

# Noise Removal Process from Label Classification using Machine Learning



Mokshada Kotwal, Shraddha Khonde

**Abstract:** Text classification and clustering approach is essential for big data environments. In supervised learning applications many classification algorithms have been proposed. In the era of big data, a large volume of training data is available in many machine learning works. However, there is a possibility of mislabeled or unlabeled data that are not labeled properly. Some labels may be incorrect resulted in label noise which in turn regress learning performance of a classifier. A general approach to address label noise is to apply noise filtering techniques to identify and remove noise before learning. A range of noise filtering approaches have been developed to improve the classifiers performance. This paper proposes noise filtering approach in text data during the training phase. Many supervised learning algorithms generates high error rates due to noise in training dataset, our work eliminates such noise and provides accurate classification system.

**Index Terms:** label noise, majority voting, unlabeled, supervised learning.

## I. INTRODUCTION

In a supervised classification approach, the standard of a dataset is characterized by data sources: the predictor attributes and also the categorical attribute that defines the categories. The standard of the predictors is set by their quality to represent the instances to be classified, and also the quality of the category attribute is set by the right assignment of every instance. The standard of a dataset is set by internal and external factors. The inner issue reveals if the predictors and also the categories has been properly chosen and area unit well outlined. The external issue measures errors introduced within the predictors or within the category assignment, either consistently or unnaturally. In explicit, associate degree instance contains noise once it causes issues because of external reasons. As the number of electronic documents are provided from resources text mining studies have gained more importance. The resources of unstructured and semistructured data integrate , governmental electronic repositories, news articles, biological databases, chat rooms, digital libraries, on-line forums, electronic message and diary repositories. Therefore, it is essential to have a correct classification and data exploration from these resources. The

noise contained in the training data are divided into two classes : 1) attribute noise 2) label noise. Attribute noise is outlined as associate degree of fault or an inaccuracy within the attribute values, whereas label noise is due to mislabeling. These forms of noise have been studied in many works that conlude that removal of attribute noise decrease the prognosticative accuracy of a classifier when same attribute values are fed to the classifier in next iteration whereas removing label noise continuously improves the predictive accuracy of classifier. The proposed work focuses on label noise. The effect of different kinds of noise on text classification performance is analyzed by performing experiments on synthetic and real life noisy datasets .

## II. LITERATURE SURVEY

According Donghai Guan et. al. [1] Proposed a system to improve label noise filtering process by effectively using untagged knowledge. System propose a completely unique noise filtering algorithmic program called as increased soft majority choice by exploring untagged knowledge (ESMVU), that is an ensemble-learning-based filter that follows a soft majority choice method. To improve noise filtering process ESMVU uses untagged knowledge by calculating the reliability of labels from a noisy training dataset. By recognizing label confidence and statistical distribution ESMVU provides an effective use of untagged knowledge. This approach uses untagged knowledge for illegal knowledge filtering. In ESMVU the vote process and confidence measurement are followed to improve the classifier accuracy.

Hongqiang Wei et. al. [2] proposed a system to improve label noise identification using untagged data. It projected a way called as MFUDCM (Multiple Filtering with the help of untagged information victimization Confidence Measurement). This methodology follows the unique multiple soft majority ballot plan to operate on untagged information. Additionally, MFUDCM is anticipated to own a better accuracy of characteristic illegal information by victimization the idea of voting. This work adopts multiple filtering technique than one level filtering as other works followed. To utilize the untagged knowledge it uses algorithmic rules. Hongjiao Guan et. al. [3] The main concept of this projected technique is to rigorously place a lot of particularize in the bulk category than the minority category throughout knowledge improvement. In depth operations on artificial and real knowledge confirms the prevalence method against alternative knowledge improvement strategies. This paper proposes personal improvement for innumerable unbalanced knowledge sets.

