

ANALYSIS OF COMPRESSION AND CLUSTERING OF DATA ANALYSIS USING GRAPHICAL PLOT GENERATION

¹Akkipalli Sowjanya, ²Shiva Kumar B, ³Ashok M, ⁴Vakidi Vineetha

^{1,2,3}Assistant Professor, ⁴UG Student, ^{1,2,3,4}Department of Information Technology, Rishi MS Institute of Engineering and Technology for Women, Kukatpally, Hyderabad.

ABSTRACT

The act of clustering is the division or distribution of a population or set of data points into different groups, or clusters, where data points in the same group are more similar to one another and differ from those in other groups. In many clustering techniques, the cluster centroids play a key role. The project's objectives are to compress the data file, visualize the clustering procedure using a few graphical charts, and finally pinpoint the positions of the cluster centroids. Prior to clustering, the data points are displayed in the first figure, and the data clusters and centroids are displayed in the second plot. Various clusters will be plotted using various colors so that they may be clearly identified from one another. This project is divided into four sections. The first module focuses on developing a user interface for saving customer data in a database and exporting a.csv file from the database. The second module is concerned with compressing the data file and then plotting the unclustered data points. The third module is concerned with the creation of clusters of various colors, as well as the centroids. The final module is concerned with creating the centroids' coordinates.

INTRODUCTION

A tough topic nowadays is assessing the behavior of complex systems, not least because of the challenges brought on by emergence and self-organization. In this sense, behavior refers to the activities and interactions of the dealers as well as how this also affects the surroundings. Often, a device can be described by using simple rules, but these rules can result in many unusual behaviors depending on the initial conditions, and frequently, the observer is unsure whether or not they may also have missed some interesting behavior even after running numerous simulations or making numerous observations in the case of a real system. It seems obvious that the ability to automatically characterize intriguing behaviors inside a complicated system would be extremely valuable. However, to discover such behaviors, it is necessary first to outline what we would possibly imply when we describe a behavior as being interesting. Within the statistics mining and know-how discovery community, there are a number of definitions of the time period "interestingness". As there is no preceding evaluation of interestingness inside complicated systems, this paper explores whether or not it is feasible to routinely classify when fascinating behavior has happened inside a complicated system. As a case study, it will mainly use one-dimensional (1D) fundamental mobile automata as an instance of a complicated machine to locate fascinating behavior inside such systems.

Cellular automata had been invented in the 1940s. There has been extensive research, given that they are nonetheless simply as applicable these days in fields ranging from its use as a pseudo-random sequence generator [1] to designing tiers in mazes [2]. This paper will study 1D cell automata which produce exceptional patterns as a end result of their behaviour relying upon the preliminary prerequisites and the regulations used. In a 1D fundamental mobile automata, the cellphone can be inactive (white) or lively (black) [3] with 256 simple guidelines available. The activation country of the cellphone and the two on both facet of it influences the cellphone on the subsequent iteration. This paper will use cell automata with a single activated cellphone in the

centre of the preliminary row (with subsequent rows being generated every time step when visualising the behaviour of the system). There are many diversifications of policies and preliminary stipulations that generate many distinctive behaviours that are hard to discover manually.

LITERATURE SURVEY

produced, many may additionally be referred to as “uninteresting” patterns. These boring patterns consist of the empty areas and diagonal strains transferring from the centre to a corner, or repeated patterns of horizontal stripes. For an preliminary circumstance of a single cell, there are fascinating patterns that are produced via some of the rules, such as those that seem to be comparable to a Sierpiński Triangle. Currently, there is no listing of cell automata policies that are labeled by means of their fascinating output. This paper will advocate such a list, and additionally a new way of mechanically classifying fascinating mobile automata patterns primarily based on compression.

model based. The approaches used in these methods are discussed with their respective states of art and applicability. The measures of similarity as well as the evaluation criteria, which are the central components of clustering, are also presented in the paper. The applications of clustering in some fields like image segmentation, object and character recognition and data mining are highlighted.

A. Saxena et al., "A review of clustering techniques and developments," *Neurocomputing*, vol. 267, pp. 664-681, 2017. J. Verrelst et al., "Machine learning regression algorithms for biophysical parameter retrieval: Opportunities for Sentinel-2 and 3," *Remote Sensing of Environment*, vol. 118, pp. 127-139, 2012.

