

A Systematic Overview of Data Mining Algorithms

Data Mining Algorithm

- A well-defined procedure that takes data as input and produces output as models or patterns
- *well-defined*: precisely encoded as a finite set of rules.
- *algorithm*: terminates after finite no of steps and produces an output
 - e.g., gradient 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 $X > x_1$ then $prob(Y > y_1) = p_1$

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

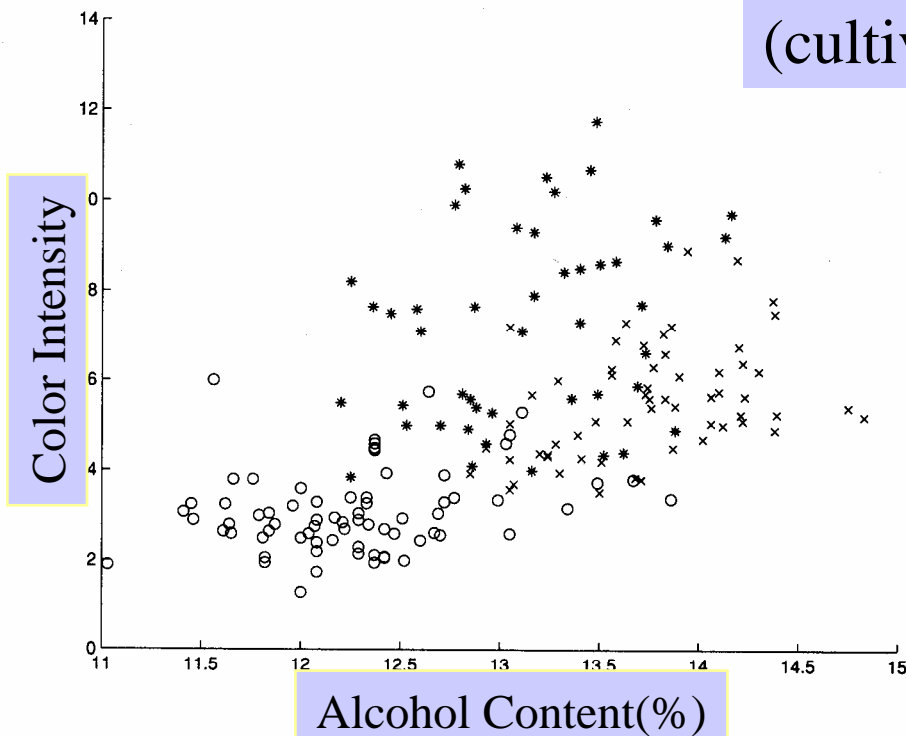
	CART (model)	Backpropagation (parameter est.)	A Priori
Task	Classification and Regression	Classification and Regression	Rule Pattern Discovery
Structure	Decision Tree	Neural Network	Association Rules
Score Function	Cross-validated Loss Function	Squared Error	Support/Accur acy
Search Methods	Greedy Search over Structures	Gradient descent on Parameters	Breadth-First with Pruning
Data Mgmt Techniques	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 label y

Example Task: Wine Classification

Two-dimensional wine data

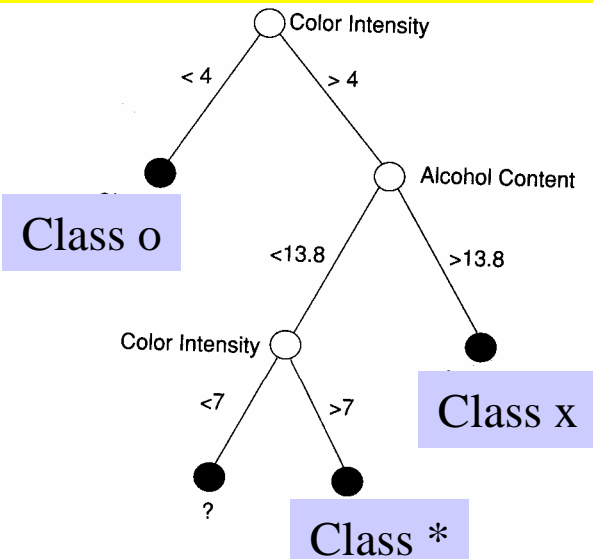


Scatterplot of Data

Data originate from a 13-dimensional data set – each variable measuring a different wine characteristic

