

A suggested approach for data mining in the education sector

Dr. Vyankat Vishnupant Munde

Department of Computer Science Vaidyanath college Parli-V.

Abstract

Education-related Data Mining Data mining is essential in many different sectors. It takes time and a high degree of precision to access a lot of data. The potential impact of data mining on student learning outcomes and procedures was recognized in higher education. Given that practically all educational institutions, both public and private, have thousands of records from students enrolled in a wide range of programs and courses, this is especially true in the sphere of education. Gaining an understanding of the advantages of data retrieval can help the educational process. Developing a student-focused approach and giving institutions the right tools to employ for quality improvement are two benefits of using data mining in education. In this paper, we will find out the benefits of applying data mining in the education sector using classification, prediction, association and clustering methods.

Keywords: Data Mining, Education, Data, Data Mining Process

Introduction

Data mining methods are becoming more and more important in the field of education. One technique that may improve the effectiveness of a certain area of education that is more concerned with digitally preserving crucial information is data mining. Keeping their information in a database will ensure security. for their personal data, which might be more susceptible if it takes the shape of tangible material. Additionally, data mining may identify, categorize, and even forecast the data's future potential, enabling an organization to take proactive measures and make decisions. EDM and educational data mining are two distinct fields used to illustrate the use and use of data mining in the education sector.

They develop a system that can continually gather, process, report, and operate on digital data in order to enhance the educational process. Numerous problems in the subject of education that often take a long time to answer can be answered using data mining technologies [1]. The EDM process converts unprocessed data from educational systems into information that may be used to inform research and study in the field of education. Data mining makes it feasible to look for forecasted data that specialists can misunderstand for a variety of reasons beyond of their control. Particularly in higher education, data mining can forecast a student's likelihood of failing or graduating.

Information from data mining itself may be used by a variety of institutions to concentrate on raising the performance of students who are most likely to fail. Nonetheless, data mining's use in education is still relatively new and requires further study. The use of data mining in the field of education is identified in this research. Every data mining technique used by different researchers will be described and contrasted in this study. This paper's primary goal is to Determine the advantages of data mining for the education industry, what can be accomplished with its assistance, and the best and most practical ways to use it.

Literature Review

Big data can generally refer to datasets that traditional IT and software tools could not perceive, collect, manage, and process. Big data has been used in many sectors, particularly business and education. In 2010, Apache Hadoop defined big data as "datasets that could not be captured, managed, processed within an acceptable scope by general computers." In May 2011, the global consultancy McKinsey & Company announced that big data was the next frontier for innovation, competition, and productivity. This definition has two meanings: first, the volumes of data sets that conform to the big data standard are changing and may increase over time or with technological advancements; second, the volumes of datasets that conform to the big data standard vary from one another in different applications.

Big data refers to datasets that are not accessible, archivable, or manageable with conventional database software [2]. The Big Data phenomena began in the 2000s when Doug Laney, a group of industry experts, presented the idea of Big Data, which is divided into three key components known as 3V [3]. Volume: Users get information from a variety of sources, including social media, business transactions, and machine or system data. These kinds of tasks are challenging to do due to the lack of sophisticated data collecting technology, but they are no longer an issue with the help of technologies like Hadoop, Spark, Google BigQuery, etc. speed,

the volume of input data that requires processing. Think of the hundreds of thousands or perhaps millions of data points that telecommunications providers create every hour of the day in the form of messages, information, data updates, and payment activities.

It is essential to use hardware and software to handle data flow correctly and swiftly. Variety: The forms and types of data that are gathered are often diverse. As more and more data becomes digital, one of the most exciting technical advancements is variety. Numerical data (date, cost, price), documents (text, fax), emails, audio, video, medical records, financial transactions, and so forth are examples of structured data. Many industries that require to store a lot of crucial data, like banking, telecommunications, healthcare, infrastructure, etc., have embraced data mining.

Universities and other educational institutions typically adopt CRISP-DM (Cross-Industry Standard Process Data Mining) as a standard process. There are six phases in the CRISP-DM process: business comprehension, data comprehension, data preparation, modeling, evaluation, and dissemination.

1) Business understanding

Setting company objectives, evaluating the present state, identifying data mining targets, and creating project plans while maintaining focus on the objectives should be the initial steps at this point.

