

# A Systematic Overview of Data Mining Algorithms

Sargur Srihari

University at Buffalo

The State University of New York

# Topics

- Data Mining Algorithm Definition
- Example of CART Classification
  - Iris, Wine Classification
- Reductionist Viewpoint
  - Data Mining Algorithm as a 5-tuple
  - Three Cases
    - MLP for Regression/Classification
    - A Priori Algorithm
    - Vector-space Text Retrieval

# Data Mining Algorithm Definition

- A data mining algorithm is a well-defined procedure
  - that takes *data* as input and
  - produces as output: *models* or *patterns*
- Terminology in Definition
  - well-defined:
    - procedure can be precisely encoded as a finite set of rules
  - algorithm:
    - procedure terminates after finite no of steps and produces an output
  - computational method (procedure):
    - has all properties of an algorithm except guaranteeing finite termination
    - e.g., search based on steepest descent is a *computational method*- for it to be an algorithm need to specify where to begin, how to calculate direction of descent, when to terminate search
  - model structure
    - a global summary of the data set,
    - e.g.,  $Y=aX+c$  where  $Y, X$  are variables;  $a, c$  are extracted parameters
  - pattern structure: statements about restricted regions of the space
    - If  $X > x_l$  then  $prob( Y > y_l) = p_l$

# Components of a Data Mining Algorithm

## 1. Task

e.g., visualization, classification, clustering, regression, etc

## 2. Structure (functional form) of model or pattern

e.g., linear regression, hierarchical clustering

## 3. Score function to judge quality of fitted model or pattern,

e.g., generalization performance on unseen data

## 4. Search or Optimization method

e.g., steepest descent

## 5. Data Management technique storing, indexing and retrieving data. ML algorithms do not specify this. Massive data sets need it.

# Components of 3 well-known Data Mining algorithms

Component/ Name	CART (model)	Backpropagation (parameter est.)	A Priori
1. Task	Classification and Regression	Classification and Regression	Rule Pattern Discovery
2. Structure	Decision Tree	Neural Network	Association Rules
3. Score Functn	Cross-validated Loss Function	Squared Error	Support/ Accuracy
4. Search Methd	Greedy Search over Structures	Gradient descent on Parameters	Breadth-First with Pruning
5. Data Mgmt Tx	Unspecified	Unspecified	Linear Scans

# CART Algorithm Task

- Classification and Regression Trees
- Widely used statistical procedure
- Produces classification and regression models with a tree-based structure
- Only classification considered here:
  - Mapping input vector  $\mathbf{x}$  to categorical (class) label  $y$

# Classification Aspect of CART

- Task = prediction (classification)
- Model Structure = Tree
- Score Function = Cross-validated Loss Function
- Search Method = greedy local search
- Data Management Method = Unspecified

# Van Gogh: Irises





# Iris Classification

*Iris versicolor*



Iris Setosa



Iris Versicolor



Iris Virginica



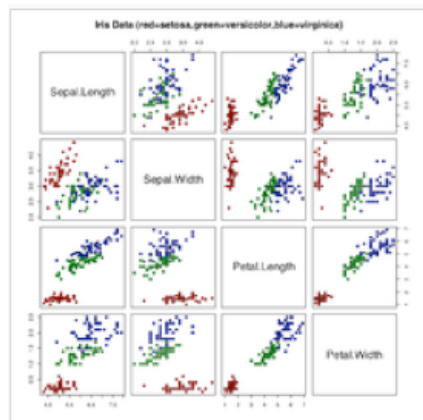
# Fisher's Iris Data Set

UCI Repository

<b>Data Set Characteristics:</b>	Multivariate	<b>Number of Instances:</b>	150	<b>Area:</b>	Life
<b>Attribute Characteristics:</b>	Real	<b>Number of Attributes:</b>	4	<b>Date Donated</b>	1988-07-01
<b>Associated Tasks:</b>	Classification	<b>Missing Values?</b>	No	<b>Number of Web Hits:</b>	121967