Goal is to classify into 3 different wine types (cultivars)

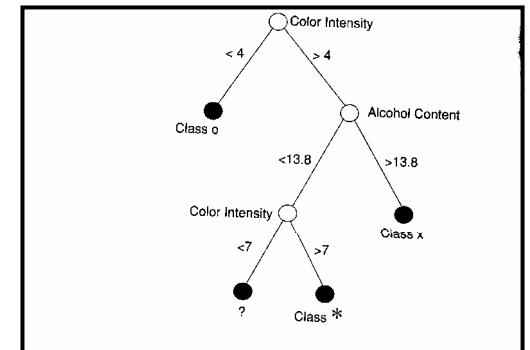
Classification Tree



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

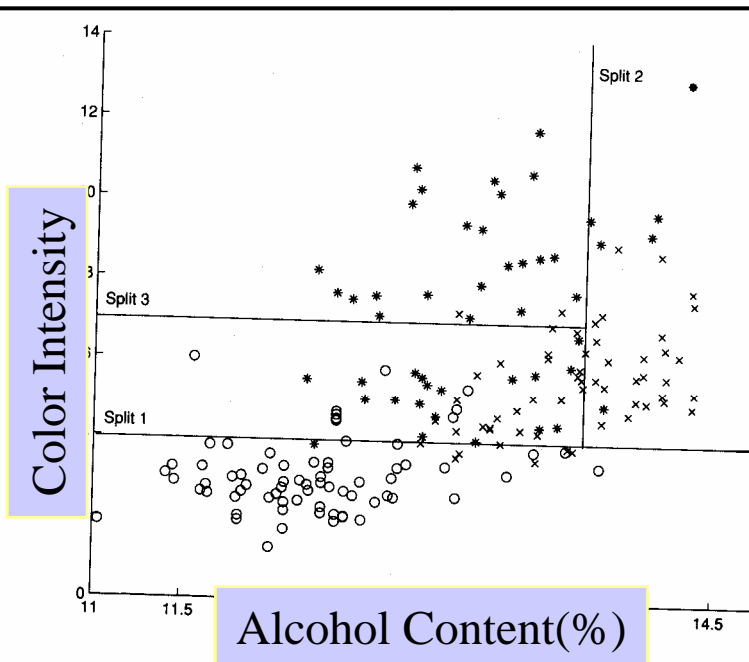
$$\sum_{i=1}^n C(y(i), \hat{y}(i))$$

- Loss incurred when class label for i 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!

