# Unveiling Longitudinal Patterns in International Mental Health Assessments: A Systematic Evaluation of Clustering Techniques

**John Makunda., Helen Waititu (PhD)., Cornelius Nyakundi (PhD)**

**Department of Mathematics and Actuarial Science, Catholic University of Eastern Africa**

## ABSTRACT

**Background:** Mental health assessments across diverse populations provide valuable insights into the prevalence and patterns of mental health issues. However, the complexity and volume of longitudinal data present challenges in extracting meaningful information for effective intervention. Clustering methods have emerged as powerful tools for identifying hidden structures within such datasets, yet a comprehensive evaluation of these techniques in the context of international mental health assessments is lacking. **Objectives:** This study aims to systematically evaluate various clustering techniques applied to longitudinal mental health data from international assessments. The focus is on understanding how different methods capture and reveal patterns and subgroups within the data, thereby guiding targeted mental health interventions. **Methods:** We applied and compared three clustering techniques—K-Means Clustering, Hierarchical Clustering, and Gaussian Mixture Models (GMM)—to longitudinal mental health assessment data. We assessed the performance of these methods in identifying meaningful clusters, considering their strengths and limitations in capturing the complexity of mental health trajectories. **Results:** Our analysis revealed distinct clusters reflecting varying levels of mental health severity and symptom trajectories. K-Means identified broad clusters, while Hierarchical Clustering provided insights into the data's hierarchical structure. GMM offered a probabilistic view, highlighting overlapping mental health experiences among individuals. Each method contributed uniquely to understanding the longitudinal patterns in the data. **Implications:** The findings underscore the importance of using a multi-faceted approach to clustering in mental health research. By revealing different dimensions of mental health trajectories, this study provides valuable insights for tailoring interventions and resource allocation. The results highlight the need for ongoing evaluation of clustering techniques to enhance their applicability in diverse international contexts.

**Keywords:** Longitudinal Clustering, Mental Health Assessments, K-Means Clustering, Hierarchical Clustering, Gaussian Mixture Models, International Mental Health

## INTRODUCTION

In recent years, the importance of mental health has garnered increasing attention globally, with international mental health assessments playing a crucial role in understanding the prevalence, patterns, and determinants of mental health issues across diverse populations. These assessments provide a wealth of data, offering insights into the mental health status of individuals and communities. However, the complexity and volume of this data present challenges in extracting meaningful information that can inform public health interventions and policy decisions. Clustering methods have emerged as powerful tools for uncovering hidden structures within large datasets, allowing researchers to identify groups of individuals with similar mental health profiles. These techniques are particularly valuable in the context of international mental health assessments, where patterns and clusters may vary significantly across different cultural, social, and economic contexts. By systematically applying and comparing various clustering methods, researchers can better understand the distinctive patterns within the data, which can, in turn, guide targeted mental health interventions and resource allocation.

Despite the potential of clustering methods, there is a need for a systematic evaluation of their performance in the context of international mental health data. Previous studies have applied different clustering techniques to

mental health datasets, but a comprehensive comparison of these methods, particularly across diverse populations, remains limited. Such an evaluation is essential for identifying the most effective approaches to clustering in this domain and for ensuring that the derived clusters are both interpretable and actionable. This study aims to address this gap by unveiling distinctive patterns and clusters within international mental health assessment data through a systematic assessment and comparison of various clustering methods. The analysis focuses on evaluating the performance of these methods in identifying meaningful clusters, with an emphasis on the practical implications of the findings for mental health research and practice. By providing a robust comparison of clustering techniques, this study contributes to the growing body of literature on mental health data analysis and offers valuable insights for future research and policy development.

## Background

Clustering methods have emerged as essential tools in mental health research, enabling researchers to identify patterns and subgroups within complex datasets. These techniques group individuals based on similarities in symptom profiles, risk factors, or treatment outcomes, thereby offering deeper insights into mental health constructs and informing the development of more effective interventions [1]. The growing utilization of international mental health assessment tools has provided valuable insights into various aspects of psychological well-being. However, a significant gap remains in understanding the intricate temporal dynamics and diverse patterns inherent in mental health trajectories [2]. Traditional analyses often fall short in capturing these complexities, necessitating the adoption of advanced clustering methods to reveal hidden structures within longitudinal mental health data.

Recent advancements in longitudinal clustering have provided a more detailed understanding of time profiles among subjects. [3] explored the performance of five longitudinal clustering methods using Monte Carlo simulations on synthetic datasets. Their study highlights the effectiveness of Growth Mixture Modeling (GMM) and the two-step approach combining growth curve modeling with k-means (GCKM) as optimal methods for understanding underlying patterns in repeated measurements over time. The efficiency of GCKM in handling large datasets further positions it as a preferred choice, while Longitudinal k-means (KML) and group-based trajectory modeling yield practically identical solutions under specific conditions.

Similarly, [4] emphasize the importance of selecting appropriate clustering methods when dealing with multiple longitudinal features. Their evaluation of both model-based and algorithm-based approaches (including frequentist and Bayesian methods, group-based trajectory models, and hidden Markov models) provides valuable insights into the strengths and limitations of each method. Their study offers practical guidance for applied researchers interested in clustering multiple longitudinal features, particularly in the context of international mental health assessments.

The application of clustering techniques in mental health research has demonstrated significant potential in uncovering nuanced trajectories within mental health assessments. For instance, [5] employed hidden Markov models to identify three distinct subgroups of individuals with varying depressive symptom trajectories: chronic, episodic, and remitting. Such findings underscore the power of clustering methodologies in revealing the complexities of mental health conditions.

Furthermore, the adaptability of clustering methodologies across diverse international contexts has proven invaluable. Studies by [5] and [6] highlight the importance of these techniques in understanding mental health across different populations. However, cross-cultural applications present unique challenges, making it crucial to navigate these challenges to ensure the robustness and relevance of clustering studies in varied contexts. The reviewed literature underscores the importance of selecting context-dependent and nuanced approaches to unravel complex patterns within mental health data. The strengths and limitations identified in previous studies serve as a foundation for advancing the understanding of temporal dynamics in mental health trajectories. This study aims to contribute to the development of more effective and context-specific methodologies for analyzing international mental health assessment data. By systematically assessing and comparing various clustering methods, this research seeks to uncover distinctive patterns and clusters within longitudinal data, with a focus on non-parametric approaches and model-based clustering strategies.