2) Understanding data

The next stage after deciding on company objectives is to gather and describe data. Gathering and describing data from several sources is crucial. Next, assess the quality of the data to determine whether an issue exists.

3) Data preparations

During this stage, data that has been collected need to be cleaned from data that is not needed. Data that has

been cleaned will be sent to the Modeling and building new data stage. In this phase, the obtained data must be filtered to remove unnecessary information. After cleaning, the data will go on to the modeling and new data construction step.

4) Modeling

At this stage, data that has been prepared will be evaluated, so that it can be used to design the test.

5) Evaluations

After determining tools and techniques, results of the modeling process will be evaluated and analyzed

according to business objectives. The impact from the model is needed to understand and review data before

it will be disseminated.

3. Research Method

It is crucial to comprehend the fundamental ideas of data mining in order to comprehend how it might improve educational performance. Classification, categorization, estimation, and visualization are the four key techniques in data mining [10]. Classification is used to separate subjects for each student and to find class differences and relationships. The function of categorization is to manage the outcomes of algorithmic induction categories like "Move" or "Stay," "Survive" or "Expelled." Predictive functions and continuous outcome variables, like GPA or value and pay levels, are related to estimates. Visualization is far more advanced than bar charts or pie charts and employs dynamic images to quantitatively illustrate rules and values. Higher education institutions can use classification methods, to analyze student characteristics, or use estimates to predict the likelihood of various outcomes such as perseverance, performance in an area and the level of graduation of a course. There are many methods in data mining that are used in various industries. Predictions, Classifications,

Associations and Clustering are the most common in the education sector.

3.1. Classification

The process of identifying a group of models that characterize and differentiate the class or notion of a data set is known as classification [11]. The outcome can be shown in a number of ways, including neural networks, decision trees, and categorization algorithms.

The class label of the data may then be predicted using the model. Re-predicting part of the missing data in comparison to the class label is still required in the majority of situations. [12] claimed that categorization techniques can aid in increasing the effectiveness of the higher education system by precisely forecasting the total grades of the students enrolled in a certain course. This method entails: Examining students' enthusiasm for the educational process, Finding pupils that lack motivation, interaction between the educational activities and the pupils, Evaluating if a student has completed an assignment, Continuing evaluate learning performance of students, Examining levels of participation to avoid students dropout from e-learning courses.

To sum up, data mining categorization may be used to group pupils according to their grades, accomplishments, and knowledge, as well as those who lack motivation. In addition to offering various guidelines for the higher education system, classification is utilized to improve the effectiveness and quality of learning activities.

It has been demonstrated that classification would give decision-makers greater latitude in assessing the performance and conduct of a group of students in order to determine how certain group members can perform well in a learning process, even if their specific knowledge or skills do not align with the task.

These techniques can be used effectively to provide early support in the form of educational assistance, in particular to encourage students who are expected to perform unsuccessfully in a particular activity or class, and to effectively calculate the pro and con responses that make up the efficiency of a classification pattern.

3.2. Predictions

The goal of prediction is to create a model that can infer one element (the predicted variable) from a number of other aspects of the data (the predictor variable). A label for the output variable for a small collection of data is necessary for prediction; this label indicates some trustworthy "Basic Trust" information regarding the output variable's value in Specific_scenario.

A few of the primary applications of higher education prediction include forecasting students' grades, conduct, performance, knowledge, and skills. Nevertheless, whether these labels are estimations or not entirely reliable, it is sometimes necessary to reevaluate them. In order to provide details about the underlying structure, prediction can be utilized to determine which model elements are crucial for predictions.

Without initially forecasting intermediate elements, the study program's general strategy can forecast students' educational results. It may be anticipated that prediction can be utilized successfully for prediction purposes, much as classification approaches. Nonetheless, the predictor in a classification is the continuous task in prediction or the categorical task, which is a numerical value. Because of this, researchers often employ a variety of prediction

methodologies to forecast students' academic performance and to pinpoint factors that might indicate whether they would succeed or fail in college courses.

3.3. Association

Association method in data mining is to find a collection of variables in a database repeatedly. Association analysis is the discovery of association rules that indicate the condition of attribute values that often occur together repeatedly in a particular data set. The association method is 'find if-then rules' which means that if the value of one variable is found, the value of another variable will have a certain value [13]. This is because it help teachers to more efficiently evaluate learning patterns for students and coordinate the course content. It could also be used to provide feedback to support decision making by teachers, suggest learning content based on student access history, encourage collaborative learning, recognize irregular learning patterns, evaluate student performance and predict student grades. For example, students who choose networking can also choose computers as a specialization. In addition, students who choose business courses will also choose the MBA program.

