



www.ijarr.org

Underwater Sonar Signals Recognition by Incremental Data Stream Mining with Conflict Analysis

Nagam Aanjaneyulu, Shaik Guntur Mahabub Subhani,

Chevula Rekha, Aremandla Sai Pujitha

Associate Professor^{1,2}, Assistant Professor^{3,4}

anji.amrexamcell@gmail.com¹, subhanimehandi@gmail.com²,
rekhavenkat16@gmail.com³, sapuja502@gmail.com⁴

Department of CSE, A.M. Reddy Memorial College of Engineering and Technology,
Petlurivaripalem, Narasaraopet, Andhra Pradesh

1. Introduction

University of Macau, Taipa Campus, Department of Computer and Information Science

Location: 2Computer Science and Engineering Department, Cambridge Institute of Technology, Ranchi
835103, India

Sydney, NSW2052, Australia; Level 3, School of Computer Science and Engineering, University of New
South Wales

4Electronic Engineering Department, Beijing University of Technology, Beijing 100022, China Address
any communication to Simon Fong at ccfong@umac.mo.

The manuscript was received on March 3, 2014; accepted on April 15, 2014; and published on May 28,
2014.

Professorial Editor: Sabah Mohammed

2014 Simon Fong et al. Copyright. This is an open-access article shared under the terms of the Creative Commons Attribution License, which allows for free use, distribution, and reproduction in any form, provided the original work is properly credited.

Recognizing sonar sounds is a crucial step in finding large items submerged in the ocean. Rather than relying on visual cues, the military often uses sonar signals to guide them through the depths of the ocean and/or detect the presence of enemy submarines. In particular, data mining's classification technique has

been used to sonar signal detection to determine the nature of reflected surfaces. By training a classification model using the whole dataset in batches, classification algorithms in the conventional data mining technique provide reasonable accuracy. It's well known that data streams are continuously acquired from sonar waves.

The prior classification techniques may not be applicable to incremental classifier learning, notwithstanding their efficacy in conventional batch training. To meet the need for fast speed, data preparation time must be minimized despite the infinitely large data streams that might result from sonar signals. To avoid having to learn everything about a dataset all at once, this study introduces a new approach to data mining that is well-suited to the gradual elimination of noisy data through quick conflict analysis of the data stream. Through rigorous simulation studies, we find strong evidence for the methodology's success.

1. Set the Scene

The acronym "sonar" refers to a kind of sound propagation technology often employed in underwater environments for navigating and measuring distance.

underwater navigating, talking, and/or finding something. Current methods in this area have been recently reviewed in [1]. The detection/classification problem of sonar sounds was singled out as one of the field's most difficult issues.

The success of underwater item detection relies heavily on selecting an appropriate classification model for sonar signals identification. Numerous uses for underwater sensor networks are mentioned in [2], including ocean sampling networks, environment monitoring, offshore explorations, disaster prevention, aided navigation, and mine reconnaissance. Easy installation, no need for wires, and minimal disruption to shipping lanes are just a few of the benefits of underwater sensor networks. Unfortunately, noise and interference are common problems for sonar signals that travel underwater, particularly over long distances. In particular, data mining classification approaches have been extensively used in sonar signal identification to identify the surface of the target item from which the sonar waves were reflected [3-5].

By utilizing the whole dataset to induce a classification model, classification algorithms in the classic data mining technique may be able to attain high levels of accuracy. Unfortunately, induction is often performed and repeated in batches, which suggests a reduction in accuracy between model updates [6]. In addition, when new data accumulates, the dataset as a whole will likely grow in size, which might lengthen the time it takes to do updates.

Sonar signals, like any other data stream, are constantly being sent and received. Although the trained model generated by batchmode classification methods is accurate, it may not be appropriate for use in streaming settings

detection via sonar For real-time sensing and reconnaissance, it is vital to make the data processing time extremely low since sonar signal data streams might potentially add up to infinity.

In this research, we provide a new approach to data stream mining that makes use of rapid conflict analysis to remove noisy data in increments from the stream-based training dataset. The abbreviation iDSM-CA stands for "incremental data stream mining with conflict analysis." The approach benefits from gradually developing a classification model from stream data. For the purpose of demonstrating the effectiveness of the suggested technique, simulation tests are conducted, focusing on the problem of eliminating noisy data from the sonar data while they are streaming.

Here is how the remainder of the paper is structured. In Section 2, we look at three well-known computational methods for eliminating unwanted background noise in data sets used for training.

The third section explains our new data stream mining strategy and the "conflict analysis" process we utilize to get rid of incorrectly categorized occurrences. In Section 4, we demonstrate the effectiveness of the stream mining method using a battery of sonar recognition trials. The paper is finished with Section 5.

2.Related Research

A variety of methods have been used by researchers in an effort to identify and eliminate noisy data, often known as random chaos in the training dataset. Essentially, these methods seek for occurrences in the data that cause the training model unnecessary confusion, hence improving classification accuracy. As a rule, they probe for anomalies in the data and examine how they influence the accuracy of classifications. Statistics-based techniques, similarity-based approaches, and classification-based strategies make up the bulk of the available options.

In this section, we will discuss noise detection techniques based on statistics, which we will refer to as 2.1. Data with very unusual values, known as outliers, are assumed to be random fluctuations when using this technique. Methods of detection suggested in the literature vary from simple normalcy checks to more involved procedures like looking for outlying values over a threshold.

Outlier detection techniques are surveyed in detail in [7, 8] for the purpose of locating sources of noise during preprocessing. The authors of [9] used a novel method of outlier identification in which they double-checked the predicted behavior based on the dataset. The information represented by a point is judged abnormal if it is sparse in a lower low-dimensional projection.

The projections are determined through trial and error or, at best, with the use of heuristics.

A comparable strategy is described in [10], which constructs a height-balanced tree using clustering characteristics on both the nodes that are not leaves and the leaves themselves. Outlier leaf nodes with an abnormally low density are removed.

To spot noise, similarity-based techniques are discussed in Section 2.2. This class of procedures often calls for some kind of reference against which data may be compared to establish a degree of similarity or dissimilarity.

