# **Data Mining**

(with many slides due to Gehrke, Garofalakis, Rastogi)

### Raghu Ramakrishnan Yahoo! Research

University of Wisconsin–Madison (on leave)

### Introduction

### Definition

Data mining is the exploration and analysis of large quantities of data in order to discover valid, novel, potentially useful, and ultimately understandable patterns in data.

Valid: The patterns hold in general.
Novel: We did not know the pattern beforehand.
Useful: We can devise actions from the patterns.
Understandable: We can interpret and comprehend the patterns.

## Case Study: Bank



- Business goal: Sell more home equity loans
- Current models:
  - Customers with college-age children use home equity loans to pay for tuition
  - Customers with variable income use home equity loans to even out stream of income
- Data:
  - Large data warehouse
  - Consolidates data from 42 operational data sources



- 1. Select subset of customer records who have received home equity loan offer
  - Customers who declined
  - Customers who signed up

Income	Number of	Average Checking	•••	Reponse
	Children	Account Balance		
\$40,000	2	\$1500		Yes
\$75,000	0	\$5000		No
\$50,000	1	\$3000		No
• • •	•••	•••	• • •	•••



- 2. Find rules to predict whether a customer would respond to home equity loan offer
- IF (Salary < 40k) and (numChildren > 0) and (ageChild1 > 18 and ageChild1 < 22) THEN YES



3. Group customers into clusters and investigate clusters





### 4. Evaluate results:

- Many "uninteresting" clusters
- One interesting cluster! Customers with both business and personal accounts; unusually high percentage of likely respondents

### Example: Bank (Contd.)



9

### Action:

• New marketing campaign

### **Result:**

Acceptance rate for home equity offers more than doubled

### **Example Application: Fraud Detection**

- Industries: Health care, retail, credit card services, telecom, B2B relationships
- Approach:
  - Use historical data to build models of fraudulent behavior
  - Deploy models to identify fraudulent instances

## Fraud Detection (Contd.)

- Examples:
  - Auto insurance: Detect groups of people who stage accidents to collect insurance
  - Medical insurance: Fraudulent claims
  - Money laundering: Detect suspicious money transactions (US Treasury's Financial Crimes Enforcement Network)
  - Telecom industry: Find calling patterns that deviate from a norm (origin and destination of the call, duration, time of day, day of week).

### **Other Example Applications**

- CPG: Promotion analysis
- Retail: Category management
- Telecom: Call usage analysis, churn
- Healthcare: Claims analysis, fraud detection
- Transportation/Distribution: Logistics management
- Financial Services: Credit analysis, fraud detection
- Data service providers: Value-added data analysis

### What is a Data Mining Model?

A data mining model is a description of a certain aspect of a dataset. It produces output values for an assigned set of inputs.

### Examples:

- Clustering
- Linear regression model
- Classification model
- Frequent itemsets and association rules
- Support Vector Machines

### **Data Mining Methods**

## Overview

- Several well-studied tasks
  - Classification
  - Clustering
  - Frequent Patterns
- Many methods proposed for each
- Focus in database and data mining community:
  - Scalability

TECS 2007, Data Mining

- Managing the process
- Exploratory analysis

### Classification

Goal:

Learn a function that assigns a record to one of several predefined classes.

Requirements on the model:

- High accuracy
- Understandable by humans, interpretable
- Fast construction for very large training databases

### Classification

### Example application: telemarketing



Copyright 3 1997 United Feature Syndicate, Inc. Redistribution in whole or in part prohibited

### Classification (Contd.)

- Decision trees are one approach to classification.
- Other approaches include:
  - Linear Discriminant Analysis
  - k-nearest neighbor methods
  - Logistic regression
  - Neural networks
  - Support Vector Machines

### **Classification Example**

#### • Training database:

- Two predictor attributes: Age and Car-type (Sport, Minivan and Truck)
- Age is ordered, Car-type is categorical attribute
- Class label indicates whether person bought product
- Dependent attribute is categorical

Age	Car	Class
20	Μ	Yes
30	Μ	Yes
25	Т	No
30	S	Yes
40	S	Yes
20	Т	No
30	Μ	Yes
25	М	Yes
40	Μ	Yes
20	S	No

TECS 2007, Data Mining

## Types of Variables

- *Numerical*: Domain is ordered and can be represented on the real line (e.g., age, income)
- *Nominal* or *categorical*: Domain is a finite set without any natural ordering (e.g., occupation, marital status, race)
- Ordinal: Domain is ordered, but absolute differences between values is unknown (e.g., preference scale, severity of an injury)

### Definitions

- Random variables X<sub>1</sub>, ..., X<sub>k</sub> (predictor variables) and Y (dependent variable)
- X<sub>i</sub> has domain dom(X<sub>i</sub>), Y has domain dom(Y)
- P is a probability distribution on dom(X<sub>1</sub>) x ... x dom(X<sub>k</sub>) x dom(Y) Training database D is a random sample from P
- A predictor d is a function
   d: dom(X<sub>1</sub>) ... dom(X<sub>k</sub>) → dom(Y)

### **Classification Problem**

- If Y is categorical, the problem is a *classification* problem, and we use C instead of Y. |dom(C)| = J, the number of classes.
- C is the *class label*, d is called a *classifier*.
- Let r be a record randomly drawn from P.
   Define the *misclassification rate* of d: RT(d,P) = P(d(r.X<sub>1</sub>, ..., r.X<sub>k</sub>) != r.C)
- <u>Problem definition</u>: Given dataset D that is a random sample from probability distribution P, find classifier d such that RT(d,P) is minimized.

