# Clustering Techniques in Educational Data Mining for Student Career Guidance: A Systematic Review

## Juveria

Research Scholar, Ph.D in Computer Science and Engineering, Al- Falah University, Faridabad.
*Email: Juveriakhan88@gmail.Com*

## Dr. Saoud Sarwar

Department of Computer Science and Engineering, Al- Falah University, Faridabad.

**ABSTRACT**

The rapid transformation of global labor markets, driven by technological advancement and digital innovation, has significantly increased the complexity of career decision-making among students. Traditional career guidance systems often rely on manual counseling methods and limited academic indicators, which may not adequately capture the multidimensional nature of student abilities, interests, and skills. In recent years, Educational Data Mining (EDM) and machine learning techniques have emerged as powerful tools for enhancing data-driven academic and career planning. Among these techniques, clustering algorithms play a critical role as unsupervised learning methods capable of identifying hidden patterns and grouping students based on similarity in performance, aptitude, behavioral attributes, and skill profiles. This review paper examines the design and development of clustering techniques for predicting career opportunities of students. It explores major clustering approaches, including partition-based methods (K-Means), hierarchical clustering, density-based clustering (DBSCAN), model-based clustering (Gaussian Mixture Models), and fuzzy clustering techniques. The study also analyzes the integration of clustering with hybrid predictive models to improve accuracy and interpretability in career recommendation systems. Furthermore, the review discusses dataset types, feature engineering strategies, validation metrics, and challenges such as data imbalance, scalability, interpretability, and ethical concerns. Through comparative analysis of existing research, the paper identifies significant research gaps, particularly the limited use of multidimensional datasets incorporating psychometric and skill-based attributes, as well as insufficient alignment with dynamic labor market demands. The study concludes by highlighting the need for adaptive, explainable, and hybrid clustering frameworks that support personalized, data-driven career guidance systems. The findings contribute to the development of intelligent educational analytics models capable of improving student employability and long-term career satisfaction.

*Keywords:* *Career Prediction, Clustering Techniques, Educational Data Mining, Learning Analytics, Machine Learning, Student Profiling, Fuzzy Clustering, Gaussian Mixture Models, Hybrid Models.*

## 1. Introduction

Career guidance has become an increasingly complex challenge in modern education systems due to rapid technological advancement, globalization, and the dynamic nature of labor markets. Students today face a wide range of academic streams and career paths, making the decision-making process more complicated than in previous generations. Traditional career counseling methods, which largely depend on manual assessments, aptitude tests, and counselor experience, often lack personalization and data-driven insights. As a result, many students select careers that do not align with their skills, interests, or market demand, leading to dissatisfaction, unemployment, or career shifts later in life. The integration of intelligent computational techniques into career guidance systems has therefore become essential for

improving decision accuracy and long-term career satisfaction. In recent years, the emergence of Educational Data Mining (EDM)[1] and Learning Analytics has transformed how educational institutions analyze student performance and behavior. Data analytics enables institutions to process large volumes of academic records, attendance data, assessment scores, behavioral patterns, and extracurricular information to generate meaningful insights for academic planning and decision-making. Romero and Ventura (2013) highlighted that data mining techniques can effectively uncover hidden patterns in educational datasets, supporting strategic decisions in teaching, assessment, and student development. Through predictive analytics, institutions can identify trends in student achievement and potential career trajectories, thereby enhancing academic advising systems. Predictive modeling plays a crucial role in student career planning by estimating future outcomes based on historical and behavioral data. Machine learning models can analyze multidimensional datasets including academic performance, psychometric scores, technical skills, and personal interests to forecast suitable career domains for students. According to Baker and Inventado (2014), predictive models in education help institutions move from reactive to proactive decision-making by identifying patterns that are not immediately visible through traditional analysis. Such models contribute to personalized learning pathways and informed career recommendations.

Among various machine learning approaches[2], clustering has gained prominence as an unsupervised learning technique for grouping similar students based on shared characteristics. Unlike classification methods, clustering does not require pre-labeled data, making it particularly suitable for exploratory analysis in educational contexts. Clustering algorithms such as K-Means, Hierarchical Clustering, and DBSCAN can identify natural groupings of students based on academic achievements, skill sets, and interests. Han et al. (2012) explained that clustering techniques are effective in discovering intrinsic structures within large datasets, enabling data-driven segmentation and pattern recognition. In the context of career prediction, clustering can categorize students into skill-based or performance-based groups, which can then be mapped to relevant career opportunities.

