# Introduction to Machine Learning (CSCI-UA.0480-007)

## David Sontag

## New York University

Slides adapted from Luke Zettlemoyer, Pedro Domingos, and Carlos Guestrin

# Logistics

- **Class webpage:**
  - http://cs.nyu.edu/~dsontag/courses/ml16/
  - Sign up for Piazza!
- **Office hours:** TBD
- **Teaching assistant**:
  Kevin Jiao <jjiao@stern.nyu.edu>
- **Graders:**
  - Yijun Xiao <ryjxiao@nyu.edu>
  - Alexandre Sablayrolles
    <alexandre.sablayrolles@gmail.com>

# Evaluation

- 6-7 homeworks (50%)
  - Both theory and programming
  - Collaboration policy:
    - First try to solve the problems on your own
    - Then, can discuss with other classmates
    - Write-up solutions on your own
    - List names of anyone you talked to
- Midterm exam (25%)
- Project (20%)
- Course participation (5%)

# Projects

- Be creative – think of new problems that you can tackle using machine learning
  - Scope: ~40 hours/person

- Logistics:
  - 2-3 students per group
  - Begins mid-March. Project proposal due week after midterm exam
  - Will still be problem sets during this period!

# Prerequisites

**REQUIRED:**

- **Basic algorithms** (CS 310)
  - Dynamic programming, algorithmic analysis
  - *Can be taken concurrently*

**STRONGLY RECOMMENDED:**

- **Linear algebra** (Math 140)
  - Matrices, vectors, systems of linear equations
  - Eigenvectors, matrix rank
  - Singular value decomposition
- **Multivariable calculus** (Math 123)
  - Derivatives, integration, tangent planes
  - Optimization, Lagrange multipliers
- **Good programming skills:** Python highly recommended

# Source Materials

**No textbook required. Readings will come from freely available online material.**

If you really want a book for an additional reference, these are OK options:

• C. Bishop, ***Pattern Recognition and Machine Learning***, Springer, 2007

• K. Murphy, ***Machine Learning: a Probabilistic Perspective***, MIT Press, 2012

• … may update this list throughout semester. I wouldn't buy anything yet.

# What is Machine Learning ?
## (by examples)

# Classification

## from data to discrete classes

# Spam filtering

**data**                                     **prediction**



⟶ Spam
vs.
Not Spam

# Face recognition



Example training images
for each orientation

# Weather prediction

# Regression

## predicting a numeric value

# Stock market

# Weather prediction revisited

# Ranking

**comparing items**

# Web search

# Given image, find similar images



http://www.tiltomo.com/

# Collaborative Filtering

# Recommendation systems

# Recommendation systems

Machine learning competition with a $1 million prize

# Clustering

**discovering structure in data**

# Clustering Data: Group similar things

# Clustering images

Set of Images



$C_1$

$C_2$

$C_3$

$C_4$

$C_5$

[Goldberger et al.]

# Clustering web search results

# Embedding

**visualizing data**

# Embedding images

- Images have thousands or millions of pixels.

- Can we give each image a coordinate, such that similar images are near each other?

[Saul & Roweis '03]

# Embedding words



[Joseph Turian]

# Embedding words (zoom in)



[Joseph Turian]

# Structured prediction

## from data to discrete classes

# Speech recognition

# Natural language processing



I need to hide a body
noun, verb, preposition, …

# Growth of Machine Learning

- Machine learning is preferred approach to
    - Speech recognition, Natural language processing
    - Computer vision
    - Medical outcomes analysis
    - Robot control
    - Computational biology
    - Sensor networks
    - …
- This trend is accelerating
    - Big data
    - Improved machine learning algorithms
    - Faster computers
    - Good open-source software

# Course roadmap

- **First half of course: supervised learning**
  - SVMs, kernel methods
  - Learning theory
  - Decision trees, boosting, deep learning

- **Second half of course: data science**
  - Unsupervised learning, EM algorithm
  - Dimensionality reduction
  - Topic models

# Supervised Learning: find $f$

- Given: Training set $\{(x_i, y_i) \mid i = 1 \ldots N\}$
- Find: A good approximation to $f : X \rightarrow Y$

Examples: what are $X$ and $Y$ ?