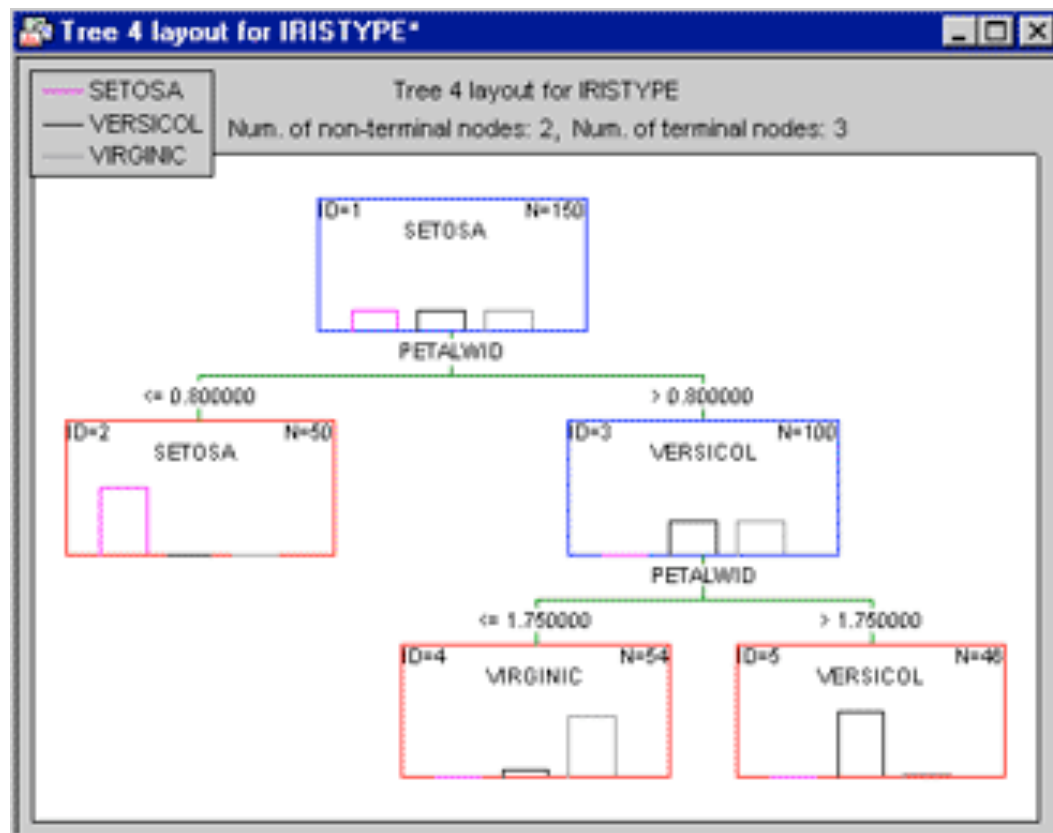
1. sepal length in cm
2. sepal width in cm
3. petal length in cm
4. petal width in cm
5. class:
  - Iris Setosa
  - Iris Versicolour
  - Iris Virginica



Fisher's Iris Data

Sepal Length	Sepal Width	Petal Length	Petal Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
4.6	3.1	1.5	0.2	setosa
5.0	3.6	1.4	0.2	setosa
5.4	3.9	1.7	0.4	setosa
4.6	3.4	1.4	0.3	setosa
5.0	3.4	1.5	0.2	setosa
4.4	2.9	1.4	0.2	setosa
4.9	3.1	1.5	0.1	setosa
5.4	3.7	1.5	0.2	setosa
4.8	3.4	1.6	0.2	setosa
4.9	3.0	1.4	0.1	setosa