First, the researchers in [11] partitioned the data into several subsets, and then they looked for the subset that caused the effect. training dataset elements whose removal would result in the largest decrease in remaining dissimilarity. Any function that produces a little difference in value when comparing two otherwise identical items may be used as the dissimilarity function, and

a significant value between contrasting factors, such as dissimilarity.

However, the authors concede that settling on a single dissimilarity function is no easy task. Applying a hyperclique-based data cleaner, Xiong et al. [12] developed theHCleaner method. A hyperclique pattern is characterized by a high degree of resemblance between any two items that are close together. As noise, the HCleaner treats instances that do not fit any hyperclique pattern.

The k-nearest neighbors (k-NN) approach was used by a different group of researchers [13], which effectively compares test data with nearby data to determine whether they are outliers. Disparate pieces of information are identified and eliminated if they do not closely resemble their closest neighbors.

Wilson's editing technique, a series of criteria that automatically pick the data to be deleted, was developed by the authors after they analyzed patterns of activity among the data.

Methods for Detecting Noise Using Classification 2.3

Methods based on classification use one or more pre-built classifiers as references to determine which data instances have been erroneously categorised and should be purged.

Mislabeled examples were discovered using an n-fold cross validation strategy by the authors of [14]. For this method, the dataset is split into a number of smaller subsets, say n. In order to classify the examples in the excluded subset, m classifiers are trained on the instances in the other n - 1 groups. If a classifier makes a mistake, it marks the occurrence as misclassified. The filtering procedure may be done using either majority vote or a consensus method. For the exclusion of anomalous data, another group of researchers [15] has developed a powerful decision tree approach. This technique involves creating a pruning tree to categorize the training data.

Erroneous examples that the trimmed tree labels are deleted from the dataset. Iterations of these steps are taken until the resulting trimmed tree accurately labels every occurrence in the training set. Ingeniously, the researchers in the work described in [16] employed a genetic algorithm (GA) to generate a sample of potentially noisy examples and then choose a prototype from among them. The GA searches for erroneously labeled examples using a generic classifier that was trained in advance as its fitness function.

3. The Data Stream Mining Model We Suggest

All of the aforementioned methods need a whole dataset before deciding which instances should be removed, making them best suited for batch-mode data preparation. Compared to the methods described in Section 2, this proposal for data pretreatment and model learning stands out.

This approach of preprocessing has often been treated as an independent procedure before model learning.

At least once, the whole dataset is examined to identify outliers that should be eliminated from the dataset because they provide a risk of incorrect classification in the future. An example of a filtered training set is

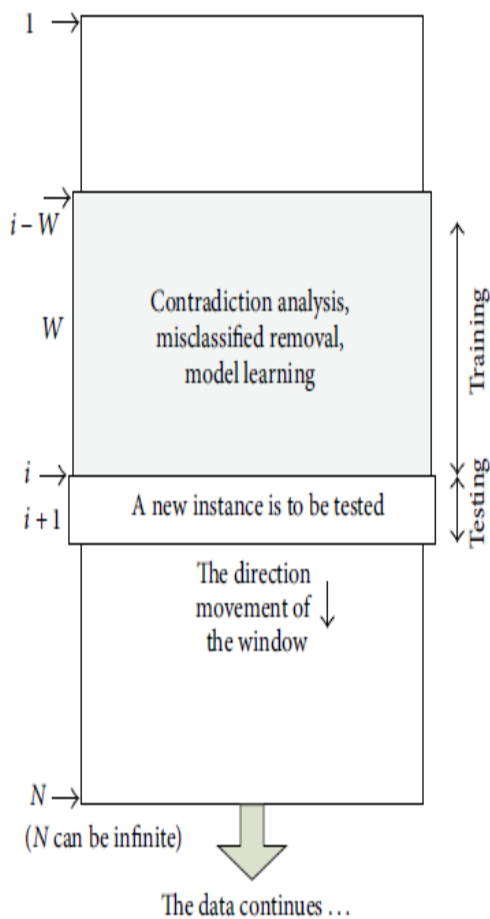


FIGURE 1: Illustration of how iDSM-CA works.

subsequently fed back into the learning process with the hope that it would help to create a more peaceful environment for learning.

On the other hand, iDSM-CA is a part of the gradual

the learning process, with each stage (noise detection, elimination of incorrectly categorized data, and learning) occurring simultaneously. In this two-pronged strategy, the incoming data is first preprocessed and trained, then put to the test. Figure 1 depicts a window of size W moving with the data flow. Conflict analysis (for noise detection), misclassified data removal, and training are all performed on the data during the time frame (model building).

After the model has been properly trained, it is put to the test on new data.

By collecting data at regular intervals, we can determine the performance level at each stage of the process and use that data to get the average performance level at the conclusion of the operation.

Model preprocessing and incremental learning 3.1 workflow. Figure 2 displays the iDSM-comprehensive CA's process flow. The window used for preprocessing and training moves with the data stream from the beginning and is not likely to use all of the data that is available. In data mining, this is referred described

as a "anytime" technique since the model may be used at any time without having to wait for all of the training data (for testing). Whenever there is a fresh influx of information, the window will gradually cover the new instances while fading away the old ones. Every time the analysis resumes, a new version of the model is built from scratch in real time.

Benefiting from the benefits of deleting misclassified instances, this method eliminates the requirement to assume the dataset is fixed and limited. New information is added to the training dataset inside window W at regular intervals, allowing the model to grow in sophistication.

3

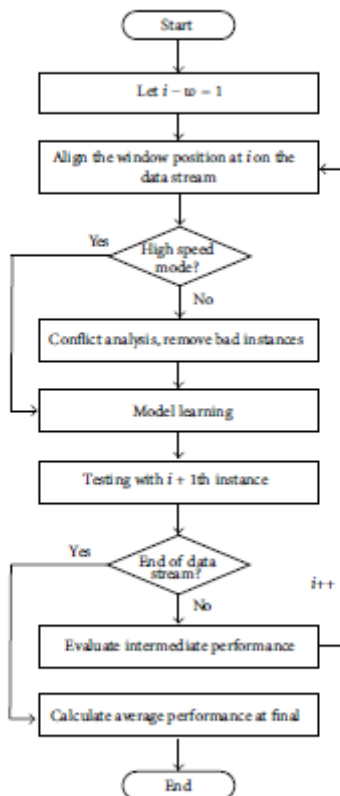


FIGURE 2: Workflow of the incremental learning method.

new information being included into W . One other advantage of the suggested method is that the data preserved by the rolling window W may be added together. Through the compilation of data sets

The properties of the data are recorded delicately from a long-run global viewpoint based on the contradiction analysis carried out inside each frame of the window as it rolls ahead.

