



Feature Selection towards Soil Classification in the context of Fertility classes using Machine Learning

Janmejy pant, Pushpa Pant, Ashutosh Bhatt, Himanshu Pant, Nirmal Pandey

Abstract: Soil is recognized as one of the most valuable entity which is responsible for sustaining life on the earth. It is clear that Indian economy is largely dependent on agriculture. In India, for a large section of the population the sole or major source of earnings is agriculture. There are many factors which are responsible in yield production and also affect agriculture. Soil is one of them and plays a crucial role in yield production. Soil nutrients are important aspects that contribute to soil fertility. In this paper we classify soil in the form of fertility index level using machine learning. We induced rule model on training data and apply this model on test data for making predictions for fertility level classes. In this paper we have used Rough Set method to classify the data based on fertility level class and calculate the accuracy of the used classifier.

Keywords: Feature Selection, Fertility Index, Rough Set, Soil nutrients, Reduct, Rule

I. INTRODUCTION

Soil nutrients are a crucial property that contributes to the soil fertility and other environment factors [1]. In agriculture it is important to enhance yield or crop production which is dependent on soil management. Enhancing the yield production depends and is directly affected by soil parameters or nutrients such as pH value, organic matters, electrical conductivity, P, K etc. Agriculture is fully dependent on soil quality. But due to heavy agriculture production the quality of soil may decrease and the nutrients present in soil may lost. In recent years machine learning has been very helpful in studying the elimination of nutrients in soil. Nowadays soil classification and prediction problems are easily handled by Machine Learning techniques [2]. Various Machine Learning techniques were used to predict the soil nutrients, soil type and soil moisture [2], [3].

Revised Manuscript Received on October 30, 2019.

* Correspondence Author

Janmejy Pant*, Dept. of Computer Science, Graphic Era Hill University, Bhimtal, India., Email: jpant@gehu.ac.in

Puspha Pant, Dept. of Geography, M.B. Govt. P.G College, Haldwani, India, Email: pushpapant01@gmail.com

Ashutosh Bhatt, Dept. of Computer Science , BIAS, Bhimtal, India, Email:ashutoshbhatt123@gmail.com

Himanshu Pant, Dept. of Computer Science, Graphic Era Hill University, Bhimtal, India, Email: himanshupantbareilly@gmail.com

Nirmal Pandey, Dept. of Computer Science and Applications, Amrapali Group of Institute Haldwani, India, Email:nirmal.bias@gmail.com

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

Soil fertility is affected by many factors like air, water, organic matter and nutrients. Due to some factors, like use of fertilizers, pesticides, insecticides, cultivation etc, and the fertility of soil is regularly decreasing in recent times [4]. In India there is deficiency of primary nutrients in agricultural soil. In this regard it would be better to utilize the different type of land such as saline, alkaline wastelands by analyzing different soil fertility parameters such as fertility index and soil depletion index [4].

II. DATA SET DESCRIPTION

In order to classify soil data based on fertility level, the used data set was collected from one of the soil testing laboratories in Almora district (Uttarakhand). We have used Rosetta tool for generating the rules through which we can further classify our data set and we make prediction for the test data. Based on the soil properties the fertility of soil is defined in three different levels. These are Level 1, Level 2 and Level 3.

We have collected the 53 samples of soil from Almora district (Uttarakhand) by random sampling. The complete data is taken from three different sites corresponding to three different varieties of oak trees growing there, namely Banj oak, Tlionj oak and Kharsu oak. Some samples are distributed and some are undistributed. The primary data was sent to soil testing laboratory for measuring the properties of soil. The data set has 11 independent attributes and 1 is dependent or target attribute which is fertility level. Table 1 describes data collected from each soil sample.

Table I: Attribute Description

Attribute	Description
Depth	in cm from sample is collected
sand	Sand particles in soil ,%
slit	Slit particles in soil ,%
clay	clay particles in soil ,%
moisture	Moisture present in soil,%
pH	pH value of soil
C	Organic carbon,%
N	Nitrogen,%
OM	Organic Matter,%
C:N	Carbon and Nitrogen Ratio in soil
Fertility Level	Dependent class

III. MATERIALS AND METHODS

Machine Learning method is used for classification. For this the following concepts are used:



A. Feature Selection

FS technique [5], [6] is used to obtain a reduced representation of the data set that is smaller in volume in comparison to the original one. The reduced features also maintain the integrity of the original data. Now the feature selection technique is expected to identify most significant soil properties specifying the soil fertility [7].

