# ECE 365: Data Science and Engineering
## Spring 2019
https://courses.engr.illinois.edu/ece365

**Instructors**: Venu Veeravalli, Subhonmesh Bose and Olgica Milenkovic.

**Course Coordinator**: Venu Veeravalli

**Prerequisites**: ECE 313 (or campus equivalent on basic undergrad probability) and some basic linear algebra. General mathematical maturity expected of engineering undergraduates.

**Textbook:** None. Relevant course notes will be handed out to the students.

**Target Audience**: Juniors or Seniors

**Outline**: Big Data is all around us. Petabytes of data is collected by Google and Facebook. 24 hours of video is uploaded on Youtube every minute. Making sense of all this data in the relevant context is a critical question. This course takes a holistic view towards understanding how this data is collected, represented and stored, retrieved and computed/analyzed upon to finally arrive at appropriate outcomes for the underlying context. The course is divided into three parts, with the first part focusing on foundations of machine learning, and the remaining two on specific application areas. Each application topic is covered at four discrete levels.

- We start with the context of where the data comes from, how it is acquired, what are the biases and noise levels in the data leading to statistical and physical models of the data acquired.
  Appropriate data representation mechanisms and distributed storage and computing architectures are discussed next. Based on the type of the data, different compression/ coding methods are appropriate. Images, videos, genomic data, medical imaging data, smart grid data, each bring their own unique characteristics which can be harnessed towards efficient representation.
- Once data is stored and represented efficiently, we look for the right statistical and algorithmic tools to analyze the data. Spectral methods (including Fourier methods and PCA), Clustering algorithms, SVM, Mining algorithms are studied in the specific context of the data.
- Finally, the analyzed data leads to appropriate inferences or visualizations as appropriate to the physical problem we started out with. This closes the loop bringing utility to the original setting and context in which the data was acquired.

For Spring 2019 the application areas will be:

- *Machine learning for power systems:* Grid operation relies on efficient processing of data and identifying patterns in them. In this module, we explore applications of machine learning in grid operations. Specifically, we explore regression and classification tasks such as those that arise in load prediction, consumer electricity usage, recognizing valid power system measurements, and virtual bidding markets.

- *Biological Data Analytics:* It may be argued that biology and medical sciences are the two disciplines with the fastest growing datasets and data repositories. It is nowadays common to refer to data being of genomic, rather than astronomic size. What is known as –Omics data gives invaluable information about the structure and composition of our genomes, our unique genetic markers, the communication activity between genes and

other molecules, the structure of our building block proteins and many other health related issues. In this part of the course we will cover diverse topics of relevance in bioinformatics, ranging from de Bruijn  graphs (used to stitch DNA sequence fragments produced by experiments into a complete DNA sequence) to suffix trees (used for efficient data representation) and community detection (used to identify cancer gene communities). You will also get acquainted with modern biological data acquisition technologies, data libraries and publicly available data processing software.

**Course Plan**

**Part 1 (Weeks 1-5): Foundations of Machine Learning**

**Lecture 1**: Introduction to the course; Review of Linear Algebra and Probability
**Lecture 2**: k-Nearest Neighbor Classifiers and Bayes Classifiers
**Lecture 3**: Linear Classifiers and Linear Discriminant Analysis
**Lecture 4**: Naïve Bayes, Kernel Tricks
**Lecture 5**: Logistic Regression, SVM and Model Selection
**Lecture 6**: K-Means Clustering and Applications
**Lecture 7**: Linear Regression and Applications
**Lecture 8**: SVD and Eigen-Decomposition
**Lecture 9**: Principal Component Analysis
**Lecture 10**: Optimization Techniques for Machine Learning, Q&A

**Labs (Weeks 1-5)**
Lab 1: Introduction to Python and the Canopy environment
Lab 2: Linear Classification: k-NN and LDA
Lab 3: Linear Classification: SVM
Lab 4: Clustering and Linear Regression
Lab 5: Eigen-Decompositions, SVD and PCA

**Grading**: 30% pre-lab quizzes (in class), 70% labs and lab reports.

**Part 2 (Weeks 6-10): Smart Grid**

**Lecture 1**: Introduction to power systems, basics of neural networks
**Lecture 2**: Neural networks and load prediction
**Lecture 3**: Power flow equations
**Lecture 4**: SVM for detecting corrupt power system measurements
**Lecture 5**: Detecting network structure
**Lecture 6**: Basics of electricity markets, virtual bidding
**Lecture 7**: Trading strategies for virtual bidding
**Lecture 8**: Wrapping up virtual bidding, understand customer data
**Lecture 9**: Logistic regression for customer data analysis
**Lecture 10**: Customer billing and cost savings from solar

**Labs**
Lab 1: Day-ahead load prediction in ERCOT markets
Lab 2: Detecting bad sensors in power system measurements
Lab 3: Virtual bidding in NYISO's markets
Lab 4: Analyze customer data from Austin, Texas.

**Grading**: 30% pre-lab quizzes (in class), 70% labs and lab reports

**Part 3 (Weeks 11-15): Biological Data Analytics**

**Lecture 1**: Introduction to bioinformatics. Biological data.
**Lecture 2**: Sequence alignment. Global vs local alignment. Dynamic programming.
**Lecture 3**: The Smith-Waterman and Needlman-Wunsch algorithms. BLAST.

**Lecture 4**. Suffix trees and the Burrows-Wheeler transform. Bowtie2.

**Lecture 5**: Dynamic programming for sequence folding prediction. Vienna and Mfold. Stochastic grammars for folding models.

**Lecture 6**: Sanger sequencing. Overview of Next Generation and Third Generation Sequencing technologies.

**Lecture 7**: Basics of graph theory. Genome assembly via de Bruijn Graphs. EULER and IDBA_UD.

**Lecture 8**: Statistical read error-correction for Illumina, PacBio and Oxford Nanopore sequencers. Quake.

**Lecture 9**: Biological data repositories and databases.

**Lecture 10:** Biological data compression. Reference-based compression. CRAM. Context-tree weighting.

**Labs**

Lab 1: Sequence alignment and applications of BLAST.

Lab 2: Bowtie and DNA forensics.

Lab 3: Genome assembly. Influence of sequencing errors on assembler accuracy.

Lab 4: -Omics data compression.

Lab 5: Genomic sequence amplification and primer selection.

**Grading**: 30% pre-lab quizzes (in class), 70% labs and lab reports.