automation

# Current Healthcare Themes (2): Smart and Connect Health

- Smart and connected health (Quantified Self)
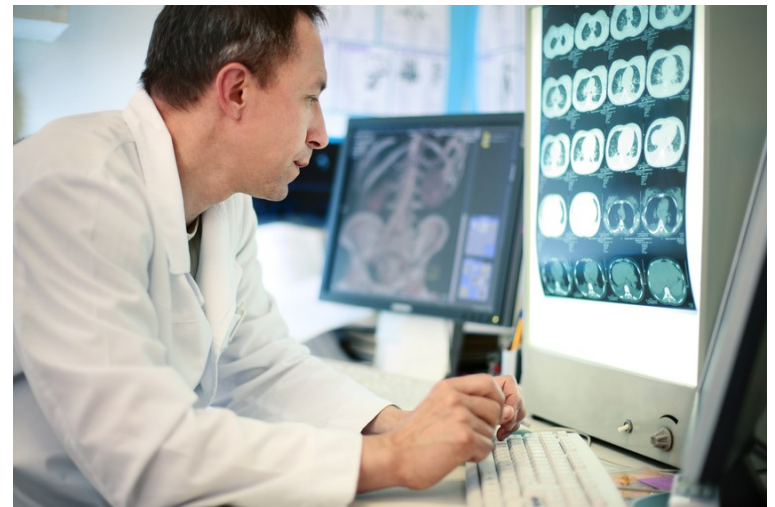


Polar: embedded with heart rate-monitoring sensors, a motion-tracking sensor to track speed, distance and acceleration



Bodytrak: a pair of earbuds equipped with an in-ear thermometer to measure core body temperature



Flow, a smart air quality tracker

# Current Healthcare Themes (3): Omics, Imaging, …

- Complex heterogeneous types of data
- Past: single experiment, small set of results
- Today: using sophisticated instruments, we can generate very large datasets, e.g. gene sequencing.

# Current Healthcare Themes (4): Information Access

- **Information access**
  - Health-related searches among the most popular
    - Most Web information are anecdotal

# Industry & Healthcare

# Healthcare Costs

- 2019: $3.8 trillion, $10K+ per person
- 1960: $27.2 billion, $146 per person

Top 10 most valuable companies combined

Net worth of

$72 billion        $58 billion

Yan Liu, Jimeng Sun, Deep Learning Models for Health Care - Challenges and Solutions, 2017

# Healthcare is Broken

## How the U.S. Can Reduce Waste in Health Care Spending by $1 Trillion

POTENTIAL SAVINGS (IN BILLIONS OF DOLLARS)

$600 — Innovations to eliminate clinical waste, fraud and abuse, administrative complexity, and excessive prices

Additional broad system interventions (e.g., reducing administrative complexity to levels in other service industries)

$130 — Incremental savings from comprehensive demand-side and aggressive supply-side reforms
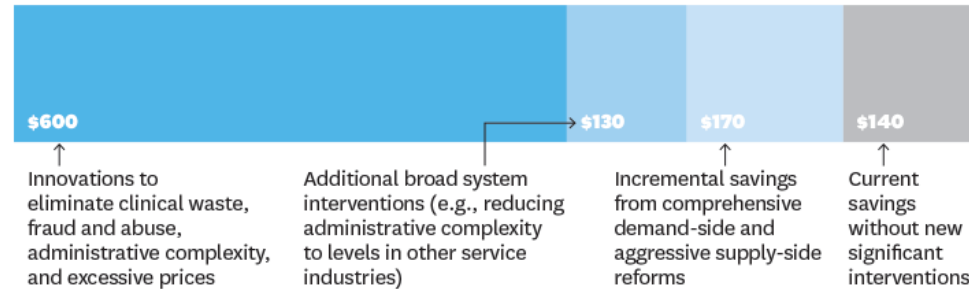
$170

