# Data Mining and Bagging and Boosting for Uncertain Numerical Data: A Fundamental Study

Santosh S Lomte[1], Sanket Gunderao Torambekar[2]

*[1]Radhai Mahavidyalay., Aurangabad, Maharashtra, India*

*[2]K.P.C.Yogeshwari Polytechnic .,Ambajogai 431517, Maharashtra, India*

## *Abstract*

Boosting and bagging are one of the most common re-sampling grouping techniques that use the same learning algorithm for the base classifiers to combine and generate a variety of classifiers. Algorithms of boosting are thought to be stronger than bagging on noise-free information. There are clear scientific principles. signs, however that bagging in noisy settings is far more resilient compared to boosting. In ambiguous data management, the Possible Universe has grown to be one of the most powerful methods for working with different forms of uncertainty concerning data. However on the basis of a hypothetical universe, a few unknown algorithms for data classification are suggested. For some data, most current unknown algorithms for data classification are clearly generalised from conventional classification algorithms. They cope with the volatility of data dependent on the reasonably ideal distribution of probability and expectations of the data form, so it is difficult to apply for different application scenarios. Classification is one of the most important techniques for data mining and Decision Tree is the most common classification structure used in many applications. The classifier of the decision tree operates on precise and known data. Traditional classifiers of the decision tree function with data whose values are known and correct. During the process of data collection, value ambiguity occurs in many applications. Measurement/quantisation errors, data staleness, and numerous repetitive measurements include examples of sources of uncertainty. The basic fundamentals of classification using different techniques and requirements are primarily applied to the Decision Tree Classifier in this paper. For the assembly of the classifiers, the job also has bagging and boosting, the unknown numerical data is regarded as the data form for the process flow.

*Keywords:* Low Density Polymers; High Density polymers; Cross-linked Polymers; Insulation; Cables; Power; Electrical; Voltage.

## 1. Introduction

In data mining and machine learning, classification plays a crucial role. Traditional algorithms for classification concentrate on certain details. In many actual applications, however, data ambiguity inevitably occurs [1-4]. When we monitor an object's position with GPS devices, for instance, the recorded location can have errors of a several meters [5]. For another illustration, because of the presence of multiple noisy variables, sensor measurements can be imprecise to any degree.[6]. In addition, the treatment of probe-level ambiguity of gene expression microarray data is also a central study factor in the biomedical research domain [7]. Uncertain knowledge has provided conventional

classification algorithms with a major challenge. Several techniques have been suggested to deal with data uncertainty [8], including the potential world model that has been displayed to be successful in coping with different forms of uncertainty concerning data in the management of ambiguous data [9-12]. Nonetheless, as far as we do, few unpredictable algorithms for classification of data are built on the basis of a possible universe.

In data mining, classification is a classic problem. The role is to create a model algorithmically which predicts the class mark of an unseen test tuple centred on the tuple's function vector, provided a set of tuples of training results, each with a class mark and described by a function vector. The decision tree model is one of the most common models of classification. Decision trees are common because they are easy to understand and are realistic. Laws can also be readily omitted from decision trees. For decision tree construction, several algorithms, for example IC4.5 and ID3, have been invented. In a broad variety of applications, such algorithms are commonly implemented and used, including image detection, medical diagnosis, loan applicant credit ranking, science research, detection of fraud, and target marketing.

A function A tuple (attribute) is either categorical or numerical in design in conventional decision tree classification. It is commonly thought that an accurate and definite point value is for the latter. Nonetheless, data instability is popular in various systems. Hence A feature/attribute worth is ideally captured, not by value of a single point value, but by a set of values that gives rise to a distribution of probability. The abstraction of distributions of odds through summary statistics or example variance and means is an easy way to manage data uncertainty. We name this process averaging. Another approach is to accept the entire knowledge carried to construct a decision tree by the probability distributions.

Current uncertain algorithms for data classification depend mainly on random or probabilistic sampling. concepts for the extension of

standard Algorithms for classification of unknown results, for example Bayes algorithms[13,14], decision tree-based algorithms[15,16], nearest neighbor-based algorithms[17,18], SVM algorithms[19,20], rule-based algorithms[21,22], FDA algorithms[21,22], For certain data, of these algorithms is merely expanded from a single conventional algorithm of classification, thereby eventually inheriting the inherent limitations of the real algorithms[26]. In addition, they struggle with the volatility of data dependent on reasonably ideal assumptions of distribution of probabilities and data form, so it is difficult to apply it to different Scenarios for deployment. A potential world-based paradigm for classification of unknown details has recently been developed (PWCLA)[27], which uses the concept of consistency to learn an affinity matrix of agreement for ambiguous data and then uses To extend the model Spectral examination in order to classify unexplained effects. In essence, however the transudative classification algorithm is a semi-supervised PWCLA, since it involves awareness of all test data throughout the training phase, i.e. it is not usable in cases where the test data during the training is totally unseen.

