


CSE 446 Machine Learning

Daniel Weld
Xiao Ling
Congle Zhang

1




What is Machine Learning ?

©2005-2009 Carlos Guestrin 2


Machine Learning

Study of algorithms that

- improve their performance
- at some task
- with experience

Data →  → Understanding

©2005-2009 Carlos Guestrin 3

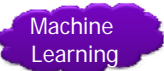


Why?

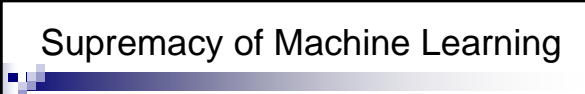
Is this topic important?

4

Exponential Growth in Data

Data →  → Understanding

©2005-2009 Carlos Guestrin 5



Supremacy of Machine Learning

- Machine learning is preferred approach to
 - Speech recognition, Natural language processing
 - Web search – result ranking
 - 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

©2005-2009 Carlos Guestrin 6

Logistics

7

Syllabus

- Covers a wide range of Machine Learning techniques – from basic to state-of-the-art
- You will learn about the methods you heard about:
 - Naive Bayes, logistic regression, nearest-neighbor, decision trees, boosting, neural nets, overfitting, regularization, dimensionality reduction, error bounds, loss function, VC dimension, SVMs, kernels, margin bounds, K-means, EM, mixture models, semi-supervised learning, HMMs, graphical models, active learning
- Covers algorithms, theory and applications
- **It's going to be fun and hard work 😊**

©2005-2009 Carlos Guestrin, D. Weld

8

Prerequisites

- Probabilities
 - Distributions, densities, marginalization...
- Basic statistics
 - Moments, typical distributions, regression...
- Algorithms
 - Dynamic programming, basic data structures, complexity...
- Programming
 - Mostly your choice of language, but Python (NumPy) + Matlab will be very useful
- We provide some background, but the class will be fast paced
- Ability to deal with "abstract mathematical concepts"

©2005-2009 Carlos Guestrin

9

Staff

- Two Great TAs:
Fantastic resource for learning, interact with them!
 - **Xiao Ling**, CSE 610, xiaoling@cs
 - Office hours: TBA
 - **Congle Zhang**, CSE 524, clzhang@cs
 - Office hours: TBA
- Administrative Assistant
 - Alicen Smith, CSE 546, asmith@cs



10

Text Books

- Required Text:
 - Pattern Recognition and Machine Learning; Chris Bishop
- Optional:
 - Machine Learning; Tom Mitchell
 - The Elements of Statistical Learning: Data Mining, Inference, and Prediction; Trevor Hastie, Robert Tibshirani, Jerome Friedman
 - Information Theory, Inference, and Learning Algorithms; David MacKay
 - Website: Andrew Ng's AI class videos
 - Website: Tom Mitchell's AI class videos

11

Grading

- 4 homeworks (55%)
 - First one goes out Fri 1/6/12
 - Start early, Start early, Start early, Start early, Start early, Start early, Start early, Start early, Start early, Start early
- Midterm (15%)
 - Circa Feb 10 in class
- Final (30%)
 - TBD by registrar

12

Homeworks

- Homeworks are hard, start early ☺
- Due at the beginning of class
 - Minus 33% credit for each day (or part of day) late
- All homeworks **must be handed in**, even for zero credit
- Collaboration
 - You may **discuss** the questions
 - Each student writes their own answers
 - Write on your homework anyone with whom you collaborate
 - Each student must write their own code for the programming part
 - Please don't search for answers on the web, Google, previous years' homeworks, etc.**
 - Ask us if you are not sure if you can use a particular reference

Communication

- Main discussion board
 - <https://catalyst.uw.edu/gqpost/board/xling/25219/>
- Urgent announcements
 - cse446@cs
 - Subscribe: <http://mailman.cs.washington.edu/mailman/listinfo/cse446>
- To email instructors, always use:
 - cse446_instructor@cs

Space of ML Problems

Type of Supervision (eg, Experience, Feedback)

What is Being Learned?