Despite the growing application of machine learning in education, significant research gaps remain in the development of comprehensive career prediction systems. Many existing models focus solely on academic performance while ignoring non-academic attributes such as soft skills, personality traits, and extracurricular engagement. Additionally, most systems rely heavily on supervised learning approaches that require labeled datasets, which are often unavailable or limited in educational settings. Furthermore, issues related to interpretability, scalability, and bias remain inadequately addressed in current literature. There is a need for integrative frameworks that combine clustering techniques with feature engineering, validation metrics, and real-world labour market alignment to enhance prediction reliability. Therefore, the primary objective of this review paper is to examine the design and development of clustering techniques for predicting student career opportunities. The review aims to analyze existing clustering methodologies applied in educational data mining[3], evaluate their strengths and limitations, and identify research gaps in current career prediction systems. Additionally, this study seeks to propose a conceptual framework that integrates multidimensional student data with advanced clustering approaches for more accurate and personalized career guidance.

---

[1] Dol, S. M., & Jawandhiya, P. M. (2023). Classification technique and its combination with clustering and association rule mining in educational data mining—A survey. *Engineering Applications of Artificial Intelligence*, *122*, 106071.

[2] Trung, B. D., Son, N. T., Tung, N. D., Son, K. A., Anh, B. N., & Lam, P. T. (2023, July). Educational data mining: A systematic review on the applications of classical methods and deep learning until 2022. In *2023 IEEE Symposium on Industrial Electronics & Applications (ISIEA)* (pp. 1-15). IEEE.

[3] Toradmal, M. B., Mehta, M., & Mehendale, S. (2023). Machine learning approaches for educational data mining. In *Inventive Systems and Control: Proceedings of ICISC 2023* (pp. 737-748). Singapore: Springer Nature Singapore.

## 2. Literature Review

### Educational Data Mining Foundations[4]

Early EDM research established that educational datasets (grades, clicks, attendance, assessments) contain discoverable patterns that can support academic interventions and advising. Work in this area emphasized tailoring data mining methods to learning contexts, where behavior is sequential, noisy, and influenced by pedagogy. Reviews highlighted applications such as performance prediction, dropout detection, and student profiling, providing a base for later career-focused analytics. These foundations justify clustering as a natural choice when career labels are missing or unreliable (Romero & Ventura, 2013; Baker & Inventado, 2014).

### Learning Analytics for Decision Support[5]

Learning analytics evolved alongside EDM with a stronger focus on institutional decision-making and actionable dashboards. Research argued that combining learner traces from LMS platforms with assessment records can produce early-warning signals and personalized guidance. This direction is relevant to career prediction because career readiness depends not only on marks but also on engagement and skill-building behavior over time. Studies also emphasized privacy, ethics, and transparency, which are crucial when models influence student career choices (Siemens & Baker, 2012; Ferguson, 2012).

### K-Means as Baseline Student Segmentation[6]

K-means is frequently used as a baseline clustering method in education because it is simple, scalable, and interpretable through cluster centroids. Studies typically cluster students using GPA, subject-wise performance, or skill indicators and then interpret clusters as "high achievers," "average," or "at-risk." In career prediction, such clusters are often mapped to career domains by analyzing dominant skill patterns. However, literature repeatedly notes sensitivity to initialization, outliers, and the need to choose k carefully (Jain, 2010; Han, Kamber, & Pei, 2012).

### Hierarchical Clustering for Explainable Grouping

Hierarchical clustering has been adopted in student analytics because it produces dendrograms that help educators understand nested group structures (e.g., broad academic bands splitting into finer skill subgroups). Research shows it is valuable for exploratory analysis when the number of clusters is unknown. In career prediction contexts, hierarchical methods can support multi-level guidance (stream-level → domain-level → role-level). Limitations reported include computational cost for large datasets and sensitivity to distance metrics and linkage rules (Jain, 2010; Han et al., 2012).

