

Data Mining Assignment 2

Cost sensitive classifiers

08305006: Prashanth K
08305023: Prashant Borole
08305028: Sriram Kashyap
08305045: Anup Kulkarni

October 10, 2008

1 Introduction

1.1 Basic Info

- Classifiers Implemented:
 - Cost sensitive Decision Tree
 - Cost sensitive Naive Bayes Classifier
- Files Created/Modified:
 - weka/classifiers/trees/CostTree.java
 - weka/classifiers/bayes/CostSensitiveNaiveBayes.java
 - weka/gui/GenericObjectEditor.props
 - costs.cfg

1.2 Instructions

Instructions to run:

- Run weka.jar from the command-line: `java -jar weka.jar`
- Edit the costs.cfg file so that the first line is the file path of the attribute cost file and the second line is the file path of the cost matrix file (The remaining lines are ignored).

- Load the dataset, and run one of these classifiers:
 - `trees.CostTree`
 - `bayes.CostSensitiveNaiveBayes`

Note: The classifiers work for data sets with numeric/nominal/missing attributes, and nominal classes (numeric classes are currently not supported).

2 Cost Sensitive Naive Bayes Classifier

2.1 Formulation

Given (D, C, T) , where:

- D is a training dataset consisting of N samples (x_1, x_2, \dots, x_N) from P classes (c_1, c_2, \dots, c_P) . Each sample x_i is described by M attributes (A_1, x_2, \dots, A_M) among whom there can be missing values.
- C is a misclassification cost matrix. Each entry $C_{ij} = C(i, j)$ specifies the cost of classifying a sample from class c_i as belonging to class c_j ($1 \leq i, j \leq P$). Usually, $C_{ii} = 0$.
- T is a test-cost vector. Each entry $T_k = T(k)$ specifies the cost of taking a test on attribute A_k ($1 \leq k \leq M$);

Build a test-cost sensitive naive Bayes classifier $csNB$ and for every test case, a test strategy with the aim to minimize the sum of the misclassification cost C_{mc} and test cost C_{test} .

2.1.1 Strategy

A sequential test strategy is as follows. During the process of classification, based on the results of previous tests, decisions are made sequentially on whether a further test on an unknown attribute should be performed, and if so, which attribute to select.

Suppose that $x = (a_1, a_2, \dots, a_M)$ is a test example. Each attribute a_i can be either known or unknown. Let \tilde{A} denote the set of known attributes among all the attributes A and \bar{A} the unknown attributes. The expected misclassification cost of classifying x as class c_j based on \tilde{A} is:

$$R(c_j|x) = R(c_j|\tilde{A}) = \sum C_{ij} \times P(c_i|\tilde{A}), 1 \leq j \leq P \quad (1)$$

c_{j*} with the minimum expected cost is predicted as the class label.

To decide which attribute $\tilde{A} \in \bar{A}$ to select, we define the *utility* of testing an unknown attribute \bar{A}_i as follows:

$$Util(\bar{A}_i) = Gain(\tilde{A}, \bar{A}_i) - T_i \quad (2)$$

$$Gain(\tilde{A}, \bar{A}_i) = C_{mc}(\tilde{A}) - C_{mc}(\tilde{A} \cup \bar{A}_i) \quad (3)$$

$$C_{mc}(\tilde{A}) = \min_j R(c_j | \tilde{A}) \quad (4)$$

$$C_{mc}(\tilde{A} \cup \bar{A}_i) = \sum_{k=1}^{|\bar{A}_i|} P(\bar{A}_i = x_k | \tilde{A}) \times \min_j R(c_j | \tilde{A}, \bar{A}_i = x_k) \quad (5)$$

Overall, an attribute \bar{A}_i is worth testing if testing it offers more gain than the cost it brings. We select the attribute A_i with the maximum utility, read its value, update \tilde{A} and \bar{A} with A_i , and continue this process until none of the remaining attributes are worth testing.

3 CostTree

The cost sensitive decision tree is based on a modification of the algorithm suggested by Ling et al. in “Decision Trees with Minimal Costs”. It has been built using weka’s REP-Tree as a base, and modifying the gain functions that it uses.