Numerous methods for the development of an ensemble of classifiers have been proposed[28]. Mechanisms needed to construct a classifier category include:

i.   Use distinct training data sub-sets of a single learning process,
ii.  Using multiple training criteria for a single form of training,
iii. Use multiple types of learning.

Open and casual rationale, from a mathematical, representational and computed point of view, is given in[28] on why ensembles can enhance performance. If the classifiers are sufficiently distinct from each other in a scheme or in other words, that the particular classifiers contain minimum of popular failures is the secret to the success of sets. If a mistake is made by one classifier, the others should not be likely to make the same mistake.

Bagging [29] and boosting [30] are the common ensemble algorithms. Boosting and bagging have two big variations. First, improving the adaptive distribution of the training set based on the output of previously generated classifiers, thereby altering the stochastic distribution of the training set by bagging. Second, boosting uses a feature of a classifier's output as a voting weight, whereas bagging utilises equal weight voting. Algorithms for boosting are deemed better compared to bagging on data without noise, but in noisy environments, Bagging is far more durable compared to boosting.

## 2. Ensembles of Classifiers

Learning algorithms attempt to search a hypothesis in a specified space H of hypotheses and, if we have sufficient evidence, they will find the best one for a given problem in certain instances. But we only have minimal data sets in actual situations, and often only there are a few samples available. Learning algorithm will find many hypotheses in these cases that seem same correct with regard to the training data available, and while we can often pick the easiest one or the lowest ability among them, to get a reasonable estimate of the uncertain genuine theory, we may prevent the issue of mixing them.

There is therefore an increasing awareness that classifier combinations can be more efficient than single classifiers. When a mixture will achieve a more accurate and specific outcome of many, why depend on the best single classifier? This is basically the logic behind the principle of multiple systems for classifiers.

### 2.1 Boosting and Bagging

Aggregation Bootstrap, or bagging, is a tool suggested by [31] which is used to decrease the uncertainty Linked with forecasting and thus enhance the phase of prediction with many classification methods and regression methods. This is a pretty clear concept: several samples of bootstraps are taken from the data available, each bootstrap sample is added to some prediction process, and the observations are then mixed. to obtain the overall prediction by averaging regression and basic classification voting, with a reduction in variance, because of the average.

Boosting is a committee-based technique, like bagging, and can be used to maximise the precision of methods of regression and classification. A weighted average of outcomes received by applying a projection process is used for boosting to different samples, in comparison to bagging, which utilises a simple average of outcomes to produce an complete prediction. Also with boosting, not all the samples used at each step are taken from the same population in the same way, but rather, during the next step, increased weight is assigned to the incorrectly expected cases of a given stage. Boosting is thus an iterative process, integrating weights, in comparison to being centred, as is the case with bagging, on a simple average of predictions. Moreover, boosting is also used in poor learners (e.g., a basic classifier such as a decision tree of two nodes), while bagging is not the case.

The predecessor was developed by Schapire (1990) [32] to subsequently improve algorithms developed by him and others. His initial solution involved two-class classifiers., and pooled the outcomes of three classifiers, generated by simple majority voting from separate learning samples. By integrating the outcomes of a greater number of weak learners, Freund (1995) [33] expanded Schapire's original system. Then the algorithm AdaBoost was developed by Schapire and Freund (1996) [34], which became quite famous quickly. Breiman (1998) [35] extended the complete boosting policy and consideration of the algorithm of Schapire and Freund as a special scenario' of the arcing algorithm class, proposing the term arcing via adaptive merging and resampling. But this chapter will concentrate on AdaBoost in the interest of brevity, and because of the algorithm success of Freund and Schapire, and refer to similar methods only momentarily.

Many of the authors have stated this raise is one of the best important ideas for statistical/machine

learning to be implemented during the nineties, and it was proposed (see for instance, Breiman (1998) argument in Olshen (2001) [37]. The application of reinforcement to classification trees contributes to classifiers that are normally equivalent of every other classifier. A especially nice thing about boosting and bagging is that simple "off-the-shelf" classifiers can be very effectively used for them (as compared to the need to tune and modify the classifiers carefully). This reality tends to balance somewhat the accusation that the expense of increased computing time comes with both boosting and bagging the enhanced output. It should also be noted that Breiman (1998) suggests that it can actually be much faster to apply boosting to CART to construct a classifier than to fit a neural net classifier. One possible downside of disrupting and merging strategies, however is the final rule of prediction could be considerably more complicated compared to that achieved using a process predictors that do not combine. It is also the case of course, that simplicity must be compromised in order to achieve greater precision.