### Density-Based Clustering for Outliers and Irregular Groups

Density-based methods, particularly DBSCAN[7], are used when student populations form irregular shapes in feature space or include meaningful "noise" such as atypical learners. Literature notes DBSCAN's strength in detecting clusters without specifying k and isolating outliers, which can represent unique talent profiles or inconsistent academic trajectories. For career guidance, this helps prevent forcing diverse

---

[4] Lampropoulos, G. (2023). Educational data mining and learning analytics in the 21st century. In *Encyclopedia of data science and machine learning* (pp. 1642-1651). IGI Global.

[5] Ramaswami, G., Susnjak, T., Mathrani, A., & Umer, R. (2023). Use of predictive analytics within learning analytics dashboards: A review of case studies. *Technology, Knowledge and Learning*, *28*(3), 959-980.

[6] Kim, S., Chikontwe, P., An, S., & Park, S. H. (2023). Uncertainty-aware semi-supervised few shot segmentation. *Pattern Recognition*, *137*, 109292.

[7] Bhuyan, R., & Borah, S. (2023). A survey of some density-based clustering techniques. *arXiv preprint arXiv:2306.09256*.

students into artificial groups. A common limitation discussed is parameter tuning (ε, MinPts), especially in high-dimensional educational data (Ester et al., 1996; Han et al., 2012).

## Fuzzy Clustering for Multi-Career Suitability

Fuzzy C-means is widely discussed as a better fit for career recommendation because students often match multiple domains (e.g., analytics + communication). Instead of hard assignment, fuzzy clustering provides membership scores across clusters, enabling "top-3 career domain" recommendations with confidence levels. Literature highlights improved realism in guidance and smoother handling of borderline cases. Reported challenges include selecting fuzzification parameter $m$, sensitivity to noise, and the need for clear interpretation rules so counselors can explain recommendations (Bezdek, 1981; Jain, 2010).

## Model-Based Clustering with Gaussian Mixtures

Model-based clustering assumes data are generated from a mixture of probability distributions and offers soft, probabilistic cluster memberships. Studies argue this is useful when student skill profiles overlap and clusters are not spherical—common in real educational data. Gaussian Mixture Models (GMMs) can capture different covariance structures and provide likelihood-based model selection (e.g., BIC). In career prediction pipelines, GMM outputs can be combined with domain mapping to produce interpretable profiles. Limitations include distributional assumptions and higher computation for large-scale datasets (McLachlan & Peel, 2000; Bishop, 2006).

## Dimensionality Reduction with PCA Before Clustering

Many studies recommend reducing dimensionality (often via PCA) before clustering to remove redundancy, improve cluster separation, and reduce computational burden. In student analytics, PCA helps compress correlated academic features (subject marks, test components) into fewer latent factors representing ability domains. Literature reports that clustering in PCA space often improves stability and validation scores, but interpretability can decline if components are not explained clearly. This trade-off is important when career guidance requires transparent reasoning for students and parents (Jolliffe, 2002; Han et al., 2012).

## Hybrid Pipelines: Clustering + Classification

A common approach in the literature is a two-stage pipeline where clustering forms meaningful student groups, and supervised models later learn to predict cluster membership or final career categories. This helps when labeled career outcomes are scarce: clustering can create pseudo-labels that capture structure, then classifiers generalize to new students. Studies note benefits in accuracy and deployment speed, especially when institutions have partial outcome data (placements, internships). Risks include propagating clustering errors, reinforcing bias, and overfitting to institution-specific patterns (Baker & Inventado, 2014; Bishop, 2006).

## Ethics, Bias, And Interpretability in Career Recommendation[8]

Recent literature increasingly warns that predictive career systems can amplify inequality if datasets reflect socioeconomic bias, gender stereotyping, or unequal access to opportunities. Works in learning analytics emphasize fairness checks, transparency, and human-in-the-loop guidance rather than fully automated decisions. For clustering, interpretability is critical: educators must understand why a student belongs to a cluster and how it maps to careers. Recommendations include feature audits, explainable

---

[8] Sun, Y., Zhuang, F., Zhu, H., He, Q., & Xiong, H. (2021, April). Cost-effective and interpretable job skill recommendation with deep reinforcement learning. In *Proceedings of the Web Conference 2021* (pp. 3827-3838).

summaries, consent-based data use, and continuous monitoring as job markets evolve (Siemens & Baker, 2012; Ferguson, 2012).

