

# Introduction to Machine Learning and Python

Shrey Gupta

*Applied Machine Learning* (HOUSECS 59-01), Duke University

August 29, 2018

# Topics

- ▶ **Classification:** naive Bayes, support vector machines, kernel methods, and neural networks.
- ▶ **Regression:** spline interpolation and linear and polynomial regression.
- ▶ **Unsupervised learning:** mixture of Gaussians clustering.
- ▶ **Computer vision:** object detection via convolutional neural networks, feature extraction, edge detection, and processing methods.
- ▶ **Dimensionality reduction:** generalized discriminant analysis.
- ▶ **Evaluation of machine learning models.**

# Prerequisites

- ▶ This is an *applied* course, but requires fundamental understanding of the algorithms and techniques being used.
- ▶ Prerequisites: basic fluency in programming and mathematics at the single-variable calculus level.

# Grading

- ▶ **Attendance:** required to attend at least 11 (of 14) classes. 11 classes are lectures, and three are office hours.
- ▶ **Programming assignments:** required to complete all three assignments, with individual scores of 70% or greater.
- ▶ **Final paper:** 1500-word write-up detailing machine learning model from final programming assignment.
- ▶ Class is graded on a *satisfactory/unsatisfactory* basis.
- ▶ See syllabus for more detailed information.

# What is Machine Learning?

- ▶ “Give computers the ability to learn without being explicitly programmed” (Arthur Samuel).
- ▶ Formal problem specification: “A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$ , and performance measure  $P$  if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ ” (Tom Mitchell).

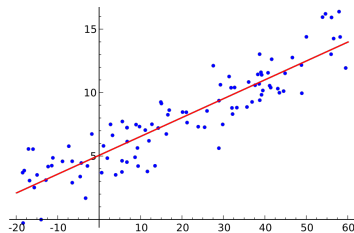
# Terminology

- ▶ Grouped into two categories: *supervised* and *unsupervised* learning.
- ▶ Input data comes in *features* that describe each data point.
- ▶ *Training data*: data used to train algorithm (i.e. create model).
- ▶ *Testing data*: data used to evaluate performance of algorithm.

# Supervised Learning

- ▶ Data (a subset from a larger distribution) is labeled, and we attempt to generalize to (predict) the larger distribution.
- ▶ Regression: predicts a continuous value output.
- ▶ Classification: predicts a discrete class output.

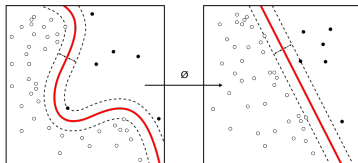
# Regression: Examples



- ▶ Given data about square footage, age, zip code, and housing demand, predict the selling price for a house.
- ▶ Predict the percentage increase or decrease in the price of an equity.



# Classification: Examples



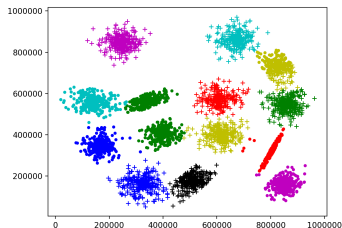
- ▶ Given data about temperature, humidity, and wind speed, predict whether it will be sunny, cloudy, or raining.
- ▶ Predict whether the price of an equity will increase or decrease.

*Image source: Wikipedia*

# Unsupervised Learning

- ▶ Data is unlabeled (no “ground truth”).
- ▶ Problems: clustering, density estimation, and pattern detection.

# Clustering: Examples



- ▶ Given consumption data, partition the consumers into market segments.
- ▶ Given several news articles (and their text), group them based on similarity.

# Python

- ▶ We'll be using *Python 3.x* throughout the course.
- ▶ Libraries and frameworks: NumPy, SciPy, Pandas, Matplotlib, SciKit, TensorFlow, NLTK, and Jupyter.
- ▶ Today's notebook will ensure all of the packages are downloaded and work through an introduction of Python.