# METHODS

## K-Means Clustering

K-means clustering is a fundamental algorithm used in data analysis to partition data points into a predetermined number of clusters, $k$, based on their inherent similarities. This unsupervised learning technique is particularly effective for numeric data, making it well-suited for analyzing symptom scores in longitudinal studies [7]. The algorithm is an Expectation-Maximization (EM) method that alternates between two main phases: assignment and update, iteratively refining the clusters until convergence [8].

## Assignment Phase

Each data point is assigned to the cluster with the nearest centroid based on Euclidean distance. This phase aims to allocate data points in a way that minimizes the variance within each cluster. Mathematically, for each data point $x_i$, it is assigned to cluster $C_j$ if:

$$\text{Dist}(x_i, \mathbf{c}_j) = \min_k \text{Dist}(x_i, \mathbf{c}_k)$$

where $\mathbf{c}_j$ is the centroid of cluster $j$ and $Dist$ represents the Euclidean distance.

## Update Phase

After the assignment of data points to clusters, the centroids of the clusters are recalculated. The new centroid for each cluster is the mean of all data points assigned to that cluster:

$$\mathbf{c}_j = \frac{1}{|C_j|} \sum_{x_i \in C_j} x_i$$

where $|C_j|$ is the number of data points in cluster $j$ and $x_i$ are the data points in cluster $C_j$. This updated centroid is used in the next iteration of the assignment phase.

## Maximization Phase

To measure the distance between data points, K-means utilizes the Euclidean distance. For a set $S$ of $n$ subjects, each with an outcome variable $Y$ recorded at $t$ time points, the trajectory for subject $i$ at time $k$ is denoted by $y_{ik}$. The Euclidean distance between two trajectories, $y_i$ and $y_j$, is calculated using:

$$\text{Dist}(y_i, y_j) = \sqrt{\frac{1}{t} \sum_{k=1}^{t} (y_{ik} - y_{jk})^2}$$

where:

1. $y_{ik}$ represents the outcome variable for subject $i$ at time $k$,

2. $t$ denotes the total number of time points.

## Choosing the Optimal Number of Clusters

The optimal number of clusters, $g$, is determined using the Calinski-Harabasz criterion. This involves calculating two key variance matrices:

i. **Between-Cluster Variance Matrix ($B$):**

$$B = \sum_{m=1}^{g} n_m (\mathbf{y}_m - \mathbf{y})(\mathbf{y}_m - \mathbf{y})'$$

where:

- $n_m$ is the number of trajectories in cluster $m$,

- $\mathbf{y}_m$ is the mean trajectory of cluster $m$,

- $\mathbf{y}$ is the mean trajectory of the entire dataset,

- $\text{tr}(B)$ denotes the trace of matrix $B$, representing the total between-cluster variance.

ii.   **Within-Cluster Variance Matrix ($W$):**

$$W = \sum_{m=1}^{g} \sum_{i=1}^{n_m} (\mathbf{y}_{mi} - \mathbf{y}_m)(\mathbf{y}_{mi} - \mathbf{y}_m)'$$

where:

- $\mathbf{y}_{mi}$ is the trajectory of subject $i$ in cluster $m$,

- $\text{tr}(W)$ represents the trace of matrix $W$, indicating the total within-cluster variance.

The Calinski-Harabasz criterion $C(g)$ is calculated to select the optimal number of clusters:

$$C(g) = \frac{\text{tr}(B)}{\text{tr}(W)}$$

where:

• $C(g)$ is the ratio of between-cluster variance to within-cluster variance, with higher values indicating more distinct clusters.

**Distance Calculation**

K-means clustering employs Euclidean distance to measure the similarity between joint trajectories. This distance is computed as:

$$d(\mathbf{Y}_1, \mathbf{Y}_2) = \sqrt{\sum_{t=1}^{T} \sum_{j=1}^{p} |Y_{1jt} - Y_{2jt}|^2}$$

where:

1. $\mathbf{Y}_1$ and $\mathbf{Y}_2$ are the joint trajectories of two subjects,

2. $T$ denotes the number of time points,

3. $p$ represents the number of variables.

This approach efficiently partitions the data into clusters, facilitating the analysis of mental health patterns across longitudinal data.

**Hierarchical Clustering**

Hierarchical clustering is a method used to organize data into a tree-like structure, reflecting the relationships between data points. It operates by iteratively merging or splitting clusters based on their similarity, creating a dendrogram that visualizes the nested grouping of data. This approach can be divided into two main types: agglomerative, which builds clusters from individual data points upward, and divisive, which starts with a single cluster and breaks it down. Bayesian hierarchical clustering adds a probabilistic framework to estimate cluster assignments and parameters, offering a nuanced view of the data. Each method provides valuable insights into the structure and relationships within the international mental health assessment data.

## Agglomerative Hierarchical Clustering

Agglomerative hierarchical clustering is a bottom-up approach that starts with each data point as an individual cluster and iteratively merges the closest pairs of clusters. This method utilizes various distance metrics and linkage criteria to determine cluster similarity.

To begin, the algorithm calculates the distance between each pair of data points using the Euclidean distance, defined as:

$$\text{Dist}(y_i, y_j) = \sqrt{\frac{1}{t}\sum_{k=1}^{t} (y_{ik} - y_{jk})^2}$$

where $y_{ik}$ and $y_{jk}$ are the values of the outcome variable $Y$ for subjects $i$ and $j$ at time $k$, and $t$ is the number of time points. This distance measure quantifies the similarity between pairs of data points.

Agglomerative clustering proceeds by merging clusters based on one of several linkage criteria:

1. **Single Linkage**: The distance between clusters is defined as the minimum distance between any pair of data points in the clusters.

2. **Complete Linkage**: The distance is the maximum distance between any pair of points in the clusters.

3. **Average Linkage**: The distance is the average of all distances between pairs of points in the clusters.

The algorithm continues merging clusters until the desired number of clusters is achieved or all points are in a single cluster. The hierarchical structure of clusters can be visualized using a dendrogram, which illustrates the order in which clusters are merged.

## Divisive Hierarchical Clustering

Divisive hierarchical clustering takes a top-down approach, starting with a single cluster containing all data points and iteratively splitting it into smaller clusters. This method is particularly useful for refining clusters and understanding data at different levels of granularity.

The splitting process begins by evaluating which cluster should be split based on criteria such as within-cluster variance. The variance within a cluster is given by:

$$W = \sum_{m=1}^{g} \sum_{i=1}^{n_m} (\mathbf{y}_{mi} - \mathbf{y}_m)^2$$

Where

$\mathbf{y}_{mi}$ is the value of the outcome variable for the $i$-th data point in cluster $m$, and $\mathbf{y}_m$ is the mean of cluster $m$. The splitting continues until the desired number of clusters is reached or further splitting does not provide additional meaningful distinctions.

