

Intelligent Data Mining Techniques For Coal Mining Data

Ravi K Jade^{#1} , L. K. Verma^{*2} , Kesari Verma^{#3}

^{#1}Department of Mining Engineering
National Institute of Technology Raipur

^{*2}Department of Computer Applications
Raipur Institute of Technology Raipur

^{#3} Department of Computer Applications
National Institute of Technology Raipur

Abstract

Classification is an important problem in data mining. Given a database of records, each with a class label, a classifier generates a concise and meaningful description for each class that can be used to classify future records whose classes are unknown. A number of popular classifier exists like naïve bays classifier, Neural Network, SVM classifier, CARD etc. In this paper we applied classification algorithm for coal mines dataset. We also discussed decision tree, nearest neighbor classifier, NN classifier for prediction of class label.

Keywords: Data Mining algorithm, C4.5 algorithm, Naïve Bayes Algorithm, Intelligent Data mining for coal data

Category and Subject Descriptors:

H.2.8 [Information System] : Database Management - database application, data Mining.

1. Introduction

Classification is an important problem in data mining. Under the guise of supervised learning, classification has been studied extensively by the machine learning community as a possible solution to the “knowledge acquisition” or “knowledge extraction” problem. The input to the classifier construction is a training set of records, each of which is tagged with a class label. A set of attribute values defines each record. Attributes with discrete domains are referred to as categorical, while those with ordered domains are referred to as numeric. The goal is to induce a concise model or description for each class in terms of the attributes. The model is then used by the classifier to

classify (i.e., assign class labels to) future records whose classes are unknown [10]. It thus, reduces outlier problem. The classification model of dataset is shown in figure 1.

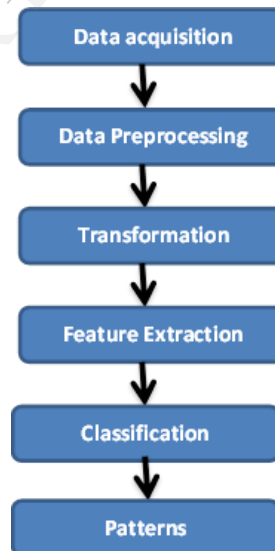


Fig . Data Mining process

Coal mine dataset is shown in Table I that comprises of Spacing(M), Burden(M), Depth (M), Stemming(M). Second part Types of Explosive that consist of Booster(Kg), Primer(Kg), Boulders (Nos.).

TABLE I.

COAL MINE DATASET

Spac	Burd.	Depth	Stemm	Boost.	Prim.	SME	Boulders (Nos.)
5	3.25	6.25	3	10.5	0	4600	1340
5.75	3.5	7	3.5	15.3	294	7500	420
5.75	3.5	7	3.5	10.5	0	5700	1200
6.25	4	8.25	3.5	35.2	0	15000	625
5	3.25	6.25	4.0	13.2	0	9000	1075

1.1 Data Acquisition process

The data is acquired from coal mines where the blasting process for mines are performed. The spacing(M) denotes the space between holes. The class label Boulder denotes after the blasting the size of boulder produced in blasting.

1.2 Data Pre-processing

Data pre-processing is an often neglected but important step in the data mining process. The phrase “Garbage in, Garbage Out” is applicable in data mining and machine learning. Data gathering methods are loosely controlled, resulting in out of range values, impossible data combination, missing values. Analyzing data that has not been carefully screen for such problems can produce misleading results. The representation and quality of data is first and foremost before running an analysis. In our dataset we filled the missing value by mean of whole data set of specified attributes. Some data values the range is given, we computed the mean value and it is replaced by mean of the data e.g. for stemming(M) the range is given 3.5 to 3.75 that data is replaced by 3.625.

1.3 Data Transformation

In data transformation, the data are transformed or consolidated into forms appropriate for mining. Data transformation process are follows.

Normalization- where the attribute data are scaled in the range such as -1.0 to +1.0.

$$v' = \frac{(v - Old_{min})}{(Old_{max} - Old_{Min})} * (New_{Max} - New_{Min}) + Min_{New}$$

1.4 Feature Selection

The main idea of feature selection is to choose a subset of input variables by eliminating features with little or no predictive information. Feature selection [2] can significantly improve the comprehensibility of the resulting classifier models and often build a model that generalizes better to unseen points. In machine learning and statistics, feature selection [2] also known as variable selection, attribute selection or variable subset selection, is the process of selecting a subset of relevant features for use in model construction. The central assumption when using a feature selection technique is that the data contains many redundant or irrelevant features. Redundant features are those which provide no more information than the currently selected features, and irrelevant features provide no useful information in any context.