# Tree for Iris Data



## Interpretation of tree:

If petal width is less than or equal to 0.8, flower classified as *Setosa*

If petal width is greater than 0.8 and less than or equal to 1.75, Then flower classified as *Virginic* else, it belongs to class *Versicol*

# CART Approach to Classification

- Model structure is a classification tree
  - Hierarchy of univariate binary decisions
  - Each node of tree specifies a binary test
    - On a single variable
    - using thresholds on real and integer variables
    - Subset membership for categorical variables
- Tree derived from data, not specified a priori
- Choosing best variable for splitting data

# Wine Classification



# Wine Data Set

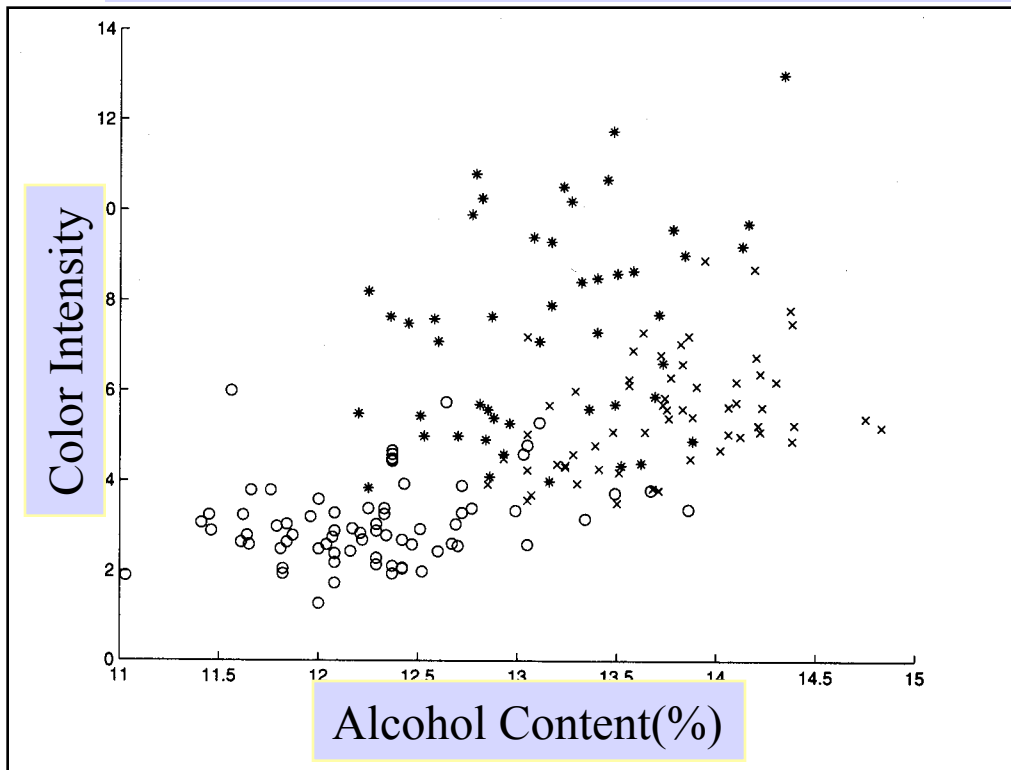
UCI Repository  
Three wine types

<b>Data Set Characteristics:</b>	Multivariate	<b>Number of Instances:</b>	178	<b>Area:</b>	Physical
<b>Attribute Characteristics:</b>	Integer, Real	<b>Number of Attributes:</b>	13	<b>Date Donated</b>	1991-07-01
<b>Associated Tasks:</b>	Classification	<b>Missing Values?</b>	No	<b>Number of Web Hits:</b>	80636

- 1) Alcohol
- 2) Malic acid
- 3) Ash
- 4) Alcalinity of ash
- 5) Magnesium
- 6) Total phenols
- 7) Flavanoids
- 8) Nonflavanoid phenols
- 9) Proanthocyanins
- 10) Color intensity
- 11) Hue
- 12) OD280/OD315 of diluted wines
- 13) Proline