*2.2 Decision Tree as a Base Classifier*

In knowledge discovery (KD) and data mining DTs [38] are effective, productive and common approaches for exploring broad and to define patterns, difficult databases that are valuable. This field is of great significance because it allows the modelling and removing of information from big data sets. In order making the approach more cost-effective, simpler, precise and reliable, both practitioners and theoreticians are continually seeking methods. Data analysis, machine learning, data retrieval, text pattern and mining detection are used in different disciplines., DTs are very powerful instruments.

### A. DT construction/induction

DTs[39] are predictive models which analyse data in a fashion-like tree. DTs are designed mainly for the supervised mining of data. DT is a structure-like flowchart, where it serves each internal node an

attribute test, the attribute test is denoted by each branch. result, and every leaf node represents the mark of the class. A root node has an incoming degree of zero, meaning it does not have any incoming edges. Each tuple are initially at the root node. By separating the branches of a tree, the tree obtains classification, where any split reflects a data attribute test. This division branches proceed to the last identified stage as the it hits terminal stage, where samples from one class are composed of all data tuples at one node. [39-41] provides the classification algorithm for the generation of DTs.
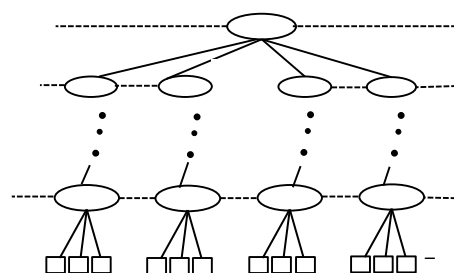


**Fig. 1**. Examples of a general (complete) balanced decision tree

Mainly, DT induction consists of two levels:

- **Tree generation:** All data tuples were initially at the root node. The splitting criterion is introduced and the best characteristic for splitting is chosen for the next step. The worth of the best attribute for splitting indicates the degree of branches belonging to the node. Partitioning continues till the sample is very small at one node and each part consists of a one class sample. Finally, full tree is formed which consist of the accuracy of training data will not be improved by any further enhancement of the leaf node.

- **Tree pruning:** By eliminating the sub-tree that represents outliers and noise, pruning reduces the size of the tree. An algorithm for checking repeated and replicated sub-trees is applied after the complete tree is created. DT is pruned whenever any sub-tree of this nature exits. Outcomes of pruning in DTs that are quicker and more accurate.

**Table I** Comparison of common algorithms for decision tree induction.

| Algorithm | Author | Splitting criteria | Pruning criteria | Implementation | Type of data | Drawbacks |
|---|---|---|---|---|---|---|
| ID3 | Quinlan (1986) | Information gain | Simple pruning | Serial | Categorical attributes | Not scalable, applicable for small datasets |
| | | | | | | Does not guarantee optimal solutions. |
| CART | Quinlan (1986) | Information gain | Simple pruning | Serial | Categorical and continuous attributes | Memory resident, applicable for small datasets. |
| | | | | | | Sorting is performed at every node |
| C4.5 | Breiman et al. (1984 | Quinlan (1993) | Gain ratio | Serial | Categorical and continuous attributes | Applicable only for small datasets. |
| | | | | | | Require skills for understanding. |
| SLIQ | Mehta et al. (1996) | Gini index | MDL principle | Serial | Categorical and continuous attributes | Attribute list is memory resident |
| | | | | | | Not applicable for parallel implementation |
| BOAT | Gehrke et al. (1999) | Various methods based on impurity | MDL principle | Serial | Categorical and continuous attributes | Require small datasets to training purpose of classifier. |
| | | | | | | |

| | | | | | | |
|---|---|---|---|---|---|---|
| | | | | | | Allow dynamic insertions and deletions which require rigorous studies. |
| **PUBLIC** | Rastogi and Shim (1998) | Entropy method | MDL principle | Parallel | | Does not concentrate on building phase |
| **SPRINT** | Shafer et al. (1996) | Gini index | MDL principle | Serial and parallel both | Categorical and continuous attributes | At each split attribute list requires to be rewritten and resorted. |
| **RainForest** | Gehrke et al. (1996) | Various methods based on impurity | MDL principle | Parallel | | Not any particular drawback as such. |

### B. Applications of DT

In several areas of real life, the DT algorithms can be implemented. DTC fields of application are summarised below:

- **Business:** In visualising probabilistic business models and even in CRM, DTCs are used (Control of client partnerships).
- **E-commerce:** DT is used to create an online catalogue that is necessary for an e-commerce website to succeed.
- **Medicine:** The relevant application fields for DT methods are disease diagnosis and medical study. In the diagnosis of multiple conditions, such as diagnosis of cardiac sounds, DT is most useful.
- **Industry:** In fault detection and non-destructive checks, DT algorithms are used.
- **Intelligent vehicles:** In designing smart vehicles, looking for road lane boundaries is a key activity. [42] addressed the framework for lane identification of intelligent vehicles that use DTs.
- **Remote sensing:** Nowadays, working with DTs for pattern recognition is emerging as a key field of interest for researchers. Categories of field cover algorithms are too suggested.