### **Regression Problem**

- If Y is numerical, the problem is a *regression problem*.
- Y is called the dependent variable, d is called a *regression function.*
- Let r be a record randomly drawn from P.
   Define mean squared error rate of d: RT(d,P) = E(r.Y - d(r.X<sub>1</sub>, ..., r.X<sub>k</sub>))<sup>2</sup>
- <u>Problem definition</u>: Given dataset D that is a random sample from probability distribution P, find regression function d such that RT(d,P) is minimized.

### **Regression Example**

#### • Example training database

- Two predictor attributes: Age and Car-type (Sport, Minivan and Truck)
- Spent indicates how much person spent during a recent visit to the web site
- Dependent attribute is *numerical*

Age	Car	Spent
20	М	\$200
30	М	\$150
25	Т	\$300
30	S	\$220
40	S	\$400
20	Т	\$80
30	М	\$100
25	Μ	\$125
40	Μ	\$500
20	S	\$420

### **Decision Trees**

### What are Decision Trees?



TECS 2007, Data Mining

Bee-Chung Chen, Raghu Ramakrishnan, Jude Shavlik, Pradeep Tamm R. Ramakrishnan, Yahoo! Research

### **Decision Trees**

- A *decision tree* T encodes d (a classifier or regression function) in form of a tree.
- A node t in T without children is called a *leaf* node. Otherwise t is called an *internal node*.

### **Internal Nodes**

- Each internal node has an associated splitting predicate. Most common are binary predicates. Example predicates:
  - Age <= 20
  - Profession in {student, teacher}
  - $-5000^{*}$ Age + 3<sup>\*</sup>Salary -10000 > 0

### Internal Nodes: Splitting Predicates

- Binary Univariate splits:
  - Numerical or ordered X: X <= c, c in dom(X)</p>
  - Categorical X: X in A, A subset dom(X)
- Binary Multivariate splits:
  - Linear combination split on numerical variables:  $\Sigma a_i X_i \le C$
- k-ary (k>2) splits analogous

### Leaf Nodes

Consider leaf node t:

- Classification problem: Node t is labeled with one class label c in dom(C)
- Regression problem: Two choices
  - Piecewise constant model:
     t is labeled with a constant y in dom(Y).
  - Piecewise linear model: t is labeled with a linear model  $Y = y_t + \Sigma a_i X_i$

### Example



#### **Encoded classifier:**

- If (age<30 and carType=Minivan) Then YES
- If (age <30 and (carType=Sports or carType=Truck)) Then NO

### **Issues in Tree Construction**

- Three algorithmic components:
  - Split Selection Method
  - Pruning Method
  - Data Access Method

### **Top-Down Tree Construction**

BuildTree(Node *n*, Training database *D*, Split Selection Method **S**)

- [(1) Apply **S** to *D* to find splitting criterion ]
- (1a) **for** each predictor attribute X
- (1b) Call **S**.findSplit(AVC-set of X)
- (1c) endfor
- (1d) S.chooseBest();
- (2) **if** (*n* is not a leaf node) ...

#### S: C4.5, CART, CHAID, FACT, ID3, GID3, QUEST, etc.

### **Split Selection Method**

 Numerical Attribute: Find a split point that separates the (two) classes



### Split Selection Method (Contd.)

• Categorical Attributes: How to group?



### Impurity-based Split Selection Methods

- Split selection method has two parts:
  - Search space of possible splitting criteria.
     Example: All splits of the form "age <= c".</li>
  - Quality assessment of a splitting criterion
- Need to quantify the quality of a split: Impurity function
- Example impurity functions: Entropy, gini-index, chi-square index
#### Data Access Method

 Goal: Scalable decision tree construction, using the complete training database

#### **AVC-Sets**

#### **Training Database**

Age	Car	Class
20	М	Yes
30	Μ	Yes
25	Т	No
30	S	Yes
40	S	Yes
20	Т	No
30	Μ	Yes
25	Μ	Yes
40	Μ	Yes
20	S	No

#### AVC-Sets

Age	Yes	No
20	1	2
25	1	1
30	3	0
40	2	0

Car	Yes	No
Sport	2	1
Truck	0	2
Minivan	5	0

#### Motivation for Data Access Methods



#### In principle, one pass over training database for each node. Can we improve?

#### First scan:



Build AVC-sets for root

TECS 2007, Data Mining

Bee-Chung Chen, Raghu Ramakrishnan, Jude Shavlik, Pradeep TammR. Ramakrishnan, Yahoo! Research



TECS 2007, Data Mining

Bee-Chung Chen, Raghu Ramakrishnan, Jude Shavlik, Pradeep Tamm R. Ramakrishnan, Yahoo! Research



Further optimization: While writing partitions, concurrently build AVC-groups of as many nodes as possible in-memory. This should remind you of Hybrid Hash-Join!



R. Ramakrishnan, Yahoo! Research

#### CLUSTERING

# Problem

- Given points in a multidimensional space, group them into a small number of clusters, using some measure of "nearness"
  - E.g., Cluster documents by topic
  - E.g., Cluster users by similar interests

# Clustering

- Output: (k) groups of records called clusters, such that the records within a group are more similar to records in other groups
  - Representative points for each cluster
  - Labeling of each record with each cluster number
  - Other description of each cluster
- This is unsupervised learning: No record labels are given to learn from
- Usage:
  - Exploratory data mining
  - Preprocessing step (e.g., outlier detection)

#### Clustering (Contd.)

- Example input database: Two numerical variables
- How many groups are here?



#### Customer Demographics



e-Chung Chen, Raghu Ramakrishnan, Jude Shavlik, Pradeep Tamm R. Ramakrishnan, Yahoo! Research 47

# Improve Search Using Topic Hierarchies

• Web directories (or topic hierarchies) provide a hierarchical classification of documents (e.g., Yahoo!)