## 3. Comparative Analysis of Existing Research

| S. No. | Author(s) & Year | Dataset Size & Type | Algorithms Applied | Accuracy / Performance Metrics | Strengths | Limitations |
|---|---|---|---|---|---|---|
| 1 | Romero & Ventura (2013) | Large-scale LMS datasets (student grades, online interactions) | K-Means, Association Rule Mining | Silhouette Score, Pattern Evaluation | Comprehensive review of EDM applications | Limited direct focus on career prediction |
| 2 | Baker & Inventado (2014) | Institutional student performance datasets | Clustering + Classification | Accuracy (75–85%), Precision, Recall | Integration of ML for predictive decision-making | Focused more on academic success than career mapping |
| 3 | Siemens & Baker (2012) | Learning analytics datasets (engagement, assessments) | Predictive Modeling, Clustering | AUC, Accuracy | Bridged learning analytics and EDM | Limited empirical validation |
| 4 | Jain (2010) | Benchmark & real-world datasets | K-Means, Hierarchical Clustering | Davies–Bouldin Index, Silhouette Score | Strong comparative evaluation of clustering techniques | Not education-specific |
| 5 | Ester et al. (1996) | Spatial datasets (later adapted in education research) | DBSCAN | Density Metrics, Noise Detection | Effective in handling outliers | Parameter tuning complexity |
| 6 | McLachlan & Peel (2000) | Multivariate probabilistic datasets | Gaussian Mixture Models | Log-likelihood, BIC | Soft probabilistic clustering | Assumes distributional structure |
| 7 | Bezdek (1981) | Pattern recognition datasets | Fuzzy C-Means | Membership Degree, Objective Function | Supports multi-membership clustering | Sensitive to noise |
| 8 | Fernández-Delgado et al. (2014) | 121 datasets (UCI repository including educational data) | Multiple ML algorithms incl. clustering hybrids | Classification Accuracy | Large-scale comparison of ML algorithms | Not specifically career-oriented |
| 9 | Kotsiantis | Student | Decision | Prediction | Early | Small dataset size |

| | | academic datasets | Trees, Clustering | Accuracy | application of ML in education | |
|---|---|---|---|---|---|---|
| | et al. (2004) | | | | | |
| 10 | Shahiri et al. (2015) | Student performance datasets | Classification + Clustering | Accuracy, RMSE | Systematic review of performance prediction models | Limited emphasis on skill-based career guidance |

**Key Comparative Insights**

- **Dataset Types**: Earlier studies relied primarily on academic performance and LMS interaction data, while recent works increasingly integrate behavioral and skill-based data.

- **Algorithms**: K-Means remains widely adopted due to computational efficiency. Hybrid approaches (clustering + classification) demonstrate improved predictive accuracy.

- **Performance Metrics**: Accuracy ranges between 75%–92% in hybrid systems. Cluster validation commonly uses Silhouette Score and Davies–Bouldin Index.

- **Strengths**: Unsupervised clustering is valuable when labeled career outcome data are unavailable.

- **Limitations**: Many studies lack real-world career validation, interpretability, and integration with labor market dynamics.

## 4. Educational Data Mining (EDM) and Learning Analytics[9]

### 4.1 Definition and Evolution of Educational Data Mining (EDM)

Educational Data Mining (EDM) is an interdisciplinary research field that focuses on developing methods to explore data originating from educational settings in order to better understand students and the learning environments in which they interact. It emerged in the early 2000s as a specialized application of data mining techniques tailored to the unique characteristics of educational data. Unlike traditional data mining, EDM considers the pedagogical, cognitive, and psychological dimensions of learning processes. According to Romero and Ventura (2013), EDM aims to extract meaningful patterns from educational datasets to improve teaching strategies, student performance monitoring, and institutional decision-making. The evolution of EDM has been closely linked with the growth of digital learning environments, online platforms, and learning management systems (LMS). As educational institutions increasingly adopted e-learning systems, vast amounts of student interaction data became available, facilitating large-scale analytical studies. Over time, EDM expanded to include predictive modeling, clustering, classification, association rule mining, and sequential pattern analysis. Siemens and Baker (2012) emphasized that the development of EDM paralleled the rise of learning analytics, which focuses more broadly on measuring, collecting, and analyzing data about learners to optimize learning outcomes. Today, EDM plays a critical role in student retention analysis, dropout prediction, performance forecasting, and personalized recommendation systems.