3. Classification

There are various classifier to classify the data like NN, SVM, Decision tree and others.

3.1 Nearest Neighbour

Nearest neighbour (NN) [11] also known as Closest Point Search is a mechanism that is used to identify the unknown data point based on the nearest neighbour whose value is already known. It has got a wide variety of applications in various fields such as Pattern recognition, mage databases, Internet marketing, Cluster analysis etc.

The NN algorithm can also be adapted for use in estimating continuous variables. One such implementation uses an inverse distance weighted average of the k-nearest multivariate neighbors. This algorithm functions as follows:

1. Compute Euclidean or Mahalanobis distance from target plot to those that were sampled.
2. Order samples taking for account calculated distances.
3. Choose heuristically optimal k nearest neighbor based on RMSE done by cross validation technique.

3.2 Decision Tree

Decision Tree Classifier [4]A decision tree is a class discriminator that recursively partitions the training set until each partition consist entirely or dominantly of examples from one class. Each leaf node of the tree contains one or more attributes and determines how the data is partitioned.

A decision tree classifier is built in two phases [4] growing phase followed by pruning phase. In growth phase the tree is built by recursively partitioning the

data until each partition is either “pure” or sufficiently small.

The algorithm for building tree [5].

Procedure buildTree (S)

- 1) Initialize root node using dataset S
- 2) Initialize queue Q to contain root node
- 3) While Q is not empty do {
- 4) dequeue the first node N in Q
- 5) if N is not pure {
- 6) for each attribute A
- 7) Evaluate splits on attributes A
- 8) Use best split to split node N into N1 and N2
- 9) Append N1 and N2 to Q
- 10) }
- 11) }

3.3 Bayesian Network

Bayesian network (BN) is also called belief networks. A BN is a graphical representation of probability distribution. It belongs to the family of probabilistic graphical models. This BN consist of two components. First component is mainly a directed acyclic graph (DAG) in which the nodes in the graph are called the random variables and the edges between the nodes or random variables represents the probabilistic dependencies among the corresponding random variables. Second component is a set of parameters that describe the conditional probability of each variable given its parents. The conditional dependencies in the graph are estimated by statistical and computational methods [6], [7].

This prior expertise about the structure of Bayesian network algorithm work as follows.

1. Declare that a node is root node.
2. Declare that a node is leaf node.
3. Declaring that a node has direct effect of another noe.
4. Declaring that a node is not directly connected to another node.
5. Declaring that two nodes are independent, giving a condition set.
6. Providing partial ordering among the nodes.

3.4 Artificial Neural Network

Artificial Neural Network (ANN) is a computational model based on biological neural network. ANN also called Neural Network [ANN]. It contains interconnected group of artificial neurons and processes the information by a connectionist approach. ANN is an adaptive system because it changes its structure based on information flow during the learning phase.

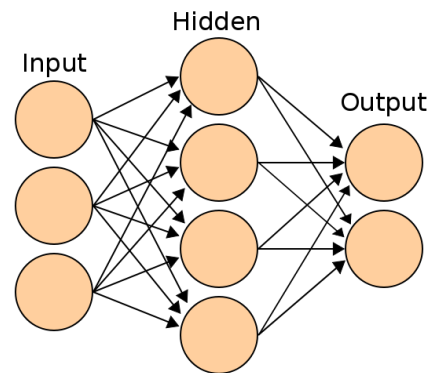


Fig. 2. Architecture of ANN

Actual algorithm for a 3-layer network (only one hidden layer) [8]:

Initialize the weights in the network (often randomly)

Do

For each example e in the training set

O = neural-net-output(network, e) ; forward pass

t = teacher output for e

- 1) Calculate output $y = x_1w_1 + x_2w_2$,
- 2) Calculate error (T - O) at the output units $E = (t - y)^2$
- 3) Compute delta_w for all weights from hidden layer to output layer ; backward pass
- 4) Compute delta_wi for all weights from input layer to hidden layer ; backward pass continued.
- 5) Update the weights in the network

The different classification techniques and their comparative charts is shown in Table 1.