$140 — Current savings without new significant interventions
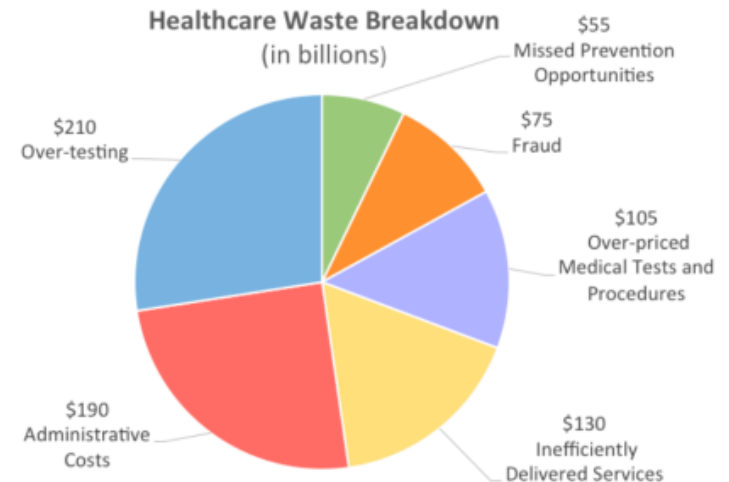
SOURCE  ANALYSIS BY NIKHIL SAHNI ET AL.; "ELIMINATING WASTE IN U.S. HEALTH CARE"
BY DONALD M. BERWICK AND ANDREW D. HACKBARTH, 2012

© HBR.ORG

**$765 billions** = 50 years budget NASA

### Healthcare Waste Breakdown (in billions)

- $55 Missed Prevention Opportunities
- $75 Fraud
- $105 Over-priced Medical Tests and Procedures
- $130 Inefficiently Delivered Services
- $190 Administrative Costs
- $210 Over-testing

Graph by MMS Analytics, Inc.
Data courtesy of the Institue of Medicine

# Data Revolution

- All the data processing we did in the last 2 years is more than all the data processing we did in the last three thousand years.

- We are now being exposed to as much information in a single day as our 15th century ancestors were exposed to in their entire lifetime.

- Every two days the human race is now generating as much data as were generated from the dawn of humanity through the year 2003.

Image credit: The IT staffing company

# Healthcare data sources



Weber, Griffin M., Kenneth D. Mandl, and Isaac S. Kohane. 2014. "Finding the Missing Link for Big Biomedical Data." *JAMA: The Journal of the American Medical Association* 311 (24): 2479–80.

# Clinical Data

- ## Adoption of Electronic Health record (EHR) data

   % non-federal acute care hospitals with EHR adoption



2008     2011     2015

NR   0-19%   20-39%   40-59%   60-79%   80-100%

MarketScan Databases

Inpatient Outpatient View

The IBM® Family of MarketScan® Research Databases is the largest of its kind in the industry, with data on over 245 million unique patients since 1995.

Example: MarketScan dataset 245 million patients

# Physiological Data



If you use MIMIC data or code in your work, please cite the following publication:

*MIMIC-III, a freely accessible critical care database. Johnson AEW, Pollard TJ, Shen L, Lehman L, Feng M, Ghassemi M, Moody B, Szolovits P, Celi LA, and Mark RG. Scientific Data (2016). DOI: 10.1038/sdata.2016.35. Available from: http://www.nature.com/articles/sdata201635*

# MoleculeNet

- A large set of dataset useful for molecular machine learning

Wu, Zhenqin, Bharath Ramsundar, Evan N. Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S. Pappu, Karl Leswing, and Vijay Pande. "MoleculeNet: a benchmark for molecular machine learning." Chemical science 9, no. 2 (2018): 513-530.

# Many Other Sources of Data


Biomedical Images


Wearables


Genomics


Social Media


Lab Tests

# Course Logistics

# Main Topics Covered

- It can be subject to some change:
  - Python Programming
  - Biomedical Ontologies (e.g. ICD-9, SNOMED, ...)
  - Conventional Machine Learning Techniques
    - kNN, SVM, RF
  - Deep Learning
    - Images, Time series, Text data
  - Useful machine learning and deep learning libraries
    - Scikit-learn, Pandas, Keras, DeepChem
  - Introduction to Natural Language Processing (NLP)

