# A Review on Artificial Neural Network Approach in Data Mining

**Mithun Mhatre[1] and Ranjeet Pawar[2]**
In-charge HOD, Department of Computer Technology[1]
In-charge HOD, Department of Information Technology[2]
Bharti Vidyapeeth Institute of Technology, Navi Mumbai, Maharashtra, India

**Abstract:** *Data mining has been used as a term describing explorative analysis of large data sets (frequently stored in data warehouses with the objective to identify hidden relationships among the variables in the set). Artificial neural network is one of many tools for data mining. This paper summarizes the state-of-the-art of the principles beyond using neural models in data mining. Artificial Neural Networks are suitable in data-rich environments and are typically used for extracting embedded knowledge in the form of rules, quantitative evaluation of these rules, clustering, self-organization, classification and regression, feature evaluation and dimensionality reduction. In this paper, we try to understand the basics of neural network modeling, some specific applications, and the process of implementing a neural network in data mining. Like shake and tap on SOS button.*
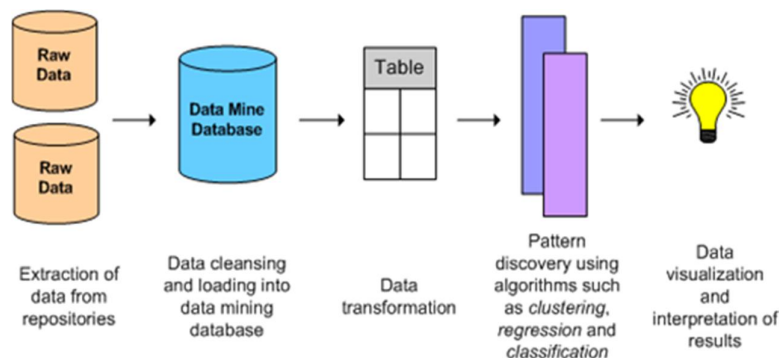
**Keywords:** ANN (Artificial Neural Network), Data pre-processing, Rules extraction, Data mining

## I. INTRODUCTION

Data mining (DM) is the nontrivial extraction of implicit, previously unknown, interesting, and potentially useful information (usually in the form of knowledge patterns or models) from data. Historically data mining has grown from large business database applications, such as finding patterns in customer purchasing activities from transactions databases.

The most vehiculated Data Mining problems are reduced to traditional statistical and machine leaning methods: classification, prediction, association rule extraction, and sequence detection. The following are the major stages in solving a DM problem.

1. Define the problem.
2. Collect and select data, such as deciding which data to collect and how to collect them
3. Prepare data, such as transform data to a certain format, or data cleansing.
4. Data preprocessing; this task is concerned mainly with enhancement of data quality.
5. Select an appropriate mining method, which consists of:
    a. Selecting a model or algorithm.
    b. Selecting model/algorithm training parameters.
6. Training/testing the data or applying the algorithm, where evaluation set of data is used in the trained architecture.
7. Final integration and evaluation of the generated model.

## II. WHAT IS ARTIFICIAL NEURAL NETWORK?

Artificial Neural Networks are relatively crude electronic models based on the neural structure of the brain. The brain basically learns from experience. It is natural proof that some problems that are beyond the scope of current computers are indeed solvable by small energy efficient packages. This brain modeling also promises a less technical way to develop machine solutions. This new approach to computing also provides a more graceful degradation during system overload than its more traditional counterparts. These biologically inspired methods of computing are thought to be the next major advancement in the computing industry. Even simple animal brains are capable of functions that are currently impossible for computers

This research shows that brains store information as patterns. Some of these patterns are very complicated and allow us the ability to recognize individual faces from many different angles. This process of storing information as patterns, utilizing those patterns, and then solving problems encompasses a new field in computing.