The feature selection method is implemented to extract the most significant soil property for soil fertility level class. In this paper we have used one machine learning techniques that is Rough Set [5], [7], [8]. Feature Selection has four major steps [7]-

- (i) Generation Procedure
- (ii) Evaluation Subset
- (iii) Stopping Criterion
- (iv) Validation Procedure.

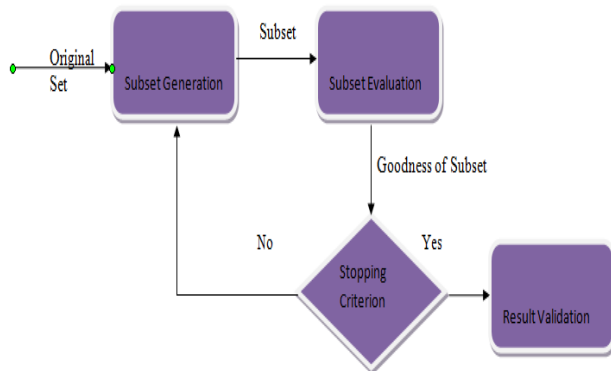


Fig. 1. The feature selection implementation [8].

In this paper, after extracting the important features we have developed an automated system for generating the rules which help in soil classification based on soil fertility.

B. Rough Set Theory

RST [12], [14] is one of the most important mathematical tools to deal with uncertainty which may arise due to granularity in domain of discourse or universe. Rough set is also useful to deal with vagueness of data. It is a mathematical framework for analyzing tabular data. There are various basic concepts defined by Z. Pawlak [12], [14]. Some of them are:

Information System (IS)

The appropriate data which is basically used for research work and is arranged in a tabular format is known as information system (IS) [9]. An information system can be formally defined as: Let $A = \{A_1, A_2, A_3, \dots, A_k\}$ be a non empty finite set of attributes and $U = \{a_1, a_2, \dots, a_k\}$ be a non empty finite set of k -tuples, known as the objects. $V(A_i)$ denotes the set of values for the attributes A_i [9]. Then information system is defined as an ordered pair $I(U, A)$ such that for all $i=1, 2, \dots, k$ there is a function f_i [9], [5]. This means, every object in the set U has an attribute value for every element in the set A . The set U is called the universe of the information system [8], [9], [12]. Thus, we can say that Information System is a data table. Let us assume that $I = (U, A, C, D)$ [5] is a decision table or decision system. Here, U is universe of finite set, A is an attribute of data table I and C and D are subsets of A . C is a conditional subset of A and D is decision attribute. V_a is the value of set of a , the elements of U are objects [5], [8], [9], [12].

Indiscernibility Relation

Let $I = (U, A)$ be an information system where $U = \{a_1 \dots a_k\}$ is the non empty finite set of k tuples known as the object and $A = \{A_1 \dots A_k\}$ is a non empty finite set of attributes [10], [12]. Let P be a subset of attributes [10], [12]. Then the set of P -indiscernible objects is defined as the set of objects having the same set of attribute values [10], [12].

The RST treats data in terms of equivalence classes this means set of objects that are indiscernible with respect to attributes [8], [12].

Reducts

The objects of an information system can be divided into a group of similar types, i. e., equivalence classes [8], [10]. An Indiscernibility relation can do this task. However, the entire set of attributes, A , may not be necessary to preserve the Indiscernibility among the set of indiscernible object [8], [9]. We can say that there may a subset that is sufficient to maintain the classification based on Indiscernibility. A minimal set of attributes required to preserve the Indiscernibility relation among the objects of an information system is called a reduct. Let Say, given an information system a reduct [8], [9], [11] is a minimum group of attributes such that [8], [9].

Decision Rules

After getting reducts, rules may be generated from reducts for classification of the objects. IF-THEN rules [9], [12] are constructed by reading of the values for each attribute in the reduct and associating them with one or more decision classes. The THEN-part will only include one decision class unless the decision class is rough with respect to the attributes in the reduct. Rules are evaluated according to how general they are and how specific they are. How general they are in term of coverage and how specific they are in terms of accuracy.