# Wine Classification

Constituents of 3 different wine types  
(cultivars)

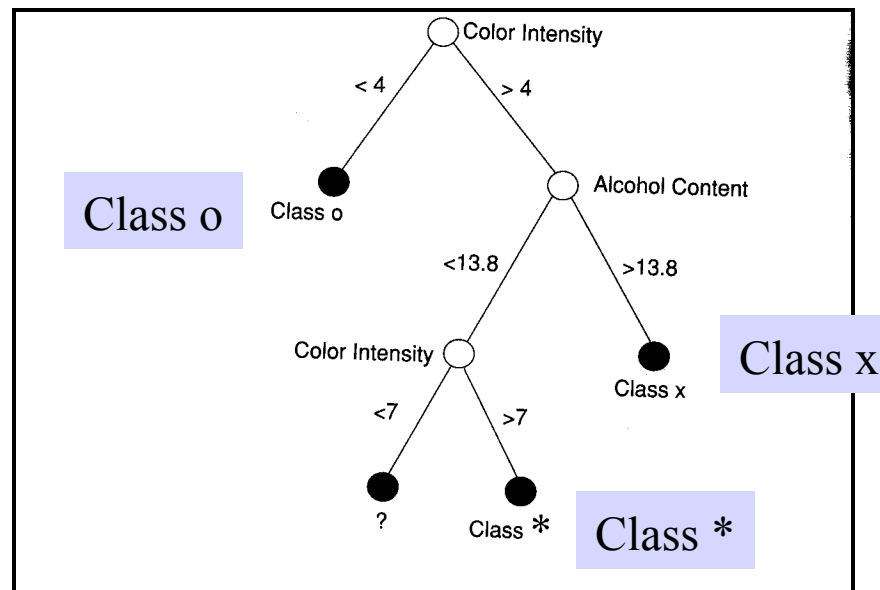


Scatterplot of two variables

- From 13 dimensional data set
- Each variable measures a particular characteristic of a specific wine

# Tree for Wine Classification

Classification into 3 different wine types  
(cultivars)



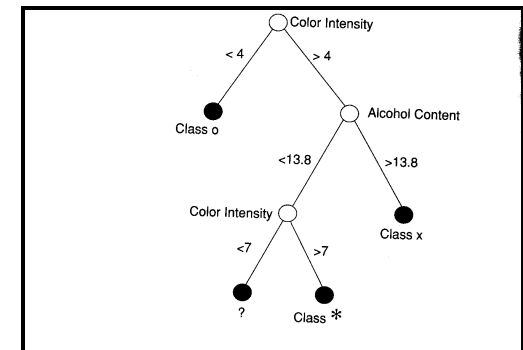
Test of Thresholds (shown beside branches)  
Uncertainty about class label at leaf node labelled as ?



# CART 5-tuple

1. **Task** = prediction (classification)
2. **Model Structure** = tree
3. **Score Function** = cross-validated loss function
4. **Search Method** = greedy local search
5. **Data Management Method** = unspecified

- Hierarchy of univariate binary decisions
- Each internal node specifies a binary test on a single variable
  - Using thresholds on real and integer valued variables
- Can use any of several splitting criteria
- Chooses best variable for splitting data