- Searches performed in the context of a topic restricts the search to only a subset of web pages related to the topic
- Clustering can be used to generate topic hierarchies

# Clustering (Contd.)

- Requirements: Need to define "similarity" between records
- Important: Use the "right" similarity (distance) function
  - Scale or normalize all attributes. Example: seconds, hours, days
  - Assign different weights to reflect importance of the attribute
  - Choose appropriate measure (e.g., L1, L2)

#### Distance Measure D

- For 2 pts x and y:
  - $D(\mathbf{x},\mathbf{x}) = 0$
  - D(x,y) = D(y,x)
  - $D(x,y) \le D(x,z)+D(z,y)$ , for all z
- Examples, for x,y in k-dim space:
  - L1: Sum of |xi-yi| over I = 1 to k
  - L2: Root-mean squared distance

#### Approaches

- Centroid-based: Assume we have k clusters, guess at the centers, assign points to nearest center, e.g., K-means; over time, centroids shift
- Hierarchical: Assume there is one cluster per point, and repeatedly merge nearby clusters using some distance threshold

# Scalability: Do this with fewest number of passes over data, ideally, sequentially

#### K-means Clustering Algorithm

- Choose *k* initial means
- Assign each point to the cluster with the closest mean
- Compute new mean for each cluster
- Iterate until the *k* means stabilize

#### Agglomerative Hierarchical Clustering Algorithms

- Initially each point is a distinct cluster
- Repeatedly merge closest clusters until the number of clusters becomes k

- Closest: dmean (Ci, Cj) = 
$$||m_i - m_j||$$
  
dmin (Ci, Cj) =  $\min_{p \in C_i, q \in C_j} ||p-q||$ 

Likewise dave (Ci, Cj) and dmax (Ci, Cj)

#### Scalable Clustering Algorithms for Numeric Attributes

CLARANS DBSCAN BIRCH CLIQUE CURE

• Above algorithms can be used to cluster documents after reducing their dimensionality using SVD

#### Birch [ZRL96]

#### Pre-cluster data points using "CF-tree" data structure



# BIRCH [ZRL 96]

- Pre-cluster data points using "CF-tree" data structure
  - CF-tree is similar to R-tree
  - For each point
    - CF-tree is traversed to find the closest cluster
    - If the cluster is within epsilon distance, the point is absorbed into the cluster
    - Otherwise, the point starts a new cluster
- Requires only single scan of data
- Cluster summaries stored in CF-tree are given to main memory clustering algorithm of choice

#### Background

*Given a cluster of instances*  $\{\vec{X}_i\}$ *, we define:* 

Centroid 
$$\vec{X0} = \frac{\sum_{i=1}^{N} \vec{X_i}}{N}$$
  
Radius  $R = (\frac{\sum_{i=1}^{N} (\vec{X_i} - \vec{X0})^2}{N})^{\frac{1}{2}}$   
Diameter  $D = (\frac{\sum_{i=1}^{N} \sum_{j=1}^{N} (\vec{X_i} - \vec{X_j})^2}{N(N-1)})^{\frac{1}{2}}$   
(Euclidean) Distance  $D0 = ((\vec{X0}_1 - \vec{X0}_2)^2)^{\frac{1}{2}}$ 

#### The Algorithm: Background

We define the **Euclidean** and **Manhattan** distance between any two clusters as:

$$D0 = ((\vec{X0}_1 - \vec{X0}_2)^2)^{\frac{1}{2}}$$
$$D1 = |\vec{X0}_1 - \vec{X0}_2| = \sum_{i=1}^d |\vec{X0}_1^{(i)} - \vec{X0}_1^{(i)}|$$

TECS 2007, Data Mining

Bee-Chung Chen, Raghu Ramakrishnan, Jude Shavlik, Pradeep TammR. Ramakrishnan, Yahoo! Research

#### Clustering Feature (CF)

Given a cluster  $\{\vec{X_1}, \vec{X_2}, \dots, \vec{X_N}\}$   $\mathbf{CF} = (N, \vec{LS}, SS)$ 

$$N$$
 is the number of data points  
 $\vec{LS} = \sum_{i=1}^{N} \vec{X_i}$   
 $SS = \sum_{i=1}^{N} \vec{X_i}^2$ 

 $\mathbf{CF_1} + \mathbf{CF_2} = (N_1 + N_2, \vec{LS_1} + \vec{LS_2}, SS_1 + SS_2)$ 

#### Allows incremental merging of clusters!

#### Points to Note

- Basic algorithm works in a single pass to condense metric data using spherical summaries
  - Can be incremental
- Additional passes cluster CFs to detect nonspherical clusters
- Approximates density function
- Extensions to non-metric data

# **CURE [GRS 98]**

- Hierarchical algorithm for dicovering arbitrary shaped clusters
  - Uses a small number of representatives per cluster
  - Note:
    - Centroid-based: Uses 1 point to represent a cluster => Too little information ... Hyper-spherical clusters
    - MST-based: Uses every point to represent a cluster =>Too much information ... Easily mislead
- Uses random sampling
- Uses Partitioning
- Labeling using representatives

#### **Cluster Representatives**

A representative set of points:

- Small in number : c
- Distributed over the cluster
- Each point in cluster is close to one representative
- Distance between clusters:

smallest distance between representatives

Market Basket Analysis: Frequent Itemsets

# Market Basket Analysis

- Consider shopping cart filled with several items
- Market basket analysis tries to answer the following questions:
  - Who makes purchases
  - What do customers buy

# Market Basket Analysis

#### • Given:

- A database of customer transactions
- Each transaction is a set of items
- Goal:
  - Extract rules

TID	CID	Date	Item	Qty
111	201	5/1/99	Pen	2
111	201	5/1/99	Ink	1
111	201	5/1/99	Milk	3
111	201	5/1/99	Juice	6
112	105	6/3/99	Pen	1
112	105	6/3/99	Ink	1
112	105	6/3/99	Milk	1
113	106	6/5/99	Pen	1
113	106	6/5/99	Milk	1
114	201	7/1/99	Pen	2
114	201	7/1/99	Ink	2
114	201	7/1/99	Juice	4

#### Market Basket Analysis (Contd.)

- Co-occurrences
  - 80% of all customers purchase items X, Y and Z together.
- Association rules
  - 60% of all customers who purchase X and Y also buy Z.
- Sequential patterns
  - 60% of customers who first buy X also purchase Y within three weeks.

