Outline

Associative Classifiers Evaluation Measures Classification with Rare Classes

Luiza Antonie PhD Candidate

• What are associative classifiers?

- How do you use them?
- How can you evaluate and compare multiple classification systems?
- What measures are better?
- How do classifiers perform when there are rare classes?

Tuesday, November 09, 2004

Luiza Antonie - Cmput 695 Presentation

Association Rules (typical)

- Association rule mining aims at discovering associations between items in a transactional database.
- Given D={T1...Tn} a set of transactions and I={i1...in} a set of items such that any Ti in D is a set of items in I.
- An association rule is an implication A→B where A and B are subsets of Ti given some support and confidence thresholds.
- The **support** of the rule is the probability that A and B hold together among all the possible presented cases.
- The **confidence** of the rule is the conditional probability that the consequent B is true under the condition of the antecedent A.

Association Rules for Classification



Associative Classifier Rule Discovery (1st stage) • Given a training set a set of association rules is • Rule discovery discovered using an ARM algorithm: - Using an ARM algorithm – Apriori, FP-tree, etc. • Pruning - Modify the algorithms to mine the form of rules that you - Discarding those rules that are redundant or not want: interesting - Mine all the association rules and filter them afterwards: • Classification • CBA, ARC-AC and ARC-BC - Apriori based; - Based on a scoring scheme, use the set of rules to • CMAR classify new, unseen instances - FP-tree: Luiza Antonie - Cmput 695 Presentation Luiza Antonie - Cmput 695 Presentation Tuesday, November 09, 2004 Tuesday, November 09, 2004 Pruning (2nd stage) **Pruning** (2nd stage) Noisy information • Database Coverage Large number of rules - select a small set of high quality rules Long classification time - A set of rules (SR) • Removing low ranked specialized rules; • The rules are ordered by confidence and support $R_1: F_1 \Longrightarrow C$ Confidence 90% - for each rule (R) $\Rightarrow R_1$ $R_2: F_1 \wedge F_2 \Longrightarrow C$ Confidence 80% • if R classifies correctly at least one example, keep R • remove the examples covered by R • Eliminate conflicting rules; - stop when there are no more examples or all the $F_1 \Longrightarrow C_1 \land F_1 \Longrightarrow C_2$ rules have been checked • Database coverage;

Tuesday, November 09, 2004

Classification (3rd stage)

- A set of rules (SR)
 - The rules are ordered by confidence and support
- A new instance to be classified
- From SR a subset of rules SR' matches the new instance
 - Divide SR' in subsets based on the class label
 - SR'C₁, SR'C₂, SR'C_n

Classification (3rd stage)

- Different strategies
 - CBA
 - Choose the first matching rule (highest confidence)
 - CMAR
 - For each SR'C₁, SR'C₂, ..., SR'C_n set
 - Computes a weighted chi-square
 - Chooses the class with the best score (best chi-square)
 - ARC-AC and ARC-BC
 - For each SR'C₁, SR'C₂, ..., SR'C_n set
 - Computes the average of the confidences
 - Chooses the class with the best score (best average confidence)

Tuesday,	November	09, 200
----------	----------	---------

Luiza Antonie - Cmput 695 Presentation

Classification Stage for ARC

Luiza Antonie - Cmput 695 Presentation

Let S be the classification system

A new object O <f1; f3; f4; f7; f9 >

 $f1 \Rightarrow C1$ confidence 0.9 f3 & f4 => C2 confidence 0.85 $f4 \Rightarrow C2$ confidence 0.8 $f7 \Rightarrow C1$ confidence 0.6 $f9 \Rightarrow C3$ confidence 0.5



Association Rules Classification with All Categories



Tuesday, November 09, 2004

Tuesday, November 09, 2004

Using the dominance factor we chose the

 $\delta = 80\%$ O is predicted to fall in C2 and C1.

winning categories. If $\delta = 100\%$ C2 is winning. If

11

ARC-AC



ARC-BC



Association Rules Classification by Category



Evaluation

- Why do we do it?
 - To study the performance of the classification systems and to compare with other algorithms
- How do we do it?
 - Accuracy/error
 - Precision, Recall, F-measure
 - Graphical methods
 - ROC
 - Lift curves, PN curves, cost curves