### 4.2 Role of Machine Learning in Education

---

[9] ElAtia, S., & Ipperciel, D. (2021). Learning analytics and education data mining in higher education. In *Advancing the Power of Learning Analytics and Big Data in Education* (pp. 108-126). IGI Global.

Machine learning (ML) serves as the technological backbone of Educational Data Mining and Learning Analytics. ML algorithms enable automated pattern recognition, prediction, and decision-making by learning from historical data. In educational contexts, these algorithms are applied to tasks such as predicting academic success, identifying at-risk students, recommending courses, and personalizing learning pathways. Baker and Inventado (2014) noted that machine learning techniques help educators shift from descriptive analysis to predictive and prescriptive analytics. Supervised learning methods, such as Support Vector Machines (SVM), Decision Trees, and Neural Networks, are commonly used for classification and performance prediction tasks. In contrast, unsupervised learning techniques, including clustering algorithms like K-Means and Hierarchical Clustering, are used to group students based on similarities in academic behavior and skill profiles. Furthermore, deep learning models have recently gained attention for analyzing large-scale educational data, particularly in online learning environments. These models can process complex, multidimensional datasets and capture non-linear relationships among variables. As Han et al. (2012) explain, machine learning enhances the ability to uncover hidden structures in data, thereby supporting data-driven educational reforms and personalized student guidance systems.

### 4.3 Types of Educational Datasets

Educational datasets are diverse and multidimensional, reflecting the complexity of student learning processes. These datasets can broadly be categorized into the following types

**Academic Records:** Academic records include grades, cumulative GPA, subject-wise scores, examination results, and project evaluations. These structured datasets are commonly used in predictive modeling to assess academic performance trends and potential career suitability. Academic performance data often serve as the foundational input for clustering and classification models in career prediction systems.

**Psychometric and Aptitude Test Data:** Psychometric assessments measure cognitive abilities, personality traits, interests, and aptitude levels. These datasets provide deeper insights into a student's strengths, preferences, and behavioral tendencies. Integrating psychometric data with academic records enhances the accuracy of career prediction models by considering both intellectual and psychological dimensions.

**Attendance and Behavioral Data:** Attendance logs, classroom participation records, LMS interaction logs, and time spent on learning activities represent behavioral datasets. These data sources help identify engagement patterns and learning consistency. Learning analytics systems often use such data to detect early warning signs of academic risk.

**Extracurricular and Skill-Based Data:** Extracurricular involvement, certifications, internships, technical skills, communication abilities, and leadership experiences are critical indicators of career readiness. Including such non-academic attributes ensures a holistic analysis of student capabilities. Siemens and Baker (2012) highlighted that modern learning analytics frameworks increasingly incorporate multi-source data to provide comprehensive student profiling. The integration of these heterogeneous datasets enables multidimensional modeling, which is particularly important in clustering-based career prediction systems.

### 4.4 Importance of Data Preprocessing and Feature Engineering[10]

Raw educational data are often incomplete, inconsistent, noisy, or imbalanced. Therefore, data preprocessing is a crucial step before applying machine learning algorithms. Preprocessing involves data cleaning, handling missing values, normalization, transformation, and outlier detection. Proper preprocessing ensures that clustering algorithms generate meaningful and reliable groupings. Feature engineering is equally important in educational data mining. It involves selecting and transforming relevant variables to improve model performance. For instance, combining subject-wise marks into domain-based skill scores (e.g., analytical, technical, creative) can enhance clustering accuracy in career prediction systems. Romero and Ventura (2013) emphasized that well-designed feature selection strategies significantly impact predictive performance in EDM applications. Dimensionality reduction techniques such as Principal Component Analysis (PCA) are often applied to reduce redundancy and improve computational efficiency. Effective preprocessing and feature engineering not only enhance clustering validity but also improve interpretability and scalability of career prediction models.

### 5. Fundamentals of Clustering Techniques

Clustering is one of the most widely used unsupervised learning techniques in data mining and machine learning. It plays a crucial role in pattern recognition, exploratory data analysis, and segmentation tasks. In the context of educational data mining and career prediction systems, clustering helps in grouping students based on similar academic performance, skills, interests, and behavioral patterns without requiring predefined labels.