hierarchical, partitional, grid, density based and

This paper presents a comprehensive study on clustering: exiting methods and developments made at various times. Clustering is defined as an unsupervised learning where the objects are grouped on the basis of some similarity inherent among them. There are different methods for clustering the objects such as Leaf area index (LAI) is a crucial crop biophysical parameter that has been widely used in a variety of fields. Five state-of-the-art machine learning regression algorithms (MLRAs), namely, artificial neural network (ANN), support vector regression (SVR), Gaussian process regression (GPR), random forest (RF) and gradient boosting regression tree (GBRT), have been used in the retrieval of cotton LAI with Sentinel-2 spectral bands. The performances of the five machine learning models are compared for better applications of MLRAs in remote sensing, since challenging problems remain in the selection of MLRAs for crop LAI retrieval, as well as the decision as to the optimal number for the training sample size and spectral bands to different MLRAs. A comprehensive evaluation was employed with respect to model accuracy, computational efficiency, sensitivity to training sample size and sensitivity to spectral bands. We conducted the comparison of five MLRAs in an agricultural area of Northwest China over three cotton seasons with the corresponding field campaigns for modeling and validation. Results show that the GBRT model outperforms the other models with

respect to model accuracy in average ($R^2 = 0.854$, $R^2 = 0.674$ and $R^2 = 0.456$). SVR

achieves the best performance in computational efficiency, which means it is fast to train, and to validate that it has great potentials to deliver near-real-time operational products for crop management. As for sensitivity to training sample size, GBRT behaves as the most robust model, and provides the best model accuracy on the average among the variations of training sample

size, compared with other models ($R^2 = 0.884$, $R^2 = 0.615$ and $R^2 = 0.452$). Spectral bands sensitivity analysis with dCor (distance correlation), combined with the backward elimination approach, indicates that SVR, GPR and RF provide relatively robust performance to the spectral bands, while ANN outperforms the other models in terms of model accuracy on the average among the reduction of spectral

bands ($R^2 = 0.881$, $R^2 = 0.625$ and $R^2 = 0.480$). A comprehensive evaluation indicates that GBRT is an appealing alternative for cotton LAI retrieval, except

for its computational efficiency. Despite the different performance of the ML models, all models exhibited considerable potential for cotton LAI retrieval, which could offer accurate crop parameters information timely and accurately for crop fields management and agricultural production decisions.

[3] E. Formisano, F. De Martino, and

G. Valente, "Multivariate analysis of fMRI time series: classification and regression of brain responses using machine learning," *Magnetic resonance imaging*, vol. 26, no. 7, pp. 921-934, 2008.

Machine learning and Pattern recognition techniques are being increasingly employed in Functional magnetic resonance imaging (fMRI) data analysis. By taking into account the full spatial pattern of brain activity measured simultaneously at many locations, these methods allow detecting subtle, non-strictly localized effects that may remain invisible to the conventional analysis with univariate statistical methods. In typical fMRI applications, pattern recognition algorithms "learn" a functional relationship between brain response patterns and a perceptual, cognitive or behavioral state of a subject expressed in terms of a label, which may assume discrete (classification) or continuous (regression) values. This learned functional relationship is then used to predict the unseen labels from a new data set ("brain reading"). In this article, we describe the mathematical foundations of machine learning applications in fMRI. We focus on two methods, support vector machines and relevance vector machines, which are respectively suited for the classification and regression of fMRI patterns. Furthermore, by means of several examples and applications, we illustrate and discuss the methodological challenges of Using machine learning algorithms in the context of fMRI data analysis.

[4] M. Shannag, "High strength concrete containing natural pozzolan and silica fume," *Cement and concrete composites*, vol. 22, no. 6, pp. 399- 406, 2000.

Various combinations of a local natural pozzolan and silica fume were used to produce workable high to very high strength mortars and concretes with a compressive strength in the range of 69- 110 MPa. The mixtures were tested for workability, density, compressive strength, splitting tensile strength, and modulus of elasticity. The results of this study suggest that certain natural pozzolan-silica fume combinations can improve the compressive and splitting tensile strengths, workability, and elastic modulus of concretes, more than natural pozzolan and silica fume alone. Furthermore, the use of silica fume at 15% of the weight of cement was able to produce relatively the highest strength increase in the presence of about 15% pozzolan than without pozzolan. This study recommends the use of natural pozzolan in combination with silica fume in the production of high strength concrete, and for providing technical and economical advantages in specific local uses in the concrete industry.

PROPOSED SYSTEM

This project consists of four modules. The first module is about creating the GUI for storing the customer data in the database, and, to create the .csv file from the DB. The second module deals with compressing the data file and then generating the plot of the unclustered data points. The third module deals with generating clusters in different colours, along with the centroids. The final module deals with generating the coordinates of the centroids.