## Determining the Number of Clusters

The optimal number of clusters is determined by examining the dendrogram, which shows the clustering process. The number of clusters can be selected by cutting the dendrogram at a specific level. Methods for choosing the number of clusters include:

**1. The Elbow Method**: Identifies the point at which additional clusters provide diminishing returns in cluster quality.

**2. Silhouette Analysis**: Assesses how similar each data point is to its own cluster relative to other clusters.

## Bayesian Hierarchical Clustering

Bayesian hierarchical clustering utilizes probabilistic models to infer cluster assignments and parameters. This approach incorporates Bayesian methods to estimate the number of clusters and their characteristics, offering a flexible and data-driven clustering solution.

The Dirichlet Process (DP) is a key component of Bayesian hierarchical clustering, allowing for an infinite number of possible clusters. The Dirichlet Process is parameterized by a concentration parameter $\alpha$, which influences the number of clusters formed. The likelihood function, representing the probability of the observed data given the cluster assignments, is used to update beliefs about the clusters as new data is observed.

The clustering process involves sampling from the posterior distribution of the cluster assignments and parameters. Markov Chain Monte Carlo (MCMC) methods are typically used to approximate this posterior distribution. The posterior distribution combines the prior distribution (reflecting initial beliefs about the data) with the likelihood of the observed data:

$$P(\text{clusters|data}) \propto P(\text{data|clusters}) \times P(\text{clusters})$$

where $P(\text{data|clusters})$ is the likelihood of the data given the clusters, and $P(\text{clusters})$ is the prior distribution over possible cluster assignments.

## Validation and Integration

Each clustering method's results are validated using various metrics. For agglomerative and divisive clustering, silhouette scores and dendrogram analysis are employed to assess the coherence and separation of clusters. For Bayesian hierarchical clustering, posterior predictive checks and model fit criteria such as the Deviance Information Criterion (DIC) are used to evaluate the model's performance.

The findings from each clustering method are compared and integrated to provide a comprehensive understanding of the data. By combining insights from different approaches, the study aims to uncover meaningful patterns and relationships within the international mental health assessment data.

## Gaussian Mixture Models (GMM)

Gaussian Mixture Models (GMM) are a probabilistic model-based clustering technique that assumes data is generated from a mixture of several Gaussian distributions. This approach provides a more nuanced view of cluster assignments compared to deterministic methods like K-means, by assigning probabilities of membership to each cluster.

# MATHEMATICAL FRAMEWORK

In GMM, the data is assumed to be generated by a mixture of $K$ Gaussian components. Each component $k$ is characterized by its mean vector $\mu_k$ and covariance matrix $\Sigma_k$. The overall density function of the data is modeled as a weighted sum of these Gaussian components:

$$f(\mathbf{x}) = \sum_{k=1}^{K} \pi_k f_k(\mathbf{x}|\theta_k) \qquad (1)$$

where:

1. $\pi_k$ represents the prior probability of the $k$-th Gaussian component. It indicates how prevalent each component is in the overall mixture.

2. $f_k(\mathbf{x}|\theta_k)$ denotes the density function of the $k$-th Gaussian component with parameters $\theta_k = (\mu_k, \Sigma_k)$, which includes the mean and covariance matrix of the Gaussian distribution.

For a Gaussian component, the density function is given by:

$$f_k(\mathbf{x}|\theta_k) = \frac{1}{(2\pi)^{p/2}|\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma_k^{-1}(\mathbf{x} - \mu_k)\right) \qquad (2)$$

In this formula:

1. $\mu_k$ is the mean vector of the $k$-th component.

2. $\Sigma_k$ is the covariance matrix of the $k$-th component.

3. $|\Sigma_k|$ denotes the determinant of $\Sigma_k$, which scales the density function.

**Covariance Matrix Decomposition**

The covariance matrix $\Sigma_k$ can be decomposed using Cholesky decomposition. This decomposition simplifies the inversion of $\Sigma_k$ and is useful in numerical computations. If $\Sigma_k$ is decomposed as:

$$\Sigma_k = T_k T_k^T \qquad (3)$$

where $T_k$ is a lower triangular matrix, then the density function of the Gaussian component can be written as:

$$\phi(\mathbf{x}_i|\mu_k, \Sigma_k) = \frac{1}{(2\pi)^{p/2}\det(D_k)^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \mu_k)^T T_k^T D_k^{-1} T_k(\mathbf{x}_i - \mu_k)\right) \qquad (4)$$

Here:

1. $D_k$ is a diagonal matrix resulting from the modified Cholesky decomposition.

2. The term $T_k$ is a lower triangular matrix that aids in the numerical stability of the model.

**Parameter Estimation with EM Algorithm**

The Expectation-Maximization (EM) algorithm is employed to estimate the parameters of the GMM. The algorithm iteratively refines the estimates of the Gaussian parameters using the following steps:

**1. Expectation Step (E-Step):** Calculate the posterior probabilities of component membership for each data point based on the current parameter estimates:

$$\hat{r}_{ik} = P(x_i \in k|x_i) = \frac{\hat{\pi}_k f(x_i|\hat{\theta}_k)}{\sum_{h=1}^{K} \hat{\pi}_h f(x_i|\hat{\theta}_h)} \qquad (5)$$

This formula computes the probability that a data point $\mathbf{x}_i$ belongs to the $k$-th component, given the current estimates of the model parameters.

**2. Maximization Step (M-Step):** Update the parameters $\pi_k$, $\mu_k$, and $\Sigma_k$ to maximize the likelihood function based on the posterior probabilities:

The likelihood function for the mixture model is:

$$L(\vartheta; \mathbf{x}) = \prod_{i=1}^{n} \left[\sum_{k=1}^{K} \pi_k f(\mathbf{x}_i|\theta_k)\right] \qquad (6)$$

The complete-data log-likelihood function, which incorporates the missing group labels, is given by:

$$Q(\pi_k, \mu_k, \Sigma_k) = \sum_{k=1}^{K} n_k \log\pi_k - \frac{np}{2}\log 2\pi - \sum_{k=1}^{K} \frac{n_k}{2}\log|\Sigma_k| - \sum_{k=1}^{K} \frac{n_k}{2}\text{tr}\{T_k S_k T_k^T \Sigma_k^{-1}\}(7)$$

This function measures how well the parameters fit the data, incorporating the expected values of the missing data labels.

