The Potential of Associative **Classifiers**



Presentation Outline

- Typical Machine Learning
- What is Classification?
- What are the Challenges?
- The Associative Classifier
- Breast Cancer Detection
- Other Examples
- Dealing with efficiency & effectiveness

The Potential of Associative Classifiers

Osmar R. Zaïar

Text & Multimedia Mining



What is Classification?

The goal of data classification is to organize and categorize data in distinct classes.

- A model is first created based on the data distribution.
- ▶ The model is then used to classify new data.
- Given the model, a class can be predicted for new data.



rishane. December 2005



Challenges

- Dealing with high dimensional spaces
- Handling missing data
- Deriving a model that can be interpreted (Transparency leads to trust for some applications)
- Deriving a model that can be edited by human experts (to inject domain knowledge)
- Dealing with very large or evolving training sets
- Allowing multi-label classification
- Dealing with uneven representation of classes in training sets



Basic Concepts

A transaction is a set of items: $T = \{i_a, i_b, \dots i_t\}$

 $T \subset I$, where *I* is the set of all possible items $\{i_1, i_2, ..., i_n\}$

D, the task relevant data, is a set of transactions.

An association rule is of the form: $P \rightarrow Q$, where $P \subset I$, $Q \subset I$, and $P \cap Q = \emptyset$



What Is Association Rule Mining?

 Association rule mining searches for relationships between items in a dataset:



The Potential of Associative Classifiers

Osmar R. Zaïan







- We want to find associations between extracted features and class labels

Osmar R. Zaïane

 Constrain the association rule mining such that the rules found are of the following form:



 We used a constrained version of frequent itemset mining.

Prichana December 2005

- The class label has to be part of any {A, B, C, Class} frequent itemset

The Potential of Associative Classifiers

- The class label is a consequent, and all other items are the antecedent of a rule A, B, C -> Class

Modeling documents



How do Associative Classifiers Work?





Brisbane, December 2005





Improving the Quality of Images

- Digitization introduces noise
- Inconsistent illumination conditions
- Inconsistent sizes and distributions

Automatic Cropping: Removes unwanted parts and artifacts. Enhancement: Diminishes the effect of over brightness and

over darkness. Histogram equalization to increase contrast range.







Histogram

Classifiers

Equalization



Original mammogram

ram

Cropping

Partitioning & orientation

Osmar R. Zaïane

۲

Brisbane, December 2005

Experimental Results







Precision over 10 splits



Brisbane, December 2005

ifiers

Recall over 10 splits

Osmar R. Zaïane

Experimental Results with Reuters

Reuters collection: ModApte version: 12,202 documents consisting of 9.603 training documents and 3,299 testing documents.

Osmar R. Zaïane

Osmar R. Zaï

BEP	ARC-BC with $\delta = 50$			Bayes	Rocchio C4.5		k-NN	bigrams	SVM	SVM
	10%	15%	20%				1970) (1990) - AM (1997) N	0	(poly)	(rbf)
acq	90.9	89.9	87.8	91.5	92.1	85.3	92.0	73.2	94.5	95.2
corn	69.6	82.3	70.9	47.3	62.2	87.7	77.9	60.1	85.4	85.2
crude	77.9	77.0	80.7	81.0	81.5	75.5	85.7	79.6	87.7	88.7
earn	92.8	89.2	86.6	95.9	96.1	96.1	97.3	83.7	98.3	98.4
grain	68.8	72.1	73.1	72.5	79.5	89.1	82.2	78.2	91.6	91.8
interest	70.5	70.1	75.3	58.0	72.5	49.1	74.0	69.6	70.0	75.4
money- fx	70.5	72.4	70.5	62.9	67.6	69.4	78.2	64.2	73.1	75.4
ship	73.6	73.2	63.0	78.7	83.1	80.9	79.2	69.2	85.1	86.6
trade	68.0	69.7	69.8	50.0	77.4	59.2	77.4	51.9	75.1	77.3
wheat	84.8	86.5	85.3	60.6	79.4	85.5	76.6	69.9	84.5	85.7
micro- avg	82.1	81.8	81.1	72.0	79.9	79.4	82.3	73.3	85.4	86.3
macro- avg	76.74	78.24	76.32	65.21	79.14	77.78	82.05	67.07	84.58	86.01

Precision/Recall-breakeven point on ten most populated Reuters categories for ARC-BC and most known classifiers

Localization of Proteins

plast

mito-





A eukaryotic cell: many compartments



KDD Cup 2002



- Motivation
- What is Classification?
- What are the Challenges?
- The Associative Classifier
- Breast Cancer Detection
- Other Examples
- **Dealing with efficiency & effectiveness**

Osmar R. Zaï

Brishane, December 2004

Brisbane, December 200

Brisbane, December 2005

Improving Efficiency (Rule Generation)



Positive & Negative Rule Generation

- It generates all positive and negative association rules with strong correlation
- minsupp, minconf user-defined
- correlation starts at ρ_{min} = 0.5

Error Rate Comparison

+R

5.5

23.3

16.3

28.7

27.4

6.6

C4.5 CBA

4.2

25.3

18.5

7.1

27.8

27.6

3.9

27.6

18.9

5.5

26.5

27.5

- The process of rules generation is apriori-like $-C_k = F_{k-1} \times F_1$
- For each pair X,Y, where X∪Y is itemset in C_k – correlation(X,Y) is computed

Some ARC-PAN Results

ARC-PAN

4.8

25.4

17.0

6.6

28.7

27.1

+R&-R

3.8

25.1

16.2

6.0

28.9

26.9

rules

17.000 5.0

200,000 24.7

4.000 21.8

140 7.3

4.000 34.3

4.000 22.0

the form: $\neg F \rightarrow C$

Only positive rules

+R&-AR

Osmar R. Zaïane

rules error

1,000 5.5

60 6.6

Osmar R. Zai

40 23.3

80 16.3

500 28.7

50 27.4

Support vs. Correlation

Strong Rules Correlated Rules

error

All positive and negative rules

ositive rules and rules of

Positive & Negative Rule Generation

 $A \rightarrow \neg B$

• if the correlation is **positive**:



• if the correlation is negative:

 $\neg A \rightarrow B$

 if the rules have high confidence they are added to the discovered set of rules

Repetition of Features

ICDE'00 and PAKDD'05

Osmar R. Zaïan

Model transactions of features not as binary

{Item_a, ...Item_k } But as enumerations of repeated features

 $\{\boldsymbol{\alpha} \text{Item}_{a}, \dots \boldsymbol{\beta} \text{Item}_{k}\}$

- Use reoccurring frequent itemset mining (ICDE 2000) to generate rules such as: $\{\alpha \text{ Item}_a, \dots \beta \text{ Item}_k\} \rightarrow \text{Class}_x$

Nothing to do with quantitative association rules

Brisbane, December 2005

Prichana December 200

Datasets

Diabetes

Breast

Heart

Led7

Pima

Iris

Brisbane, December 2005

Osmar R. Zaïane 📢