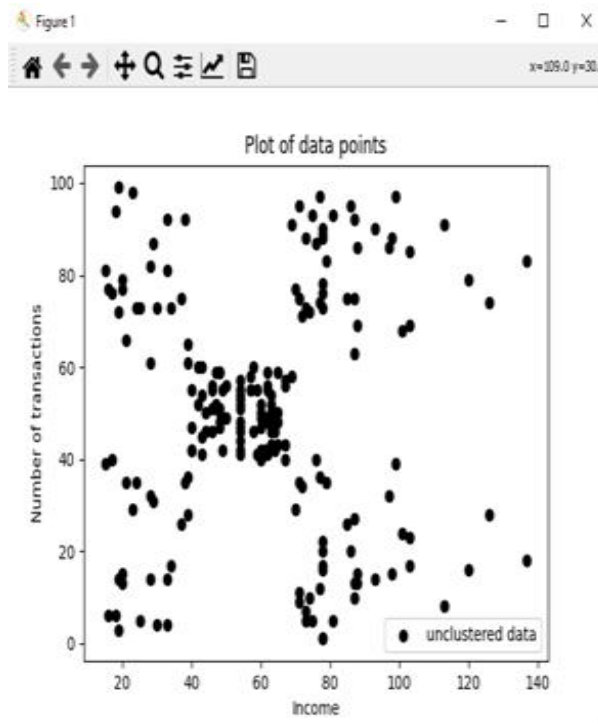
To create the necessary Graphical plots, this project employs the matplotlib.pyplot tool. Matplotlib and its computational mathematics extension NumPy are used to behave like MATLAB. Each individual pyplot function makes some change to the figure: e.g., creates a figure, creates a plotting area, it plots some lines in a plotting area, decorates the plot with labels, etc.

In matplotlib.pyplot various states are preserved across function calls, so that it keeps track of things like the current figure and plotting area, and the plotting functions are directed to the current axes. Generating visualizations with pyplot is a very quick process.

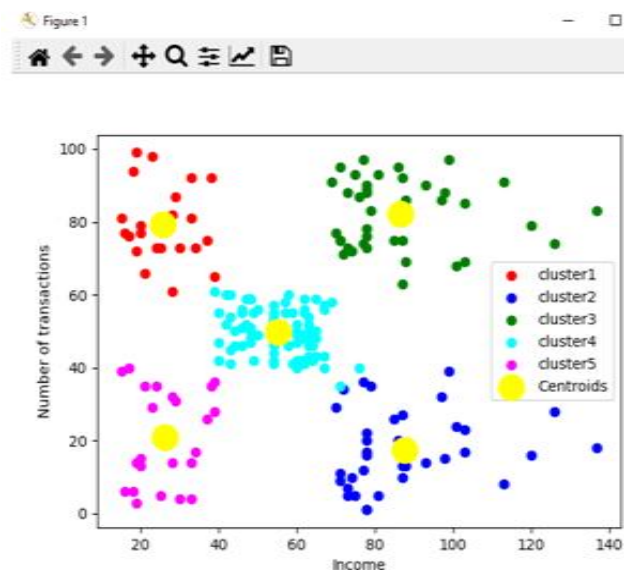
plotting libraries used in Python. It offers an object-oriented API for embedding plots into applications using GUI toolkits such as Tkinter, wxPython, and Qt. Pyplot is a Matplotlib module that aids in the creation of a MATLAB-like user interface. It is intended to be used in the same way as MATLAB, but with the added benefit of being free and open-source. matplotlib.pyplot is a set of command-style functions that make the matplotlib that we're

RESULTS AND DISCUSSIONS

UNCLUSTERED DATA PLOT



CLUSTERED DATA PLOT WITH CENTROID



CONCLUSION

This project is entitled “Data Compression and Clustering through Graphical Plots generation.” is useful to the data analysts for clustering the data and to identify the centroids of those clusters. The project gives clear and colorful pictures of the clustering process, so that the clustering concepts can be easily understood by the beginners. This project finally leads to the improvement of quality of the clustering process.

FUTURE SCOPE

As of now, the clustering process is performed by using the clustering module of the matplotlib. The feasibility of carrying out this project, without using this module needs to be further explored.

REFERENCES

1. M. Kreuseler, T. Nocke, H. Schumann: Integration of Cluster Analysis and Visualization Techniques for Visual Data Analysis
2. AJongwan Kim: Grid-Based Spatial Data Compression Scheme for Clustering in Artificial Intelligence, Indian Journal of Computer Science and Engineering
3. <https://www.python.org/> [4] <https://github.com/baoboa/pyqt5/blob/master/pyuic/uic/pyuic.py>
4. H. Liu, D. Orban, in 2008 Eighth IEEE International Symposium on Cluster Computing and the Grid (CCGRID) (IEEE, 2008), pp. 295–305
5. X. Huang, T. Hu, C. Ye, G. Xu, X. Wang, L. Chen, Electric load data compression and classification based on deep stacked auto-encoders. *Energies* 12(4), 653 (2019)
6. B. Ghose, Z. Rehena, A mechanism for air health monitoring in smart city using context aware computing. *Procedia Computer Science* 171, 2512–2521 (2020)
7. B. Ghose, Z. Rehena, in 2022 International Conference on IoT and Blockchain Technology (ICIBT) (IEEE, 2022), pp. 1–6
8. A. Shehabi, S. Smith, D. Sartor, R. Brown, M. Herrlin, J. Koomey, E. Masanet, N. Horner, I. Azevedo, W. Lintner, United states data center energy usage report (2016)
9. J. Zhang, Real-time lossless compression of soc trace data (2015)
10. F. Marcelloni, M. Vecchio, An efficient lossless compression algorithm for tiny nodes of monitoring wireless sensor networks. *the computer journal* 52(8), 969–987 (2009)