# An Introduction to Machine Learning

LING-7800 Computational Lexical Semantics Presented by Lee Becker and Shumin Wu

### What is Machine Learning?

• AKA

- Pattern Recognition
- Data Mining

### What is Machine Learning?



### What is Machine Learning?

- Programming computers to do tasks that are (often) easy for humans to do, but hard to describe algorithmically.
- Learning from observation
- Creating models that can predict outcomes for unseen data
- Analyzing large amounts of data to discover new patterns

### What is Machine Learning?

- Isn't this just statistics?
  - Cynic's response: Yes
  - CS response: Kind of
    - Unlike in statistics, machine learning is also concerned with the complexity, optimality, and tractability in learning a model
    - Statisticians are often dealing with much smaller amounts of data.

### Ok, so where do we start?

- Observations
  - Data! The more the merrier (usually)
- Representations
  - Often raw data is unusable, especially in natural language processing
  - Need a way to represent observations in terms of its properties (features)

fn

• Feature Vector

f<sub>0</sub> f<sub>1</sub>



### Machine Learning Paradigms

- Supervised Learning
  - Deduce a function from labeled training data to minimize labeling error on future data
- Unsupervised Learning
  - Learning with unlabeled training data
- Semi-supervised Learning
  - Learning with (usually small amount of) labeled training data and (usually large amount of)
- Active Learning
  - Actively query for specific labeled training data
- Reinforcement Learning
  - Learn actions in environment to maximize (often long-term) reward



### Supervised Learning

- Stages
  - Train model on data
  - Tune parameters of the model
  - Select best model
  - Evaluate



# Evaluation But what are we comparing against? Typically the data is divided into three parts Training Development Test / Validation Typically accuracy on the validation set is reported Why all this extra effort? The goal in machine learning is to select the model that does the best on unseen data This divide is an attempt to keep our experiment honest Avoids overfitting







### Naïve Bayes Dog vs Cat Classifier

### • 2 features: weight & how frequent it chases mouse

mouse chase	weight	label
0.7	55	dog
0.05	15	dog
0.2	100	dog
0.25	42	dog
0.2	32	dog
0.6	25	cat
0.2	15	cat
0.55	8	cat
0.15	12	cat
0.4	15	cat

Given an animal that weighs no more than 20 lbs and chases mouse at least 21% of time, is it a cat or dog?

 $f(dog, w \le 20, m \ge .21) = p(dog)p(w \le 20 | dog)p(m \ge 0.21 | dog) = 0.5 \times 0.2 \times 0.4 = 0.04$ 

 $\begin{array}{l} f(cat, w \leq 20, m \geq .21) = \\ p(cat)p(w \leq 20 \,|\, cat)p(m \geq 0.21 \,|\, cat) = \\ 0.5 \times 0.4 \times 0.6 = 0.12 \end{array}$ 

So, it's cat! In fact, naïve Bayes is 75% certain it's a cat over a dog.

### Linear Classifier

• Features have linear relationships with each other:

$$g(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots + \beta_m x_m$$
$$f(\mathbf{x}) = \begin{pmatrix} \text{class 1 if } g(\mathbf{x}) \ge 0\\ \text{class 2 if } g(\mathbf{x}) < 0 \end{pmatrix}$$















### **Decision Trees**

- A better approach
  - Find the most important attribute first
  - Prune the tree based on these decision
  - Lather, Rinse, and Repeat as necessary

### Decision Trees

- Choosing the best attribute
  - Measuring Information (Entropy):

$$I(P(v_1),...,P(v_2)) = \sum_{i=1}^{n} -P(v_i)\log_2 P(v_i)$$

Examples:

• Tossing a fair coin  $I(P(heads), P(tails)) = I(\frac{1}{2}, \frac{1}{2}) = -\frac{1}{2}\log_2 \frac{1}{2} - \frac{1}{2}\log_2 \frac{1}{2} = 1$  bits

- Tossing a biased coin
  - $I(P(heads), P(tails)) = I(\frac{1}{100}, \frac{99}{100}) = -\frac{1}{100}\log_2 \frac{1}{100} \frac{99}{100}\log_2 \frac{99}{100} = 0.08$  bits
- Tossing a fair die
  - $I(P(1), P(2), ..., P(6)) = I(\frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}, \frac{1}{6}) = 2.58$  bits

### **Decision Trees**

- Choosing the best attribute cont'd
  - New information requirement due to an attribute

Remainder (A) = 
$$\sum_{i=1}^{\nu} I\left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i}\right)$$