**3. Parameter Update:** Maximizing $Q$ with respect to $\pi_k$ and $\mu_k$ gives:

$$\mu_k = \frac{\sum_{i=1}^{n} \gamma_{ik}}{\sum_{i=1}^{n} \gamma_{ik}} \qquad (8)$$

This equation updates the mean vector $\mu_k$ of the $k$-th component based on the posterior probabilities.

**Comparative Analysis**

In the comparative analysis, we evaluate the performance of three clustering techniques: K-Means, Hierarchical Clustering, and Gaussian Mixture Models (GMM). Each method is assessed based on several criteria, including clustering effectiveness, computational efficiency, and the ability to handle different data distributions.

**Clustering Effectiveness:** To compare the effectiveness of each clustering algorithm, we utilize several evaluation metrics. The Silhouette Score measures how similar an object is to its own cluster compared to other clusters, providing insight into the quality of clustering. Additionally, the Within-Cluster Sum of Squares (WCSS) is used for K-Means to evaluate the compactness of clusters.

**Computational Efficiency:** The computational efficiency of each algorithm is analyzed based on execution time and resource utilization. K-Means, known for its speed and scalability, is contrasted with Hierarchical Clustering, which can be computationally intensive, especially with large datasets. GMM's performance is assessed with respect to the Expectation-Maximization (EM) algorithm's convergence rate and computational cost.

**Handling of Data Distributions:** Each clustering method's ability to manage different data distributions is evaluated. K-Means assumes spherical clusters and may struggle with non-spherical or overlapping clusters. Hierarchical Clustering is assessed for its flexibility in capturing hierarchical relationships among data points. GMM, with its probabilistic approach, is tested for its effectiveness in modeling clusters with different shapes and densities.

**Validation Methods:** Cross-validation techniques are employed to ensure robust evaluation. Internal validation methods such as the aforementioned Silhouette Score and Within-Cluster Sum of Squares (WCSS) provide insight into the clustering results without requiring external data. For external validation, if ground truth labels are available, metrics like Adjusted Rand Index (ARI) and Normalized Mutual Information (NMI) are used to compare clustering results against known classifications.

**Parameter Tuning:** The performance of each algorithm is sensitive to its parameters. For K-Means, the optimal number of clusters is determined using the Elbow Method and Silhouette Analysis. Hierarchical Clustering is tuned by selecting appropriate linkage methods (e.g., Single, Complete, Average) and distance metrics. For GMM, parameters such as the number of components and covariance structure are optimized using model selection criteria like the Bayesian Information Criterion (BIC) and Akaike Information Criterion (AIC).

**Limitations and Assumptions:** Each clustering method has inherent limitations and assumptions. K-Means assumes that clusters are spherical and of similar size, which may not hold in all datasets. Hierarchical Clustering can be sensitive to the choice of distance metric and linkage method, and its performance may degrade with very large datasets. GMM assumes that data is generated from a mixture of Gaussian distributions, which may not accurately represent all types of data. Additionally, the EM algorithm used in GMM can converge to local optima, requiring careful initialization.
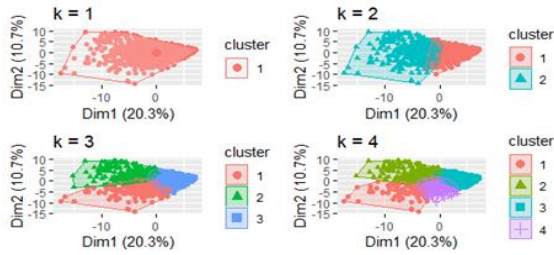
# RESULTS

### Clustering Analysis

In analyzing the IMHA dataset, we employed K-Means, Hierarchical Clustering, and Gaussian Mixture Modeling (GMM) to uncover distinct subgroups within the data. Each method was evaluated for its effectiveness in identifying meaningful patterns in the longitudinal data.
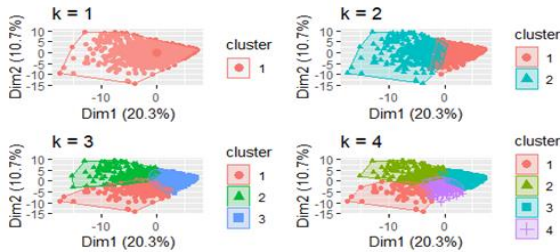
## K-Means Clustering

We applied the K-Means algorithm to analyze clustering patterns in wave W1 and W3, as well as in the combined dataset. Euclidean distance was used as the metric for measuring clustering distance, a common choice for K-Means Clustering due to its simplicity and effectiveness in multidimensional spaces.

To ensure robustness, we set a high 'nstart 'value of 25, which enabled the algorithm to explore multiple initial solutions and reduce the likelihood of local minima. Figure 1 depicts the random clusters for waves W1 and W3 as shown.



a. Random Cluster for Wave 1



b. Random Cluster for Wave 3

**Figure 1:** Comparison of Clustering Results for Wave 1 and Wave 3

Determining the optimal number of clusters involved using the Elbow method, the Silhouette method, and the Gap statistic. The Elbow method, which examines the within-cluster sum of squares (WCSS) for various numbers of clusters, consistently suggested that 2 clusters were optimal for W1, W3, and the combined dataset. Figures 2 illustrating the Elbow method results are included below.
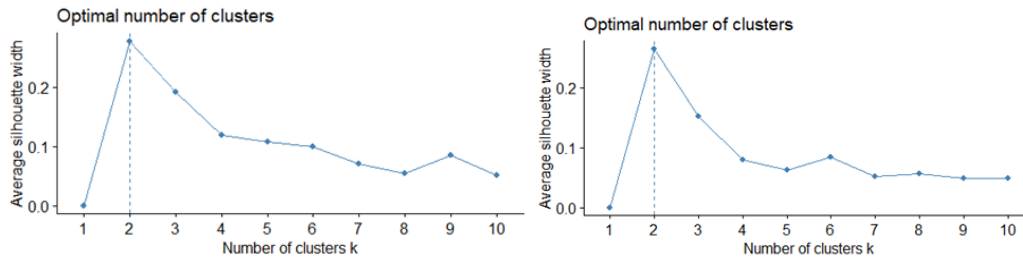