#### **Confidence and Support**

We prune the set of all possible association rules using two interestingness measures:

• Confidence of a rule:

- X => Y has confidence c if P(Y|X) = c

• Support of a rule:

- X => Y has support s if P(XY) = s

We can also define

• Support of a co-ocurrence XY:

- XY has support s if P(XY) = s

### Example

- Example rule: {Pen} => {Milk}
   Support: 75%
   Confidence: 75%
- Another example: {Ink} => {Pen} Support: 100% Confidence: 100%

TID	CID	Date	Item	Qty
111	201	5/1/99	Pen	2
111	201	5/1/99	Ink	1
111	201	5/1/99	Milk	3
111	201	5/1/99	Juice	6
112	105	6/3/99	Pen	1
112	105	6/3/99	Ink	1
112	105	6/3/99	Milk	1
113	106	6/5/99	Pen	1
113	106	6/5/99	Milk	1
114	201	7/1/99	Pen	2
114	201	7/1/99	Ink	2
114	201	7/1/99	Juice	4

#### Exercise

Can you find all itemsets
 with

support >= 75%?

TID	CID	Date	Item	Qty
111	201	5/1/99	Pen	2
111	201	5/1/99	Ink	1
111	201	5/1/99	Milk	3
111	201	5/1/99	Juice	6
112	105	6/3/99	Pen	1
112	105	6/3/99	Ink	1
112	105	6/3/99	Milk	1
113	106	6/5/99	Pen	1
113	106	6/5/99	Milk	1
114	201	7/1/99	Pen	2
114	201	7/1/99	Ink	2
114	201	7/1/99	Juice	4

#### Exercise

 Can you find all association rules with support >= 50%?

TID	CID	Date	Item	Qty
111	201	5/1/99	Pen	2
111	201	5/1/99	Ink	1
111	201	5/1/99	Milk	3
111	201	5/1/99	Juice	6
112	105	6/3/99	Pen	1
112	105	6/3/99	Ink	1
112	105	6/3/99	Milk	1
113	106	6/5/99	Pen	1
113	106	6/5/99	Milk	1
114	201	7/1/99	Pen	2
114	201	7/1/99	Ink	2
114	201	7/1/99	Juice	4

# Extensions

- Imposing constraints
  - Only find rules involving the dairy department
  - Only find rules involving expensive products
  - Only find rules with "whiskey" on the right hand side
  - Only find rules with "milk" on the left hand side
  - Hierarchies on the items
  - Calendars (every Sunday, every 1<sup>st</sup> of the month)

### Market Basket Analysis: Applications

- Sample Applications
  - Direct marketing
  - Fraud detection for medical insurance
  - Floor/shelf planning
  - Web site layout
  - Cross-selling
## **DBMS Support for DM**

#### Why Integrate DM into a DBMS?



#### **Integration Objectives**

- Avoid isolation of querying from mining
  - Difficult to do "ad-hoc" mining
- Provide simple programming approach to creating and using DM models

- Make it possible to add new models
- Make it possible to add new, scalable algorithms

#### Analysts (users)

DM Vendors

# SQL/MM: Data Mining

- A collection of classes that provide a standard interface for invoking DM algorithms from SQL systems.
- Four data models are supported:
  - Frequent itemsets, association rules
  - Clusters
  - Regression trees
  - Classification trees

#### DATA MINING SUPPORT IN MICROSOFT SQL SERVER \*

#### \* Thanks to Surajit Chaudhuri for permission to use/adapt his slides

TECS 2007, Data Mining

Bee-Chung Chen, Raghu Ramakrishnan, Jude Shavlik, Pradeep Tamm R. Ramakrishnan, Yahoo! Research 77

# **Key Design Decisions**

#### • Adopt relational data representation

 A Data Mining Model (DMM) as a "tabular" object (externally; can be represented differently internally)

#### • Language-based interface

- Extension of SQL
- Standard syntax

# **DM Concepts to Support**

- Representation of input (cases)
- Representation of *models*
- Specification of *training step*
- Specification of *prediction step*

#### Should be independent of specific algorithms

#### What are "Cases"?

- DM algorithms analyze "cases"
- The "case" is the entity being categorized and classified
- Examples
  - Customer credit risk analysis: Case = Customer
  - Product profitability analysis: Case = Product
  - Promotion success analysis: Case = Promotion
- Each case encapsulates all we know about the entity

#### Cases as Records: Examples

Cust ID	Age	Marital Status	Wealth
1	35	М	380,000
2	20	S	50,000
3	57	М	470,000

Age	Car	Class
20	Μ	Yes
30	Μ	Yes
25	Т	No
30	S	Yes
40	S	Yes
20	Т	No
30	Μ	Yes
25	Μ	Yes
40	Μ	Yes
20	S	No

# **Types of Columns**

Cust ID	Ane	Marital	Wealth	Pro	duct Purch	ases
	Age	Status	Weatth	Product	Quantity	Туре
1	35	М	380,000	TV	1	Appliance
				Coke	6	Drink
				Ham	3	Food

