Machine Learning Intro

CompSci 370 Ronald Parr Department of Computer Science Duke University



- Considered a hallmark of intelligence
- Viewed as way to reduce programming burden
 - Not enough programmers in the world to produce custom solutions to all problems – even if we knew how
 - Programmers are expensive!
- Many algorithms assume parameters that are difficult to determine exactly a priori
 - What is the right formula to filter spam?
 - When should your smart thermostat turn on the heat?

Examples

- SPAM classification
- Computational Biology/medicine
 - Distinguish healthy/diseased tissue (e.g., skin/colon cancer)
 - Find structure in biological data (regulatory pathways)
- Financial events
 - Predict good/bad credit risks
 - Predict price changes
 - Response to marketing
- Object/person recognition
- Natural language processing
- Document categorization and user preferences
- Recommend products to users
- Learn to play games, e.g., go, chess, etc.
- Learn to control systems, e.g., robots or helicopters
- Public database of (old) benchmark learning problems:
 - http://www.ics.uci.edu/~mlearn/MLSummary.html



What is Machine Learning?

- Learning Element
 - The thing that learns
- Performance Element
 - Objective measure of progress
- Learning is simply an increase in the ability of the learning element over time (with data) to achieve the task specified by the performance element



ML vs. Data Mining

- Machine Learning is:
 - (Arguably) more formal
 - (Arguably) more task driven/decision theoretic
- Data Mining is:
 - More constrained by size of data set
 - More closely tied to database techniques







Features can be anything Images, sounds, text Real values (height, weight) Integers, or binaries Targets can be discrete classes: Safe mushrooms vs. poisonous Malignant vs. benign Good credit risk vs. bad Label of image Or numbers Selling price of house Life expectancy

How Most Supervised Learning Algorithms Work

- Main idea: Minimize error on training set
- How this is done depends on:
 - Hypothesis space
 - Type of data
- Big Question: What is the "right" hypothesis space?
- The following example for *regression* (continuous targets) is from Chris Bishop







































Why Neural Networks?

Maybe computers should be more brain-like:

	Computers	Brains
Computational Units	10 ¹⁰ transistors/CPU	10 ¹¹ neurons/brain
Storage Units	10 ¹¹ bits RAM 10 ¹³ bits HD	10 ¹¹ neurons 10 ¹⁴ synapses
Cycle Time	10 ⁻⁹ S	10 ⁻³ S
Bandwidth	10 ¹⁰ bits/s*	10 ¹⁴ bits/s
Compute Power	10 ¹⁰ Ops/s	10 ¹⁴ Ops/s





Artificial Neural Networks

- Develop abstraction of function of actual neurons
- Simulate large, massively parallel artificial neural networks on conventional computers note that even supercomputers have very low connectivity compared to a brain
- Some have tried to build the hardware too
- Try to approximate human learning, robustness to noise, robustness to damage, etc.

Neural Network Lore

- Neural nets have been adopted with an almost religious fervor within the AI community several times
 - First coming: Perceptron
 - Second coming: Multilayer networks
 - Third coming (present): Deep networks
- Sound science behind neural networks: gradient descent
- Unsound social phenomenon behind neural networks: HYPE!





















Differentiating h

• Recall the logistic sigmoid:

$$h(x) = \frac{e^{x}}{1 + e^{x}} = \frac{1}{1 + e^{-x}}$$
$$1 - h(x) = \frac{e^{-x}}{1 + e^{-x}} = \frac{1}{1 + e^{x}}$$
• Differentiating:

$$h'(x) = \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1}{(1+e^{-x})} \frac{e^{-x}}{(1+e^{-x})} = h(x)(1-h(x))$$



Summary of Gradient Update

- Gradient calculation, parameter updates have recursive formulation
- Decomposes into:
 - Local message passing
 - No transcendentals:
 - h'(x)=1-h(x)² for tanh(x)
 - H'(x)=h(x)(1-h(x)) for logistic sigmoid
- Highly parallelizable
- Biologically plausible(?)
- Celebrated *backpropagation* algorithm















Good News

- Can represent any continuous function with two layers (1 hidden)
- Can represent essentially any function with 3 layers
- (But how many hidden nodes?)
- Multilayer nets are a universal approximation architecture with a highly parallelizable training algorithm



Backprop Issues

- Backprop = gradient descent on an error function
- Function is nonlinear (= powerful)
- Function is nonlinear (= local minima)
- Big nets:
 - Many parameters
 - Many optima
 - Slow gradient descent
 - Risk of overfitting
 - Biological plausibility ≠ Electronic plausibility
- Many NN experts became experts in numerical analysis (by necessity)



Deep Networks

- Not a learning algorithm, but a family of techniques
 - Improved training techniques (though still essentially gradient descent)
 - Clever crafting of network structure convolutional nets
 - Some new activation functions
- Exploit massive computational power
 - Parallel computing
 - GPU computing
 - Very large data sets (can reduce overfitting)

Conclusions

- Supervised learning = successful way to take training (input, output pairs) and induce functions that generalize to test data drawn from the same distribution as the training data.
- Methods for learning linear functions are well understood and perform well with good features
- Non-linear methods, such as neural networks are more powerful and require less feature engineering but are more computationally expensive and less predictable in practice
 - Historically wild swings in popularity
 - Currently on upswing due to clever changes in training methods, use of parallel computation, and large data sets