

Classification and Feature Selection Techniques in Data Mining

Sunita Beniwal*, Jitender Arora

Department of Information Technology, Maharishi Markandeshwar University, Mullana,
Ambala-133203, India

Abstract

Data mining is a form of knowledge discovery essential for solving problems in a specific domain. Classification is a technique used for discovering classes of unknown data. Various methods for classification exist like bayesian, decision trees, rule based, neural networks etc. Before applying any mining technique, irrelevant attributes need to be filtered. Filtering is done using different feature selection techniques like wrapper, filter, embedded technique. This paper is an introductory paper on different techniques used for classification and feature selection.

Keywords: KDD, Preprocessing, Neural Networks, Decision trees

1. Introduction

As the world grows in complexity, overwhelming us with the data it generates, data mining becomes the only hope for elucidating the patterns that underlie it [1]. The manual process of data analysis becomes tedious as size of data grows and the number of dimensions increases, so the process of data analysis needs to be computerised.

The term Knowledge Discovery from data (KDD) refers to the automated process of knowledge discovery from databases. The process of KDD is comprised of many steps namely data cleaning, data integration, data selection, data transformation, data mining, pattern evaluation and knowledge representation.

Data mining is a step in the whole process of knowledge discovery which can be explained as a process of extracting or mining knowledge from large amounts of data [2]. Data mining is a form of knowledge discovery essential for solving problems in a specific domain. Data mining can also be explained as the non trivial process that automatically collects the useful hidden information from the data and is taken on as forms of rule, concept, pattern and so on [3]. The knowledge extracted from data mining, allows the user to find interesting patterns and regularities

deeply buried in the data to help in the process of decision making.

The data mining tasks can be broadly classified in two categories: descriptive and predictive. Descriptive mining tasks characterize the general properties of the data in the database. Predictive mining tasks perform inference on the current data in order to make predictions. According to different goals, the mining task can be mainly divided into four types: class/concept description, association analysis, classification or prediction and clustering analysis [4].

This paper provides a survey of various feature selection techniques and classification techniques used for mining.

2. Data Preprocessing

Data available for mining is raw data. Data may be in different formats as it comes from different sources, it may consist of noisy data, irrelevant attributes, missing data etc. Data needs to be pre processed before applying any kind of data mining algorithm which is done using following steps [5]:

Data Integration – If the data to be mined comes from several different sources data needs to be integrated which involves removing inconsistencies in names of attributes or attribute value names between data sets of different sources.

Data Cleaning – This step may involve detecting and correcting errors in the data, filling in missing values, etc. Some data cleaning methods are discussed in [6,7].

Discretization – When the data mining algorithm cannot cope with continuous attributes, discretization needs to be applied. This step consists of transforming a continuous attribute into a categorical attribute, taking only a few discrete values. Discretization often improves the comprehensibility of the discovered knowledge [8, 9].

Attribute Selection – not all attributes are relevant so for selecting a subset of attributes relevant for mining, among all original attributes, attribute selection is required.

3. Feature Selection

Many irrelevant attributes may be present in data to be mined. So they need to be removed. Also many mining algorithms don't perform well with large amounts of features or attributes. Therefore feature selection techniques need to be applied before any kind of mining algorithm is applied. The main objectives of feature selection are to avoid overfitting and improve model performance and to provide faster and more cost-effective models.