- <u>Keys</u>: Columns that uniquely identify a case
- Attributes: Columns that describe a case
  - Value: A state associated with the attribute in a specific case
  - Attribute Property: Columns that describe an attribute
    - Unique for a specific attribute value (TV is always an appliance)
  - Attribute Modifier: Columns that represent additional "meta" information for an attribute
    - Weight of a case, Certainty of prediction

# More on Columns

- Properties describe attributes
  - Can represent generalization hierarchy
- Distribution information associated with attributes
  - Discrete/Continuous
  - Nature of Continuous distributions
    - Normal, Log\_Normal
  - Other Properties (e.g., ordered, not null)



- Special parameters
- Model is represented as a (nested) table
  - Specification = Create table
  - Training = Inserting data into the table
  - Predicting = Querying the table



#### CREATE MINING MODEL

CREATE MINING MODEL [Age Prediction]

```
[Customer ID] LONG
                       KEY,
[Gender]
                      TEXT
                              DISCRETE
                                         ATTRIBUTE,
[Age]
                      DOUBLE CONTINUOUS ATTRIBUTE PREDICT,
[ProductPurchases] TABLE (
[ProductName]
              TEXT
                     KEY,
[Quantity]
                      DOUBLE NORMAL CONTINUOUS,
[ProductType] TEXT DISCRETE RELATED TO [ProductName]
USING [Microsoft Decision Tree]
```

Note that the ProductPurchases column is a nested table. SQL Server computes this field when data is "inserted".

# Training a DMM

- Training a DMM requires passing it "known" cases
- Use an INSERT INTO in order to "insert" the data to the DMM
  - The DMM will usually not retain the inserted data
  - Instead it will analyze the given cases and build the DMM content (decision tree, segmentation model)
    - INSERT [INTO] <mining model name>
       [(columns list)]
       <source data guery>

#### **INSERT INTO**

```
INSERT INTO [Age Prediction]
(
[Gender],[Hair Color], [Age]
)
OPENQUERY([Provider=MSOLESQL...,
`SELECT
    [Gender], [Hair Color], [Age]
FROM [Customers]'
)
```

# **Executing Insert Into**

- The DMM is trained
  - The model can be retrained or incrementally refined
- Content (rules, trees, formulas) can be explored
- Prediction queries can be executed

## What are Predictions?

- Predictions apply the trained model to estimate missing attributes in a data set
- Predictions = Queries
- Specification:
  - Input data set
  - A trained DMM (think of it as a truth table, with one row per combination of predictor-attribute values; this is only conceptual)
  - Binding (mapping) information between the input data and the DMM

## **Prediction Join**

```
SELECT [Customers].[ID],
MyDMM.[Age],
PredictProbability(MyDMM.[Age])
FROM
MyDMM PREDICTION JOIN [Customers]
ON MyDMM.[Gender] = [Customers].[Gender] AND
MyDMM.[Hair Color] =
[Customers].[Hair Color]
```

#### Exploratory Mining: Combining OLAP and DM

### **Databases and Data Mining**

- What can database systems offer in the grand challenge of understanding and learning from the flood of data we've unleashed?
  - The plumbing
  - Scalability

## **Databases and Data Mining**

- What can database systems offer in the grand challenge of understanding and learning from the flood of data we've unleashed?
  - The plumbing
  - Scalability
  - Ideas!
    - Declarativeness
    - Compositionality
    - Ways to conceptualize your data

## **Multidimensional Data Model**

- One fact table  $\Delta = (\mathbf{X}, \mathbf{M})$ 
  - $X = X_1, X_2, \dots$  Dimension attributes
  - $M=M_1, M_2, \dots$  <u>Measure attributes</u>
- Domain hierarchy for each dimension attribute:
  - Collection of domains  $Hier(X_i) = (D_i^{(1)}, \dots, D_i^{(k)})$
  - The extended domain:  $EX_i = \bigcup_{1 \le k \le t} DX_i^{(k)}$
- Value mapping function:  $\gamma_{D1 \rightarrow D2}(x)$ 
  - − e.g.,  $\gamma_{month \rightarrow year}(12/2005) = 2005$
  - Form the value hierarchy graph
  - Stored as dimension table attribute (e.g., week for a time value) or conversion functions (e.g., month, quarter)

#### **Multidimensional Data**



## **Cube Space**

- Cube space:  $C = EX_1 \times EX_2 \times ... \times EX_d$
- Region: Hyper rectangle in cube space  $-c = (v_1, v_2, ..., v_d)$ ,  $v_i \in EX_i$
- Region granularity:
  - $\text{gran}(c) = (d_1, d_2, ..., d_d), d_i = \text{Domain}(c.v_i)$
- Region coverage:

 $- \operatorname{coverage}(c) = \operatorname{all} \operatorname{facts} \operatorname{in} c$ 

• Region set: All regions with same granularity

#### **OLAP Over Imprecise Data**

with Doug Burdick, Prasad Deshpande, T.S. Jayram, and Shiv Vaithyanathan In VLDB 05, 06 joint work with IBM Almaden

#### **Imprecise** Data



## **Querying Imprecise Facts**



How do we treat p5?