Manuscript published on 30 September 2019

\* Correspondence Author

**Mokshada Kotwal**, Computer Department, M.E.S College of Engineering, Pune, India.

**Prof. Shraddha Khonde**, Computer Department, M.E.S College of Engineering, , Pune, India.

© The Authors. Published by Blue Eyes Intelligence Engineering and Sciences Publication (BEIESP). This is an [open access](https://creativecommons.org/licenses/by-nc-nd/4.0/) article under the CC-BY-NC-ND license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

Current systems for knowledge improvement having complications with unique cases and outliers in minority categories, specially in extremely unequal knowledge. The downside gets inadequate and incorrect cases to remove in order to enhance the hardiness and functions accurate knowledge improvement, system adopts a weighted edited nearest neighbor (WENN) that identifies and discards noise instances from both classes intelligently.

Kaio Rodrigues et. al. [4] In this work it projected a method to remove DUST by using an outsized range of URLs gathered by net crawlers results in webpage's having replicated contents in nearby areas. To drag, accumulate and use such a replicated knowledge results in a wastage of resources, low end users experiences and foundation of quality based ranks. Many research works have been done to change above downside to explore and take away replicated work while not taking their contents. In order to achieve this, the predicted strategies study normalization rules to redesign all replicate URLs into a similar authorised way. A difficult side of this approach is etymologizing a collection of common and actual rules. In this task, system gift a DUSTER a replacement way to obtain quality rules that profit of a multiple sequence alignment strategy. System shows a complete multiple sequence alignment of URLs with replicated content before the formation of foundations which may leads to the foundation of powerful rules.

Turki et. al. [5] proposed a system of gene regulatory networks by combination of supervised and unattended ways. Supervised ways for inferring factor regulative networks (GRN) functions well with sensible coaching information. Yet, once coaching information is absent, these ways aren't applicable. Unattended ways don't want coaching information however their accuracy is low. In this work, system mix supervised and unattended mechanism to deduce GRNs blaxploitation space organic phenomenon information. Particularly, system use outcomes returned from unattended ways to coach supervised ways. However there is noise contained within results, system generate a rule for knowledge clearing to get rid of noise, improving the standard of the coaching information. These classy coaching information measure usually guide classifiers as well as support vector machines and deep learning tools to deduce GRNs via link prediction.

Victoria Chayes et. al. [6] proposed Pre-processing and classification of hyper spectral imaging representational process via selective in painting. System planned semi-supervised rule for process and classification of hyper spectral imagery. For data formatting, system keeps hundredth of the information complete, and use Principle Element Analysis to eliminate volume elements from noisy bands and pixels. After that system uses an Accelerated Proximal Gradient Rule (AGPL) or a changed AGPL rule cost for distance between painted elements and finish members on the began information cube to inpaint the missing information. APGL and APGL Hyper differentiated by performance on datasets either all pixels removed or noise removed. This in painting method ends up with one by one band information cube pointing and noise removal from each element.

Chun Lung Philip Chen et. al. [7] proposed a system a weighted couple sparse representation. Within the projected work, the sophisticated exchange in the recreation and also the pictures square measure used to form the secret writing

coefficients applicable to repair the noiseless image. Additionally, the image square measure divided as clear, low impaired and high impaired. Total data-integrity regularizations square measure applied to different pixels to raise the de-noising work. In projected methodology, the lexicon is explicitly skilled on the crying knowledge by sending a weighted rank one diminution drawback, that may capture more features of the original knowledge.

Chao Shang et. al. [8] proposed a system that initially treats every view as a separate domain and associate a domain to domain mappings. To regenerate the removed read from the GAN outputs, it employs a multi-modal denoising motor vehicle encoder that supports paired information over the views. By enhancing DAE and GAN together design permits the data mixture for domain mappings that corresponds to the efficient recovery of removed knowledge. Results shows the validation of this research work operating on standard datasets.

Hao S et. al. [9] proposed a system for noisy data detection supported by DBSCAN and SVDD which is a replacement technique in order to improve the standard of real information sets by taking noisy data as input. Density-Based abstraction Agglomeration of Applications with Noise (DBSCAN) and support vector information description (SVDD) was developed in this paper. First DBSCAN formula is applied to gather information and take away the exceptions. Secondly, SVDD track the classified information in step with the gathered result, and achieve discriminate design for every cluster. These discrimination designs were utilized in entity dataset to distribute the information. The data doesn't fit to any category is known as noisy information therefore removed.

