The University of Texas ARLINGTON.

Data Mining Classification: Decision Trees (chapter 3.3)

Instructor: Sharma Chakravarthy sharmac@cse.uta.edu The University of Texas at Arlington

Classification: preamble

Lab

- > It can be used as a predictive model
- For previously unlabeled instances
- It can be used as a descriptive model (rules)
 Explanation of prediction in terms of attribute values
- Explanator of prededon in terms of attribute values
 Especially critical for some applications (e.g., medical diagnosis)
- Remember, one of the complaints about the Neural Networks (ANN, CNN, ...) is that it is non-descriptive!
- The class attribute needs to be nominal (not necessarily binary) or converted to finite ranges!
- Other attributes can be of any types (discrete, continuous, qualitative, quantitative)
- Not all attributes may be relevant; finding optimal combination of attributes that best discriminates instances is challenging!

Classification: preamble

- Classifier and model are often used synonymously
- A classification technique may build a single model, multiple models, or an ensemble of models
- ✤ Decision tree builds a single model
- ✤ k-nearest neighbors does not build an explicit model
- * Ensemble classifiers combine the output of a collection of models
- The induction (building the model) and deduction (predicting on unseen instances) steps are done separately
- Confusion matrix is used for assessing the performance (accuracy and error rate) of a classifier

Classification: Definition

- > Given a collection of records (training set)
- Each record contains a set of *attributes*, one of the attributes being the *class*.
- Find a model for class attribute as a function of the values of other attributes.
- Goal: previously unseen records should be assigned a class as accurately as possible.
 - ♦ A *test set* is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.



Examples of Classification Task

- o Predicting tumor cells as benign or malignant
- Classifying credit card transactions as legitimate or fraudulent
- Classifying secondary structures of protein as alpha-helix, beta-sheet, or random coil
- Categorizing news stories as finance, weather, entertainment, sports, etc.

Classification vs. Prediction

Classification

- o Predicts categorical class labels
- o Most suited for nominal attributes
- o Less effective for ordinal attributes (why?)

Prediction

- models continuous-valued functions or ordinal attributes, i.e., predicts unknown or missing values
- o E.g., Linear regression

Supervised vs. Unsupervised Learning

Supervised learning (e.g., classification)

Supervision: The training data (observations, measurements, etc.) are accompanied by labels indicating the class of the observations
New data is classified based on the training set

Unsupervised learning (e.g., clustering)

- The class labels of training data is unknown
- Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

Classification Techniques

- o Decision Tree based Methods
- o Rule-based Methods
- o Memory based reasoning
- o Neural Networks
- o Naïve Bayes classification
- o Bayesian Belief Networks
- o Tf-idf or Support Vector Machines (SVM)







© 2012 Microsoft Corporation. All rights reserved. Microsoft, Windows, Windows Vista and other product names are or may be registered trademarks and/or trademarks in the U.S. and/or other countries. The information herein is for informational purposes only and represents the current view of Microsoft Corporation as of the date of this presentation. Because Microsoft must respond to changing market conditions, it should not be interpreted to be a commitment on the part of Microsoft. and Microsoft corporation as of the date of this presentation. MICROSOFT MAKES NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, AS TO THE INFORMATION IN THIS PRESENTATION.

 \mathbf{k}















Decision Tree problem
Large search space (why?) what does it depend on?
of tuples in training set? # of attributes?
 Exponential in size, with respect to the set of attributes
 Finding the optimal decision tree is computationally infeasible (why?)
Hence, need an Efficient algorithm for suboptimal decision tree
 Greedy strategy (top-down, recursive, divide-and-conquer) Grow the tree by making locally optimal decisions in selecting the attributes
 Generating an optimal tree is NP-complete
> What roles does the attribute values play?
What roles does the attribute values play? Greened CA072001 The University of Texa at Adments Al Rober Reserved.

