# Introduction to Machine Learning (CSCI-UA.0480-002)

## David Sontag

## New York University

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

# Logistics

- **Class webpage:**
  - http://cs.nyu.edu/~dsontag/courses/ml12/
  - Sign up for mailing list!

- **Office hours:**
  - Tuesday 3:30-4:30pm and by appointment.
  - 715 Broadway, 12th floor, Room 1204

- **Grader:** Jinglun Dong
  - Email: jinglundong@gmail.com

# Evaluation

- About 7 homeworks (50%)
  - Both theory and programming
  - See collaboration policy on class webpage

- Midterm & final exam (45%)

- Course participation (5%)

# Prerequisites

- **Basic algorithms** (CS 310)
  - Dynamic programming, algorithmic analysis
- **Linear algebra** (Math 140)
  - Matrices, vectors, systems of linear equations
  - Eigenvectors, matrix rank
  - Singular value decomposition
- **Multivariable calculus** (Math 123)
  - Derivatives, integration, tangent planes
  - Lagrange multipliers
- **Probability** (Math 233 or 235)

# Source Materials

Optional textbooks:

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

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

# A Few Quotes

- "A breakthrough in machine learning would be worth ten Microsofts" (Bill Gates, Chairman, Microsoft)

- "Machine learning is the next Internet" (Tony Tether, former director, DARPA)

- "Machine learning is the hot new thing" (John Hennessy, President, Stanford)

- "Web rankings today are mostly a matter of machine learning" (Prabhakar Raghavan, former Dir. Research, Yahoo)

- "Machine learning is going to result in a real revolution" (Greg Papadopoulos, former CTO, Sun)

- "Machine learning is today's discontinuity" (Jerry Yang, former CEO, Yahoo)

# What is Machine Learning ?
# (by examples)

# Classification

## from data to discrete classes

# Spam filtering

☆ **Osman Khan** to Carlos                     show details Jan 7 (6 days ago)   ↩ Reply | ▼

sounds good
+ok

Carlos Guestrin wrote:
> Let's try to chat on Friday a little to coordinate and more on Sunday in person?
>
> Carlos

**Welcome to New Media Installation: Art that Learns**

☆ **Carlos Guestrin** to 10615-announce, Osman, Michel  show details 3:15 PM (8 hours ago)  ↩ Reply | ▼

Hi everyone,

Welcome to New Media Installation:Art that Learns

The class will start tomorrow.
***Make sure you attend the first class, even if you are on the Wait List.***
The classes are held in Doherty Hall C316, and will be Tue, Thu 01:30-4:20 PM.

By now, you should be subscribed to our course mailing list: 10615-announce@cs.cmu.edu.
You can contact the instructors by emailing: 10615-instructors@cs.cmu.edu

**Natural _LoseWeight SuperFood Endorsed by Oprah Winfrey, Free Trial 1 bottle,
pay only $5.95 for shipping mfw rlk**  Spam | X

☆ **Jaquelyn Halley** to nherrlein, bcc: thehorney, bcc: ang  show details 9:52 PM (1 hour ago)  ↩ Reply | ▼

=== Natural WeightL0SS Solution ===

Vital Acai is a natural WeightL0SS product that Enables people to lose wieght and cleansing their bodies
faster than most other products on the market.

Here are some of the benefits of Vital Acai that You might not be aware of. These benefits have helped
people who have been using Vital Acai daily to Achieve goals and reach new heights in there dieting that
they never thought they could.

* Rapid WeightL0SS
* Increased metabolism - BurnFat & calories easily!
* Better Mood and Attitude
* More Self Confidence
* Cleanse and Detoxify Your Body
* Much More Energy
* BetterSexLife
* A Natural Colon Cleanse

## Spam
## vs.
## Not Spam

# Object detection



Example training images
for each orientation

# Weather prediction

# Regression

**predicting a numeric value**

# Stock market

# Weather prediction revisted



Temperature

72° F

# 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



[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
  - Improved machine learning algorithms
  - Improved data capture, networking, faster computers
  - Software too complex to write by hand
  - New sensors / IO devices
  - Demand for self-customization to user, environment

# Supervised Learning: find $f$

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

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)

# Example: Spam Filter

- Input: email
- Output: spam/ham
- Setup:
  - Get a large collection of example emails, each labeled "spam" or "ham"
  - Note: someone has to hand label all this data!
  - Want to learn to predict labels of new, future emails

- Features: The attributes used to make the ham / spam decision
  - Words: FREE!
  - Text Patterns: $dd, CAPS
  - Non-text: SenderInContacts
  - …

❌
> Dear Sir.
>
> First, I must solicit your confidence in this transaction, this is by virture of its nature as being utterly confidencial and top secret. …

❌
> TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.
>
> 99  MILLION EMAIL ADDRESSES
>   FOR ONLY $99

✅
> Ok, Iknow this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

# Example: Digit Recognition

- Input: images / pixel grids

- Output: a digit 0-9

- Setup:
  - Get a large collection of example images, each labeled with a digit
  - Note: someone has to hand label all this data!
  - Want to learn to predict labels of new, future digit images

- Features: The attributes used to make the digit decision
  - Pixels: (6,8)=ON
  - Shape Patterns: NumComponents, AspectRatio, NumLoops
  - …

0

1

2

1

??

# Important Concepts

- Data: labeled instances, e.g. emails marked spam/ham
    - Training set
    - Held out set (sometimes call Validation set)
    - Test set

- Features: attribute-value pairs which characterize each x

- Experimentation cycle
    - Select a hypothesis $f$ to best match training set
    - (Tune hyperparameters on held-out or *validation* set)
    - Compute accuracy of test set
    - Very important: never "peek" at the test set!

- Evaluation
    - Accuracy: fraction of instances predicted correctly

- Overfitting and generalization
    - Want a classifier which does well on *test* data
    - Overfitting: fitting the training data very closely, but not generalizing well
    - We'll investigate overfitting and generalization formally in a few lectures

Training Data

Held-Out Data

Test Data

# A Supervised Learning Problem

- Consider a simple, Boolean dataset:
  - $f : X \rightarrow Y$
  - $X = \{0,1\}^4$
  - $Y = \{0,1\}$

- 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?

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 |

# 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 |