In short, the association method can be used to open a new tertiary institution, promote new courses and specializations based on some rules. Association is used to define interactions between the behaviors, learning materials and performance disparity characteristics of students.

3.4. Clustering

It means grouping similar objects or clustering. Rajshree et al described clustering as the process of grouping a series of physical or abstract objects into the same object class. The grouping is not to predict, classify, or estimate target variables number but segment entire data inside a homogeneous subgroup [14]. Furthermore, the task of grouping in the education sector especially universities are based on registration, student transfer, re-admission, selected course, gender, specialization and student behavior. Clustering in higher education is mainly used to support the interaction of students in different learning situations, to recommend similar users' activities and resources, to examine student performance and participation in the education process, to find groups of students with similar learning characteristics based on their behavior knowledge and skills.

4. Data Mining Techniques

There are many techniques discussed in order to understand data mining. [15] indicate 7 types of data mining models: Association, Classification, Clustering, Forecasting, Regression, Sequence Discovery and Visualization. Data mining techniques are used to obtain useful or important information in a particular sector. Of the many data mining techniques in the education sector, there are two very important techniques: decision tree and neural network.

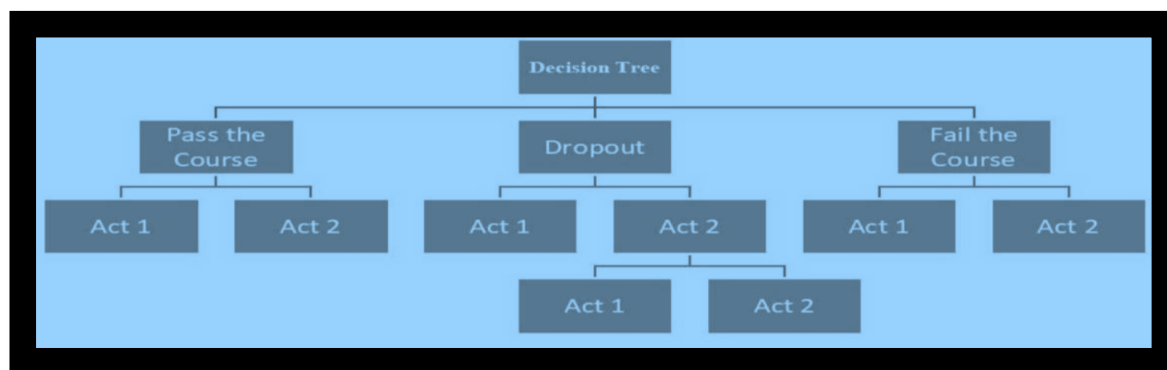


Figure. 1. Decision tree concept

Decision trees illustrated on figure 1 above, are data mining techniques that can be used to classify and predict big data. Decision tree is used to create a customer profile. In addition, decision trees make classifications that are easy to understand and also results-oriented. Many sectors use it to predict and classify customer behavior, release and retention. Likewise, the education sector can also use this technique to predict or classify student performance and behavior.

4.1. Neural Networks

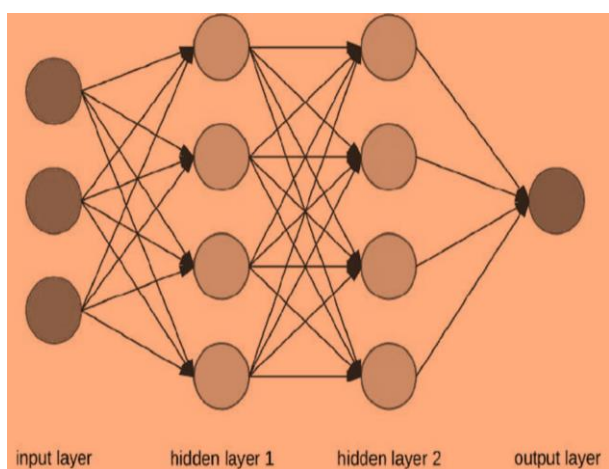


Figure. 2. Neural network concept

Neural network Figure 2, also known as a distributed parallel processing network, is a computing paradigm that is broadly styled on neuronal brain structures. It consists of interconnected processing elements, called nodes or neurons, which work together to create a function of output. This is a technique that can be used to classify large complex data. It is usually used to learn student selection of courses, determine student satisfaction on course and their grades, and decide their specialization selection. Data input is also represented by a neurons that connected to prototype neurons. Each of these connections has a weight that learned the ability to adapt while learning.