- **Web applications:** Problems of internet site failures are listed in [43] and the use of DTC for intelligent web caching is addressed in [44].

### C. DT benefits

A forum or mechanism is given by the tree structure to evaluate all possible alternatives for a decision to be made. Several advantages provided by DTs are:

- Self-described and when compacted, very quick to follow
- A variety of input data forms may be interacted with by DTCs: textual, numerical and categorical.
- Capability to process mistaken missing values or datasets
- Relative to limited computing effort, high performance
- A variety of data mining packages are accessible on a selection of platforms.
- useful for different functions for example clustering, collection of functionality, regression and grouping.

### D. Limitations of DT

DTs are reasonably easy to interpret and comprehend for humans, their creation requires

little information of the domain. But despite of many uses, there were many limitations in earlier DTs. Few of the restrictions are:

- **DT size:** The dynamics of the tree has a crucial influence on the precision of prediction. The DT size increases when the scale of the data for training grows, it's complicated and time consuming to prepare such big trees.
- **Instability:** Changing any element, with the exception of data on replication, or differing information the midway sequence may cause significant tree changes. This will cause the tree to be redrawn.
- **Scalability:** In data mining applications, large training sets of several million attributes are popular. For certain intricate datasets, earlier DTs were unable to do so.
- **Missing values:** Noisy data and missing values could not be managed by initial DTs.

### 3. Uncertain Data

Information includes noise that causes it to drift from the right, expected or real values is uncertain data. One of the distinguishing aspects of data in the era of big data is complexity or veracity of details. In length, variety, velocity and uncertainty (1/veracity), knowledge is continuously increasing. Uncertain knowledge on the internet, it is contained in abundance today, Both in their organised and unstructured origins in sensor networks within enterprises. In an enterprise dataset, for example, there might be confusion about the address of a consumer or The reading of the temperature collected by a sensor owing to ageing of the sensor.

In 2012, in its global technology outlook report, IBM called for handling unpredictable data on a scale that provides a systematic study looking into the future for three to ten years, aiming to recognise important, disruptive developments that will change the world. Analyses would certainly take into consideration the for several There are various forms of volatility of very large quantities of data on the basis of real-world data to make confident decisions about companies. The quality of

subsequent judgments may be compromised by analyses centred on unclear evidence, so it is not possible to disregard in this unknown evidence, the amount and type of inaccuracies.

The abstraction of Distributions of odds through summary statistics for example variances and means an easy way to manage data uncertainty. Data instability occurs naturally for different reasons in many applications. Three categories are briefly discussed here: staleness of results, measurement errors, and repetitive measurements.

a) *Measurement Errors:* Due to measurement errors, data collected from physical system measurements is often imprecise. As an example, by calculating the ear drum's temperature, through the infrared sensor, a tympanic (ear) thermometer measures body temperature. A standard ear thermometer has a quoted ±0.2 ?? C calibration error, that is around 6.7% of the usual operating frame, noting that the temperature of the human body varies average from 37?? C (normal) to 40?? C (severe fever). The calculation error compound can be high with other variables such as placement and technique. For instance, around 24% of the scales are greater than 0.5 C, or around 17% of the operating spectrum., it is recorded in [45]. The quantization errors introduced by the digitization method are another source of error. These errors can be treated correctly by assuming an effective error model, for example the distribution of Gaussian errors for random noise or the distribution of uniform errors for failures in quantization.

b) *Data Staleness:* Values of data are constantly changing in some applications and reported data is often stale. One example is a tracking system based on position. Only by implementing a model of uncertainty on its last recorded position [46] can the location of a mobile device be approximated. A traditional model of uncertainty involves

awareness of the device's moving speed as well as what the movement is constrained (A vehicle driving through a road network for example,) or unregulated (Like an elephant that moves on the plateau). To model such uncertainty, usually a 2D probability density function over a bounded area is described.

c) ***Repeated Measurements:*** Maybe the most prevalent cause of misunderstanding arises from repetitive observations. For example, during a day, the temperature of body of a patient multiple occasions should be taken; Wind speed could be measured per minute by an anemometer; a large number of heat sensors were mounted around the surface of the space shuttle. When talking about a patient's temperature, or wind direction, or the temperature of a certain section of the shuttle, what values do we use? Or Through consideration of the distribution given by the collected data values, would it be better to use all the information?