Classification Tree

# Score Function of CART

- Quality of Tree structure

$$\sum_{i=1}^n C(y(i), \hat{y}(i))$$

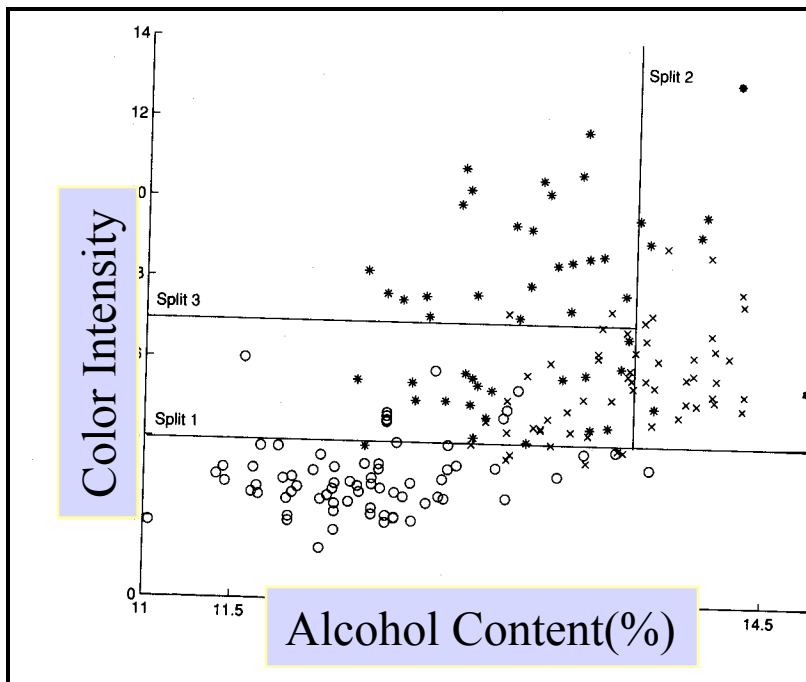
– A misclassification function

- Loss incurred when class label for  $i^{\text{th}}$  data vector  $y(i)$  is predicted by the tree to be  $\hat{y}(i)$
- Specified by an  $m \times m$  matrix, where  $m$  is the number of classes

# CART Search

- Greedy local search to identify candidate structures
- Recursively expands from root node
- Prunes back specific branches of large tree
- Greedy local search is most common method for practical tree learning!

# Classification Tree for Wine



Representational power is coarse:

Decision regions are constrained to be hyper-rectangles with boundaries parallel to input variable axes

Decision Boundaries of Classification Tree Superposed on Data.  
Note parallel nature of boundaries

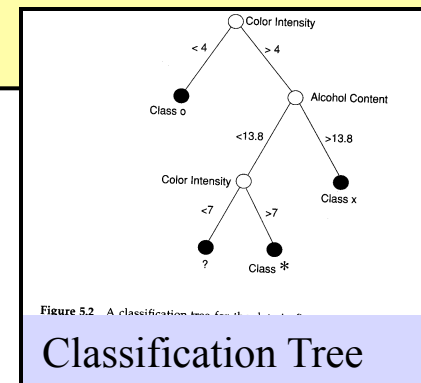


Figure 5.2 A classification tree for wine data.