C. Fundamental Characteristics of Rough Set Theory[8]

1. This theory doesn't require any prelude and additional information about data. Other similar theories need the same information as preliminary, for example, statistics needs probability, fuzzy set theory needs grade of membership as prior information. But Rough set theory does not require these types of information as prior.
2. For extracting the hidden patterns form the data rough set provides efficient and intelligent methods, algorithms and tools.
3. By using rough set we can reduce the original data to find the minimal sets which contain the identical knowledge as in the initial data, i.e. we can reduce the dimension or size of original data without losing the information which is derived by the original data. For this we can use feature selection.
4. It allows calculating the importance of data.
5. Rough set allows producing the set of decision rules in automatic way.

IV. METHODOLOGY

In this paper we have tried to describe how Rough Set theory is useful for feature selection of soil data.

The whole data, which is stored in tabular form, is split in two parts one is training set and the other is test set. The division of data set in two parts is useful for classification and prediction. The training set is classified using rough set to find reducts and rules. Once we get rules, we then apply the rules in test data to make some predictions and classify the test data. Researchers are using Rough set theory as a classifier to extract the relevant attributes to achieve the accurate and exact classification results. Feature selection is basically finding the reducts from data to reduce the volume of the data without losing the information. In this paper we find the reducts, generate rules and obtain classification results for the soil sample based on fertility levels.

V. EXPERIMENT AND RESULTS

We have used rough set in ROSTTA tool to apply our soil samples of different locations of hilly areas in Almora district in Uttarakhand. As we have already discussed in methodology section, we split the data in training and testing parts. For the data we induced a model that is rule model on one part of the data, i.e., induce rule model on training part. We used the same rule model to classify the objects in remaining data, i.e., used the generated rule model on test part. In this section we describe the data used in our experiments to generate results using Rough Set theory.

Fig. 2 Demonstration of Soil Data

As stated above we have collected 53 soil samples from the different location of hilly area in Almora district, Uttarakhand. A snapshot of this data set is in Fig 2. In this data the different sites from where we collected the samples were renamed: Banjoak site renamed as S1, Tlionj oak site renamed as S2 and Kharsu oak site renamed as S3. The above data is split out in two parts with spilt factor 0.5. This means 50 percent of data will spilt in training part and rest 50 percent will split into test part.

Fig. 3 Training Set

Fig. 4 Test Set

Once the data is divided into two parts, we can extract reducts and induce Classification rules on training set by Johnson’s algorithm and standard voting method respectively.

Fig. 5 Extracted Reducts

For feature selection we get 28 reducts followed by their support and length.

Rule	LHS Support	RHS Support	RHS Accuracy	LHS Coverage	RHS Coverage	RHS Stability	LHS Length	RHS Length
1	pH_Value(4) => F(Level2)	2	2	1.0	0.074074	0.089657	1.0	1
2	pH_Value(3) => F(Level2)	2	2	1.0	0.074074	0.089657	1.0	1
3	pH_Value(5) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	1
4	Moisture(52) AND C%(2) => F(Level2)	2	2	1.0	0.074074	0.089657	1.0	2
5	Moisture(48) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
6	Moisture(39) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
7	Moisture(39) AND C%(2) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
8	Moisture(34) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
9	Moisture(31) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
10	Moisture(29) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
11	Moisture(31) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
12	Moisture(27) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
13	Moisture(42) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
14	Moisture(48) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
15	Moisture(29) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
16	Moisture(32) AND C%(2) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
17	Moisture(28) AND C%(1) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
18	Moisture(39) AND C%(2) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
19	Moisture(62) AND C%(2) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
20	Moisture(56) AND C%(2) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	2
21	Clay(24) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	1
22	Clay(15) => F(Level2)	1	1	1.0	0.037037	0.049478	1.0	1

Fig. 6 Generated Sorted Rules

713 classification rules have been generated. These rules are very important for classifying the test set of our soil data. Based on these rules we can predict the exact fertility level class for the test part objects. In the above Fig 6 the generated rules are sorted in by IF and THEN part of the rule. There are 713 rules that have been generated. The following statistics can be defined for each rule.

- **LHS Support:** the total number of objects in the training set matching the IF part.
- **RHS Support:** Number of objects in the training set matching the IF-part and the THEN-part (LHS and RHS support is the same unless the THEN-part contains several decisions).
- **RHS Accuracy:** RHS support divided by LHS support (Accuracy is 1.0 unless the THEN-part contains several decisions).
- **LHS Coverage:** LHS support divided by the number of objects in the training set.
- **RHS Coverage:** RHS Support divided by the number of objects in the decision class listed in the THEN part of the rule.
- **RHS Stability:** Not applicable for the Johnson algorithm (always 1.0).
- **LHS Length:** Number of attributes in the IF-part of the rule.
- **RHS Length:** Number of decisions in the THEN-part of the rule.

