

ANN: A Novel Technique in Data Mining

Aleem Ali¹, Naresh Kumar¹, Sanyogita Chouhan²

1). F/O Engg, University Polytechnic, JMI New Delhi-110025

2). M.Tech(CS) Scholar, Department Of Computer Science, JamiaHamdard(Hamdard University)-110065

Abstract

Data mining is a multidisciplinary field, can be described as "the nontrivial extraction of implicit, previously unknown, and potentially useful information from data and the science of extracting useful information from large data sets or databases. We present technique for the discovery of patterns hidden in large data sets, focusing on issues relating to their feasibility, usefulness, effectiveness, and scalability. Predictive data mining is the most common type of data mining and one that has the most direct business applications. This paper is an overview of artificial neural networks and questions their position as a preferred tool by data mining practitioners and the key technology and ways to achieve the data mining based on neural networks are also researched.

Keywords: Artificial Neural Network (ANN), Advantages, Back propagation algorithm, Data mining, Data modeling.

1. INTRODUCTION

Data Mining is an analytic process designed to explore data (usually large amounts of data - typically business or market related) in search of consistent patterns and systematic relationships between variables, and then to validate the findings by applying the detected patterns to new subsets of data[1]. The ultimate goal of data mining is knowledge discovery and prediction. Knowledge discovery provides explicit information that has a readable form and can be understood by a user[2]. Forecasting, or predictive modeling provides predictions of future events and may be transparent and readable in some approaches (e.g. rule based systems) and opaque in others such as neural networks. Moreover, some data mining systems such as neural networks are inherently geared towards prediction and pattern recognition, rather than knowledge discovery[3].

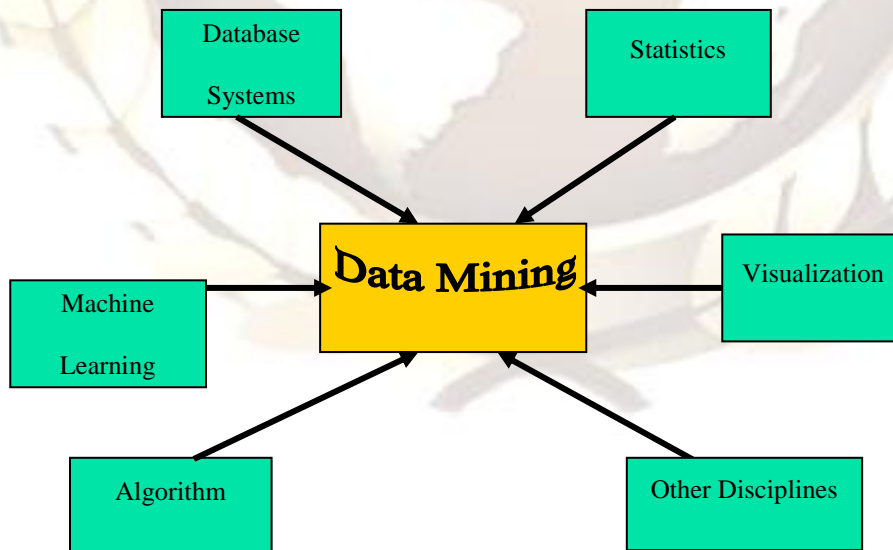


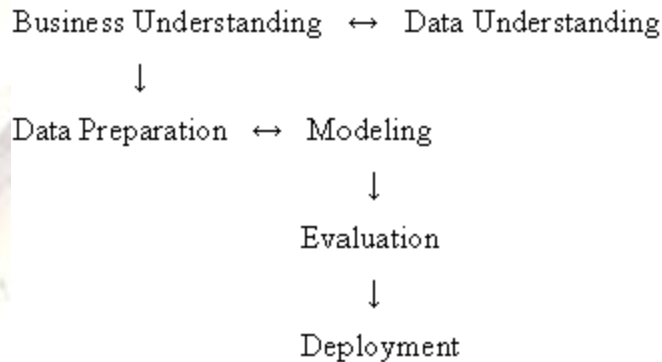
Fig. 1

Data mining relies on the use of real world data (typically in business applications). This data is extremely vulnerable to collinearity precisely because data from the real world may have unknown interrelations.