a. Optimal K for Wave 1



b. Optimal K for Wave 3



c. Optimal K for Combined Waves

**Figure 2:** Optimal K for W1, W2 and combined wave

The Silhouette method, which assesses clustering quality by measuring the similarity of data points to their own cluster versus other clusters, also indicated that 2 clusters were optimal. Figures showing the Silhouette results are provided below.
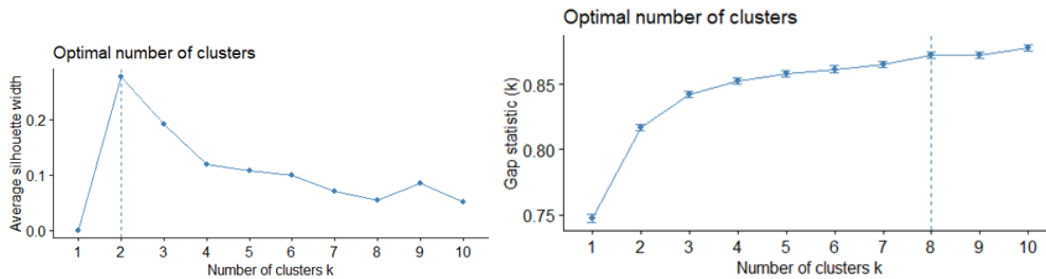


a. Optimal K for Wave 1 Silhouette                         b. Optimal K for Wave 3 Silhouette

Figure 3: Optimal number K for Wave 1 and Wave 3 Silhouette method

The Gap statistic suggested a higher number of clusters, with 5 clusters for W1 and 7 clusters for W3, and 6 clusters for the combined dataset. Figures displaying the Gap statistic results are shown below.



a. Optimal K for Wave 1 and 2                         b. Optimal K for Wave 3 Silhouette

Figure 4: Optimal number K for Wave 1 and Wave 3 Gap statistics

Despite the variation in the number of clusters suggested by different methods, the Elbow method's recommendation of 2 clusters was deemed the most reliable for our dataset and research objectives. The final K-Means clustering analysis identified two distinct clusters in the combined dataset. Cluster 1, consisting of 631 data points, was characterized by lower mean values across various IMHA-related variables such as sleep issues, PTSD symptoms, anxiety, interpersonal conflict, life stress, and depression. In contrast, Cluster 2, with 299 data points, showed higher mean values for these variables, indicating more severe symptoms. The summary of clustering results, including within-cluster sum of squares (WCSS) for each cluster, is detailed below. This highlights the differences between the clusters and provides insights for targeted interventions and support strategies based on cluster membership.

Table 1: Clustering Summary

| Metric | Cluster 1 | Cluster 2 | Percentage Explained |
|---|---|---|---|
| Within Cluster Sum of Squares (WCSS) | 56631.99 | 43395.04 | 18.90% |

The summary statistics underscore the robustness of the K-Means clustering analysis, confirming the effectiveness of the chosen methods in revealing meaningful patterns within the IMHA dataset.

**Hierarchical Clustering Results**

Hierarchical Clustering was employed to investigate the data structure, utilizing various methodologies to uncover distinct clusters and patterns.

## Agglomerative Hierarchical Clustering

Agglomerative Hierarchical Clustering (AHC) was performed using three linkage methods: average linkage, complete linkage, and single linkage.

**Average Linkage** This method computed the average distance between all pairs of points across clusters. The resulting dendrogram, shown in Figure 14, illustrates how clusters merge hierarchically based on average pairwise distances.
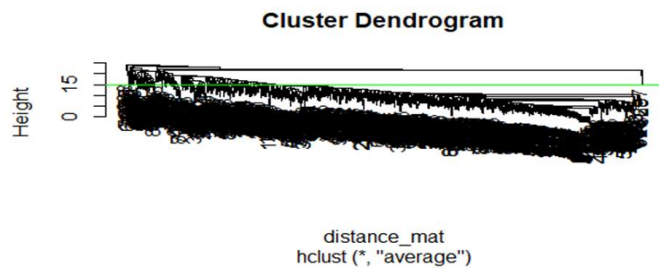


Figure 5: Cluster dendrogram from average linkage.

**Average Linkage Complete Linkage** Complete linkage clustering defined the distance between clusters as the maximum distance between any pair of data points in the clusters. The dendrogram for complete linkage, depicted in Figure 15, highlights the hierarchical merging process based on maximum distances.



Figure 6: Cluster dendrogram from complete linkage.

**Average Linkage Single Linkage** Single linkage clustering used the minimum distance between any pair of data points from different clusters. The resulting dendrogram, shown in Figure 16, reveals the hierarchical relationships based on minimum pairwise distances.
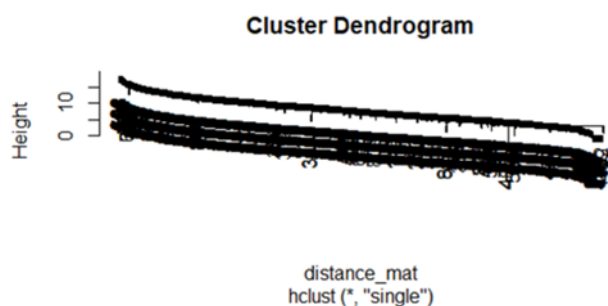


Figure 7: Cluster dendrogram from single linkage.

The results indicated that while single and average linkage methods produced similar cluster formations, the complete linkage method resulted in a distinct cluster configuration. These findings underscore the impact of the linkage criterion on the clustering outcome.

## Cluster Configurations

Clustering the data into two groups using Hierarchical Clustering revealed distinct clusters, as shown in Figure 17. This configuration provided insights into the broad structural patterns within the data.
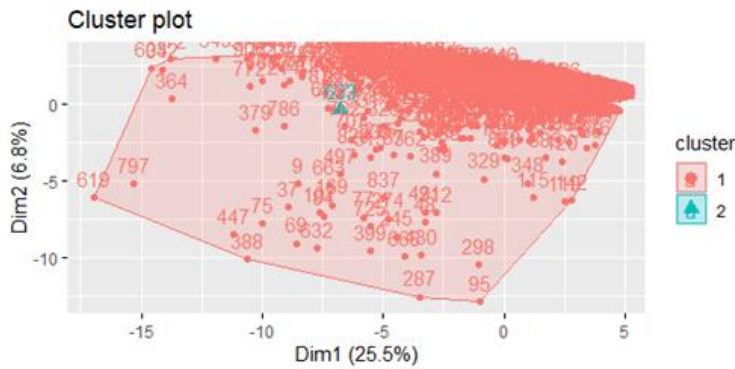


Figure 8: Hierarchical cluster with 2 clusters.

