



Outline

- Motivation: why should agents learn?
- Different models of learning
- Learning from observation
 - classification and regression
- Learning decision trees
- Linear regression

Machine Learning is everywhere!

- Every time you speak on the phone to an automated program (travel, UPS, Fedex, ...)
- Google, Amazon, Facebook all extensively use machine learning to predict user behavior (they hire many of our PhD's)
- Machine learning is one of the most sought after disciplines for prospective employers





The Future: ML everywhere!

- Hand-held devices will have terabytes of RAM and petabytes of disk space
- Massive use of machine learning across all smartphones, web software, OS, desktops
- Cars will increasingly use machine learning to drive autonomously
- Hard to overestimate the impact of ML

Human Learning

- Learning is a hallmark of intelligence
- Human abilities depend on learning
 - Learning a language (e.g., English, French)
 - Learning to drive
 - Learning to recognize people (faces)
 - Learning in the classroom







Types of Learning

- There are many types of learning
 - Supervised learning
 - Unsupervised learning
 - Reinforcement learning
 - Evolutionary (genetic) learning

Supervised Learning

- Simplest model of learning
- An agent is given positive and negative examples of some concept or function
- The goal is to learn an approximation of the desired concept or function
- Classification: discrete concept spaces
- Regression: real-valued functions

Character Recognition

Apple Apple Apple Apple Apple 383 383 383 383 383

- Humans can effortlessly recognize complex visual patterns (characters, faces, text)
- This apparently simple problem is formidably difficult for machines







Attribute-Value Data

Examples described by attribute values (Boolean, discrete, continuous, etc.) E.g., situations where I will/won't wait for a table:

Example	Attributes										Target
r	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	WillWait
X_1	Т	F	F	Т	Some	\$\$\$	F	Т	French	0–10	Т
X_2	Т	F	F	Т	Full	\$	F	F	Thai	30–60	F
X_3	F	Т	F	F	Some	\$	F	F	Burger	0–10	Т
X_4	Т	F	Т	Т	Full	\$	F	F	Thai	10–30	Т
X_5	Т	F	Т	F	Full	\$\$\$	F	Т	French	>60	F
X_6	F	Т	F	Т	Some	\$\$	Т	Т	Italian	0–10	Т
X_7	F	Т	F	F	None	\$	Т	F	Burger	0–10	F
X_8	F	F	F	Т	Some	\$\$	Т	Т	Thai	0–10	Т
X_9	F	Т	Т	F	Full	\$	Т	F	Burger	>60	F
X_{10}	Т	Т	Т	Т	Full	\$\$\$	F	Т	Italian	10–30	F
X_{11}	F	F	F	F	None	\$	F	F	Thai	0–10	F
X_{12}	Т	Т	Т	Т	Full	\$	F	F	Burger	30–60	Т

Classification of examples is positive (T) or negative (F)







How many distinct decision trees with n Boolean attributes??

- = number of Boolean functions
- = number of distinct truth tables with $2^n \mbox{ rows} = 2^{2^n}$

E.g., with 6 Boolean attributes, there are 18,446,744,073,709,551,616 trees

 \odot

How many purely conjunctive hypotheses (e.g., $Hungry \land \neg Rain$)??

Each attribute can be in (positive), in (negative), or out

 \Rightarrow 3^n distinct conjunctive hypotheses

More expressive hypothesis space

- increases chance that target function can be expressed

– increases number of hypotheses consistent $w/\ training\ set$

 \Rightarrow may get worse predictions \bigotimes

Learning Decision Trees

Aim: find a small tree consistent with the training examples

Idea: (recursively) choose "most significant" attribute as root of (sub)tree





Entropy

- We can apply ideas from the field of information theory to attribute selection
- Given a set of events, each of which occurs with probability p_i, entropy measures the surprise associated with a particular outcome
- Low frequency events are more surprising

$$H(P) = -\sum_{i} p_i \log_2 p_i$$

Information Theory

Suppose we have p positive and n negative examples at the root $\Rightarrow \ H(\langle p/(p+n), n/(p+n)\rangle) \text{ bits needed to classify a new example E.g., for 12 restaurant examples, } p = n = 6 \text{ so we need 1 bit}$

An attribute splits the examples E into subsets E_i , each of which (we hope) needs less information to complete the classification

Let E_i have p_i positive and n_i negative examples

 $\Rightarrow H(\langle p_i/(p_i+n_i), n_i/(p_i+n_i)\rangle) \text{ bits needed to classify a new example} \\\Rightarrow expected number of bits per example over all branches is$

$$\sum_{i} \frac{p_i + n_i}{p + n} H(\langle p_i / (p_i + n_i), n_i / (p_i + n_i) \rangle)$$

For Patrons?, this is 0.459 bits, for Type this is (still) 1 bit

 $\Rightarrow~$ choose the attribute that minimizes the remaining information needed





Regression

- Another common type of learning involves making continuous real-valued predictions
 - How much does it cost to fly to Europe?
 - How long to drive to Northampton?
 - How much money will I make when I graduate?











Minimize Squared Error

- Process of curve fitting is based on minimizing a loss function
- One example is minimizing sum of squared errors (between predicted and actual values)

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{t_n - \mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_n)\}^2$$

Widrow-Hoff Algorithm

- Incremental algorithm that modifies the weights based on the gradient of the error
- For each example i in a dataset:

$$w^{t+1} \leftarrow w^t + \alpha_t (t_i - \phi(x_i)^T w^t) \phi(x_i)$$

• until the error is small enough



Summary

- Learning is a fundamental component of intelligence
- Classification is a way of discriminating among categories
- Decision trees: simple classification method
- Regression is the estimation of functions
- Least-squares is a standard method for regression