#### Lecture Slides for

### **INTRODUCTION TO Machine Learning 2nd Edition**



ETHEM ALPAYDIN, modified by Leonardo Bobadilla and some parts from http://www.cs.tau.ac.il/~apartzin/MachineLearning/ and

www.cs.princeton.edu/courses/archive/fall01/cs302 /notes/11.../EM.ppt © The MIT Press, 2010 alpaydin@boun.edu.tr

http://www.cmpe.boun.edu.tr/~ethem/i2m

### Outline

Previous class Ch 8: Clustering This class: Ch 9: Decision Trees

### CHAPTER 7: Clustering

### Example:

- Data came from mix of Gaussians  $\hat{p}_i(\mathbf{x}^t | \Phi) \sim \mathcal{N}(\mathbf{m}_i, \mathbf{S}_i)$ ,
- Maximize likelihood assuming we know latent "indicator variable"

$$\boldsymbol{m}_{i}^{l+1} = \frac{\sum_{t} h_{i}^{t} \boldsymbol{x}^{t}}{\sum_{t} h_{i}^{t}}$$
$$\boldsymbol{S}_{i}^{l+1} = \frac{\sum_{t} h_{i}^{t} (\boldsymbol{x}^{t} - \boldsymbol{m}_{i}^{l+1}) (\boldsymbol{x}^{t} - \boldsymbol{m}_{i}^{l+1})^{T}}{\sum_{t} h_{i}^{t}}$$

• E-step: expected value of indicator variables

$$h_i^t = \frac{\pi_i |\mathbf{S}_i|^{-1/2} \exp[-(1/2)(\mathbf{x}^t - \mathbf{m}_i)^T \mathbf{S}_i^{-1}(\mathbf{x}^t - \mathbf{m}_i)]}{\sum_j \pi_j |\mathbf{S}_j|^{-1/2} \exp[-(1/2)(\mathbf{x}^t - \mathbf{m}_j)^T \mathbf{S}_j^{-1}(\mathbf{x}^t - \mathbf{m}_j)]}$$



### EM for Gaussian mixtures

- Assume all groups/clusters are Gaussians
- Multivariate Uncorrelated
- Same Variance
- Harden indicators
  - EM: expected values are between 0 and 1
  - K-means: 0 or 1
- Same as k-means

### Dimensionality Reduction vs. Clustering

- Dimensionality reduction methods find correlations between features and group features
  - Age and Income are correlated
- Clustering methods find similarities between instances and group instances
  - Customer A and B are from the same cluster

## Clustering: Usage for supervised learning

- Describe data in terms of cluster
  - Represent all data in cluster by cluster mean
  - Range of attributes
- Map data into new space(preprocessing)
  - d- dimension original space
  - k- number of clusters
  - Use indicator variables as data representations
  - k might be larger then d

### Mixture of Mixtures

- In classification, the input comes from a mixture of classes (supervised).
- If each class is also a mixture, e.g., of Gaussians, (unsupervised), we have a mixture of mixtures:

$$p(\mathbf{x} | \mathbf{C}_i) = \sum_{j=1}^{k_i} p(\mathbf{x} | \mathbf{G}_{ij}) P(\mathbf{G}_{ij})$$
$$p(\mathbf{x}) = \sum_{i=1}^{K} p(\mathbf{x} | \mathbf{C}_i) P(\mathbf{C}_i)$$

### **Hierarchical Clustering**

- Probabilistic view
  - Fit mixture model to data
  - Find codewords minimizing reconstruction error
- Hierarchical clustering
  - Group similar items together
  - No specific model/distribution
  - Items in groups is more similar to each other than instances in different groups

### **Hierarchical Clustering**

Minkowski ( $L_p$ ) (Euclidean for p = 2)  $d_m(\mathbf{x}^r, \mathbf{x}^s) = \left[\sum_{j=1}^d (x_j^r - x_j^s)^p\right]^{1/p}$ 

City-block distance  $d_{cb}(\mathbf{x}^r, \mathbf{x}^s) = \sum_{j=1}^d \mathbf{x}_j^r - \mathbf{x}_j^s$ 