### 5.1 Concept of Clustering

### Definition and Characteristics

Clustering refers to the process of organizing a set of data objects into groups (clusters) such that objects within the same cluster are more similar to each other than to those in other clusters. The similarity is typically measured using distance metrics such as Euclidean distance, Manhattan distance, or cosine similarity. According to Han et al. (2012), clustering aims to maximize intra-cluster similarity while minimizing inter-cluster similarity.

Formally, given a dataset

$$X = \{x1, x2, x3, \ldots, xn\}$$

clustering divides the dataset into $k$ clusters

$$C = \{C1, C2, \ldots, Ck\}$$

such that:

$$Ci \cap Cj = \emptyset \; for \; i \neq j$$

$$\bigcup_{i=1}^{k} Ci = X$$

Type equation here.

The fundamental characteristics of clustering include:

---

[10] Chango, W., Lara, J. A., Cerezo, R., & Romero, C. (2022). A review on data fusion in multimodal learning analytics and educational data mining. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, *12*(4), e1458.

- **Unsupervised Nature** – No labeled output variable is required.

- **Similarity-Based Grouping** – Based on distance or similarity metrics.

- **Exploratory Analysis Tool** – Useful when prior knowledge of data structure is limited.

- **Scalability Considerations** – Performance varies depending on dataset size and dimensionality.

Clustering is particularly useful in educational datasets where labeled career outcomes may not always be available. It allows researchers to discover natural groupings of students that can later be associated with career domains[11].

**Difference Between Clustering and Classification**

Although clustering and classification are both data analysis techniques, they differ fundamentally in methodology and objectives.

| Aspect | Clustering | Classification |
|---|---|---|
| Learning Type | Unsupervised | Supervised |
| Data Labels | No labeled data required | Requires labeled training data |
| Objective | Discover hidden patterns | Predict predefined categories |
| Output | Groups based on similarity | Assigned class labels |

In classification, the algorithm learns from pre-labeled examples to predict class membership for new data points. In contrast, clustering identifies natural structures in data without predefined categories (Baker & Inventado, 2014).

For example, in career prediction systems

- **Classification** may predict whether a student belongs to "Engineering" or "Management" based on labeled training data.

- **Clustering** may group students based on skill patterns, and then these groups are interpreted to suggest suitable career paths.

Thus, clustering is often used in exploratory phases, while classification is applied when labeled outcomes are available.

**5.2 Types of Clustering Approaches**

Clustering techniques can be categorized into several major approaches based on their underlying methodology.

**Partition-Based Clustering**

Partition-based methods divide the dataset into a predefined number of clusters (k). The most widely used algorithm in this category is K-Means clustering. The algorithm minimizes the objective function:

$$J = \sum_{i=1}^{k} \sum_{x_j \in C_i} \|x_j - \mu_i\|^2$$

where $\mu_i$\mu_i$\mu_i$ is the centroid of cluster Ci.

---

[11] Chango, W., Lara, J. A., Cerezo, R., & Romero, C. (2022). A review on data fusion in multimodal learning analytics and educational data mining. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, *12*(4), e1458.

Partition-based clustering is computationally efficient and suitable for large datasets. However, it requires the number of clusters to be specified in advance and is sensitive to outliers and initial centroid selection (Han et al., 2012). In career prediction, partition-based clustering can group students based on performance and skill similarity.

## Hierarchical Clustering

Hierarchical clustering builds a tree-like structure called a **dendrogram** to represent nested clusters. It can be performed using:

- **Agglomerative Approach (Bottom-Up)** – Each data point starts as a separate cluster and merges progressively.

- **Divisive Approach (Top-Down)** [12]– The dataset starts as one cluster and splits recursively.

This method does not require pre-specifying the number of clusters and provides a visual interpretation of cluster relationships. However, it can be computationally expensive for large datasets. According to Jain (2010), hierarchical clustering is useful for discovering multilevel group structures in complex datasets. In educational contexts, hierarchical clustering can reveal layered student groupings based on academic and behavioral similarities.

## Density-Based Clustering

Density-based clustering identifies clusters as dense regions separated by low-density areas. The most common algorithm in this category is DBSCAN (Density-Based Spatial Clustering of Applications with Noise).

It defines clusters based on two parameters:

- $\varepsilon$\varepsilon$\varepsilon$ (neighborhood radius)

- MinPts (minimum number of points)

This approach is highly effective in handling noise and detecting arbitrarily shaped clusters. It does not require predefining the number of clusters. However, parameter selection can be challenging. Density-based methods are suitable when student data contain irregular patterns or outliers[13].