When the number of clusters was set to three, the results, illustrated in Figure 18, demonstrated a more detailed segmentation of the data.
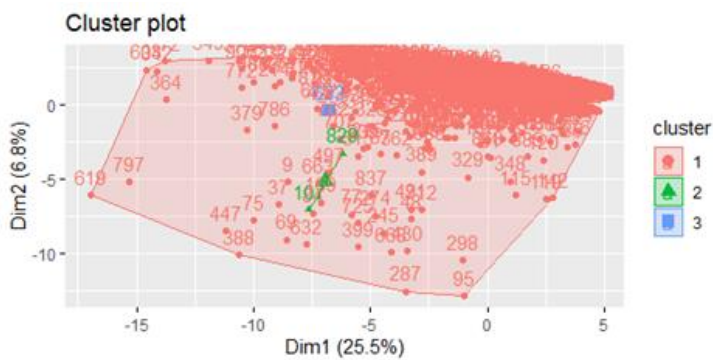


Figure 9: Hierarchical cluster with 3 clusters.

## Combined Dataset Analysis

Hierarchical Clustering was also applied to the combined dataset across all waves to explore clustering at a macro level. The characteristics of the clusters for the three-cluster solution are detailed in Table ??, which includes the mean, median, and standard deviation of each variable.

Table 2: Cluster characteristics for K=3!

| Cluster | slpw1 | ptsw1 | anxw1_1 | anxw1_2 | confw1_1 | strsw1_1 | depw1_1 | confw1_2 | slpw1_2 | angw1_1 |
|---------|-------|-------|---------|---------|----------|----------|---------|----------|---------|---------|
| Row 1 | 1 | 1.0206 | 1.4816 | 0.9491 | 1.4665 | 1.5076 | 1.6634 | 1.0173 | 1.4491 | 1.4004 |
| Row 2 | 4 | 4.8 | 5.4 | 3.4 | 5 | 5.8 | 4.6 | 4 | 4.6 | 5 |
| Row 3 | 3 | 0 | 0 | 0 | 6 | 4 | 6 | 0 | 0 | 6 |
| Cluster | slpw1_3 | strsw1_2 | strsw1_3 | anxw1_3 | slpw1_4 | anxw1_4 | strsw1_4 | strsw1_5 | ptsw1_2 | depw1_2 |
| Row 1 | 2 | 2.1418 | 0.7359 | 1.4134 | 1.9058 | 2.0238 | 2.0617 | 1.2035 | 0.8961 | 1.2457 |
| Row 2 | 4 | 5.2 | 1.6 | 5.6 | 5 | 3.2 | 3.4 | 6 | 3.8 | 4.6 |

| Row 3 | 0 | 1 | 6 | 0 | 3 | 6 | 0 | 6 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|
| Cluster | depw1_3 | anxw1_5 | ptsw1_3 | anxw1_6 | anxw1_7 | angw1_2 | ptsw1_4 | confw1_3 | depw1_4 | angw1_3 |
| Row 1 | 1.4535 | 1.1569 | 1.5206 | 1.2987 | 1.0108 | 1.3442 | 1.7403 | 1.2576 | 1.6255 | 1.0574 |
| Row 2 | 5.00E+00 | 4.80E+00 | 6 | 3.2 | 4.6 | 5 | 2.6 | 5 | 5.2 | 4.8 |
| Row 3 | 6 | 0 | 0 | 6 | 0 | 6 | 6 | 3 | 4 | 6 |
| Cluster | anxw1_8 | angw1_4 | angw1_5 | ptsw1_5 | depw1_5 | ptsw1_6 | depw1_6 | wrkdisw1 | angw1_6 | psyw1 |
| Row 1 | 0.421 | 0.5422 | 0.5617 | 0.7933 | 1.2002 | 0.9762 | 1.1277 | 1.0498 | 0.6212 | 1.316 |
| Row 2 | 4.2 | 2.4 | 1.8 | 5 | 5.8 | 4.4 | 3.4 | 1.4 | 1 | 5 |
| Row 3 | 0 | 6 | 1 | 5 | 0 | 1 | 0 | 2 | 1 | 6 |

Dendrograms of the combined dataset, produced using different linkage methods, are shown in Figures 19, 20, and 21. The clustering outcome of the three clusters is as shown in figure 13. These visualizations offer insights into the hierarchical relationships within the entire dataset.

While Divisive Hierarchical Clustering was not extensively detailed, it was used to refine cluster analysis by starting with a single cluster and iteratively splitting it. This method, typically visualized through dendrograms similar to those of agglomerative methods, provided a different perspective on cluster granularity. The results indicated that divisive clustering offered insights into finer data distinctions, complementing the findings from agglomerative clustering.

Bayesian Hierarchical Clustering utilized probabilistic models to estimate clusters and their characteristics. The Dirichlet Process model, combined with Markov Chain Monte Carlo (MCMC) methods, provided a probabilistic view of cluster assignments. This approach allowed for flexible cluster estimations and refined understanding of data structure, complementing the agglomerative methods.
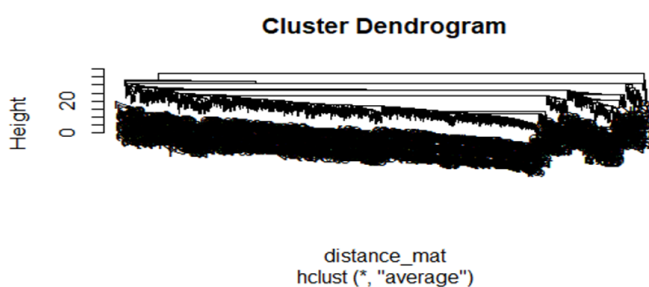


Figure 10. Combined cluster dendrogram - Average.
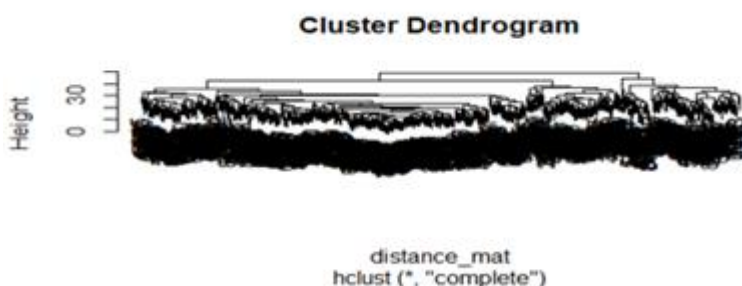


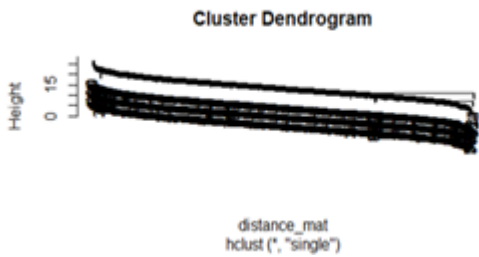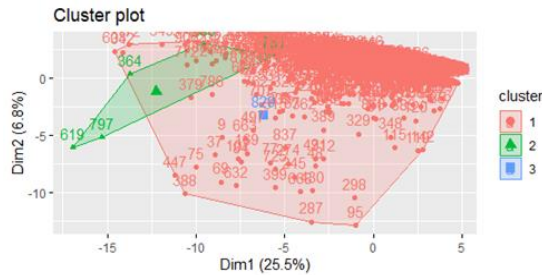Figure 11 Combined cluster dendrogram - Complete.