Classification Tree

# CART Scoring/Stopping Criterion

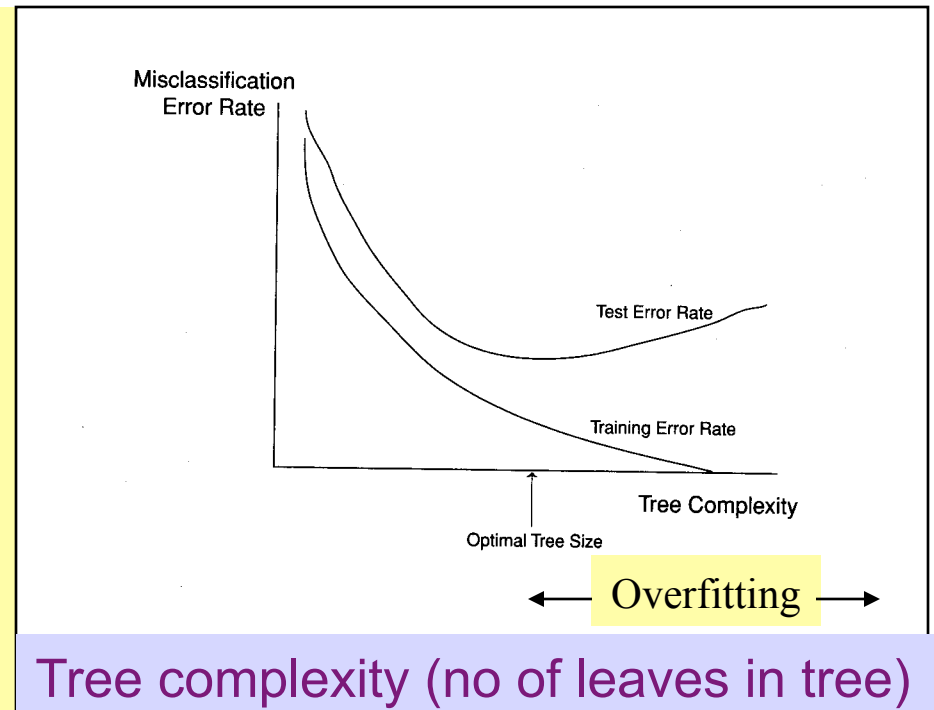
## Cross Validation

to estimate misclassification:

Partition sample into *training* and *validation* sets

Estimate misclassification on validation set

Repeat with different partitions and average results for each tree size



# CART Data Management

- Assumes that all the data is in main memory
- For tree algos data management non-trivial
  - Since it recursively partitions the data set
  - Repeatedly find different subsets of observations in database
  - Naïve implementation involves repeated scans of secondary storage medium leading to poor time performance

# Reductionist Viewpoint of Data Mining Algorithms

- A Data Mining Algorithm is a tuple:  
{model structure, score function, search method, data management techniques}
- Combining different model structures with different score functions, etc will yield a potentially infinite number of different algorithms

# Reductionist Viewpoint applied to 3 algorithms

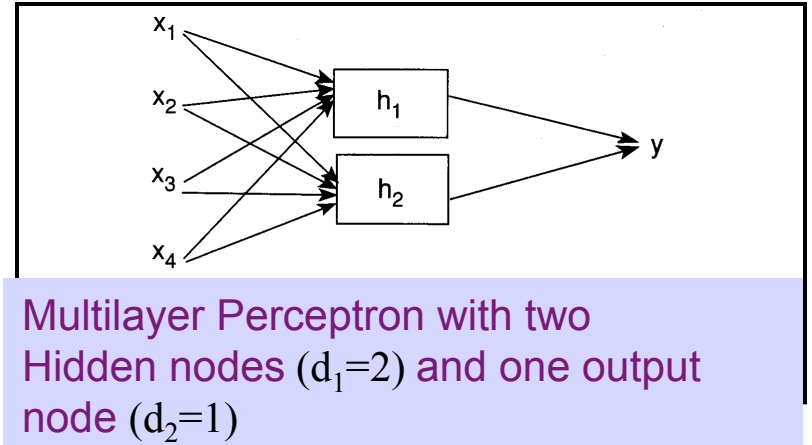
1. *Multilayer Perceptron (MLP)* for Regression and Classification
2. *A Priori* Algorithm for Association Rule Learning
3. *Vector Space Algorithms* for Text Retrieval



# Multilayer Perceptron (MLP)

- Artificial Neural Network
- Non-linear mapping from real-valued input vector  $\mathbf{x}$  to real-valued output vector  $\mathbf{y}$
- Thus MLP can be used as a nonlinear model for regression as well as for classification