2. MODELS FOR DATA MINING

In the business environment, complex data mining projects may require the coordinate efforts of various experts, stakeholders, or departments throughout an entire organization. In the data mining literature, various "general frameworks" have been proposed to serve as blueprints for how to organize the process of gathering data, analyzing data, disseminating results, implementing results, and monitoring improvements.

One such model, CRISP (Cross-Industry Standard Process for data mining) was proposed in the mid-1990s by a European consortium of companies to serve as a non-proprietary standard process model for data mining. This general approach postulates the following (perhaps not particularly controversial) general sequence of steps for data mining projects:



Another approach - the Six Sigma methodology - is a well-structured, data-driven methodology for eliminating defects, waste, or quality control problems of all kinds in manufacturing, service delivery, management, and other business activities. This model has recently become very popular (due to its successful implementations) in various American industries, and it appears to gain favor worldwide. It postulated a sequence of, so-called, DMAIC steps -

Define → Measure → Analyze → Improve → Control

that grew up from the manufacturing, quality improvement, and process control traditions and is particularly well suited to production environments (including "production of services," i.e., service industries).

Another framework of this kind (actually somewhat similar to Six Sigma) is the approach proposed by SAS Institute called SEMMA -

Sample → Explore → Modify → Model → Assess

which is focusing more on the technical activities typically involved in a data mining project. All of these models are concerned with the process of how to integrate data mining methodology into an organization, how to "convert data into information," how to involve important stakeholders, and how to disseminate the information in a form that can easily be converted by stakeholders into resources for strategic decision making.

3. ARTIFICIAL NEURAL NETWORKS

Neural Networks are analytic techniques modeled after the (hypothesized) processes of learning in the cognitive system and the neurological functions of the brain and capable of predicting new observations (on specific variables) from other observations (on the same or other variables) after executing a process of so-called *learning* from existing data[4]. Neural networks is one of the Data mining techniques.

Some authors stress the fact that *neural networks* use, or we should say are expected to use, massively parallel computation models. For example Haykin (1994) defines *neural network* as: "a massively parallel distributed processor that has a natural propensity for storing experiential knowledge and making it available for use. It resembles the brain in two respects: (1) Knowledge is acquired by the network through a learning process, and (2) Interneuron connection strengths known as synaptic weights are used to store the knowledge".

The first step is to design a specific network architecture (that includes a specific number of "layers" each consisting of a certain number of "neurons") [5]. The size and structure of the network needs to match the nature (e.g., the formal complexity) of the investigated phenomenon. Because the latter is obviously not known very well at this early stage, this task is not easy and often involves multiple "trials and errors." (Now, there is, however, neural network software that applies artificial intelligence techniques to aid in that tedious task and finds "the best" network architecture.)

The new network is then subjected to the process of "training." In that phase, neurons apply an iterative process to the number of inputs (variables) to adjust the weights of the network in order to optimally predict (in traditional terms, we could say find a "fit" to) the sample data on which the "training" is performed. After the phase of learning from an existing data set, the new network is ready and it can then be used to generate predictions.

One of the major advantages of *neural networks* is that, theoretically, they are capable of approximating any continuous function, and thus the researcher does not need to have any hypotheses about the underlying model, or even to some extent, which variables matter. An important disadvantage, however, is that the final solution depends on the initial conditions of the network, and, as stated before, it is virtually impossible to "interpret" the solution in traditional, analytic terms, such as those used to build theories that explain phenomena.