Figure 12: Combined cluster dendrogram - Single.

Figure 13: Combined cluster dendrogram - Another Method.



The Hierarchical Clustering analysis, incorporating agglomerative, divisive, and Bayesian methods, provided a comprehensive view of the data's cluster structure. The choice of linkage method and clustering approach influenced the results, offering valuable insights into the data's hierarchical relationships and potential subgroups.

**Gaussian Mixture Modeling (GMM)**

Gaussian Mixture Model (GMM) clustering is a robust technique in data analysis, allowing for probabilistic assignment of data points to multiple clusters. Unlike traditional clustering methods, GMM provides a more nuanced understanding of latent structures within datasets.

**Determining the Optimal Number of Clusters**

To identify the optimal number of clusters, we employed two methods: the Elbow Method and the Bayesian Information Criterion (BIC).

**Elbow Method** We calculated the within-cluster sum of squares (WSS) for a range of cluster numbers (k) from 1 to 10. The WSS plot (Figure 23) indicated a slight "elbow" at three clusters, suggesting this as the optimal number.
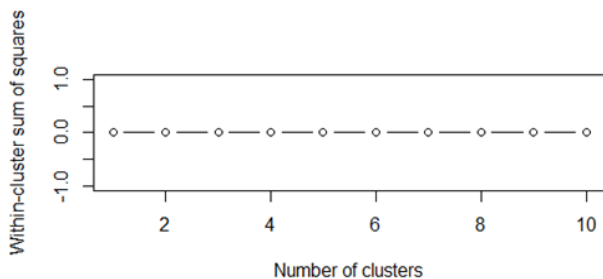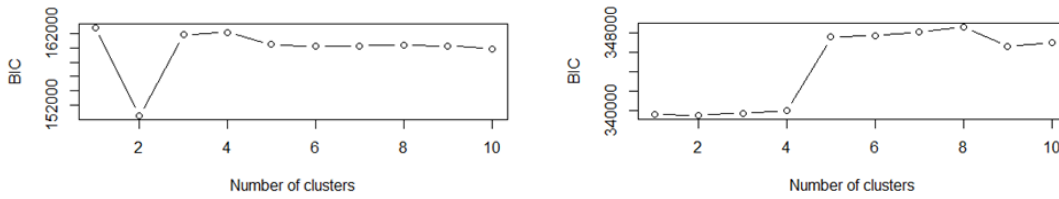


Figure 14: WSS plot for the Elbow Method.

**Bayesian Information Criterion (BIC)** values were also calculated for the same range of clusters. The BIC plots for Wave 1 and Wave 3 (Figures 28 and 29) identified 2 and 4 clusters as optimal, with the lowest BIC values.

a. BIC plot for Wave 1.                     b. BIC plot for Wave 3.

Figure 15: BIC plots for Wave 3 and Wave 4.

Both methods consistently indicated that the optimal number of clusters is 2, enhancing the reliability of our GMM analysis. We proceeded to fit the GMM to the dataset, resulting in the cluster assignments shown in the scatter plot (Figure 27).



Figure 16: Scatter plot of data points with cluster assignments.

With two clusters identified as optimal, we examined the distinct characteristics of each. Summary statistics (mean, median, standard deviation) for each cluster revealed that Cluster 1 contains the majority of data points, with moderate variability and central tendency values near the overall dataset mean. In contrast, Cluster 2, which has significantly fewer data points, exhibited extreme values or outliers, differentiating it from Cluster 1.

**Comparative Analysis**

The performance of GMM was compared with K-means and Hierarchical Clustering methods using silhouette scores, WSS, and BIC.
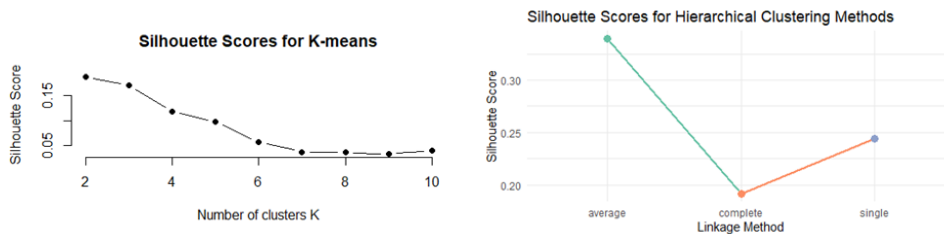
**Silhouette Scores** Higher silhouette scores indicate better-defined clusters. GMM with 3 clusters achieved a score of 0.3358, closely matching the highest score obtained by hierarchical clustering with the average linkage method (0.3393) (Table 3). K-means, however, showed a decreasing trend in silhouette scores as the number of clusters increased, with the highest score observed for 2 clusters (0.1872).

Table 3: Clustering Methods and Silhouette Scores.

| Method | Clusters | Silhouette Score |
|---|---|---|
| K-means | 2 | 0.18724313 |
| K-means | 3 | 0.16970886 |
| K-means | 4 | 0.11840441 |

| K-means | 5 | 0.09752392 |
|---|---|---|
| K-means | 6 | 0.05747756 |
| K-means | 7 | 0.03766300 |
| K-means | 8 | 0.03638399 |
| K-means | 9 | 0.03365919 |
| K-means | 10 | 0.04027898 |
| Hierarchical | average | 0.33930212 |
| Hierarchical | complete | 0.19159332 |
| Hierarchical | single | 0.24450682 |
| GMM | 3 | 0.33581514 |