### Agglomerative Clustering

- Start with clusters each having single point
- At each step merge similar clusters
- Measure of similarity
  - Minimal distance(single link)
    - Distance between closest points in 2 groups
  - Maximal distance(complete link)
    - Distance between most distant points in 2 groups
  - Average distance
    - Distance between group centers

### Example: Single-Link Clustering



Based on for E ALPAYDIN 2004 Introduction to Machine Learning © The MIT Press (V1.1)

### Choosing k

- Defined by the application, e.g., image quantization
- Plotting data in two dimensions using PCA
- Incremental (leader-cluster) algorithm: Add one at a time until "elbow" (reconstruction error/log likelihood/intergroup distances)

### CHAPTER 9: Decision Trees

Slides from: Blaž Zupan and Ivan Bratko magix.fri.uni-lj.si/predavanja/uisp

### Motivation

- Parametric Estimation
  - Assume model for class probability or regression
  - Estimate parameters from all data

### Motivation

- Pre-split training data into region using small number of simple rules organized in hierarchical manner
- Decision Trees
  - Internal decision nodes have splitting rule
  - Terminal leaves have class labels for classification problem or values for regression problem

### Tree Uses Nodes, and Leaves



Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

# An Example Data Set and Decision Tree

3	sunny	med	big	yes
4	sunny	no	small	yes
5	sunny	big	big	yes
6	rainy	no	small	no
7	rainy	med	small	yes
8	rainy	big	big	yes
9	rainy	no	big	no
10	rainy	med	big	no



### Classification

#	Attribute			Class
	Outlook	Company	Sailboat	Sail?
1	sunny	no	big	?
2	rainy	big	small	?



### Learning Decision Trees

- Data Set (Training Set)
  - Each example = Attributes + Class
- Description = Decision tree
- Recursive Partitioning

### Analysis of Severe Trauma Patients Data



PH\_ICU and APPT\_WORST are exactly the two factors (theoretically) advocated to be the most important ones in the study by Rotondo et al., 1997.

### Prostate cancer recurrence



### ID3 Algorithm

To construct decision tree T from learning set X:

If all examples in X belong to some class C
Then

make leaf labeled C

#### - Otherwise

- select the "most informative" attribute A
- partition S according to A's values
- recursively construct subtrees T1, T2, ..., for the subsets of S

### ID3 Algorithm

• Resulting tree T is:



### Another Example

#	Attribute			Class	
	Outlook	Temperature	Humidity	Windy	Play
1	sunny	hot	high	no	Ν
2	sunny	hot	high	yes	Ν
3	overcast	hot	high	no	Р
4	rainy	moderate	high	no	Р
5	rainy	cold	normal	no	Р
6	rainy	cold	normal	yes	Ν
7	overcast	cold	normal	yes	Р
8	sunny	moderate	high	no	Ν
9	sunny	cold	normal	no	Р
10	rainy	moderate	normal	no	Р
11	sunny	moderate	normal	yes	Р
12	overcast	moderate	high	yes	Р
13	overcast	hot	normal	no	Р
14	rainy	moderate	high	yes	Ν

### Simple Tree



### **Complicated Tree**



### Attribute Selection Criteria

- Main principle
  - Select attribute which partitions the learning set into subsets as "pure" as possible
- Various measures of purity
  - Information-theoretic
  - Gini index
  - ....

### Information-Theoretic Approach

- To classify an object, a certain information is needed
  - I, information
- After we have learned the value of attribute A, we only need some remaining amount of information to classify the object
  - Ires, residual information
- Gain Gain(A) = I Ires(A)
- The most 'informative' attribute is the one that minimizes Ires, *i.e.*, maximizes Gain

### Entropy

• The average amount of information *I* needed to classify an object is given by the entropy measure

$$I = -\sum_{c} p(c) \log_2 p(c)$$

• For a two-class problem:



### **Residual Information**

- After applying attribute A, S is partitioned into subsets according to values v of A
- *Ires* is equal to weighted sum of the amounts of information for the subsets

$$I_{res} = -\sum_{v} p(v) \sum_{c} p(c|v) \log_2 p(c|v)$$

### **Triangles and Squares**

#	Attribute			Shape
	Color	Outline	Dot	_
1	green	dashed	no	triange
2	green	dashed	yes	triange
3	yellow	dashed	no	square
4	red	dashed	no	square
5	red	solid	no	square
6	red	solid	yes	triange
7	green	solid	no	square
8	green	dashed	no	triange
9	yellow	solid	yes	square
10	red	solid	no	square
11	green	solid	yes	square
12	yellow	dashed	yes	square
13	yellow	solid	no	square
14	red	dashed	yes	triange

Data Set: A set of classified objects



### Entropy



- 5 triangles
- 9 squares
- class probabilities

$$p(\Box) = \frac{9}{14}$$

$$p(\Delta) = \frac{5}{14}$$

entropy

$$I = -\frac{9}{14}\log_2\frac{9}{14} - \frac{5}{14}\log_2\frac{5}{14} = 0.940$$
 bits





I(yellow) = 0.0 bits

$$I_{res}(\text{Color}) = \sum p(v)I(v) = \frac{5}{14}0.971 + \frac{5}{14}0.971 + \frac{4}{14}0.0 = 0.694 \text{ bits}$$





I(yellow) = 0.0 bits

 $Gain(Color) = I - I_{res}(Color) = 0.940 - 0.694 = 0.246 \ bits$ 

### Information Gain of The Attribute

- Attributes
  - Gain(Color) = 0.246
  - Gain(Outline) = 0.151
  - Gain(Dot) = 0.048
- Heuristics: attribute with the highest gain is chosen
- This heuristics is local (local minimization of impurity)



Gain(Outline) = 0.971 - 0 = 0.971 bits Gain(Dot) = 0.971 - 0.951 = 0.020bits





### **Decision Tree**



### **Decision Trees**

- Start from univariate decision trees
  - Each node looks only at single input feature
- Want smaller decision trees
  - Less memory for representation
  - Less computation for a new instance
- Want smaller generalization error

### Decision and Leaf Node

- Implement simple test function f<sub>m</sub>(x)
- Output: labels of branches
- $f_m(x)$  discriminant in d-dimensional space
- Complex discriminant is broken down into hierarchy of simple decisions
- Leaf node describes a region in d-dimensional space with same value
  - Classification label
  - Regression value

### Classification Trees

- What is the good split function?
- Use Impurity measure
- Assume  $N_m$  training samples reach node m
- $N_m^i$  of  $N_m$  belong to class  $C_i$ , with  $\sum_i N_m^i = N_m$ .

• 
$$\hat{P}(C_i|\mathbf{x},m) \equiv p_m^i = \frac{N_m^i}{N_m}$$

• Node m is pure if for all classes either 0 or 1

### Entropy

- Measure amount of uncertainty on a scale from 0 to 1
- Example: 2 events
- If p1=p2=0.5, entropy is 1 which is maximum uncertainty
- If p1=1=1-p0, entropy is 0, which is no uncertainty  $K_{m} = -\sum p_{m}^{i} \log_{2} p_{m}^{i}$

i=1

### Entropy



### Best Split

- Node is impure, need to split more
- Have several split criteria (coordinates), have to choose optimal
- Minimize impurity (uncertainty) after split
- Stop when impurity is small enough
  - Zero stop impurity=>complex tree with large variance
  - Larger Stop impurity=>\$mail tress but

### Best Split

- Impurity after split:  $N_{mj}$  of  $N_m$  take branch *j*.
- *N*<sup>*i*</sup><sub>*mj*</sub> belong to C<sub>*i*</sub>

$$\hat{P}(C_i \mid \boldsymbol{X}, m, j) \equiv p_{mj}^i = \frac{N_{mj}^i}{N_{mj}}$$

$$\mathbf{I'}_{m} = -\sum_{j=1}^{n} \frac{N_{mj}}{N_{m}} \sum_{i=1}^{K} p_{mj}^{i} \log_{2} p_{mj}^{i}$$

- Find the variable and split that min impurity
  - among all variables
  - split positions for numeric variables