5. Data Mining Methods in the Education Sector

There are many ways to apply data mining in education. But some of them have very helpful advantages, such as:

1. Prediction of Student Registration

Aksenova et al [16] developed a predictive model for new, current and returning students at the undergraduate and graduate level. This model builds on the region's population, the unemployment rate in the region, the tuition fees of an institution, household income, and the recording of the institution's past registration data. Data is mined or obtained by help of Cubist tools. The conclusion is data mining is applied and has a very big impact on higher education.

2. Curriculum Development

According to [17] predicting completion rates, study program preferences, and professional registrants used data mining algorithms such as decision trees, decision forests and link analysis. The institution can find correlations between course categories and applicant's profession or occupation. They stressed on how important data mining in developing marketing and curriculum in the higher education. This helps an institute improve the quality of their registrants and meet the institute's business targets themselves.

3. Subject Completion

The university can arrange students according to their loyalty, degree of complaints, and course satisfaction in order to comprehend how they tend to finish their courses.

4. Targeting Students

Woo and associates [18] kids are targeted when they are "the process of building strategies against specific students." According to them, the customer map is a visualization technique used to reach pupils. Customer maps aid in the development of student-focused tactics. This is a "new method for identifying the ideal target student with their needs, values, and character." This is predicated on three aspects of consumer targeting: customer attributes (demographic and psychological), customer values (use and behavior), and customer wants (complaints and satisfaction).

5. Student Course Selection

Factors such as student workload, student characteristics, grades, type of course, course duration, final exam and student needs can influence student choice in subjects using neural networks [19]. These factors act as input for neural network modeling.

6. Teacher's performance in teaching

The instructor's attitude, employee status, student attendance, and student feedback are factors that can influence the teacher or instructor teaching performance at university by using stepwise regression and decision trees in data mining techniques [20].

7. Student performance

Student performance can be determined based on student IDs and grades obtained in the course using datamining techniques, specifically classification techniques like Decision trees and

Naïve Bayes [21][22]. Even the data mining process can also be used for classify teachers performance which helps in improving education system.

6. Discussion

Learning institutions may be able to make wise decisions, provide better advanced planning to support students, more precisely predict possible patterns and behaviors, and allocate staff and resources more precisely thanks to the benefits obtained from various data mining techniques. In our opinion, data mining may be used in education to identify trends, forecast student behaviors and academic success, and enhance the learning experience and grades of students. based on the classifications, predictions, associations, and clustering techniques that we have covered. Every technique has a certain function and usefulness. The instructor can use classification to group pupils according to their grades, accomplishments, or even those who are less motivated.

Information, data, or values can be predicted using prediction techniques. Generally speaking, researchers utilize it to forecast students' final scores from a given course or their passing percentages. In order to determine the correlations between two variables, associations are employed. For instance, a student enrolled in a computer course will also select a programming course. This approach is frequently used by academics to assess student performance and spot anomalous learning patterns. Lastly, pupils with particular patterns can be grouped using the Clustering approach. This is helpful for spotting undesirable habits or forecasting how well certain student groups would learn.

7. Conclusion

Prediction techniques can be used to forecast values, information, or data. In general, researchers use it to predict students' passing percentages or final grades in a particular course. Associations are used to ascertain the relationships between two variables. A student taking a computer course, for example, will also choose to take a programming course. Scholars regularly employ this method to evaluate student performance and identify unusual learning trends. Finally, the Clustering technique may be used to group students who exhibit specific patterns. This is useful for identifying bad behaviors or predicting the learning outcomes of particular student groupings.

References

- [1] A. A. Kardan, H. Sadeghi, S. S. Ghidary, and M. R. F. Sani, "Prediction of student course selection in online higher education institutes using neural network," *Comput. Educ.*, vol. 65, pp. 1–11, 2013, doi: 10.1016/j.compedu.2013.01.015.
- [2] L. M. Caesarius and J. Hohenthal, "Searching for big data: How incumbents explore a possible adoption of big data technologies," *Scand. J. Manag.*, vol. 34, no. 2, pp. 129–140, 2018, doi: 10.1016/j.scaman.2017.12.002.