FactID	Auto	Loc	Repair	
p1	F150	NY	100	
p2	Sierra	NY	500	
р3	F150	MA	100	
p4	Sierra	MA	200	
p5	Truck	MA	100	

### Allocation (1)



FactID	Auto	Loc	Repair
p1	F150	NY	100
p2	Sierra	NY	500
р3	F150	MA	100
p4	Sierra	MA	200
p5	Truck	MA	100

## Allocation (2)

#### (Huh? Why 0.5 / 0.5? - Hold on to that thought)



ID	FactID	Auto	Loc	Repair	Weight
1	p1	F150	NY	100	1.0
2	p2	Sierra	NY	500	1.0
3	р3	F150	MA	100	1.0
4	p4	Sierra	MA	200	1.0
5	р5	F150	MA	100	0.5
6	р5	Sierra	MA	100	0.5

## Allocation (3)



## **Allocation Policies**

- The procedure for assigning allocation weights is referred to as an allocation policy:
  - Each allocation policy uses different information to assign allocation weights
  - Reflects assumption about the correlation structure in the data
    - Leads to EM-style iterative algorithms for allocating imprecise facts, maximizing likelihood of observed data

#### Allocation Policy: Count



#### Allocation Policy: Measure



ID	Sales
p1	100
p2	150
р3	300
p4	200
p5	250
p6	400

#### **Allocation Policy Template**



Bee-Chung Chen, Raghu Ramakrishnan, Jude Shavlik, Pradeep Tamm R. Ramakrishnan, Yahoo! Research 107

## What is a Good Allocation Policy?


#### **Desideratum I: Consistency**



Consistency specifies the relationship between answers to related queries on a fixed data set

#### **Desideratum II: Faithfulness**



 Faithfulness specifies the relationship between answers to a fixed query on related data sets

# **Results on Query Semantics**

- Evaluating queries over extended data model yields expected value of the aggregation operator over all possible worlds
- Efficient query evaluation algorithms available for SUM, COUNT; more expensive dynamic programming algorithm for AVERAGE
  - Consistency and faithfulness for SUM, COUNT are satisfied under appropriate conditions
  - (Bound-)Consistency does not hold for AVERAGE, but holds for E(SUM)/E(COUNT)
    - Weak form of faithfulness holds
  - Opinion pooling with LinOP: Similar to AVERAGE

# **Allocation Policies**

- Procedure for assigning allocation weights is referred to as an allocation policy
  - Each allocation policy uses different information to assign allocation weight
- Key contributions:
  - Appropriate characterization of the large space of allocation policies (VLDB 05)
  - Designing efficient algorithms for allocation policies that take into account the correlations in the data (VLDB 06)



# **Query Semantics**

- Given all possible worlds together with their probabilities, queries are easily answered using expected values
  - But number of possible worlds is exponential!
- Allocation gives facts weighted assignments to possible completions, leading to an extended version of the data
  - Size increase is linear in number of (completions of) imprecise facts
  - Queries operate over this extended version

#### Exploratory Mining: Prediction Cubes

#### with Beechun Chen, Lei Chen, and Yi Lin In VLDB 05; EDAM Project

# The Idea

- Build OLAP data cubes in which cell values represent decision/prediction behavior
  - In effect, build a tree for each cell/region in the cube observe that this is not the same as a collection of trees used in an ensemble method!
  - The idea is simple, but it leads to promising data mining tools
  - Ultimate objective: Exploratory analysis of the entire space of "data mining choices"
    - Choice of algorithms, data conditioning parameters ...

### Example (1/7): Regular OLAP

**Goal:** Look for patterns of unusually high numbers of applications:

Z: Dimensions	Y: Measure
---------------	------------

Location	Time	# of App
AL, <b>USA</b>	Dec, 04	2
WY, <b>USA</b>	Dec, 04	3



Bee-Chung Chen, Raghu Ramakrishnan, Jude Shavlik, Pradeep TammR. Ramakrishnan, Yahoo! Research 117

### Example (2/7): Regular OLAP

**Goal:** Look for patterns of unusually high numbers of applications:

Cell value: Number of loan applications

Z: Dimensions Y: Measure

Location	Time	# of App	
AL, <b>USA</b>	Dec, 04	2	
WY, <b>USA</b>	Dec, 04	3	



#### Finer regions

# Example (3/7): Decision Analysis

Goal: Analyze a bank's loan decision process w.r.t. two dimensions: *Location* and *Time* 

#### Fact table **D**



Bee-Chung Chen, Raghu Ramakrishnan, Jude Shavlik, Pradeep TammR. Ramakrishnan, Yahoo! Research 119

# Example (3/7): Decision Analysis

- Are there branches (and time windows) where approvals were closely tied to sensitive attributes (e.g., race)?
  - Suppose you partitioned the training data by location and time, chose the partition for a given branch and time window, and built a classifier. You could then ask, "Are the predictions of this classifier closely correlated with race?"
- Are there branches and times with decision making reminiscent of 1950s Alabama?
  - Requires comparison of classifiers trained using different subsets of data.

### Example (4/7): Prediction Cubes



- 1. Build a model using data from USA in Dec., 1985
- 2. Evaluate that model

Measure in a cell:

- Accuracy of the model
- Predictiveness of Race measured based on that model
- Similarity between that model and a given model

Model  $h(\mathbf{X}, \sigma_{[\text{USA, Dec 04}]}(\mathbf{D}))$ E.g., decision tree

## Example (5/7): Model-Similarity



TECS 2007, Data Mining

# Example (6/7): Predictiveness

#### Given: Data table **D** Location Time Sex Race Approval - Data table **D** ... - Attributes V AL, USA Dec, 04 White Μ Yes ... ... - Test set $\Delta$ w/o labels . . . ... ... WY, USA Dec. 04 Black No . . . 2004 2003 ... Dec Jan Dec Jan ... ... ... 0.2 0.3 0.6 0.5 CA 0.4 ·... ... Yes Yes **Build models** 0.9 **USA** 0.2 0.3 No No ... ... ... <u>×</u> ... ... ... ... ... ... ... No Yes $h(\boldsymbol{X})$ h(X-V)Level: [Country, Month] Race Sex ... F Predictiveness of V White ... ... ... ... Μ Black Race was an important predictor of loan ... approval decision in USA during Dec 04