# MLP Formulas



- From first layer of weights

$$s_1 = \sum_{i=1}^4 \alpha_i x_i \quad s_2 = \sum_{i=1}^4 \beta_i x_i$$

- Non-linear Transformation at hidden nodes

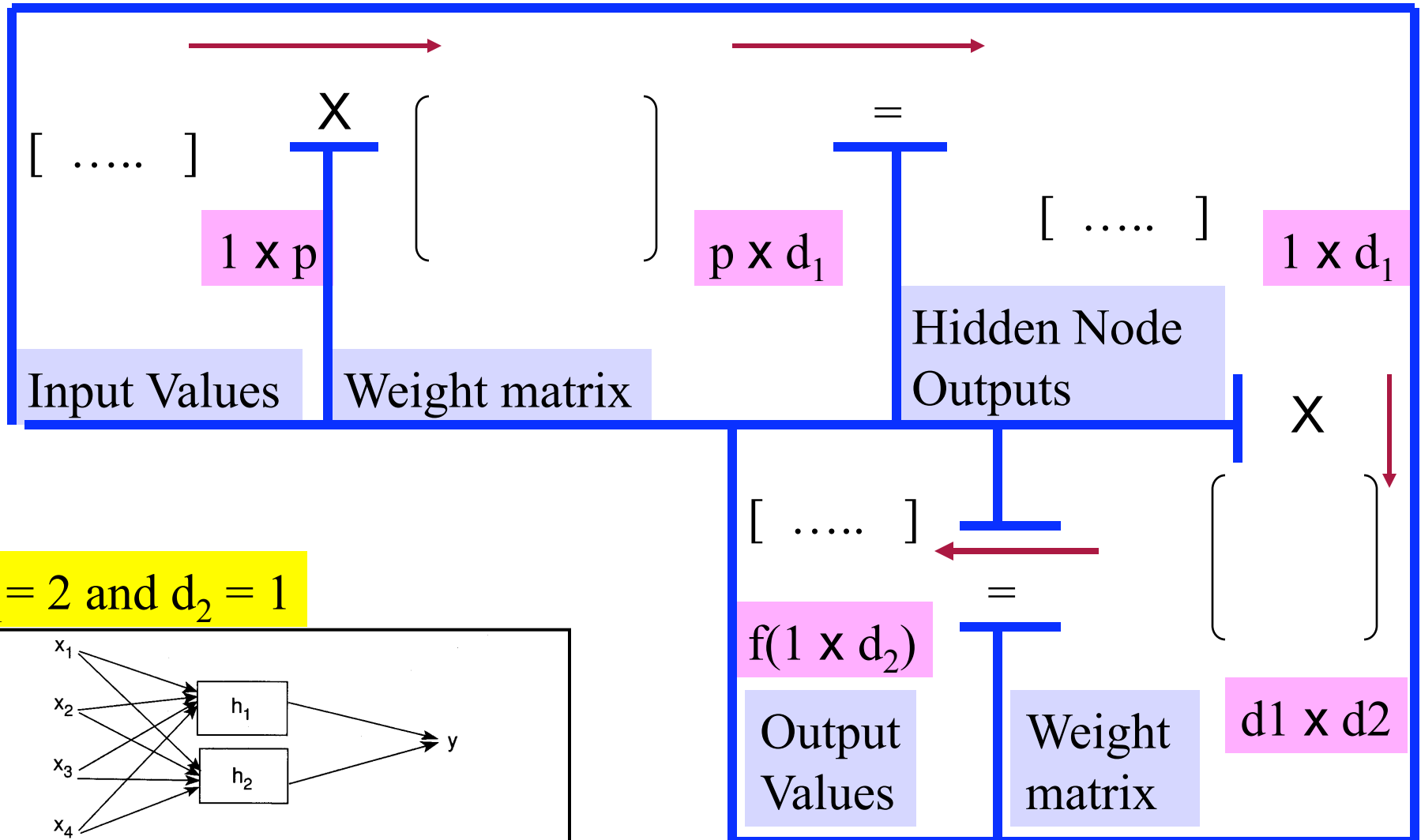
$$h_1 = h(s_1) = \frac{1}{1 + e^{-s_1}}$$