The selection of optimal features adds an extra layer of complexity in the modelling as instead of just finding optimal parameters for full set of features, first optimal feature subset is to be found and the model parameters are to be optimised [10]. Attribute selection methods can be broadly divided into filter and wrapper approaches. In the filter approach the attribute selection method is independent of the data mining algorithm to be applied to the selected attributes and assess the relevance of features by looking only at the intrinsic properties of the data. In most cases a feature relevance score is calculated, and low-scoring features are removed. The subset of features left after feature removal is presented as input to the classification algorithm. Advantages of filter techniques are that they easily scale to high-dimensional datasets are computationally simple and fast, and as the filter approach is independent of the mining algorithm so feature selection needs to be performed only once, and then different classifiers can be evaluated. Disadvantages of filter methods are that they ignore the interaction with the classifier and that most proposed techniques are univariate which means that each feature is considered separately, thereby ignoring feature dependencies, which may lead to worse classification performance when compared to other types of feature selection techniques. In order to overcome the problem of ignoring feature dependencies, a number of multivariate filter techniques were introduced, aiming at the incorporation of feature dependencies to some degree. Wrapper methods embed the model hypothesis search within the feature subset search. In the wrapper approach the attribute selection method uses the result of the data mining algorithm to determine how good a given attribute subset is. In this setup, a search procedure in the space of possible feature subsets is defined, and various subsets of features are generated and evaluated. The major characteristic of the wrapper approach is that the quality of an attribute subset is directly

measured by the performance of the data mining algorithm applied to that attribute subset. The wrapper approach tends to be much slower than the filter approach, as the data mining algorithm is applied to each attribute subset considered by the search. In addition, if several different data mining algorithms are to be applied to the data, the wrapper approach becomes even more computationally expensive [11]. Advantages of wrapper approaches include the interaction between feature subset search and model selection, and the ability to take into account feature dependencies. A common drawback of these techniques is that they have a higher risk of overfitting than filter techniques and are very computationally intensive. Another category of feature selection technique was also introduced, termed embedded technique in which search for an optimal subset of features is built into the classifier construction, and can be seen as a search in the combined space of feature subsets and hypotheses. Just like wrapper approaches, embedded approaches are thus specific to a given learning algorithm. Embedded methods have the advantage that they include the interaction with the classification model, while at the same time being far less computationally intensive than wrapper methods [12].

3. Classification

Data mining algorithms can follow three different learning approaches: supervised, unsupervised, or semi-supervised.

In supervised learning, the algorithm works with a set of examples whose labels are known. The labels can be nominal values in the case of the classification task, or numerical values in the case of the regression task.

In unsupervised learning, in contrast, the labels of the examples in the dataset are unknown, and the algorithm typically aims at grouping examples according to the similarity of their attribute values, characterizing a clustering task.

Finally, semi-supervised learning is usually used when a small subset of labeled examples is available, together with a large number of unlabeled examples.

The classification task can be seen as a supervised technique where each instance belongs to a class, which is indicated by the value of a special goal attribute or simply the class attribute. The goal attribute can take on categorical values,

each of them corresponding to a class. Each example consists of two parts, namely a set of predictor attribute values and a goal attribute value. The former are used to predict the value of the latter. The predictor attributes should be relevant for predicting the class of an instance. In the classification task the set of examples being mined is divided into two mutually exclusive and exhaustive sets, called the training set and the test set. The classification process is correspondingly divided into two phases: training, when a classification model is built from the training set, and testing, when the model is evaluated on the test set. In the training phase the algorithm has access to the values of both predictor attributes and the goal attribute for all examples of the training set, and it uses that information to build a classification model. This model represents classification knowledge – essentially, a relationship between predictor attribute values and classes – that allows the prediction of the class of an example given its predictor attribute values. For testing, the test set the class values of the examples is not shown. In the testing phase, only after a prediction is made is the algorithm allowed to see the actual class of the just-classified example. One of the major goals of a classification algorithm is to maximize the predictive accuracy obtained by the classification model when classifying examples in the test set unseen during training. The knowledge discovered by a classification algorithm can be expressed in many different ways like rules, decision trees, Bayesian network etc. Various techniques used for classification are explained in the following section.

4. Classification Techniques

4.1 Rule Based Classifiers

Rule based classifiers deals with the the discovery of high-level, easy-to-interpret classification rules of the form if-then. The rules are composed of two parts mainly rule antecedent and rule consequent. The rule antecedent, is the if part, specifies a set of conditions referring to predictor attribute values, and the rule consequent, the then part, specifies the class predicted by the rule for any example that satisfies the conditions in the rule antecedent. These rules can be generated using different classification algorithms, the most well known being the decision tree induction algorithms and sequential covering rule induction algorithms [13].

