

Syllabus for:  
Python for Machine Learning and Data Analysis  
Spring 2019

Time and place: Thursday 6:00pm – 9:00pm      TBD  
Instructor: Handan Liu      [h.liu@northeastern.edu](mailto:h.liu@northeastern.edu)  
Office hours:  
TA:

Course Rationale:

This course we will cover is broadly applicable, and has led to significant advances in many fields. Once you understand the basics of the technology of machine learning and data analysis, and the close connection between theory and practice, it's a very open field, where lots of progress can be made quickly. More important, you will transform your theoretical knowledge into practical skill using many hands-on labs in this course.

Prerequisites:

You should be comfortable programming in Python. Basic knowledge of Unix/Linux will be helpful.

Grading:

Homework:	40%	(HW1 10%, HW2 15% and HW3 15%)
Project:	30%	
Attendance:	10%	
Final exam:	20%	

Assignments:

Lectures are complemented by homeworks (programming assignments) to bridge the theory with the practice. The homeworks are associated with the three main parts of the course that mostly consist of programming assignments to exercise a technology or programming model.

Homework is due electronically. All assignments have a specific due date and time.

Submissions will be accepted up to one day after the deadline with a 50% penalty. For example, an on-time submission might receive a grade of 80 points. The same assignment submitted after the deadline would receive 40 points ( $80 \times 0.5$ ).

Submission will be done through Blackboard.

Make-up Policy

Students who miss the presentation and the final exam will not, as a matter of course, be able to make up it. If there is a legitimate reason why a student will not be able to complete an assignment on time or not be present for the presentation or the exam, then they should contact the instructor beforehand. Under extreme circumstances, as decided on a case-by-case basis by the instructor, students may be allowed to make up assignments or exam without first informing the instructor.

#### Course Schedule:

Week	Topics	Assignments
1	<ul style="list-style-type: none"> <li>• Course/Syllabus</li> <li>• Python Infrastructure and Development Tools: Anaconda, an enhanced interactive science-centric console; and the Jupyter Notebook, a web-based application that mixes code, plots, and rich media, making it ideal for sharing and publishing analyses with peers, etc.</li> <li>• Anaconda, conda and pip</li> <li>• <u>Lab 1</u>: install anaconda <ul style="list-style-type: none"> <li>○ Custom environment settings</li> <li>○ Jupyter Notebook and JupyterLab</li> <li>○ Script editor for Python</li> </ul> </li> </ul>	Install Software on your computer
2	<p>Python Language Essentials: Introduction to Python's core language features and packages: Python's built-in data structures, including how and where each might be used and what trade-offs are present, and cover Python's looping and control flow constructs, etc.</p> <ul style="list-style-type: none"> <li>• Fundamental data types and data structures; Organizing code with functions, modules and packages; Loading packages, namespaces; Reading and writing data; Control flow</li> <li>• <u>Lab 2</u>: samples - writing python code and run it</li> </ul>	
3	<p>Numerical Analysis and Data Exploration with <b>NumPy</b> Arrays: NumPy is a critical tool for rapidly manipulating and processing large data sets.</p> <ul style="list-style-type: none"> <li>• The NumPy array; Selecting data using slicing and logical indexing; Efficient numerical processing with multi-dimensional arrays; Expressive array operations and manipulations; Access larger-than-RAM data using memory mapped arrays</li> <li>• <u>Lab 3</u>: run samples</li> </ul>	

4	<p>Data Visualization with <b>Matplotlib</b></p> <ul style="list-style-type: none"> <li>• 2D plotting with Matplotlib: line plots, scatter plots, histograms, labeling, and more.</li> <li>• <u>Lab 4</u>: run samples</li> </ul>	HW1
5	<p><b>Pandas</b>, the Python Data Analysis Library is a powerful package for working with tabular data for data aggregation and reorganization, such as:</p> <ul style="list-style-type: none"> <li>• Loading from CSV and other structured text formats; Accessing data stored in SQL databases; 1D and 2D data structures; Stripping out extraneous information; Normalizing data; Dealing with missing data; Data manipulation (alignment, aggregation, and summarization); Group-based operations: split-apply-combine; Statistical analysis; Date and time series analysis with Pandas; Visualizing data</li> <li>• <u>Lab 5</u>: run samples</li> </ul>	
6	<p><u>Lab 6</u>: Hand-on Lab:</p> <ul style="list-style-type: none"> <li>• Setup computers to be able to access Discovery Cluster</li> <li>• Learn how to submit and manage jobs on the HPC cluster</li> <li>• Write our first program and run on HPC cluster</li> </ul>	
7	<p>Parallel programming in Python:</p> <ul style="list-style-type: none"> <li>• Python multiple threads: <b>multiprocessing</b></li> <li>• Python MPI: <b>mpi4py</b></li> </ul> <p><u>Lab 7</u>: run python <b>parallel</b> code on HPC cluster</p>	HW2
8	<p>Introduction to Machine Learning:</p> <p>This section starts with a short conceptual introduction to machine learning and explain how it works, and what kinds of problems it's best suited to solve. And the section covers the frameworks and tools provided by <b>scikit-learn</b>, a widely used library for machine learning, such as:</p> <ul style="list-style-type: none"> <li>• Linear and nonlinear models; Constant and variable learning-rates; Cost functions, regularization methods, and other constraints; Fitting, transforming, and predicting</li> <li>• <u>Lab 8</u>: run samples</li> </ul>	
9		

10	<b>Machine Learning: Numeric Data:</b> <ul style="list-style-type: none"> <li>Logarithmic and curvilinear transforms; Data scaling; Outliers; Linear regressors; l1 and l2 normalization; Support vector machines (SVM)</li> </ul> <b>Machine Learning: Categorical Data:</b> <ul style="list-style-type: none"> <li>Contrast encoding; Missing values; Categorical rebinning; Linear classifiers; Tree-based classifiers; Ensemble methods; Boosting methods; Unbalanced designs</li> </ul>	HW3
11	<b>Machine Learning: Image Data:</b> <ul style="list-style-type: none"> <li>Image storage formats; Scikit-image; Smoothing and denoising; Edge detection; Feature-based segmentation; K-means clustering</li> </ul> <u>Lab 9</u> : run samples	Project
12	Brief Introduction to the Python Deep Learning Library <b>TensorFlow and Keras</b> ; <u>Lab 10</u> : run python code of ML and DL on HPC cluster for parallel on CPU and GPU	
13	Final exam; debrief	
14	Presentation of the project by every student	

Note: This schedule and contents will be adjusted as needed throughout the semester.

#### Resources:

Required Textbook: None

Other Material:

- Python: <https://www.python.org/>
- Anaconda: <https://www.anaconda.com/>
- Numpy: <http://www.numpy.org/>
- Matplotlib: <https://matplotlib.org/>
- Pandas: <https://pandas.pydata.org/>
- Scikit-learn: <https://scikit-learn.org/>
- TensorFlow: <https://www.tensorflow.org/>
- Keras: <https://keras.io/>

Other information will be completed later.