## 4. Related Work

There has been a surge in curiosity in the mining of unknown data. For unknown data clustering, the well-known k-means algorithm for clustering in [47] is expanded to the UK-means algorithm. As we've outlined,, pdf's, those are commonly defined via collections of sample values, typically capture data uncertainty. Mining unknown data is also As a consequence of the blast, computationally cost of information (sets of samples vs. single values). Pruning techniques have been suggested to boost the efficiency of UK-means. Min-max-dist pruning [48] and CK-means [49] are examples. In addition to studies of partition-based uncertain data clustering, density-based clustering (for example FDBSCAN [50]), regular itemset mining [51] and density-based classification [52] are other directions for uncertain data mining. Density-based classification involves the distribution of the attributes of knowledge' joint probability to be remembered. In [52], an error model is given for each data point. Every test tuple

is a data valued by a point after testing. These are somewhat distinct from our data model, since we don't include the awareness of the data attributes' joint probability distribution. Each attribute is separately treated and could have its own model and of error. In addition, the test tuples can contain uncertainty in our model, like the training tuples.

For decades, classification of the decision tree on unclear data was discussed in the missing value's form [53]. Missing values occur If, during data processing or due to mistakes in data entry, certain attribute values are not available. Solutions involve using the attribute classifier to approximate missing values containing the majority value or to infer missing values (either by precise or probabilistic values) (e.g. the organised tree attribute and the probabilistic tree attribute) [54]. Missing values in training data are treated by the use of fractional tuples in C4.5 [55] and probabilistic decision trees. Every missing value is substituted during testing Many values dependent on the training tuples, with probabilities, thereby enabling the effects of the probabilistic classification.

[56] Data uncertainty resulting from human interpretation and comprehension is modelled by fuzzy knowledge. The complexity here is the vagueness and uncertainty of hypotheses, e.g. If temperature measurements are taken into account, it is difficult to understand how hot when "hot" is the available data value. Attributes as well as class labels may be fuzzy in fuzzy classifications and are defined in fuzzy terms. In these models, a decision tree node doesn't deliver A crisp exam that deterministically determines which branch is sent down the tree training or test tuple.

[57] For decades, decision tree classification on unknown data was discussed in the form of missing values. When certain attribute values are not accessible during data collection or data entry errors, missing values appear. Solutions involve using the attribute classifier to approximate missing values with the majority value or to infer missing values (either by probabilistic or exact values) (The

organised tree attribute and the probabilistic tree attribute).

[58] The uncertainty of data has been generally categorised into existential uncertainty and uncertainty of value. There is existential uncertainty when it is unclear if an entity or a tuple of data exists. On the other hand, value uncertainty occurs when a tuple exists, but its values are not precisely defined. Unprecise query processing is one well-studied subject on value uncertainty. A probabilistic promise of its correctness is associated with the response to such a question. Always been there a rising interest in the mining of unknown data.

[59] This approach uses the Rule rule-based prediction algorithm to tackle data uncertainty. In view of the unknown data interval and probability distribution function for generating pruning and optimization, this algorithm considers new steps computed. Rules extracted using uRule show the relationship between the label attribute and class.

The rule coverage lists the number of instances that satisfy the requirement. The consistency of a rule is the fraction of instances that satisfy the condition and assign a rule's output, normalised by condition, to the class name. These procedures are used to support the uRule algorithm: uLearnOneRule(), uGrow(), splitUncertain(). The best rule for the class is created by uLearnOneRule() from an uncertain training set. It has two increasing and pruning components. The splitUncertain() function returns part of the instance that the rule protects. UGrow(initial )'s rule is that the left hand is vacant. and the right hand side includes the current class. For attribute and split point collection, probabilistic data benefit is used. If the rule covers an instance, it is removed from the dataset as the rule expands.

[60] The outcomes of data mining efficiency and quality are largely based upon data uncertainty. It needs to be correctly modelled and stored. This technique focuses on one kind of data uncertainty frequently encountered. If the exact data value is not available and the distribution of the likelihood of the data is known, then the value of the data is replaced by the estimated value. While simple and straightforward, this technique can trigger valuable loss of knowledge. To address this problem, the traditional neural networks classifier is expanded so that some data and unknown distribution of probabilities can be taken as the input. For this method, the Gaussian distribution of probability is considered. UNN would appropriately implement the classification because according to the probability distribution information, it measures the probability of P belonging to classes it is supposed to be in class that has a greater probability. Therefore, greater classification accuracy can be achieved via the uncertain neural network.

## 5. Problem Formulation

The adaptive classification technique (AdaBoost) algorithm was implemented by Schapire and Freund. In the AdaBoost algorithm, the main term is the repeated use of instead of the weights for same training data, randomly choosing new ones, since, in comparison with other classification algorithms, the Adaboost is not a complex and broad technique.

Instead of drawing a set of independent bootstrap samples from the original instances, the higher the weight, the more the instance affects the classifier learned, the greater the weight for each instance. The weight vector is modified at each trial to represent the output of the corresponding classifier, thereby raising the weight of misclassified instances. By voting, the final classifier also aggregates the learned classifiers, but each vote of the classifiers is a function of its precision.