By using this kind of worldwide data, contradiction analysis may be able to recognize noisy data with greater precision. That is to say, the function's sensitivity to noise improves as it gains experience (through use of accumulated data). Obviously, noise is something that can only be identified in relation to something else.

Dissecting the Conflict, Section 3.2. As attribute values and class labels are interdependent, a modified pair-wise based classifier (PWC) is utilized for contradiction analysis. When an instance comes, PWC is triggered like an instance-based classifier or lazy classifier and trains the classifier in increments of no more than one round.

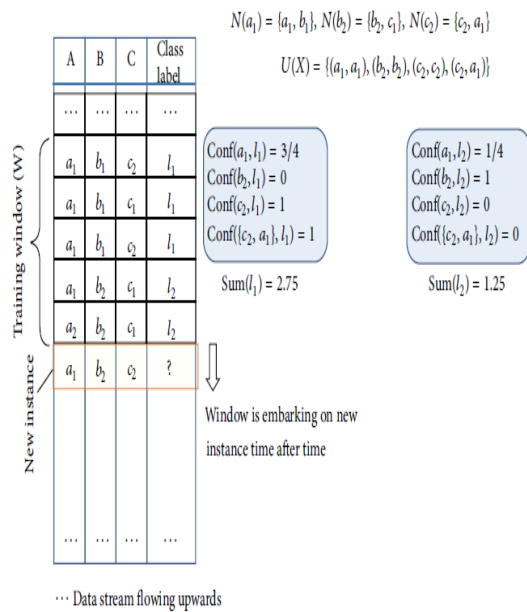


FIGURE 3: Illustration of a dynamic rolling window and conflict analysis.

Along with its high processing speed, which is necessary for lightweight preprocessing, PWC offers other advantages over other approaches. Among the benefits is the ease of just

When calculating the supports and confidence values for estimating which target label an instance should be classified into, no persistent tree structure or trained model is required beyond small registers for statistics, and the number of samples (reference) needed for noise detection can scale elastically to any amount (W). Example of a PWC with weights, based on [17]. Figure 3.

With each incremental model update, the sliding window comprises W possible training samples across three characteristics (A, B, and C) and one target class, where I is the current position. A new instance X with the values a 1, b 2, and c 2 is created at position I + 1, which is right before the previous end of the window. The neighborhood sets for each attribute value of X are shown in the upper-right part of the picture, assuming $k = W/2$, which rounds up to 2. For a 1, the greatest two $\text{conf}(a_1, a_1)$ values are 1 and $\text{conf}(a_1, b_1)$ values are 0.75, therefore $N(a_1, b_1) = a_1, b_1$. You can see the final U(X) set down below. For example, the pair (a 1, a 1) is formed when a 1 and a 1 are both present in U(X). Although b 1 is smaller than N(a 1), it does not belong in X and must be left out of U(x).

U(X) associated with c 1 also includes c 2, which is part of N(c 2), and a 1, which is part of N(c 2), so PWC checks the confidence levels for each member of U(X) against l 1 and l 2 to determine which is more likely. For (a 1, l 1), for example, we verify for first-order dependence by computing a $\text{conf}(a_1, l_1) = \text{support}(a_1, l_1) / \text{support}(a_1) = 3/4 = 0.75$. For the other pair (c 2, a 1), we look at the second-order dependence by computing $\text{conf}(c_2, a_1, l_1) = \text{support}(a_1, c_2, l_1) / \text{support}(a_1, c_2) = 2/2 = 1$. By adding up the values of certainty in each group, we get $\text{Sum}(l_1) = 2.75$ and $\text{Sum}(l_2) = 1.25$; hence, the

The new instance must be a member of the l 1 class. In this method, we can verify whether there is a disagreement by seeing if the computed membership

of the class coincides with the new instance's labeled class. If the new instance matches the PWC computation, then there is considered to be no conflict, and the window advances by one row, omitting the last row, and adding the new instance to the dataset. Whenever a new instance is created, its class label is compared to the result of the calculated class label. If there is a discrepancy, the new instance is marked as a conflict and removed.

One thing we changed was how we handle neighbor sets. Only the most recent (and hence most reliable) confidence values for couples inside the current window (the "Local Sets") are kept in the neighbor sets at any given time. Each time the window shifts to a new location and a new instance is added, the previous information in the Local Sets is discarded and replaced with freshly calculated data. We use a similar buffer, termed Global Sets, which instead of being replaced by freshly calculated confidence values for each pair in the window frame, just accumulates them. Such algorithms may naturally be used to implement conflict analysis. The use of PWC allows for the minimum amount of data and computation to be collected and processed, resulting in a streamlined process.

Iterate and test

The experiment aims to test how well the suggested iDSM-CA technique performs for recognizing sonar signals from below the surface. The recognition of sonar signal data in a data stream mining environment is of special interest, and we want to test how well iDSM-CA performs in contrast to more conventional data mining techniques in this context.

Using iDSM-CA, we tested six different classification algorithms for sonar recognition. Traditional batch-based learning may make use of two different algorithms.

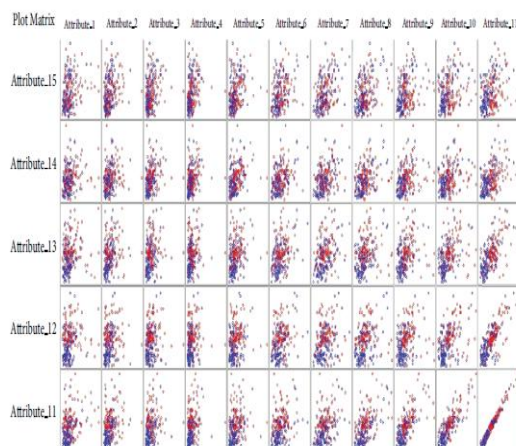


FIGURE 4: Matrix visualization of data over multiple attributes of the sonar dataset.

are a kind of neural network that uses back-propagation of its (SVM). When it comes to iDSM-CA, the best option is to use instance-based classifiers, which gradually improve with each new set of inputs.

flood in and include locally weighted learning, K-nearest neighbors classifiers, and decision tables (LWL). For the sake of curiosity, we also include an incremental version of the method, updateable naïve Bayesian (NBup), which is a modification of the original naïve Bayesian. All six methods may be used in either batch learning or incremental learning mode. Using conflict analysis to ensure iterative training

and updating, the incremental learning mode treats the data as a data stream. A trained model is put to the test on the next incoming data instance as the window advances. This method allows for the continuous collection of data on performance, beginning with the first measurement and continuing until the last. Typically, while training a model, we utilize all of the available data, and then we use 10-fold validation to assess how well the model performed.