## Model Accuracy

 A probabilistic view of classifiers: A dataset is a random sample from an underlying pdf p\*(X, Y), and a classifier

 $h(\mathbf{X}; \mathbf{D}) = \operatorname{argmax}_{y} p^{*}(Y=y | \mathbf{X}=\mathbf{x}, \mathbf{D})$ 

- i.e., A classifier approximates the pdf by predicting the "most likely" y value
- Model Accuracy:
  - $E_{\mathbf{x},y}[I(h(\mathbf{x}; \mathbf{D}) = y)]$ , where  $(\mathbf{x}, y)$  is drawn from  $p^*(\mathbf{X}, Y | \mathbf{D})$ , and  $I(\Psi) = 1$  if the statement  $\Psi$  is true;  $I(\Psi) = 0$ , otherwise
  - In practice, since p\* is an unknown distribution, we use a set-aside test set or cross-validation to estimate model accuracy.

# Model Similarity

The prediction similarity between two models, h1(X) and h2(X), on test set Δ is

$$\frac{1}{|\Delta|} \sum_{\mathbf{x} \in \Delta} I(h_1(\mathbf{x}) = h_2(\mathbf{x}))$$

 The KL-distance between two models, h1(X) and h2(X), on test set Δ is

$$\frac{1}{|\Delta|} \sum_{\mathbf{x} \in \Delta} \sum_{y} p_{h_1}(y \mid x) \log \frac{p_{h_1}(y \mid x)}{p_{h_2}(y \mid x)}$$

### **Attribute Predictiveness**

Intuition: V ⊆ X is not predictive if and only if V is independent of Y given the other attributes X – V; i.e.,

$$p^*(Y \mid \boldsymbol{X} - \boldsymbol{V}, \, \boldsymbol{D}) = p^*(Y \mid \boldsymbol{X}, \, \boldsymbol{D})$$

- In practice, we can use the distance between h(X; D) and h(X V; D)
- Alternative approach: Test if h(X; D) is more accurate than h(X V; D) (e.g., by using cross-validation to estimate the two model accuracies involved)

#### Example (7/7): Prediction Cube

	2004		2003				
	Jan		Dec	Jan		Dec	
CA	0.4	0.1	0.3	0.6	0.8		
USA	0.7	0.4	0.3	0.3			



	04	03	
CA	0.3	0.2	
USA	0.2	0.3	

#### Cell value: Predictiveness of Race



s of Race		2004		2003		•••		
		Jan		Dec	Jan		Dec	
	AB	0.4	0.2	0.1	0.1	0.2		
СА		0.1	0.1	0.3	0.3			
	ΥT	0.3	0.2	0.1	0.2			
	AL	0.2	0.1	0.2				
JSA		0.3	0.1	0.1				
	WY	0.9	0.7	0.8				

# **Efficient Computation**

- Reduce prediction cube computation to data cube computation
  - Represent a data-mining model as a distributive or algebraic (bottom-up computable) aggregate function, so that data-cube techniques can be directly applied

# Bottom-Up Data Cube Computation



#### Cell Values: Numbers of Ioan applications

# **Scoring Function**

- Represent a model as a function of sets
- Conceptually, a machine-learning model h(X; σ<sub>z</sub>(D)) is a scoring function Score(y, x; σ<sub>z</sub>(D)) that gives each class y a score on test example x
  - $h(\mathbf{x}; \sigma_{\mathbf{z}}(\mathbf{D})) = \operatorname{argmax}_{y} \operatorname{Score}(y, \mathbf{x}; \sigma_{\mathbf{z}}(\mathbf{D}))$
  - Score( $y, \mathbf{x}; \sigma_{\mathbf{z}}(\mathbf{D})$ )  $\approx p(y \mid \mathbf{x}, \sigma_{\mathbf{z}}(\mathbf{D}))$
  - $-\sigma_z(\mathbf{D})$ : The set of training examples (a cube subset of **D**)

# **Machine-Learning Models**

- Naïve Bayes:
  - Scoring function: algebraic
- Kernel-density-based classifier:
  - Scoring function: distributive
- Decision tree, random forest:
  - Neither distributive, nor algebraic
- PBE: Probability-based ensemble (new)
  - To make any machine-learning model distributive
  - Approximation

### **Efficiency Comparison**



#### Bellwether Analysis: Global Aggregates from Local Regions

#### with Beechun Chen, Jude Shavlik, and Pradeep Tamma In VLDB 06

TECS 2007, Data Mining

Bee-Chung Chen, Raghu Ramakrishnan, Jude Shavlik, Pradeep TammR. Ramakrishnan, Yahoo! Research 133

# Motivating Example

- A company wants to predict the first year worldwide profit of a new item (e.g., a new movie)
  - By looking at features and profits of previous (similar) movies, we predict expected total profit (1-year US sales) for new movie
    - Wait a year and write a query! If you can't wait, stay awake ...
  - The most predictive "features" may be based on sales data gathered by releasing the new movie in many "regions" (different locations over different time periods).
    - Example "region-based" features: 1<sup>st</sup> week sales in Peoria, week-toweek sales growth in Wisconsin, etc.
    - Gathering this data has a cost (e.g., marketing expenses, waiting time)
- **Problem statement:** Find the most predictive region features that can be obtained within a given "cost budget"

# Key Ideas

- Large datasets are rarely labeled with the targets that we wish to learn to predict
  - But for the tasks we address, we can readily use OLAP queries to generate features (e.g., 1<sup>st</sup> week sales in Peoria) and even targets (e.g., profit) for mining
- We use data-mining models as building blocks in the mining process, rather than thinking of them as the end result
  - The central problem is to find data subsets
    ("bellwether regions") that lead to predictive features which can be gathered at low cost for a new case

# Motivating Example

- A company wants to predict the first year's worldwide profit for a new item, by using its historical database
- Database Schema:



• The combination of the underlined attributes forms a key

# A Straightforward Approach

• Build a regression model to predict item profit



• There is much room for accuracy improvement!