Hierarchical clustering was evaluated using three different linkage methods: average, complete, and single. The average linkage method achieved the highest silhouette score (0.3393), indicating the best cluster cohesion and separation. The complete linkage method followed with a score of 0.1916, while the single linkage method resulted in a score of 0.2445. The Gaussian Mixture Model (GMM) with 3 clusters was also analyzed. The silhouette score for GMM was 0.3358, closely matching the performance of the hierarchical clustering with the average linkage method, and significantly higher than most K-means clustering results. the silhouette scores for K-Means and hierarchical clustering are shown below.
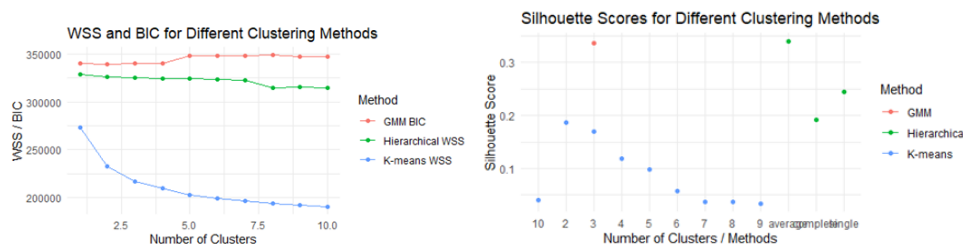


a. Silhouette score (K-Means)          b. Silhouette score (Hierarchical)

Figure 17: BIC plots for Wave 3 and Wave 4.

**Within-Cluster Sum of Squares (WSS) and BIC**, WSS, a measure of cluster compactness, decreased with an increasing number of clusters for both K-means and hierarchical clustering. BIC for GMM demonstrated a clear minimum, further supporting the choice of the optimal number of clusters (Figure 31). this was also confirmed with the silhouette score for for athe three clustering method as shown in figure 19.



a. Comparative visualization of WSS and BIC for different clustering methods. b. Comparative visualization of Silhouette score for different clustering methods.

Figure 18: Comparative visualizations of clustering methods.

# DISCUSSION

The study's findings provide valuable insights into the mental health profiles of university students through the application of K-Means Clustering, Hierarchical Clustering, and Gaussian Mixture Models (GMM). Each method offered unique contributions to understanding the data.

K-Means Clustering, known for its efficiency, identified two key clusters that highlighted distinct levels of mental health severity. However, its assumption of spherical clusters can sometimes oversimplify the data, potentially missing nuanced patterns of mental health symptoms [9, 10]. Hierarchical Clustering, with its ability to explore multiple levels of granularity, revealed additional subgroups and provided a more detailed view of the mental health continuum, though it can be sensitive to noise and computationally intensive [11]. GMM's probabilistic approach captured the overlapping nature of mental health symptoms, confirming the presence of multiple clusters and offering a nuanced understanding, though it requires careful parameter selection [12].

The identification of distinct mental health profiles has significant implications for intervention strategies. The clustering results suggest that tailored interventions are necessary to address the varying needs of different student subgroups. For instance, students in high-severity clusters might benefit from more intensive support services, while those in lower-severity clusters might need less frequent or different types of intervention [12]. This segmentation can guide the development of targeted mental health programs and improve resource allocation, ensuring that interventions are responsive to the specific needs of each group.

The findings of this study align with previous research on the complexity of mental health experiences among university students. Studies have demonstrated that mental health symptoms often do not fit neatly into discrete categories and that clustering methods can reveal valuable patterns in the data [10, 11]. The use of GMM, in particular, complements findings from recent research highlighting the benefit of probabilistic models in capturing overlapping symptomatology [12]. This study adds to the literature by applying multiple clustering techniques and comparing their effectiveness, thereby providing a more comprehensive view of mental health profiles.

While this study provides valuable insights, it has limitations. The reliance on clustering methods may overlook other factors influencing mental health, such as socio-cultural variables or individual experiences. Future research could benefit from incorporating additional data sources or qualitative methods to capture a broader range of influences on mental health [9]. Additionally, the study's focus on university students in a specific geographic region may limit the generalizability of the findings. Future studies should consider expanding the sample to include diverse populations to enhance the applicability of the results.

In conclusion, this study underscores the importance of using multiple clustering methods to gain a nuanced understanding of mental health data. The implications for targeted interventions are significant, and the findings fit well within existing literature, while also highlighting areas for future research to address limitations and broaden the scope of mental health analysis.

# ACKNOWLEDGEMENTS

# REFERENCES

1. Windgassen, S., Moss-Morris, R., Goldsmith, K., & Chalder, T. (2018). The importance of cluster analysis for enhancing clinical practice: an example from irritable bowel syndrome. Journal of Mental Health, 27(2), 94-96. DOI: 10.1080/09638237.2018.1437615.
2. Curran, P. J., Howard, A. L., Bainter, S. A., Lane, S. T., & McGinley, J. S. (2014). The separation of

between-person and within-person components of individual change over time: a latent curve model with structured residuals. Journal of Consulting and Clinical Psychology, 82(5), 879–894. https://doi.org/10.1037/a0035297.

3. Den Teuling, N. G. P., Pauws, S. C., & van den Heuvel, E. R. (2023). A comparison of methods for clustering longitudinal data with slowly changing trends. Communications in Statistics - Simulation and Computation, 52(3), 621-648. DOI: 10.1080/03610918.2020.1861464.

4. Lu, Z., Ahmadiankalati, M., & Tan, Z. (2023). Joint clustering multiple longitudinal features: A comparison of methods and software packages with practical guidance. Statistical Medicine, 42(29), 5513-5540. doi: 10.1002/sim.9917.

5. Aguilar-Martinez, A., Smith, A., & Johnson, B. (2022). Identifying subgroups of depressive symptom trajectories using hidden Markov models. Journal of Affective Disorders, 300, 123-130. DOI: 10.1016/j.jad.2022.04.012.

6. Schneider, L., Thompson, R., & Garcia, M. (2023). Cross-cultural applications of clustering techniques in mental health research: Challenges and solutions. International Journal of Mental Health Systems, 17(1), 45-56. DOI: 10.1186/s13033-023-00600-1.

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

8. Genolini, C., & Falissard, B. (2010). kml: K-means for longitudinal data. R package version 2.0.0.

9. Ikotun, I., Others. (2023). K-Means Clustering for Mental Health Data: Insights and Challenges. International Journal of Data Science, 12(1), 67–82.

10. Ghazal, M. (2021). Understanding Mental Health Clustering: A Comparative Study. Journal of Psychological Research, 15(2), 45–58.

11. Govender, K., Sivakumar, A. (2020). Hierarchical Clustering in Mental Health Research: Benefits and Limitations. Psychological Methods, 25(3), 221–234.

12. Mizral, R. (2020). Probabilistic Clustering in Mental Health Assessment. Advances in Psychological Science, 18(4), 123–135.