For every case, the estimation for the training data is given by learning algorithms. and can get the optimal prediction solution if we have appropriate results. Since the Boosting and Bagging algorithm produces various predictions, many of them look very similar and precise when considering the training dataset. One is chosen as the final for that example from all available predictions, but by integrating the available classifiers in the training data, the question may be more likely to be solved.

There is therefore an increasing awareness that classifier combinations can be more efficient than single classifiers. When a mixture will achieve a more accurate and specific result of many, why depend on the best single classifier? This is basically the logic behind the principle of multiple systems for classifiers.

Boosting and Bagging is applied to the for a specified base learner, the same training data that we would use the CART because of its classification simplicity and the issue of overfitting is also the decision tree (A slight shift in the training pattern allows the established model to change tremendously). For each case the various classifiers obtained from the different basic students are using the stacking method then clustered.

### 6. Tentative Research Methodology

In our work, we are expending the adaptive classification algorithms, the Adaboost is not a complex and broad technique. It is noted that the number of base level classifiers is not greatly affected when one prepares the ensemble, and typically researchers randomly choose 3 or 7 depending on the type of applications. As decision trees are built, if the concepts of linear regression are also followed, m regression equations are created for each of the m target groups. In an algorithm, namely M5 by Quinlan, this definition is adopted.

Ensembling the classifiers is performed in the suggested approach where two separate approaches are used as classifiers such as Decision Tree for Bagging and ANNN (Artificial Neural Network). The key part of the methodology is that the bagging and boosting are independently treated using two different methods for data classification or for learning the dataset. The measures involved in performing the bagging and boosting are below:

**Step by Step processing of the proposed work:**

- **Step 1:** Pre-processing of dataset which is to be trained,

- **Step 2:** Bagging and Boosting techniques are applied separately on the dataset,
- **Step 3:** Different learning algorithms are used differently for Bagging and Boosting,
- **Step 4:** For bagging the classifier used is Decision tree and for Boosting the classifier used is the Artificial neural network,
- **Step 5:** The obtained results are then ensembled using the MDTs (Meta Decision Tree),
- **Step 6:** Final prediction is being done on the basis of the ensembled data.

### 7. Expected Outcome

Certain criteria are considered for the final review of the work that will determine the dataset classification and the utility of the discussed process. The situation is treated as the sole one datapoint in the evaluation portion parameters and parameters are evaluated in the dataset for a single case. Some of the considered evaluation parameters include Relative absolute error, Root relative squared error, mean absolute error, Root mean squared error etc. For the evaluation of the job, a confusion matrix is being built and reflects the classification of the dataset. A uncertainty matrix is a table commonly used to describe the output on a set of test data that is assumed to be valid values of a classification model. The confusion matrix itself is reasonably simple to grasp, although the words relevant can be ambiguous.

### 8. Conclusion

In solving many pattern recognition issues, DTs provide a great deal of promise and are well known for their visualisation of performance data. The key feature of DTC is to offer versatility, i.e. the ability to use different features subsets and decision rules at different classification points, and the ability for maintaining a balance between precision of classification and efficiency of time/space. In the current work the weighted learning technique is used for the dataset, the work processes over the uncertain numerical dataset. So as to ensemble the classifiers bagging and boosting technique is used

over the uncertain numerical data. The work formatted in the paper also presents the fundamentals of the uncertain dataset, bagging and boosting. Also the work provides the fundamentals about the decision as base classifier, where the limitations, advantages, benefits of the consideration of the DT as base classifiers are depicted. The research conducted can be directly picked for the practical evaluation of the research methodology represented in the work. In the future part the work have the hint about the research direction which is needed to be evaluated further to be picked as complete and appropriate research methodology and the same is needed to be verified after experimental evaluation of the research over the defined uncertain dataset. Various parameters which are to be used for the work evaluation are defined in the work to have the basics about the evaluation of the research proposal.

## References

[1] Aggarwal, Charu C., and S. Yu Philip. "A survey of uncertain data algorithms and applications." *IEEE Transactions on knowledge and data engineering* 21.5 (2008): 609-623.

[2] Aggarwal, Charu C., and CK Reddy Data Clustering. "Algorithms and Applications." (2014).

[3] Zhang, Xianchao, Han Liu, and Xiaotong Zhang. "Novel density-based and hierarchical density-based clustering algorithms for uncertain data." *Neural Networks* 93 (2017): 240-255.

[4] Liu, Han, et al. "Self-adapted mixture distance measure for clustering uncertain data." *Knowledge-Based Systems* 126 (2017): 33-47.

[5] Trajcevski, Goce, et al. "Managing uncertainty in moving objects databases." *ACM Transactions on Database Systems (TODS)* 29.3 (2004): 463-507.