 – Gain = Original Information Requirement – New Information Requirement

$$Gain(A) = I\left(\frac{p}{p+n}, \frac{n}{p+n}\right) - Remainder(A)$$

### **Decision Trees**

arks	Chase Mice (Freq)	Chase Ball (Freq)	Weight (Pounds)	Matching Eye C	olor Category
TRUE	0.7	1 1	55	TRUE	Dog
TRUE	0.2	. 0.9	22	TRUE	Dog
TRUE	0.1	0.8	38	TRUE	Dog
TRUE	0.8	0.1	17	TRUE	Dog
TRUE	0.2	. 0	100	TRUE	Dog
FALSE	0.1	0.7	27	TRUE	Dog
FALSE	0.25	i 0.6	42	TRUE	Dog
FALSE	0.4	0.5	25	TRUE	Dog
FALSE	0.2	. 0.3	32	TRUE	Dog
FALSE	0.3	0.2	10	TRUE	Dog
FALSE	0.6	i 0.5	25	TRUE	Cat
FALSE	0.6	i 0.4	22	TRUE	Cat
FALSE	0.2	0.6	15	TRUE	Cat
FALSE	0.2	. 0.2	10	TRUE	Cat
FALSE	0.55	i 0.1	8	TRUE	Cat
FALSE	0.8	. 0	11	TRUE	Cat
FALSE	0.15	0.25	12	TRUE	Cat
FALSE	0.7	0.3	9	TRUE	Cat
FALSE	0.4	ı 0	15	FALSE	Cat
FALSE	0.3	. 0	13	TRUE	Cat



Decision Trees									
	Cats and Dogs								
	- Step 2: Information Requirement $I\left(\frac{p}{p+n}, \frac{n}{p+n}\right) = I\left(\frac{5}{15, 15}\right) = .918$ bits								
	<ul> <li>Information gain by attributes</li> </ul>								
	Attribute	P(Dog A)	P(Cat A)	P(Dog ~A)	P(Cat ~A)	Remainder	Gain		
<	Chases Mice	0	1	.5	.5	.667	.252	$\geq$	
	Chases Ball	.667	.333	.25	.75	.832	.086		
	Weight > 30	1	0	.231	.769	.675	.242		
	Eye Color Matches	.357	.642	.357	.643	.877	.041		





### **Other Popular Classifiers**

- Support Vector Machines (SVM)
- Maximum Entropy
- Neural Networks
- Perceptron

## Machine Learning for NLP (courtesy of Michael Collins)

- The General Approach:
  - Annotate examples of the mapping you're interested in
  - Apply some machinery to learn (and generalize) from these examples
- The difference from classification
  - Need to induce a mapping from one complex set to another (e.g. strings to trees in parsing, strings in machine translation, strings to database entries in information extraction)
- Motivation for learning approaches (as opposed to "hand-built" systems
  - Often, a very large number of rules is required.
  - Rules interact in complex and subtle ways.
  - Constraints are often not "categorical", but instead are "soft" or violable.
  - A classic example: Speech Recognition





### **K-Mean Clustering**

- Aims to partition *n* observations into *k* clusters. Wherein each observation is in the cluster with the nearest mean.
- Iterative 2-stage process
  - Assignment Step
  - Update Step

### Hierarchical Clustering

- Build a hierarchy of clusters
- Find successive clusters using previously established clusters
- Paradigms
  - Agglomerative: Bottom-up
  - Divisive: Top-down













### **Cluster Evaluation**

• Normalized Mutual Information

 $NMI(\Omega,C) = \frac{I(\Omega,C)}{[H(\Omega) + H(C)]/2}$ where the set of clusters  $\Omega = \{\omega_1, \omega_2, ..., \omega_k\}$ and the set of classes  $C = \{c_1, c_2, ..., c_j\}$ 

- Drawbacks

• Requires members to have labels

### Application: Automatic Verb Class Identification

- Feature Representation:
  - Want features that provide clues to the sense used
    - Word co-occurrence
      - The <u>ball</u>rolled down the hill.
      - The wheel rolled away.
      - The <u>ball</u> bounced.
    - Selectional Preferences
      - Part of Speech
      - Semantic Roles
    - Construction
      - Passive
      - Active
    - Other
      - Is the verb also a noun?

### Application: Automatic Verb Class Identification

- Goal: Given significant amounts of sentences, discover verb classes
  - Example:
    - Steal-10.5: Abduct, Annex, Capture, Confiscate, ...
    - Butter-9.9: Asphalt, Butter, Brick, Paper, ...
    - Roll-51.3.1: Bounce, Coil, Drift, Drop....
- Approach:
  - Determine meaningful feature representation for each verb
  - Extract set of observations from a corpus
  - Apply clustering algorithm
  - Evaluate



Disturb ... .004 .003 .000 .33 .21 .000 .005 0