	Labeled Examples	Reward	Nothing
Discrete Function	Classification		Clustering
Continuous Function	Regression		
Policy	Apprenticeship Learning	Reinforcement Learning	

Classification

from data to discrete classes

LAST UPDATED: 12:54 AM ET

Top Candidate		Bichelle Bachmann	David Geppert	Jim Heftman	Ron Paul	Rick Perry	Mitt Romney	Rick Santorum
Sex	Male	4%	14%	1%	24%	10%	23%	23%
	Female	8%	12%	0%	19%	10%	25%	27%
Describes self as born again or evangelical Christian	Yes	6%	14%	1%	19%	14%	14%	32%
	No	3%	12%	1%	26%	8%	38%	14%
Tea Party	Support	6%	15%	0%	19%	11%	19%	29%
	Neutral	2%	10%	1%	28%	9%	32%	17%
	Oppose	3%	9%	3%	21%	8%	43%	13%
Age	17-29	2%	0%	1%	48%	8%	13%	23%
	30-44	4%	13%	0%	26%	6%	20%	30%
	45-64	6%	14%	1%	18%	11%	20%	26%
	65 and older	6%	17%	0%	11%	13%	33%	20%
Income	Less than \$30,000	6%	13%	0%	37%	14%	15%	13%
	\$30,000 - \$49,999	7%	12%	1%	26%	14%	18%	24%
	\$50,000 - \$99,999	5%	13%	1%	21%	9%	21%	29%
	\$100,000 or more	3%	15%	0%	14%	7%	36%	24%

Spam filtering

data → prediction

Example email content:

David Geppert to Carlos Guzman (2009) Jan 7 (8 hours ago) - No Reply - 1

Hi Carlos,

Carlos Guzman writes:

Let's try to find out if there is a link to contribute and meet on Sunday in person?

Carlos Guzman

Welcome to New Media Installation: Art that Learns

Carlos Guzman to 10015-announce, Carlos Guzman (2009) Jan 9 (8 hours ago) - No Reply - 1

Hi everyone,

Welcome to New Media Installation: Art that Learns

Our goal is to build a machine learning system that can learn to recognize faces and objects in images.

There are some of the benefits of this tool that you might not be aware of. These benefits have helped people who have been using this tool to learn more about their own faces and how they change over time.

----- Natural Weight Loss Solution -----

Use Anus as a natural Weight Loss product that enables people to lose weight and clean up their bodies from their most annoying problems in the world.

Here are some of the benefits of this tool that you might not be aware of. These benefits have helped people who have been using this tool to learn more about their own faces and how they change over time.

----- Natural Weight Loss Solution -----

Use Anus as a natural Weight Loss product that enables people to lose weight and clean up their bodies from their most annoying problems in the world.

Here are some of the benefits of this tool that you might not be aware of. These benefits have helped people who have been using this tool to learn more about their own faces and how they change over time.

Text classification



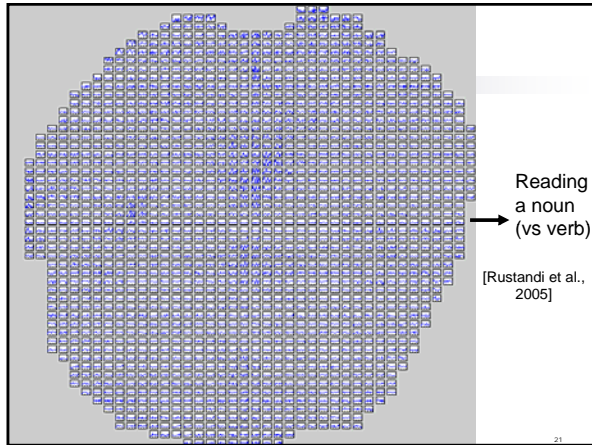
- Company home page
- vs
- Personal home page
- vs
- Univeristy home page
- vs
- ...

Object detection

(Prof. H. Schneiderman)



Example training images for each orientation



Reading a noun (vs verb)

[Rustandi et al., 2005]

Weather prediction



The classification pipeline

Training



Testing



Regression

predicting a numeric value

Stock market



©2009 Carlos Guestrin

25

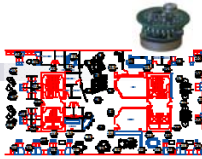
Weather prediction revisited



©2009 Carlos Guestrin

26

Modeling sensor data



- Measure temperatures at some locations
- Predict temperatures throughout the environment

[Guestrin et al. '04]