### A. Objective of System

- To identify and detect the noisy label set during the system training.
- Using Ensemble-Learning base voting approach can be detect the noise unlabeled data and measure the system performance before as well as after applying the filters.
- To design and implement a system of multi label classification approach for synthetic high dimensional data and analyze the issues of redundancy in runtime classification.
- To implement a system which can carried out the clustering as well as micro-clustering according to similarity weight.
- To classifies unlabeled data into predefined categories according to text contents with maximum accuracy and highest similarity.
- To implement a micro cluster classification approach on high dimensional data using density base approach.

## III. PROPOSED METHODOLOGY

The proposed algorithm identifies the noisy elements and differentiate them from the elements that are in class boundary. The aim of the algorithm is to identify and eliminate the noisy instances, preserving the class distribution and class boundaries such that neither separate the classes nor the discriminate power of the classification algorithm is changed.

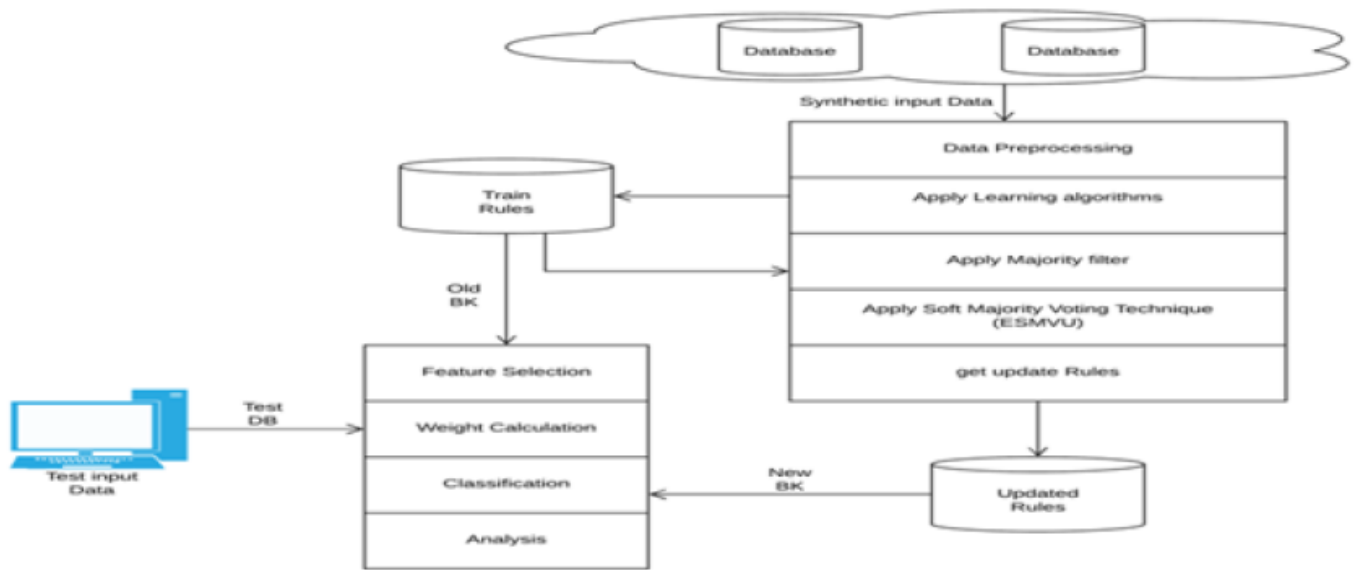


Fig: 1 Proposed System Architecture

In the proposed system we described deep learning base Machine Learning (RNN), basically the system contains two different phases like training as well as testing. In the training phase the system initially uses Natural Language Processing (NLP) to extract the best feature. During the training phase features are extracted by system and these features known as data pre-processing and normalization. Once training has done system stores relevant features into the behalf of respective domain, actually system works like a supervised learning and this extracted feature known as background knowledge of desired domain. After the completion of training phase, it moves for evaluation of the specific test object with the help of proposed classification algorithm. The system uses text feature evaluation technique based on similarity index, and RNN has used to classify the respective test object in testing. When multiple objects are given as testing module, system first extract the Background Knowledge (BK) using NLP and select top features according to their weights. TF-IDF has been used to generate the respective terms weight using NLP.