| Week | Date | Tuesday Lecture | Thursday Lecture | Code Examples | Paper Reading | HW |
|------|------|-----------------|------------------|---------------|---------------|-----|
| 1 | 08/19 - 08/23 | Course Overview | Data Sci. Introduction, Python Programming | • Intro. To Python<br>• Setting up Google Colab | - | - |
| 2 | 08/26- 08/30 (Add/Drop) | EHR, NumPy | ML Metrics | • NumPy Intro<br>• Toxicity Prediction<br>• Solubility Prediction | Watson Oncology | |
| 3 | 09/02 - 09/06 | Pandas, Preprocessing | Version Control | • Pandas Intro<br>• Preprocessing (5)<br>• Setting up Git | | HW1 |
| 4 | 09/09 - 09/13 | RF, kNN | Biophysical Modeling | • kNN<br>• RF<br>• Biophysical Modeling | Septic Shock Prediction | |
| 5 | 09/16 - 09/20 | Evaluation | SVM | • Evaluation | | HW2 |
| 6 | 09/23 - 09/27 | Neural Network | Neural Network, Deep Learning | • FC Networks | Melanoma Detection | |
| 7 | 09/30 - 10/04 | Review | Exam 1 | | | |
| 8 | 10/07 - 10/11 | ConvNet | Convnet (cont.) | • CNN | Adversarial Attacks | |
| 9 | 10/14 – 10/18 | CNN Architectures | Microscopy | • Cell Counting<br>• Cell Segmentation | | HW3 |
| 10 | 10/21 – 10/25 | RNN | RNN | • Physiological Signals (MIMIC) | Drug Discovery | |
| 11 | 10/28 – 11/01 | Intro to NLP | NLP Representations | • Clinical Notes | | HW4 |
| 12 | 11/04 – 11/08 | Generative Models | Genomics | • Generating Drug Molecules<br>• TF Binding<br>• RNA Interference<br>• Chromatic Accessibility | | |
| 13 | 11/11 – 11/15 | Feature Selection | Guest Lecture | | | HW5 |
| 14 | 11/18 – 11/22 | Review | Exam 2 | | | |
| 15 | 11/25 – 11/29 | Ethical AI | Holiday | | | |
| 16 | 12/02 – 12/06 | | Reading Day | | | |

# Course Objectives

- At the end of the semester, students are expected to be able to:
  - Understand what biomedical data science is,
  - Identify different techniques used to solve biomedical data science problems,
  - Identify when and why a certain library or platform should be used,
  - Demonstrate the ability to apply methods from each of the major domains to solve practical problems.

# Course Website

- Everything will be available on Canvas
- Canvas should be considered as the reference
  - All announcements, project postings, schedule changes, ..
  - Check your email!
  - Upload your photo

# Class

- We will meet on the following days:
  - Tuesdays
    - One session: 1:55 – 2:45 pm
  - Thursdays
    - First session: 1:55 - 2:45 pm
    - 15 minutes break
    - Second session: 3:00 – 3:50 pm

# Office Hours

- Office location: NEB 459
- Office hours:  by appointment
- E-mail address: parisa.rashidi@ufl.edu
- <span style="color:red">NOTE</span>: When contacting by email include "<span style="color:red">Course – 4760/6938</span>" in the subject line to ensure delivery.

# Supervised Teaching Student

- Subhash Nerella, [subhashnerella@ufl.edu](mailto:subhashnerella@ufl.edu),
- office location: TBA, office hours: TBA

- Programming/HW questions should be directed to Subhash.

# Textbook

- *Recommended*, not required
- Your main source: Lecture Notes,
  - You still have to take your own notes!

# Undergraduate Grading

- Final grade is calculated according to:
  - Homework (5) = 35%
    - Equal weight
  - Exams (2) = 45%
    - Equal weight
  - Paper Discussions (5) = 10%
    - Undergraduate students are not asked to present in class, but are expected to contribute to discussions in class and write up a paragraph on each paper (strengths, weakness, suggestions).
  - Quiz (5) = 10%
    - Equal weight

# Graduate Grading

- Final grade is calculated according to:
  - Homework (5) = 35%
    - Equal weight
  - Exams (2) = 40%
    - Equal weight
  - Paper Presentation (1) = 10%
    - Graduate students are expected to present in class
  - Paper Discussions (5) = 5%
    - Graduate students are also expected to contribute to discussions in class and write up a paragraph on each paper (strengths, weakness, suggestions).
  - Quiz (5) = 10%
    - Equal weight

# Exam

- Exam 1: Thursday, Oct. 3
- held during regular class hours


- Exam 2: Thursday, November 21
- held during regular class hours

# Quick Poll

- Programming experience
    - Python
    - Matlab
    - R
    - Other
    - No experience
- Prior machine learning experience

- Introduce Yourself
  - Your name
  - Your major
  - Your research topic if a graduate student
  - Why you enrolled in this class