4.2 Bayesian Networks

A Bayesian network (BN) consists of a directed, acyclic graph and a probability distribution for each node in that graph given its immediate predecessors [14]. A Bayes Network Classifier is based on a bayesian network which represents a joint probability distribution over a set of categorical attributes. It consists of two parts, the directed acyclic graph G consisting of nodes and arcs and the conditional probability tables. The nodes represent attributes whereas the arcs indicate direct dependencies. The density of the arcs in a BN is one measure of its complexity. Sparse BNs can represent simple probabilistic models (e.g., naive Bayes models and hidden Markov models), whereas dense BNs can capture highly complex models. Thus, BNs provide a flexible method for probabilistic modeling [15].

4.3 Decision Tree

A Decision Tree Classifier consists of a decision tree generated on the basis of instances. The decision tree has two types of nodes: a) the root and the internal nodes, b) the leaf nodes. The root and the internal nodes are associated with attributes, leaf nodes are associated with classes. Basically, each non-leaf node has an outgoing branch for each possible value of the attribute associated with the node. To determine the class for a new instance using a decision tree, beginning with the root, successive internal nodes are visited until a leaf node is reached. At the root node and at each internal node, a test is applied. The outcome of the test determines the branch traversed, and the next node visited. The class for the instance is the class of the final leaf node [16].

4.4 Nearest Neighbour

A Nearest Neighbor Classifier assumes all instances correspond to points in the n -dimensional space. During learning, all instances are remembered. When a new point is classified, the k -nearest points to the new point are found and are used with a weight for determining the class value of the new point. For the sake of increasing accuracy, greater weights are given to closer points [17].

4.5 Artificial Neural Network

An artificial neural network, often just called a neural network is a mathematical model or

computational model based on biological neural networks, in other words, is an emulation of biological neural system. In most cases an ANN is an adaptive system that changes its structure based on external or internal information that flows through the network during the learning phase [18]. A Neural Network Classifier is based on neural networks consisting of interconnected neurons. From a simplified perspective, a neuron takes positive and negative stimuli (numerical values) from other neurons and when the weighted sum of the stimuli is greater than a given threshold value, it activates itself. The output value of the neuron is usually a non-linear transformation of the sum of stimuli. In more advanced models, the non-linear transformation is adapted by some continuous functions.

4.6 Support vector machines

Support Vector Machines [19] are basically binary classification algorithms. Support Vector Machines (SVM) is a classification system derived from statistical learning theory. It has been applied successfully in fields such as text categorisation, hand-written character recognition, image classification, biosequences analysis, etc. The SVM separates the classes with a decision surface that maximizes the margin between the classes. The surface is often called the optimal hyperplane, and the data points closest to the hyperplane are called support vectors. The support vectors are the critical elements of the training set. The mechanism that defines the mapping process is called the kernel function. The SVM can be adapted to become a nonlinear classifier through the use of nonlinear kernels. SVM can function as a multiclass classifier by combining several binary SVM classifiers. The output of SVM classification is the decision values of each pixel for each class, which are used for probability estimates. The probability values represent "true" probability in the sense that each probability falls in the range of 0 to 1, and the sum of these values for each pixel equals 1. Classification is then performed by selecting the highest probability. SVM includes a penalty parameter that allows a certain degree of misclassification, which is particularly important for nonseparable training sets. The penalty parameter controls the trade-off between allowing training errors and forcing rigid margins. It creates a soft margin that permits some misclassifications, such as it allows some training points on the wrong side of the hyperplane. Increasing the value of the

penalty parameter increases the cost of misclassifying points and forces the creation of a more accurate model that may not generalize well [20].