Result of CART



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

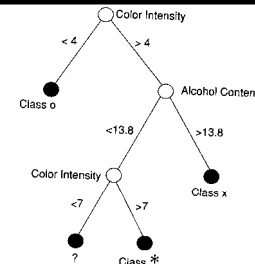


Figure 5.2 A classification tree for the data in Figure 5.1

Classification Tree

CART Scoring/Stopping Criterion

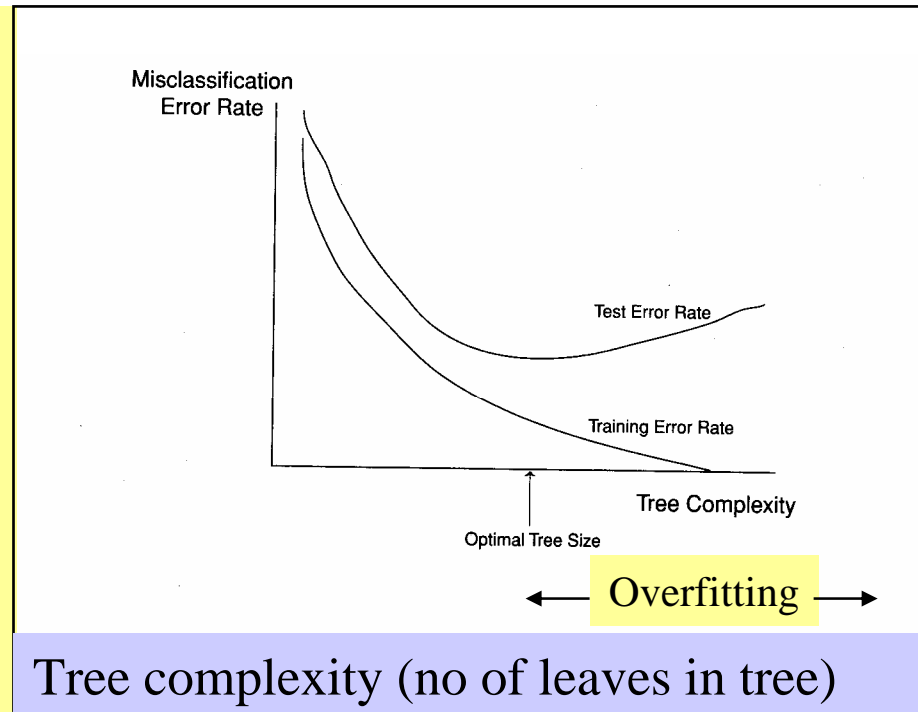
Cross Validation

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 algorithms the data management non-trivial
 - Since it recursively partitions the data set
 - Repeatedly find different subsets of observations in database
 - Naïve implementation will involve 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

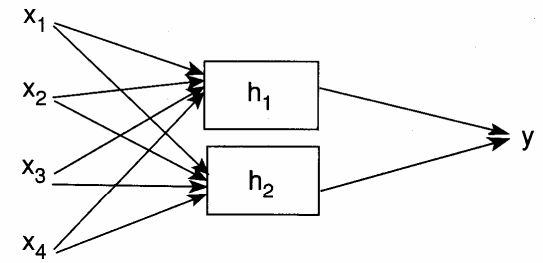
Reductionist Viewpoint applied to 3 algorithms

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 x to real-valued output vector y
- Thus MLP can be used as a nonlinear model for regression as well as for classification

MLP Formulas



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

- From first layer of weights

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

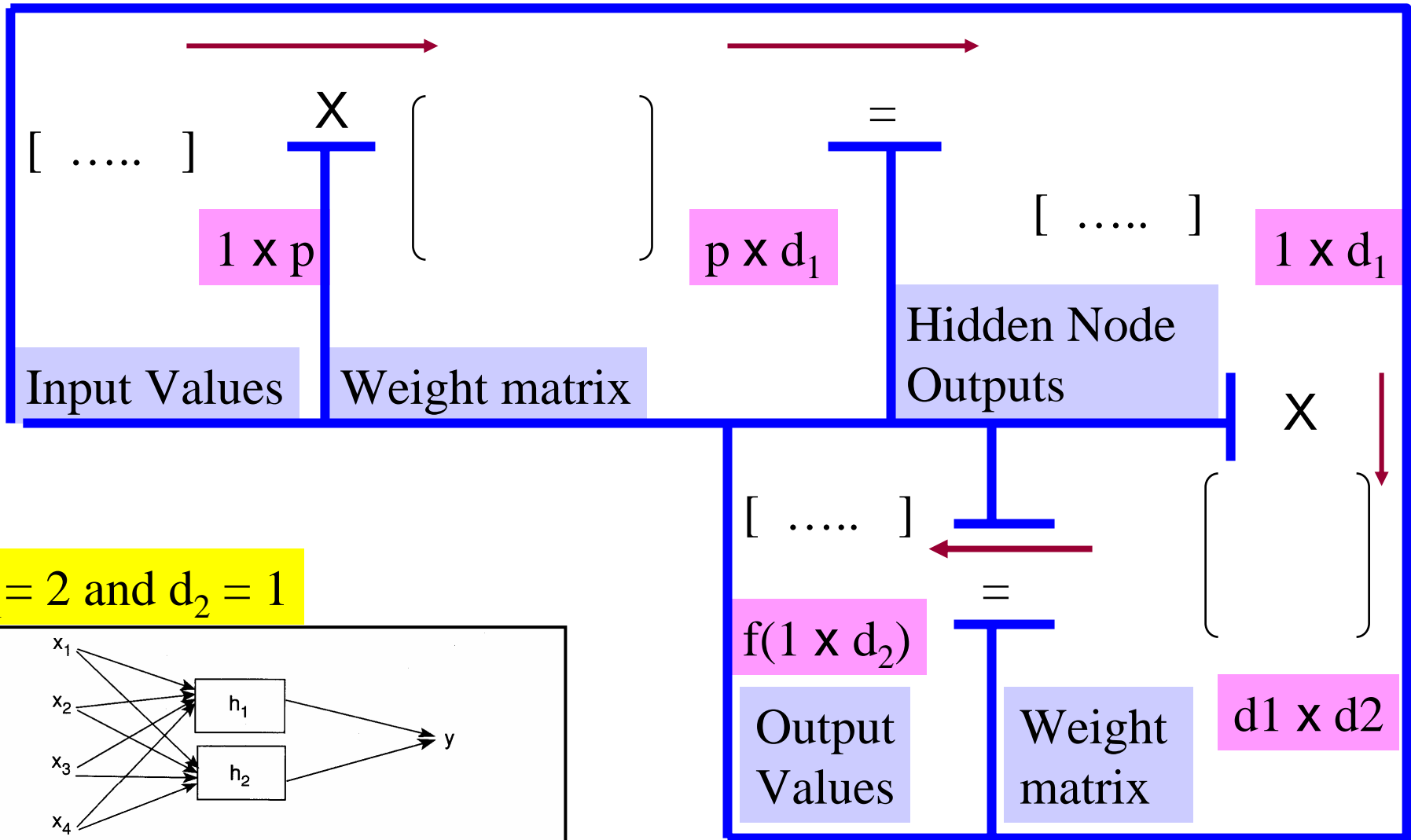
- Non-linear Transformation at hidden nodes