2-class Confusion Matrix

	Predicted class			
True class	positive	negative		
positive (P)	ТР	FN=P - TP		
negative (N)	FP	TN=N - FP		

$$accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

$$error = \frac{FP + FN}{TP + FP + TN + FN}$$

Tuesday, November 09, 2004	Luiza Antonie – Cmput 695 Presentation

2-class Confusion Matrix

	Predicted class			
True class	positive	negative		
positive (P)	ТР	FN=P - TP		
negative (N)	FP	TN=N - FP		

precision =
$$\frac{TP}{TP + FN}$$
 $F - measure = \frac{(1 + \beta^2) * P * R}{P + \beta^2 * R}$

$$F_1 - measure = \frac{2*P*R}{P+R}$$

 $recall = \frac{TP}{TP + FP}$

Luiza Antonie - Cmput 695 Presentation

18

R

ROC curve

	Predicted class			
True class	positive	negative		
positive (P)	ТР	P - TP		
negative (N)	FP	N - FP		

- Reduce the 4 numbers to two rates true positive rate = TPR = TP/Pfalse positive rate = FPR = FP/N
- Rates are independent of class ratio

Example: 3 classifiers

	Pred	icted		Pred	icted			Pred	icted
True	pos	neg	True	pos	neg	1	True	pos	neg
pos	40	60	pos	70	30		pos	60	40
neg	70	30	neg	50	50		neg	80	20
Cla	assifie	<u>r 1</u>	Cla	assifie	<u>r 2</u>		Cl	assifie	ar 3

19

Source: Rob Holte



Dominance



Linear Interpolation





Operating Range

Combining Classifiers – Convex Hull



Rare Classes

- Class imbalance occurs when some classes have many examples, while others are represented by just a few;
- Small classes are difficult to classify for existing classification algorithms;
- Applications medical data; text data; biological data; detection of intrusions;

Why is it difficult? Class imbalance Class overlapping Noisy data Small disjuncts

Methods and Solutions

Luiza Antonie - Cmput 695 Presentation

• Data manipulation

Tuesday, November 09, 2004

- balancing data by sampling
- data segmentations
- learning only the minority class
- Cost-sensitive classifiers
 - MetaCost [Domingos '99]
 - AdaCost [Fan et al. '99]

27

Methods and Solutions

- Creating new algorithms
 - PN-rule [Joshi et al. '01]
 - P-rules (rules that predict the presence of the target class)
 - N-rules (rules that imply the absence of the target class)
 - two phases
 - 1st focuses on recall; completeness retrieving the results
 - 2nd improves precision; quality retaining only the desired examples

Methods and Solutions

- Creating new algorithms
 - Boosting
 - AdaBoost [Schapire '99]
 - starts with a weak classifier and boosts is performance
 - iterative processes
 - At each iteration the weights that are attached to a training example are refined
 - misclassified weight increase
 - correctly classified weight decrease

- ARC-BC

Tuesday, November 09, 2004	Luiza Antonie – Cmput 695 Presentation	29	Tuesday, November 09, 2004	Luiza Antonie – Cmput 695 Presentation	30

Methods and Solutions

• Associative Classifier (by category)







Sampling

• Under-sampling

Tuesday, November 09, 2004

- random eliminate at random examples of the majority class
- Tomek links [Tomek '76]
- Condensed Nearest Neighbour Rule (CNN) [Hart '68]
- One-sided Selection (OSS) [Kubat and Matwin '97]
- Neighbourhood Cleaning Rule (NCL) [Laurikkala '01]
- drawbacks it can discard good examples

Sampling

- Over-sampling
 - random duplicate at random some examples belonging to the minority class;
 - Smote [Chawla '02]– create artificial examples for the minority class by interpolating between existing examples;
 - drawbacks

Tuesday, November 09, 2004

- overfitting
- execution time increase
- Under-sampling + Over-sampling

Under-Sampling vs. Over-Sampling

- Under-sampling
 - [Drummond and Holte '03]
 - [Domingos '99]
- Over-sampling

Tuesday, November 09, 2004

- [Japkowicz et al. '02] artificial datasets
- [Batista et al. '04]
- Is it dataset dependant?

Discussion

Luiza Antonie - Cmput 695 Presentation

- classification errors occur near class boundaries
- difficult to find good boundaries when classes are overlapping
- take into account the overlapping of classes
- application dependant

Discussion

Luiza Antonie - Cmput 695 Presentation

- develop new algorithms to deal with rare classes
- incorporate in the existing algorithms some knowledge about the classification with rare classes
- use some graphical evaluation measures to improve the classifier

35

		Thank You!	
--	--	------------	--