3.1 Formulation

The decision tree is made cost-sensitive by selecting those attributes that have highest gain, at each stage of the tree building process. The gain is defined as:

$$Gain = priorCost - cCost - attribCost \times N \quad (6)$$

Here, *priorCost* is the cost of misclassification before the split, and *cCost* is the cost of misclassification after the split. *attribCost* is the cost of evaluating the attribute over which the split is taking place. This cost is multiplied by the number of instances (N) for which this attribute should be evaluated.

$$currentCost = \sum_{i=0}^n \sum_{j=0}^d (N * dist_j) * C_{jk} \quad (7)$$

where: n is the number of values that the attribute can take,
 N is the number of instances,
 d is the number of attributes,
 $dist_j$ is the probability of class value j ,
 C_{jk} is the cost of misclassifying an instance of class j as that of class k , where k is the dominating class of the split.

Given a distribution for c classes, the dominating class i for that node is calculated as follows:

$$\arg \min_i \text{cost} = \sum_{j=0}^c dist_j * C_{ji} \quad (8)$$

It has been shown that if a tree is built using these rules, the best attributes to evaluate are at the top of the tree, and the most expensive attributes move to the bottom of the tree. This means that evaluating the class of a new instance using this tree is just the matter of traversing the tree, till the leaf node is reached. Any attributes which are un-economical to expand are not included in the tree building process.

The class of an instance is calculated from the distributions by weighing each class distribution by the misclassification costs of that class, as shown above.

4 Results

Pima Indians

cs dataset 1

	a	b
a	0	10
b	70	0

CS Naive Bayes

	a	b	
a	151	349	0.68
b	9	259	

SimpleNaiveBayes

	a	b	
a	422	78	0.42
b	104	164	

cs dataset 2

	a	b
a	0	1
b	1	0

	a	b	
a	500	0	0.59
b	268	0	

Confusion Matrix
RMS Error
cost matrix

cs dataset 3

	a	b
a	0	70
b	10	0

	a	b	
a	486	14	0.59
b	249	19	

cs dataset 4

	a	b
a	0	0
b	70	0

	a	b	
a	0	500	0.81
b	0	268	

cs dataset

	a	b
a	0	70
b	70	0

	a	b	
a	426	74	0.51
b	126	142	

Cancer dataset

cs dataset 1

	a	b
a	0	100
b	400	0

CS Naive Bayes

	a	b	
a	330	14	0.19
b	5	174	

SimpleNaiveBayes

	a	b	
a	324	20	0.21
b	4	175	

cs dataset 2

	a	b
a	0	1
b	1	0

	a	b	
a	344	0	0.59
b	179	0	

Confusion Matrix
RMS Error
cost matrix

cs dataset 3

	a	b
a	0	400
b	100	0

	a	b	
a	332	12	0.26
b	22	157	

cs dataset 4

	a	b
a	0	0
b	300	0

	a	b	
a	0	344	0.81
b	0	179	

cs dataset

	a	b
a	0	400
b	400	0

	a	b	
a	331	13	0.2
b	7	172	

Pima Indians

CostTree

REPTree

cs dataset 1

	a	b
a	0	10
b	70	0

	a	b
a	0	500
b	0	268

0.65

	a	b
a	422	78
b	113	155

0.43

cs dataset 2

	a	b
a	0	1
b	1	0

	a	b
a	500	0
b	268	0

0.48

Confusion Matrix
RMS Error
cost matrix

cs dataset 3

	a	b
a	0	70
b	10	0

	a	b
a	500	0
b	268	0

0.56

cs dataset 4

	a	b
a	0	0
b	70	0

	a	b
a	0	500
b	0	268

0.81

cs dataset

	a	b
a	0	70
b	70	0

	a	b
a	441	59
b	197	71

0.47

Cancer dataset

Cost Tree

REP Tree

cs dataset 1

	a	b
a	0	100
b	400	0

	a	b
a	315	29
b	4	175

0.28

	a	b
a	332	12
b	19	160

0.22

cs dataset 2

	a	b
a	0	1
b	1	0

	a	b
a	344	0
b	179	0

0.47

Confusion Matrix
RMS Error
cost matrix

cs dataset 3

	a	b
a	0	400
b	100	0

	a	b
a	338	6
b	57	122

0.34

cs dataset 4

	a	b
a	0	0
b	300	0

	a	b
a	0	344
b	0	179

0.81

cs dataset

	a	b
a	0	400
b	400	0

	a	b
a	327	17
b	19	160

0.25