The above rules can be used to classify the objects or tuples of test set. Standard voting is used as classification method. In the present study we predict the fertility index class of test set based on the rules.

The confusion matrix [7], [8], [12] in Fig.7 shows the overall accuracy (i.e. 0.961538), as well as the sensitivity [8], [13] and accuracy for each class. For example, the Level 2 decision class has a sensitivity of 1.0 (i.e. of 23+0+0 = 23 objects actually belonging to Level 2, 23 was correctly classified as Level 2: 23/23 = 1.0) and an accuracy of 0.95 (i.e. of 23+1 = 24 objects predicted to Level 2, 23 were actually belonging to this class: 23/24 = 0.95).

		Predicted			
		Level2	Level3	Level1	
Actual	Level2	23	0	0	1.0
	Level3	1	1	0	0.5
	Level1	0	0	1	1.0
		0.958333	1.0	1.0	0.961538
ROC	Class	Undefined			
	Area	3.402820e+038			
	Std. error	3.402820e+038			
	Thr. (0, 1)	3.402820e+038			
	Thr. acc.	3.402820e+038			

Fig. 7 Classification Rule

Through the confusion matrix [7], [8] designed above we can define classifier's performance in terms of accuracy for the given soil data. In our experiments on soil samples the accuracy of rough set is approx 96% which defines 96 percent of tuples are correctly classify by the classifier in terms of fertility level class that is the target class for the data samples.

VI. CONCLUSION

In this paper a feature selection is used towards soil classification in terms of soil fertility level classes using machine learning. The target class fertility level is considered in three different levels based on the soil properties. Rough set theory is applied for classification and prediction results and finally calculated the performance of classifier in terms of accuracy where predicted values and actual values also play an important role to calculate the performance of classifier

VII. FUTURE SCOPE

In current research machine learning is widely used in soil science. We can extend this assignment using other machine learning techniques like ANN, Fuzzy etc. Deep learning is another important tool that can be used for implementing a classification model to predict fertility.

ACKNOWLEDGMENT

The first author is thankful to Prof R P Pant, Dept. of Mathematics, DSB Campus, Kumaun University Nainital for his continuous guidance and valuable suggestions.

REFERENCES

1. Hao Li, Weijia Leng, Yibing Zhou and Zhilong Xiu., "Evaluation Models for Soil Nutrient Based on Support Vector Machine and Artificial Neural Networks", Hindawi Publishing Corporation, The Scientific World Journal, 2014, Volume. Article ID 478569. 7 pages.
2. M.S. Suchithra, Maya L. Pai, "Improving the prediction accuracy of soil nutrient classification by optimizing extreme learning machine parameters", Information Processing in Agriculture, 2019, <https://doi.org/10.1016/j.inpa.2019.05.003>.
3. Reashma SJ, Pillai AS, "Edaphic factors and crop growth using machine learning—A review", International conference on intelligent sustainable systems (ICISS), 2017 IEEE 2017, p. 270–74.
4. Nikhita Awasthi, Abhay Bansal, "Application of Data Mining Classification Techniques on Soil Data Using R", International Journal of Advances in Electronics and Computer Science, 2017, ISSN: 2393-2835, Volume-4. Issue.
5. Janmejy Pant, Kamlesh Padaliya, Himanshu Pant, "Rough Set Approach for Feature Selection in IDS", International Journal of Innovations & Advancement in Computer Science, 2015, Volume 4. Special Issue September. ISSN 2347 – 8616.

