# STA 221: LECTURE 1

KRISHNA BALASUBRAMANIAN

(UNIVERSITY OF CALIFORNIA, DAVIS)

- ▷ Instructor: Krishna Balasubramanian

- ▷ Office: Zoom

- ▷ Office Hours: Monday 4:30-5:30PM

- ▷ Email: kbala@ucdavis.edu

- ▷ TA: Si Teng Hao

▷ Learning Objecives:

   ▷ Learning to use Python for big-data analytics

   ▷ Learning popular Machine Learning Techniques

▷ Basics of python ?

▷ Basic statistics and machine learning ?

▷ Basics of Linear Algebra ?

Rate yourself out of 10 (10 being expert) in each of the above.

▷ 3 homeworks: 60% of total points

    ▷ Each homework is worth 20%.

▷ Final project: 40% of total points

    ▷ May 4th and 6th proposal (individual meetings)

    ▷ June 2nd and 3rd demonstration (individual meetings)

    ▷ Final report due June 10th.

▷ Follow instructions for each homework and project carefully!

▷ Late submission **WILL NOT** be accepted.

▷ Follow UC Davis Code of Academic Conduct carefully!

▷ Violations **WILL NOT** be tolerated!

▷ Two master algorithms:

  ▷ Power method for computing eigenvectors

  ▷ Gradient descent for optimization

  ▷ Note: Both methods are kind of related.

▷ Bonus algorithm:

  ▷ Randomized matrix multiplication - for learning/testing basics of Python

▷ Computing eigenvectors: Unsupervised learning. We will use scikit learn python package.

▷ Gradient descent: Supervised learning. We will use scikit learn and Pytorch.

▷ Given large amounts of unlabelled data:

▷ How to group/cluster them accordingly ?

▷ How to visualize them for exploratory data analysis ?

▷ Some techniques (all based on eigenvalue computation):

▷ k-means, mixture models, spectral clustering, page rank.

▷ PCA, Kernel PCA, Manifold Learning

▷ Predict/Classify new data based on learning from labelled training data:

   ▷ How to efficiently learn from labeled data ?

▷ Techniques (all based on gradient descent (or variations)):

   ▷ Discriminant analysis, Logistic Regression

   ▷ Support vector machine, Deep neural networks

**NO FREE LUNCH !**

# Lecture 1: Python Basics

**Growth of major programming languages**

Based on Stack Overflow question views in World Bank high-income countries

**Python compared to smaller, growing technologies**

Based on question traffic in World Bank high-income countries

14

**Traffic by industry to Python**

Comparing Jan-Aug of each year, in the United States and United Kingdom.



15

IEEE spectrum

KDnuggets Analytics, Data Science, Machine Learning Software Poll, top tools share, 2015-2017