- [3] M. S. Mahmud, J. Z. Huang, S. Salloum, T. Z. Emara, and K. Sadatdiynov, "A survey of data partitioning and sampling methods to support big data analysis," *Big Data Min. Anal.*, vol. 3, no. 2, pp. 85–101, 2020, doi:10.26599/bdma.2019.9020015.
- [4] X. Chen and M. Xie, "A split-and-conquer approach for analysis of extraordinarily large data," *Statistica Sinica*, vol. 24, no.4, pp. 1655–1684, 2014.
- [5] Y. Baashar et al., "Customer relationship management systems (CRMS) in the healthcare environment: A systematic literature review," *Comput. Stand. Interfaces*, vol. 71, no. March, p. 103442, 2020, doi: 10.1016/j.csi.2020.103442.
- [6] C. da S. Chagas, H. S. K. Pinheiro, W. de Carvalho Junior, L. H. C. dos Anjos, N. R. Pereira, and S. B. Bhering, "Data mining methods applied to map soil units on tropical hillslopes in Rio de Janeiro, Brazil," *Geoderma Reg.*, vol. 9, pp. 47–55, 2017, doi: 10.1016/j.geodrs.2017.03.004.
- [7] J. S. Challa, P. Goyal, S. Nikhil, A. Mangla, S. S. Balasubramaniam, and N. Goyal, "DD-Rtree: A dynamic distributed data structure for efficient data distribution among cluster nodes for spatial data mining algorithms," *Proc. - 2016 IEEE Int. Conf. Big Data, Big Data 2016*, pp. 27–36, 2016, doi: 10.1109/BigData.2016.7840586.
- [8] C. S. Ishikiriya, D. Miro, and C. F. S. Gomes, "Text mining business intelligence: A small sample of what words can say," *Procedia Comput. Sci.*, vol. 55, no. Itqm, pp. 261–267, 2015, doi: 10.1016/j.procs.2015.07.044.
- [9] J. R. Saura, "Using Data Sciences in Digital Marketing: Framework, methods, and performance metrics," *J. Innov. Knowl.*, 2020, doi: 10.1016/j.jik.2020.08.001.
- [10] G. M. Borkar, L. H. Patil, D. Dalgade, and A. Hutke, "A novel clustering approach and adaptive SVM classifier for intrusion detection in WSN: A data mining concept," *Sustain. Comput. Informatics Syst.*, vol. 23, pp. 120–135, 2019, doi: 10.1016/j.suscom.2019.06.002.
- [11] Y. Sato, K. Izui, T. Yamada, and S. Nishiwaki, "Data mining based on clustering and association rule analysis for knowledge discovery in multiobjective topology optimization," *Expert Syst. Appl.*, vol. 119, pp. 247–261, 2019, doi:10.1016/j.eswa.2018.10.047.
- [12] L. Siguenza-Guzman, V. Saquicela, E. Avila-Ordóñez, J. Vandewalle, and D. Cattrysse, "Literature Review of Data Mining Applications in Academic Libraries," *J. Acad. Librariansh.*, vol. 41, no. 4, pp. 499–510, 2015, doi:10.1016/j.acalib.2015.06.007.
- [13] Siemens, George, and Phil Long. "Penetrating the fog: Analytics in learning and education." *EDUCAUSE review*, vol. 46, no. 5, pp. 30, 2011. doi: 10.17471/2499-4324/195.
- [14] S. Natek and M. Zwilling, "Student data mining solution-knowledge management system related to higher education institutions," *Expert Syst. Appl.*, vol. 41, no. 14, pp. 6400–6407, 2014, doi: 10.1016/j.eswa.2014.04.024.

[15] Q. Liu, F. Xiao, and Z. Zhao, "Grouting knowledge discovery based on data mining," *Tunn. Undergr. Sp. Technol.*, vol. 95, no. December 2018, p. 103093, 2020, doi: 10.1016/j.tust.2019.103093.