Neural networks are essentially comprising three pieces: the architecture or model; the learning algorithm; and the activation functions. (Fausett (1994)) Neural networks are programmed or "trained" to "store, recognize, and associatively retrieve patterns or database entries; to solve combinatorial optimization problems; to filter noise from measurement data; to control ill-defined problems; in summary, to estimate sampled functions when we do not know the form of the functions" (Kosko (1992), p.13). It is precisely have these two abilities pattern recognition and function.

4. NEURAL NETWORK IN DATA MINING

Neural network method is used for classification, clustering, feature mining, prediction and pattern recognition [6]. It imitates the neurons structure of animals, based on the M-P model and Hebb's learning rule, so in essence it is a distributed matrix structure. Through training data mining, the neural network method gradually calculates (including repeated iteration or cumulative calculation) the weights the neural network connected. The neural network model can be broadly divided into the following three types:

- (1) Feed-forward networks: it regards the perception back-propagation model and the function network as representatives, and mainly used in the area such as prediction and pattern recognition;
- (2) Feedback network: it regards Hopfield discrete model and continuous model as representatives, and mainly used for associative memory and optimization calculation;
- (3) Self-organization networks: it regards adaptive resonance theory (ART) model and Kohonen model as representatives, and mainly used for cluster analysis.

At present, the neural network most commonly used in data mining is BP network [7]. Of course, artificial neural network is the developing science, and some theories have not really taken shape, such as the problems of convergence, stability, local minimum and parameters adjustment. For the BP network the frequent problems it encountered are that the training is slow, may fall into local minimum and it is difficult to determine training parameters. Aiming at these problems some people adopted the method of combining artificial neural networks and genetic algorithm to achieve better results.

Artificial neural network has the characteristics of distributed information storage, parallel processing, information reasoning, and self-organization learning, and has the capability of rapid fitting the non-linear data, so it can solve many problems which are difficult for other methods to solve.

4.1 FEEDFORWARD NEURAL NETWORK

One of the simplest feedforward neural networks (FFNN), such as in Figure, consists of three layers: an input layer, hidden layer and output layer [8]. In each layer there are one or more processing elements (PEs). PEs is meant to simulate the neurons in the brain and this is why they are often referred to as neurons or nodes. A PE receives inputs from either the outside world or the previous layer. There are connections between the PEs in each layer that have a weight (parameter) associated with them. This weight is adjusted during training. Information only travels in the forward direction through the network - there are no feedback loops.

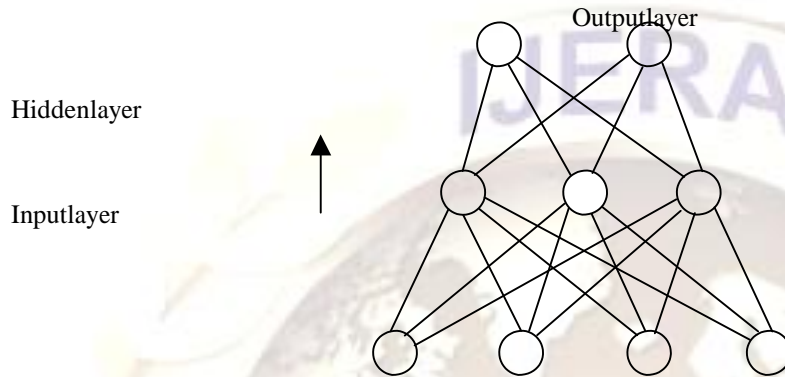


Fig 2. Multilayered feed-forward neural network (ANN)

The simplified process for training a FFNN is as follows:

1. Input data is presented to the network and propagated through the network until it reaches the output layer. This forward process produces a predicted output.
2. The predicted output is subtracted from the actual output and an error value for the network is calculated.
3. The neural network then uses supervised learning, which in most cases is back propagation, to train the network. Back propagation is a learning algorithm for adjusting the weights. It starts with the weights between the output layer PE's and the last hidden layer PE's and works backwards through the network.
4. Once back propagation has finished, the forward process starts again, and this cycle is continued until the error between predicted and actual outputs is minimized.