TABLE I.

DIFFERENT CLASSIFICATION TECHNIQUES

Method	Generative /Discriminative	Loss Functions	Parameter estimation Algorithm
K-Nearest Neighbour	Discriminative	-log P(X,Y) or Zero one	All data are unsuperised
Decision tree	Discriminative	Zero –One loss	C4.5
Bayesian Network	Generative	-log O(X,Y)	Variable Elimination
Neural Network	Discriminative	Sum-Squared Error	Forward Propagation

4. Experimental Result

The experiments were performed in Weka machine learning software [12] in Pentium IV machine with 1 GB RAM. No other application is running while performance computation. The dataset contains 47 instances, 6 attributes i.e. Number_of_holes, Burden, Depth, Primer, SME and Class. In test mode: 10-fold cross-validation is performed. All the numeric data is converted to categorical data all <50 values denotes 'Effective' Boulders and remaining are converted as 'Not-Effective' Boulders. The data is taken from Jan- December 2009. The experimental results their accuracy shown in Table II, III, IV and V.

Table II.
Result with Decision Tree Model

Correctly Classified Instances	41	87.234 %
Incorrectly Classified Instances	6	12.766 %
Mean absolute error		0.2365
Root mean squared error		0.337

Table III
Confusion Matrix for Decision tree

	Effective	Not Effective
Effective	0	6
Not Effective	0	41

Table IV
Result with Perceptron NN Model

Correctly Classified Instances	40	85.1064
Incorrectly Classified Instances	7	14.8936
Mean absolute error		0.1568
Root mean squared error		0.345

Table V
Confusion Matrix

	Effective	Not Effective
Effective	0	4
Not Effective	3	38

5. Conclusion

In this paper we applied data mining process in coal mining data. The mining process starts with data preparation, which is an important issue for data mining, as real world data tends to be incomplete, noisy and inconsistent. Data preparation includes data cleaning, data integration, data transformation and data reduction. In order to mine the effective size of Boulders (Nos.) the attribute effective are analysed. Decision tree model gives 87% of accuracy to correctly predict the class label whereas Neural Network model predicts 84% correct class label. In future we try to implement support vector machine classifier in order to improve the accuracy of classifier.

5. References

- [1] YongSeog Kim, W. Nick Street, and Filippo Menczer. Feature Selection in Data Mining.
- [2]. http://en.wikipedia.org/wiki/Feature_selection
- [3]. Ms. Aparna Raj, Mrs. Bincy G., Mrs. T.Mathu. Survey on Common Data Mining Classification Techniques. International Journal of Wisdom Based Computing, Vol. 2(1), April 2012
- [4] J. Ross Quinlan. C4.5: Programs for Machine Learning. Morgan Kaufman, 1993. [4] J. Ross Quinlan. C4.5: Programs for Machine Learning. Morgan Kaufman, 1993.
- [5] Keshri Verma, O.P. Vyas. An Approach for Decision tree classifier for Temporal data", Science Journal of Pt. Ravishankar Shukla University Raipur.,2004 ,ISSN –NO 0170 –5910 Vol. 18, No B (Science) 2005, pp 07-22.
- [6] Charniak, E. 1991, .Bayesian Networks without tears. AI Magazine, Winter 1991.
- [7]. Ben-Gal I., Bayesian Networks, in Ruggeri F., Faltin F. & Kenett R Encyclopedia of Statistics in Quality & Reliability, Wiley & Sons (2007).
- [8] Artificial Neural Networks Ajith Abraham Oklahoma State University, Stillwater, OK, USA.
- [9] <http://en.wikipedia.org/wiki/Backpropagation>
- [10] Rajeev and Kyuseok Shim. PUBLIC: A Decision Tree Classifier that Integrates Building and Pruning". In Proceedings of the 24th International Conference on Very Large Data Bases, pages 404{415, New York, USA, August 1998.
- [11] A. Djouadi and E. Bouktached, "A Fast Algorithm for the Nearest Neighbor Classifier," IEEE Transaction Pattern Analysis and Machine Intelligence, Vol. 19 no. 3 pp. 277-282,1997.
- [12] <http://www.cs.waikato.ac.nz/ml/weka/>