### I. SYSTEM ANALYSIS

#### Algorithms

##### 1: Weight calculation Algorithm (NN)

**Input :** Training Rules Tr[], Test Instances Ts[], Threshold T.

**Output :** Weight w=0.0

**Step 1 :** Read each test instance from (TsInstnace from Ts)

**Step 2 :**  $TsIns = \sum_{k=0}^n \{Ak \dots An\}$

**Step 3 :** Read each train instance from (TrInstnace from Tr)

**Step 4 :**  $TrIns = \sum_{j=0}^n \{Aj \dots Am\}$

**Step 5 :**  $w = WeightCalc(TsIns, TrIns)$

**Step 6 :** if (w >= T)

**Step 7 :** Forward feed layer to input layer for feedback  
FeedLayer[] ← {Tsf,w}

**Step 8 :** optimized feed layer weight, Cweight ← FeedLayer[0]

**Step 9 :** Return C weight

##### 2 : Stop word Removal Approach

**Input:** Stop words list L[], String Data D for remove the stop words.

**Output:** Verified data D with removal all stop words

**Step 1:** Initialize the data string S[].

**Step 2:** initialize a=0,k=0

**Step 3:** for each(read a to L)

If(a.equals(L[i]))

Then Remove S[k]

End for

**Step 4:** add S to D.

**Step 5:** End Procedure

##### 3 Stemming Algorithm

**Input :** Word w

**Output :** w with removing past participles as well.

**Step 1:** Initialize w

**Step 2:** Intialize all steps of Porter stemmer

**Step 3:** for each (Char ch from w)

If(ch.count==w.length()) && (ch.equals(e))

Remove ch from(w)

**Step 4:** if(ch.endsWith(ed))

Remove 'ed' from(w)

**Step 5:** k=w.length()

If(k (char) to k-3 .equals(tion))

Replace w with te.

**Step 6:** end procedure

##### 4 TF-IDF

**Input :** Each word from vector as Term T, All vectors V[i...n]

**Output :** TF-IDF weight for each T

**Step 1 :** Vector = {c1, c2, c3....cn}

**Step 2 :** Aspects available in each comment

**Step 3 :** D = {cmt1, cmt2, cmt3, cmtn}

and comments available in each document

Calculate the Tf score as

**Step 4 :**  $tf(t,d) = (t,d)$

t=specific term

d= specific document

in a term is to be found.



**Step 5 :**  $idf = t \rightarrow \text{sum}(d)$

**Step 6:** Return  $tf * idf$

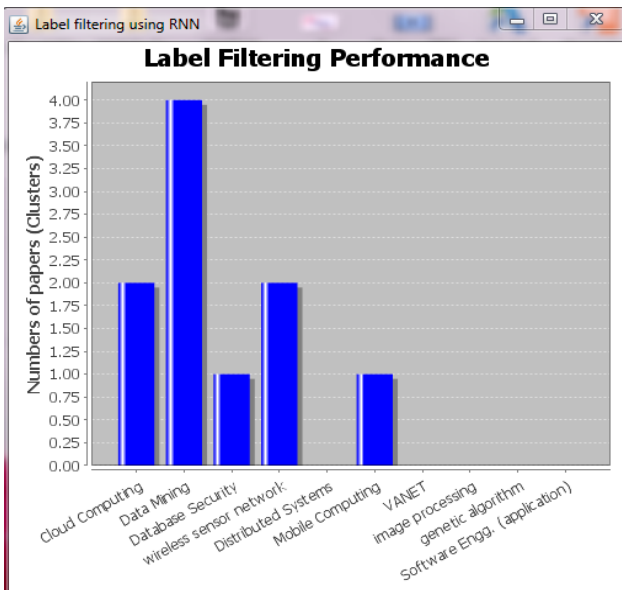
## II. RESULTS AND DISCUSSION

The proposed learning scheme explicitly models the integration of via label graph learning, which is enhanced with multilabel classification. The label integration graph have capability to work well with multilabel classification simultaneously reflects the descriptive architecture between labels. Also we presented a well established familiar center to capture the context dependent inter-label interaction information. Using classification and NN weight calculation approach the projected work strongly classifies the test instance with powerful labels. Experimental results have validate the performance of our work over number of standard datasets.