## III. WORKING OF ANN

The other parts of using neural networks revolve around the myriad of ways these individual neurons can be clustered together. This clustering occurs in the human mind in such a way that information can be processed in a dynamic, interactive, and self-organizing way. Biologically, neural networks are constructed in a three-dimensional world from microscopic components. These neurons seem capable of nearly unrestricted interconnections.
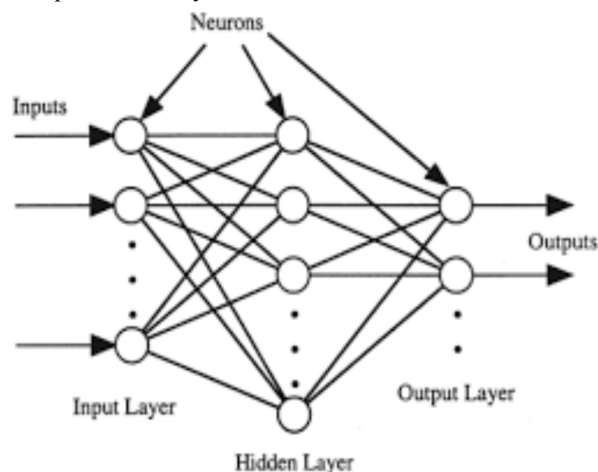


**Figure 1:** A Simple Neural Network Diagram

Although there are useful networks which contain only one layer, or even one element, most applications require networks that contain at least the three normal types of layers - Input, Hidden, and Output.

- The layer of input neurons receives the data either from input files or directly from electronic sensors in real-time applications.
- The output layer sends information directly to the outside world, to a secondary computer process, or to other devices such as a mechanical control system.
- Between these two layers can be many hidden layers. These internal layers contain many of the neurons in various interconnected structures. The inputs and outputs of each of these hidden neurons simply go to other neurons.

These lines of communication from one neuron to another are important aspects of neural networks. There are two types of these connections. One causes the summing mechanism of the next neuron to add while the other causes it to subtract. In more human terms one excites while the other inhibits. Some networks want a neuron to inhibit the other neurons in the same layer. This is called lateral inhibition. The most common use of this is in the output layer. This concept is also called competition. Another type of connection is feedback. This is where the output of one-layer routes back to a previous layer.

## IV. TRAINING OF ARTIFICIAL NEURAL NETWORKS

- **A neural network** has to be configured such that the application of a set of inputs produces (either 'direct' or via a relaxation process) the desired set of outputs. Various methods to set the strengths of the connections exist. One way is to set the weights explicitly, using a priori knowledge. Another way is to **'train' the neural network** by feeding it teaching patterns and letting it change its weights according to some learning rule. We can categorize the learning situations as follows:

- **Supervised learning** or Associative learning in which the network is trained by providing it with input and matching output patterns. These input-output pairs can be provided by an external teacher, or by the system which contains the neural network (self-supervised).

- **Unsupervised learning** or Self-organization in which an (output) unit is trained to respond to clusters of patterns within the input. In this paradigm the system is supposed to discover statistically salient features of the input population. Unlike the supervised learning paradigm, there is no a priori set of categories into which the patterns are to be classified; rather the system must develop its own representation of the input stimuli.

- **Reinforcement Learning** This type of learning may be considered as an intermediate form of the above two types of learning. Here the learning machine does some action on the environment and gets a feedback response from the environment. The learning system grades its action good (rewarding) or bad(punishable) based on the environmental response and accordingly adjusts its parameters.

- **Neural networks for data mining** soft computing methodologies (involving fuzzy sets, NN, genetic algorithms, and rough sets) are most widely applied in the DM. Fuzzy sets provide a natural framework for the process in dealing with uncertainty. NN and rough sets are used for classification and rule generation. Genetic algorithms are involved in various optimization and search processes, like query optimization and template selection. It is presently hard to separate NN as a distinct tool. The main contribution of NN toward DM stems from rule extraction and from clustering.

