

An Iterative Refinement Approach for Mining in Large Sized Data

Dursun Delen, Merylin Kletke, Jinhwa Kim
School of Business
Oklahoma State University

Introduction

- Motivation
- Data mining in Massive Data Set
- Iterative Refinement Method for Massive Data
- Tests and Analysis
- Future Studies

2 Jinhwa Kim

Motivation

Q1. What is a new opportunity in data mining?

- Information overload
- Massive data
- Stream of data

Q2. How to handle this massive data efficiently?

3 Jinhwa Kim

2. Applications of Data Mining

- Stock market prediction
- Bankruptcy prediction
- Feature detection
- Dynamic flexible manufacturing
- Criminal management
- Weather forecasting

4 Jinhwa Kim

Data mining in Massive Data Set

1. What is Massive Data Set?

A data set that is so large and complex that existing methodologies can't cope with it readily.

2. What are the problems with Massive Data Set?

- Size causes aggravation.
- : As the size increases, the accuracy goes down.
- What if data size is keep increasing?

3. Examples of Massive Data Sets

- Prediction with stock market data
- Predicting income with US Census Data
- Earth Observation Systems
- Semiconductor manufacturing
- Crime statistics
- Network traffic
- Health care
- Geographic Information Systems

4. Techniques in Data Mining

- Nonparametric Regression
- Classification
- Clustering
- Neural Networks
- Genetic Algorithms
- CART
- CHAID
- SEES

5. Possible Approaches to Massive Data Sets

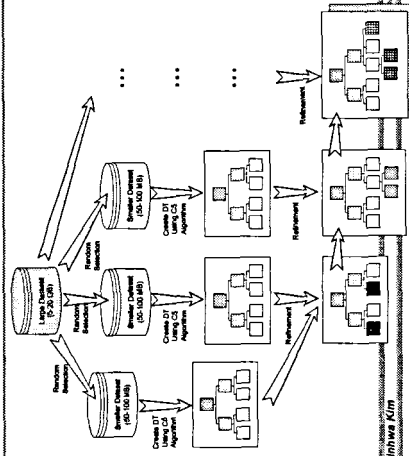
- Sampling Techniques
- Parallel Neural Networks (Takahash and Hiramatsu 1990, Chiel etc.1992, Owens 2000)
- Merging Knowledge Bases (Gehrke etc. 1998, Hall etc. 1998,Utgoff etc. 1998)
- Distributed Systems (Frieder and Kjell 1991, Paul etc., 1994, Shiroshita etc.1996, George and Knobbe 1998)

Iterative Refinement Method for Massive Data

- Conceptually it iteratively explores huge data ,collect information and add this information to grow a decision tree.
- Physically it cumulates knowledge from each iteration into a knowledge base.

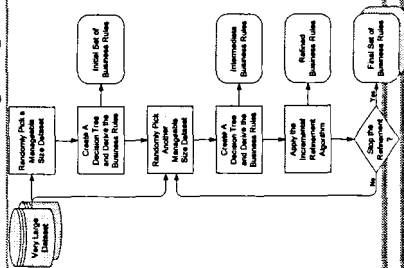
10
Jinhwa Kim

The Repetitive Methodology of Merging New Rules into the Domain Knowledge Base



10
Jinhwa Kim

A Visual Display of IRM



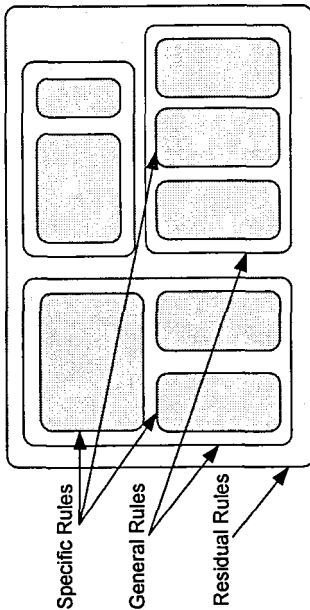
11
Jinhwa Kim

Two Steps in IRM

- Step 1: Build the domain knowledge base with a repetitive specialization process.
- Step 2: Refine the domain knowledge base with a repetitive refinement process.

12
Jinhwa Kim

Increasing Coverage Via Specializing Rules



13
Jinhwa Kim

Tests & Analysis

- Data: Census Data 1990
- Independent variables (15)
age, education, marital status, capital gain, capital loss .. etc.,
- Dependent variable
income > = \$50,000
income < \$50,000