4.7 Rough Sets

Any set of all indiscernible (similar) objects is called an elementary set. Any union of some elementary sets is referred to as a crisp or precise set - otherwise the set is rough (imprecise, vague). Each rough set has boundary-line cases, i.e., objects which cannot be with certainty classified, by employing the available knowledge, as members of the set or its complement. Obviously rough sets, in contrast to precise sets, cannot be characterized in terms of information about their elements. With any rough set a pair of precise sets - called the lower and the upper approximation of the rough set is associated. The lower approximation consists of all objects which surely belong to the set and the upper approximation contains all objects which possible belong to the set. The difference between the upper and the lower approximation constitutes the boundary region of the rough set. Rough set approach to data analysis has many important advantages like provides efficient algorithms for finding hidden patterns in data, identifies relationships that would not be found using statistical methods, allows both qualitative and quantitative data, finds minimal sets of data (data reduction), evaluates significance of data, easy to understand [21].

4.8 Fuzzy Logic

Fuzzy logic is a multivalued logic different from "crisp logic", where binary sets have two valued logic. Fuzzy logic variables have truth value in the range between 0 and 1. Fuzzy logic is a superset of conventional Boolean logic that has been extended to handle the concept of partial truth. A membership function (MF) is a curve that defines how each point in the input space is mapped to a membership value (or degree of membership) between 0 and 1. Fuzzy Logic consists of Type 1 and Type 2 fuzzy logic. Type 1 fuzzy contains the constant values. A Type-2 Fuzzy Logic is an extension of Type 1 Fuzzy Logic in which the fuzzy sets comes from Existing Type 1 Fuzzy. A type-2 fuzzy set contains the grades of membership that are themselves fuzzy. A Type-2 membership grade can be any subset in the primary membership. For each primary membership there exists a secondary membership that defines the

possibilities for the primary membership. Type-1 Fuzzy Logic is unable to handle rule uncertainties. Type-2 Fuzzy Logic can handle rule uncertainties effectively and efficiently [22]. Type 2 Fuzzy sets are again characterized by IF–THEN rules [23]. Type-2 Fuzzy is computationally intensive because type reduction is very intensive. Type-2 fuzzy is used for modeling uncertainty and imprecision in a better way. The type-2 fuzzy sets are called as “fuzzy fuzzy” sets where the fuzzy degree of membership is fuzzy itself that results from Type 1 Fuzzy [24].

4.9 Genetic algorithms

Genetic Algorithms (GA) are search algorithms based on natural genetics that provide robust search capabilities in complex spaces, thereby offering a valid approach to problems requiring efficient and effective search processes [25]. GA is an iterative process that operates on a population, i.e., a set of candidate solutions. Each solution is obtained by means of an encoding/decoding mechanism, which enables us to represent the solution as a chromosome and vice versa. Initially, the population is randomly generated. Every individual in the population is assigned, by means of a fitness function, a fitness value that reflects its quality with respect to solving the particular problem. A chromosome is evaluated by a fitness function to determine the quality of the solution, i.e., how effective it is in solving the problem. The input of the fitness function is the chromosome and the output is the fitness value of this chromosome. In each cycle, fitness of each candidate solution is determined. The next stage is selection, where a temporary population is created in which the fittest individuals are likely to have a higher number of chances than less fit individuals, to be used as parents for the next generation. The reproductive operators like crossover and mutation are applied to the individuals in this population yielding a new population [26].

References:

- [1] I.H. Witten, E. Frank and M.A. Hall, *Data mining practical machine learning tools and techniques*, Morgan Kaufmann publisher, Burlington 2011
- [2] J. Han and M. Kamber, *Data mining concepts and techniques*, Morgan Kaufmann, San Francisco 2006
- [3] T.J. Shan, H. Wei and Q. Yan, “Application of genetic algorithm in data mining”, *1st Int Work Educ Technol Comput Sci, IEEE 2*, 2009, pp. 353- 356
- [4] Z.Z. Shi, *Knowledge discovery*, Tsinghua University Press, Beijing, 2001
- [5] D. Pyle, *Data preparation for data mining*, 1st Vol., Morgan Kaufmann publisher, San Francisco, 1999
- [6] I. Guyon, N. Matic and V. Vapnik, “Discovering informative patterns and data cleaning”, In: *Fayyad UM, Piatetsky-Shapiro G, Smyth P and Uthurusamy R. (ed) Advances in knowledge discovery and data mining*, AAAI/MIT Press, California, 1996, pp. 181-203
- [7] E. Simoudis, B. Livezey B and R. Kerber R , “Integrating inductive and deductive reasoning for data mining”, In: *Fayyad UM, Piatetsky-Shapiro G, Smyth P and Uthurusamy R. (Eds.) Advances in knowledge discovery and data mining*, AAAI/MIT Press, California, 1996, pp. 353-373
- [8] B. Pfahringer, “Supervised and unsupervised discretization of continuous features”, *Proc. 12th Int. Conf. Machine Learning*, 1995, pp. 456-463.
- [9] J. Catlett, “On changing continuous attributes into ordered discrete attributes”, In *Y. Kodratoff (ed), Machine Learning—EWSL-91*, Springer-Verlag, New York, 1991, pp 164-178
- [10] W. Daelemans, V. Hoste, F.D. Meulder and B. Naudts, “Combined Optimization of Feature Selection and Algorithm Parameter Interaction in Machine Learning of Language”, *Proceedings of the 14th European Conference on Machine Learning (ECML-2003)*, Lecture Notes in Computer Science 2837, Springer-Verlag, Cavtat-Dubrovnik, Croatia, 2003, pp. 84-95
- [11] M.L. Raymer, W.F. Punch, E.D. Goodman, L.A. Kuhn and A.K. Jain, “Dimensionality Reduction Using Genetic Algorithms”, *IEEE Transactions On Evolutionary Computation*, Vol. 4, No. 2, 2000
- [12] Y. Saeys, I. Inza and P. Larranaga, “A review of feature selection techniques in bioinformatics”, *Bioinformatics-19*, 2007, pp. 2507–17.
- [13] G.L. Pappa and A.A. Freitas, *Automating the Design of Data Mining Algorithms. An Evolutionary Computation Approach*, Natural Computing Series, Springer, 2010
- [14] A. Darwiche, *Modeling and Reasoning with Bayesian Networks*, Cambridge University Press, 2009
- [15] G.F. Cooper, P. Hennings-Yeomans, S. Visweswaran and M. Barmada, “An Efficient Bayesian Method for Predicting Clinical Outcomes from Genome-Wide Data”, *AMIA 2010 Symposium Proceedings*, 2010, pp. 127-131
- [16] M. Garofalakis, D. Hyun, R. Rastogi and K. Shim, “Building Decision Trees with Constraints”, *Data Mining and Knowledge Discovery*, vol. 7, no. 2, 2003, pp. 187 – 214
- [17] T.M. Mitchell, *Machine Learning*, McGraw-Hill Companies, USA, 1997
- [18] Y. Singh Y, A.S. Chauhan, “Neural Networks in Data Mining”, *Journal of Theoretical and Applied Information Technology*, 2005, pp. 37-42

- [19] V. N. Vapnik, *Statistical Learning Theory*, Wiley New York., 1998
- [20] C.W. Hsu, C.C. Chang and C.J. Lin, "A practical guide to support vector classification", <http://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>, 2003
- [21] Z. Pawlak, "Rough sets", *International Journal of Computer and Information Sciences*, 1982, pp. 341-356
- [22] L. Tari, C. Baral and S. Kim, "Fuzzy c-means clustering with prior biological knowledge", *Journal of Biomedical Informatics*, 42(1), 2009, pp. 74-81
- [23] N.N. Karnik, J.M. Mendel and Q. Liang, "Type-2 Fuzzy Logic Systems", *IEEE Transactions on Fuzzy Systems*, Vol. 7, No. 6, 1999, 643-658
- [24] J.R. Castro, O. Castillo and L.G. Martínez, "Interval Type-2 Fuzzy Logic Toolbox", *Engineering Letters* 15(1), 2007, pp. 89-98
- [25] J.H. Holland, *Adaptation in Natural and Artificial Systems*, University of Michigan Press, Ann Arbor, MI, 1975
- [26] D.E. Goldberg, *Genetic algorithms in search, optimization and machine learning*, Addison-Wesley, New York, 1989