- **Rule Extraction and Evaluation** Typically, a network is first trained to achieve the required accuracy rate. Redundant connections of the network are then removed using a pruning algorithm. The link weights and activation values of the hidden units in the network are analysed, and classification rules are generated

- **Clustering and Dimensionality Reduction:** Korhonen's SOM proved to be an appropriate tool for handling huge data bases. Korhonen have demonstrated the utility of a SOM with more than one million nodes to partition a little less than seven million patent abstracts where the documents are represented by 500-dimensional feature vectors.

- **Incremental Learning:** When designing and implementing data mining applications for large data sets, we face processing time and memory space problems. e. The fundamental issue in incremental learning is: how can a learning system adapt to new information without corrupting or forgetting previously learned information

## V. ADVANTAGES OF NEURAL NETWORKS

1. **High Accuracy:** Neural networks are able to approximate complex non-linear mappings
2. **Noise Tolerance:** Neural networks are very flexible with respect to incomplete, missing and noisy data.
3. **Independence from prior assumptions:** Neural networks do not make a priori assumptions about the distribution of the data, or the form of interactions between factors.
4. **Ease of maintenance:** Neural networks can be updated with fresh data, making them useful for dynamic environments.
5. Neural networks can be implemented in parallel hardware
6. When an element of the neural network fails, it can continue without any problem by their parallel nature.

The following are examples of where neural networks have been used.

### 5.1 Accounting

- Identifying tax fraud
- Enhancing auditing by finding irregularities

### 5.2 Finance
- Signature and bank note verification
- Risk Management
- Foreign exchange rate forecasting
- Bankruptcy prediction
- Customer credit scoring
- Credit card approval and fraud detection
- Forecasting economic turning points
- Bond rating and trading
- Loan approvals
- Economic and financial forecasting

### 5.3 Marketing
- Classification of consumer spending pattern
- New product analysis
- Identification of customer characteristics
- Sale forecasts

### 5.4 Human resources
- Predicting employee's performance and behaviour
- Determining personnel resource requirements

## VI. CONCLUSION

Neural Networks are suitable for problems, whose inputs are both categorical and numeric, and where the relationships between inputs and outputs are not linear or the input data are not normally distributed. In such cases, classical statistical methods may not be reliable enough. Because neural network does not make any assumptions about the data distribution, their power is less affected than traditional statistical methods when data are not properly distributed. Finally, there are cases in which the neural networks simply provide one more way of building a predictive model for the situation at hand. Given the ease of experimentation using the available software tools, it is certainly worth exploring the power of neural networks in any data modelling situation.

## REFERENCES

[1]. Wang L, Sui TZ. Application of data mining technology based on neural network in the engineering. Proceedings of International Conference on Wireless Communications, Networking and Mobile Computing; Shanghai, China. 21–25 September 2007; pp. 5544–5547

[2]. Nirkhi S. Potential use of artificial neural network in data mining. Proceedings of International Conference on Computer and Automation Engineering; Singapore. 26–28 February 2010; pp. 339–343

[3]. Biryulev C, Yakymiv Y, Selemonavichus A. Research of artificial neural networks usage in data mining and semantic integration. Proceedings of International Conference on Perspective Technologies and Methods in MEMS Design; Lviv, Ukraine. 20–23 April 2010; pp. 144–149.

[4]. Setiono R, Baesens B, Mues C. A note on knowledge discovery using neural networks and its application to credit screening. Eur. J. Operation. Res. 2009; 192:326–332.

[5]. Craven M, Shavlik J. Using neural networks for data mining. Future Gener. Comput. Syt. 1997; 13:211–229.

[6]. Fu L. Rule learning by searching on adapted nets. Proceedings of National Conference on Artificial Intelligence; Anaheim, CA, USA. 1991. pp. 590–595.

[7]. Parekh R, Yang J, Honavar V. Constructive neural network learning algorithms for pattern classification. IEEE Trans. Neural Netw. 2000; 11:436–451.