4.2. THE BACK PROPAGATION ALGORITHM:

Back propagation, or propagation of error, is a common method of teaching artificial neural network how to perform a given task. The back propagation algorithm is used in layered feed-forward ANNs. This means that the artificial neurons are organized in layers, and send their signals "forward", and then the errors are propagated backwards [9]. The back propagation algorithm uses supervised learning, which means that we provide the algorithm with examples of the inputs and outputs we want the network to compute, and then the error (difference between actual and expected results) is calculated. The idea of the back propagation algorithm is to reduce this error, until the ANN *learns* the training data [10]. **Summary of the technique:**

1. Present a training sample to the neural network.
2. Compare the network's output to the desired output from that sample. Calculate the error in each output neuron.
3. For each neuron, calculate what the output should have been, and a *scaling factor*, how much lower or higher the output must be adjusted to match the desired output. This is the local error.
4. Adjust the weights of each neuron to lower the local error.

Actual Algorithm:

1. Initialize the weights in the network (often randomly)
2. repeat
 - *for each example e in the training set do
 - 1. $O = \text{neural-net-output}(\text{network}, e)$;
 - forward pass

2. $T = \text{teacher output for } e$
 3. Calculate error $(T - O)$ at the output units
 4. Compute Δ_{wi} for all weights from hidden layer to output layer ; backward pass
 5. Compute Δ_{wi} for all weights from input layer to hidden layer ; backward pass continued
 6. Update the weights in the network
- *end
3. until all examples classified correctly or stopping criterion satisfied
 4. return(network)

5. ADVANTAGES

Artificial neural networks (ANN) are just one of the tools used to find patterns in the data and to infer rules from them. Neural networks are useful in providing information on associations, classifications, clusters, and forecasting. Both the IRS and Wrangler have used neural networks in a data mining situation with good success. More examples would have been given, but the Internet search of data mining and neural networks revealed only these cases. We anticipate as time passes, and data mining grows more case studies will become available. We have also seen that for best results with neural networks a working knowledge of statistical models is desired. With all the common material between the two disciplines, neural networks and statistics, better communication between them would be advantageous to both. Computers have a long way to go before they can rival the human brain on the same parallel scale but neural networks are a start in the right direction.

REFERENCES

- [1] Aleem Ali, A Concise Artificial Neural Network in Data Mining, *International Journal of Research in Engineering & Applied Sciences*, Feb 12, 2(2), p. 418-428, 2012.
- [2] O. Maimon and Rokach (eds.), Data Mining and Knowledge Discovery Handbook: A Complete Guide for Practitioners and Researchers (Kluwer Academic Publishers, 2007).
- [3] Agrawal R and Srikant R, Fast algorithms for mining association rules, *Proceedings of the 20th VLDB Conference*, 1994, 487-499.
- [4] <http://en.wikipedia.org>
- [5] Fausett, Laurene, *Fundamentals of Neural Networks: Architectures, Algorithms and Applications* (Prentice-Hall, New Jersey, USA, 1994).
- [6] G Towell, J W Shavlik. The extraction of refined rules from knowledge based neural networks, *Machine Learning*, 1993(13):71-101.
- [7] Rumelhart, D. E and Zipser, D, Feature discovery by competitive learning. *Cognitive Science*, 1985, 9:75-112.
- [8] Zurada J.M, *An introduction to artificial neural systems* (St. Paul: West Publishing, 1992).
- [9] Aleksander I. and Morton H, *An introduction to Neural Computing*, (2nd edition, 1998).
- [10] Maulik V. Dhamecha and Bhavesh R. Akbari, COMPREHENSIVE STUDY OF MODIFIED ARTIFICIAL NEURAL NETWORK ALGORITHM AND COMPARE EFFICIENCY AND PERFORMANCE WITH BASIC ALGORITHM OF NEURAL NETWORK IN DATA MINING, *International Journal of Research in Engineering & Applied Sciences*, 2(2), 2012, pp. 29-39.