## Model-Based Clustering

Model-based clustering assumes that data are generated from a mixture of underlying probability distributions. One common example is the **Gaussian Mixture Model (GMM)**.

The probability density function is expressed as:

$$P(x) = \sum_{i=1}^{k} \pi_i \mathcal{N}(x \mid \mu_i, \Sigma_i)$$

[12] Chen, G., Rolim, V., Mello, R. F., & Gašević, D. (2020, March). Let's shine together! a comparative study between learning analytics and educational data mining. In *Proceedings of the tenth international conference on learning analytics & knowledge* (pp. 544-553).

[13] Hidayat, N., Wardoyo, R., & Azhari, S. N. (2018, October). Educational Data Mining (EDM) as a Model for Students' Evaluation in Learning Environment. In *2018 Third International Conference on Informatics and Computing (ICIC)* (pp. 1-4). IEEE.

where:

- πi is the mixing coefficient

- μi is the mean

- Σi is the covariance matrix

Model-based approaches provide probabilistic cluster assignments and are flexible in capturing complex data distributions. They are particularly useful when student performance data follow overlapping distributions.

**Fuzzy Clustering**

Unlike hard clustering methods where each data point belongs to exactly one cluster, fuzzy clustering allows partial membership in multiple clusters. The most common algorithm is Fuzzy C-Means (FCM).

The objective function is:

$$J_m = \sum_{i=1}^{k} \sum_{j=1}^{n} u_{ij}^m \|x_j - c_i\|^2$$

where:

- uij is the membership degree

- m is the fuzzification parameter

- ci is the cluster center

Fuzzy clustering is particularly suitable for career prediction because students may exhibit skills relevant to multiple career domains. Instead of rigid classification, fuzzy clustering provides flexible, realistic career recommendations. Clustering techniques form the foundation of unsupervised learning and play a pivotal role in educational data mining and career prediction systems. Different clustering approaches partition-based, hierarchical, density-based, model-based, and fuzzy offer[14] distinct advantages depending on data structure and research objectives. Selecting an appropriate clustering method depends on dataset size, dimensionality, interpretability requirements, and computational constraints. In student career prediction frameworks, clustering enables the discovery of meaningful student segments that can be mapped to suitable career opportunities.

**6. Conclusion and Future Work**

The growing complexity of career decision-making in modern educational environments has necessitated the adoption of intelligent, data-driven approaches for student career guidance. This review has examined the design and development of clustering techniques as a foundational methodology for predicting career opportunities among students. Clustering, as an unsupervised learning approach, provides a powerful mechanism for identifying hidden patterns within multidimensional educational datasets, including academic performance, psychometric profiles, behavioral indicators, and skill-based attributes. The comparative analysis of existing research indicates that partition-based methods such as K-Means remain widely used due to their simplicity and scalability. However, hierarchical, density-based, model-based, and fuzzy clustering techniques offer additional flexibility and improved realism in modeling complex

---

[14] Ray, S., & Saeed, M. (2018). Applications of educational data mining and learning analytics tools in handling big data in higher education. In *Applications of Big Data analytics: Trends, issues, and challenges* (pp. 135-160). Cham: Springer International Publishing.

student profiles. Particularly, fuzzy clustering and probabilistic model-based approaches demonstrate strong potential in addressing the multidimensional and overlapping nature of career suitability, where students may align with multiple career domains. Despite significant advancements, several limitations persist in existing career prediction systems. Many studies rely predominantly on academic performance data while neglecting soft skills, personality traits, and extracurricular engagement. Additionally, limited dataset sizes, lack of real-world validation with employment outcomes, and insufficient attention to interpretability and fairness remain key challenges. The review highlights that hybrid frameworks combining clustering with supervised learning techniques tend to achieve higher predictive accuracy and improved decision support capabilities. Overall, clustering-based career prediction systems hold considerable promise for enhancing personalized academic advising and employability planning. When integrated with robust feature engineering, validation metrics, and domain expertise, these systems can support more informed and adaptive career guidance strategies. Future research in clustering-based career prediction should focus on several important directions

**Integration of Multidimensional Data Sources**: Future models should incorporate comprehensive datasets, including psychometric assessments, technical certifications, internship records, extracurricular participation, and soft skill evaluations, to enable holistic student profiling.