## III. CONCLUSION

The system proposed noise removal approach in supervised learning algorithms. Different machine learning algorithms have been used to generate the training rules and evaluate it with different test datasets. Proposed ensemble technique also remove error rate during testing of system. Above experimental analysis finally conclude the below statements.

- Using proposed ensemble learning approach it eliminate the noise during the system training.
- It improves the accuracy of testing phase.
- Minimized the error rate, increased the efficiency of system.



**Fig. 2: System Performance Measures proposed vs. existing approaches**

## REFERENCES

1. Guan D, Wei H, Yuan W, Han G, Tian Y, Al-Dhelaan M, Al-Dhelaan A. Improving Label Noise Filtering by Exploiting Unlabeled Data. *IEEE Access*. 2018;6:11154-65
2. W. H, Z. Q, A. Khattak and C. F, ". Improved label noise identification by exploiting unlabeled data", in *International Conference on Security, Pattern Analysis, and Cybernetics (SPAC)*, 2017, pp. 284-289.

3. G. H, T. X, Z. Y and X. M, "WENN for individualized cleaning in imbalanced data", in *23rd International Conference on Pattern Recognition (ICPR)*, 2016, pp. 456-461.
4. R. K, d. A, C. M and M. ES, "Removing DUST using multiple alignment of sequences.", 2015, pp. 2261-74. *IEEE Transactions on Knowledge and Data Engineering*.
5. T. T, R. I and W. JT, "Inferring gene regulatory networks by combining supervised and unsupervised methods", in *Machine Learning and Applications (ICMLA)*, 2016, pp. 140-145.
6. B. AL, L. W, C. V, B. R and L. J, "Pre-processing and classification of hyperspectral imagery via selective inpainting.", in *In Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 6195-6199.
7. L. L, Z. Y, C. L and T. YY, "Weighted couple sparse representation with classified regularization for impulse noise removal.", in *IEEE Transactions on Image Processing*, 2015, pp. 4014-26.
8. P. A, C. KS, S. C and S. J, "VIGAN: Missing view imputation with generative adversarial networks", in *IEEE International Conference*, 2017, pp. 766-775.
9. H. S, S. H and Z. X, "A new method for noise data detection based on DBSCAN and SVDD", in *IEEE International Conference on Cyber Technology in Automation, Control, and Intelligent Systems (CYBER)*, 2015, pp. 784-789.

## AUTHORS PROFILE



**Mokshada Kotwal** , Received her Bachelor of Engineering degree in computer field from SSBT College of Engineering, North Maharashtra University in 2011. She is currently pursuing Masters degree in Engineering in Computers from M. E. S. College of Engineering, University of Pune. Currently working on a machine learning project.

She published a survey paper based on the current research in IJRCCCE, a UGC approved journal and a research paper published in National Level Conference on "Recent Advances in Computer Engineering" held at MESCoE, Pune.



**Prof. (Mrs.) Shradha Khonde** Received her Bachelor in Computer Engineering and Master in Computer Engineering. She is pursuing her Ph. D. in Computer Sc. & Engineering. She has more than 14 years of teaching experience. At present, she is working as Assistant Professor, Department of Computer Engineering at M. E. S. College of Engineering (Wadia College). Her area of

specialization includes network security and machine learning. She owned memberships of ISTE, IETE and CSI. She has published more than 30 papers in national, international conferences and journals.

Awards and Achievements:

1. Received Best paper award in International Conference on Modern Technologies in Engineering and Science (ICMTES) on "An improved Technique of extracting frequent itemsets from massive data using mapreduce", on 2nd and 3rd June 2017.
2. Received Best paper award in International Conference on Intelligent Computing and Communication (ICICC) on "A Novel Approach of Frequent Itemset Mining using HDFS Framework".