Decision Tree size

- Can we say anything about depth, breadth of the tree?
 Worst case?
- What is the effect of the number of attributes
 Determines depth
- What is the effect of attribute valuesBreadth of the tree
- In addition, what else plays a role in complexity?
 Permutation order of attributes to choose and split

Algorithm for Decision Tree Induction

wn, recursive, divide-and

Basic algorithm (a greedy algorithm)

- Tree is constructed in a top-do conquer manner
- At start, all the training examples are at the root
- Attributes are categorical (if continuous-valued, they are discretized in advance)
- Examples are partitioned recursively based on selected attributes
- Test attributes are selected on the basis of a heuristic or attributes are selected on the basis of a heuristic or
- statistical measure (we will discuss a couple of them) Conditions for stopping partitioning
- All samples for a given node belong to the same class
 There are no remaining attributes for further partitioning –
- There are no remaining attributes for further partitioning majority voting is employed for classifying the leaf
- There are no samples left

Decision Tree Induction

> Many Algorithms:

- ♦ Hunt's Algorithm (one of the earliest) (CH 3) -- 1966
- CART (uses only binary splits) (1984) uses GINI
- ◆ ID3 (1986), C4.5 (1993, successor) uses entropy
- ✤ SLIQ, SPRINT



21





Tree Induction

- ➤ Greedy strategy.
- Split the records based on an attribute test that optimizes certain criterion.
- ➤ Issues
- * Determine how to split the records
- How to specify the attribute test condition?
- How to determine the **best split**?
- Determine when to stop splitting

How to Specify Test Condition?

- > Depends on attribute types
- * Categorical vs. Numeric
- Categorical attributes: Nominal, Ordinal
- Numeric attributes: Interval, Ratio
- Discrete vs. Continuous
- > Depends on number of ways to split
- ✤ 2-way split
- ✤ Multi-way split



Splitting Based on Continuous Attribute

Different ways of handling

- * Discretization to form an ordinal categorical attribute
- Static discretize once at the beginning
- Dynamic ranges can be found by equal interval bucketing, equal frequency bucketing (remember histograms) (percentiles), or clustering
- Binary split: $(A \le v)$ or $(A \ge v)$
- consider all possible splits and find the best cut
- can be more compute intensive













Measures of Node Impurity

- 1. Gini Index
- 2. Entropy (information gain)
- 3. Misclassification error
- > Both GINI and Entropy use impurity.
- Both give a zero impurity value if the node contains instances from a single class and maximum impurity value if the node contains equal proportion of instances from multiple classes











Continuous Attributes: Computing Gini Index • Use Binary Decisions based on one value • Several Choices for the splitting value • Number of possible splitting values = Number of distinct values • Each splitting value has a count matrix associated with it • Class counts in each of the partitions, A < v and $A \ge v$ • Simple method to choose best v • For each v, scan the database to gather count matrix and compute its Gini index • Computationally Inefficient! Repetition of twork.



































Tree Induction

Greedy strategy.

Split the records based on an attribute test that optimizes certain criterion.

➤ Issues

- Determine how to split the records
 How to specify the attribute test condition?
 How to determine the best split?
- Determine when to stop splitting

Stopping Criteria for Tree Induction

- Stop expanding a node when all the records belong to the same class
- Stop expanding a node when all the records have similar attribute values
 - □ What to do? majority voting
- Early termination, e.g., when the information gain is below a threshold.

Decision Tree Based Classification

- ➤ Advantages:
- Inexpensive to construct
- Extremely fast at classifying unknown records
- Easy to interpret for small-sized trees
- Accuracy is comparable to other classification techniques for many data sets

Example: C4.5

- Simple depth-first construction.
- Uses Information Gain
- > Sorts Continuous Attributes at each node.
- Needs entire data to fit in memory.
- Unsuitable for Large Datasets.
- Needs out-of-core sorting.
- You can download the software from: <u>http://www.cse.unsw.edu.au/~quinlan/c4.5r8.tar.gz</u>