[16] Ngai, Eric WT, Li Xiu, and Dorothy CK Chau. "Application of data mining techniques in customer relationship management: A literature review and classification." *Expert systems with applications*, vol. 36, no. 2, pp. 2592- 2602, 2009. doi: 10.1016/j.eswa.2008.02.021

[17] Aksenova, Svetlana S., Du Zhang, and Meiliu Lu. "Enrollment prediction through data mining." 2006 IEEE International Conference on Information Reuse & Integration. IEEE, pp. 510-515, 2006. doi: 10.1109/IRI.2006.252466

[18] T. C. Hsia, A. J. Shie, and L. C. Chen, "Course planning of extension education to meet market demand by using data mining techniques - an example of Chinkuo technology university in Taiwan," *Expert Syst. Appl.*, vol. 34, no. 1, pp. 596–602, 2008, doi: 10.1016/j.eswa.2006.09.025.

[19] Woo, Ji Young, Sung Min Bae, and Sang Chan Park. "Visualization method for customer targeting using customer map." *Expert Systems with Applications*, vol. 28, no. 4, pp. 763-772, 2005. doi: 10.1016/j.eswa.2004.12.041

[20] Kardan, Ahmad A., et al. "Prediction of student course selection in online higher education institutes using neural network." *Computers & Education*, Vol. 65, pp. 1-11, 2013. doi: 10.1016/j.compedu.2013.01.015

[21] Mardikyan, Sona, and Bertain Badur. "Analyzing Teaching Performance of Instructors Using Data Mining Techniques." *Informatics in Education.*, vol. 10, no. 2, pp. 245-257, 2011. doi: 10.15388/infedu.2011.17

[22] N. Tomasevic, N. Gvozdenovic, and S. Vranes, "An overview and comparison of supervised data mining techniques for student exam performance prediction," *Comput. Educ.*, vol. 143, no. February 2019, p. 103676, 2020, doi:10.1016/j.compedu.2019.103676.

[23] V. Carter, "Do media influence learning? Revisiting the debate in the context of distance education," *Open Learn.*, vol. 11, no. 1, pp. 31–40, 1996, doi: 10.1080/0268051960110104.

[24] W. Doherty, "An analysis of multiple factors affecting retention in Web-based community college courses," *Internet High.Educ.*, vol. 9, no. 4, pp. 245–255, 2006, doi: 10.1016/j.iheduc.2006.08.004.

[25] R. Damadian, "Abnormal phosphorus metabolism in a potassium transport mutant of *Escherichia coli*," *BBA - Biomembr.*, vol. 135, no. 2, pp. 378–380, 1967, doi: 10.1016/0005-2736(67)90137-X.

[26] B. Holder, "An investigation of hope, academics, environment, and motivation as predictors of persistence in higher education online programs," *Internet High. Educ.*, vol. 10, no. 4, pp. 245–260, 2007, doi: 10.1016/j.iheduc.2007.08.002.

[27] M. Xenos, "Prediction and assessment of student behaviour in open and distance education in computers using Bayesian networks," *Comput. Educ.*, vol. 43, no. 4, pp. 345–359, 2004, doi: 10.1016/j.compedu.2003.09.005.

[28] J. Culpeper and M. Gillings, "Pragmatics: Data trends," *J. Pragmat.*, vol. 145, pp. 4–14, 2019, doi:10.1016/j.pragma.2019.01.004.

[29] M. F. Jefferson, N. Pendleton, S. Lucas, M. A. Horan, and L. Tarassenko, "Neural networks," *Lancet*, vol. 346, no.8991–8892, p. 1712, 1995, doi: 10.1016/S0140-6736(95)92880-4.

[30] N. H. Farhat, "Photonit neural networks and learning mathines the role of electron-trapping materials," *IEEE Expert. Syst.their Appl.*, vol. 7, no. 5, pp. 63–72, 1992, doi: 10.1109/64.163674.

[31] C. K. Loo and M. V. C. Rao, "Accurate and reliable diagnosis and classification using probabilistic ensemble simplified fuzzy ARTMAP," *IEEE Trans. Knowl. Data Eng.*, vol. 17, no. 11, pp. 1589–1593, 2005, doi: 10.1109/TKDE.2005.173.

[32] X. Lin, S. Yacoub, J. Burns, and S. Simske, "Performance analysis of pattern classifier combination by plurality voting," *Pattern Recognit. Lett.*, vol. 24, no. 12, pp. 1959–1969, 2003, doi: 10.1016/S0167-8655(03)00035-7.