7. Janmejy Pant, Amit Juyal , Shivani Bahuguna,"Soft set, a soft Computing Approach for Dimensionality Reduction", International Journal of Innovative Science, Engineering & Technology,2015,Vol. 2, Issue 4.
8. Marzieh Mokarram, Mehran Shaygan, George ch. Miliarezi,"Balancing Soil Parameters and Farmers Budget by Feature Selection and Ordered Weighted Averaging" Geographia Technica,2018,Vol.13, Issue 1,pp 73 to 84.
9. Janmejy Pant, R.P Pant, Amit Juyal, Himanshu Pant,"Rule Generation Methods for Medical Data (Heart Disease) Using Rough Set Approach",International Conference on Advances in Engineering Science Management & Technology .Availableat SSRN: <https://ssrn.com/abstract=3384652>.
10. Samir Roy, Udit Chaudhary."Introduction to Soft Computing: Neuro-Fuzzy and Genetic Algorithms", pp 145-171.
11. Priyanka Suyal, Janmejy Pant, Akhilesh Dwivedi & Manoj Chandra Lohani,"Performance Evaluation of rough set based classification models to intrusion detection system".2nd International Conference on Advances in Computing, Communication, & Automation (ICACCA), 2016.
12. Janmejy Pant, Amit Juyal, Himanshu Pant, Akhilesh Dwivedi, "A Real Time Application of Soft Set in Parameterization Reduction for Decision Making Problem ".International Journal of Electrical and Computer Engineering (IJECE), 2017, Vol. 7, No. 1, pp. 324~329 ISSN: 2088-8708, DOI: 10.11591/ijece.v7i1.pp324-329.
13. Pawlak, Z., "Rough Set", In International Journal of Computer and Information Sciences,Vol. 11, pp.341-356, 1982.
14. Jensen, J., "Combining rough set and fuzzy sets for feature selection". Ph.D Thesis from Internet, 2005.
15. Maryam Zavareha and Viviana Maggioni,"Application of Rough Set Theory to Water Quality Analysis: A Case Study", 2018, MPDI. www.mdpi.com/journal/data

of Uttarakhand Science Education & Research Centre (USERC). He had served as State Student Coordinator of Region I Uttarakhand of CSI (for three year) and also served as Secretary of Uttarakhand ACM Professional Chapter.



Himanshu Pant is pursuing Doctorate from Graphic Era Hill University. He received his MCA from UPTU Lucknow. His areas of interest are Machine Learning, Soft computing and data mining. He has published various papers in the different journals. Currently he is working as Assistant professor in Graphic Era Hill University Bhimtal. He has an experience of around 07 years of teaching



Nirmal Pandey is an Assistant Professor at Faculty of Computer Science and Applications in Amrapali Group of Institute Haldwani, Nainital. Presently he is pursuing Ph.D. from Department of Information Technology, Kumaun University, and Nainital. He has done his Master degree of Computer Applications from Birla Institute of Applied Sciences, Bhimtal; Nainital in year 2000. His research area includes Machine Learning, Soft computing. He has 6 years of software industry experience specialization in core Java and 12 years of professional teaching experience.

AUTHORS PROFILE



Janmejy Pant is pursuing his Doctorate in Information Technology from Kumaun University Nainital and done M.tech in Information Technology from Graphic Era deemed to be University, Dehradun. Currently he is working as Assistant Professor in Dept. of computer science at Graphic Era Hill University, Bhimtal Campus. He has total Academic teaching experience of more than 8 years with more than 25 publications in reputed National and International SCOPUS, UGC Approve Journals and conferences. His research area includes Machine Learning, Soft Computing, and Data Mining and Deep Learning. He is an active member of CSI. He secured First Rank in Uttarakhand State Eligibility Test (USET) in Computer Science and Applications in 2018. He also received Research Award at University Level in September 2019.



Pushpa Pant obtained her Doctorate in Geography from Kumaun University Nainital in 1992. She has Post Graduate teaching experience of more than 18 years with more than 10 publications in reputed journals, including SCI and SCOPUS indexed journals, and conferences. Currently she is working as Assistant Professor in M.B. Govt. P.G College Haldwani, Uttarakhand. She was awarded Senior Research Fellowship (SRF) by Council of Scientific and Industrial Research (CSIR) during her doctoral research work (July 1990 – June 1993) and Post Doctoral Fellowship by Indian Council of Social Sciences Research (ICSSR) during 1999 – 2001.



Ashutosh Bhatt obtained his Doctorate in Computer Science from Kumaun University Nainital in 2009. He has total Academic teaching experience of more than 16 years with more than 40 publications in SCI, SCOPUS journals and conferences. His research area includes Deep Learning, Machine Learning, Soft Computing and Data Science. Currently he is working as Associate Professor in BIAS Bhimtal, Uttarakhand. He is also associated with many renowned National/International Journals as Lead Guest Editor/reviewer/editorial board member. He is life member of CSI(Computer Society of India) and senior member of IEEE. Presently he is also serving as District Coordinator