Weka, an open-source Java platform developed at the University of Waikato, is a well-known software tool for experimenting with machine learning. Weka has comprehensive documentation for all of the aforementioned algorithms in its repository of documentation files (which may be downloaded by the general public at <http://www.cs.waikato.ac.nz/ml/weka/>). Therefore, we will not go over the same ground again. Lenovo laptop has Intel Pentium Dual-Core T3200 2GHz CPU, 8GB RAM, and 64-bit Windows 7 operating system.

The "connectionist bench (sonar, mines vs rocks) data set," or Sonar for short, is a prominent test dataset for evaluating classification systems.

You may get the dataset from UC Irvine's Machine Learning Repository (<http://archive.ics.uci.edu/ml/datasets>). To classify sonar waves, Gorman and Sejnowski [18] were the first to use this dataset in an experiment with various neural network parameters. This does the same thing.

sonar waves reflected from a metal cylinder and those reflected from a roughly cylinder-shaped rock. A graphical representation of the distribution of data points inside

Figure 4 shows the distribution of the population according to the two classifications (mine or rock) in blue and red. Figure 4 depicts a ship using sonar to distinguish between underwater mines and rocks. There are several overlapping sets of data in every attribute combination, indicating that the underlying mapping pattern is very nonlinear. This suggests a challenging categorization task, one where achieving a high degree of accuracy is problematic.

There are two distinct sorts of patterns within this sonar sample. One hundred and eleven of them are gathered experimentally by sending out sonar sounds and seeing how they reflect off a metal cylinder at various angles and under varying circumstances. Signals reflected off rocks in a similar environment account for the remaining 97 patterns. Typically, an auditory chirp of increasing frequency is used to send the sonar signal. Between the rock's 180 degrees and the metal cylinder's 90 degrees, there's a large range of transmission angles available. The qualities or characteristics of the reflected signal are represented by a vector of 60 decimal values [0, 1].

Each parameter is a measure of the total energy during a certain time period that may be assigned to a specific frequency range.

Each record's target class is binary; if the obstruction is a rock, it will be described as such, and if the item that reflected the signal is a metal cylinder, it will be described as such. Although the precise values of the angles are not stored, the characteristics in the dataset are sorted in ascending order of angles.

Figure 5 shows a visualization of the attribute values, and it is clear that the characteristics representing the various acoustic beam angles may take on a broad variety of values.

Spreads of the signal strengths are greatest at the center angles.

Perhaps the sensors collect the reflected acoustic signals by first scanning a large number of surfaces. This

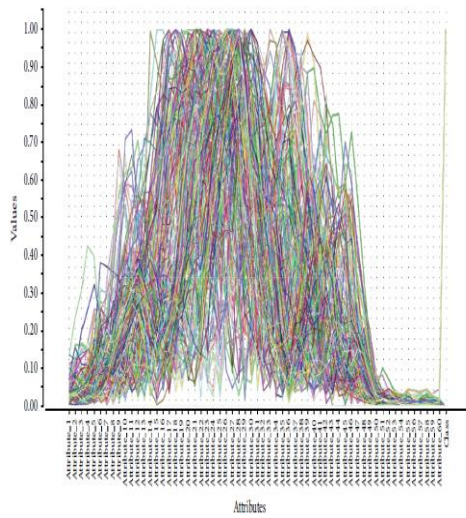


FIGURE 5: Visualization of the sensor values on a parallel coordinate plot.

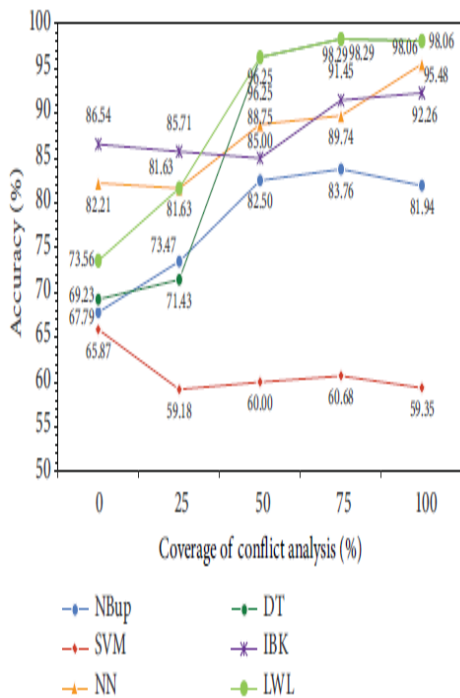


FIGURE 6: Accuracy of sonar classification by various algorithms under iDSM-CA.

indicates intricate and highly nonlinear connections between qualities and classes. Therefore, it would be difficult for a classifier to get high recognition.

accuracy.

We utilize this information to evaluate the training duration, classification model prediction accuracy, and ROC indices in our experiment.

In the case of batch learning, the length of time it takes to construct a fully operational classification model from scratch using all available data serves as a proxy for the duration of the training process. The model learning time in incremental learning is equal to the mean time per each step sliding from the beginning to the end of the data stream, taking into account the time needed to analyze the data, resolve any conflicts, and train the model. The proportion of properly labeled instances relative to the total number of examples is what we call the accuracy.

The receiver operating characteristic (ROC) curve displays the upper and lower bounds of a test's ability (the discrimination power in classification) to separate between alternative states of target objects across the entire spectrum of operating conditions, providing a single, consistent measure of accuracy between 0 and 1. As soon as the ROC value drops below 0.5, the model's predictions are no better than those obtained by chance.