14
Jinhwa Kim

Test results with mid-size (32MB) dataset

- Each training set = 6032 records
- Test set = 15K+ records

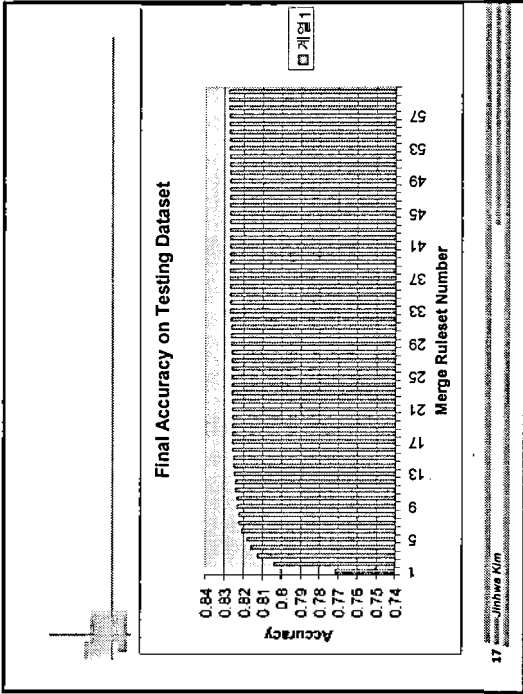
Samples	CART (160 Cross Val)			Neural Nets (MLP)			CHAD			Discriminant Analysis		
	c50K	>50K	Overall	c50K	>50K	Overall	c50K	>50K	Overall	c50K	>50K	Overall
Set #1	85.45%	80.03%	82.22%	77.02%	76.97%	77.01%	86.24%	7.35%	75.91%	93.01%	49.19%	82.24%
Set #2	90.33%	55.09%	81.65%	77.90%	77.89%	77.89%	89.05%	5.33%	76.04%	84.45%	42.65%	81.72%
Set #3	89.51%	53.22%	81.06%	79.27%	78.01%	79.05%	88.41%	48.19%	78.53%	93.40%	50.22%	82.79%
Set #4	89.31%	56.81%	81.33%	79.51%	78.41%	79.46%	82.77%	24.95%	76.14%	83.81%	40.19%	80.64%
Set #5	86.72%	54.51%	81.10%	80.02%	79.09%	79.77%	89.45%	55.55%	81.07%	93.57%	40.84%	82.62%
Mean	89.67%	56.32%	81.48%	78.74%	78.37%	78.65%	93.58%	28.24%	77.34%	93.05%	44.82%	81.69%
St. Dev.	0.40%	2.29%	0.46%	1.24%	0.87%	1.16%	4.91%	22.94%	2.35%	0.35%	4.74%	0.97%
Median	89.51%	55.22%	81.33%	79.27%	78.01%	79.05%	82.77%	24.95%	76.14%	93.40%	42.65%	81.72%
Min	89.31%	54.51%	81.06%	77.02%	76.97%	77.01%	86.41%	4.32%	75.91%	93.01%	49.19%	80.65%
Max	90.33%	60.03%	82.22%	80.02%	79.41%	79.77%	95.03%	55.55%	81.07%	94.45%	50.22%	82.79%

15
Jinhwa Kim

A Test Example - Iterations in IRM

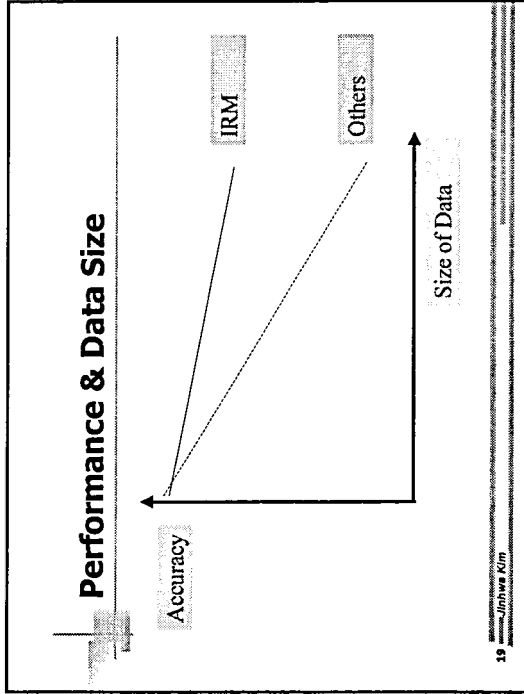
Iterations	Accuracy	Accuracy Improved	Number of Rules	No. New Rules Added
5	82.76	0	167	167
10	84.62	2.36	286	119
15	85.05	0.43	409	123
20	85.16	0.11	511	102
25	85.16	0	637	126
30	85.16	0	637	0
35	85.16	0	637	0
40	85.16	0	637	0
45	85.18	0.02	760	123
50	85.18	0	760	0
55	85.18	0	760	0
60	85.18	0	760	0

14
Jinhwa Kim



Performance Tests

C&RT	Neural Nets	CHAID	Discriminant Analysis	SEES	SEES Separated & Merged	IRM
81.48	76.65	77.54	81.60	81.82	83.18	85.18



- ### Future Studies
- Huge Data in Congress Library
 - Real Time Data Mining
 - On-line stock market data
 - Government Applications
 - Transportation
 - Criminal Management