# **Using Regional Features**

- Example region: [1<sup>st</sup> week, HK]
- Regional features:
  - Regional Profit: The 1<sup>st</sup> week profit in HK
  - Regional Ad Expense: The 1<sup>st</sup> week ad expense in HK
- A possibly more accurate model:

 $\begin{aligned} \textit{Profit}_{[1\text{yr, All}]} &= \beta_0 + \beta_1 \textit{ Laptop } + \beta_2 \textit{ Desktop } + \beta_3 \textit{ RdExpense } + \\ \beta_4 \textit{ Profit}_{[1\text{wk, KR}]} + \beta_5 \textit{ AdExpense}_{[1\text{wk, KR}]} \end{aligned}$ 

- Problem: Which region should we use?
  - The smallest region that improves the accuracy the most
  - We give each candidate region a cost
  - The most "cost-effective" region is the bellwether region

## **Basic Bellwether Problem**

- Historical database: DB
- Training item set: I
- Candidate region set: R
  - E.g., { [1-n week, Location] }

#### Location domain hierarchy



- Target generation query: τ<sub>i</sub>(DB) returns the target value of item
  i ∈ I
  - E.g.,  $\alpha_{sum(Profit)} \sigma_{i, [1-52, AII]}$  ProfitTable
- Feature generation query:  $\phi_{i,r}(DB)$ ,  $i \in I_r$  and  $r \in R$ 
  - I<sub>r</sub>: The set of items in region r
  - E.g., [ Category<sub>i</sub>, RdExpense<sub>i</sub>, Profit<sub>i, [1-n, Loc]</sub>, AdExpense<sub>i, [1-n, Loc]</sub>]
- Cost query:  $\kappa_r(DB)$ ,  $r \in R$ , the cost of collecting data from r
- Predictive model:  $h_r(\mathbf{x}), r \in \mathbf{R}$ , trained on  $\{(\phi_{i,r}(\mathbf{DB}), \tau_i(\mathbf{DB})) : i \in \mathbf{I}_r\}$ 
  - E.g., linear regression model

### **Basic Bellwether Problem**



For each region **r**, build a predictive model  $h_r(\mathbf{x})$ ; and then choose bellwether region:

- Coverage(**r**) = fraction of all items in region ≥ minimum coverage support
- Cost(**r**, **DB**) ≤ cost threshold
- *Error*(*h*<sub>*r*</sub>) is minimized

### Experiment on a Mail Order Dataset

#### **Error-vs-Budget Plot**



- Bel Err: The error of the bellwether region found using a given budget
- Avg Err: The average error of all the cube regions with costs under a given budget
- Smp Err: The error of a set of randomly sampled (non-cube) regions with costs under a given budget

(RMSE: Root Mean Square Error)

#### Experiment on a Mail Order Dataset

#### **Uniqueness Plot**



- **Y-axis:** Fraction of regions that are as good as the bellwether region
  - The fraction of regions that satisfy the constraints and have errors within the 99% confidence interval of the error of the bellwether region
- We have 99% confidence that that [1-8 month, MD] is a quite unusual bellwether region

### Subset-Based Bellwether Prediction

- Motivation: Different subsets of items may have different bellwether regions
  - E.g., The bellwether region for laptops may be different from the bellwether region for clothes
- Two approaches:

#### **Bellwether Tree**

#### **Bellwether Cube**



			<b>L</b>					
	-		Low	Medium	High			
٨	Software	OS	[1-3,CA]	[1-1,NY]	[1-2,CA]			
or								
eg	Hardware	Laptop	[1-4,MD]	[1-1, NY]	[1-3,WI]			
)at								

#### **R&D** Expenses

### Conclusions
## Related Work: Building models on OLAP Results

- Multi-dimensional regression [Chen, VLDB 02]
  - Goal: Detect changes of trends
  - Build linear regression models for cube cells
- Step-by-step regression in stream cubes [Liu, PAKDD 03]
- Loglinear-based quasi cubes [Barbara, J. IIS 01]
  - Use loglinear model to approximately compress dense regions of a data cube
- NetCube [Margaritis, VLDB 01]
  - Build Bayes Net on the entire dataset of approximate answer count queries

## Related Work (Contd.)

- Cubegrades [Imielinski, J. DMKD 02]
  - Extend cubes with ideas from association rules
  - How does the measure change when we rollup or drill down?
- Constrained gradients [Dong, VLDB 01]
  - Find pairs of similar cell characteristics associated with big changes in measure
- User-cognizant multidimensional analysis [Sarawagi, VLDBJ 01]
  - Help users find the most informative unvisited regions in a data cube using max entropy principle
- Multi-Structural DBs [Fagin et al., PODS 05, VLDB 05]

## **Take-Home Messages**

- Promising exploratory data analysis paradigm:
  - Can use models to identify interesting subsets
  - Concentrate only on subsets in cube space
    - Those are meaningful subsets, tractable
  - Precompute results and provide the users with an interactive tool
- A simple way to plug "something" into cube-style analysis:
  - Try to describe/approximate "something" by a distributive or algebraic function

## **Big Picture**

- Why stop with decision behavior? Can apply to other kinds of analyses too
- Why stop at browsing? Can mine prediction cubes in their own right
- Exploratory analysis of mining space:
  - Dimension attributes can be parameters related to algorithm, data conditioning, etc.
  - Tractable evaluation is a challenge:
    - Large number of "dimensions", real-valued dimension attributes, difficulties in compositional evaluation
    - Active learning for experiment design, extending compositional methods