Prior to beginning the experiment, the dataset is put through a battery of six algorithms in the calibration phase to determine the best value for W . Calibration might be repeated at set intervals or if the performance of the incremental learning lowers to fine tune the window size, although in reality only a small sample would be utilized. As the experiment progresses, we increase the window size from 49 to 80 to 117 to 155. The different window widths are often referred to as 25%, 50%, 75%, and 100% of the whole dataset for ease of reference. If you use $W = 0\%$ for incremental learning, you'll get the same results as if you used a complete batch of data.

Figures 6 and 7 show the classification accuracy and model induction time in seconds achieved by the iDSM-CA, while Figure 8 displays the ROC index achieved by the aforementioned approach. The proportion of correct classifications provided by the classifier is a direct measure of how well it does at separating the two classes of rocks and metals. The efficiency with which a model may be induced is an indication of how well the techniques being used scale to the data stream mining setting. The ideal methods for real-time model update and refreshment would take next to no time at all.

If the window size W is set to 0, then no noisy instances will be weeded out of the data set as a result of a conflict analysis.

As the methods and the size of the sliding window change, so do the classification models' accuracies. Figure 6 shows that the classic classifiers are not as effective as the incremental classifiers.

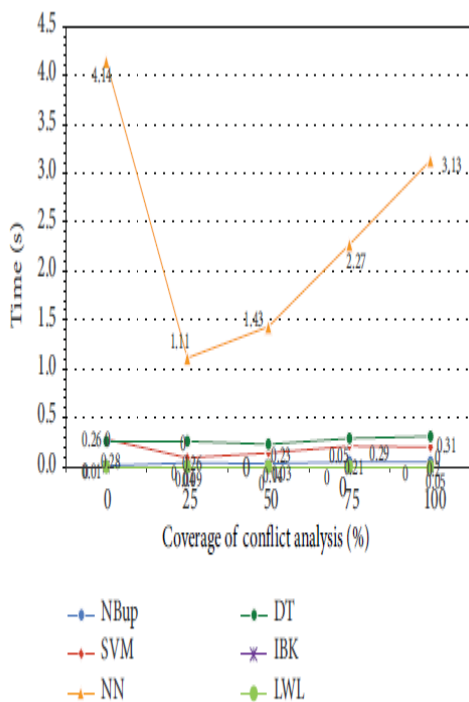


FIGURE 7: Model induction time of sonar classification by various algorithms under iDSM-CA.

other than NN. Accuracy improvements of several W% are possible for NN in general. When compared to other incremental algorithms, such as LWL and IBK, it may get better results. The accuracy of NN is better than

that of LWL for low values of W, but degrades with increasing W. Stream mining is an example where a small value for W is preferred. In a short window, IBK has the greatest accuracy, followed by LWL, NBup, and DT. When $W = 0$, there is no preprocessing to remove noise, hence IBK and NN both do well (82%-90%). Rates of accuracy for the remaining algorithms range from 73% to 66%. With the exception of SVM and IBK, all algorithms benefit from conflict analysis's noise reduction. As far as I can tell, the most likely candidates to take on the improvement are LWL, NBup, and DT.

If, on the other hand, $W = 100\%$, we treat the whole dataset as if it were a conflict analysis, we get something like to a thorough conflict analysis. Approximately 98% accuracy is achieved by LWL and DT, with NN and IBK coming in a close second and third, respectively. The benefits of comprehensive noise reduction were clearly not exploited by NBup and SVM.

Here, speed is interpreted as model induction time, and it is mostly determined by the amount of time spent on each model update, making it an essential requirement for data stream mining. In incremental learning, it is anticipated that model updates will occur at a constant rate; that is, whenever a new data instance is received, the model will be updated once to account for the new data. Figure 7 displays the typical duration required by each technique to complete a model update. Except for NN, all of the other algorithms can complete a model update in less than 0.4 seconds.

When data is not preprocessed for noise reduction ($W = 0$), NN consumes the most time. Figure 6 demonstrates that as W grows, the time required by NN steadily increases from its minimum at $W = 0$. One last performance indicator is the ROC index, which suggests the reliability of the classification

system. Figure 8 shows that, on average, all algorithms can achieve a high ROC level, with the exception of support vector machine (SVM), which fares very badly.

7

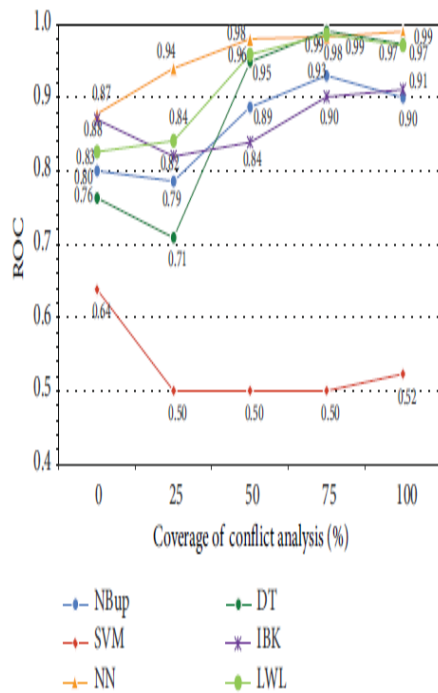


FIGURE 8: ROC of sonar classification by various algorithms under iDSM-CA.

Whenever the iDSM-CA is relevant. In terms of ROC, NN achieves the highest value, while LWL is a promising incremental algorithm that benefits from using conflict analysis.