[6] Deshpande, Amol, et al. "Model-based approximate querying in sensor networks." *The VLDB journal* 14.4 (2005): 417-443.

[7] Liu, Xuejun, et al. "A tractable probabilistic model for Affymetrix probe-level analysis across multiple chips." *Bioinformatics* 21.18 (2005): 3637-3644.

[8] Sarma, Anish Das, et al. "Representing uncertain data: models, properties, and algorithms." *The VLDB Journal* 18.5 (2009): 989.

[9] Jampani, Ravi, et al. "MCDB: a monte carlo approach to managing uncertain data." *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*. 2008.

[10] Zhang, Wenjie, et al. "Managing uncertain data: Probabilistic approaches." *2008 The Ninth International Conference on Web-Age Information Management*. IEEE, 2008.

[11] Wang, Yijie, et al. "A survey of queries over uncertain data." *Knowledge and information systems* 37.3 (2013): 485-530.

[12] Dallachiesa, Michele, Themis Palpanas, and Ihab F. Ilyas. "Top-k nearest neighbor search in uncertain data series." *Proceedings of the VLDB Endowment* 8.1 (2014): 13-24.

[13] Ren, Jiangtao, et al. "Naive bayes classification of uncertain data." *2009 Ninth IEEE International Conference on Data Mining*. IEEE, 2009.

[14] Qin, Biao, et al. "A novel Bayesian classification for uncertain data." *Knowledge-Based Systems* 24.8 (2011): 1151-1158.

[15] Qin, Biao, Yuni Xia, and Fang Li. "DTU: a decision tree for uncertain data." *Pacific-Asia conference on knowledge discovery and*

*data mining*. Springer, Berlin, Heidelberg, 2009.

[16] Tsang, Smith, et al. "Decision trees for uncertain data." *IEEE transactions on knowledge and data engineering* 23.1 (2009): 64-78.

[17] Angiulli, Fabrizio, and Fabio Fassetti. "Nearest neighbor-based classification of uncertain data." *ACM Transactions on Knowledge Discovery from Data (TKDD)* 7.1 (2013): 1-35.

[18] Tavakkol, Behnam, Myong Kee Jeong, and Susan L. Albin. "Object-to-group probabilistic distance measure for uncertain data classification." *Neurocomputing* 230 (2017): 143-151.

[19] Bi, Jinbo, and Tong Zhang. "Support vector classification with input data uncertainty." *Advances in neural information processing systems*. 2005.

[20] Yang, Jianqiang, and Steve Gunn. "Exploiting uncertain data in support vector classification." *International Conference on Knowledge-Based and Intelligent Information and Engineering Systems*. Springer, Berlin, Heidelberg, 2007.

[21] Qin, Biao, et al. "A rule-based classification algorithm for uncertain data." *2009 IEEE 25th International Conference on Data Engineering*. IEEE, 2009.

[22] Gao, Chuancong, and Jianyong Wang. "Direct mining of discriminative patterns for classifying uncertain data." *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2010.

[23] Pedagadi, Sateesh, James Orwell, Sergio Velastin, and Boghos Boghossian. "Local fisher discriminant analysis for pedestrian re-identification." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 3318-3325. 2013.

[24] Ge, Jiaqi, Yuni Xia, and Chandima Nadungodage. "UNN: a neural network for uncertain data classification." *Pacific-Asia Conference on Knowledge Discovery and Data Mining*. Springer, Berlin, Heidelberg, 2010.

[25] Cao, Keyan, et al. "An algorithm for classification over uncertain data based on extreme learning machine." *Neurocomputing* 174 (2016): 194-202.

[26] Kotsiantis, Sotiris B., I. Zaharakis, and P. Pintelas. "Supervised machine learning: A review of classification techniques." *Emerging artificial intelligence applications in computer engineering* 160.1 (2007): 3-24.

[27] Liu, Han, Xianchao Zhang, and Xiaotong Zhang. "Possible world based consistency learning model for clustering and classifying uncertain data." *Neural Networks* 102 (2018): 48-66.

[28] Dietterich, Thomas G. "Ensemble methods in machine learning." *International workshop on multiple classifier systems*. Springer, Berlin, Heidelberg, 2000.

[29] Breiman, Leo. "Bagging predictors." *Machine learning* 24.2 (1996): 123-140.

[30] Freund, Yoav, and Robert E. Schapire. "Experiments with a new boosting algorithm." *icml*. Vol. 96. 1996.

[31] Breiman, Leo. "Bagging predictors." *Machine learning* 24.2 (1996): 123-140..

[32] Schapire, Robert E. "The strength of weak learnability." *Machine learning* 5.2 (1990): 197-227.

[33] Freund, Yoav. "Boosting a weak learning algorithm by majority." *Information and computation* 121.2 (1995): 256-285.