- Spam Detection
  - Map email to {Spam, Not Spam}
- Digit recognition
  - Map pixels to {0,1,2,3,4,5,6,7,8,9}
- Stock Prediction
  - Map new, historic prices, etc. to $\Re$ (the real numbers)

# A Supervised Learning Problem

Dataset:

| Example | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|---------|-------|-------|-------|-------|-----|
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

- Our goal is to find a function $f : X \to Y$
  - $X = \{0,1\}^4$
  - $Y = \{0,1\}$

- Question 1: How should we pick the *hypothesis space*, the set of possible functions $f$?

- Question 2: How do we find the best $f$ in the hypothesis space?

# Most General Hypothesis Space

Consider all possible boolean functions over four input features!

- $2^{16}$ possible hypotheses

- $2^9$ are consistent with our dataset

- How do we choose the best one?

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|-------|-------|-------|-------|-----|
| 0 | 0 | 0 | 0 | ? |
| 0 | 0 | 0 | 1 | ? |
| 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | ? |
| 1 | 0 | 0 | 0 | ? |
| 1 | 0 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | ? |
| 1 | 0 | 1 | 1 | ? |
| 1 | 1 | 0 | 0 | 0 |
| 1 | 1 | 0 | 1 | ? |
| 1 | 1 | 1 | 0 | ? |
| 1 | 1 | 1 | 1 | ? |

Dataset:

| Example | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|---------|-------|-------|-------|-------|-----|
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

# A Restricted Hypothesis Space

Consider all conjunctive boolean functions.

- 16 possible hypotheses

- None are consistent with our dataset

- How do we choose the best one?

| Rule | Counterexample |
|------|----------------|
| $\Rightarrow y$ | 1 |
| $x_1 \Rightarrow y$ | 3 |
| $x_2 \Rightarrow y$ | 2 |
| $x_3 \Rightarrow y$ | 1 |
| $x_4 \Rightarrow y$ | 7 |
| $x_1 \wedge x_2 \Rightarrow y$ | 3 |
| $x_1 \wedge x_3 \Rightarrow y$ | 3 |
| $x_1 \wedge x_4 \Rightarrow y$ | 3 |
| $x_2 \wedge x_3 \Rightarrow y$ | 3 |
| $x_2 \wedge x_4 \Rightarrow y$ | 3 |
| $x_3 \wedge x_4 \Rightarrow y$ | 4 |
| $x_1 \wedge x_2 \wedge x_3 \Rightarrow y$ | 3 |
| $x_1 \wedge x_2 \wedge x_4 \Rightarrow y$ | 3 |
| $x_1 \wedge x_3 \wedge x_4 \Rightarrow y$ | 3 |
| $x_2 \wedge x_3 \wedge x_4 \Rightarrow y$ | 3 |
| $x_1 \wedge x_2 \wedge x_3 \wedge x_4 \Rightarrow y$ | 3 |

Dataset:

| Example | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|---------|-------|-------|-------|-------|-----|
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

# Occam's Razor Principle

- William of Occam: Monk living in the 14th century
- Principle of parsimony:

"One should not increase, beyond what is necessary, the number of entities required to explain anything"

- When many solutions are available for a given problem, we should select the simplest one
- But what do we mean by simple?
- We will use prior knowledge of the problem to solve to define what is a simple solution

*Example of a prior: smoothness*

[Samy Bengio]

# Key Issues in Machine Learning

- How do we choose a hypothesis space?
  - Often we use **prior knowledge** to guide this choice
- How can we gauge the accuracy of a hypothesis on unseen data?
  - **Occam's razor:** use the *simplest* hypothesis consistent with data! This will help us avoid overfitting.
  - ***Learning theory*** will help us quantify our ability to **generalize** as a function of the amount of training data and the hypothesis space
- How do we find the best hypothesis?
  - This is an **algorithmic** question, the main topic of computer science
- How to model applications as machine learning problems? (engineering challenge)