$$h_2 = h(s_2) = \frac{1}{1 + e^{-s_2}}$$

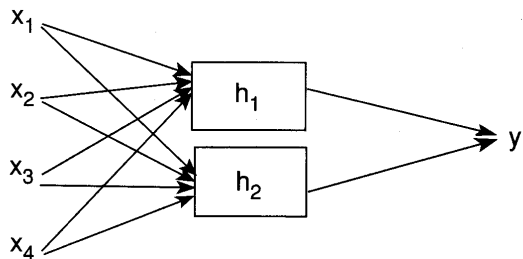
- Output Value

$$y = \sum_{i=1}^2 w_i h_i$$

# MLP in Matrix Notation

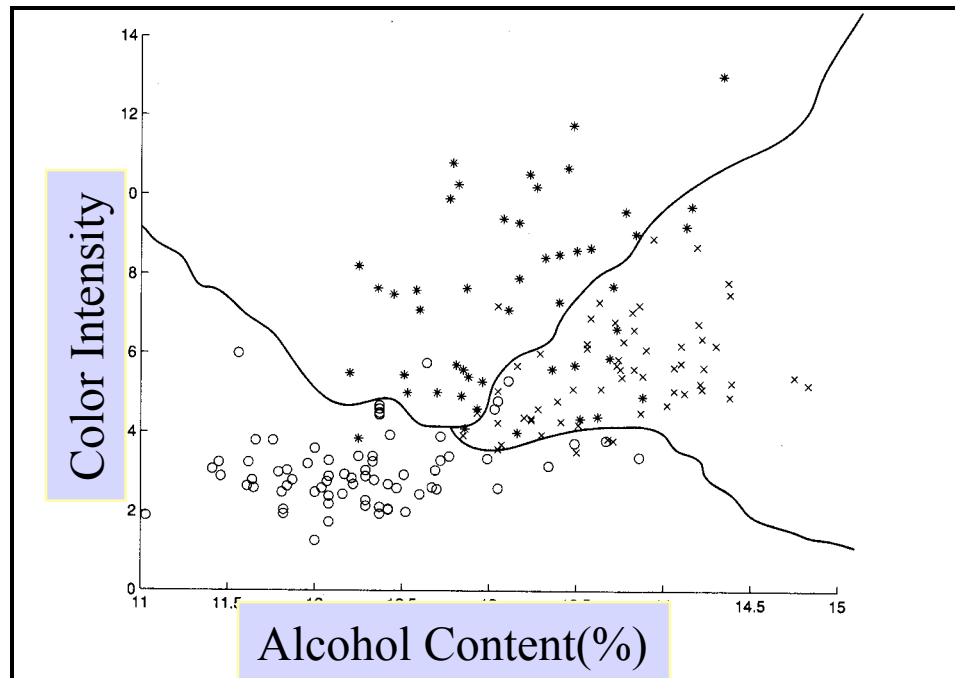


$d_1 = 2$  and  $d_2 = 1$



Multilayer Perceptron with two Hidden nodes ( $d_1=2$ ) and one output node ( $d_2=1$ )

# MLP Result on Wine Data



Highly non-linear  
decision boundaries  
Unlike CART,  
no simple summary  
form to describe  
workings of neural  
network  
model

Type of decision boundaries produced  
by a neural network on wine data

# MLP “algorithm-tuple”

1. Task = prediction: classification or regression
2. Structure = Layers of nonlinear transformations of weighted sums of inputs
3. Score Function = Sum of squared errors
4. Search Method = Steepest descent from random initial parameter values
5. Data Management Technique = online or batch

# MLP Score, Search, Data Mgmt

- Score function

$$S_{SSE} = \sum_{i=1}^n (y(i) - \hat{y}(i))^2$$

True Target Value

Output of Network

- Search
  - Highly nonlinear multivariate optimization
  - Backpropagation uses steepest descent to local minimum
- Data Management
  - On-line (update one data point at a time)
  - Batch mode (update after seeing all data points)