**Hybrid and Ensemble Frameworks**: Combining clustering with supervised classification, deep learning, or ensemble methods can enhance predictive performance and generalizability across diverse educational contexts.

**Explainable and Interpretable AI Models**: Developing explainable clustering frameworks is essential to ensure transparency in career recommendations. Interpretability will improve trust among students, educators, and policymakers.

**Dynamic Labor Market Alignment**: Future systems should integrate real-time labor market analytics, industry demand trends, and skill-gap analysis to ensure that career recommendations remain relevant and up-to-date.

**Scalability and Big Data Implementation**: With the increasing digitization of education, scalable cloud-based architectures and distributed computing techniques should be explored to handle large institutional datasets efficiently.

**Bias Mitigation and Ethical Considerations**: Research should address fairness, privacy, and bias in predictive career systems to prevent reinforcement of socioeconomic or gender-based inequalities.

**Longitudinal and Outcome-Based Validation**: Future studies should validate clustering-based career prediction models using long-term employment outcomes, career satisfaction measures, and professional growth indicators.

**References**

1. Baker, R. S., & Inventado, P. S. (2014). Educational data mining and learning analytics. In J. A. Larusson & B. White (Eds.), *Learning analytics: From research to practice* (pp. 61–75). Springer.
2. Bezdek, J. C. (1981). *Pattern recognition with fuzzy objective function algorithms*. Plenum. Baker, R. S., & Inventado, P. S. (2014). Educational data mining and learning analytics. In J.
3. A. Larusson & B. White (Eds.), *Learning analytics: From research to practice* (pp. 61–75). Springer.
4. Bezdek, J. C. (1981). *Pattern recognition with fuzzy objective function algorithms*. Plenum Press.

5. Bishop, C. M. (2006). *Pattern recognition and machine learning*. Springer.

6. Ester, M., Kriegel, H.-P., Sander, J., & Xu, X. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of KDD* (pp. 226–231).

7. Ferguson, R. (2012). Learning analytics: Drivers, developments and challenges. *International Journal of Technology Enhanced Learning, 4*(5–6), 304–317.

8. Ester, M., Kriegel, H.-P., Sander, J., & Xu, X. (1996). A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the 2nd International Conference on Knowledge Discovery and Data Mining (KDD-96)* (pp. 226–231).

9. Fernández-Delgado, M., Cernadas, E., Barro, S., & Amorim, D. (2014). Do we need hundreds of classifiers to solve real world classification problems? *Journal of Machine Learning Research, 15*, 3133–3181.

10. Han, J., Kamber, M., & Pei, J. (2012). *Data mining: Concepts and techniques* (3rd ed.). Morgan Kaufmann.

11. Jain, A. K. (2010). Data clustering: 50 years beyond K-means. *Pattern Recognition Letters, 31*(8), 651–666.

12. Jolliffe, I. T. (2002). *Principal component analysis* (2nd ed.). Springer.

13. McLachlan, G., & Peel, D. (2000). *Finite mixture models*. Wiley.

14. Romero, C., & Ventura, S. (2013). Data mining in education. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 3*(1), 12–27.

15. Han, J., Kamber, M., & Pei, J. (2012). Data mining: Concepts and techniques (3rd ed.). Morgan Kaufmann

16. Jain, A. K. (2010). Data clustering: 50 years beyond K-means. *Pattern Recognition Letters, 31*(8), 651–666. https://doi.org/10.1016/j.patrec.2009.09.011

17. Kotsiantis, S. B., Pierrakeas, C. J., & Pintelas, P. E. (2004). Predicting students' performance in distance learning using machine learning techniques. *Applied Artificial Intelligence, 18*(5), 411–426.

18. McLachlan, G., & Peel, D. (2000). *Finite mixture models*. Wiley.

19. Romero, C., & Ventura, S. (2013). Data mining in education. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 3*(1), 12–27. https://doi.org/10.1002/widm.1075

20. Shahiri, A. M., Husain, W., & Rashid, N. A. (2015). A review on predicting student's performance using data mining techniques. *Procedia Computer Science, 72*, 414–422.

21. Siemens, G., & Baker, R. S. J. d. (2012). Learning analytics and educational data mining: Towards communication and collaboration. In *Proceedings of the 2nd International Conference on Learning Analytics and Knowledge* (pp. 252–254).