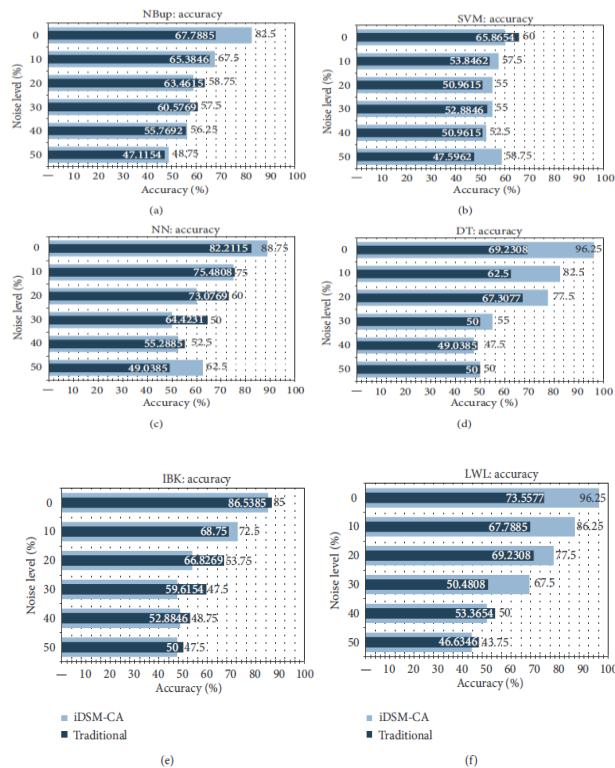
The second section of the experiment compares the performance of standard and iDSM-CA classification modes using six widely used algorithms to determine the effect of noise on sonar signal detection. In this experiment, we introduce spurious, incorrect target class values in order to amplify the background noise. An arbitrary proportion of target class values in the dataset may be fabricated with incorrect values to regulate the noise level. This would lead to a decrease in overall algorithm performance. Underwater sonar is very susceptible to significant noise in this circumstance due to factors including background inference, defective sensors, or a target of detection that is almost out of range.

In this study, we limit our manipulation of noise to a range of 50% or less.

When the noise level goes above 50%, the classifiers begin to incorrectly consider the erroneous values in the noise as true values due to majority rule, resulting in an assessment that is nonsensical.

Figure 9(a) through (f) depict the relative accuracy of both learning modes of the classifiers under consideration while the noise is amplified. Figures 10(a)–10(e) show time expenditures, whereas Figures 11(a)–11(e) show return on investment (f). In this instance, we are maintaining a constant window size of 80 for incremental learning.

It is clear from comparing the results of the batch and incremental learning modes in Figure 9(a)–(f) that, over a range of noise levels, the incremental learning mode algorithms perform better overall. As the noise level is low, NBup in incremental mode performs well; however, when the noise level rises, NBup in incremental mode lags behind that of batch mode. Both SVM and NN show the same behavior. There is a little but noticeable increase when using the incremental learning mode, with the exception of the region around the 50% noise level (the oblivion state), where it is impossible to discern between noise and real cases. Contrarily, incremental learning mode algorithms, such as DT, IBK, and LWL, demonstrate superior accuracy. To focus on LWL specifically, it demonstrates



Accuracy of NBup with additional noise in batch and incremental learning modes (a) is shown in Figure 9. Specifically, (b) the precision of SVM in batch and incremental learning modes when subjected to additional noise. (c) The precision of NN in batch and incremental learning modes when subjected to additional noise. (d) The precision of DT when exposed to additional noise in both batch and incremental learning modes. (e) The precision of IBK in batch and incremental learning modes when subjected to additional noise. (f)

Whether LWL is accurate in batch or incremental learning modes when subjected to additional noise.

Figure 9: Accuracy that excels in low- to medium-level noise environments (f). In conclusion, incremental learning mode benefits most algorithms under noise, with the exception of NN, which only performs poorly when the noise level is equivalent to 50%.

It has been observed that algorithms run in batch learning mode require much more time than those run in incremental learning mode. For instance, LWL requires very little time to learn in either incremental or batch modes, and this holds true regardless of the quantity of noise introduced into the input stream.

The time differences between NBup (10% and 20% noise) and SVMandNN are striking. When these traditional classifiers are trained incrementally, using just a subset of data at a time, but still expected to provide acceptable results, they learn far more quickly. Figure 10(a) demonstrates, however, that NBup is very unstable because of the algorithm's

International Journal of Distributed Sensor Networks

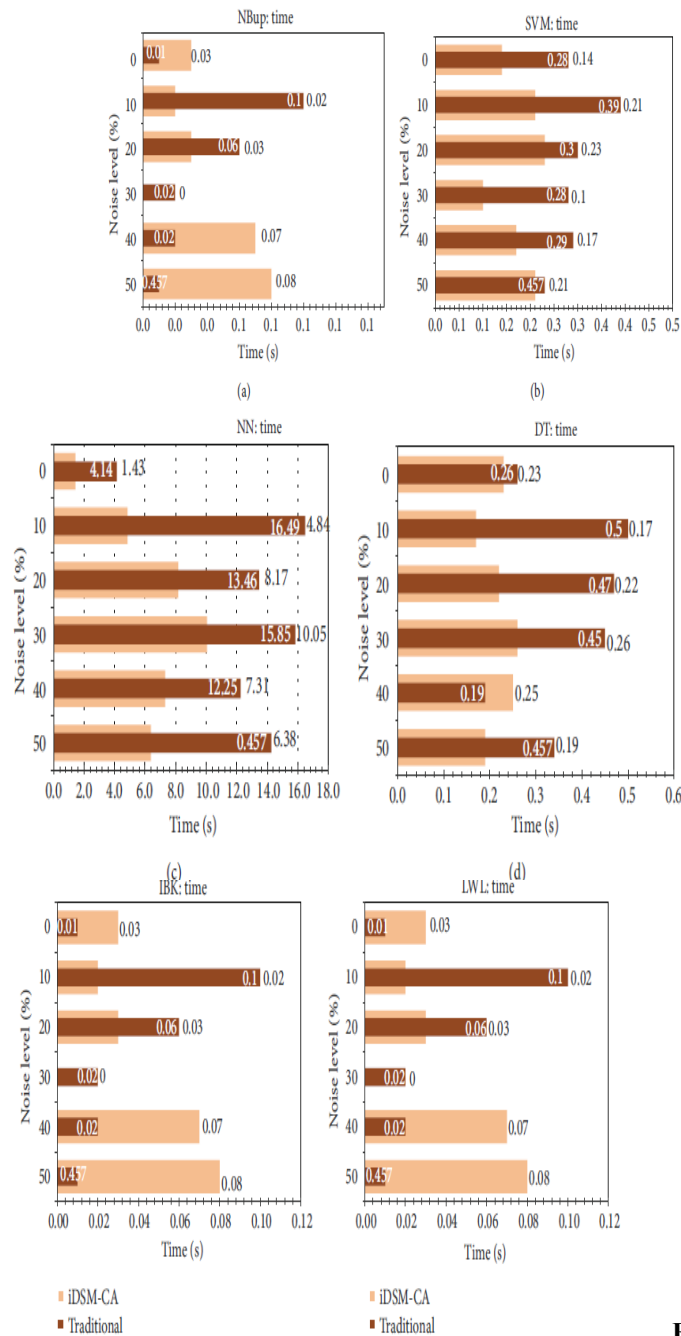
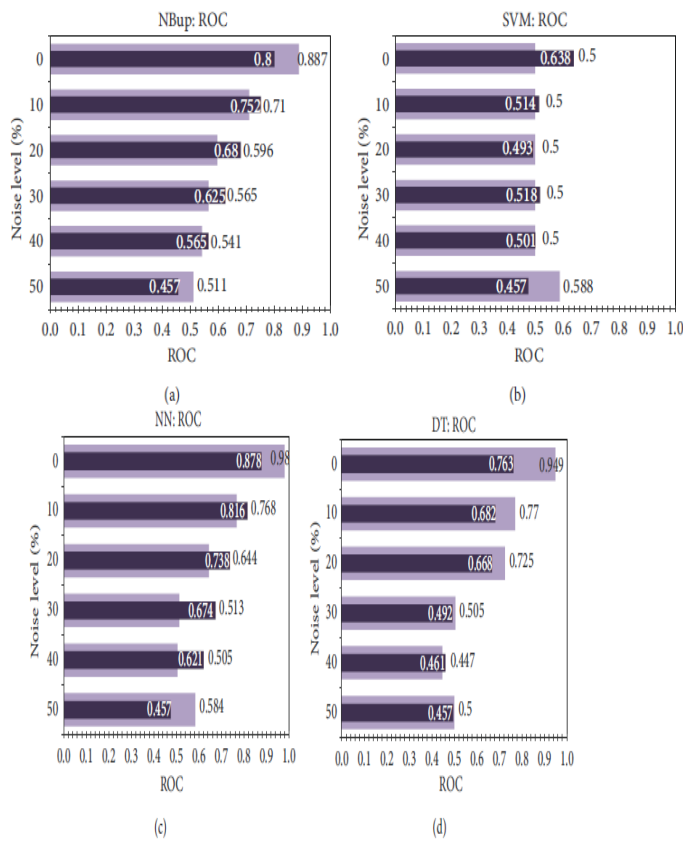