[34] Freund, Yoav, and Robert E. Schapire. "Experiments with a new boosting algorithm." *icml*. Vol. 96. 1996.

[35] Breiman, Leo. "Arcing classifier (with discussion and a rejoinder by the author)." *The annals of statistics* 26.3 (1998): 801-849.

[36] Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.

[37] Olshen, Richard, and Leo Breiman. "A conversation with Leo Breiman." *Statistical Science* (2001): 184-198.

[38] Rokach, Lior, and Oded Z. Maimon. *Data mining with decision trees: theory and applications*. Vol. 69. World scientific, 2008.

[39] Payne, Harold J., and William S. Meisel. "An algorithm for constructing optimal binary decision trees." *IEEE Transactions on Computers* 9 (1977): 905-916.

[40] Barros, Rodrigo Coelho, et al. "A survey of evolutionary algorithms for decision-tree induction." *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)* 42.3 (2011): 291-312.

[41] Srivastava, Anurag, et al. "Parallel formulations of decision-tree classification algorithms." *High Performance Data Mining*. Springer, Boston, MA, 1999. 237-261.

[42] Gonzalez, Juan Pablo, and Umit Ozguner. "Lane detection using histogram-based segmentation and decision trees." *ITSC2000. 2000 IEEE Intelligent Transportation Systems. Proceedings (Cat. No. 00TH8493)*. IEEE, 2000.

[43] Chen, Mike, et al. "Failure diagnosis using decision trees." *International Conference on Autonomic Computing, 2004. Proceedings.*. IEEE, 2004.

[44] Bonchi, Francesco, et al. "Data mining for intelligent web caching." *Proceedings International Conference on Information Technology: Coding and Computing*. IEEE, 2001.

[45] FREED, GARY L., and J. KENNARD FRALEY. "25%" Error rate" in ear temperature sensing device." *Pediatrics* 87.3 (1991): 414-415.

[46] Wolfson, Ouri, and Huabei Yin. "Accuracy and resource consumption in tracking and location prediction." *International Symposium on Spatial and Temporal Databases*. Springer, Berlin, Heidelberg, 2003.

[47] Chau, Michael, et al. "Uncertain data mining: An example in clustering location data." *Pacific-Asia conference on knowledge discovery and data mining*. Springer, Berlin, Heidelberg, 2006.

[48] Ngai, Wang Kay, et al. "Efficient clustering of uncertain data." *Sixth International Conference on Data Mining (ICDM'06)*. IEEE, 2006.

[49] Lee, Sau Dan, Ben Kao, and Reynold Cheng. "Reducing UK-means to K-means." *Seventh IEEE International Conference on Data Mining Workshops (ICDMW 2007)*. IEEE, 2007.

[50] Kriegel, Hans-Peter, and Martin Pfeifle. "Density-based clustering of uncertain data." *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*. 2005.

[51] Chui, Chun-Kit, Ben Kao, and Edward Hung. "Mining frequent itemsets from uncertain data." *Pacific-Asia Conference on knowledge*

*discovery and data mining*. Springer, Berlin, Heidelberg, 2007.

[52] Aggarwal, Charu C. "On density based transforms for uncertain data mining." *2007 IEEE 23rd International Conference on Data Engineering*. IEEE, 2007.

[53] Quinlan, J. Ross. "Induction of decision trees." *Machine learning* 1.1 (1986): 81-106.

[54] Hawarah, Lamis, Ana Simonet, and Michel Simonet. "A probabilistic approach to classify incomplete objects using decision trees." *International Conference on Database and Expert Systems Applications*. Springer, Berlin, Heidelberg, 2004.

[55] Quinlan, J. Ross. *C4. 5: programs for machine learning*. Elsevier, 2014.

[56] Vazirgiannis, Michalis, and Maria Halkidi. "Uncertainty handling in the data mining process with fuzzy logic." *Ninth IEEE International Conference on Fuzzy Systems.*

*FUZZ-IEEE 2000 (Cat. No. 00CH37063)*. Vol. 1. IEEE, 2000.

[57] Ngai, Wang Kay, et al. "Efficient clustering of uncertain data." *Sixth International Conference on Data Mining (ICDM'06)*. IEEE, 2006.

[58] Dalvi, Nilesh, and Dan Suciu. "Efficient query evaluation on probabilistic databases." *The VLDB Journal* 16.4 (2007): 523-544.

[59] Kesavaraj, Gopalan, and Sreekumar Sukumaran. "A study on classification techniques in data mining." *2013 Fourth International Conference on Computing, Communications and Networking Technologies (ICCCNT)*. IEEE, 2013.

[60] Phadatare, Manasi M., and Sushma S. Nandgaonkar. "Uncertain data mining using decision tree and bagging technique." *Int. J. Comput. Sci. Inf. Technol* 5 (2014): 3069-3073