$$h_1 = h(s_1) = \frac{1}{1 + e^{-s_1}} \quad h_2 = h(s_2) = \frac{1}{1 + e^{-s_2}}$$

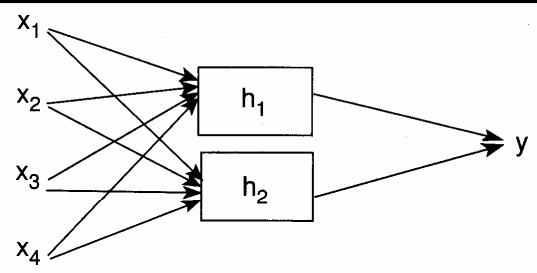
- Output Value

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

MLP in Matrix Notation

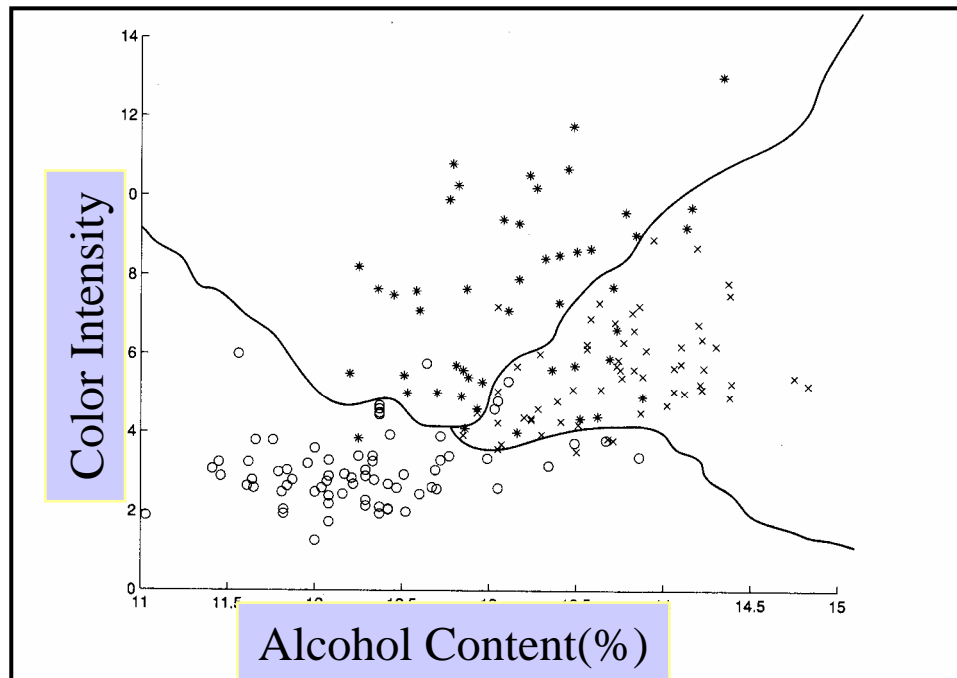


$d_1 = 2$ and $d_2 = 1$



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)