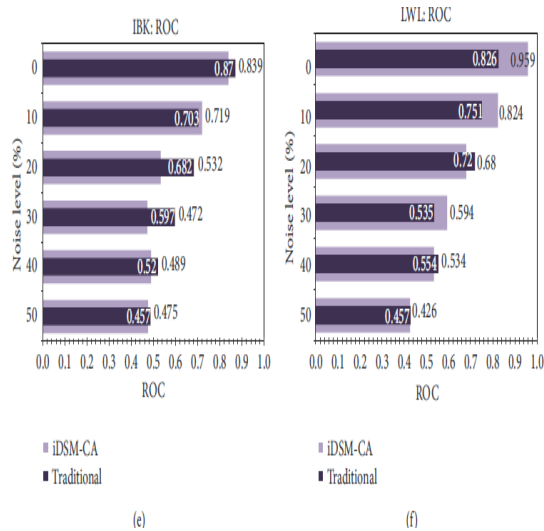
Figure 10(a) compares the NBup's time in batch and incremental learning modes when subjected to additional noise. (b) The length of time it takes for SVM to train in batch and incremental modes when subjected to additional noise. (c) The time required for NN to train under conditions of high background noise, both in batch and incremental modes. (d) The amount of time it takes DT to learn in batch and incremental modes with additional noise. (e) The length of time it takes IBK to learn in batch and incremental modes when there is additional noise. Time of LWL under (f)

increased background noise while doing incremental or batch-based training.

is largely reliant on probabilities, where ambiguity is more likely to develop.

The ROC measures how well a classifier can distinguish between similar examples, and NBup, SVM, and NN in batch learning mode all show rather stable performance over a wide range of noise levels. However, while operating in incremental learning mode, these algorithms fall short of ROC performance due to a high number of false alarms and missed targets. This is mostly because incremental learning only allows for a subset of data to be utilized for training at a time, hence reducing the overall amount of data that must be processed. Model induction in batch learning might take its time and use all available data. In most cases, however, the incremental learning method of DT and LWL proves superior in terms of ROC. Lightweight algorithms, it should be noted, seem to execute quickly and





(a) ROC of NBup with additional noise in batch and incremental learning modes. Batch and incremental ROC of SVM with additional noise. (c) Relative Operating Characteristic (ROC) of NN in the presence of supplemental noise in both batch and incremental learning modes. (d) ROC of DT in batch and incremental learning modes with additional noise. (e) ROC of IBK in batch and incremental learning modes with additional noise. Under (f), the ROC of LWL is

increased background noise while doing incremental or batch-based training.

superior performance (accuracy and ROC) in the incremental mode (iDSM-CA). To evaluate how batch learning mode and incremental learning mode fare when applied to the job of sonar signal detection, we take an average of the performance characteristics across all methods. Table 1 provides a quick overview of the mean findings. While incremental learning outperforms batch learning on average when dealing with mild sounds, its performance immediately declines when dealing with more severe noises because of the inherent limitations of inducing a model from incomplete data. Similar results are shown when examining the ROC and Kappa statistics, which are often used to represent the consistency and generalizability of the datasets.

Compared to batch learning, incremental learning is much more efficient.

TABLE 1: The averaged performance indicator values for batch learning mode and incremental learning mode over the six classification algorithms.

Noise %	0	10	20	30	40	50
Accuracy %						
iDSM-CA	85.17857	74.10714	63.75	55.89286	51.42857	51.60714
Traditional	75.68683	66.6209	65.79669	56.86814	53.15934	48.07693
Kappa						
iDSM-CA	0.66743	0.4507	0.254329	0.082214	0.011529	0.032143
Traditional	0.508057	0.323071	0.3093	0.134429	0.0585	-0.03846
Time (sec.)						
iDSM-CA	0.261429	0.752857	1.235714	1.487143	1.114286	0.98
Traditional	0.67	2.497143	2.041429	2.371429	1.821429	2.121429
ROC						
iDSM-CA	0.865	0.733286	0.615571	0.526571	0.505429	0.503714
Traditional	0.817	0.715714	0.675143	0.585286	0.543429	0.472571

5, Final Remarks

Although sonar has made important contributions, its accurate identification is a difficult challenge.

using them for war purposes.