©2009 Carlos Guestrin

27

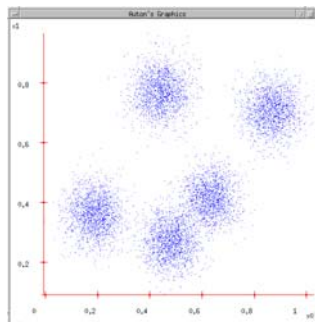
Clustering

discovering structure in data

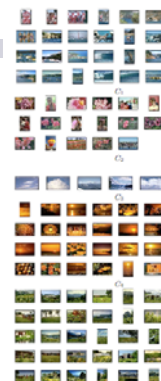
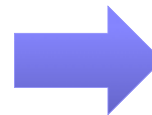
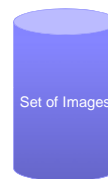
©2009 Carlos Guestrin

28

Clustering Data: Group similar things



Clustering images



©2009 Carlos Guestrin

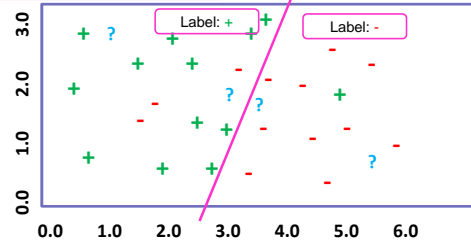
[Goldberger et al.]₃₀

Key Concepts

38

Classifier

Hypothesis:
Function for labeling examples



Generalization

- Hypotheses must **generalize** to correctly classify instances not in the training data.
- Simply memorizing training examples is a consistent hypothesis **that does not generalize**.

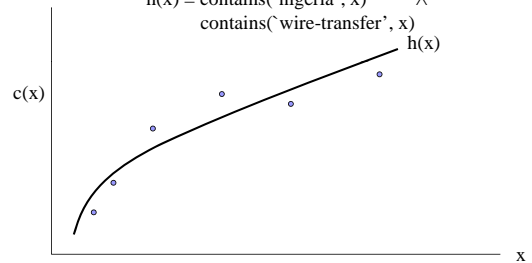
40

ML = Function Approximation

May not be any perfect fit

Classification ~ discrete functions

$h(x) = \text{contains}(\text{'nigeria'}, x) \wedge$
 $\text{contains}(\text{'wire-transfer'}, x)$



41

Why is Learning Possible?

Experience alone never justifies any conclusion about any unseen instance.

Learning occurs when
PREJUDICE meets **DATA!**

Learning a "Frobnitz"

© Daniel S. Weld

42

Bias

- The nice word for prejudice is "bias".
 - Different from "Bias" in statistics
- What kind of hypotheses will you **consider**?
 - What is allowable **range** of functions you use when approximating?
- What kind of hypotheses do you **prefer**?

© Daniel S. Weld

43

Some Typical Biases

- Occam's razor
 - "It is needless to do more when less will suffice"
 - William of Occam,
died 1349 of the Black plague
- MDL – Minimum description length
- Concepts can be approximated by
- ... **conjunctions** of predicates
 - ... by **linear** functions
 - ... by **short** decision trees

Frobnitz?

© Daniel S. Weld

44

ML as Optimization

- Specify Preference Bias
 - aka "Loss Function"
- Solve using optimization
 - Combinatorial
 - Convex
 - Linear
 - Nasty

©2005-2009 Olexa Group

45

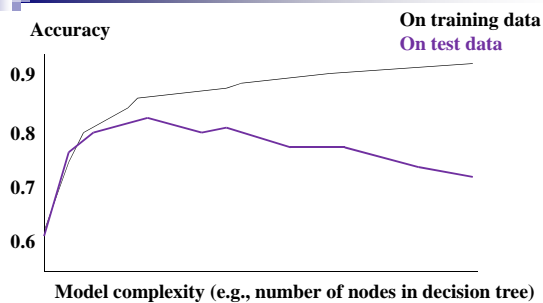
Overfitting

- Hypothesis H is *overfit* when $\exists H'$ and
 - H has **smaller** error on training examples, but
 - H has **bigger** error on test examples

Overfitting

- Hypothesis H is *overfit* when $\exists H'$ and
 - H has **smaller** error on training examples, but
 - H has **bigger** error on test examples
- Causes of overfitting
 - Training set is too small
 - Large number of features
- Big problem in machine learning
 - One solution: Validation set

Overfitting



© Daniel S. Weld

48

The Road Ahead

49