Underwater noise is a substantial contributor to the accuracy degradation.

Classification model building is complicated by noise. In the context of training datasets, "noisy data" refers to instances that are at odds with the rest of the data and hence corrupt the training patterns, leading to incorrect classification rules. The accuracy of the classification model may be improved by removing noise, which is also known as outliers, misclassified cases, or misfits. Even though this has been a subject of study for over twenty years, the methods that have been developed to date for eliminating this kind of noise all make the assumption that batch processes are necessary, with the whole dataset being utilized for noise identification.

This work presented iDSM-CA, an innovative preprocessing technique for mining incremental data streams via the use of conflict analysis. One of the key benefits of the iDSM-CA is its ability to mine continuously evolving data streams using a lightweight sliding window method. When compared to more involved methods like those described in Section 2, the iDSMCA model is a breeze to implement. Using empirical sonar data to differentiate between metal and rock items, our experiment verifies its merits in terms of its fast speed and its usefulness in delivering a noise-resilient streamlined training dataset for incremental learning.

Experiments show that iDSM-CA is efficient and useful for mining streamdata. Many big data applications, such as data mining sporting events [19] and social media data feeds, rely on real-time analysis of data streams.

real-time [20] among more examples.

An Interests Conflict Exists

This paper's authors report no conflicts of interest that would prevent them from freely sharing their findings with the public.

Acknowledgments

The authors acknowledge the funding provided by the University of Macau, FST, and RDAO via the Research Grant "Adaptive OVFD with Incremental Pruning and ROC Corrective Learning for Data Stream Mining," Grant no. MYRG073(Y3-L2)-FST12-FCC.

References

SONAR systems and underwater signal processing: traditional and new techniques, by H. Peyvandi, M. Farrokhrooz, Roufarshbaf, and S.-J. Park, in *Sonar Systems*, N. Z. Kolev, Ed., pp. 173- 206, InTech, Hampshire, UK, 2011.

International Journal of Distributed Sensor Networks, 2014, vol. 2014, Article ID 374028, 10 pages; C. Lv, S.Wang, M. Tan, and L.Chen, "UA-MAC: an underwater acoustic channel access approach for dense mobile underwater sensor networks."

Sonar sensor for precise 3D target localisation and classification. H. Akbarally and L. Kleeman. *Proceedings of the 1995 IEEE International Conference on Robotics and Automation*.

In the *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, volume 3, pages 1446-1451, November 2001, A.Heale and L.Kleeman wrote about "fast target classification using sonar."

Balleri, A., "Biologically Inspired Radar and Sonar Target Classification," Ph.D. thesis, University College London, London, UK, 2010.

Reference: [6] "Stream-based biomedical classification algorithms for evaluating biosignals," *Journal of Information Processing Systems*, Korea Information Processing Society, vol. 7, no. 4, pp. 717-732, 2011.

Artificial Intelligence Review, volume 22, issue 3, pages 177-210, 2004; X. Zhu and X.Wu, "Class noise vs. attribute noise: a quantitative investigation."

To wit: [8] V. J. Hodge and J. Austin, "A study of outlier identification approaches," *Artificial Intelligence Review*, vol. 22, no. 2, pages 85-126, 2004.

C. Aggarwal and P. Yu, "Outlier detection for high dimensional data," in *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 37-46, May 2001. According to "BIRCH: an effective data clustering approach for extremely big databases," written by T. Zhang, R. Ramakrishnan, and M. Livny and published in the *Proceedings of the Conference on Management of Data (ACM SIGMOD '96)*, pp. 103-114, 1996.

According to

[11] "A linear technique for deviation identification in big databases," by A. Arning, R. Agrawal, and P. Raghavan, published in the *Proceedings of the 1996 International Conference on Knowledge Discovery and Data Mining (KDD '96)*, Portland, Oregon, USA, pages 164-169.

"Enhancing data analysis using noise reduction," *IEEE Transactions on Knowledge and Data Engineering*, volume 18, issue 3, pages 304-319, 2006.

[12] H. Xiong, G. Pandey, M. Steinbach, and V. Kumar.

Specifically,

[13] "Advances in instance selection for instance-based learning algorithms," by H. Brighton and C. Mellish, was published in 2002 in *Data Mining and Knowledge Discovery*, volume 6, issue 2, pages 153-172.

Reference:

[14] C. E. Brodley and M. A. Friedl, "Identifying and removing mislabeled training instances," in *Proceedings of the 1996 13th National Conference on Artificial Intelligence (AAAI '96)*, pp. 799-805, AAAI Press, Portland, Ore., USA, August 1996.

The following piece of research is based on the work of G. H. John: "Robust decision tree: eliminating outliers from databases," published in the 1995 volume of the *Proceedings of the First International Conference on Knowledge Discovery and Data Mining*.

Proceedings of the 2008 International Conference on Artificial Intelligence (ICAI '08), pp. 821-827, Las Vegas, Nevada, USA, July 2008. [16] B. Byeon, K. Rasheed, and P. Doshi, "Enhancing the quality of noisy training data using a genetic algorithm and prototype selection."

International Journal of Data Warehousing and Mining, volume 3, issue 3, pages 14-27, 2007. Authors: A. Nanopoulos, A. N. Papadopoulos, Y. Manolopoulos, and T. Welzer-Druzovec.

Based on the work of R. P. Gorman and T. J. Sejnowski, "Analysis of hidden units in a multilayer network trained to categorize sonar targets," published in 1988's *Neural Networks*, volume 1, issue 1, pages 75-89.

In the Proceedings of the 2013 International Symposium on Computational and Business Intelligence (ISCBI '13), pages 88-91, 2013; I. Fister Jr., I. Fister, D. Fister, and S. Fong; "Data mining in athletic activities provided by sports trackers."

"Mining twitterspace for information: classifying sentiments programmatically using Java," by J. Fiaidhi, O. Mohammed, S. Mohammed, S. Fong, and T. H. Kim, was published in the proceedings of the IEEE 7th International Conference on Digital Information Management (ICDIM '12), held in Taipa, Macau, in August 2012.

.